

The Formal Synthesis of Control Signals for  
Systolic Arrays

Jingling Xue

PhD  
University of Edinburgh  
1992





# Abstract

The distinguishing features characteristic of systolic arrays are synchrony of computations, local and regular connections between processors and massive decentralised parallelism. The potential of the systolic array lies in its suitability for VLSI fabrication and its practicality for a variety of application areas such as signal or image processing and numeric analysis. With the increasing possibilities promised by advances in VLSI technology and computer architecture, more and more complex problems are now solvable by systolic arrays.

This thesis describes a systematic method for the synthesis of control signals for systolic arrays that are realised in hardware. Control signals ensure that the right computations are executed at the right processors at the right time. The proposed method applies for iterative algorithms defined over a domain that can be expressed by a convex set of integer coordinates. Algorithms that can be implemented as systolic arrays can be expressed this way; a large subclass can be phrased as affine (or uniform) recurrence equations in the functional style and as nested loops in the imperative style. The synthesis of control signals from a program specification is a process of program transformation and construction. The basic idea is to replace the domain predicates in the initial program specification which constitute the abstract specification of control signals by a system of uniform recurrence equations by means of data pipelining. Then, systolic arrays with a description of both data and control signals can be obtained by a direct application of the standard space-time mapping technique.



# Acknowledgements

I would like to thank first of all my supervisor, Christian Lengauer, for his help, encouragement and advice in the development of this thesis. His concern for clarity and mathematical elegance has greatly influenced my work. Also, I would like to thank Patrice Quinton and Björn Lisper for helpful discussions over the electronic network. Finally, I would like to thank my wife Lili for her support and patience in making this thesis possible.

My research work was supported financially by an Overseas Research Studentship, a University of Edinburgh Studentship, and the Science and Engineering Research Council, Grant no. GR/G55457.



# Table of Contents

<b>1. Introduction</b>	<b>1</b>
1.1 Outline of the Thesis . . . . .	6
1.2 Notation and Terminology . . . . .	8
<b>2. Data Flow Synthesis for <math>(n-1)</math>-Dimensional Systolic Arrays</b>	<b>15</b>
2.1 Introductory Remarks . . . . .	15
2.2 Affine and Uniform Recurrence Equations . . . . .	16
2.3 The Systolic Array Model . . . . .	23
2.4 The Space-Time Mapping . . . . .	24
2.5 Input and Output Extension . . . . .	29
2.6 Uniformisation . . . . .	33
2.7 Conclusion . . . . .	35
<b>3. Control Flow Synthesis for <math>(n-1)</math>-Dimensional Systolic Arrays</b>	<b>38</b>
3.1 Introductory Remarks . . . . .	38
3.2 The Requirements on Control Flow . . . . .	40
3.3 The Synthesis of Control Flow . . . . .	42
3.3.1 The Computation Control Flow . . . . .	48
3.3.2 The Propagation Control Flow . . . . .	58



3.4	The Optimisation of Control Flow . . . . .	61
3.5	The Extension of the Space-Time Mapping . . . . .	64
3.6	Related Work . . . . .	66
3.7	Examples . . . . .	67
3.7.1	Dynamic Programming . . . . .	67
3.7.2	LU-Decomposition . . . . .	74
3.8	Conclusion . . . . .	78
<b>4.</b>	<b>Data Flow Synthesis for One-Dimensional Systolic Arrays</b>	<b>82</b>
4.1	Introductory Remarks . . . . .	82
4.2	One-Dimensional Systolic Array Models . . . . .	83
4.3	The Mapping Conditions . . . . .	88
4.4	The Space-Time Diagram . . . . .	92
4.5	The Extension of the Index Space . . . . .	96
4.6	A Synthesis Procedure . . . . .	100
4.7	Related Work . . . . .	104
4.8	Conclusion . . . . .	105
<b>5.</b>	<b>Control Flow Synthesis for One-Dimensional Systolic Arrays</b>	<b>108</b>
5.1	Introductory Remarks . . . . .	108
5.2	The Synthesis of Propagation Control Flow . . . . .	112
5.3	The PCUREs for Three-Dimensional UREs . . . . .	115
5.3.1	The Evolution Control Flow . . . . .	115
5.3.2	The Initialisation Control Flow . . . . .	121
5.3.3	The Termination Control Flow . . . . .	126



5.3.4	The Space-Time Mapping . . . . .	131
5.4	The PCUREs for $n$ -Dimensional UREs . . . . .	133
5.5	Related Work . . . . .	139
5.6	Examples . . . . .	139
5.6.1	Matrix Product . . . . .	139
5.6.2	Dynamic Programming . . . . .	141
5.7	Conclusion . . . . .	142
<b>6.</b>	<b>The Elimination of Propagation Control Flow</b>	<b>145</b>
6.1	Introductory Remarks . . . . .	145
6.2	Example: 1-D Convolution . . . . .	147
6.3	Systolic Arrays of $n-1$ Dimensions . . . . .	151
6.4	Systolic Arrays of One Dimension . . . . .	154
6.5	Conclusion . . . . .	155
<b>7.</b>	<b>The Loading, Recovery and Access of Stationary Data</b>	<b>156</b>
7.1	Introductory Remarks . . . . .	156
7.2	The Loading and Recovery of Stationary Data . . . . .	157
7.3	The Access of Stationary Data . . . . .	163
7.4	Conclusion . . . . .	167
<b>8.</b>	<b>Conclusion</b>	<b>168</b>
<b>A.</b>	<b>The Normalisation of Domain Predicates</b>	<b>170</b>



## List of Tables

5-1	The host program for three-dimensional UREs (to be refined). . . .	119
5-2	The cell program for three-dimensional UREs. . . . .	120
5-3	The specification of initialisation control variables. . . . .	123
5-4	The host program for three-dimensional UREs (the final version). . . .	129
5-5	The host program for $n$ -dimensional UREs ( $n > 2$ ). . . . .	137
5-6	The cell program for $n$ -dimensional UREs ( $n > 2$ ). . . . .	138
A-1	The abstract syntax of the domain predicates. . . . .	171



## List of Figures

1-1	The synthesis of control signals. . . . .	4
1-2	Illustration of some concepts in $\mathbb{R}^2$ in convex analysis. . . . .	13
2-1	The data dependence graph of the AREs for matrix product . . . .	22
2-2	The data dependence graph of the UREs for matrix product . . . .	22
2-3	S. Y. Kung's systolic array for matrix product. . . . .	28
2-4	Kung-Leiserson's systolic array for matrix product. . . . .	31
2-5	Part of the extended data dependence graph for matrix product. .	32
3-1	The construction of the quotient set $\Phi/\mathbb{T}$ . . . . .	43
3-2	Domain predicate replacement. . . . .	52
3-3	The optimisation of pipelined UREs. . . . .	63
3-4	The data dependence graph for dynamic programming . . . . .	69
3-5	Four systolic arrays for dynamic programming. . . . .	73
3-6	The data dependence graph for LU-decomposition . . . . .	76
3-7	Two systolic arrays for LU-decomposition. . . . .	79
4-1	Two one-dimensional systolic array models. . . . .	84
4-2	A systolic array in which the neighbouring communication cannot be enforced. . . . .	90
4-3	The space-time diagrams for one-dimensional arrays. . . . .	93



4-4	Part of the extended index space for matrix product. . . . .	97
5-1	The space-time diagram for LU-decomposition. . . . .	111
5-2	The notation for specifying the cell program. . . . .	113
5-3	The evolution of the evolution control variable . . . . .	117
5-4	The evolution control variable for LU-decomposition. . . . .	121
5-5	The interplay between the evolution control variable and the initialisation control variables. . . . .	122
5-6	The initialisation control variables for LU-decomposition. . . . .	127
5-7	The interplay between the evolution control variable and the termination control variables. . . . .	128
5-8	The termination control variables for LU-decomposition. . . . .	130
5-9	A hierarchical construction of the PCUREs for $n$ -dimensional UREs. . . . .	135
5-10	Adaptation of the index space for the construction of the PCUREs. . . . .	143
6-1	The elimination of the propagation control signals for 1-D convolution. . . . .	150
7-1	The data dependence graph of the UREs with a load-recovery variable for matrix product. . . . .	161
7-2	The systolic array for matrix product that handles the loading and recovery of stationary data. . . . .	162
7-3	The data dependence graph of the UREs with an address variable for dynamic programming. . . . .	166
A-1	The points $(i, j, 1)$ satisfying predicate $2\lfloor i/2 \rfloor - j = 0$ . . . . .	176



# Chapter 1

## Introduction

The potential of VLSI for MIMD parallelism was first recognised by H. T. Kung and C. E. Leiserson [37], when they introduced the term *systolic array* to describe a processor network that was suitably restrained to meet the requirements of VLSI technology. The characteristic features of a systolic array are synchronous and decentralised parallelism and local and regular interconnections between processors (or cells). The proceedings of recent yearly or bi-yearly conferences and workshops on regular array processors attest to the variety of areas in which systolic solutions have been proposed since then, e.g., numerical analysis, signal and image processing, pattern recognition and combinatorial theory. This suggests that systolic computation has an important, enduring rôle to play in concurrent computation.

The special appeal of systolic arrays is that they can be derived systematically, and often mechanically, by provably correct and (in a sense) optimal synthesis methods. These methods transform algorithmic descriptions that do not specify concurrency or communication – usually functional programs or imperative programs – into functions that distribute the computations prescribed by the program over space and time. This process is called *systolic design*. The challenge, as stated by H. T. Kung [35], is to ensure that *the right data arrive at the right cells at the right time*. The distribution functions essentially describe the velocities and distribution of data, i.e., the flow of data over space and time. They can then be refined further and translated into a description for fabrication of a VLSI chip or into a distributed program for execution on a programmable processor array.



There have been a plethora of synthesis methods proposed in the literature. These methods focus mainly on the synthesis of data flow. They are based on geometry, linear algebra and convex analysis. The basic idea to describe a systolic array by two distribution functions: a timing function that specifies the temporal distribution of the computations and an allocation function that specifies their spatial distribution such that concurrent computations are allocated to different processors. The combination of the timing function and allocation function is called a *space-time mapping* or *space-time transformation*. As has been proved by S. K. Rao [64], a systolic array is a program resulting from an application of the space-time mapping to the initial program specification. Of course, not every space-time mapping describes a systolic array. A valid space-time mapping must preserve the behaviour of the source program in some sense and must respect the constraints imposed on resource such as channel connections in systolic arrays. As pointed out by P. Quinton [55], the space-time mapping also needs to meet conditions that allow a full mechanisation of the underlying synthesis method.

The timing and allocation function constitute a complete description for fabrication of a VLSI chip if every processor in the systolic array performs the same computation at all time steps. The description is incomplete, however, if some processor performs different computations at different time steps. In this case, a control mechanism is called for that instructs the processors in the systolic array when to perform which computation. L. J. Guibas, H. T. Kung and C. D. Thompson [26] suggested a decade ago that this control mechanism may be implemented by communicating control signals throughout the systolic array in much the same way as the data are communicated. Like in the synthesis of data flow, the challenge in the synthesis of control flow is to ensure that *the right control signals arrive at the right cells at the right time*. So far, the systematic synthesis of control flow has not received adequate attention. In most systolic solutions that have been proposed in the literature, the derivation of the control flow is conducted after that of the data flow and in an informal and problem-specific manner.

With the advent of powerful VLSI design and fabrication techniques and of programmable processor networks, the development of formal methods for the

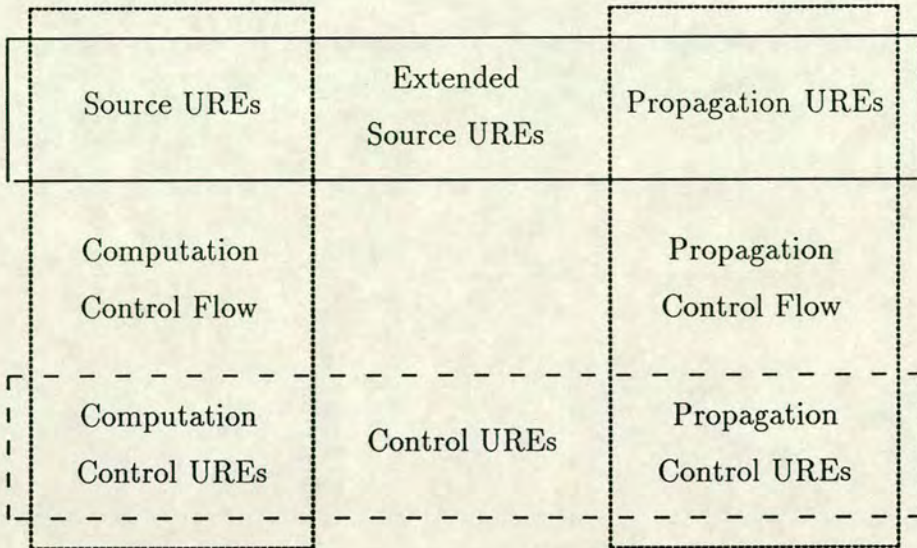


synthesis of control signals has become increasingly important. First, complex applications, like the algebraic path problem [8,18,67], that were not considered for a systolic design in the past are now solvable by systolic arrays. A manual derivation of control signals for these applications is intractable and error-prone. Second, increasing control complexity makes it necessary to evaluate and further optimise the quality of the systolic array with respect to not only the data flow but also the control flow. This is made possible by methods that allow the systematic synthesis of control signals. Third, the synthesis of systolic arrays from algorithmic descriptions is a process of program transformation. At any stage, many transformations are applicable. The choice of different transformations leads to systolic arrays with differing quality. One should be able to assess a transformation by its impact on the quality of the systolic array in terms of not only the data flow but also the control flow. Methods that allow the systematic synthesis of control signals accomplish this by making the specification of control signals explicit in the early stages of the systolic design.

This treatise is concerned with the synthesis of control signals for systolic arrays that are realised in hardware. The proposed method applies for iterative algorithms defined over a domain that can be expressed by a convex set of integer coordinates (or points) called the *index space*. Algorithms that are amenable to systolic design can be expressed in this way. They are often informally referred to as *systolic algorithms*; many can be expressed in the form of affine recurrence equations (AREs) in the functional style [14,19,47,57,60,64] and nested loops in the imperative style [23,27,41,50]. A recurrence equation defines recursively the computation of a target variable indexed by one point in terms of variables indexed by other points. The set of all index vectors associated with the target variable of a recurrence equation is specified by a predicate called a *domain predicate*. A domain predicate specifies the domain of an equation at which the target variable of the equation is defined. Informally, control signals must ensure that the computation of a recurrence equation takes place within the domain specified by the domain predicate of the equation.

A subclass of AREs, called *uniform recurrence equations* (UREs), which was





**Figure 1-1:** The synthesis of control signals.

proposed by R. M. Karp and R. E. Miller and S. Winograd for representing numerical algorithms for a parallel implementation [32], has been used extensively for the specification of systolic algorithms in systolic design. We shall present our method based on input that is in the form of UREs. We refer to these initial UREs as the *source UREs*. A generalisation of the method to certain other types of recurrence equations is straightforward and will be alluded to where appropriate.

To conform to current VLSI fabrication technology, the exchange of input and output data with the external environment is restricted to the border cells of a systolic array. The source UREs may have to be extended to include the specification of the propagation of input data from border cells to internal cells and output data from internal cells to border cells. This added specification is called *propagation UREs*. The extended specification, which consists of the source UREs and propagation UREs, is called the *extended source UREs*. The domain for which the extended source UREs are defined is called the *extended index space*. (The source UREs might be called *computation UREs* since they prescribe the computations to be performed by systolic arrays.)

The synthesis of control signals is a process of program transformation and construction. The domain predicates of the extended source UREs constitute the



initial specification of control signals. We obtain systolic arrays with a description of both data flow and control flow in two successive steps:

1. Transform the domain predicates of the extended source UREs to a system of UREs called the *control UREs* (Fig. 1-1). The control UREs can be divided into two systems of UREs: the *computation control UREs* (CCUREs), which specify the control signals for the appropriate computations prescribed by the source UREs, and the *propagation control UREs* (PCUREs), which specify the control signals for the propagation of input and output data prescribed by the propagation UREs. Then replace the domain predicates in the source UREs by predicates in computation control variables and those in the propagation UREs by predicates in propagation control variables.
2. Obtain systolic arrays with a description of both data and control flow from the previously transformed UREs by applying different space-time mappings.

In the first step, the challenge lies in the specification of control signals in terms of the control UREs. Once this is accomplished, the replacement of domain predicates that arises due to the introduction of the control UREs is straightforward. In the second step, the issues that need to be addressed are (1) the definition of systolic array models, (2) the corresponding mapping conditions that ensure the validity of the space-time mapping and (3) the procedures by which valid space-time mappings, i.e., systolic arrays, are systematically generated. These three issues will not be emphasised in the thesis; they have been studied extensively in the realm of systolic design. Because the specification of control signals in the first step depends on issues addressed in the second step, we shall conduct our presentation in the reverse of the order in which systolic arrays are derived. That is, we first describe the systolic array model under consideration, then the mapping conditions for the validity of the space-time mapping, then some procedures for the generation of systolic arrays, and finally the specification of control signals.

Next, we give an outline of this thesis. Each chapter contains a section of introductory remarks to prepare the reader for its technical context.



## 1.1 Outline of the Thesis

In Chap. 2, we review the basic technique of the synthesis of data flow for  $(n-1)$ -dimensional systolic arrays from  $n$ -dimensional UREs. We start with the definition of  $n$ -dimensional AREs and the systolic array model that has been used extensively in systolic design. Next, we describe the standard space-time mapping technique, which delivers systolic arrays of  $n-1$  dimensions with maximal parallelism, i.e., a shortest execution (or latency) derivable from the  $n$ -dimensional source UREs. This is followed by a brief description of uniformisation techniques that transform AREs to UREs. This chapter concludes with a brief review of several current research issues in systolic design.

In Chap. 3, we present our method for the specification of control signals for  $(n-1)$ -dimensional systolic arrays from  $n$ -dimensional UREs. The specification of control signals relies on a partition of the extended index space such that all the computations associated with one block in the partition are specified by a common set of control signals. We provide necessary and sufficient conditions for the correctness of control flow and a mechanisable procedure that constructs the specification of control signals from the source UREs.

In Chap. 4, we review the synthesis of data flow for one-dimensional systolic arrays and extend and improve previous results. We first define the two most frequently adopted one-dimensional systolic array models. We then investigate the corresponding mapping conditions for the validity of the space-time mapping. In addition, we provide a range of equivalent mapping conditions, which increase our understanding of various properties of one-dimensional systolic arrays and provide insight into the construction of the control UREs for one-dimensional systolic arrays in Chap. 5. Finally, we describe a procedure for the systematic generation of one-dimensional systolic arrays.

In Chap. 5, we focus on the specification of control signals for one-dimensional systolic arrays. We begin by demonstrating that the control UREs that are constructed in Chap. 3 for  $(n-1)$ -dimensional systolic arrays may result in very



inefficient one-dimensional systolic arrays. The problem lies in the specification of the PCUREs and is caused by the non-injectivity of the space-time mapping for one-dimensional systolic arrays. We solve this problem by providing an alternative construction of the PCUREs. To complete our method for the specification of control signals for systolic arrays, we show that the PCUREs constructed in this chapter and the CCUREs constructed in the previous chapter apply for systolic arrays of any  $r$  dimensions with  $0 < r < n$ .

In Chap. 6, we are concerned with the elimination of control signals for systolic arrays. We consider a special class of UREs that do not require the CCUREs. We show how algebraic properties of some operators in the source UREs, like the satisfaction of the Unit Law and the Zero Law (Sect. 1.2), can be exploited to eliminate the PCUREs. We provide necessary and sufficient conditions for the elimination of the PCUREs. Although it is difficult to give a general treatment for arbitrary UREs, the underlying idea of the optimisation of control signals can be generalised straightforwardly to other systolic algorithms.

In Chap. 7, we describe a scheme for loading, recovering and accessing stationary variables. It is sometimes beneficial to make certain variables stationary in order to improve some aspects of the systolic array, e.g., to reduce the latency or size of the array. Unfortunately, the space-time mapping does not provide any help in handling stationary variables. We present a systematic method that modifies the source UREs such that the specification for loading, recovering and accessing stationary variables is included. Then, specifying the control signals for systolic arrays synthesised from this new system of UREs is just a matter of constructing the control UREs based on the methods described in Chaps. 3 and 5.

Finally, in Chap. 8, we give an overview of the results presented in the thesis and comment on some open problems and suggestions for future research.



## 1.2 Notation and Terminology

**Logic** The logic connectives are  $\neg$  (not),  $\vee$  (or),  $\wedge$  (and),  $\implies$  (implies),  $\iff$  (if and only if),  $\forall$  (for all) and  $\exists$  (there exists). We use `true` and `false` to denote the propositional constants. The notation `iff` is also used for  $\iff$ .

The notation for a quantified expression is  $(quant : range : term)$ , where *quant* specifies a quantifier and a list of dummy variables for the quantification, *range* specifies the range of the dummies and *term* is some function or predicate on the dummies [20]. Any binary, commutative, associative operator that has an identity element may be used as a quantifier (for the definition of identity element of an operator, see the subsequent paragraph of this section on abstract algebra). Examples of quantifiers (and the corresponding interpretation in the case of an empty range) include:  $\forall$  (`true`),  $\exists$  (`false`),  $\max$  ( $-\infty$ ) and  $\min$  ( $+\infty$ ); the values `true`, `false`,  $-\infty$  and  $+\infty$  are the identity elements of the operators  $\forall$ ,  $\exists$ ,  $\max$  and  $\min$ , respectively.

Some proofs follow the notation of [20], in which a proof step has the following layout:

$$\begin{array}{l}
 \dots \\
 \neg(P \wedge \neg P) \wedge (Q \vee \neg Q) \\
 \iff \{ \text{De Morgan's law; excluded-middle law} \} \\
 (\neg P \vee \neg \neg P) \wedge \text{true} \\
 \dots
 \end{array}$$

in place of  $\iff$  to the left of the brace may be  $\implies$ ,  $=$ ,  $\subseteq$ , etc. The brace may contain a number of hints separated by conjunctive semicolons.

The symbol  $\square$  marks the end of theorems, definitions, examples, and so forth.

**Programming Languages** A conditional command, i.e., `if`-statement is written as follows [20]:



```

if   $B_1 \rightarrow S_1$ 
     $\square$   $B_2 \rightarrow S_2$ 
        ...
     $\square$   $B_n \rightarrow S_n$ 
fi

```

or, if it is short and simple enough, on one line as:

$$\mathbf{if} \ B_1 \rightarrow S_1 \ \square \ B_2 \rightarrow S_2 \ \cdots \ \square \ B_n \rightarrow S_n \ \mathbf{fi}$$

$B_i$  is a Boolean expression and  $S_i$  a command.  $B_i \rightarrow S_i$  is a *guarded command*;  $B_i$  acts as a *guard*; it has to be validated before  $S_i$  can be executed. One statement  $S_i$  whose guard  $B_i$  is true is executed; if no guard is true, the **if**-statement aborts. If **else** appears in the place of  $B_n$ , it stands for  $(\forall i : 0 < i < n : \neg B_i)$ . **skip** denotes the empty statement. All **if**-statements that appear in the thesis are deterministic.

A **for**-loop has the usual meaning and is written as follows:

```

for  $i$  from  $lb$  to  $rb$  do
     $S_1$ 
     $S_2$ 
    ...
     $S_n$ 

```

We indicate scoping by indenting:  $S_1, S_2, \dots, S_n$  are the statements executed at each loop step.

**Sets** The notation  $\{x \mid \mathcal{P}(x)\}$  denotes the set of elements  $x$  that satisfy the condition  $\mathcal{P}(x)$ . The set whose members are all the objects appearing in the list  $x_1, x_2, \dots, x_m$  and no others is denoted by  $\{x_1, x_2, \dots, x_m\}$ . The empty set is denoted by  $\emptyset$ . If two sets have no element in common, they are called *disjoint*. The number of elements in a set  $S$  is denoted by  $|S|$ . Let  $A$  and  $B$  denote sets. We say that  $A$  is a *subset* of  $B$  (or that  $A$  is *contained in*  $B$ ) iff every member of  $A$  is also a member of  $B$ . Our notation for “ $A$  is a subset of  $B$ ” is  $A \subseteq B$ . When  $A \subseteq B$  and  $A \neq B$ , we say that  $A$  is a *proper subset* of  $B$  and we write



$A \subset B$ . If  $A$  and  $B$  are sets, the *relative complement* of  $B$  in  $A$ , denoted by  $A \setminus B$ , is the set  $\{x \mid x \in A \wedge x \notin B\}$ . Let  $A$  be a subset of a given set  $S$ . The *complement* of  $A$  in  $S$  is the subset  $\{x \mid x \in S \setminus A\}$ . It is denoted by  $\complement_S A$ . For the set-theoretic notions of *union* and *intersection*, we use the symbols  $\cup$  and  $\cap$ . For sets  $S_1, S_2, \dots, S_m$ , where the elements of every set  $S_i$  are themselves sets, we define  $(\bigcap_{i: 0 < i \leq m} S_i) = \{(\bigcap_{i: 0 < i \leq m} X_i) \mid (\forall i: 0 < i \leq m : X_i \in S_i)\}$ .

$\mathbf{Z}$ ,  $\mathbf{Q}$  and  $\mathbf{R}$  denote the set of integers, rationals and reals, respectively. Let  $S$  be  $\mathbf{Z}$ ,  $\mathbf{Q}$  or  $\mathbf{R}$ .  $S^+$  denotes the positive subset,  $S_0^+$  the non-negative subset (i.e.,  $S^+ \cup \{0\}$ ) and  $S^n$  the  $n$ -fold Cartesian product of  $S$ .

**Functions**  $f : D \rightarrow R$  indicates that  $f$  is a function with *domain*  $D$  and *range* (or *codomain*)  $R$ . The set  $\{f(x) \mid x \in S\}$  for a set  $S \subseteq D$  is called the *image* of the set  $S$  under  $f$  and is denoted by  $f(S)$ . The *composite*  $f \circ g$  of two functions is the function obtained by applying them in succession – first  $g$ , then  $f$  – provided the domain of  $f$  is the image of  $g$ . The notation  $f : D \mapsto R$  denotes that function  $f$  is injective from domain  $D$  to range  $R$ . Two sets  $A$  and  $B$  are called *isomorphic* if there exists a bijection  $f : A \rightarrow B$ . Let  $f : S \rightarrow \{a, b\}$  be a surjective function. If  $T \subseteq S$  and  $f(T) \cap f(S \setminus T) = \emptyset$ , then  $f$  is called the *characteristic function* of  $T$  in  $S$  and is denoted by  $\chi_T : S \rightarrow \{a, b\}$ . The function **sign** is a mapping from  $\mathbf{Q}$  to the set  $\{-1, 0, 1\}$ : **sign**( $x$ ) = if  $x < 0 \rightarrow -1$   $\square$   $x = 0 \rightarrow 0$   $\square$   $x > 0 \rightarrow 1$  fi.

**Linear Algebra** Let  $V$  be a vector space over a field  $F$ . A set of vectors  $x_1, x_2, \dots, x_m$  in  $V$  is called *linearly dependent* if there are coefficients  $\lambda_1, \lambda_2, \dots, \lambda_m$  in  $F$ , but not all zero, such that  $(\sum_{i: 0 < i \leq m} \lambda_i x_i) = 0$ . Otherwise, it is called *linearly independent*. A set  $B$  of vectors in  $V$  is a *basis* for a subspace  $L$  of  $V$  if  $B$  is a maximal linearly independent subset of  $L$ . Every linear subspace  $L$  has a basis, and all bases of  $L$  have the same number of elements. This common number of elements is called the *dimension* of  $L$ . The vector space *spanned* by the vectors  $x_1, x_2, \dots, x_m$  in  $V$  is the set  $\text{span}(x_1, x_2, \dots, x_m) = \{(\sum_{i: 0 < i \leq m} \lambda_i x_i) \mid (\forall i: 0 < i \leq m : \lambda_i \in F)\}$ . For an  $m \times n$  matrix  $A$  over a field  $F$ , the set of solutions to  $Ax = 0$  is a vector space called the *null space* of  $A$ . The *rank* of



the matrix  $A$ , denoted by  $\text{rank}(A)$ , is the number of maximal linearly independent column (or row) vectors of  $A$ . Often we do not distinguish whether a vector is a row or column vector and assume that this is deducible from the context. When we write  $xy$  for  $x, y \in V$ , for example, we mean that  $x$  is a row vector and  $y$  is a column vector with the same number of components. For an  $n$ -vector  $x$  in  $F$ ,  $i \cdot x$  ( $0 < i \leq n$ ) denotes its  $i$ -th component.

**Abstract Algebra** For any two sets  $A$  and  $B$ , a subset  $R \subseteq A \times B$  is called a *binary relation* between  $A$  and  $B$ .  $(a, b) \in R$  is often written as  $aRb$ . Let  $X$  be a set. A *partition*  $P$  of a set  $X$  consists of a set of disjoint subsets (called  $P$ -blocks) of  $X$  whose union is the set  $X$ . A relation  $R$  between  $X$  and  $X$  is called an *equivalence relation* on  $X$  when it has the following three properties: (1) *Reflexivity*:  $xRx$  for all  $x \in X$ . (2) *Symmetry*:  $xRy$  implies  $yRx$  for all  $x, y \in X$ . (3) *Transitivity*:  $xRy$  and  $yRz$  imply  $xRz$  for all  $x, y, z \in X$ . Given an equivalent relation  $E$  on a set  $X$ , the *equivalent class* under  $E$  of any element  $x \in X$  is the set  $C_E(x)$  of all the elements  $y$  of  $X$  that bear the relation  $E$  to  $x$ :  $C_E(x) = \{y \mid y \in X \wedge yEx\}$ . Any subset  $C$  of  $X$  that has the form  $C = C_E(x)$ , for some  $x$ , is called an *equivalence class* of  $E$  (or an  $E$ -class). The set of all possible  $E$ -classes is  $X/E = \{C \mid C \subseteq E \wedge C = C_E(x) \wedge x \in X\}$  and is called the *quotient set* of  $X$  by  $E$ . When restricting  $E$  to a subset  $S \subseteq X$ , the quotient set of  $S$  by  $E$  becomes:  $S/E = \{S\} \cap X/E$ , i.e.,  $S/E = \{S \cap C \mid C \in X/E\}$ . Let  $P$  and  $Q$  be partitions of a set  $X$ .  $P$  is *finer* than  $Q$  if  $(\forall x, y : x, y \in X : (\exists S : S \in P : x, y \in S) \implies (\exists T : T \in Q : x, y \in T))$ . Let  $E$  and  $F$  be equivalence relations on a set  $X$ .  $E$  is *finer* than  $F$  if  $X/E$  is finer than  $X/F$ , i.e., if  $(\forall x, y : x, y \in X : xEy \implies xFy)$ .

Let  $S$  be a non-empty set. A *binary operation* on  $S$  is any function  $* : S \times S \rightarrow S$ . An element  $u \in S$  is called a *unit element* for operation  $*$  if

$$(\forall s : s \in S : u * s = s = s * u) \quad (\text{Unit Law})$$

A unit element is often called an *identity element* (or *neutral element*). We call  $u$  a *left (right) unit* for  $*$  if  $(\forall s : s \in S : u * s = s)$  ( $(\forall s : s \in S : s * u = s)$ ). An element  $z \in S$  is called a *zero element* for operation  $*$  if

$$(\forall s : s \in S : z * s = z = s * z) \quad (\text{Zero Law})$$



A zero element is also called an *annihilator*. We call  $z$  a *left (right) zero* if  $(\forall s : s \in S : z * s = z)$  ( $(\forall s : s \in S : s * z = z)$ ).

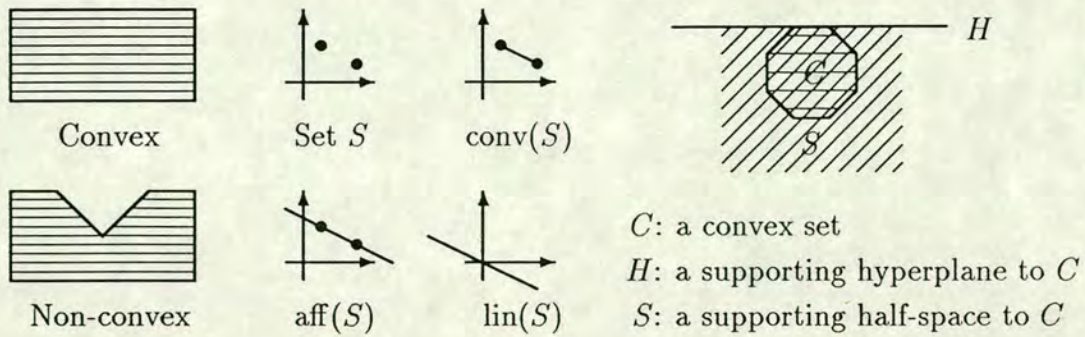
**Convex Analysis [66]** If  $S \subset \mathbf{R}^n$  and  $x \in \mathbf{R}$ , the set  $x + S = \{x + y \mid y \in S\}$  is called a *translate* of  $S$ . For any two sets  $A$  and  $B$ ,  $A + B = \{a + b \mid a \in A \wedge b \in B\}$ . A translate of a subspace of  $\mathbf{R}^n$  is called an *affine set*. Two affine sets are *parallel* if one is a translate of the other. All non-empty parallel affine sets have a unique subspace. The *dimension* of an affine set is the dimension of the corresponding parallel subspace. The *dimension* of a set  $S$  is the dimension of the smallest affine set containing it, and is denoted by  $\dim(S)$ . By convention,  $\dim(\emptyset) = -1$ . A set  $S \subset \mathbf{R}^n$  is *full-dimensional* if  $\dim(S) = n$ . An affine set of dimension 1 is a *line*. An affine set of dimension  $n-1$  is a *hyperplane*. Let  $\{x \mid \pi x = \delta\}$  be a hyperplane, where  $\pi \in \mathbf{R}^n \setminus \{0\}$ ,  $x \in \mathbf{R}^n$ , and  $\delta \in \mathbf{R}$ . The sets  $\{x \mid \pi x \leq \delta\}$  and  $\{x \mid \pi x \geq \delta\}$  are called *closed half-spaces*. The sets  $\{x \mid \pi x < \delta\}$  and  $\{x \mid \pi x > \delta\}$  are called *open half-spaces*. We may speak unambiguously of the open and closed half-spaces associated with a given hyperplane. We denote the hyperplane  $\{x \mid \pi x = \delta\}$  by  $[\pi : \delta]$ . We use  $\square$  to denote a fixed but arbitrary relational operator in  $\{<, \leq, >, \geq\}$ . The notation  $[\pi \square \delta]$  stands for the half-space  $\{x \mid \pi x \square \delta\}$ .

The *affine hull* of a set  $S$  is the intersection of all affine sets that contain  $S$ . It is denoted by  $\text{aff}(S)$ . The unique linear subspace parallel to the affine hull of a set  $S$  is denoted by  $\text{lin}(S)$ . A subset  $S$  of  $\mathbf{R}^n$  is a *convex set* if  $(1-\lambda)x + \lambda y \in S$  for every  $x \in S$ ,  $y \in S$  and  $0 < \lambda < 1$  in  $\mathbf{R}$  (it is an affine set if the range of  $\lambda$  is extended to  $\mathbf{R}$ ). The *convex hull* of a set  $S$  is the intersection of all convex sets that contain  $S$ . It is denoted by  $\text{conv}(S)$ .

The hyperplane  $[\pi : \delta]$  is said to *bound* the set  $S$  if either  $(\forall x : x \in S : \pi x \leq \delta)$  or  $(\forall x : x \in S : \pi x \geq \delta)$ . A hyperplane  $H$  is said to *support* a set  $S$  at a point  $x \in S$  if  $x \in H$  and  $H$  bounds  $S$ . A *supporting half-space* to  $S$  is a half-space containing  $S$  and bounded by a supporting hyperplane to  $S$ .

Let  $A$  and  $B$  be two non-empty sets in  $\mathbf{R}^n$ . A hyperplane  $H$  is said to *separate*  $A$  and  $B$  if (1)  $A$  is contained in one closed half-space and  $B$  is contained in the opposite closed half-space and (2)  $H$  is not a supporting hyperplane for both  $A$  and





**Figure 1-2:** Illustration of some concepts in  $\mathbf{R}^2$  in convex analysis.

$B$ . The standard definition of separation only requires that  $A$  and  $B$  be contained in the opposite closed half-spaces (i.e., (2) is not required). For the purpose of this thesis, our definition of separation deliberately excludes the case when  $H$  supports both  $A$  and  $B$ . A hyperplane  $H$  is said to separate  $A$  and  $B$  *strictly* if  $A$  and  $B$  are contained in the two opposite open half-spaces.

Fig. 1-2 gives a pictorial illustration of some concepts defined previously.

A *polyhedral convex set* (or, *convex polyhedron*) is the intersection of finitely many closed half-spaces. A *polytope* (or, *convex polytope*) is the convex hull of finitely many points, i.e., a bounded convex polyhedron. Let  $S$  be a convex polyhedron in  $\mathbf{R}^n$ . A subset  $F$  of  $S$  is called a *face* of  $S$  if either  $F = \emptyset$  or  $F = S$  or if there exists a supporting hyperplane to  $S$  such that  $F = S \cap H$ . The faces  $\emptyset$  and  $S$  are called *improper*. All other faces are called *proper*. If the dimension of  $F$  is  $k$ , then  $F$  is called a  $k$ -face of  $S$ . It is customary to refer to the 0-faces of  $S$  as *vertices*, the 1-faces as *edges* and the  $(n-1)$ -faces as *facets*. The notation  $\text{facets}(S)$  stands for the set of all facets of  $S$ . For a facet  $F$  of  $S$ ,  $\text{sup}(S, F)$  denotes the supporting half-space to  $S$  that contains  $F$ .

A subset  $S$  of  $\mathbf{R}^n$  is called a *cone* if it is closed under positive multiplication, i.e.,  $\lambda x \in S$  for all  $x \in S$  and  $\lambda \in \mathbf{R}^+$ . A *convex cone* is a cone that is a convex set. The cone *finitely generated* by the vectors  $x_1, x_2, \dots, x_m$  is the set  $\text{cone}(x_1, x_2, \dots, x_m) = \{(\sum_{i=1}^m \lambda_i x_i) \mid (\forall i : 0 < i \leq m : \lambda_i \in \mathbf{R}_0^+)\}$ , i.e, it is the smallest convex cone containing  $x_1, x_2, \dots, x_m$ . A cone that is a proper subset of a line is called a *ray* (or *half-line*). A 1-face of a convex polyhedron that is a ray is called an *extreme ray*.



Let  $x_1, x_2, \dots, x_m$  be linearly independent vectors in  $\mathbf{R}^n$  and  $o \in \mathbf{R}^n$  be the reference point or (the origin). Let  $L_i$  be the line segment from the origin  $o$  to the point  $\alpha_i x_i$  for some non-zero  $\alpha_i$  in  $\mathbf{R}$ . Then  $(\sum_{i: 0 < i \leq m} L_i)$  is called a *k-parallelepiped*. It is called a *hypercube* if the vectors  $x_1, x_2, \dots, x_m$  are mutually orthogonal, i.e., if  $(\forall i, j : 0 < i \leq m \wedge 0 < j \leq m : i \neq j \implies x_i \cdot x_j = 0)$ .

Let  $A$  be an  $m \times n$  real matrix and  $a \in \mathbf{R}^m$ . The function from  $\mathbf{R}^n$  to  $\mathbf{R}^m$  that maps every element  $x \in \mathbf{R}^n$  to the element  $Ax \in \mathbf{R}^m$  is called a *linear transformation*. The function from  $\mathbf{R}^n$  to  $\mathbf{R}^m$  that maps every element  $x \in \mathbf{R}^n$  to the element  $Ax + a \in \mathbf{R}^m$  is called an *affine transformation*.

All concepts introduced so far in  $\mathbf{R}^n$  have parallels in  $\mathbf{Q}^n$  ( $\mathbf{Z}^n$ ). Let  $S = \{x \mid \mathcal{P}(x)\}$  be any of the sets defined previously in  $\mathbf{R}^n$ , the corresponding set in  $\mathbf{Q}^n$  ( $\mathbf{Z}^n$ ) is the set of all the rational (integral) points in  $S$ , and is given by  $\{x \mid \mathcal{P}(x) \wedge x \in \mathbf{Q}^n\}$  ( $\{x \mid \mathcal{P}(x) \wedge x \in \mathbf{Z}^n\}$ ). For example, a convex polyhedron in  $\mathbf{Z}^n$  (or an integral convex polyhedron) is the set of all integral points in the intersection of finitely many closed half-spaces.

**Arithmetic** As is customary, we write  $\lfloor x \rfloor$  for the *floor* of  $x \in R$ , namely for the greatest integer not greater than  $x$ , and we write  $\lceil x \rceil$  for the *ceiling* of  $x$ , i.e., for  $-\lfloor -x \rfloor$ , the smallest integer not smaller than  $x$ . We write  $|x|$  for the absolute value of  $x \in R$ , i.e.,  $|x| = \text{if } x \geq 0 \rightarrow x \text{ [] else } \rightarrow -x \text{ fi}$ . The operator **div** denotes integer division. The operator **mod** denotes the modulo operation. For  $x, y \in \mathbf{Z}$ ,  $x \mid y$  iff  $x$  divides  $y$ , i.e., iff  $y \bmod x = 0$ .

**Symbols for Depicting Systolic Arrays** Unless otherwise stated, we use the following symbols in the depiction of systolic arrays. Boxes represent processors, lines represent the communication channels between processors. Solid lines are *data channels*; they carry the data flow. Dashed lines are *control channels*; they carry the control flow. The delay buffers associated with data and control channels are represented by fat dots and circles, respectively. The direction of data or control signals that move along a channel is represented by an arrow tip at one end of the channel.



## Chapter 2

# Data Flow Synthesis for $(n-1)$ -Dimensional Systolic Arrays

### 2.1 Introductory Remarks

The synthesis of systolic arrays from a program specification proceeds in two successive steps:

1. the specification is refined and transformed to a systolisable source, and
2. the systolisable source is mapped to a systolic array.

Based on the work of [32], researchers have shown that a program in the form of UREs [55,64] (or in the form of nested loops [50]) serves well as the systolisable source. This has spurred research in systolic design in two major directions. One is to develop methods for the mapping of UREs to systolic arrays. The main design tool is a space-time mapping that delivers systolic arrays by means of index transformations of the UREs. The other focuses on the transformation of program specifications of a more general form to UREs. The main result achieved so far is a uniformisation technique for transforming AREs to UREs.

This chapter provides the technical background for the thesis. We review the basic techniques used in the previous two steps, i.e., the basic techniques



underlying the synthesis of data flow for systolic arrays. We only describe in detail the techniques that are necessary for the understanding of the results of this thesis. Other related issues are discussed in the conclusion.

The rest of this chapter is organised as follows. Sect. 2.2 presents the definition of affine and uniform recurrence equations and some related concepts. Sect. 2.3 defines the most commonly used systolic array model. Sect. 2.4 describes the standard space-time mapping technique for the synthesis of the data flow with respect to the systolic array model. In addition, it defines the validity of the space-time mapping and the mapping conditions that ensure this validity. Sect. 2.5 describes a technique for ensuring that input and output data are always handled at the border cells of the systolic array, as required by the systolic array model defined in Sect. 2.3. Sect. 2.6 presents three uniformisation methods for the transformation of AREs to UREs. Sect. 2.7 concludes the chapter by outlining some other related issues and commenting on some current research directions.

## 2.2 Affine and Uniform Recurrence Equations

**Definition 2.1** A *system of affine recurrence equations* is a finite collection of equations each of which has the form [57]

$$I \in D \rightarrow V(I) = f(W(\rho_W(I)), \dots) \quad (2.1)$$

where

- $D \subset \mathbf{Z}^n$  is a union of disjoint convex polytopes in  $\mathbf{Z}^n$  and is called the *domain of the equation*. All defining equations with the same target variable must be distinct. This simplifies the presentation with no loss of generality. Domains of equations that have the same target variable must be disjoint. This ensures that the system of AREs is well-formed. The *index space*, denoted  $\Phi$ , is the convex hull of the union of all the domains of equations. The *domain of variable*  $V$ , denoted  $\Phi_V$ , is the convex hull of the union of the domains of all equations with target variable  $V$ . (deq stands for a domain of equation)



$\text{deq}(\Phi)$  denotes the set of all domains of equations.  $\text{deq}(V)$  denotes the set of the domains of all equations with target variable  $V$ .

- $I$  is called an *index vector* or *point*. The predicate  $I \in D$  is called a *domain predicate*. Each domain predicate can be put into disjunct normal form, where each disjunct consists of a conjunction of predicates called *conditionals* that are generally linear (or affine) functions of the index vector  $I$ . Conditionals that are in one of the following forms:

$$\pi I = \delta, \pi I \neq \delta, \pi I < \delta, \pi I \leq \delta, \pi I > \delta, \pi I \geq \delta \quad (2.2)$$

are called *basic conditionals*, where  $\pi \in \mathbf{Z}^n$  is a *coefficient vector* and  $\delta \in \mathbf{Z}$  is a *constant*.  $\pi I = \delta$  is called an *equality*.  $\pi I \neq \delta$  is called an *inequality*. (Note that  $\pi I = \delta \iff \pi I \leq \delta \wedge \pi I \geq \delta$  and  $\pi I \neq \delta \iff \pi I < \delta \vee \pi I > \delta$ .) A domain predicate is said to be in *normal form* if its constituent conditionals are basic conditionals.

- $V$  and  $W$  are variable names belonging to a finite set  $\mathcal{V}$ .
- The “...” stands for an arbitrary but fixed number of similar arguments.
- $f$  is a function that has time complexity  $O(1)$ .
- $\rho_W$  is an affine mapping from  $\mathbf{Z}^n$  to  $\mathbf{Z}^\ell$  ( $0 \leq \ell \leq n$ ) called an *index mapping*:

$$\rho_W(I) = \Delta_W I - \vartheta_W \quad (2.3)$$

where  $\Delta_W \in \mathbf{Z}^\ell \times \mathbf{Z}^n$  and  $\vartheta_W \in \mathbf{Z}^\ell$ .  $\Delta_W$  is called the *linear part* of  $\rho_W$ . If  $\ell = n$ ,  $n$ -vector  $I - \rho_W(I)$  is called a *data dependence (vector)*.

- The symbol  $\rightarrow$  separates the domain predicate from the corresponding defining equation. (Note that  $\rightarrow$  is used in [57]. We choose  $\rightarrow$  to avoid notational confusion with the use of  $\rightarrow$  in the **if**-statement (Sect. 1.2).)  $\square$

The elements of a variable  $V$  are partitioned into three subsets:

- *Input data* are only read. They are identified by adding the following equa-



tions, called *input equations*, to the source program:

$$I \in D_V \rightarrow V(I) = v_{in_V(I)} \quad (2.4)$$

where  $in_V$  is a mapping from  $\mathbf{Z}^n$  to  $\mathbf{Z}^{h_{in_V}}$  ( $0 \leq h_{in_V} \leq n$ ). That is, input element  $V(I)$  has the value  $v_{in_V(I)}$ .

- *Output data* are only assigned. They are identified by adding the following equations, called *output equations*, to the source program:

$$I \in E_V \rightarrow v_{out_V(I)} = V(I) \quad (2.5)$$

where  $out_V$  is a mapping from  $\mathbf{Z}^n$  to  $\mathbf{Z}^{h_{out_V}}$  ( $0 \leq h_{out_V} \leq n$ ). That is, output element  $V(I)$  has the value  $v_{out_V(I)}$ .

- *Intermediate data* are read and assigned. Equations that are neither input nor output equations are called *computation equations*. In other words, the computation equations are those introduced by Def. 2.1.

**Remark** Input and output equations are not taken into account in the definitions of  $\text{deq}(V)$ ,  $\text{deq}(\Phi)$ ,  $\Phi_V$  and  $\Phi$  (Def. 2.1). They only serve to declare which elements of variables are input and output data. It is the computation equations that specify the computations to be performed in the systolic array.  $\square$

Following [14,58], we sometimes write the set of all defining equations with the same target variable  $V$  in the following abbreviated form

$$V(I) = \begin{cases} I \in D_1 & \rightarrow f_1(W(\rho_1(I)), \dots) \\ I \in D_2 & \rightarrow f_2(W(\rho_2(I)), \dots) \\ & \dots \\ I \in D_p & \rightarrow f_p(W(\rho_p(I)), \dots) \end{cases} \quad (2.6)$$

When specifying a system of AREs, we shall adopt one of the previous two styles of notation, whichever is more convenient.

The concept of a data dependence graph plays an important rôle in systolic design. A *data dependence graph* has one node for each point of the index space and a directed arc from node  $J$  to node  $I$  iff a variable indexed by  $J$  is an argument in the equation for a variable indexed by  $I$ .



UREs are a subclass of AREs. They differ in the format of the index mapping. In UREs, the index mapping  $\rho_W$  of (2.3) is of the form:

$$\rho_W(I) = I - \vartheta_W \quad (2.7)$$

That is, the linear part of the index mapping is the identity matrix. The data dependence vector  $I - \rho_W(I)$  in the AREs becomes  $\vartheta_W$  in the UREs.  $\vartheta_W$  is a constant, i.e., is independent of the index vector  $I$ . This exposes the uniformity and regularity of the computations prescribed by UREs, as opposed to AREs that are not UREs. Data dependence vector  $I - \rho_W(I)$  in the AREs depends on the index vector  $I$  if the linear part  $\Delta_W$  is not the identity matrix. If  $\Delta_W$  is singular, the data dependence is not injective. Non-injective dependences are called *broadcast dependences*. Injective dependences are called *pipelining dependences*. Broadcast dependences allow the same variable to appear as arguments in the defining equations of an unbounded number of target variables. This represents the broadcast of a datum from one processor to an unbounded number of processors. In addition, irregular and non-constant data dependences result in irregular and non-constant channel connections. The advantage of UREs is that they enforce pipelining and regular and local channel connections.

By convention, a system of UREs is called *unconditional* if all its equations are defined for the index space, i.e., if  $I \in \Phi$  is the domain predicate for all the equations. Otherwise, it is *conditional*.

We now introduce some concepts that are related to UREs. The *data dependence matrix*  $\mathcal{D}$  is a matrix whose column vectors are the data dependence vectors [64]. For notational convenience, we assume that, for each variable name  $V$ , there is only one associated data dependence vector, denoted  $\vartheta_V$ . When we write  $\vartheta \in \mathcal{D}$ , we mean that  $\vartheta$  is a data dependence vector, that is, a column of  $\mathcal{D}$ .

The points of the index space are called the *computation points*. If  $I$  is inside and  $I - \vartheta_V$  is outside the domain  $\Phi_V$  of variable  $V$ ,  $I$  is called a *first computation point* of  $V$ . If  $I$  is inside and  $I + \vartheta_V$  is outside the domain  $\Phi_V$  of variable  $V$ ,  $I$  is called a *last computation point* of  $V$ . The set of first computation points,



$\text{fst}(\Phi_V, \vartheta_V)$ , and the set of last computation points,  $\text{lst}(\Phi_V, \vartheta_V)$ , of  $V$  are given by

$$\begin{aligned}\text{fst}(\Phi_V, \vartheta_V) &= \{I \mid I \in \Phi_V \wedge I - \vartheta_V \notin \Phi_V\} \\ \text{lst}(\Phi_V, \vartheta_V) &= \{I \mid I \in \Phi_V \wedge I + \vartheta_V \notin \Phi_V\}\end{aligned}\quad (2.8)$$

The set  $\text{in}(\Phi_V, \vartheta_V)$  is the translate of  $\text{fst}(\Phi_V, \vartheta_V)$  by  $-\vartheta_V$ :

$$\text{in}(\Phi_V, \vartheta_V) = \text{fst}(\Phi_V, \vartheta_V) - \vartheta_V \quad (2.9)$$

Since  $\text{fst}(\Phi_V, \vartheta_V)$  and  $\text{in}(\Phi_V, \vartheta_V)$  are translates of each other, we shall use either of them in our analysis, as is convenient. By convention, the input data are supplied at the points of  $\text{in}(\Phi_V, \vartheta_V)$  and used for the first time at the points of  $\text{fst}(\Phi_V, \vartheta_V)$ . The output data are assumed to be available at the points of  $\text{lst}(\Phi_V, \vartheta_V)$ . If this is not the case, they can be made so by adding pipelining equations [64]; a *pipelining equation* is of the form

$$I \in D \rightarrow V(I) = V(I - \vartheta_V) \quad (2.10)$$

A variable is called a *pipelining variable* if it is specified by pipelining equations.

Having explicitly defined the domains at which input data are supplied and the domains at which output data are defined, we can make the input equations of (2.4) and the output equations of (2.5) in the source UREs more explicit. Replacing  $D_V$  of (2.4) by  $\text{in}(\Phi_V, \vartheta_V)$  yields the input equations:

$$(\forall V : V \in \mathcal{V} : I \in \text{in}(\Phi_V, \vartheta_V) \rightarrow V(I) = v_{\text{in}_V(I)}) \quad (2.11)$$

Replacing  $E_V$  of (2.5) by  $\text{lst}(\Phi_V, \vartheta_V)$  yields the output equations:

$$(\forall V : V \in \mathcal{V} : I \in \text{lst}(\Phi_V, \vartheta_V) \rightarrow v_{\text{out}_V(I)} = V(I)) \quad (2.12)$$

where  $v_{\text{in}_V(I)}$  and  $v_{\text{out}_V(I)}$  are input and output data of variable  $V$ , respectively.

Let us introduce some notation. For  $I, \vartheta \in \mathbb{Q}^n$ ,  $\text{ray}(I, \vartheta)$  denotes the set of the elements  $I, I + \vartheta, I + 2\vartheta, \dots$ , i.e.,

$$\text{ray}(I, \vartheta) = \{J \mid J = I + m\vartheta \wedge m \in \mathbb{Z}_0^+\}$$

We write  $\text{ray}(I, \pm\vartheta)$  for  $\text{ray}(I, \vartheta) \cup \text{ray}(I, -\vartheta)$ . For convenience, we write  $I \xrightarrow{\vartheta} J$  if  $J \in \text{ray}(I, \vartheta)$  and  $I \not\xrightarrow{\vartheta} J$  if  $J \notin \text{ray}(I, \vartheta)$ . For  $D \subset \mathbb{Q}^n$ , we define

$$\text{rays}(D, \vartheta) = \left( \bigcup I : I \in D : \text{ray}(I, \vartheta) \right)$$

We write  $\text{rays}(D, \pm\vartheta)$  for  $\text{rays}(D, \vartheta) \cup \text{rays}(D, -\vartheta)$ .



Let us use the multiplication of  $n \times n$  matrices as an example for illustration. This example will be used for illustration throughout the paper.

**Example 2.1**  $m \times m$  Matrix Product

*Specification:*  $(\forall i, j : 0 < i \leq m \wedge 0 < j \leq m : c_{i,j} = (\sum k : 0 < k \leq m : a_{i,k} b_{k,j}))$

*AREs:*

$$\begin{aligned}
 A(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 = j \wedge 0 < k \leq m & \rightarrow a_{i,k} \\ 0 = i \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow b_{k,j} \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 = k & \rightarrow 0 \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow C(i, j, k-1) \\ & + A(i, 0, k) B(0, j, k) \end{cases} \\
 c_{i,j} &= \begin{cases} 0 < i \leq m \wedge 0 < j \leq m \wedge k = m & \rightarrow C(i, j, k) \end{cases}
 \end{aligned}$$

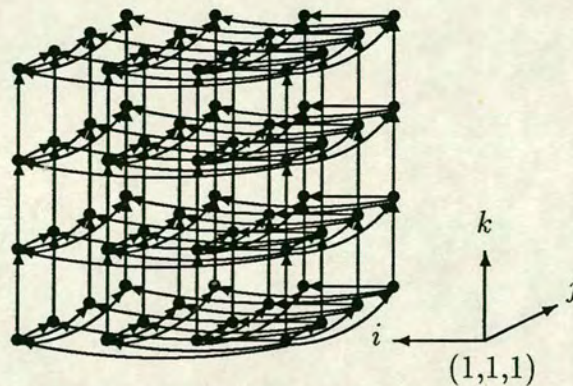
*Index Space:*  $\Phi = \{(i, j, k) \mid 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m\}$

*Domains of Variables:*  $\Phi_A = \emptyset \quad \Phi_B = \emptyset \quad \Phi_C = \Phi$

*Index Mappings:*

$$\begin{aligned}
 \Delta_A &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \Delta_B &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \Delta_C &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\
 \vartheta_A &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} & \vartheta_B &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} & \vartheta_C &= \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}
 \end{aligned}$$

*Data Dependence Graph (m=4):*



**Figure 2-1:** The data dependence graph of the previous AREs ( $m=4$ ).



UREs:

$$\begin{aligned}
 A(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 = j \wedge 0 < k \leq m & \rightarrow a_{i,k} \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow A(i, j-1, k) \end{cases} \\
 B(i, j, k) &= \begin{cases} 0 = i \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow b_{k,j} \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow B(i-1, j, k) \end{cases} \\
 C(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 < j \leq m \wedge 0 = k & \rightarrow 0 \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow C(i, j, k-1) \\ & + A(i, j-1, k)B(i-1, j, k) \end{cases} \\
 c_{i,j} &= \begin{cases} 0 < i \leq m \wedge 0 < j \leq m \wedge k = m & \rightarrow C(i, j, k) \end{cases}
 \end{aligned}$$

Index Space:  $\Phi = \{(i, j, k) \mid 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m\}$

Domains of Variables:  $\Phi_A = \Phi_B = \Phi_C = \Phi$

Data Dependence Matrix:  $D = [\vartheta_A, \vartheta_B, \vartheta_C] = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

First Computation Points:  $\text{fst}(\Phi, \vartheta_A) = \{(i, 1, k) \mid 0 < i \leq m \wedge 0 < k \leq m\}$   
 $\text{fst}(\Phi, \vartheta_B) = \{(1, j, k) \mid 0 < j \leq m \wedge 0 < k \leq m\}$   
 $\text{fst}(\Phi, \vartheta_C) = \{(i, j, 1) \mid 0 < i \leq m \wedge 0 < j \leq m\}$

Last Computation Points:  $\text{lst}(\Phi, \vartheta_A) = \{(i, m, k) \mid 0 < i \leq m \wedge 0 < k \leq m\}$   
 $\text{lst}(\Phi, \vartheta_B) = \{(m, j, k) \mid 0 < j \leq m \wedge 0 < k \leq m\}$   
 $\text{lst}(\Phi, \vartheta_C) = \{(i, j, m) \mid 0 < i \leq m \wedge 0 < j \leq m\}$

Data Dependence Graph ( $m=4$ ):

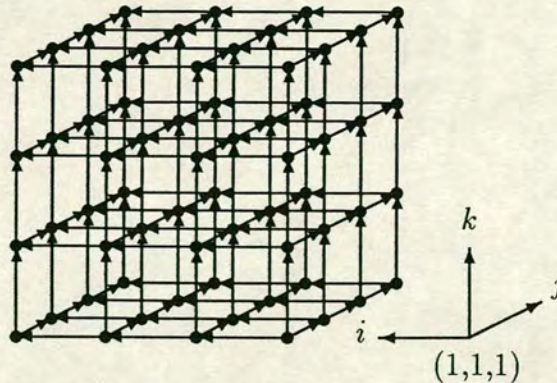


Figure 2-2: The data dependence graph of the previous UREs ( $m=4$ ). □



## 2.3 The Systolic Array Model

There is one systolic array model that has been the basis of many synthesis methods [19,27,50,55,64].

**Definition 2.2** A *systolic array* is a network of cells that are placed at the grid points of an  $r$ -dimensional set  $\mathcal{L} \in \mathbf{Z}^r$  ( $0 < r < n$ ) and that satisfy the following properties:

Prop. 1. (Synchrony of Computation) The array is driven by a global clock that ticks in unit time. A cell is active at every clock cycle.

Prop. 2. (Uniqueness of Channel Connections) For every variable  $V$ , postulate the existence of a unique, directed connection from the cell at location  $p$  to the cell at location  $p+d_V$ , for some constant vector  $d_V \in \mathbf{Z}^r$ . This postulate is either true for all  $p \in \mathcal{L}$  or false for all  $p \in \mathcal{L}$ . A directed connection is also called a *channel*; it is an *input channel* to the cell at its destination and an *output channel* to the cell at its source. If  $p \in \mathcal{L}$  and  $p-d_V \notin \mathcal{L}$ , cell  $p$  is called an *input cell* of variable  $V$ . If  $p \in \mathcal{L}$  and  $p+d_V \notin \mathcal{L}$ , cell  $p$  is called an *output cell* of variable  $V$ . Both input and output cells of  $V$  are also called *border cells* of variable  $V$ ; they are connected to the external environment. The cells that are not border cells of  $V$  are called *internal cells* of  $V$ .

Prop. 3. (Linearity of Velocity) All channels represented by  $d_V$ , for a fixed variable  $V$ , are associated with the same number of delay buffers. Thus, all elements of  $V$  move with the same constant velocity.  $\square$

Prop. 1 assumes that the computation at a point in the index space takes unit time. The sentence “a cell is active at every clock cycle” in Prop. 1 is interpreted as follows. If a cell receives a value on an input channel at a given time step, then it will receive a value on the same input channel and send a value on the corresponding output channel at the next time step. As a special case, if the values



on all the input channels of a cell at a given time step are undefined (or  $\perp$ ), then the cell is conventionally regarded as inactive. However, it is sometimes convenient to regard such a cell as active. In this case, the cell at the given step performs an undefined computation; it will send the undefined value,  $\perp$ , at the next time step on all its output channels. By Prop. 2, only border cells are connected with the external environment. Depending on the length of a channel, a border cell may be a processor inside the array – a few processors away from the real border of the array. Prop. 2 requires that a cell always receives (sends) the elements of a variable via a fixed input (output) channel. It also enforces the regularity and locality of channel connections. Prop. 2 and Prop. 3 ensure that all elements of the same variable move with a fixed velocity.

## 2.4 The Space-Time Mapping

The space-time mapping maps every computation (or point) in the source UREs to a time step and a processor allocation. It is a linear (or affine) transformation of the index space. The advantage of linear (or affine) index transformations is that they enforce easily the linearity of variables' velocities and the regularity of channel connections.

**Definition 2.3** A space-time mapping consists of two components: **step** and **place**. Useful functions defined in terms of **step** and **place** are **flow** and **pattern**.

**step** :  $\Phi \rightarrow \mathbb{Z}$ ,  $\text{step}(I) = \lambda I$ ,  $\lambda \in \mathbb{Z}^n$ . **step** specifies the temporal distribution.  $\lambda$  is the *scheduling vector*.  $I$  is computed at step  $\lambda I$ .

**place** :  $\Phi \rightarrow \mathbb{Z}^r$ ,  $\text{place}(I) = \sigma I$ ,  $\sigma \in \mathbb{Z}^{r \times n}$ . **place** specifies the spatial distribution.  $\sigma$  is the *allocation matrix*.  $\mathcal{P} = \{\sigma I \mid I \in \Phi\}$  is the *processor space* of  $r$  dimensions.  $I$  is computed at cell  $\sigma I$ .

**flow** :  $\mathcal{V} \rightarrow \mathbb{Q}^r$ ,  $\text{flow}(V) = \sigma \vartheta_V / \lambda \vartheta_V$ ,  $\vartheta_V \in \mathcal{D}$ . **flow** specifies the velocity with which elements of a variable travel at each step. Variable  $V$  is called *moving*



if  $\text{flow}(V) \neq 0$  and *stationary* if  $\text{flow}(V) = 0$ .  $\sigma\vartheta_V$  represents the direction and length of a connecting channel at a cell for variable  $V$ ; the number of delay buffers associated with that channel is  $\lambda\vartheta_V - 1$ .

$\text{pattern} : \mathcal{V} \rightarrow \Phi_{\text{fst}} \rightarrow \mathbf{Z}^{r \times n}$ ,  $\text{pattern}(V(I)) = \text{place}(I) - (\text{step}(I) - t_{\text{fst}})\text{flow}(V)$ .  
 $\Phi_{\text{fst}} = (\cup V : V \in \mathcal{V} : \text{fst}(\Phi_V, \vartheta_V))$ .  $t_{\text{fst}}$  is the first step number.  $\text{pattern}$  specifies the location of variables in the processor space at the first step.

$\text{step}$  is called the *step function* or *timing function*.  $\text{place}$  is called the *place function* or *allocation function*. □

This chapter considers only space-time mappings with  $r = n - 1$ . They describe systolic arrays of  $n - 1$  dimensions. The motivation for employing an  $(n - 1)$ -dimensional array to implement a system of  $n$ -dimensional UREs is to obtain a time-minimal systolic array, an array with the smallest number of execution steps (called *latency*) derivable from the source UREs.

The *space-time matrix*,  $\Pi$ , is the matrix formed by the scheduling vector in the first row and the allocation matrix in the remaining rows:

$$\Pi = \begin{bmatrix} \lambda \\ \sigma \end{bmatrix} \quad (2.13)$$

A space-time matrix uniquely determines a space-time mapping. For convenience, we sometimes refer to the space-time matrix as the space-time mapping. We shall denote the image of some  $x$  under a given space-time mapping with an overbar:  $\bar{x}$ . Here,  $x$  may be a point, a set of points, a data dependence vector, and so on. For example, the image  $\bar{\vartheta}_V$  of data dependence vector  $\vartheta_V$  is given by

$$\bar{\vartheta}_V = \begin{bmatrix} \lambda\vartheta_V \\ \sigma\vartheta_V \end{bmatrix} \quad (2.14)$$

The space-time mapping can be interpreted geometrically. All points that are scheduled concurrently belong to a hyperplane called a *temporal hyperplane*. The set of points that are scheduled at step  $t$  is in the temporal hyperplane  $[\lambda : t]$ . The processor space is obtained by a projection, namely, by eliminating one dimension



from the index space along a chosen direction called the *projection vector*, and denoted by  $u$ . Here,  $u \in \mathbb{Z}^n$  is a normalised vector that satisfies  $\sigma u = 0$  [64]. The points that are in a line parallel to the projection vector are assigned the same processor location. The set of points that are allocated to cell  $p$  is in the line ray( $I, \pm u$ ), where  $\sigma I = p$ . In what follows, we shall use either the projection vector or the allocation matrix, as is convenient.

Data dependence vectors are projected onto channels. For variable  $V$ , dependence vector  $\vartheta_V$  is projected onto channel  $\bar{\vartheta}_V$ . The first component of  $\bar{\vartheta}_V$  decremented by 1 represents the number of delay buffers associated with the channel and the remaining components represent the length and direction of the channel.

The following mapping rules characterise the validity of a space-time mapping.

**Definition 2.4** A space-time mapping is *valid* for the source UREs in the systolic array model (Def. 2.2) iff

*Precedence Rule:* The step assigned to a target variable is smaller than the steps assigned to its arguments. (This ensures that the data dependences prescribed in the source UREs are preserved.)

*Computation Rule:* Concurrent computations are mapped to different processors.

*Communication Rule:* At most one element of a variable is injected to a fixed input cell at any time step.

*Delay Rule:* The number of buffers associated with a communication channel is a non-negative integer. □

The following mapping conditions ensure the validity of a space-time mapping [57,64].

**Theorem 2.1** A space-time mapping  $\Pi$  is valid if

- $(\forall V : V \in \mathcal{V} : \lambda \vartheta_V > 0)$  (*Precedence Constraint*)
- $\text{rank}(\Pi) = n$ , i.e.,  $\lambda u \neq 0$  (*Computation and Communication Constraint*)



It is straightforward to see that the precedence constraint is equivalent to the precedence rule. We sometimes write  $\lambda\vartheta_V \geq 1$  instead of  $\lambda\vartheta_V > 0$  in the precedence constraint if we want to emphasise the fact that the evaluation of a point takes unit time. The computation and communication constraint implies the computation and communication rule. This can be understood in two different ways. First, the non-singularity of  $\Pi$  (i.e.,  $\text{rank}(\Pi) = n$ ) ensures that the space-time mapping is a bijection from  $\mathbb{Q}^n$  to  $\mathbb{Q}^n$ . Therefore, points that are mapped to the same step must differ in their processor coordinates. Second,  $\lambda u \neq 0$  means that the projection vector is not orthogonal to the scheduling vector. Hence, any line parallel to the projection vector can only intersect at most one point in any temporal hyperplane. To see that the delay rule is satisfied, recall that data dependence vectors are projected onto channels. The number of delay buffers associated with a channel for variable  $V$  is the non-negative integer  $\lambda\vartheta_V - 1$ .

We may view the synthesis of systolic arrays as a process of program transformations [14,56]. It amounts to finding a suitable linear index transformation, i.e., a space-time mapping of the index space. In the transformed UREs, one index represents time and the remaining indices represent processor coordinates. Finding different systolic arrays means just finding different valid space-time mappings.

We use matrix product to illustrate both the geometrical and transformational view of the synthesis of systolic arrays.

### Example 2.2 $m \times m$ Matrix Product

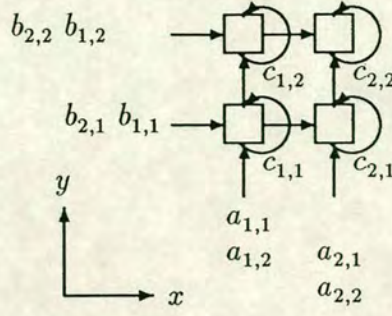
The following space-time mapping describes S. Y. Kung's two-dimensional arrays for matrix product [64] (Fig. 2-3):

$$\Pi = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad \Pi I = \bar{I} = (t, x, y)$$

The array consists of  $m^2$  cells and runs in  $3m-2$  steps.

First, we consider the geometrical view. The processor space is obtained by projecting the index space along the projection vector  $u = (0, 0, 1)$ . The points





**Figure 2–3:** S. Y. Kung's systolic array for matrix product with data distribution at the first step ( $m=2$ ).

that are scheduled at step  $t$  are in the hyperplane  $[\lambda : t]$ . The points that are mapped to cell  $p$  are in the line  $\text{ray}(I, \pm u)$ , where  $\sigma I = p$ . Elements of  $A$  travel along channel  $\sigma\vartheta_A = (0, 1)$  with a speed of one unit per cycle. Elements of  $B$  travel along channel  $\sigma\vartheta_B = (1, 0)$  with a speed of one unit per cycle. Elements of  $C$  are stationary because  $\sigma\vartheta_C = 0$ .

Second, we consider the transformational view. This array is described precisely by the following UREs called the *space-time UREs* obtained from the source UREs and the space-time mapping by replacing (1) the indices  $i, j$ , and  $k$  of all the variables by  $t, x$  and  $y$ , respectively, (2) all the data dependence vectors by their images under the space-time mapping as defined in (2.14), and (3) the indices  $i, j$  and  $k$  in the domain predicates of the source UREs by the expressions in the corresponding components of vector  $\Pi^{-1}(t, x, y)$ .

*The Space-Time UREs:*

$$\begin{aligned}
 A(t, x, y) &= \begin{cases} 0 < x \leq m \wedge 0 = y \wedge 0 < t - x - y \leq m & \rightarrow a_{x, t-x-y} \\ 0 < x \leq m \wedge 0 < y \leq m \wedge 0 < t - x - y \leq m & \rightarrow A(t-1, x, y-1) \end{cases} \\
 B(t, x, y) &= \begin{cases} 0 = i \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow b_{t-x-y, y} \\ 0 < x \leq m \wedge 0 < y \leq m \wedge 0 < t - x - y \leq m & \rightarrow B(t-1, x-1, y) \end{cases} \\
 C(t, x, y) &= \begin{cases} 0 < x \leq m \wedge 0 < y \leq m \wedge 0 = t - x - y & \rightarrow 0 \\ 0 < x \leq m \wedge 0 < y \leq m \wedge 0 < t - x - y \leq m & \rightarrow C(t-1, x, y) \\ & + A(t-1, x, y-1)B(t-1, x-1, y) \end{cases} \\
 c_{x,y} &= \begin{cases} 0 < x \leq m \wedge 0 < y \leq m \wedge t - x - y = m & \rightarrow C(t, x, y) \end{cases}
 \end{aligned}$$



From the space-time UREs, we can directly extract the following information: the latency, the number of channels required, the number of processors required, the input/output characteristics, the storage requirement, etc.  $\square$

## 2.5 Input and Output Extension

There are space-time mappings that do not project first and last computation points to border cells. The corresponding input and output data are mapped to internal cells. This is costly because the I/O pins of a chip are restricted to the boundary of the chip and, consequently, non-local communication channels may be needed. To impose the restriction of border communication, the index space can be extended in such a way that the resulting first and last computation points are all projected to border cells.

The basic idea is to replicate a data dependence vector in the data dependence graph forward and backward until the images of points so created are at the boundary of the processor space. This extension depends on the projection vector and works only for moving variables. We give a formal definition of such an extension, based on an informal idea described in [64]. Note that  $\text{rays}(\Phi, \pm u)$  denotes the set of all the points in  $\mathbf{Z}^n$  whose images under the space-time mapping are in the processor space.

**Definition 2.5** The *extended index space*  $\Psi$  of the index space  $\Phi$  with respect to the projection vector  $u$  is defined as follows:

$$\begin{aligned} \Psi_V^s &= \begin{cases} \text{if } \text{flow}(V)=0 & \rightarrow \emptyset \\ \text{if } \text{flow}(V) \neq 0 & \rightarrow \text{rays}(\Phi, \pm u) \cap \text{rays}(\text{fst}(\Phi_V, \vartheta_V) - \vartheta_V, -\vartheta_V) \\ \text{fi} \end{cases} \\ \Psi_V^d &= \begin{cases} \text{if } \text{flow}(V)=0 & \rightarrow \emptyset \\ \text{if } \text{flow}(V) \neq 0 & \rightarrow \text{rays}(\Phi, \pm u) \cap \text{rays}(\text{lst}(\Phi_V, \vartheta_V) + \vartheta_V, \vartheta_V) \\ \text{fi} \end{cases} \\ \Psi_V &= \Phi_V \cup \Psi_V^s \cup \Psi_V^d \\ \Psi^P &= (\cup V : V \in \mathcal{V} : \Psi_V^s \cup \Psi_V^d) \end{aligned}$$



$$\begin{aligned}\Psi &= \Phi \cup \Psi^P \\ \Psi_V^\perp &= \Psi \setminus \Psi_V\end{aligned}$$

The points of  $\Psi_V^s$  are called the *soaking points* of  $V$ ; they are generated in the extension of  $V$  along direction  $-\vartheta_V$ . The points of  $\Psi_V^d$  are called the *draining points* of  $V$ ; they are generated in the extension of  $V$  along direction  $\vartheta_V$ . The points of  $\Psi_V^\perp$  are called the *undefined points* of  $V$ ; they are generated in the extension of the other variables.  $\Psi_V$  is the *extended domain of variable  $V$* .  $\Psi^P$  contains the soaking and draining points of all variables; they are called the *pipelining points*. The data dependence graph resulting from the extension is called the *extended data dependence graph*.  $\square$

The following concepts are analogous to those defined for the index space. The set of *new* first computation points of variable  $V$  over the extended index space is  $\text{fst}(\Psi_V, \vartheta_V)$ . The set of *new* last computation points of variable  $V$  over the extended index space is  $\text{lst}(\Psi_V, \vartheta_V)$ .

The two sets  $\text{in}(\Psi_V, \vartheta_V)$  and  $\text{in}(\Phi_V, \vartheta_V)$  are isomorphic by the bijection:

$$\mathcal{I}_V : \text{in}(\Psi_V, \vartheta_V) \rightarrow \text{in}(\Phi_V, \vartheta_V), \quad \mathcal{I}_V(I) = J \quad \text{where } I \xrightarrow{\vartheta_V} J \quad (2.15)$$

( $\text{fst}(\Psi_V, \vartheta_V)$  and  $\text{fst}(\Phi_V, \vartheta_V)$  are also isomorphic because they are translates of  $\text{in}(\Psi_V, \vartheta_V)$  and  $\text{in}(\Phi_V, \vartheta_V)$ , respectively.) Similarly, the two sets  $\text{lst}(\Psi_V, \vartheta_V)$  and  $\text{lst}(\Phi_V, \vartheta_V)$  are isomorphic by the bijection:

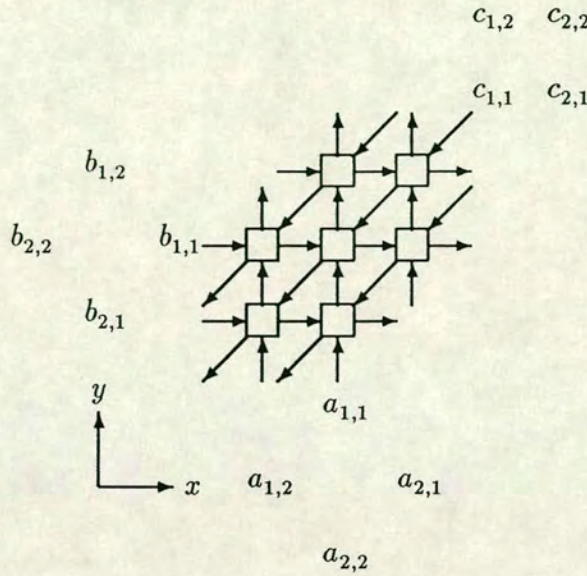
$$\mathcal{O}_V : \text{lst}(\Psi_V, \vartheta_V) \rightarrow \text{lst}(\Phi_V, \vartheta_V), \quad \mathcal{O}_V(I) = J \quad \text{where } J \xrightarrow{\vartheta_V} I \quad (2.16)$$

The motivation underlying the extension of the index space is to have the input data of  $V$  supplied at the points of  $\text{in}(\Psi_V, \vartheta_V)$  and the output data of  $V$  defined at the points of  $\text{lst}(\Psi_V, \vartheta_V)$ . The *extended source UREs* defined over the extended index space are a refinement of the source UREs; they are obtained as follows:

- Replace the input equations as defined in (2.11) by

$$(\forall V : V \in \mathcal{V} : I \in \text{in}(\Psi_V, \vartheta_V) \rightarrow V(I) = v_{\text{in}_V(\mathcal{I}_V(I))}) \quad (2.17)$$





**Figure 2-4:** Kung-Leiserson's systolic array for matrix product with data distribution at the first step ( $m=2$ ).

- Replace the output equations as defined in (2.12) by

$$(\forall V : V \in \mathcal{V} : I \in \text{lst}(\Psi_V, \vartheta_V) \rightarrow v_{\text{out}_V(\mathcal{O}_V(I))} = V(I)) \quad (2.18)$$

- Add the following pipelining equations

$$(\forall V : V \in \mathcal{V} : I \in \Psi_V \setminus \Phi_V \rightarrow V(I) = V(I - \vartheta_V)) \quad (2.19)$$

Let us illustrate how Kung-Leiserson's two-dimensional array for matrix product [37] can be obtained by extending the index space.

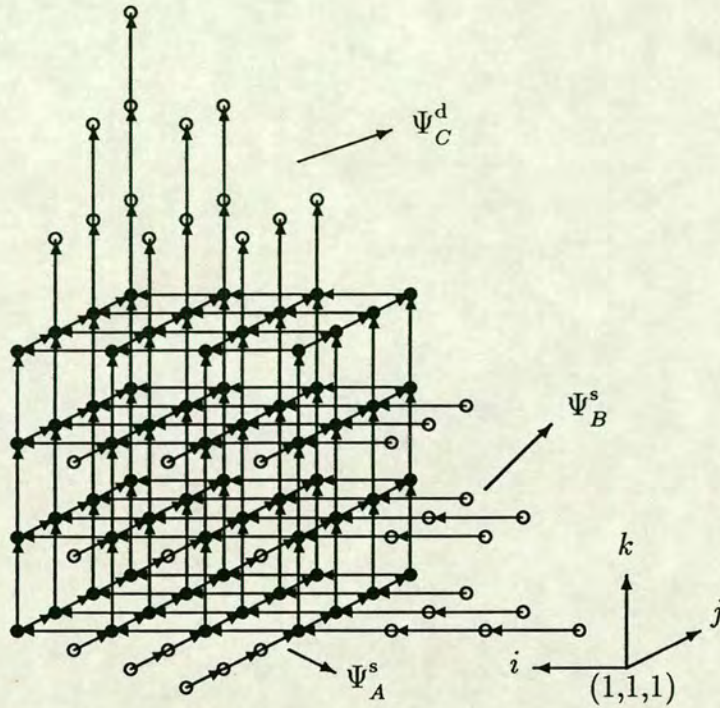
### Example 2.3 $m \times m$ Matrix Product

The following space-time mapping describes Kung-Leiserson's two-dimensional array for matrix product (Fig. 2-4):

$$\Pi = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & -1 \\ 0 & 1 & -1 \end{bmatrix} \quad \Pi I = \bar{I} = (t, x, y)$$

The array consists of  $3m^2 - 3m + 1$  and runs in  $5m - 4$  time steps. Compared with S. Y. Kung's array,  $2m - 2$  extra time steps are needed for the propagation of





**Figure 2-5:** Part of the extended data dependence graph for matrix product ( $m=4$ ). The extension along  $-\vartheta_C$  is not depicted; it is symmetric to that along  $\vartheta_C$  with respect to the projection vector. Extensions along  $\vartheta_A$  and  $\vartheta_B$  are not necessary because the output of  $A$  and  $B$  is of no interest.

input data from border cells to internal cells and output data from internal cells to border cells.

The processor space is obtained by projecting the index space along the projection vector  $u = (1,1,1)$ . Let us look at the case where the index space is not extended. The set of first (last) computation points for each of the three variables  $A$ ,  $B$  and  $C$  are contained in a hyperplane (Fig. 2-2). Since the projection vector is not orthogonal to the normal of any of these hyperplanes, some input and output data need to be handled at internal cells. For example, point  $(1,1,1)$  is mapped to the cell at the center of the array. As a result, data elements  $a_{1,1}$ ,  $b_{1,1}$  and  $c_{1,1}$  must be injected at this cell, which by definition is not a border cell. The extension of the index space enforces border communication (Fig. 2-5), because all input and output data are now defined at the boundary of the extended index space, which are projected to the border cells.  $\square$



## 2.6 Uniformisation

The projection vector completely determines the size and shape of the processor space and the channel interconnection pattern between processors. (The number of buffers associated with a channel remains to be determined by the scheduling vector.) Unless uniform, a system of recurrence equations cannot be mapped to a systolic array. For example, no systolic array can be obtained from a projection of the data dependence graph depicted in Fig. 2-1, since channels of unbounded length would be generated.

*Uniformisation* refers to the process by which AREs are transformed to UREs [57,58,79]. It aims at replacing non-constant data dependence vectors (i.e., broadcast and non-constant pipelining dependences) with constant ones (i.e., constant pipelining dependences). For this reason, uniformisation is also referred to as *data pipelining*. Uniformisation proceeds in two steps. In the first step, one identifies a set of constant data dependence vectors called *pipelining vectors* for the UREs to be constructed. In the second step, one expresses each data dependence vector in the AREs as a linear combination of the pipelining vectors and then replaces the data dependence vector by the pipelining vectors. Uniformisation is a difficult problem because the choice in neither of these two steps is unique.

There are basically three different uniformisation methods. We distinguish them by the way the pipelining vectors are introduced and used. We only describe the idea underlying each method and refer details to the corresponding references.

The method of [57,58] uses the null bases of the linear parts of the index mappings as the pipelining vectors. Consider the AREs and the corresponding data dependence graph depicted in Fig. 2-1.  $I-\rho_A(I)$  and  $I-\rho_B(I)$  are broadcast dependences. Every element of  $A$  or  $B$  needs to be broadcast to  $m$  points. Let us consider variable  $A$  first. Element  $a_{i,k}$  must be broadcast to the points  $(i, 1, k)$ ,  $(i, 2, k)$ ,  $\dots$ ,  $(i, m, k)$ . We want to find a pipelining vector such that  $a_{i,k}$  can be routed incrementally to the corresponding  $m$  points along the pipelining vector. This pipelining vector may be chosen from the null basis of the linear part  $\Delta_A$  of the index



mapping for variable  $A$ . Let us choose  $(0, 1, 0)$ , since  $\Delta_A(0, 1, 0) = 0$ . Similarly, we choose  $(1, 0, 0)$  for the pipelining vector for variable  $B$ , since  $\Delta_B(1, 0, 0) = 0$ . Consequently, we have transformed the AREs to the UREs described in Sect. 2.2. This method cannot transform non-constant pipelining dependences to constant pipelining dependences. The reason is that the null basis of a non-singular matrix is  $\emptyset$ , which does not contribute pipelining vectors. The method also breaks down if the point at which a datum is produced has a non-constant distance from the point at which the datum is used for the first time in the pipeline.

The method of [79] uses canonical or non-canonical bases of the vector space  $\mathbb{Q}^n$  as pipelining vectors. Let there be a data dependence between two points  $I$  and  $J$ . Choose  $n$  linearly independent integral vectors as canonical or non-canonical bases of the vector space  $\mathbb{Q}^n$ . We then express the dependence vector  $I - J$  as a linear combination of the  $n$  pipelining vectors. If the linear parts of the index mappings are either the identity matrix or singular, all dependences of this form can be transformed to constant pipelining dependences. The resulting UREs may not allow the existence of a step function, but it is possible to find multi-step functions; one step function for each variable (Sect. 2.7; for details, see [64]).

The method of [57] uses the extreme rays of the cone generated by the data dependence vectors in the AREs as the pipelining vectors. Consider the set of data dependence vectors  $\{(i-k, 0, 1) \mid k < i \leq m \wedge k \leq j \leq m+1 \wedge 0 < k \leq m\}$  for an ARE specification of Gauss-Jordan elimination for square matrices [57]. There is one extreme ray  $(1, 0, 0)$  and one vertex  $(1, 0, 1)$  in the cone generated by these data dependence vectors. We can express the original data dependence vector  $(i-k, 0, 1)$  as a linear combination of the two pipelining vectors  $(1, 0, 0)$  and  $(1, 0, 1)$ :

$$(i-k, 0, 1) = (i-k-1)(1, 0, 0) + (1, 0, 1)$$

Thus, the non-constant dependences are transformed to constant pipelining dependences. Assuming that there exists a step function for the original AREs (i.e., that there exists a mapping of integers to the points in the index space such that the integer mapped to a target variable is always smaller than those mapped to its arguments) and that the original data dependence vectors generate a cone, all AREs can be transformed to UREs that permit a step function.



## 2.7 Conclusion

We have presented the basic techniques for the synthesis of systolic arrays from UREs and for the transformation of AREs to UREs. We conclude this chapter with a brief description of some other related issues in systolic design and comments on some further research topics.

The step function described previously is linear. A more general case is an affine step function:  $\text{step}(V, I) = \lambda I + \alpha_V$ , where  $\alpha_V \in \mathbf{Z}$ . There is one step function for each variable  $V \in \mathcal{V}$  in the source UREs [64]. The affine step function is more general than the linear step function because it may apply in the presence of zero data dependence vectors. The precedence constraint must be changed accordingly:

$$(\forall V, W, I, J : V, W \in \mathcal{V} \wedge I, J \in \Phi : W(I) \dashrightarrow V(J) \implies \text{step}(V, J) > \text{step}(W, I))$$

where  $W(I) \dashrightarrow V(J)$  denotes that  $W(I)$  is an argument appearing in the defining equation of target variable  $V(J)$ . For simplicity, we adopt linear rather than affine step functions. The existence of an affine step function in the source UREs guarantees the existence of a linear step function in the UREs resulting from a pre-index translation to the source UREs [64]. A pre-index translation can be obtained from the additive constants in  $\{\alpha_V \mid V \in \mathcal{V}\}$ . It serves to eliminate zero dependence vectors in the initial UREs. A search for more general forms of step and place that are capable of describing more general systolic arrays has been one of the research topics in systolic design [55,74].

A systolic array that is described by the space-time mapping as defined here is often referred to as a *parameterised systolic array*, since its size grows with the size of the problem. The *partitioning* of a parameterised systolic array to reduce it to some prescribed shape and size has been an active research topic recently [11,17,52,82]. Partitioning is achieved by first dividing the processor space into a number of blocks such that no two processors in the same block are assigned concurrent computations, and then merging the processors in the same block into one processor. Partitioning is possible with an extension of the standard space-time mapping [17]. An alternative approach to partitioning takes place at the



program level. The idea is to partition the problem into smaller problems that can be mapped directly to individual systolic arrays [29,53].

As a special case, partitioning also solves the problem of designing space-time minimal systolic arrays [9,12,16]. A systolic array is *space-time minimal* when it uses as few processors as any systolic array that has a minimal latency. It is also a parameterised systolic array since we need to retain the latency of the original array. The step and place functions for the partitioned systolic arrays are no longer linear functions (of the indices); they are piecewise linear and generally contain **mod** or **div** operators [22,74].

The space-time mapping technique delivers a systolic array of  $n-1$  dimensions. If the systolic array postulates more dimensions than are available, a *projection* of the processor space becomes necessary. Projection can also be considered as a special case of partitioning, but it deserves a solution in its own right. We shall address this topic in Chap. 4, when we deal with the synthesis of data flow for one-dimensional arrays.

There have been some attempts to automate the generation of step and place. Early work on this subject reported in [46] only handles some specific problems. The minimisation of the step function with respect to the computation of one specific variable (rather than all variables) can be formulated as a linear programming problem [69,70]. The step function can also be minimised by combinatorial optimisation. A branch-and-bound method has been proposed and demonstrated to be effective on practical problems [80]. A method for minimising the number of processors has also been reported [81]. The solution space of projection vectors is bounded for a given step function; the minimising place function can be obtained by enumeration.

Systolic arrays can also be emulated in software [6,7,21,63,65]. The space-time UREs can be coded in programming languages, which can then be executed in processor networks like Warp [1] and Transputers [25]. Work on the parallelisation of nested loops is relevant here [5,40,77]. The formerly distinct areas of the parallelisation of nested loops and the synthesis of systolic arrays from algorithmic descriptions have reached closer proximity [31,76].



Two different methods for the synthesis of systolic arrays deserve particular attention. One is based on graph theory. It transforms a synchronous circuit to a systolic array by redistributing the delays buffers associated with channels using the retiming theorem [43]. The retiming theorem is often transferred to the context of systolic automata [15,30]. The other is based on the theory of functional program transformations. It captures the regularity of algorithmic descriptions by the regularity exhibited in function compositions (so-called circuit combinators) rather than the regularity exhibited in data dependence graphs [49]. This method is appealing for its simplicity, elegance and facility to reason about various properties of systolic arrays in a concise manner.

A lot of issues in systolic design are now well-understood, but much remains to be studied. Each issue mentioned previously needs further study. In addition, we describe several research topics that have not received adequate attention. One topic is the transformation of program specifications of a more general form to UREs. A desirable specification should only allow the existence of data dependences that enforce the correctness of the program, as is possible in GAMMA [3,4] and Unity [13]. But these languages are too general to support an implemented scheme of systolic design. They offer too many implementation choices because they are targeted at a much wider range of architectures than systolic arrays. To support the systematic development of systolic implementations, we need to develop a set of transformation rules that are specifically targeted at systolic arrays.

A second topic is the synthesis of systolic arrays with multi-level pipelining processors [48,73]. A large amount of computations prescribed by systolic algorithms have efficient pipelined implementations [36]. A third topic is the analysis and synthesis of programmable systolic arrays [28,42]. A further topic is the synthesis of control hardware for systolic arrays. This is becoming increasingly important if we are dealing with partitioned, projected or programmable systolic arrays. In subsequent chapters, we shall present a method for the systematic derivation of control signals for parameterised systolic arrays.



## Chapter 3

# Control Flow Synthesis for $(n - 1)$ -Dimensional Systolic Arrays

### 3.1 Introductory Remarks

When the notion of a systolic array was introduced a decade ago [34,35], a systolic array was defined to be a network of simple cells, which always perform the same computation. Therefore, no control mechanism was required. Later, the notion of a systolic array was extended to permit a cell to perform different computations at different steps [26]. This made a control mechanism necessary that instructs the cell when to perform what computation. An example in [26] suggested that this control mechanism could be implemented by pipelining control signals, in addition to the data, though.

In Chap. 2, we have reviewed the space-time mapping technique for the synthesis of data flow for  $(n - 1)$ -dimensional arrays. The space-time mapping provides a description of the layout and velocities of data with a guarantee that the right data arrive at the right cells at each step. But it fails to provide the control signals necessary to guarantee that the right computations are executed at each step.

The formal derivation of control flow has not received adequate attention for several reasons. First, there was an initial focus on the synthesis of data flow.



Application problems considered for systolic design were very simple. A manual derivation of control signals was adopted. This is no longer feasible when complex problems are implemented as systolic arrays. Second, a systolic array can be viewed as a system of space-time UREs. Rather than in hardware, it can be realised in software on processor networks like Warp or Transputers [7,21,65], where the issue of the derivation of control flow becomes irrelevant. Third, the space-time mapping technique does not provide a specification of the control flow. Alternative means must be sought. The necessity and complexity of developing a formal method for the derivation of control flow were not well understood and therefore research on this topic was neglected.

This chapter describes a method for the systematic synthesis of control signals for  $(n-1)$ -dimensional systolic arrays synthesised from iterative algorithms defined over a convex set of integer coordinates. We first present the method with respect to the  $n$ -dimensional UREs and then generalise it to algorithms of a more general form in Sect. 3.8.

Recall that a variable  $V$  in the source UREs is of the following form (see (2.6)):

$$V(I) = \begin{cases} I \in D_1 & \rightarrow f_1(W(I-\vartheta_1), \dots) \\ I \in D_2 & \rightarrow f_2(W(I-\vartheta_2), \dots) \\ & \dots \\ I \in D_p & \rightarrow f_p(W(I-\vartheta_p), \dots) \end{cases} \quad (3.1)$$

We assume that  $\{D_1, D_2, \dots, D_p\}$  is a partition of the index space  $\Phi$ . If not, we can always enforce it by adding the equation  $I \in D_{p+1} \rightarrow \perp$ , where  $D_{p+1} = \Phi \setminus (\cup_{i: 0 < i \leq p} D_i)$ . Therefore,  $\text{deq}(V)$ , which is the set of all domains of equations with target variable  $V$  (Def. 2.1), is a partition of the index space.

Just like the derivation of the data flow, the derivation of the control flow is a process of program transformation and refinement. In the derivation of the data flow, a problem specification is first transformed to a system of UREs and then the data flow is derived from the UREs by means of a space-time mapping. The derivation of the control flow can proceed along the same lines. The domain predicates in the source UREs constitute the initial specification of the control



mechanism for the systolic array. We are motivated to transform these domain predicates to a system of UREs called the *control UREs*. We then replace the domain predicates by predicates in the variables of the control UREs. We say that the domain predicates are *pipelined*. To derive the control flow, we simply apply the same space-time mapping that was previously applied to the source UREs, now to the control UREs. Of course, the space-time mapping must be valid for both the source UREs and the control UREs.

To make the presentation more precise, we duplicate the terminology of source UREs for control UREs and prefix the words “data” and “control”, respectively. That is, we speak of data variables vs. control variables, and so on.

The rest of this chapter is organised as follows. Sect. 3.2 characterises control signals. Sect. 3.3 focuses on the construction of the control UREs and the subsequent replacement of the domain predicates of the source UREs by predicates in control variables. Sect. 3.4 discusses some techniques for optimising the specification of control signals before and after the space-time mapping is chosen. Sect. 3.5 extends the standard space-time mapping technique so that systolic arrays with a description of both data and control flow are directly synthesised from both the source and control UREs. Sect. 3.6 contains a survey of related work. Sect. 3.7 illustrates our method with two examples. Sect. 3.8 contains the conclusion of this chapter. It discusses the generalisation of our method to recurrence equations of a more general form.

## 3.2 The Requirements on Control Flow

We aim at implementing the control mechanism for a systolic array by pipelining control signals across the array. The control signals received at the input channels of a cell determine the computation to be performed by that cell. The communication of control signals in systolic arrays needs to satisfy certain requirements. Different requirements may lead to different specifications of control signals. We adopt the requirements on the communication of data defined in the systolic ar-



ray model (Def. 2.2) for the communication of control signals. When viewed as electronic signals, control signals and data need not be distinguished.

The control signals in  $(n-1)$ -dimensional arrays must satisfy four restrictions:

Rst. 1. Control signals are both input and output at the border cells of the array.

Rst. 2. The number of different control values that travel along a control channel is a constant, i.e., is independent of the size and dimension of the source UREs.

Rst. 3. A control signal moves with a constant velocity.

Rst. 4. Once input, a control signal remains unchanged.

Rst. 1 keeps the number of pins required in the systolic array to a minimum. Rst. 2 ensures that the bit-width of a channel is a constant. Both restrictions cater to the constraints imposed by the current VLSI technology. Rst. 3 stems from the fact that the control signals are specified by UREs and that the space-time mapping is linear. Rst. 4 is not imperative but useful. It happens to be a property of the control UREs constructed for  $(n-1)$ -dimensional arrays, but will not hold for the control UREs constructed for one-dimensional arrays in Chap. 5.

Rsts. 1 and 3 are imposed on the data in the systolic array model (Def. 2.2). Rst. 2 is guaranteed by the implicit assumption that a data channel has constant bit-width.

**Remark** The domain predicates of the source UREs constitute the initial specification of control signals. As far as the synthesis of control signals is concerned, the defining equations of the variables in the source UREs are irrelevant. When we use the source UREs for illustration, it is sufficient to describe them by the sets  $\text{deq}(V)$  for all the variables  $V \in \mathcal{V}$ . We shall not specify these sets in terms of predicates in indices. Rather, we specify them by means of diagrams; one diagram per variable. The diagram associated with variable  $V$  depicts the partition of the index space specified by  $\text{deq}(V)$ . □



### 3.3 The Synthesis of Control Flow

Our systolic array model requires that the data are exchanged with the external environment only through border cells. The restriction to border communication can be enforced by an extension of the index space with respect to the projection vector. The source UREs defined over the index space  $\Phi$  can be extended to the extended index space  $\Psi$  by adding pipelining equations for all variables at the newly generated points, i.e., in  $\mathbb{C}_\Psi\Phi$ . For notational convenience, we write the extension of (3.1) over the extended index space as follows:

$$V(I) = \begin{cases} I \in \Phi \wedge I \in D_1 & \rightarrow f_1(W(I-\vartheta_1), \dots) \\ I \in \Phi \wedge I \in D_2 & \rightarrow f_2(W(I-\vartheta_2), \dots) \\ & \dots \\ I \in \Phi \wedge I \in D_p & \rightarrow f_p(W(I-\vartheta_p), \dots) \\ I \in \mathbb{C}_\Psi\Phi & \rightarrow V(I-\vartheta) \end{cases} \quad (3.2)$$

The last defining equation of  $V$  corresponds to the pipelining equations in (2.19). In other words, it specifies that the elements of  $V$  at the cells of  $\text{place}(\mathbb{C}_\Psi\Phi)$  are propagated at the steps  $\text{step}(\mathbb{C}_\Psi\Phi)$ . This suffices for the purpose of the control flow derivation, because it is the domain  $\mathbb{C}_\Psi\Phi$  rather than the dependence vector  $\vartheta$  that determines the specification of control signals.

We shall construct the control UREs for the extended source UREs. The construction is based on a characterisation of the points in the extended index space. This characterisation aims at partitioning the extended index space into blocks such that all points in one block may be specified by a common set of control signals. For this purpose, we introduce the concept of a type for points. It relies on the following equivalence relation on  $\Phi$ .

**Definition 3.1**  $\mathbb{T}$  is the following equivalence relation on  $\Phi$ :

$$I \mathbb{T} J \iff (\forall D : D \in \text{deq}(\Phi) : I \in D \iff J \in D)$$

We say that  $I$  and  $J$  are of the same *type* if  $I \mathbb{T} J$ . Points that are not of the same type are said to have different types.  $\square$



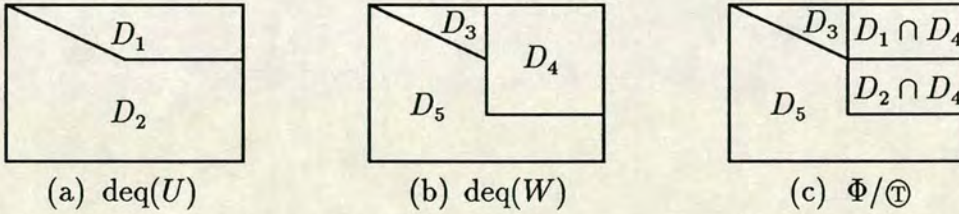
Recall that the defining equations of all domains with the same target variable are distinct (Def. 2.1). Therefore, the defining equations for a variable of the source UREs agree at points of the same type and do not agree at points of different types. The definition of  $\mathbb{T}$  does not take the input and output equations in the source UREs into account. Remember that  $\text{deq}(\Phi)$  does not contain the domains of input and output equations (Sect. 2.2). Input and output equations only serve to declare which elements of variables are input and output data. They do not specify the computations to be performed in the systolic array.

The quotient set of  $\Phi$  by  $\mathbb{T}$  is given by

$$\Phi/\mathbb{T} = (\mathbb{m} V : V \in \mathcal{V} : \text{deq}(V)) \quad (3.3)$$

To see why this is so, we note that  $\text{deq}(V)$  for a fixed variable  $V$  is a partition of the index space. So a subset of  $\Phi$  is a  $\mathbb{T}$ -class iff it is the intersection of  $|\Phi/\mathbb{T}|$  blocks, one from every partition  $\text{deq}(V)$ .

**Example 3.1** Let the source UREs be specified as in Figs. 3-1(a) and (b).



**Figure 3-1:** The construction of the quotient set  $\Phi/\mathbb{T}$ .

The quotient set  $\Phi/\mathbb{T}$  depicted in Fig. 3-1(c) follows from (3.3):

$$\begin{aligned} & \Phi/\mathbb{T} \\ = & \{ \Phi/\mathbb{T} \text{ of (3.3)} \} \\ & (\mathbb{m} V : V \in \mathcal{V} : \text{deq}(V)) \\ = & \{ \mathcal{V} = \{U, W\} \} \\ & \text{deq}(U) \mathbb{m} \text{deq}(W) \\ = & \{ \text{Definition of } \mathbb{m} \text{ in Sect. 1.2} \} \\ & \{ D_1 \cap D_3, D_1 \cap D_4, D_1 \cap D_5, D_2 \cap D_3, D_2 \cap D_4, D_2 \cap D_5 \} \\ = & \{ D_1 \cap D_3 = D_3; D_1 \cap D_5 = \emptyset; D_2 \cap D_3 = \emptyset; D_2 \cap D_5 = D_5 \} \\ & \{ D_3, D_1 \cap D_4, D_2 \cap D_4, D_5 \} \end{aligned}$$

□



Once the control UREs are constructed, we need to replace the domain predicates of the extended source UREs by predicates in control variables. Here, the quotient set  $D/\mathbb{T}$  of a domain  $D \in \text{deq}(\Phi)$  of equation by  $\mathbb{T}$  comes into play.

The quotient sets  $D/\mathbb{T}$  for all domains  $D \in \text{deq}(\Phi)$  in Ex. 3.1 are:

$$\begin{aligned} D_1/\mathbb{T} &= \{D_3, D_1 \cap D_4\} \\ D_2/\mathbb{T} &= \{D_2 \cap D_4, D_5\} \\ D_3/\mathbb{T} &= \{D_3\} \\ D_4/\mathbb{T} &= \{D_1 \cap D_4, D_2 \cap D_4\} \\ D_5/\mathbb{T} &= \{D_5\} \end{aligned}$$

They are subsets of  $\Phi/\mathbb{T}$ . The following lemma states that this is always so.

**Lemma 3.1** ( $\forall D : D \in \text{deq}(\Phi) : D/\mathbb{T} \subseteq \Phi/\mathbb{T}$ ).

**Proof** The domain  $D$  of an equation is a subset of  $\Phi$ , i.e.,  $D \subseteq \Phi$ .  $D/\mathbb{T}$  is the quotient set  $\Phi/\mathbb{T}$  restricted to  $D$ . That is,  $D/\mathbb{T} = \{D\} \pitchfork \Phi/\mathbb{T}$ . By the hypothesis,  $D \in \text{deq}(\Phi)$ . An application of (3.3) yields  $D/\mathbb{T} = \{D\} \pitchfork \Phi/\mathbb{T} \subseteq \Phi/\mathbb{T}$ .  $\square$

We sometimes write  $\Phi_i$  ( $0 < i \leq r \wedge r = |\Phi/\mathbb{T}|$ ) for a fixed but arbitrary  $\mathbb{T}$ -class if we are not concerned with the exact elements in the class. This enables us to write  $\Phi/\mathbb{T} = \{\Phi_i \mid 0 < i \leq r\}$ .  $\{\mathbb{L}_\Psi \Phi, \Phi_1, \Phi_2, \dots, \Phi_r\}$  is a partition of the extended index space  $\Psi$ . All points in one block of the partition can be specified by a common set of control signals. The specification of control signals proceeds in two steps:

1. Construct the control UREs from the domain predicates in the extended source UREs. The control UREs are divided into two systems:
  - *Computation control UREs* (CCUREs) distinguish different types of computation points, (i.e., different  $\mathbb{T}$ -classes). They specify the control signals for the evaluation of the computation points in the index space.
  - *Propagation control UREs* (PCUREs) distinguish pipelining points (of  $\mathbb{L}_\Psi \Phi$ ) from computation points (of  $\Phi$ ). They specify the control signals for the pipelining of input (output) data that are mapped to the internal cells of the array from (to) the border cells.



The dichotomy exhibited by the control UREs is natural since  $\{\mathcal{C}_\Psi\Phi, \Phi\}$  is a partition of the extended index space  $\Psi$ . The points in  $\mathcal{C}_\Psi\Phi$  are specified by pipelining equations, and the points in  $\Phi$  are specified by the source UREs.

2. Replace the domain predicates in the extended source UREs by equivalent predicates in the control variables: replace those that are in the source UREs by predicates in computation control variables and those that are in the extended source UREs but not in the source UREs, i.e.,  $I \in \Phi$  and  $I \in \mathcal{C}_\Psi\Phi$ , by predicates in propagation control variables. The resulting extended source UREs and the control UREs are called *pipelined UREs*.

If the source UREs are unconditional, i.e., if  $I \in \Phi$  is the domain predicate of all equations, then  $\Phi/\mathbb{T} = \{\Phi\}$ . In this case, computation control is not needed. If all input and output data prescribed by the source UREs are mapped to border cells, propagation control is not needed.

The specification of a control variable  $V$  proceeds in three steps:

- Choose a non-zero normalised vector  $\vartheta_V \in \mathbf{Z}^n$  as the control dependence vector. (Each control variable is associated with one control dependence vector.)
- Define the set of control values, denoted  $\text{sig}(V)$ , and the input equation that specifies the initialisation of  $V$  with control values in  $\text{sig}(V)$ . The domain of the input equation is  $\text{in}(\Psi, \vartheta_V)$ .
- Define the computation equation as  $I \in \Psi \rightarrow V(I) = V(I - \vartheta_V)$ .

The control signals of control variable  $V$  are initialised at  $\text{in}(\Psi, \vartheta_V)$  and pipelined along  $\vartheta_V$  across the extended index space  $\Psi$ . We shall address the minimisation of the domains of control variables in Sect. 3.4. The control UREs are unconditional. Otherwise, we would need to add a second level of control UREs that governs the computations prescribed by the initial control UREs, and so on.

Let us discuss how the control UREs specified this way satisfy the four restrictions on control signals described in Sect. 3.2. The sets  $\text{fst}(\Psi, \vartheta_V)$  and  $\text{lst}(\Psi, \vartheta_V)$



of control variable  $V$  are part of the boundary of the extended index space. The extension of the index space ensures that all first and last computation points of data variables are mapped to border cells. That is, only points of  $(\bigcup V : V \in \mathcal{V} : \text{fst}(\Psi_V, \vartheta_V) \cup \text{lst}(\Psi_V, \vartheta_V))$  are mapped to the border cells. The sets  $\text{fst}(\Psi, \vartheta_V)$  and  $\text{lst}(\Psi, \vartheta_V)$  may contain points that are mapped to internal cells. To enforce the border communication of control signals, i.e., Rst. 1, we extend the extended index space (with respect to the projection vector) for the control UREs in the same way as we extended the index space for the source UREs in Def. 2.5 (see Figs. 6–1(a) and (c)). Rst. 2 is satisfied if we restrict  $|\text{sig}(V)|$  to be a constant, Rst. 3 is satisfied since the space-time mapping is linear, and Rst. 4 is satisfied since the computation equations of control variables are pipelining equations.

The construction of the control UREs relies on the following theorem about the separation of two convex sets in  $\mathbf{R}^n$  by hyperplanes [66].

**Theorem 3.1** (*Separation Theorem*) *If  $C_1$  and  $C_2$  are non-empty disjoint convex polytopes in  $\mathbf{R}^n$ , there always exists a hyperplane that strictly separates  $C_1$  and  $C_2$ .*

To apply the Separation Theorem, each  $\mathbb{T}$ -class must be a convex polytope in  $\mathbf{Z}^n$ . If not, it must be partitioned into a union of convex polytopes. We refer to a partition of the index space as a  $\mathbb{P}$ -partition if it is finer than the quotient set  $\Phi/\mathbb{T}$  and if every block of the partition is a convex polytope. The quotient set  $\Phi/\mathbb{T}$  is a  $\mathbb{P}$ -partition if all  $\mathbb{T}$ -classes are convex.

Let us address the problem of obtaining  $\mathbb{P}$ -partitions from the source UREs. We can approach this problem in two different but equivalent ways. In the first approach, we obtain  $\mathbb{P}$ -partitions by partitioning each  $\mathbb{T}$ -class into a union of convex polytopes. We refer to a partition of a  $\mathbb{T}$ -class  $C \in \Phi/\mathbb{T}$  as a *convex class partition*, denoted  $\text{clapar}(C)$ , if the blocks of the partition are convex polytopes. This gives rise to a  $\mathbb{P}$ -partition, denoted  $\mathbb{P}_{\text{clapar}}$ :

$$\mathbb{P}_{\text{clapar}} = \left( \bigcup C : C \in \Phi/\mathbb{T} : \text{clapar}(C) \right) \quad (3.4)$$

In the second approach, we obtain  $\mathbb{P}$ -partitions by partitioning each domain of equations into a union of convex polytopes. We refer to a partition of a domain



$D \in \text{deq}(\Phi)$  as a *convex domain partition*, denoted  $\text{dompar}(D)$ , if the blocks of the partition are convex polytopes. We define

$$\text{condeq}(V) = \left( \bigcup D : D \in \text{deq}(V) : \text{dompar}(D) \right) \quad (3.5)$$

$$\text{condeq}(\Phi) = \left( \bigcup V : V \in \mathcal{V} : \text{condeq}(V) \right) \quad (3.6)$$

This gives rise to a  $\mathbb{P}$ -partition, denoted  $\mathbb{P}_{\text{dompar}}$ :

$$\mathbb{P}_{\text{dompar}} = \left( \mathbb{m} V : V \in \mathcal{V} : \text{condeq}(V) \right) \quad (3.7)$$

$\mathbb{P}_{\text{dompar}}$  is a  $\mathbb{P}$ -partition because  $\text{condeq}(V)$  is finer than  $\text{deq}(V)$  and its blocks are convex polytopes.

The previous two ways of obtaining  $\mathbb{P}$ -partitions are equivalent.

**Lemma 3.2** *Let  $\mathbb{P}_c$  be the set of all  $\mathbb{P}$ -partitions defined by (3.4), and  $\mathbb{P}_d$  be the set of all  $\mathbb{P}$ -partitions defined by (3.7). Then,  $\mathbb{P}_c = \mathbb{P}_d$ .*

**Proof** Every  $\mathbb{P}_{\text{clapar}}$  in  $\mathbb{P}_c$  is also in  $\mathbb{P}_d$  if we choose  $\text{dompar}(D) = \{D\} \mathbb{m} \mathbb{P}_{\text{clapar}}$  for every domain of equation  $D \in \text{deq}(\Phi)$ . Every  $\mathbb{P}_{\text{dompar}}$  in  $\mathbb{P}_d$  is also in  $\mathbb{P}_c$  if we choose  $\text{clapar}(C) = \{C\} \mathbb{m} \mathbb{P}_{\text{dompar}}$  for every  $\mathbb{T}$ -class  $C \in \Phi/\mathbb{T}$ .  $\square$

To find  $\mathbb{P}$ -partitions, we do not need to construct the quotient set  $\Phi/\mathbb{T}$  in order to obtain convex class partitions for non-convex  $\mathbb{T}$ -classes. Lemma 3.2 states that finding different  $\mathbb{P}$ -partitions amounts to finding different convex domain partitions for all domains of the equations in the source UREs. Recall the definition of AREs in Def. 2.1. Each domain of an equation must be expressed as a union of disjoint convex polytopes. This means that the set of these convex polytopes is a convex domain partition. Therefore, the domain predicates in the source UREs uniquely determine a  $\mathbb{P}$ -partition, which is defined by (3.5), (3.6) and (3.7), where convex domain partitions  $\text{dompar}(D)$  for every  $D \in \text{deq}(\Phi)$  are specified in the source UREs. To distinguish this partition from the others, we denote it by  $\mathbb{P}_{\text{src}}$ .

**Remark** From now on, when using the source UREs for illustration, we shall describe them by the sets  $\text{condeq}(V)$  for all variables  $V \in \mathcal{V}$  in the source UREs. We continue to use diagrams to depict these sets. In the diagram for a variable  $V$ , all



solid lines form the partition  $\text{deq}(V)$ , all dashed lines inside the region that depicts a non-convex polytope  $D \in \text{deq}(V)$  form the convex domain partition  $\text{dompar}(D)$  and, consequently, all solid and dashed lines form the partition  $\text{condeq}(V)$ .  $\square$

### 3.3.1 The Computation Control Flow

The construction of the CCUREs imposes a  $\textcircled{P}$ -partition by means of the computation control variables. To do so, we must find a set of separating hyperplanes that mutually separate all  $\textcircled{P}$ -blocks. Each of these separating hyperplanes divides  $\mathbf{Z}^n$  into two halves. We write  $\mathfrak{S}$  for the set of half-spaces associated with these hyperplanes, one half-space per hyperplane, and define  $\mathfrak{S}^c = \{\complement_{\mathbf{Z}^n} S \mid S \in \mathfrak{S}\}$ , i.e.,  $\mathfrak{S}^c$  contains the complements of the half-spaces of  $\mathfrak{S}$ .

**Definition 3.2** A set  $S$  of half-spaces is called a set of *defining half-spaces* of an integral convex polytope  $C$  if  $C$  is the intersection of these half-spaces, i.e., if  $C = (\bigcap H : H \in S : H)$ .  $\square$

**Definition 3.3** A set  $\mathfrak{S}$  of half-spaces is called a *separation set* (of a  $\textcircled{P}$ -partition) if  $\mathfrak{S} \cup \mathfrak{S}^c$  contains a set of defining half-spaces for every  $\textcircled{P}$ -block.  $\square$

Since  $\textcircled{P}$ -blocks are disjoint convex polytopes, the existence of a separation set is guaranteed by the Separation Theorem. In Sect. 3.3.1.1, we present the specification of computation control flow from a separation set for a fixed but arbitrary  $\textcircled{P}$ -partition. In Sect. 3.3.1.2, we describe the construction of a separation set for the  $\textcircled{P}$ -partition  $\textcircled{P}_{\text{src}}$ . The way that this separation set is constructed permits the development of a mechanisable procedure for the specification of computation control flow directly from the source UREs.

#### 3.3.1.1 Specifying Computation Control Flow from a Separation Set

The construction of the CCUREs from a separation set  $\mathfrak{S}$  is straightforward. For  $S \in \mathfrak{S}$ , the notation  $S^\equiv$  stands for the corresponding hyperplane. We associate a distinct computation control variable,  $C_i$ , with  $S_i^\equiv = [\pi_i : \delta_i]$  for every  $S_i \in \mathfrak{S}$ ,



where  $\pi_i \in \mathbf{Z}^n$  is normalised and  $\delta_i \in \mathbf{Z}$ . The associated control dependence vector  $\vartheta_{C_i}$  of  $C_i$  is a solution of  $\vartheta_{C_i} \pi_i = 0$ .  $C_i$  takes on two different values:  $c_i$  in the half-space  $S_i$  and  $\bar{c}_i$  in the half-space  $\mathcal{C}_{\mathbf{Z}^n} S_i$ . That is,  $\text{sig}(C_i) = \{c_i, \bar{c}_i\}$ . For convenience, we define  $c_i = \mathbf{1}$  and  $\bar{c}_i = \mathbf{0}$ . (This is data pipelining (Sect. 2.6). Instead of performing tests  $I \in S_i$  and  $I \in \mathcal{C}_{\mathbf{Z}^n} S_i$  for every point in the extended index space, the results of these tests at the corresponding half-space are shared.)

The CCUREs are defined as follows:

$$\left( \forall S_i : S_i \in \mathfrak{S} : \begin{cases} \vartheta_{C_i} \pi_i = 0 \\ C_i(I) = \begin{cases} I \in \text{in}(\Psi, \vartheta_{C_i}) \cap S_i & \rightarrow C_i(I) = c_i \\ I \in \text{in}(\Psi, \vartheta_{C_i}) \cap \mathcal{C}_{\mathbf{Z}^n} S_i & \rightarrow C_i(I) = \bar{c}_i \\ I \in \Psi & \rightarrow C_i(I) = C_i(I - \vartheta_{C_i}) \end{cases} \end{cases} \right) \quad (3.8)$$

**Remark** We abuse the notion of UREs slightly here. The equation  $\vartheta_{C_i} \pi_i = 0$  is not a URE; it specifies that  $\vartheta_{C_i}$  is any vector that is orthogonal to  $\pi_i$ .  $\square$

Once the CCUREs are constructed, we can proceed to replace the domain predicates in the extended source UREs that are also in the source UREs by predicates in the computation control variables.

**Definition 3.4** Let  $S \subseteq \Psi$ .  $\mathcal{P}(S)$  denotes the predicate in computation (pipelining) control variables that replaces predicate  $I \in S$ .  $S'$  contains the points of  $\Psi$  at which  $\mathcal{P}(S)$  holds:

$$S' = \{I \mid I \in \Psi \wedge \mathcal{P}(S)\} \quad \square$$

The definition of  $\mathcal{P}(S)$  aims at the establishment of the equality of  $S$  and  $S'$ .

**Definition 3.5** The computation control flow is *correct* if  $(\forall D : D \in \text{deq}(\Phi) : D = D')$ .  $\square$

The following lemma permits a hierarchical definition of predicates in control variables for the purpose of replacing the domain predicates of the source UREs.

**Lemma 3.3** Let  $S_1, S_2, \dots, S_s$  be subsets of a set  $S \subseteq \Psi$ . Let  $\mathcal{P}(S) = (\exists i : 0 < i \leq s : \mathcal{P}(S_i))$ . Then (1)  $S' = (\bigcup i : 0 < i \leq s : S'_i)$ . (2)  $S = (\bigcup i : 0 < i \leq s : S'_i) \iff S = S'$ .



**Proof** Since (1) implies (2), it suffices to prove (1).

$$\begin{aligned}
& S' \\
= & \{\text{Def. 3.4}\} \\
& \{I \mid I \in \Psi \wedge \mathcal{P}(S)\} \\
= & \{\mathcal{P}(S) = (\exists i : 0 < i \leq s : \mathcal{P}(S_i)) \text{ by hypothesis}\} \\
& \{I \mid I \in \Psi \wedge (\exists i : 0 < i \leq s : \mathcal{P}(S_i))\} \\
= & \{\text{Algebraic manipulation}\} \\
& (\bigcup i : 0 < i \leq s : \{I \mid I \in \Psi \wedge \mathcal{P}(S_i)\}) \\
= & \{\text{Def. 3.4}\} \\
& (\bigcup i : 0 < i \leq s : S'_i) \quad \square
\end{aligned}$$

We need to define predicates  $\mathcal{P}(D)$  in order to replace domain predicate  $I \in D$  for  $D \in \text{deq}(\Phi)$ . By Lemma 3.1,  $D/\mathbb{T} \subseteq \Phi/\mathbb{T}$ , i.e., all the elements of  $D/\mathbb{T}$  are  $\mathbb{T}$ -classes.  $\mathcal{P}(D)$  is defined in terms of predicates  $\mathcal{P}(S)$  for the  $\mathbb{T}$ -classes  $S \in D/\mathbb{T}$ :

$$\mathcal{P}(D) = (\exists S : S \in D/\mathbb{T} : \mathcal{P}(S)) \quad (3.9)$$

**Lemma 3.4**  $(\forall D : D \in \text{deq}(\Phi) : D = D') \iff (\forall C : C \in \Phi/\mathbb{T} : C = C')$ .

**Proof** Sufficiency follows from Lemma 3.3 by using (3.9) and the hypothesis. Let us prove necessity. For every  $C \in \Phi/\mathbb{T}$ , either of the following must be true by the definition of  $\mathbb{T}$ : (1)  $C \in \text{deq}(\Phi)$ , and (2) there are at least two elements  $X$  and  $Y$  in  $\text{deq}(\Phi)$  such that  $C \subset X$ ,  $C \subset Y$ ,  $C = X \cap Y$ ,  $X \setminus Y \neq \emptyset$  and  $Y \setminus X \neq \emptyset$ . In (1),  $C = C'$  by the hypothesis. In (2), we assume that  $C \neq C'$ . By (1) of Lemma 3.3 and (3.9),  $X' = (\bigcup S : S \in X/\mathbb{T} : S')$  and  $Y' = (\bigcup S : S \in Y/\mathbb{T} : S')$ . Since  $C \subset X$  and  $C \subset Y$ , we obtain  $C \in X/\mathbb{T}$  and  $C \in Y/\mathbb{T}$ . We consider two subcases: (a)  $C' \setminus C \neq \emptyset$ , and (b)  $C' \subset C$ . In (a),  $C'$  cannot be contained in both  $X$  and  $Y$ , since  $X \setminus Y \neq \emptyset$  and  $Y \setminus X \neq \emptyset$ . Hence, either  $X \neq X'$  or  $Y \neq Y'$ . Let us consider (b).  $C' \subset C$  implies  $C \setminus C' \neq \emptyset$ . In order for  $X = X'$  to hold, there must be some  $S \neq C$  in  $X/\mathbb{T}$  such that  $S'$  contain some elements in  $C \setminus C'$ . This implies that  $S' \setminus S \neq \emptyset$ .  $S$  is a  $\mathbb{T}$ -class. Regarding  $S$  as the previous  $C$ , we are back to (a) via (2).  $\square$

Next, we define predicates  $\mathcal{P}(\Phi_i)$  for all  $\mathbb{T}$ -classes in  $\Phi/\mathbb{T} = \{\Phi_i \mid 0 < i \leq r\}$ .  $\mathbb{T}$ -class  $\Phi_i$  is the union of the convex polytopes in partition  $\text{clapar}(\Phi_i)$ . Hence,



$\mathcal{P}(\Phi_i)$  is defined in terms of predicates  $\mathcal{P}(S)$  for the blocks  $S \in \text{clapar}(\Phi_i)$ :

$$\mathcal{P}(\Phi_i) = (\exists S : S \in \text{clapar}(\Phi_i) : \mathcal{P}(S)) \quad (3.10)$$

**Lemma 3.5**  $(\forall C : C \in \Phi/\oplus : C = C') \iff (\forall C : C \in \Phi/\oplus : C = (\cup S : S \in \text{clapar}(C) : S'))$ .

**Proof** Lemma 3.3 and (3.10). □

Finally, we define the predicates  $\mathcal{P}(S)$  for all the blocks in  $\text{clapar}(\Phi_i)$ . All these blocks are convex polytopes. Let  $X \subseteq \Phi$  be a  $\oplus$ -block. We define  $\mathcal{H}(X, \mathfrak{S})$  to be any subset of  $\mathfrak{S} \cup \mathfrak{S}^c$  that is a set of defining half-spaces of  $X$ . If no such set exists,  $\mathcal{H}(X, \mathfrak{S}) = \emptyset$ .  $\mathcal{P}(X)$  is given by

$$\mathcal{P}(X) = \left( \forall S_i : S_i \in \mathcal{H}(X, \mathfrak{S}) : \begin{cases} \text{if } S_i \in \mathfrak{S} \rightarrow C_i(I) = c_i \\ \square S_i \in \mathfrak{S}^c \rightarrow C_i(I) = \bar{c}_i \\ \text{fi} \end{cases} \right) \quad (3.11)$$

The next lemma establishes the containment of a set of defining half-spaces of  $X$  in  $\mathfrak{S} \cup \mathfrak{S}^c$  as a necessary and sufficient condition for the equality of  $X$  and  $X'$ .

**Lemma 3.6** *Let  $X \subseteq \Phi$  be a convex polytope. Then  $\mathcal{H}(X, \mathfrak{S}) \neq \emptyset \iff X = X'$ .*

**Proof**

$$\begin{aligned} & X \\ \iff & \{ \mathcal{H}(X, \mathfrak{S}) \text{ is a set of defining half-spaces of } X; \text{ Def. 3.2} \} \\ & (\cap S_i : S_i \in \mathcal{H}(X, \mathfrak{S}) : S_i) \\ \iff & \{ \mathcal{P}(X) \text{ of (3.11); the CCUREs of (3.8)} \} \\ & \{ I \mid I \in \Psi \wedge \mathcal{P}(X) \} \\ \iff & \{ \text{Def. 3.4} \} \\ & X' \end{aligned} \quad \square$$

We have completed the definition of predicates  $\mathcal{P}(D)$  for every  $D \in \text{deq}(\Phi)$ . Let us have a look at the situation when these definitions can be simplified. If  $S$  in  $\text{deq}(\Phi)$  ( $\Phi/\oplus$ ) is a convex polytope, predicate  $\mathcal{P}(S)$  defined by (3.9) (by (3.10)) can be defined directly by (3.11).



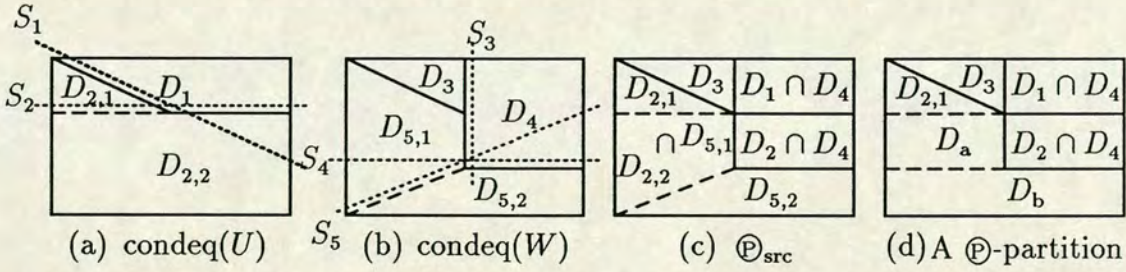


**Lemma 3.7** Let  $S$  in  $\text{deq}(\Phi)$  or  $\Phi/\mathbb{T}$  be a convex polytope. If  $\mathfrak{S}$  is a separation set, then  $\mathcal{H}(S, \mathfrak{S}) \neq \emptyset$ .

**Proof** If  $S \in \text{deq}(\Phi)$ , then  $S = (\cup X : X \in S/\mathbb{T} : (\cup Y : Y \in \text{clapar}(X) : Y))$  by Lemma 3.1 and the definition of a  $\mathbb{T}$ -partition. If  $\mathfrak{S}$  is a separation set,  $\mathcal{H}(P, \mathfrak{S}) \neq \emptyset$  when  $P$  is a  $\mathbb{T}$ -block. Thus,  $(\forall X : X \in S/\mathbb{T} : (\forall Y : Y \in \text{clapar}(X) : \mathcal{H}(Y, \mathfrak{S})))$  contains a set of defining half-spaces of  $S$ . The proof is similar when  $S \in \Phi/\mathbb{T}$ .  $\square$

By definition, the blocks in  $\text{dompar}(D)$  for all non-convex domains  $D \in \text{deq}(\Phi)$ , i.e., in  $\text{condeq}(\Phi) \setminus \text{deq}(\Phi)$ , are convex polytopes. However, Lemma 3.7 may not hold for every one of them, as the following example demonstrates. This example also provides insight to the formulation of the sufficient condition for Lemma 3.7 to hold for all these blocks.

**Example 3.2** Figs. 3–1(a) and (b) depict the source UREs as specified in Ex. 3.1 with  $\text{dompar}(D_2) = \{D_{2,1}, D_{2,2}\}$  and  $\text{dompar}(D_5) = \{D_{5,1}, D_{5,2}\}$ .



**Figure 3–2:** Domain predicate replacement.  $\text{deq}(\Phi) = \{D_1, D_2, D_3, D_4, D_5\}$ .  $\text{condeq}(\Phi) = \{D_1, D_{2,1}, D_{2,2}, D_3, D_4, D_{5,1}, D_{5,2}\}$ .  $S_1$  and  $S_2$  are supporting half-spaces to  $D_1$ ,  $S_3$  and  $S_4$  to  $D_4$ , and  $S_5$  to  $D_{5,1}$ .

Fig. 3–2(c) depicts the  $\mathbb{T}$ -partition  $\mathbb{T}_{\text{src}}$ . For simplicity, we shall not include the supporting half-spaces of the index space in a separation set.  $\mathfrak{S} = \{S_1, S_2, S_3, S_4\}$  is a separation set of the  $\mathbb{T}$ -partition depicted in Fig. 3–2(d), since

$$\begin{aligned}
 \mathcal{H}(D_{2,1}, \mathfrak{S}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_1, S_2\} \\
 \mathcal{H}(D_3, \mathfrak{S}) &= \{S_1, \mathbb{C}_{\mathbb{Z}^n} S_3\} \\
 \mathcal{H}(D_1 \cap D_4, \mathfrak{S}) &= \{S_2, S_3\} \\
 \mathcal{H}(D_2 \cap D_4, \mathfrak{S}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_2, S_3, S_4\} \\
 \mathcal{H}(D_a, \mathfrak{S}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_2, \mathbb{C}_{\mathbb{Z}^n} S_3, S_4\} \\
 \mathcal{H}(D_b, \mathfrak{S}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_4\}
 \end{aligned} \tag{3.12}$$



Lemma 3.7 holds for convex domains  $D_1$ ,  $D_3$  and  $D_4$  in  $\text{deq}(\Phi)$ , since

$$\begin{aligned}\mathcal{H}(D_1, \mathfrak{S}) &= \{S_1, S_2\} \\ \mathcal{H}(D_3, \mathfrak{S}) &= \{S_1, \mathbb{L}_{\mathbb{Z}^n} S_3\} \\ \mathcal{H}(D_4, \mathfrak{S}) &= \{S_3, S_4\}\end{aligned}\tag{3.13}$$

$D_2$  and  $D_5$  are non-convex domains in  $\text{deq}(\Phi)$ .  $\text{condeq}(\Phi) \setminus \text{deq}(\Phi) = \{D_{2,1}, D_{2,2}, D_{5,1}, D_{5,2}\}$ . Lemma 3.7 holds for  $D_{2,1}$  and  $D_{2,2}$ , since

$$\begin{aligned}\mathcal{H}(D_{2,1}, \mathfrak{S}) &= \{\mathbb{L}_{\mathbb{Z}^n} S_1, S_2\} \\ \mathcal{H}(D_{2,2}, \mathfrak{S}) &= \{\mathbb{L}_{\mathbb{Z}^n} S_2\}\end{aligned}\tag{3.14}$$

But it does not hold for  $D_{5,1}$  and  $D_{5,2}$  since  $\mathcal{H}(D_{5,1}, \mathfrak{S}) = \emptyset$  and  $\mathcal{H}(D_{5,2}, \mathfrak{S}) = \emptyset$ . Thus, we cannot define predicates  $\mathcal{P}(D_{5,1})$  and  $\mathcal{P}(D_{5,2})$  by (3.11).

The set  $\mathfrak{S}$  is not a separation set of  $\mathbb{P}_{\text{src}}$ , since  $\mathcal{H}(D_{2,2} \cap D_{5,1}, \mathfrak{S}) = \emptyset$  and  $\mathcal{H}(D_{5,2}, \mathfrak{S}) = \emptyset$ . Adding  $S_5$  to it gives rise to a separation set of  $\mathbb{P}_{\text{src}}$ :  $\mathfrak{S} = \{S_1, S_2, S_3, S_4, S_5\}$ , since, in addition to the first four equations of (3.12), we have

$$\begin{aligned}\mathcal{H}(D_{2,2} \cap D_{5,1}, \mathfrak{S}) &= \{\mathbb{L}_{\mathbb{Z}^n} S_2, \mathbb{L}_{\mathbb{Z}^n} S_3, S_5\} \\ \mathcal{H}(D_{5,2}, \mathfrak{S}) &= \{\mathbb{L}_{\mathbb{Z}^n} S_4, \mathbb{L}_{\mathbb{Z}^n} S_5\}\end{aligned}$$

For the new separation set, Lemma 3.7 is valid for every element in  $\text{condeq}(\Phi)$ , since, in addition to (3.13) and (3.14), we have

$$\begin{aligned}\mathcal{H}(D_{5,1}, \mathfrak{S}) &= \{\mathbb{L}_{\mathbb{Z}^n} S_1, \mathbb{L}_{\mathbb{Z}^n} S_3, S_5\} \\ \mathcal{H}(D_{5,2}, \mathfrak{S}) &= \{\mathbb{L}_{\mathbb{Z}^n} S_4, \mathbb{L}_{\mathbb{Z}^n} S_5\}\end{aligned}$$

**Lemma 3.8** *If  $\mathfrak{S}$  is a separation set of  $\mathbb{P}_{\text{src}}$ ,  $\mathcal{H}(S, \mathfrak{S}) \neq \emptyset$  for every  $S \in \text{condeq}(\Phi)$ .*

**Proof** By (3.7), every  $S \in \text{condeq}(\Phi)$  is a union of  $\mathbb{P}_{\text{src}}$ -blocks. The rest of proof proceeds exactly as for Lemma 3.7.  $\square$

This lemma allows us to define predicate  $\mathcal{P}(D)$  for every  $D \in \text{condeq}(\Phi)$  by (3.11). So, predicate  $\mathcal{P}(D)$ , which replaces domain predicate  $I \in D$  for every  $D \in \text{deq}(\Phi)$ , is given by

$$\mathcal{P}(D) = (\exists S : S \in \text{dompar}(D) : \mathcal{P}(S))\tag{3.15}$$

Lemma 3.8 is the basis for the procedure, presented in Sect. 3.3.1.2, for the synthesis of computation control from the source UREs.



Replacing domain predicates  $I \in D_i$  in (3.2) by the new predicates  $\mathcal{P}(D_i)$  yields

$$V(I) = \{ I \in \Psi \rightarrow \left\{ \begin{array}{ll} \text{if } I \in \Phi \wedge \mathcal{P}(D_1) & \rightarrow f_1(W(I-\vartheta_1), \dots) \\ \square I \in \Phi \wedge \mathcal{P}(D_2) & \rightarrow f_2(W(I-\vartheta_2), \dots) \\ & \dots \\ \square I \in \Phi \wedge \mathcal{P}(D_p) & \rightarrow f_p(W(I-\vartheta_p), \dots) \\ \square I \in \mathbb{C}_\Psi \Phi & \rightarrow V(I-\vartheta) \\ \text{fi} & \end{array} \right. \quad (3.16)$$

The new predicates are part of the guards in this equation. There are as many guards as there are domain predicates in the original equation. From now on, when we speak of the guards in the pipelined UREs, we mean the predicates  $\mathcal{P}(D)$ . Note that the resulting extended source UREs after domain predicate replacement and the CCUREs have the same domain predicate  $I \in \Psi$ . They are unconditional. Therefore, deriving control signals for the source UREs amounts to transforming the source UREs to a system of unconditional UREs such that the input and output data prescribed by the new system are mapped to border cells.

**Theorem 3.2** *If  $\mathfrak{S}$  is a separation set of some  $\mathbb{P}$ -partition, then the computation control flow is correct.*

**Proof** We consider only predicates  $\mathcal{P}(D)$  for  $D \in \text{deq}(\Phi)$  that are defined by (3.9). The proof is similar if they are defined by (3.15).

$$\begin{aligned} & \mathfrak{S} \text{ is a separation set of some } \mathbb{P}\text{-partition} \\ \iff & \{\text{Def. 3.2; Def. 3.3}\} \\ & (\forall X : X \in \mathbb{P} : (\exists S : S \subseteq \mathfrak{S} \cup \mathfrak{S}^c \wedge (S \text{ is a set of defining half-spaces of } X))) \\ \iff & \{\text{Choose } \mathcal{H}(X, \mathfrak{S}) \text{ as } S; X \text{ is convex and } \mathcal{H}(X, \mathfrak{S}) \neq \emptyset; \text{ Lemma 3.6}\} \\ & (\forall X : X \in \mathbb{P} : X = X') \\ \implies & \{\Phi/\mathbb{T} \text{ of (3.3); } \mathcal{P}(\Phi_i) \text{ of (3.10); Def. 3.4; Lemma 3.5}\} \\ & (\forall i : 0 < i \leq r : \Phi_i = \Phi'_i) \\ \iff & \{(\forall V : V \in \mathcal{V} : (\forall D : D \in \text{deq}(V) : D/\mathbb{T} \subseteq \Phi/\mathbb{T})); \text{ Def. 3.4; Lemma 3.4}\} \\ & (\forall V : V \in \mathcal{V} : (\forall D : D \in \text{deq}(V) : D = D')) \\ \iff & \{\text{Def. 3.5}\} \end{aligned}$$

The computation control flow is correct

□



The definition of  $\mathcal{P}(D_i)$  of (3.9) amounts to redefining the domain  $D_i$  of the domain predicate  $I \in D_i$  to:

$$\left( \bigcup X : X \in D_i / \textcircled{\text{P}} : \left( \bigcup Y : Y \in \text{clapar}(X) : \left( \bigcap H : H \in \mathcal{H}(Y, \mathfrak{S}) : H \right) \right) \right)$$

A search for different separation sets is a search for different  $\textcircled{\text{P}}$ -partitions, which amounts to searching for source UREs with different specifications of domain predicates by the virtue of Lemma 3.2. Different separation sets may lead to computation control of varying quality. Take Ex. 3.2 for example. The second separation set is inferior to the first, since it contains one more element (i.e.,  $S_5$ ) and thus results in one more control variable in the CCUREs.

This thesis does not provide systematic solutions to finding the best control UREs. We shall at the end of Sect. 3.3.1.2 describe some heuristics that guide the user in the optimisation of domain predicates. The user is responsible for the choice of domain predicates that yields efficient systolic arrays. Systolic design systems [2,54,56,75] generally rely on the user in difficult decision-making and performance optimisation. In the synthesis of data flow, optimisations of step and place [46,70,80,81] are carried out with respect to the source UREs (rather than the initial problem specification). Also for their choice, the user is responsible.

In Sect. 3.3.1.2, we present a mechanisable procedure for the construction of a separation set of  $\textcircled{\text{P}}_{\text{src}}$ , denoted  $\mathfrak{S}_{\text{src}}$ , from the source UREs. The construction of  $\mathfrak{S}_{\text{src}}$  completely depends on the domain predicates of the source UREs. So does the subsequent domain predicate replacement.

### 3.3.1.2 Specifying Computation Control Flow from the Source UREs

This section describes a procedure for the synthesis of computation control from the source UREs. We assume that the domain predicates of the source UREs are in normal form (Sect. 2.2), i.e., their constituent conditionals are in one of the following forms:

$$\pi I = \delta, \pi I \neq \delta, \pi I < \delta, \pi I \leq \delta, \pi I > \delta, \pi I \geq \delta \quad (3.17)$$



For simplicity and without loss of generality, we further assume that all conditionals are neither equalities nor inequalities. If not, we can always replace every equality  $\pi I = \delta$  by  $\pi I \leq \delta \wedge \pi I \geq \delta$ , and every inequality  $\pi I \neq \delta$  by  $\pi I < \delta \vee \pi I > \delta$ . App. A discusses the normalisation of domain predicates.

Each conditional,  $\pi \sqsupset \delta$ , in one of the rightmost four forms depicted in (3.17) corresponds to half-space  $[\pi \sqsupset \delta]$ . The idea underlying the construction of separation set  $\mathfrak{S}_{\text{src}}$  is to just to collect the half-spaces corresponding the conditionals in the domain predicates. If two half-spaces are complements of each other, it is only necessary to put one of them in  $\mathfrak{S}_{\text{src}}$  ( $\mathfrak{S}_{\text{src}}^c$  will contain the other). Once  $\mathfrak{S}_{\text{src}}$  is obtained, the CCUREs follow from (3.8) with  $\mathfrak{S} = \mathfrak{S}_{\text{src}}$ . The replacement of the domain predicates of the source UREs relies on (3.15). That is, we define predicate  $\mathcal{P}(D)$  for every  $D \in \text{condeq}(\Phi)$  by finding a set of defining half-spaces of  $D$  from  $\mathfrak{S}_{\text{src}} \cup \mathfrak{S}_{\text{src}}^c$ . For every conditional  $\pi \sqsupset \delta$  in the domain predicates, the corresponding half-space  $[\pi \sqsupset \delta]$  is contained in  $\mathfrak{S}_{\text{src}} \cup \mathfrak{S}_{\text{src}}^c$ . By the specification of the CCUREs, any two complement half-spaces in  $\mathfrak{S}_{\text{src}} \cup \mathfrak{S}_{\text{src}}^c$  are associated with a unique control variable; the variable takes on the value **1** at the points in the half-space in  $\mathfrak{S}_{\text{src}}$  and the value **0** at the points in the half-space in  $\mathfrak{S}_{\text{src}}^c$ . Let  $C$  be the control variable associated with  $[\pi \sqsupset \delta]$ . We replace conditional  $\pi \sqsupset \delta$  by predicate  $C(I) = \mathbf{1}$  if  $[\pi \sqsupset \delta]$  is in  $\mathfrak{S}$  and predicate  $C(I) = \mathbf{0}$  if  $[\pi \sqsupset \delta]$  is in  $\mathfrak{S}_{\text{src}}^c$ . This amounts to defining  $\mathcal{H}(D, \mathfrak{S}_{\text{src}})$  for  $D \in \text{condeq}(\Phi)$  as the set of half-spaces that correspond to the constituent conditionals of  $I \in D$ . Trivially,  $\mathcal{H}(D, \mathfrak{S}_{\text{src}})$  is a set of half-spaces of  $C$ .

From now on, we shall use  $C.S$  to denote the control variable associated with a half-space  $S$  in  $\mathfrak{S} \cup \mathfrak{S}^c$ . Conversely, we use  $H.C$  ( $H^c.C$ ) to denote the half-space in  $\mathfrak{S}$  ( $\mathfrak{S}^c$ ) associated with a control variable  $C$ .

**Procedure 3.1** (Construction of computation control flow from the source UREs)

INPUT:           The source UREs.

OUTPUT:          The specification of computation control flow.

1.  $\mathfrak{S}_{\text{src}} := \emptyset$ .



2. Replace every conditional  $\pi I < \delta$ ,  $\pi I \leq \delta$ ,  $\pi I > \delta$ , or  $\pi I \geq \delta$  in the domain predicates by  $-\pi I > -\delta$ ,  $-\pi I \geq -\delta$ ,  $-\pi I < -\delta$ , or  $-\pi I \leq -\delta$  if the first non-zero component of  $\pi$  is not positive. (For fixed  $\pi$  and  $\delta$ , this prevents simultaneous occurrences of  $\pi I \leq \delta$  and  $-\pi I \geq -\delta$  or  $\pi I < \delta$  and  $-\pi I > -\delta$ .)
3. For every conditional  $\pi \square \delta$  in the domain predicates:

$$\mathfrak{S}_{\text{src}} := \begin{cases} \text{if } \square = < \rightarrow \mathfrak{S}_{\text{src}} + \{I \mid \pi I < \delta\} \\ \square \square = \leq \rightarrow \mathfrak{S}_{\text{src}} + \{I \mid \pi I \leq \delta\} \\ \square \square = > \rightarrow \mathfrak{S}_{\text{src}} + \{I \mid \pi I < \delta\} \\ \square \square = \geq \rightarrow \mathfrak{S}_{\text{src}} + \{I \mid \pi I \leq \delta\} \\ \text{fi} \end{cases}$$

4. The CCUREs are given by (3.8) with  $\mathfrak{S} = \mathfrak{S}_{\text{src}}$ .
5. Replace the domain predicates of the source UREs by predicates in computation control variables. A conditional  $\pi \square \delta$  is replaced by  $C.[\pi \square \delta](I) = 1$  if  $[\pi \square \delta] \in \mathfrak{S}$  and by  $C.[\pi \square \delta](I) = 0$  otherwise.  $\square$

The choice of  $\mathfrak{S}_{\text{src}}$  is unique because the first two steps of Proc. 3.1 are deterministic. So is the specification of computation control flow.

**Lemma 3.9**  $\mathfrak{S}_{\text{src}}$  is a separation set of  $\mathbb{P}$ -partition  $\mathbb{P}_{\text{src}}$ .

**Proof** We need to show that  $\mathfrak{S}_{\text{src}} \cup \mathfrak{S}_{\text{src}}^c$  contains a set of defining half-spaces for every  $S \in \mathbb{P}_{\text{src}}$ . By construction,  $\mathfrak{S}_{\text{src}} \cup \mathfrak{S}_{\text{src}}^c$  contains  $[\pi I \square \delta]$  for every conditional  $\pi I \square \delta$  in the domain predicates. Choose  $\mathcal{H}(S, \mathfrak{S}_{\text{src}})$  as the half-spaces that are in  $\mathfrak{S}_{\text{src}} \cup \mathfrak{S}_{\text{src}}^c$  and that correspond to the constituent conditionals of  $I \in S$ .  $\square$

**Theorem 3.3** The computation control flow specified by Proc. 3.1 is correct.

**Proof** Lemma 3.6, Lemma 3.8 and Thm. 3.2.  $\square$

Each conditional of the domain predicates corresponds to a distinct control variable. The more conditionals there are in the domain predicates, the more control variables there are in the CCUREs. We describe some heuristics that



guide the user in the minimisation of the number of conditionals when he/she is searching for URE specifications from the initial problem specification.

**Procedure 3.2** (Construction of the domain predicates of the source UREs)

INPUT: The set  $\text{deq}(\Phi)$  of the domains of all equations of the source UREs.

OUTPUT: The specification of the domain predicates of the source UREs.

1.  $H := \emptyset$ .
2. For every convex domain  $D$  in  $\text{deq}(\Phi)$ , choose a set of defining half-spaces of  $D$  from  $H$  whenever possible, and add to  $H$  the new half-spaces and their complements.
3. For every non-convex domain  $D$  in  $\text{deq}(\Phi)$ , use the half-spaces in  $H$  or their translates to form a convex domain partition of  $D$  whenever possible, and add to  $H$  the new half-spaces and their complements.  $\square$

The motivation for using translates of the half-spaces in convex domain partitions is that we can always combine the control variables whose associated hyperplanes are parallel to one control variable (Sect. 3.4).

Let us use this procedure to obtain convex domain partitions for non-convex domains  $D_2$  and  $D_5$  in Ex. 3.2. Assume that  $H = \{S_1, S_2, S_3, S_4\}$  when Step 2 of Proc. 3.2 completes. The choice of  $\text{dompar}(D_2) = \{D_{2,1}, D_{2,2}\}$  in Ex. 3.2 is permitted by the procedure. It amounts to partitioning  $D_2$  to  $\{\mathbb{C}_{\mathbf{Z}^n} S_1 \cap S_2, \mathbb{C}_{\mathbf{Z}^n} S_2\}$ . But, the choice of  $\text{dompar}(D_5) = \{D_{5,1}, D_{5,2}\}$  is not permitted. This is because we can use the half-spaces in  $H$  to obtain a convex domain partition for  $D_5$ . One possible solution is  $\text{dompar}(D_5) = \{D_{2,1} \cup D_a, D_b\} = \{\mathbb{C}_{\mathbf{Z}^n} S_1 \cap \mathbb{C}_{\mathbf{Z}^n} S_3 \cap S_4, \mathbb{C}_{\mathbf{Z}^n} S_4\}$ . We thus get rid of the conditional corresponding to  $S_5$ .

### 3.3.2 The Propagation Control Flow

The construction of the PCUREs proceeds along similar lines. We redefine the partition  $\{\Phi, \mathbb{C}_{\Psi} \Phi\}$  of the extended index space  $\Psi$ . To do so, we must find a set of



defining half-spaces of  $\Phi$ :  $\mathfrak{B} = \{\text{sup}(\Phi, F) \mid F \in \text{facets}(\Phi)\}$ .  $\mathfrak{B}$  plays the same rôle as  $\mathfrak{S}$  previously. We associate a distinct propagation control variable,  $P_i$ , with  $B_i = [\pi_i : \delta_i]$  for every  $B_i \in \mathfrak{B}$  ( $0 < i \leq |\mathfrak{B}|$ ), where  $\pi_i \in \mathbb{Z}^n$  is normalised and  $\delta_i \in \mathbb{Z}$ . The associated control dependence vector  $\vartheta_{P_i}$  of  $B_i$  is a solution of  $\vartheta_{P_i} \pi_i = 0$ . Each propagation control variable takes on two different values:  $p_i$  in the half-space  $B_i$  and  $\bar{p}_i$  in the half-space  $\mathbb{C}_{\mathbb{Z}^n} B_i$ . That is,  $\text{sig}(P_i) = \{p_i, \bar{p}_i\}$ . For convenience, we define  $p_i = 1$  and  $\bar{p}_i = 0$ .

The PCUREs are defined as follows:

$$\left( \forall B_i : B_i \in \mathfrak{B} : \begin{cases} \vartheta_{P_i} \pi_i = 0 \\ P_i(I) = \begin{cases} I \in \text{in}(\Psi, \vartheta_{P_i}) \cap B_i & \rightarrow P_i(I) = p_i \\ I \in \text{in}(\Psi, \vartheta_{P_i}) \cap \mathbb{C}_{\mathbb{Z}^n} B_i & \rightarrow P_i(I) = \bar{p}_i \\ I \in \Psi & \rightarrow P_i(I) = P_i(I - \vartheta_{P_i}) \end{cases} \end{cases} \right) \quad (3.18)$$

Predicate  $\mathcal{P}(\Phi)$  and  $\mathcal{P}(\mathbb{C}_{\Psi} \Phi)$  replace predicate  $I \in \Phi$  and  $I \in \mathbb{C}_{\Psi} \Phi$ , respectively:

$$\begin{aligned} \mathcal{P}(\Phi) &= (\forall i : 0 < i \leq |\mathfrak{B}| : P_i(I) = p_i) \\ \mathcal{P}(\mathbb{C}_{\Psi} \Phi) &= (\exists i : 0 < i \leq |\mathfrak{B}| : P_i(I) = \bar{p}_i) \end{aligned} \quad (3.19)$$

$\{\Phi', \mathbb{C}_{\Psi} \Phi'\}$  is the partition defined by the PCUREs:

$$\begin{aligned} \Phi' &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi)\} \\ \mathbb{C}_{\Psi} \Phi' &= \{I \mid I \in \Psi \wedge \mathcal{P}(\mathbb{C}_{\Psi} \Phi)\} \end{aligned} \quad (3.20)$$

Replacing the domain predicates  $I \in \Phi$  and  $I \in \mathbb{C}_{\Psi} \Phi$  in the extended source UREs by predicates  $\mathcal{P}(\Phi)$  and  $\mathcal{P}(\mathbb{C}_{\Psi} \Phi)$  changes (3.16) to

$$V(I) = \{ I \in \Psi \rightarrow \begin{cases} \text{if } \mathcal{P}(\Phi) \wedge \mathcal{P}(D_1) & \rightarrow f_1(W(I - \vartheta_1), \dots) \\ \square \mathcal{P}(\Phi) \wedge \mathcal{P}(D_2) & \rightarrow f_2(W(I - \vartheta_2), \dots) \\ & \dots \\ \square \mathcal{P}(\Phi) \wedge \mathcal{P}(D_p) & \rightarrow f_p(W(I - \vartheta_p), \dots) \\ \square \mathcal{P}(\mathbb{C}_{\Psi} \Phi) & \rightarrow V(I - \vartheta) \\ \text{fi} \end{cases} \quad (3.21)$$

**Definition 3.6** The propagation control flow is *correct* if  $\Phi = \Phi'$  and  $\mathbb{C}_{\Psi} \Phi = \mathbb{C}_{\Psi} \Phi'$ .



**Theorem 3.4** *The propagation control flow as defined by (3.18), (3.19) and (3.21) is correct iff  $\mathfrak{B}$  is a set of defining half-spaces of  $\Phi$ .*

**Proof** Similar to the proof of Thm. 3.2. □

**Definition 3.7** The control flow (the specification of the pipelined UREs) is *correct* if the computation control and the propagation control flow are correct.

The PCUREs serve to propagate data from the internal cells to (from) the border cells. The half-space  $\text{sup}(\Phi, F)$  in  $\mathfrak{B}$  is superfluous if all data variables at the pipelining points in  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} \text{sup}(\Phi, F)$  are undefined. That is, control signals are unnecessary if no input/output data are pipelined across domain  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} \text{sup}(\Phi, F)$ .

The fact that valid space-time mappings are bijections from  $\Psi$  to  $\bar{\Psi}$  gives rise to the following lemma, which allows us to regard CCUREs also as PCUREs. The image of  $X'$  under the space-time mapping is given by

$$\bar{X}' = \{\bar{I} \mid \bar{I} \in \bar{\Psi} \wedge \overline{\mathcal{P}(X)}\}$$

$\overline{\mathcal{P}(X)}$  denotes the image of predicate  $\mathcal{P}(X)$ ; it is the predicate with all the occurrence of index vectors  $I$  appearing in  $\mathcal{P}(X)$  substituted by  $\bar{I}$ .

**Lemma 3.10**  $X = X' \implies \bar{X} = \bar{X}'$ .

**Proof** The bijectivity of valid space-time mappings from  $\mathbb{Q}^n$  to  $\mathbb{Q}^n$ . □

This lemma allows us to redefine  $\mathcal{P}(\Phi)$  and  $\mathcal{P}(\mathcal{L}_{\Psi}\Phi)$  of (3.19) using computation control variables only.

$$\begin{aligned} \mathcal{P}(\Phi) &= \text{true} \\ \mathcal{P}(\mathcal{L}_{\Psi}\Phi) &= (\forall S : S \in \Phi / \oplus : \neg \mathcal{P}(S)) \end{aligned} \tag{3.22}$$

**Lemma 3.11** *The propagation control flow as defined by (3.18), (3.21) and (3.22) is correct iff  $(\forall i : 0 < i \leq r : \bar{\Phi}_i = \bar{\Phi}'_i)$ .*

**Proof** To prove necessity, we assume that  $\bar{\Phi}_k \neq \bar{\Phi}'_k$  for some  $k$ . By Lemma 3.4, there is some  $D \in \text{deq}(\Phi)$  such that  $\bar{D} \neq \bar{D}'$ . Since a valid space-time mapping is



bijjective over  $\Phi$ ,  $\overline{D} \subset \overline{D'}$ . So, some pipelining point is treated as a computation point. Let us prove sufficiency. Since  $(\forall i : 0 < i \leq r : \overline{\Phi}_i = \overline{\Phi'_i})$ , we have  $(\forall D : D \in \text{deq}(\Phi) : \overline{D} = \overline{D'})$  by Lemma 3.4. Thus.  $\overline{\mathcal{P}(D)} \implies \overline{\mathcal{P}(\Phi)}$  and  $\overline{\mathcal{P}(\mathcal{C}_\Psi\Phi)} \iff (\forall S : S \in \Phi/\mathbb{T} : \neg\overline{\mathcal{P}(S)})$ .  $\square$

**Theorem 3.5** *The propagation control flow as defined by (3.18), (3.21) and (3.22) is correct.*

**Proof** Lemmata 3.10 and 3.11.  $\square$

In the synthesis of systolic arrays of reduced dimension, a valid space-time mapping is non longer a bijection. In Chap. 5, we shall see that Lemma 3.10 and, consequently, Thm. 3.10 are not valid for one-dimensional arrays.

If the propagation control flow is specified by (3.18), (3.19) and (3.21), we can disregard the half-spaces of  $\mathfrak{S}$  that support the index space. That is, we just need to find a set  $\mathfrak{S}$  such that  $\mathfrak{S} \cup \mathfrak{S}^c$  contains a set of half-spaces for every  $\mathbb{P}$ -block such that the intersection of these half-space contains that block and does not intersect the remaining blocks. This is due to the fact that  $\Psi$  and  $\mathcal{C}_\Psi\Phi$  can be identified by propagation control.

The notations  $C.S$ ,  $H.C$  and  $H^c.C$  that are defined in Sect. 3.3.1.1 for computation control flow also carry over to propagation control flow.

## 3.4 The Optimisation of Control Flow

This section describes three techniques for the optimisation of control flow (i.e., pipelined UREs). Our objective is to eliminate redundant control hardware in the systolic array. The first two optimisations are independent of the space-time mapping, while the third is dependent on the projection vector.

The first optimisation is to merge control variables in order to reduce the number of control signals. Assume that  $\mathfrak{S} \cup \mathfrak{B}$  contains  $p$  half-spaces  $S_1, S_2, \dots, S_p$  whose corresponding hyperplanes are parallel. Instead of introducing  $p$  control



variables, one per hyperplane, we introduce only one control variable,  $C$ , for the  $p$  hyperplanes. The control dependence vector  $\vartheta_C$  is perpendicular to the normal of these hyperplanes. The  $p$  hyperplanes partition  $\text{in}(\Psi, \vartheta_C)$  into  $p+1$  subsets; each is initialised with a distinct control value, which is pipelined along the pipelining vector  $\vartheta_C$  across the extended index space  $\Psi$ . Let  $\text{sig}(C) = \{c_1, c_2, \dots, c_{p+1}\}$  and  $S_1 \subset S_2 \subset \dots \subset S_p$ . The specification of  $C$  is as follows:

$$C(I) = \begin{cases} I \in \text{in}(\Psi, \vartheta_C) \cap S_1 & \rightarrow c_1 \\ I \in \text{in}(\Psi, \vartheta_C) \cap \mathbb{C}_{\mathbf{Z}^n} S_1 \cap S_2 & \rightarrow c_2 \\ \dots & \dots \\ I \in \text{in}(\Psi, \vartheta_C) \cap \mathbb{C}_{\mathbf{Z}^n} S_p & \rightarrow c_{p+1} \end{cases} \quad (3.23)$$

This optimisation reduces the number of control signals from  $2p$  for the  $p$  variables to  $p+1$  for the variable that replaces them.

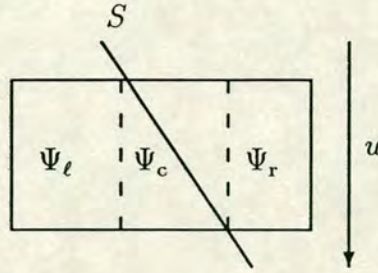
The second optimisation is to minimise the domain of each control variable. The values of a control variable  $C$  at some points may not appear in the guards the pipelined UREs. Thus,  $C$  need not be defined at these points. Let  $D_C$  be the set of all points  $I$  such that  $C(I)$  appears in the guards of the pipelined UREs.  $D_C$  is the union of all domains of equations  $D_1, D_2, \dots, D_d$  such that every predicate  $\mathcal{P}(D_i)$  refers to  $C(I)$  (see (3.9)). Thus, we can substitute  $\text{rays}(\Phi, \pm u) \cap \text{ray}(D_C, \pm \vartheta_C)$  for  $\Psi$  in the defining equation of  $C$ . The domain of  $C$  is given by

$$\Psi_C = \text{rays}(\Phi, \pm u) \cap \text{ray}(D_C, \pm \vartheta_C) \quad (3.24)$$

The CCUREs serve to distinguish different types of computation points; they are defined for the index space. If  $C$  is a computation control variable,  $\text{ray}(D_C, \pm \vartheta_C) \subseteq \text{ray}(\Phi, \pm \vartheta_C)$ . The PCUREs serve to distinguish pipelining from computation points. If  $C$  is a propagation control variable,  $\text{ray}(D_C, \pm \vartheta_C) = \text{rays}(\Psi \cup \text{H}^c.C, \pm \vartheta_C) \cup \text{ray}(\Phi, \pm \vartheta_C)$ .

**Remark** To apply Proc. 3.1, we must replace  $\pi I = \delta$  by  $\pi I \leq \delta \wedge \pi I \geq \delta$ . This implies that the corresponding control variable must be defined at the two half-spaces. If the control variable is only referenced at the points of the hyperplane  $[\pi : \delta]$ , the domain of the control variable can be reduced to the intersection of the hyperplane and the index space.  $\square$





**Figure 3–3:** The optimisation of pipelined UREs. The solid rectangle depicts the index space.  $S \in \mathfrak{S}$  supports  $\Psi_\ell$ .

The third optimisation can be applied after the projection vector  $u$  has been chosen (Fig. 3–3). It removes conditionals in the space-time UREs that become redundant due to the chosen projection vector. For control variable  $C$  (its associated half-space in  $\mathfrak{S}$  is  $S$ ), we can partition the extended index space into three blocks:

$$\begin{aligned}
 \Psi_\ell &= \{I \mid \text{ray}(I, \pm u) \subseteq S \wedge \text{ray}(I, \pm u) \not\subseteq \mathfrak{L}_{\mathbf{Z}^n} S\} \\
 \Psi_r &= \{I \mid \text{ray}(I, \pm u) \not\subseteq S \wedge \text{ray}(I, \pm u) \subseteq \mathfrak{L}_{\mathbf{Z}^n} S\} \\
 \Psi_c &= \Psi \setminus (\Psi_\ell \cup \Psi_r)
 \end{aligned} \tag{3.25}$$

Imagine that we draw lines parallel to the projection vector  $u$  in the extended index space. The three sets defined in (3.25) are the union of lines that only contain points in  $S$ , the union of lines that only contain points in  $\mathfrak{L}_{\mathbf{Z}^n} S$  and the union of the remaining lines. Thus, we can substitute  $\text{true}$  for  $C(\bar{I}) = \mathbf{1}$  at the cells of  $\text{place}(\Psi_\ell)$  and for  $C(\bar{I}) = \mathbf{0}$  at the cells of  $\text{place}(\Psi_r)$ , but we must retain predicates  $C(\bar{I}) = \mathbf{1}$  and  $C(\bar{I}) = \mathbf{0}$  at the cells of  $\text{place}(\Psi_c)$ .

If after this optimisation some guard of an equation with domain  $D$  is true, the processors at  $\text{place}(D)$  can be pre-designed to apply that equation. An extreme case occurs when the projection vector is orthogonal to the normal of the separating hyperplane  $S^\perp$  associated with control variable  $C$ . In this case,  $\Psi_c = \emptyset$ . Then the corresponding control variable is said to be *redundant*. This is because  $\Psi \cap S$  and  $\Psi \cap \mathfrak{L}_{\mathbf{Z}^n} S$  are projected to two disjoint sets of processors. They are already separated spatially and do not have to be separated temporally.

**Theorem 3.6** *If all computation (propagation) control variables are redundant under the space-time mapping, then the CCUREs (PCUREs) are unnecessary.*



**Proof** If all computation control variables are redundant, all  $\oplus$ -classes are projected to disjoint sets of processors. The CCUREs are unnecessary because a cell can be pre-designed to compute the computation points mapped to it. The PCUREs are unnecessary, if all propagation control variables are redundant. Then there are no pipelining points, i.e.,  $\Psi = \Phi$ .  $\square$

Unfortunately, the advantage of making control variables stationary does not carry over to data variables. Complications arise when some data variables become stationary. A local processor memory is required to store the stationary data allocated to the processor. The control UREs do not provide a specification for handling stationary variables. In Chap. 8, we shall present a systematic method for handling stationary variables. The basic idea is to rewrite the source UREs to encompass the specification for handling stationary variables.

### 3.5 The Extension of the Space-Time Mapping

The mapping conditions for the space-time mapping stated in Thm. 2.1 must be augmented such that the specification of the control UREs is taken into account.

**Theorem 3.7** *A space-time mapping is valid if*

- *The scheduling vector  $\lambda$  satisfies:*

$$- (\forall \vartheta : \vartheta \in \mathcal{D} : \lambda\vartheta > 0) \tag{3.26}$$

$$- (\forall S : S \in \mathfrak{S} \cup \mathfrak{B} : (\exists \vartheta : \vartheta \in \text{lin}(S^\#) : \lambda\vartheta > 0)) \tag{3.27}$$

- *The projection vector  $u$  satisfies  $\lambda u \neq 0$ .*

**Proof** The pipelined UREs are UREs. Apply Thm. 2.1.  $\square$

(3.27) means that any vector  $\vartheta$  in  $\text{lin}(S^\#)$  can be taken as the control dependence vector for the control variable associated with  $S^\#$ , provided it satisfies  $\lambda\vartheta > 0$ . If we choose control dependence vectors from the set of data dependence vectors (the control dependence vectors must be normalised), the latency of the systolic array can be retained when the PCUREs are not needed.



Let  $C^*$  be the set of scheduling vectors satisfying (3.26), and let  $C^{**}$  be the set of scheduling vectors satisfying (3.26) and (3.27). Both  $C^*$  and  $C^{**}$  are cones. First, we show that, if there is a valid scheduling vector for the source UREs, there must be valid scheduling vectors for the pipelined UREs.

**Theorem 3.8**  $C^* \neq \emptyset \implies C^{**} \neq \emptyset$ .

**Proof** Let  $\Theta^*$  be the cone generated by the data dependence vectors in  $\mathcal{D}$ . To prove that  $C^{**}$  is not empty, it suffices to prove that we can obtain a cone,  $\Theta^{**}$ , which contains  $\Theta^*$  and one vector from  $\text{lin}(S^-)$  for every  $S \in \mathfrak{S} \cup \mathfrak{B}$ . Take any non-zero vector  $\vartheta \in \text{lin}(S^-)$ . Either  $\Theta^* \cup \{\vartheta\}$  or  $\Theta^* \cup \{-\vartheta\}$  must be a cone.  $\square$

Thm. 3.8 indicates that all domain predicates can be pipelined. The proof of the theorem has the following implication. If we choose all control dependence vectors from  $\Theta^*$ , then a space-time mapping that is valid for the source UREs must also be valid for the pipelined UREs.

Next, we show that the set of scheduling vectors that are valid for the source UREs but are not valid for the pipelined UREs is just the set  $N(\mathfrak{S} \cup \mathfrak{B})$  of the normals (up to some scalars) of the corresponding hyperplanes of some half-spaces in  $\mathfrak{S} \cup \mathfrak{B}$ :

$$N(\mathfrak{S} \cup \mathfrak{B}) = \{\alpha\pi \mid [\pi I \square \delta] \in \mathfrak{S} \cup \mathfrak{B} \wedge \alpha \in \mathbb{Z}\} \cap C^*$$

**Theorem 3.9**  $C^* = C^{**} \cup N(\mathfrak{S} \cup \mathfrak{B})$  and  $C^{**} \cap N(\mathfrak{S} \cup \mathfrak{B}) = \emptyset$ .

**Proof** Clearly,  $C^* \supseteq C^{**}$ . For  $S = [\pi I \square \delta]$  in  $\mathfrak{S} \cup \mathfrak{B}$ ,  $(\forall \vartheta : \vartheta \in \text{lin}(S^-) : \lambda\vartheta = 0)$  iff  $\lambda = \alpha\pi$  for some  $\alpha \in \mathbb{Z}$ .  $\square$

Thm. 3.9 indicates that the price paid for the pipelining of domain predicates is that the set of valid scheduling vectors may be reduced. Thus, in the formulation of the domain predicates, we should try to avoid conditionals that exclude desirable scheduling vectors from  $C^*$ . Finally, if we change “>” to “=” in (3.27), then  $C^* = C^{**}$ . This implies that all control signals must be broadcast across the array. It is also possible to broadcast some control signals selectively to optimise some aspects of the array, e.g., its latency or number of delay buffers.



### 3.6 Related Work

Chen [14] proposed a method that replaces a time-dependent conditional of a space-time equation by a control signal pipelined from the boundary of the array. We explain the basic idea by considering the pipelining of a generic time-dependent conditional  $bt \sqcap (\sum i : 0 < i < n : a_i p_i)$ , where  $t$  represents time and  $p_1, p_2, \dots, p_{n-1}$  represent processor coordinates. If the conditional evaluates to true at processor  $(x_1, \dots, x_j, \dots, x_{n-1})$  at step  $t_0$ , it also evaluates to true at processors  $(x_1, \dots, x_j + mb, \dots, x_{n-1})$  at steps  $t_0 + ma_j$  ( $m \in \mathbf{Z}$ ). Instead of performing the test  $bt \sqcap (\sum i : 0 < i < n : a_i p_i)$  at each of these processors, the results of these tests can be shared. This is achieved by injecting a one-bit control signal at the appropriate border cell and pipelining it along direction  $(0, \dots, 0, b/a_j, 0, \dots, 0)$  across the array at a velocity of  $b/a_j$ , i.e., passing  $b$  processors per  $a_j$  clock ticks.

This method does not apply for time-dependent conditionals of the form  $bt \sqcap a_i$ ; they do not consist of processor coordinates. In this case, the corresponding control signals must be broadcast (Thm. 3.9). This happens when the control dependence vectors associated with these conditionals are orthogonal to the scheduling vector. But, Chen's method provides no means for preventing these conditionals with the choice of the space-time mapping. Our method does away with these conditionals by the introduction of control dependence vectors. Chen's method cannot enforce design constraints on the control signals since there is no notion of control dependence vector. In the previous example conditional,  $(\forall i : 0 < i < n : b/a_i > 1)$  will result in non-neighbouring channels. The question whether we can get away with only neighbouring channels cannot be answered. Chen's method also does not address control for the propagation of input/output data that are mapped to internal cells of the array. Our method specifies this control with the PCUREs. Finally, Chen's method does not generalise to systolic arrays of reduced dimension because it relies on space-time mappings for  $(n-1)$ -dimensional arrays.

The method reported in [71] is similar to Chen's method. The method reported in [58] only deals with conditionals that are equalities.



## 3.7 Examples

This section illustrates the synthesis of control signals with two examples: dynamic programming and LU-decomposition. The construction of the PCUREs is similar to that of the CCUREs. We shall therefore restrict ourselves mainly to the latter. For the former, we only present a set  $\mathfrak{B}$  of defining half-spaces of the index space. So, the presentation of the extended source UREs is unnecessary. Proc. 3.1 is completely mechanical. We go through its steps to highlight the rôles that the various concepts and theorems introduced so far play in the construction of the control UREs. For each example, we present the  $\textcircled{1}$ -classes and the separation set  $\mathfrak{S}_{\text{src}}$  (with optimisations) returned by Proc. 3.1. We then describe the pipelined UREs (without the propagation UREs and the PCUREs). We shall omit the half-spaces in a separation set which are the supporting half-spaces to the index space. We shall depict every point of the index space by a graphical symbol (e.g.,  $\blacksquare$ ,  $\blacktriangle$ , etc). When we write  $\Phi_s$  for a graphical symbol  $s$ , we mean the set of all the points of the index space depicted by that symbol. When we label a computation equation of the source UREs by a number of graphical symbols, we mean that the equation is only defined for the points depicted by those symbols.

We use the first example to show how control complexity should be taken into account in the appraisal of the systolic array, and the second example to show how the CCUREs or PCUREs can be avoided by an appropriate choice of space-time mapping. These two examples will also be used for illustration in Chap. 5.

### 3.7.1 Dynamic Programming

We apply dynamic programming to solve the optimal string parenthesisation problem [26]. Let a string of items be indexed by integers 1 through  $m$  from left to right. A parenthesisation of a string of  $m$  items has  $m-1$  pairs; each parenthesis pair encloses two elements each of which is either an item or another parenthesis



pair. The optimal cost, denoted  $c_{i,j}$ , of parenthesising the substring consisting items  $i$  through  $j-1$  is defined recursively as follows [14,24,58]:

$$\text{Specification: } c_{i,j} = (\min k : i < k < j : c_{i,k} + c_{k,j}) + w_{i,j}.$$

where  $w_{i,j}$  is the additional cost for the outermost parenthesis pair.  $c_{1,m+1}$  is the optimal cost for parenthesising a string of  $m$  items.

### 3.7.1.1. The Source UREs

This specification can be transformed into a system of UREs using the methods described in [57,58,79]. The following UREs are adapted from [58].

UREs:

$$\begin{aligned}
 A(i,j,k) &= \begin{cases} j-i=2k \rightarrow D(i,j,k) & \blacktriangleleft \blacklozenge \\ j-i \neq 2k \rightarrow A(i,j-1,k) & \blacktriangleright \bullet \blacksquare \star \end{cases} \\
 B(i,j,k) &= \begin{cases} k=1 \rightarrow C(i+1,j,k) & \blacktriangleleft \blacktriangleright \blacksquare \\ k \neq 1 \rightarrow B(i+1,j,k-1) & \blacklozenge \bullet \star \end{cases} \\
 D(i,j,k) &= \begin{cases} k=1 \rightarrow C(i,j-1,k) & \blacktriangleleft \blacktriangleright \blacksquare \\ k \neq 1 \rightarrow D(i,j-1,k-1) & \blacklozenge \bullet \star \end{cases} \\
 E(i,j,k) &= \begin{cases} j-i=2k \rightarrow B(i,j,k) & \blacktriangleleft \blacklozenge \\ j-i \neq 2k \rightarrow E(i+1,j,k) & \blacktriangleright \bullet \blacksquare \star \end{cases} \\
 c_{1,m+1} &= C(1,m+1,1) \\
 C(i,j,k) &= \begin{cases} k=1 \wedge j-i=2k-1 \rightarrow w_{i,j} \blacksquare \\ k=1 \wedge j-i \geq 2k \rightarrow \min \begin{pmatrix} A(i,j,k) + B(i,j,k) \\ C(i,j,k+1) \\ D(i,j,k) + E(i,j,k) \end{pmatrix} + w_{i,j} \blacktriangleleft \blacktriangleright \\ k \neq 1 \wedge j-i \geq 2k \rightarrow \min \begin{pmatrix} A(i,j,k) + B(i,j,k) \\ C(i,j,k+1) \\ D(i,j,k) + E(i,j,k) \end{pmatrix} \blacklozenge \bullet \\ k \neq 1 \wedge j-i < 2k \wedge j-i \geq 2k-2 \rightarrow \infty \star \end{cases}
 \end{aligned}$$

$$\text{Index Space: } \Phi = \{(i,j,k) \mid 0 < i < j \leq m+1 \wedge 0 < k \leq (j-i)/2+1\}$$



Data Dependence Matrix:  $D = [\vartheta_A, \vartheta_B, \vartheta_C, \vartheta_D, \vartheta_E] = \begin{bmatrix} 0 & -1 & 0 & 0 & -1 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & -1 & 1 & 0 \end{bmatrix}$

Data Dependence Graph ( $m=7$ ):

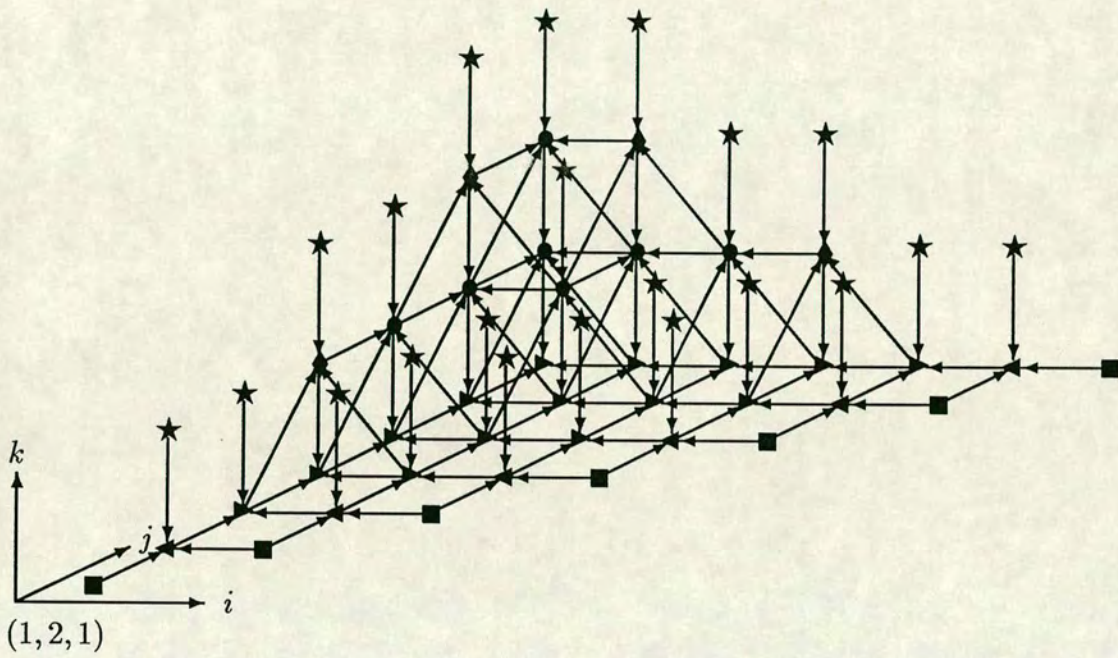


Figure 3-4: The data dependence graph for dynamic programming ( $m=7$ ).

### 3.7.1.2. The Computation Control Flow

$\Phi/\mathbb{T} = \{\Phi_{\blacksquare}, \Phi_{\blacktriangleleft}, \Phi_{\blacktriangleright}, \Phi_{\star}, \Phi_{\blacklozenge}, \Phi_{\bullet}\}$ ; the  $\mathbb{T}$ -classes are:

$$\begin{aligned} \Phi_{\blacksquare} &= \Phi \cap \{(i, j, k) \mid j-i=k \wedge k=1\} \\ \Phi_{\blacktriangleleft} &= \Phi \cap \{(i, j, k) \mid j-i=2k \wedge k=1\} \\ \Phi_{\blacktriangleright} &= \Phi \cap \{(i, j, k) \mid j-i > 2k \wedge k=1\} \\ \Phi_{\star} &= \Phi \cap \{(i, j, k) \mid j-i < 2k \wedge j-i \geq 2k-2\} \\ \Phi_{\blacklozenge} &= \Phi \cap \{(i, j, k) \mid j-i=2k \wedge k \neq 1\} \\ \Phi_{\bullet} &= \Phi \cap \{(i, j, k) \mid j-i > 2k \wedge k \neq 1\} \end{aligned}$$

All six  $\mathbb{T}$ -classes are convex polytopes. Thus,  $\Phi/\mathbb{T}$  is a  $\mathbb{T}$ -partition.

The set  $\mathfrak{S}_{\text{src}} = \{S_1, S_2, S_3\}$  given by



$$\begin{aligned}
S_1 &= \{(i, j, k) \mid k \leq 1\} \\
S_2 &= \{(i, j, k) \mid -i + j - 2k \leq -1\} \\
S_3 &= \{(i, j, k) \mid -i + j - 2k \leq 0\}
\end{aligned}$$

is a separation set of  $\Phi/\mathbb{T}$ , since

$$\begin{aligned}
\mathcal{H}(\Phi_{\blacksquare}, \mathfrak{S}_{\text{src}}) &= \{S_1, S_2\} \\
\mathcal{H}(\Phi_{\blacktriangleleft}, \mathfrak{S}_{\text{src}}) &= \{S_1, \mathbb{C}_{\mathbb{Z}^n} S_2, S_3\} \\
\mathcal{H}(\Phi_{\blacktriangleright}, \mathfrak{S}_{\text{src}}) &= \{S_1, \mathbb{C}_{\mathbb{Z}^n} S_3\} \\
\mathcal{H}(\Phi_{\star}, \mathfrak{S}_{\text{src}}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_1, S_2\} \\
\mathcal{H}(\Phi_{\blacklozenge}, \mathfrak{S}_{\text{src}}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_1, \mathbb{C}_{\mathbb{Z}^n} S_2, S_3\} \\
\mathcal{H}(\Phi_{\bullet}, \mathfrak{S}_{\text{src}}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_1, \mathbb{C}_{\mathbb{Z}^n} S_3\}
\end{aligned}$$

$S_2^=$  and  $S_3^=$  are parallel. The two control variables associated with them can be merged (Sect. 3.4). Thus, the CCUREs must specify two computation control variables:  $P$  and  $Q$ .  $P$  is associated with  $S_1$ .  $\text{sig}(P) = \{p_1, p_2\}$ .  $Q$  is associated with  $S_2$  and  $S_3$ .  $\text{sig}(Q) = \{q_1, q_2, q_3\}$ . An application of (3.8) and (3.23) yields the following CCUREs:

$$\begin{aligned}
(0, 0, 1)\vartheta_P &= 0 \\
P(I) &= \begin{cases} I \in \text{in}(\Psi, \vartheta_P) \cap S_1 & \rightarrow p_1 \\ I \in \text{in}(\Psi, \vartheta_P) \cap \mathbb{C}_{\mathbb{Z}^n} S_1 & \rightarrow p_2 \\ I \in \Psi & \rightarrow P(I) = P(I - \vartheta_P) \end{cases} \\
(-1, 1, -2)\vartheta_Q &= 0 \\
Q(I) &= \begin{cases} I \in \text{in}(\Psi, \vartheta_Q) \cap S_2 & \rightarrow q_1 \\ I \in \text{in}(\Psi, \vartheta_Q) \cap \mathbb{C}_{\mathbb{Z}^n} S_2 \cap S_3 & \rightarrow q_2 \\ I \in \text{in}(\Psi, \vartheta_Q) \cap \mathbb{C}_{\mathbb{Z}^n} S_3 & \rightarrow q_3 \\ I \in \Psi & \rightarrow Q(I) = Q(I - \vartheta_Q) \end{cases}
\end{aligned}$$

The previous six  $\mathbb{T}$ -classes are redefined in the computation control variables by (3.10) and (3.11):



$$\begin{aligned}
\Phi'_{\blacksquare} &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi_{\blacksquare})\} = \{I \mid I \in \Psi \wedge P(I)=p_1 \wedge Q(I)=q_1\} \\
\Phi'_{\blacktriangleleft} &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi_{\blacktriangleleft})\} = \{I \mid I \in \Psi \wedge P(I)=p_1 \wedge Q(I)=q_2\} \\
\Phi'_{\blacktriangleright} &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi_{\blacktriangleright})\} = \{I \mid I \in \Psi \wedge P(I)=p_1 \wedge Q(I)=q_3\} \\
\Phi'_{\star} &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi_{\star})\} = \{I \mid I \in \Psi \wedge P(I)=p_2 \wedge Q(I)=q_1\} \\
\Phi'_{\blacklozenge} &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi_{\blacklozenge})\} = \{I \mid I \in \Psi \wedge P(I)=p_2 \wedge Q(I)=q_2\} \\
\Phi'_{\bullet} &= \{I \mid I \in \Psi \wedge \mathcal{P}(\Phi_{\bullet})\} = \{I \mid I \in \Psi \wedge P(I)=p_2 \wedge Q(I)=q_3\}
\end{aligned}$$

By Lemma 3.1,  $D/\mathbb{T} \subseteq \Phi/\mathbb{T}$  for every  $D \in \text{deq}(\Phi)$ .  $D/\mathbb{T}$  for the domain  $D$  of an equation in the source UREs in Sect. 3.7.1.1 contains the  $\mathbb{T}$ -classes whose subscripts are those labeling the equation. The calculation of predicate  $\mathcal{P}(D)$  follows from (3.9). Replacing the domain predicates in the source UREs by predicates in the computation control variables as in (3.16) with optimisations yields:

$$\begin{aligned}
A(i, j, k) &= \begin{cases} Q(i, j, k)=q_2 \rightarrow D(i, j, k) & \blacktriangleleft \blacklozenge \\ Q(i, j, k) \neq q_2 \rightarrow A(i, j-1, k) & \blacktriangleright \bullet \blacksquare \star \end{cases} \\
B(i, j, k) &= \begin{cases} P(i, j, k)=p_1 \rightarrow C(i+1, j, k) & \blacktriangleleft \blacktriangleright \blacksquare \\ P(i, j, k)=p_2 \rightarrow B(i+1, j, k-1) & \blacklozenge \bullet \star \end{cases} \\
D(i, j, k) &= \begin{cases} P(i, j, k)=p_1 \rightarrow C(i, j-1, k) & \blacktriangleleft \blacktriangleright \blacksquare \\ P(i, j, k)=p_2 \rightarrow D(i, j-1, k-1) & \blacklozenge \bullet \star \end{cases} \\
E(i, j, k) &= \begin{cases} Q(i, j, k)=q_2 \rightarrow B(i, j, k) & \blacktriangleleft \blacklozenge \\ Q(i, j, k) \neq q_2 \rightarrow E(i+1, j, k) & \blacktriangleright \bullet \blacksquare \star \end{cases}
\end{aligned}$$

$$c_{1, m+1} = C(1, m+1, 1)$$

$$\begin{aligned}
C(i, j, k) &= \\
&\left\{ \begin{array}{l} P(i, j, k)=p_1 \wedge Q(i, j, k)=q_1 \rightarrow w_{i,j} \blacksquare \\ P(i, j, k)=p_1 \wedge Q(i, j, k) \neq q_1 \rightarrow \min \begin{pmatrix} A(i, j, k) + B(i, j, k) \\ C(i, j, k+1) \\ D(i, j, k) + E(i, j, k) \end{pmatrix} + w_{i,j} \blacktriangleleft \blacktriangleright \\ P(i, j, k)=p_2 \wedge Q(i, j, k) \neq q_1 \rightarrow \min \begin{pmatrix} A(i, j, k) + B(i, j, k) \\ C(i, j, k+1) \\ D(i, j, k) + E(i, j, k) \end{pmatrix} \blacklozenge \bullet \\ P(i, j, k)=p_2 \wedge Q(i, j, k)=q_1 \rightarrow \infty \star \end{array} \right.
\end{aligned}$$



### 3.7.1.3. The Propagation Control Flow

There are four supporting half-spaces to the index space:

$$\begin{aligned} B_1 &= \{(i, j, k) \mid k \geq 1\} \\ B_2 &= \{(i, j, k) \mid i \geq 1\} \\ B_3 &= \{(i, j, k) \mid j \leq m+1\} \\ B_4 &= \{(i, j, k) \mid -i+j-2k \geq -2\} \end{aligned}$$

No data variables are defined at the points of  $\Psi \cap \mathbb{C}_{\mathbb{Z}^n} B_2$  and  $\Psi \cap \mathbb{C}_{\mathbb{Z}^n} B_3$ . Hence  $\mathfrak{B} = \{B_1, B_4\}$ .

### 3.7.1.4. The Space-Time Mapping

The following space-time matrices describe the systolic arrays depicted in Fig. 3-5;  $\Pi_x$  ( $x \in \{a, b, c, d\}$ ) describes the array shown in Fig. 3-5(x).

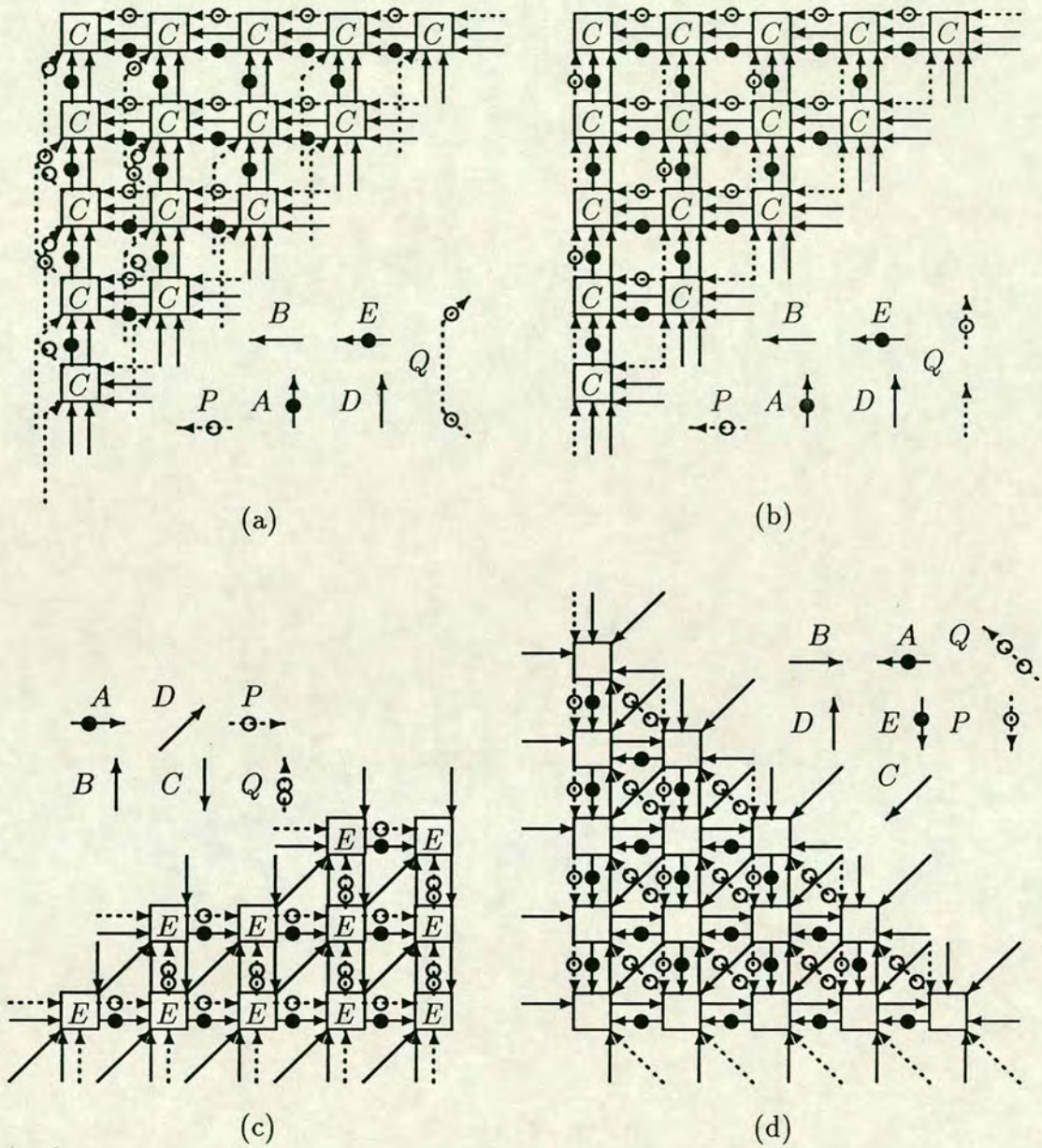
$$\begin{aligned} \Pi_a &= \begin{bmatrix} -2 & 2 & -1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} & \Pi_b &= \begin{bmatrix} -2 & 2 & -1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \\ \Pi_c &= \begin{bmatrix} -2 & 2 & -1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \Pi_d &= \begin{bmatrix} -2 & 2 & -1 \\ 0 & -1 & 1 \\ 1 & 0 & 1 \end{bmatrix} \end{aligned}$$

The four arrays each run in  $2m-1$  steps. Array (c) has  $\lfloor m/2 \rfloor (\lfloor m/2 \rfloor + 1) + m$  cells if  $m$  is odd and  $m^2/4 + m - 1$  cells if  $m$  is even. The other three arrays each have  $m(m+1)/2$  cells.

The space-time mapping of array (a) is chosen with respect to the source UREs only; it restricts the choice of the control dependence vectors. The best choice for control variable  $Q$  results in non-neighbouring channels, each connecting every other processors along the direction of these channels. A control signal of  $Q$  passes two processors in three time steps.

Array (b) can be regarded as a modification of array (a) as follows. First, eliminate all channels that connect processors  $(x, y)$ , where  $x+y$  is even. Second,





**Figure 3-5:** Four systolic arrays for dynamic programming ( $m = 5$ ). (a) Guibas-Kung-Thompson's array [26] is described by  $\lambda = (-2, 2, -1)$  and  $u = (0, 0, 1)$  with  $\vartheta_P = (-1, 0, 0)$  and  $\vartheta_Q = (0, 2, 1)$ . (b) Chen's array [14] is a modification of Guibas-Kung-Thompson's array (see text). (c) Gachet-Joinnault-Quinton's array [24] is described by  $\lambda = (-2, 2, -1)$  and  $u = (1, 0, 0)$  with  $\vartheta_P = (0, 1, 0)$  and  $\vartheta_Q = (-2, 0, 1)$ . (d) A new systolic array that we propose is described by  $\lambda = (-2, 2, -1)$  and  $u = (-1, 1, 1)$  with  $\vartheta_P = (-1, 0, 0)$  and  $\vartheta_Q = (0, 2, 1)$ .



remove all buffers from the remaining channels, make the remaining channels connect neighbours and then place one buffer on each output channel of the processors  $(x, y)$ , where  $x+y$  is even. This adaptation is a significant improvement, but it is hard to generalise. Note that the control flow of  $Q$  is not constant. Optimisations of this kind are problem-specific.

The motivations for array (c) is to reduce the number of processors. Again, the space-time mapping is chosen with respect to the source UREs only. The control channels of control variables  $P$  and  $Q$  are neighbouring channels.

Array (d) is synthesised from the pipelined UREs. We have two reasons for proposing it. First, we want to enforce neighbouring communications for both data and control variables. Second, we do not allow the existence of stationary variables. The concept of pipelined UREs enables all these design constraints to be taken into account straightforwardly. We derived this array by hand, but it could be derived using the methods in [23,51,83]. These methods allow the synthesis of systolic arrays with prescribed channel interconnections.

Let us review the control complexity of the four arrays. The control complexity of array (b) is similar to that of array (a). In arrays (a) and (d), different weights  $w_{i,j}$  reside in different cells. Accessing these weights is simple. In array (c), weights  $w_{i,j}$  with fixed  $j$  reside in the same cell. Additional access control is needed to choose the correct  $i$ . Array (a) needs to recover  $C$ , which is stationary. The stationary  $E$  in array (c) does not cause any problem because  $E$  needs not be recovered.

### 3.7.2 LU-Decomposition

LU-decomposition is the unique decomposition of a non-singular  $m \times m$  matrix  $C$  into a lower-triangular matrix  $A$  and an upper-triangular matrix  $B$  such that  $AB=C$ . The elements of the upper triangle of  $A$  and the elements of the lower triangle of  $B$  (excluding the diagonal) are 0; the diagonal elements of  $A$  are 1.

*Specification:*  $(\forall i, j : 0 < i \leq m \wedge 0 \leq j \leq m : \sum k : 0 < k \leq m : a_{i,k} b_{k,j} = c_{i,j})$



### 3.7.2.1. The Source UREs

UREs:

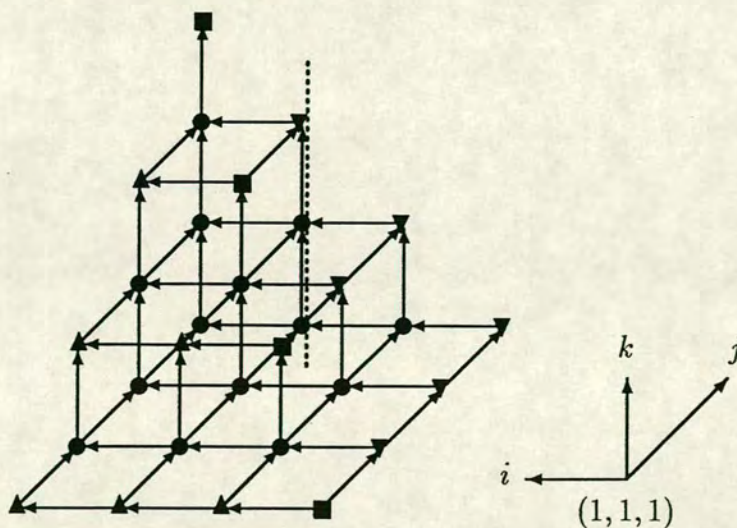
$$\begin{aligned}
 A(i, j, k) &= \begin{cases} k < i \leq m \wedge k = j \wedge 0 < k \leq m & \rightarrow C(i, j, k-1)B(i, j, k) \quad \blacktriangle \\ k < i \leq m \wedge k < j \leq m \wedge 0 < k \leq m & \rightarrow A(i, j-1, k) \quad \bullet \end{cases} \\
 B(i, j, k) &= \begin{cases} k = i \wedge k < j \leq m \wedge 0 < k \leq m & \rightarrow C(i, j, k-1) \quad \blacktriangledown \\ k = i \wedge k = j \wedge 0 < k \leq m & \rightarrow C(i, j, k-1)^{-1} \quad \blacksquare \\ k < i \leq m \wedge k \leq j \leq m \wedge 0 < k \leq m & \rightarrow B(i-1, j, k) \quad \bullet \blacktriangle \end{cases} \\
 C(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 < j \leq m \wedge 0 = k & \rightarrow c_{i,j} \\ k \leq i \leq m \wedge k \leq j \leq m \wedge 0 < k \leq m & \rightarrow C(i, j, k-1) - A(i, j-1, k) \\ & B(i-1, j, k) \quad \blacksquare \bullet \blacktriangle \blacktriangledown \end{cases} \\
 a_{i,k} &= \begin{cases} k < i \leq m \wedge m = j \wedge 0 < k \leq m & \rightarrow A(i, j, k) \end{cases} \\
 b_{k,j} &= \begin{cases} m = i \wedge k \leq j \leq m \wedge 0 < k \leq m & \rightarrow B(i, j, k) \end{cases}
 \end{aligned}$$

**Remark** Element  $b_{i,i}$  for every  $i$  in the source UREs is the reciprocal of its corresponding element in the previous specification.  $\square$

*Index Space:*  $\Phi = \{(i, j, k) \mid k \leq i \leq m \wedge k \leq j \leq m \wedge 0 < k \leq m\}$

*Data Dependence Matrix:*  $\mathcal{D} = [\vartheta_A, \vartheta_B, \vartheta_C] = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

*Data Dependence Graph ( $m=4$ ):*



**Figure 3-6:** The data dependence graph for LU-decomposition ( $m=4$ ). The dotted line will be referred to in Chap. 5.



### 3.7.2.2. The Computation Control Flow

$\Phi/\mathbb{T} = \{\Phi_{\blacksquare}, \Phi_{\blacktriangle}, \Phi_{\blacktriangledown}, \Phi_{\bullet}\}$ , the  $\mathbb{T}$ -classes are:

$$\begin{aligned}\Phi_{\blacksquare} &= \{(i, j, k) \mid k=i \wedge k=j \wedge 0 < k \leq m\} \\ \Phi_{\blacktriangle} &= \{(i, j, k) \mid k < i \leq m \wedge k=j \wedge 0 < k \leq m\} \\ \Phi_{\blacktriangledown} &= \{(i, j, k) \mid k=i \wedge k < j \leq m \wedge 0 < k \leq m\} \\ \Phi_{\bullet} &= \{(i, j, k) \mid k < i \leq m \wedge k < j \leq m \wedge 0 < k \leq m\}\end{aligned}$$

All four  $\mathbb{T}$ -classes are convex polytopes. Thus,  $\Phi/\mathbb{T}$  is a  $\mathbb{P}$ -partition.

The set  $\mathfrak{S}_{\text{src}} = \{S_1, S_2\}$  given by

$$\begin{aligned}S_1 &= \{(i, j, k) \mid i - k \leq 0\} \\ S_2 &= \{(i, j, k) \mid j - k \leq 0\}\end{aligned}$$

is a separation set, since

$$\begin{aligned}\mathcal{H}(\Phi_{\blacksquare}, \mathfrak{S}_{\text{src}}) &= \{S_1, S_2\} \\ \mathcal{H}(\Phi_{\blacktriangle}, \mathfrak{S}_{\text{src}}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_1, S_2\} \\ \mathcal{H}(\Phi_{\blacktriangledown}, \mathfrak{S}_{\text{src}}) &= \{S_1, \mathbb{C}_{\mathbb{Z}^n} S_2\} \\ \mathcal{H}(\Phi_{\bullet}, \mathfrak{S}_{\text{src}}) &= \{\mathbb{C}_{\mathbb{Z}^n} S_1, \mathbb{C}_{\mathbb{Z}^n} S_2\}\end{aligned}$$

The CCUREs specify two computation control variables:  $P$  and  $Q$ .  $P$  is associated with  $S_1$ .  $\text{sig}(P) = \{p_1, p_2\}$ .  $Q$  is associated with  $S_2$ .  $\text{sig}(Q) = \{q_1, q_2\}$ . An application of (3.8) yields the following CCUREs:

$$\begin{aligned}(1, 0, -1)\vartheta_P &= 0 \\ P(I) &= \begin{cases} I \in \text{in}(\Psi, \vartheta_P) \cap S_1 & \rightarrow p_1 \\ I \in \text{in}(\Psi, \vartheta_P) \cap \mathbb{C}_{\mathbb{Z}^n} S_1 & \rightarrow p_2 \\ I \in \Psi & \rightarrow P(I) = P(I - \vartheta_P) \end{cases} \\ (0, 1, -1)\vartheta_Q &= 0 \\ Q(I) &= \begin{cases} I \in \text{in}(\Psi, \vartheta_Q) \cap S_2 & \rightarrow q_1 \\ I \in \text{in}(\Psi, \vartheta_Q) \cap \mathbb{C}_{\mathbb{Z}^n} S_2 & \rightarrow q_2 \\ I \in \Psi & \rightarrow Q(I) = Q(I - \vartheta_Q) \end{cases}\end{aligned}$$



The previous four  $\oplus$ -classes are redefined in the computation control variables by (3.10) and (3.11):

$$\begin{aligned}\Phi'_{\blacksquare} &= \{I \in \Psi \wedge \mathcal{P}(\Phi_{\blacksquare})\} = \{I \mid I \in \Psi \wedge P(I) = p_1 \wedge Q(I) = q_1\} \\ \Phi'_{\blacktriangle} &= \{I \in \Psi \wedge \mathcal{P}(\Phi_{\blacktriangle})\} = \{I \mid I \in \Psi \wedge P(I) = p_2 \wedge Q(I) = q_1\} \\ \Phi'_{\blacktriangledown} &= \{I \in \Psi \wedge \mathcal{P}(\Phi_{\blacktriangledown})\} = \{I \mid I \in \Psi \wedge P(I) = p_1 \wedge Q(I) = q_2\} \\ \Phi'_{\bullet} &= \{I \in \Psi \wedge \mathcal{P}(\Phi_{\bullet})\} = \{I \mid I \in \Psi \wedge P(I) = p_2 \wedge Q(I) = q_2\}\end{aligned}$$

The calculation of predicates  $\mathcal{P}(D)$  for every  $D \in \text{deq}(\Phi)$  proceeds exactly as in Sect. 3.7.1.2. Replacing the domain predicates of the source UREs by predicates in the computation control variables as in (3.16) yields:

$$\begin{aligned}A(i, j, k) &= \begin{cases} P(i, j, k) = p_2 \wedge Q(i, j, k) = q_1 \rightarrow C(i, j, k-1)B(i, j, k) & \blacktriangle \\ P(i, j, k) = p_2 \wedge Q(i, j, k) = q_2 \rightarrow A(i, j-1, k) & \bullet \end{cases} \\ B(i, j, k) &= \begin{cases} P(i, j, k) = p_1 \wedge Q(i, j, k) = q_2 \rightarrow C(i, j, k-1) & \blacktriangledown \\ P(i, j, k) = p_1 \wedge Q(i, j, k) = q_1 \rightarrow C(i, j, k-1)^{-1} & \blacksquare \\ P(i, j, k) = p_2 \rightarrow B(i-1, j, k) & \bullet \blacktriangle \end{cases} \\ C(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 < j \leq m \wedge 0 = k \rightarrow c_{i,j} \\ \text{true} \rightarrow C(i, j, k-1) - A(i, j-1, k)B(i-1, j, k) & \blacksquare \bullet \blacktriangle \blacktriangledown \end{cases} \\ a_{i,k} &= \begin{cases} k < i \leq m \wedge m = j \wedge 0 < k \leq m \rightarrow A(i, j, k) \end{cases} \\ b_{k,j} &= \begin{cases} m = i \wedge k \leq j \leq m \wedge 0 < k \leq m \rightarrow B(i, j, k) \end{cases}\end{aligned}$$

### 3.7.2.3. The Propagation Control Flow

There are five supporting half-spaces to the index space:

$$\begin{aligned}B_1 &= \{(i, j, k) \mid i \leq m\} \\ B_2 &= \{(i, j, k) \mid j \leq m\} \\ B_3 &= \{(i, j, k) \mid k \geq 1\} \\ B_4 &= \{(i, j, k) \mid i - k \geq 0\} \\ B_5 &= \{(i, j, k) \mid j - k \geq 0\}\end{aligned}$$

The input of  $A$  and  $B$  is undefined and the output of  $C$  is not of interest. That is, all data variable are undefined at the points in  $\Psi \cap \mathcal{C}_{\mathbb{Z}^n} B_4$  and  $\Psi \cap \mathcal{C}_{\mathbb{Z}^n} B_5$ . Hence,  $\mathfrak{B} = \{B_1, B_2, B_3\}$ .



### 3.7.2.4. The Space-Time Mapping

Choose the following space-time mapping (Fig. 3-7(a)):

$$\Pi_0 = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

The projection vector is  $u = (1, 1, 1)$ . By Thm. 3.6, all computation control variables are redundant. Three propagation control variables,  $P_1$ ,  $P_2$  and  $P_3$  (associated with  $B_1$ ,  $B_2$  and  $B_3$ , respectively), are necessary for propagating the input data of  $C$  to the array and the output data of  $A$  and  $B$  from the array. Let us choose the data dependence vectors as control dependence vectors:  $\vartheta_{P_1} = (0, 1, 0)$ ,  $\vartheta_{P_2} = (1, 0, 0)$  and  $\vartheta_{P_3} = (0, 1, 0)$ .

Choose the following space-time mapping (Fig. 3-7(b)):

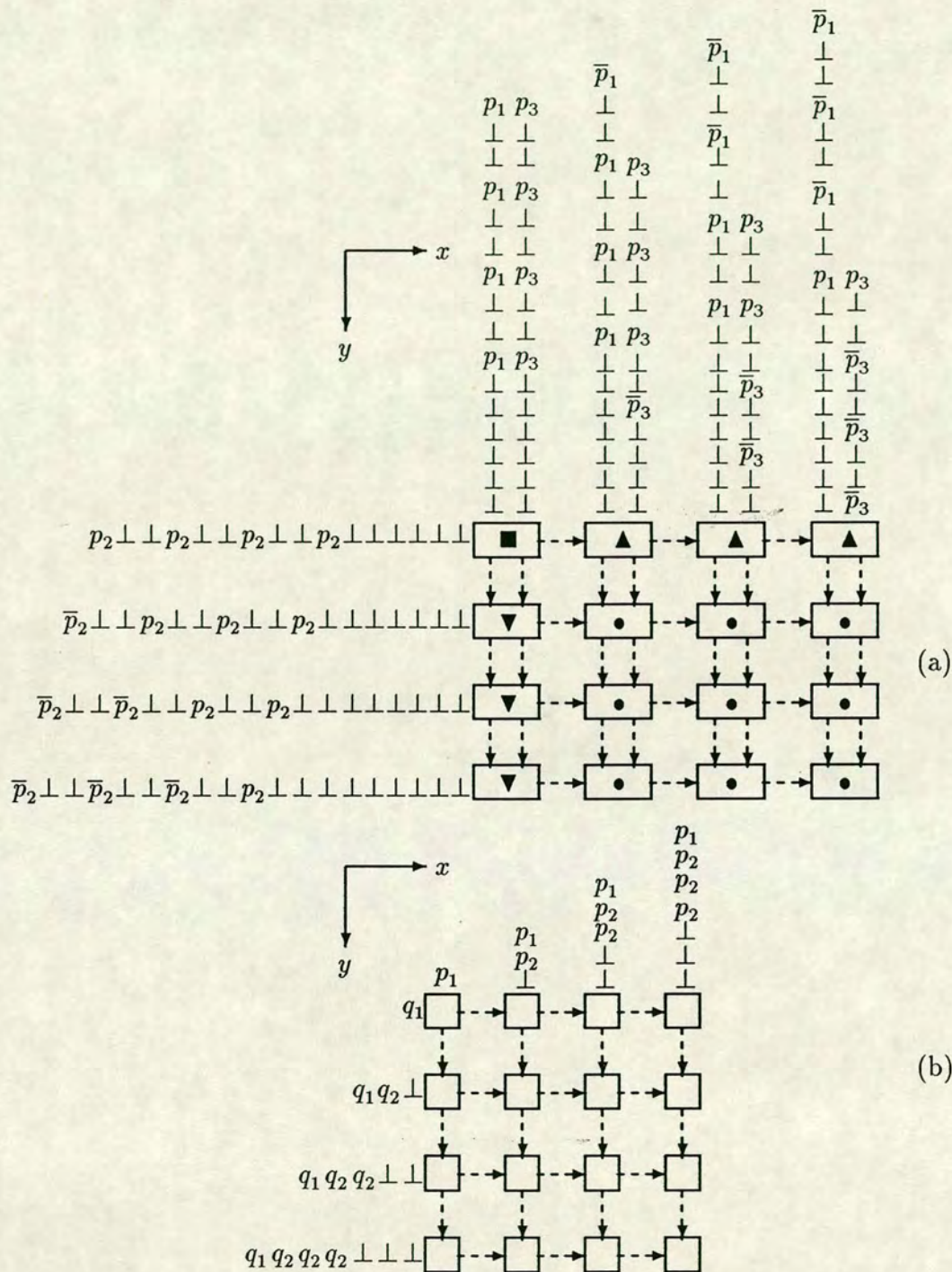
$$\Pi_1 = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

The projection vector is  $u = (0, 0, 1)$ . By Thm. 3.6, the propagation control variables  $P_1$  and  $P_2$  are redundant.  $P_3$  serves to propagate the input data of variable  $C$  to the array; it is unnecessary since  $C$  is stationary variable, i.e.,  $\Psi \cap \mathbb{C}_{\mathbb{Z}^n} B_3 = \emptyset$ . Following the previous line of reasoning, we set  $\vartheta_P = (0, 1, 0)$  and  $\vartheta_Q = (1, 0, 0)$ .

## 3.8 Conclusion

We have presented a method that enables a specification of control signals for systolic arrays in terms of control UREs. We have given necessary and sufficient conditions for the correctness of the pipelined UREs. These conditions allow the user to rewrite the domain predicates in order to obtain a better control flow. We have also shown how the pipelined UREs can be optimised to eliminate redundant control hardware in the systolic array.





**Figure 3-7:** Two systolic arrays for LU-decomposition with control distribution at the first step ( $m=4$ ).  $\perp$  represents an arbitrary choice from the two values in the respective input sequence. Computations at cells receiving  $\perp$  are undefined and therefore can be interpreted as either pipelining or computation points. (a) The space-time mapping is  $\Pi_0$ . The highlighting of a processor by a symbol in  $\{\blacksquare, \blacktriangle, \blacktriangledown, \bullet\}$  indicates that only computation points depicted by that symbol in Fig. 3-6 are mapped to the processor. (b) The space-time mapping is  $\Pi_1$ .



Our method has the following advantages. First, by specifying the control signals in terms of control UREs, we can choose space-time mappings following prescribed design criteria with respect to both data and control flow. For example, we can easily enforce neighbouring communication for both data and control signals, if possible, because both data and control dependence vectors are explicit in the pipelined UREs [83]. Second, the correctness of the pipelined UREs is independent of the space-time mapping. The pipelined UREs can be viewed as a refinement of the source UREs with the domain predicates substituted by predicates in variables of UREs. Therefore, the pipelined UREs can be manipulated like the source UREs by available synthesis methods [56]. For example, they can be easily mapped to fixed-size arrays [11,17]. Of course, such a partitioning calls for extra control signals due to the change in the specification of the pipelined UREs. These extra control signals can be specified to preserve the behaviour of the control UREs.

Our method for the construction of control UREs applies for any recurrence equations that are defined over a domain in  $\mathbf{Z}^n$ . This class of recurrence equations is an extension of AREs (Def 2.1); each equation is of the following form:

$$I \in D \rightarrow V(I) = f(W(\rho_W(I)), \dots) \quad \text{where } \rho_W(I) : \mathbf{Z}^n \rightarrow \mathbf{Z}^\ell \quad (3.28)$$

$\rho_W$  is any function from  $\mathbf{Z}^n$  to  $\mathbf{Z}^\ell$ . The domain predicate  $I \in D$  is defined as in Tab. A-1, if we apply Proc. 3.1 to derive control signals. Otherwise it can be any predicate in the index vector  $I \in \mathbf{Z}^n$ , provided that one can construct a separation set. Recall that only UREs map directly to systolic arrays. If we are concerned with the design of systolic arrays, we must first transform the specifying recurrence equations to UREs and then add control UREs (say by applying Proc. 3.1) in order to obtain pipelined UREs.

Equivalence relation  $\textcircled{\text{T}}$  plays an important rôle in the construction of the control UREs. This relation can be replaced by the equivalence relation  $\textcircled{\text{T}}_\perp$ :

$$I \textcircled{\text{T}}_\perp J \iff (\forall D : D \in \text{deq}(\Phi) : I \in D \implies J \in D \vee V(J) = \perp)$$

That is, two points are of the same type if, for every data variable of the source UREs, either the defining equations at both points agree or the defining equation



at one of the two points is undefined.  $\mathbb{T}$  is finer than  $\mathbb{T}_\perp$ . The use of  $\mathbb{T}_\perp$  is feasible because the defining equation of a variable at a computation point  $\bar{I} = (t, p)$ , if undefined, is of free interpretation. Cell  $p$  can be instructed at step  $t$  to perform any computation. Thus, this relaxation increases the space of separation sets at no extra cost.

If the index space itself is the only  $\mathbb{T}$ -class ( $\Phi/\mathbb{T} = \{\Phi\}$ ), the CCUREs are unnecessary. In this case, the PCUREs may be eliminated by exploiting some algebraic properties of operators in the source UREs. This is the topic of Chap. 6.

The notion of an extended index space serves us well in the construction of the control UREs by (3.8) and (3.18). But our method does not depend on the space-time mapping. This follows from the fact that the construction of  $\mathfrak{S}$  in Sect. 3.3.1 and  $\mathfrak{B}$  in Sect. 3.3.2 are independent of the space-time mapping. The knowledge of the projection vector does allow us to remove redundant elements in  $\mathfrak{S}$  and  $\mathfrak{B}$  (Thm. 3.6), though. Both the CCUREs and the PCUREs are defined conceptually for the extended index space. An explicit extension of the index space is unnecessary. To see which of two control values of a control variable,  $C$ , needs to be input at border cell  $p = (p_1, p_2, \dots, p_{n-1})$  at step  $t$ , we only need to perform the test  $(t, p_1, p_2, \dots, p_{n-1}) \in \overline{H.C}$  ( $\overline{H.C}$  denotes the image of  $H.C$  under the space-time mapping). We input control value **1** if the test succeeds and **0** if it fails.

There have been a number of publications on the optimisation of the latency and processor count of a systolic array [69,80,81]. In these papers, latency is defined with respect to data only. An accurate definition should also include control signals. The concept of pipelined UREs enables us to synthesise systolic arrays that satisfy prescribed design criteria with respect to both data and control signals.



## Chapter 4

# Data Flow Synthesis for One-Dimensional Systolic Arrays

### 4.1 Introductory Remarks

The space-time mapping technique described in Chap. 2 is for the synthesis of time-minimal  $(n-1)$ -dimensional arrays from  $n$ -dimensional UREs. The  $(n-1)$ -dimensional processor space is obtained by a projection of the data dependence graph along the projection vector. If the resulting systolic array requires more dimensions than are available, further projections of the processor space become necessary. A sequence of projections is called a *multi-projection* [38,64]. The dimensions that are eliminated by the multiprojection are traded to time.

In practice, one-dimensional arrays have the following advantages [36]: 100% utilisation of non-faulty cells on a wafer, a constant I/O bandwidth that can be achieved by restricting external communication to the two boundary cells, and a clock rate that is independent of the size of the array. This chapter is concerned with the synthesis of data flow for one-dimensional arrays. There are two aspects to this study. The first one is an extension and improvement of previous results, and the second is a characterisation of some properties of one-dimensional arrays that will be used in the synthesis of control flow for one-dimensional arrays presented in Chap. 5.



The rest of this chapter is organised as follows. Sect. 4.2 defines the two most frequently used one-dimensional systolic array models: one allows non-neighbouring channel connections, the other does not. Sect. 4.3 formulates conditions for the validity of a space-time mapping in both models. Sect. 4.4 discusses the space-time behaviour of one-dimensional arrays. In particular, it characterises the distribution of pipelining points. Sect. 4.5 discusses the space-time behaviour of one-dimensional arrays further (in the extended index space). This leads to the realisation that the non-injectivity of the space-time mapping plays a crucial rôle in the analysis and synthesis of one-dimensional arrays. Sect. 4.6 describes a procedure for the synthesis of data flow for both models. Sect. 4.7 reviews related work. Sect. 4.8 contains the conclusion of the chapter.

## 4.2 One-Dimensional Systolic Array Models

Recall that a space-time mapping that describes an  $r$ -dimensional systolic array consists of two components: **step** and **place**.

- **step** :  $\Phi \rightarrow \mathbf{Z}$ ,  $\text{step}(I) = \lambda I$ ,  $\lambda \in \mathbf{Z}^n$ . **step** specifies the temporal distribution.
- **place** :  $\Phi \rightarrow \mathbf{Z}$ ,  $\text{place}(I) = \sigma I$ ,  $\sigma \in \mathbf{Z}^{r \times n}$ . **place** specifies the spatial distribution.

This chapter considers only space-time mappings that describe one-dimensional arrays ( $r=1$ ). In this case,  $\sigma$  is a vector, the *allocation vector*. It determines the coordinate of the leftmost cell  $p_{\min}$  and the rightmost cell  $p_{\max}$  of the array:

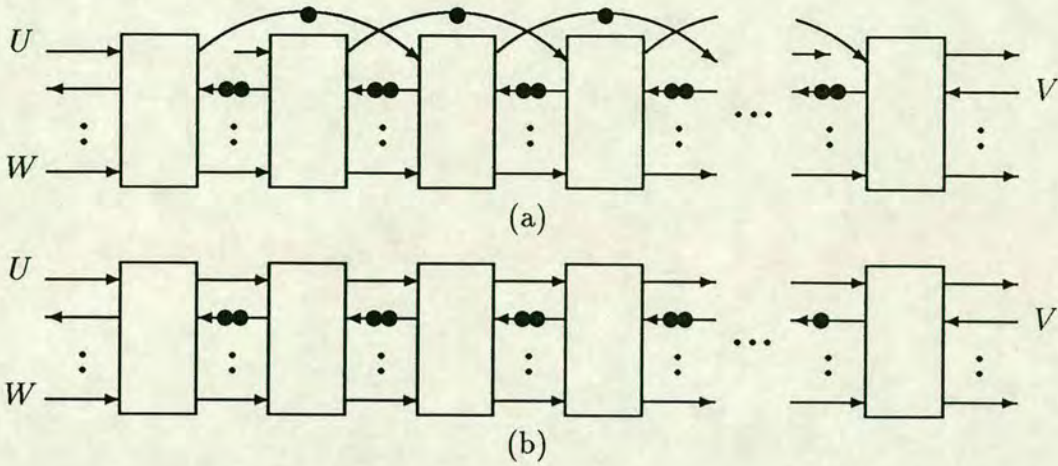
$$\begin{aligned} p_{\min} &= (\min I : I \in \Phi : \sigma I) \\ p_{\max} &= (\max I : I \in \Phi : \sigma I) \end{aligned}$$

By convention, variable  $V$  moves to the right if  $\text{flow}(V) > 0$ , to the left if  $\text{flow}(V) < 0$ .

To distinguish them from space-time mappings  $\Pi$  that describe  $(n-1)$ -dimensional arrays, we denote the space-time mappings for one-dimensional arrays by  $\pi$ :

$$\pi = \begin{bmatrix} \lambda \\ \sigma \end{bmatrix} = \begin{bmatrix} \lambda_1, \lambda_2, \dots, \lambda_n \\ \sigma_1, \sigma_2, \dots, \sigma_n \end{bmatrix}$$





**Figure 4-1:** Two one-dimensional systolic array models.  $\lambda\vartheta_U = 2$ ,  $\sigma\vartheta_U = 2$ ,  $\lambda\vartheta_V = 3$ ,  $\sigma\vartheta_V = -1$ ,  $\lambda\vartheta_W = 1$  and  $\sigma\vartheta_W = 1$ . (a) The  $\pi$ -model. (b) The  $\chi$ -model.

The space-time mapping can be interpreted geometrically. The interpretation of the scheduling vector  $\lambda$  is the same as for  $(n-1)$ -dimensional arrays. That is, the points that are scheduled concurrently belong to a hyperplane whose normal is  $\lambda$ . The processor space is obtained by projecting the index space along the  $(n-1)$ -dimensional space that is orthogonal to the allocation vector  $\sigma$ . The points that are mapped to the same cell belong to a hyperplane whose normal is  $\sigma$ .

We study space-time mappings with respect to two one-dimensional array models. They differ in channel topology. One model is called the  $\pi$ -model; in it both the direction and length of a channel depend on the space-time mapping  $\pi$ . This model is an instance of the  $r$ -dimensional systolic array model of Def. 2.2 when  $r=1$ . The other model is called the  $\chi$ -model; it allows only neighbouring communication. In the  $\chi$ -model, only the direction of a channel depends on the space-time mapping; its length is constant.

**Definition 4.1** ( $\pi$ -model) A one-dimensional systolic array consists of a finite sequence of cells with the following properties (Fig. 4-1(a)):

Prop. 1 (Synchrony of Computation) The array is driven by a global clock that ticks in unit time. Each cell is active at every clock cycle.



Prop. 2 (Uniqueness of Channel Connections) The image  $\bar{\vartheta}_V$  of data dependence vector  $\vartheta_V$  is given by:

$$\bar{\vartheta}_V = \text{if flow}(V)=0 \rightarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix} \square \text{if flow}(V) \neq 0 \rightarrow \begin{bmatrix} \lambda\vartheta_V \\ \sigma\vartheta_V \end{bmatrix} \text{fi}$$

For variable  $V$ , every cell  $p$  has an unique input (output) channel, which connects to cell  $p - \sigma\vartheta_V$  ( $p + \sigma\vartheta_V$ ). The number of buffers associated with the channel is the first component of  $\bar{\vartheta}_V$  decremented by 1. If  $p - \sigma\vartheta_V \notin \mathcal{P}$ , cell  $p$  is an *input cell* of  $V$ . If  $p + \sigma\vartheta_V \notin \mathcal{P}$ , cell  $p$  is an *output cell* of  $V$ . Both input and output cells of  $V$  are called *border cells* of  $V$ ; they are connected to the external environment. The cells that are not border cells of  $V$  are called *internal cells* of  $V$ .

Prop. 3 (Linearity of Velocity) A variable moves with a constant velocity.  $\square$

Prop. 1 is the standard assumption that the evaluation of a point takes unit time. Prop. 2 ensures that a cell receives elements of a fixed variable from a fixed channel. Depending on the length of the channel, a border cell may be a processor inside the array – a few processors away from the real border. (In Fig. 4-1(a), the two cells on the left (right) are the input (output) cells of  $U$ . All four cells are the border cells of  $U$ .) The definition of  $\bar{\vartheta}_V$  for stationary variable  $V$  implies that each element of  $V$  travels along a loop channel without delay buffers. This definition conforms to the space-time diagram that will be introduced in Sect. 4.4. (In the systolic array model of Def. 2.2, the definition of  $\bar{\vartheta}_V$  does not distinguish stationary from moving variables. There, each element of stationary variable  $V$  can be viewed as travelling along a loop channel associated with  $\lambda\vartheta_V - 1$  delay buffers.) Prop. 3 is guaranteed by the linearity of space-time mappings.

We refer to the sequence of all the channels associated with a variable in which every two adjacent channels connect to the same cell as a *link* for that variable. There are  $|\sigma\vartheta_V|$  links for variable  $V$ . Therefore, there are  $|\sigma\vartheta_V|$  input (output) cells for variable  $V$ , one per link. By the  $i$ -th input (output) cell of  $V$  ( $i$  is numbered from 1) we mean the input (output) cell that has distance  $i - 1$  from the input (output) cell of  $V$  which is the leftmost or rightmost cell.



Next, we determine where and when an input (output) value is injected (ejected). Recall that **pattern** specifies the distribution of the input data in the processor space at the first step. For one-dimensional arrays, it is more convenient to reason about the temporal and spatial distribution of input (output) data separately.

Function **pi** specifies the coordinates of the input cells:

$$\begin{aligned} \text{pi} : \mathcal{V} \rightarrow \Phi_{\text{fst}} \rightarrow \mathbf{Z}, \quad \Phi_{\text{fst}} &= (\bigcup V : V \in \mathcal{V} : \text{fst}(\Phi_V, \vartheta_V)) \\ \text{pi}(V(I)) &= \begin{cases} \text{if } \text{flow}(V) > 0 & \rightarrow p_{\text{min}} + (\text{place}(I) - p_{\text{min}}) \bmod |2 \cdot \bar{\vartheta}_V| \\ \square \text{ flow}(V) < 0 & \rightarrow p_{\text{max}} + (\text{place}(I) - p_{\text{max}}) \bmod |2 \cdot \bar{\vartheta}_V| \\ \text{fi} \end{cases} \end{aligned}$$

We write  $\textcircled{1}_V$  for the equivalence relation on  $\text{fst}(\Phi_V, \vartheta_V)$  such that two points  $I$  and  $J$  are in the same  $\textcircled{1}_V$ -class iff  $V(I)$  and  $V(J)$  are injected at the same input cell, i.e.,

$$(\forall I, J : I, J \in \text{fst}(\Phi_V, \vartheta_V) : I \textcircled{1}_V J \iff \text{pi}(V(I)) = \text{pi}(V(J)))$$

There are  $|\sigma \vartheta_V|$   $\textcircled{1}_V$ -classes, as many as there are links for variable  $V$ .

Function **po** specifies the coordinates of the output cells:

$$\begin{aligned} \text{po} : \mathcal{V} \rightarrow \Phi_{\text{lst}} \rightarrow \mathbf{Z}, \quad \Phi_{\text{lst}} &= (\bigcup V : V \in \mathcal{V} : \text{lst}(\Phi_V, \vartheta_V)) \\ \text{po}(V(I)) &= \begin{cases} \text{if } \text{flow}(V) > 0 & \rightarrow p_{\text{max}} + (\text{place}(I) - p_{\text{max}}) \bmod |2 \cdot \bar{\vartheta}_V| \\ \square \text{ flow}(V) < 0 & \rightarrow p_{\text{min}} + (\text{place}(I) - p_{\text{min}}) \bmod |2 \cdot \bar{\vartheta}_V| \\ \text{fi} \end{cases} \end{aligned}$$

We write  $\textcircled{0}_V$  for the equivalence relation on  $\text{lst}(\Phi_V, \vartheta_V)$  given by

$$(\forall I, J : I, J \in \text{lst}(\Phi_V, \vartheta_V) : I \textcircled{0}_V J \iff \text{po}(V(I)) = \text{po}(V(J)))$$

There are  $|\sigma \vartheta_V|$   $\textcircled{0}_V$ -classes, as many as there are links for variable  $V$ .

Having defined functions that determine the cells that perform input (output), we next define functions that determine the steps at which input (output) is performed.

Function **input** specifies the steps at which input data are injected into the array:

$$\begin{aligned} \text{input} : \mathcal{V} \rightarrow \Phi_{\text{fst}} \rightarrow \mathbf{Z}, \quad \Phi_{\text{fst}} &= (\bigcup V : V \in \mathcal{V} : \text{fst}(\Phi_V, \vartheta_V)) \\ \text{input}(V(I)) &= \text{step}(I) - (\text{place}(I) - \text{pi}(V(I))) / \text{flow}(V) \end{aligned}$$



Function **output** specifies the steps at which output data are ejected from the array.

$$\text{output} : \mathcal{V} \rightarrow \Phi_{\text{lst}} \rightarrow \mathbb{Z}, \quad \Phi_{\text{lst}} = (\cup V : V \in \mathcal{V} : \text{lst}(\Phi_V, \vartheta_V))$$

$$\text{output}(V(I)) = \text{step}(I) - (\text{place}(I) - \text{po}(V(I)) / \text{flow}(V))$$

The first step  $t_{\text{fst}}$  at which an input value is injected and the last step  $t_{\text{lst}}$  at which an output value is ejected are given by

$$t_{\text{fst}} = (\min V, I : V \in \mathcal{V} \wedge I \in \Phi_{\text{fst}} : \text{input}(V(I)))$$

$$t_{\text{lst}} = (\max V, I : V \in \mathcal{V} \wedge I \in \Phi_{\text{lst}} : \text{output}(V(I)))$$

**Definition 4.2** ( $\chi$ -model) A one-dimensional systolic array consists of a finite sequence of cells with the following properties (Fig. 4-1(b)):

Prop. 1 (Synchrony of Computation) The array is driven by a global clock that ticks in unit time. Each cell is active at every clock cycle.

Prop. 2 (Uniqueness of Channel Connections) The image  $\bar{\vartheta}_V$  of data dependence vector  $\vartheta_V$  is given by:

$$\bar{\vartheta}_V = \text{if } \text{flow}(V) = 0 \rightarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix} \square \text{flow}(V) \neq 0 \rightarrow \begin{bmatrix} \lambda \vartheta_V / |\sigma \vartheta_V| \\ \text{sign}(\text{flow}(V)) \end{bmatrix} \text{fi}$$

For variable  $V$ , there is a unique channel between every two neighbouring cells. The number of buffers associated with the channel is the first component of  $\bar{\vartheta}_V$  decremented by 1.  $p_{\text{min}}$  is the *input (output) cell* of  $V$  and  $p_{\text{max}}$  is the *output (input) cell* of  $V$  if  $\text{flow}(V) > 0$  ( $\text{flow}(V) < 0$ ). Both input and output cells of  $V$  are called *border cells* of  $V$ ; they are connected to the external environment. The cells that are not border cells of  $V$  are called *internal cells* of  $V$ .

Prop. 3 (Linearity of Velocity) A variable moves with a constant velocity.  $\square$

Prop. 2 requires that, for variable  $V$  travelling along non-neighbouring channels in the  $\pi$ -model, neighbouring communication must be enforced in the  $\chi$ -model.



This is achieved by breaking a non-neighbouring channel into  $|\sigma\vartheta_V|$  neighbouring channels, making the new channels connect  $|\sigma\vartheta_V|-1$  intermediate cells, and evenly distributing the  $\lambda\vartheta_V-1$  buffers associated with the original channel over the new channels (see variable  $U$  in Fig. 4-1). Since each intermediate cell acts as a delay buffer, the number of buffers associated with a new channel is  $(\lambda\vartheta_V-1-(|\sigma\vartheta_V|-1))/|\sigma\vartheta_V|$ , which simplifies to  $\lambda\vartheta_V/|\sigma\vartheta_V|-1$ . The definition of  $\bar{\vartheta}_V$  for stationary variable  $V$  is consistent with the previous interpretation for moving variables. Each element of  $V$  travels along a loop channel without delay buffers.

**Remark** The definitions of  $\text{pi}$ ,  $\text{po}$ ,  $\text{input}$ ,  $\text{output}$ ,  $\textcircled{I}_V$  and  $\textcircled{O}_V$  for the  $\pi$ -model carry over to the  $\chi$ -model. Since the  $\chi$ -model enforces neighbouring communication,  $\text{fst}(\Phi_V, \vartheta_V)/\textcircled{I}_V = \{\text{fst}(\Phi_V, \vartheta_V)\}$  and  $\text{lst}(\Phi_V, \vartheta_V)/\textcircled{O}_V = \{\text{lst}(\Phi_V, \vartheta_V)\}$ .

The functions  $\text{pi}$ ,  $\text{po}$ ,  $\text{input}$  and  $\text{output}$  are undefined for stationary variables. Whenever we apply these functions, it is understood that the variable they are applied to is moving. For example, when we write  $(\forall V : V \in \mathcal{V} : \sigma\vartheta_V | \lambda\vartheta_V)$  we mean  $(\forall V : V \in \mathcal{V} \wedge \text{flow}(V) \neq 0 : \sigma\vartheta_V | \lambda\vartheta_V)$ .  $\square$

### 4.3 The Mapping Conditions

We cannot ensure the computation and communication rules of Def. 2.4 by simply choosing non-singular space-time matrices, because these matrices are of size  $2 \times n$  and are not square when  $n > 2$ . The previous mapping conditions for  $(n-1)$ -dimensional arrays must be modified and extended.

**Theorem 4.1** *A space-time mapping  $\pi$  is valid for the source UREs in the  $\pi$ -model iff*

- $(\forall V : V \in \mathcal{V} : \lambda\vartheta_V \geq 1)$  *(Precedence Constraint)*
- $\pi : \Phi \mapsto \mathbf{Z}^2$  *(Computation Constraint)*
- $(\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V)/\textcircled{I}_V : \text{input} : S \mapsto \mathbf{Z}))$  *(Communication Constraint)*



**Proof** The precedence, computation and communication constraint are equivalent to the respective rules of Def. 2.4. The delay rule is enforced by Prop. 2 of the  $\pi$ -model.  $\square$

**Theorem 4.2** *A space-time mapping  $\pi$  is valid for the source UREs in the  $\chi$ -model iff*

- $(\forall V : V \in \mathcal{V} : \lambda \vartheta_V \geq 1)$  *(Precedence Constraint)*
- $\pi : \Phi \mapsto \mathbf{Z}^2$  *(Computation Constraint)*
- $(\forall V : V \in \mathcal{V} : \sigma \vartheta_V | \lambda \vartheta_V)$  *(Delay Constraint)*
- $(\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \mathbb{1}_V : \text{input} : S \mapsto \mathbf{Z}))$   
*(Communication Constraint)*

**Proof** By Thm. 4.1, we only need to consider the delay constraint. By Prop. 2 of the  $\chi$ -model, the delay constraint is equivalent to the delay rule of Def. 2.4.  $\square$

The communication constraint in Thm. 4.2 can be written equivalently as  $(\forall V : V \in \mathcal{V} : \text{input} : \text{fst}(\Phi_V, \vartheta_V) \mapsto \mathbf{Z})$ , since  $\text{fst}(\Phi_V, \vartheta_V) / \mathbb{1}_V = \{\text{fst}(\Phi_V, \vartheta_V)\}$ . The fact that the communication constraint can be phrased equivalently in both models allows us to talk about it without an explicit reference to either model.

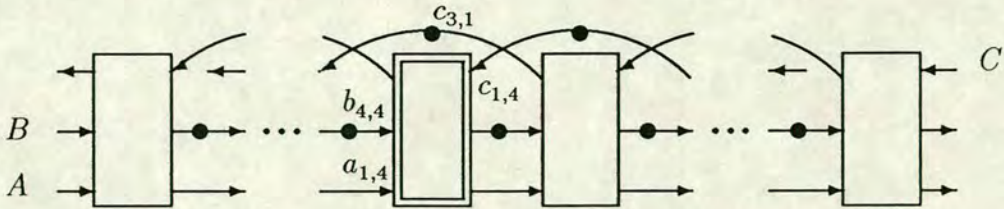
It is easy to see that the set of mapping conditions is stronger for the  $\chi$ -model than for the  $\pi$ -model. Put another way, space-time mappings that are valid in the  $\chi$ -model are also valid in the  $\pi$ -model, but the converse is not true. Let us use matrix product as example to illustrate this fact and the essential rôle that the communication constraint plays in ensuring the validity of the space-time mapping.

**Example 4.1**  $4 \times 4$  Matrix Product.

Choose the following space-time mapping:

$$\pi = \begin{bmatrix} 2 & 1 & 2 \\ 1 & 1 & -2 \end{bmatrix}$$





**Figure 4-2:** A systolic array in which the neighbouring communication cannot be enforced. The cell highlighted by a doubly bounded box is computing point  $(1, 4, 4)$ ; it is performing the last accumulation of element  $c_{1,4}$ . The delay buffer above the cell is propagating element  $c_{3,1}$ , which is the value represented by  $C(3, 1, 3)$ . If the two links associated with  $C$  were merged, the cell would have to both accumulate  $C_{1,4}$  and propagate  $c_{3,1}$  at the same step.

The precedence, delay and computation constraint are satisfied. The communication constraint is satisfied for variables  $A$  and  $B$ . Let us consider variable  $C$  with  $\lambda\vartheta_C = 2$  and  $\sigma\vartheta_C = -2$ . In the  $\chi$ -model, only one link per variable is allowed. The communication constraint is violated since  $\text{input}(C(1, 4, 0)) = \text{input}(C(3, 1, 0)) = 5$  and  $\text{input}(C(2, 4, 0)) = \text{input}(C(4, 1, 0)) = 8$  (Fig. 4-2). In the  $\pi$ -model, there are two links associated with variable  $C$ ; each link consist of channels that connect every other cells. The communication constraint is satisfied because  $\text{pi}(C(1, 4, 0)) = \text{pi}(C(4, 1, 0)) = p_{\max} - 1$  and  $\text{pi}(C(3, 1, 0)) = \text{pi}(C(2, 4, 0)) = p_{\max}$ . That is, elements  $C(1, 4, 0)$  and  $C(3, 1, 0)$  are injected at the same step but at different input cells. They travel along different links. The same is true of elements  $C(2, 4, 0)$  and  $C(4, 1, 0)$ .

Choose the following space-time mapping:

$$\pi = \begin{bmatrix} 6 & 2 & 2 \\ 1 & -2 & 1 \end{bmatrix}$$

The precedence, delay and computation constraint are satisfied. The communication constraint is satisfied for variables  $A$  and  $B$ . Let us consider variable  $C$  with  $\lambda\vartheta_C = 2$  and  $\sigma\vartheta_C = 1$ . There is only one link consisting of neighbouring channels associated with variable  $C$  in both models. Since  $\text{input}(C(1, 3, 0)) = \text{input}(C(4, 1, 0)) = 34$  and  $\text{input}(C(1, 4, 0)) = \text{input}(C(4, 2, 0)) = 40$ , the communication constraint for variable  $C$  is violated in both models.  $\square$



**Remark** There are two one-dimensional array models. Some of the following analysis applies to one model but not to the other. For example, the definitions of  $\Gamma_V$  in (4.3) and  $\delta_V$  in (4.4) are intended for the  $\chi$ -model only. It would be confusing to switch constantly between the two models in the presentation. We avoid this confusion by giving the definition of a concept for two array models even if the concept is meaningful for only one of the two models. So  $\Gamma_V$  in (4.3) and  $\delta_V$  in (4.4) are also defined for the  $\pi$ -model. All our results apply for both array models unless specified otherwise.  $\square$

The next lemma provides conditions under which the communication constraint implies the computation constraint. (The converse is generally not true as has been demonstrated by Ex. 4.1.) Intuitively, if two distinct points are computed simultaneously at the same cell, the two elements of a fixed variable indexed by the two points must travel along a common link and, consequently, must be input simultaneously at the corresponding input cell of the link.

**Lemma 4.1** *Let  $\Phi = \Phi_V$  and  $\lambda\vartheta_V \neq 0$ .*

$$(\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \mathbb{I}_V : \text{input} : S \mapsto \mathbf{Z}) \implies \pi : \Phi \mapsto \mathbf{Z}^2 \quad (4.1)$$

**Proof** Assume  $\pi I = \pi J$  for two distinct points  $I \neq J \wedge I, J \in \Phi$ .

true  
 $\iff$  {Assumption}  
 $\pi I = \pi J \wedge I \neq J$   
 $\implies$  {Definition of input}  
 $\text{input}(V(I)) = \text{input}(V(J)) \wedge I \neq J$   
 $\iff$   $\{\lambda\vartheta_V \neq 0; \lambda I = \lambda J\}$   
 $\text{input}(V(I)) = \text{input}(V(J)) \wedge I \overset{\vartheta_V}{\neq} J$   
 $\iff$  {Choose  $K, L \in \text{fst}(\Phi_V, \vartheta_V)$  such that  $K \overset{\vartheta_V}{\neq} I \wedge L \overset{\vartheta_V}{\neq} J$ }  
 $\text{input}(V(K)) = \text{input}(V(L)) \wedge K \overset{\vartheta_V}{\neq} L$   
 $\iff$  {Definition of  $\mathbb{I}_V$ }  
 $(\exists S : S \in \text{fst}(\Phi_V, \vartheta_V) / \mathbb{I}_V : K, L \in S)$   
 $\iff$  {Antecedent of (4.1)}  
 false

$\square$



## 4.4 The Space-Time Diagram

In this section, we discuss the space-time behaviour of one-dimensional arrays and describe a different formulation of the communication constraint.

The set  $\Upsilon$  of *space-time points* for a one-dimensional array is given by:

$$\Upsilon = \{(t, p) \mid t_{\text{fst}} \leq t \leq t_{\text{lst}} \wedge p_{\text{min}} \leq p \leq p_{\text{max}} \wedge t \in \mathbf{Z} \wedge p \in \mathbf{Z}\}$$

where  $t$  represents time and  $p$  represents space. The space-time points in the systolic array are divided into two categories.

- The set  $\bar{\Phi}$  of *computation points* is the image of  $\Phi$ . The defining equations of variables at these points are those at the inverse images of the points.
- The set  $\mathcal{C}_{\Upsilon}\bar{\Phi}$  of *pipelining points* is the complement of  $\bar{\Phi}$  in  $\Upsilon$ . The defining equations of variables at these points are pipelining equations.

Note that we refer to the points in both  $\Phi$  and  $\bar{\Phi}$  as computation points and the points in both  $\Psi^{\text{P}}$  and  $\mathcal{C}_{\Upsilon}\bar{\Phi}$  as pipelining points ( $\Psi^{\text{P}}$  now contains the set of points generated in the extension of the index space with respect to the allocation vector; it will be formally defined in Def. 4.3). No confusion should arise since  $\bar{\Phi}$  and  $\bar{\Psi}^{\text{P}}$  are disjoint under valid space-time mappings.

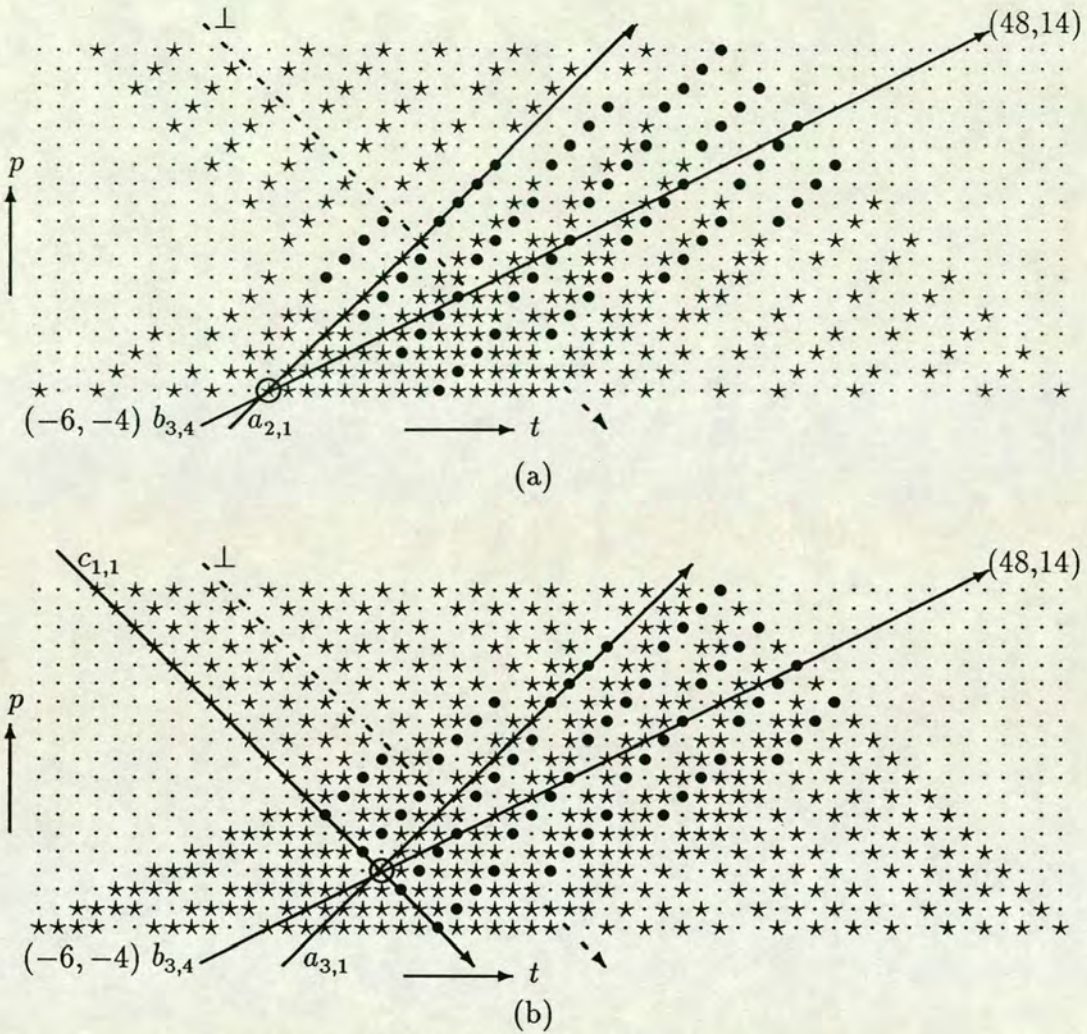
We introduce the concept of a path, which contains the points of the index space at which a given element of a variable is computed.  $p(S, \vartheta, I)$ , called a  $\vartheta$ -*path* or *path*, is the intersection of the set  $S$  and the set  $\text{ray}(I, \pm\vartheta)$ :

$$p(S, \vartheta, I) = S \cap \text{ray}(I, \pm\vartheta)$$

So  $p(S, \vartheta, I) = p(S, \vartheta, J)$  iff  $I$  and  $J$  are in the same  $\vartheta$ -path. We write  $\text{paths}(S, \vartheta)$  for the set of all  $\vartheta$ -paths associated with a set  $S$ :

$$\text{paths}(S, \vartheta) = \{p(S, \vartheta, I) \mid I \in S\}$$





**Figure 4-3:** The space-time diagrams with respect to (a) the  $\pi$ -model and (b) the  $\chi$ -model that are described by the following space-time mapping [41] ( $m=4$ ):

$$\pi_{\text{even}} = \begin{bmatrix} 2m-2 & 1 & m/2 \\ m-1 & 1 & -m/2 \end{bmatrix} \quad \pi_{\text{odd}} = \begin{bmatrix} 2m & 1 & (m+1)/2 \\ m & 1 & -(m+1)/2 \end{bmatrix}$$

$\pi_{\text{even}}$  applies for even  $m$  and  $\pi_{\text{odd}}$  for odd  $m$ . The point at the bottom-left corner is  $(t_{\text{fst}}, p_{\text{min}}) = (-6, -4)$ , that at the top-right corner is  $(t_{\text{lst}}, p_{\text{max}}) = (48, 14)$ . The fat dots represent computation points. The regular dots and stars represent pipelining points; the stars depict the pipelining points at which at least one element is propagated but exclude those that only propagate data elements of  $A$  and  $B$  which are no longer used, since the output of  $A$  and  $B$  is of no interest. The pipelining points highlighted by circles are referred to in this section and in Sect. 4.5. Arrows depict related paths. Their tails are labelled with the corresponding input data.



It is convenient to represent a one-dimensional array as a *space-time diagram*, which is obtained by viewing the space-time points as a two-dimensional lattice, where the horizontal axis represents time and the vertical axis space (Fig. 4-3). From the space-time diagram, we can directly extract the following information: the latency, the number of cells required, the velocities of variables, and input and output characteristics. To find out how the input and output data are handled in the systolic array, we note that the elements of variable  $V$  must be injected at the first points of  $\bar{\vartheta}_V$ -paths and ejected at the last points of  $\bar{\vartheta}_V$ -paths. Some  $\bar{\vartheta}_V$ -paths of  $V$  may not start at input cells in the following two cases:

- $V$  is a moving. The input data for these paths must be  $\perp$ ; they can be disregarded; i.e., they need not be supplied from the external environment.
- $V$  is a stationary. The input data for these paths must be loaded before the computation starts, i.e., before step  $t_{fst}$ .

The two sets  $\text{fst}(\Phi_V, \vartheta_V)$  and  $\text{paths}(\Phi_V, \vartheta_V)$  are isomorphic by the bijection

$$\eta_V : \text{fst}(\Phi_V, \vartheta_V) \rightarrow \text{paths}(\Phi_V, \vartheta_V), \quad \eta_V(I) = p(\Phi_V, \vartheta_V, I) \quad (4.2)$$

The communication constraint states that at most one element of a variable can be injected at any input cell at a given step. Expressed in terms of the mapping from  $\text{paths}(\Phi_V, \vartheta_V)$  to  $\text{paths}(\Upsilon, \bar{\vartheta}_V)$ , this means that different  $\vartheta_V$ -paths must have different images, i.e., different  $\bar{\vartheta}_V$ -paths in the space-time diagram. This is not true if  $V$  is stationary, because more than one input value, i.e., more than one  $\vartheta_V$ -path may be projected to the same cell.

**Lemma 4.2** *For a moving variable  $V$ , let  $\tau_V : \text{paths}(\Phi_V, \vartheta_V) \rightarrow \text{paths}(\Upsilon, \bar{\vartheta}_V)$ ,  $\tau_V(p(\Phi_V, \vartheta_V, I)) = p(\Upsilon, \bar{\vartheta}_V, \bar{I})$ . The communication constraint is satisfied iff*

$$(\forall V : V \in \mathcal{V} : \tau_V : \text{paths}(\Phi_V, \vartheta_V) \rightarrow \text{paths}(\Upsilon, \bar{\vartheta}_V))$$

**Proof**

$$\begin{aligned} & (\forall V : V \in \mathcal{V} : \tau_V : \text{paths}(\Phi_V, \vartheta_V) \rightarrow \text{paths}(\Upsilon, \bar{\vartheta}_V)) \\ \iff & \{ \eta_V \text{ of (4.2)} \} \end{aligned}$$



$$\begin{aligned}
& \tau_V \circ \eta_V : \text{fst}(\Phi_V, \vartheta_V) \mapsto \text{paths}(\Upsilon, \bar{\vartheta}_V) \\
\iff & \{\text{Definition of } \bar{\vartheta}_V\text{-path}\} \\
& (\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \textcircled{I}_V : \text{input} : S \mapsto \mathbf{Z})) \quad \square
\end{aligned}$$

Next, we characterise the distribution of pipelining points in the space-time diagram. There are two types of  $\bar{\vartheta}_V$ -paths (Fig. 4-3):

- Paths that contain computation points. These paths can be viewed as consisting of the following three consecutive segments:

- The first segment consists of pipelining points called *oaking points*.
- The second segment consists of both computation and pipelining points. These pipelining points are called *relaying points*. There are  $\Gamma_V - 1$  relaying points between every two neighbouring computation points:

$$\begin{aligned}
\Gamma_V &= \text{if } \text{flow}(V)=0 \rightarrow \lambda \vartheta_V \parallel \text{flow}(V) \neq 0 \rightarrow 1 \text{ fi } (\pi\text{-model}) \\
\Gamma_V &= \text{if } \text{flow}(V)=0 \rightarrow \lambda \vartheta_V \parallel \text{flow}(V) \neq 0 \rightarrow |\sigma \vartheta_V| \text{ fi } (\chi\text{-model})
\end{aligned} \tag{4.3}$$

The relaying points arise in the  $\chi$ -model due to the conversion of non-neighbouring channels to neighbouring channels. (E.g., see the  $\vartheta_B$ -path depicted by the solid arrow in Fig. 4-3(b); in this case,  $\Gamma_B = 3$ .)

- The third segment consists of pipelining points called *draining points*.
- Paths that do not contain computation points. The points of these paths are called *undefined points*. (E.g., see the  $\vartheta_C$ -paths depicted by the dashed arrows in Figs. 4-3(a) and (b).)

**Remark** If  $V$  is stationary variable, the definitions of  $\bar{\vartheta}_V$  in both models are consistent with the definitions of  $\Gamma_V$ . In this case, each element of  $V$  is accessed every  $\Gamma_V$  steps. For any  $\bar{\vartheta}_V$ -path in the space-time diagram, there are  $\Gamma_V - 1$  relaying points between every two neighbouring computation points.  $\square$

The soaking points of  $V$  serve to propagate the input data of  $V$  from input cells to internal cells; the draining points of  $V$  serve to propagate the output data of  $V$  from internal cells to output cells; the relaying points of  $V$  serve to



employ intermediate cells as delay buffers for relaying the data of  $V$  between non-neighbouring cells. (E.g., see the  $\vartheta_B$ -path highlighted in Fig. 4-3(b).)

## 4.5 The Extension of the Index Space

This section describes how pipelining points are generated by means of a two-step extension of the index space. The first step follows from Def. 2.5 and the second step does not apply for the  $\pi$ -model. The extension of the index space leads to the realisation that the non-injectivity of the space-time mapping causes complications in the synthesis of control flow for one-dimensional arrays. In addition, it allows several different formulations of the communication constraint.

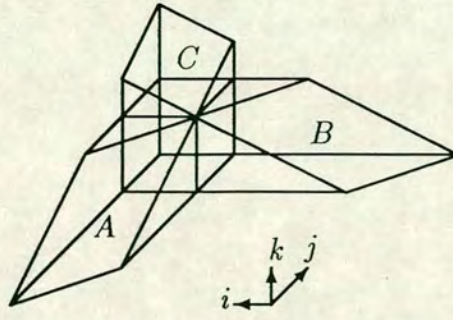
The generation of the extended index space depends on the allocation vector. In the first step, we impose the restriction of border communication, just as for  $(n-1)$ -dimensional arrays (Def. 2.5). In the second step, we impose the restriction of neighbouring communication for the  $\chi$ -model. The new data dependence vector  $\delta_V$  of  $\vartheta_V$  is defined by:

$$\begin{aligned} \delta_V &= \text{if flow}(V)=0 \rightarrow \vartheta_V/\lambda\vartheta_V \ \square \ \text{flow}(V)\neq 0 \rightarrow \vartheta_V \ \text{fi} \quad (\pi\text{-model}) \\ \delta_V &= \text{if flow}(V)=0 \rightarrow \vartheta_V/\lambda\vartheta_V \ \square \ \text{flow}(V)\neq 0 \rightarrow \vartheta_V/|\sigma\vartheta_V| \ \text{fi} \quad (\chi\text{-model}) \end{aligned} \quad (4.4)$$

The directed arc represented by  $\vartheta_V$  between any two neighbouring points in the same  $\vartheta_V$ -path that was created in the previous step is sliced into  $\Gamma_V$  consecutive directed arcs  $\delta_V$ . This creates  $\Gamma_V-1$  points in the original arc. Consider one such point  $K$ ; cell  $\text{place}(K)$  serves at step  $\text{step}(K)$  as a delay buffer to propagate an element of  $V$  from the cell at location  $\text{place}(I)$  to the cell at location  $\text{place}(J)$ .

**Remark** Recall the interpretation of  $\bar{\vartheta}_V$  in the two one-dimensional array models. Its first component decremented by 1 represents the number of delay buffers associated with a channel of  $V$ . Its second component represents the direction and length of that channel. Such an interpretation is consistent with the standard interpretation of channel connections and buffer distributions as in (2.14), since  $\bar{\vartheta}_V$  is the image of  $\delta_V$  under the space-time mapping as defined in (2.14).  $\square$





**Figure 4-4:** Part of the extended index space with respect to the allocation vector of Fig. 4-3 for matrix product ( $m=4$ ). The cube depicts the original index space. The three polytopes labelled  $A$ ,  $B$  and  $C$  represent the points created in the extension along direction  $-\vartheta_A$ ,  $-\vartheta_B$  and  $\vartheta_C$ , respectively. The extension along  $-\vartheta_C$  is not depicted. Extensions along  $\vartheta_A$  and  $\vartheta_B$  are not necessary because the output of  $A$  and  $B$  is of no interest.

To formalise the extension of the index space, we need the following notation. For  $S \subset \mathbb{Q}^n$  and  $\pi \in \mathbb{Q}^n$ ,  $\text{planes}(S, \pi)$  denotes the union of all parallel hyperplanes that intersect the set  $S$  and whose normal is  $\pi$ :

$$\text{planes}(S, \pi) = \{I \mid I \in \mathbb{Q}^n \wedge \pi I = \pi J \wedge J \in S\} \quad (4.5)$$

Note that  $\text{planes}(\Phi, \sigma)$  plays the rôle in one-dimensional arrays that  $\text{rays}(\Phi, \pm u)$  plays in  $(n-1)$ -dimensional arrays. That is, it is the set of the points in  $\mathbb{Q}^n$  that are mapped to the space-time diagram.

**Definition 4.3** The *extended index space*  $\Psi$  of the index space  $\Phi$  is defined as follows (Fig. 4-4):

$$\Psi_V^s = \begin{cases} \text{if } \text{flow}(V)=0 & \rightarrow \emptyset \\ \square \text{ flow}(V) \neq 0 & \rightarrow \text{planes}(\Phi, \sigma) \cap \text{rays}(\text{fst}(\Phi_V, \vartheta_V) - \delta_V, -\delta_V) \\ \text{fi} & \end{cases}$$

$$\Psi_V^d = \begin{cases} \text{if } \text{flow}(V)=0 & \rightarrow \emptyset \\ \square \text{ flow}(V) \neq 0 & \rightarrow \text{planes}(\Phi, \sigma) \cap \text{rays}(\text{lst}(\Phi_V, \vartheta_V) + \delta_V, \delta_V) \\ \text{fi} & \end{cases}$$

$$\Psi_V^r = (\text{rays}(\Phi_V, -\delta_V) \cap \text{rays}(\Phi_V, \delta_V)) \setminus \Phi_V$$

$$\Psi_V = \Phi_V \cup \Psi_V^s \cup \Psi_V^r \cup \Psi_V^d$$



$$\begin{aligned}
\Psi^P &= (\forall V : V \in \mathcal{V} : \Psi_V^s \cup \Psi_V^r \cup \Psi_V^d) \\
\Psi &= \Phi \cup \Psi^P \\
\Psi_V^\perp &= \Psi \setminus \Psi_V
\end{aligned}$$

The portion attributed to variable  $V$  is  $\Psi_V$  (compare Def. 2.5). □

We have given a formal definition of the extended index space for  $(n-1)$ -dimensional arrays in Def. 2.5. Every set defined in Def. 4.3, except  $\Psi_V^r$ , has a parallel in Def. 2.5. The interpretation of the points of the two corresponding sets is identical. The points of  $\Psi_V^r$  are mapped to the relaying points of  $V$ . Because of the existence of relaying points, the second step in the extension of the index space depends on the scheduling vector in the presence of stationary variables.

The points of  $\Psi \setminus \Phi$  may be rational rather than integral:  $\Psi \subset \mathbb{Q}^n$ . The images of the points in  $\Psi \setminus \Phi$  are the pipelining points at which data may be propagated. The computations at the remaining pipelining points are undefined. They can be interpreted as propagating the undefined value  $\perp$ .

The extended source UREs over the extended index space are as defined in Sect. 2.5 except that all data dependence vectors  $\vartheta_V$  are replaced by  $\delta_V$ .

The following concepts are analogous to those defined for the index space. The set of *new* first computation points of variable  $V$  over the extended index space is  $\text{fst}(\Psi_V, \delta_V)$ . The set of *new* last computation points of variable  $V$  over the extended index space is  $\text{lst}(\Psi_V, \delta_V)$ . The points of  $\text{fst}(\Psi_V, \delta_V)$  are mapped to input cells. We write  $\mathbb{I}'_V$  for the equivalent relation on  $\text{fst}(\Psi_V, \delta_V)$  given by

$$(\forall I, J : I, J \in \text{fst}(\Psi_V, \delta_V) : I \mathbb{I}'_V J \iff \text{pi}(V(I)) = \text{pi}(V(J))) \quad (4.6)$$

The points in a  $\mathbb{I}'_V$ -class are mapped to the same input cell. Thus, the communication constraint ensures the injectivity of `step` for each  $\mathbb{I}'_V$ -class. This formulation of the communication constraint is useful in the validity proof of the space-time mappings generated by the procedure described in Sect. 4.6.

**Lemma 4.3** *The communication constraint is satisfied iff*

$$(\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Psi_V, \delta_V) / \mathbb{I}'_V : \text{step} : S \mapsto \mathbf{Z}))$$



**Proof**

$$\begin{aligned}
& (\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Psi_V, \delta_V) / \mathbb{D}'_V : \text{step} : S \mapsto \mathbf{Z})) \\
\iff & \{(\forall I, J : I, J \in S : \text{place}(I) = \text{place}(J))\} \\
& (\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Psi_V, \delta_V) / \mathbb{D}'_V : \text{input} : S \mapsto \mathbf{Z})) \\
\iff & \{\text{Define } \rho_V : \text{fst}(\Phi_V, \vartheta_V) \rightarrow \text{fst}(\Psi_V, \delta_V), \rho_V(I) = J, \text{ where } J \xrightarrow{\vartheta_V} I. \\
& \text{Then } \rho_V \text{ is a bijection and } \text{input}(V(I)) = \text{input}(V(J))\} \\
& (\forall V : V \in \mathcal{V} : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \mathbb{D}_V : \text{input} : S \mapsto \mathbf{Z})) \quad \square
\end{aligned}$$

The next lemma rephrases the communication constraint by viewing the space-time mapping as a mapping from  $\Psi_V$  to  $\Upsilon$ .

**Lemma 4.4** *The communication constraint is satisfied iff*

$$(\forall V : V \in \mathcal{V} : \pi : \Psi_V \mapsto \mathbf{Z}^2)$$

**Proof** Similar to the proof of Lemma 4.1. □

This lemma indicates that the communication constraint is to ensure that the space-time mapping is injective over  $\Psi_V$ . However, valid space-time mappings from  $\Psi$  to  $\Upsilon$  are generally not injective; one example is the mapping displayed in Fig. 4-3. Take the pipelining point  $\bar{I} = (12, -1)$  highlighted in Fig. 4-3(b) as an example. Elements  $a_{3,1}$ ,  $b_{3,4}$ , and  $c_{1,1}$  are propagated at this point, since the paths  $p(\Upsilon, \bar{\vartheta}_A, \bar{I})$ ,  $p(\Upsilon, \bar{\vartheta}_B, \bar{I})$  and  $p(\Upsilon, \bar{\vartheta}_C, \bar{I})$  intersect there. In the extended source UREs,  $A(3, -8, 1) = a_{3,1}$ ,  $B(1/3, 4, 3) = b_{3,4}$  and  $C(1, 1, 5/2) = c_{1,1}$ , where  $\pi(3, -8, 1) = \pi(1/3, 4, 3) = \pi(1, 1, 5/2) = \bar{I}$ . Point  $\bar{I}$  has three inverse images in the extended index space. In general, a space-time point may have as many inverse images (in the extended index space) as there are data variables – one inverse image per variable. We also note that the inverse images of the three paths, i.e.,  $p(\Psi_A, \delta_A, (3, -8, 1))$ ,  $p(\Psi_B, \delta_B, (1/3, 4, 3))$ , and  $p(\Psi_C, \delta_C, (1, 1, 5/2))$  do not intersect. A similar analysis applies for the pipelining point  $\bar{I} = (6, -4)$  highlighted in Fig. 4-3(a). Elements  $a_{2,1}$  and  $b_{3,4}$  are propagated at this point.

**Theorem 4.3** *A valid space-time mapping may not be injective from  $\Psi$  to  $\Upsilon$ .*

**Proof** The space-time mapping in Fig. 4-3 is not injective from  $\Psi$  to  $\Upsilon$ . □



This theorem is not valid for  $(n-1)$ -dimensional arrays because a valid space-time mapping, in this case, is a bijection from  $\mathbb{Q}^n$  to  $\mathbb{Q}^n$ .

In general, the non-injectivity of the space-time mapping increases the importance of control signals in one-dimensional arrays. Take the array shown in Fig. 4-3(b). For example, cell  $p = -1$  must be instructed at step  $t = 12$  to propagate rather than accumulate element  $c_{1,1}$ . However, the non-injectivity of the space-time mapping is a virtue rather than an evil. It permits the parallelism in the propagation of data between cells to be fully exploited. If  $(t, p)$  is a pipelining point, cell  $p$  at step  $t$  is allowed to propagate up to  $|\mathcal{V}|$  elements; one element for every variable in  $\mathcal{V}$ . If the injectivity of the space-time mapping from  $\Psi$  to  $\Upsilon$  is enforced, a cell can propagate at most one element at a given step. This will certainly increase the latency of the array. On the other hand, propagation control signals may be eliminated for a special class of systolic algorithms; this is the topic of Chap. 6.

## 4.6 A Synthesis Procedure

Based on the results developed in the previous sections, we now present a synthesis procedure for the generation of valid space-time mappings for both models. The procedure first constructs  $n \times n$  matrices and then converts them to  $2 \times n$  matrices, which are valid space-time matrices.

First, we construct a space-time matrix, denoted  $\Pi$ , that is an  $n \times n$  integer matrix, whose first  $n-1$  rows, denoted by  $\Lambda$ , form a *scheduling matrix* and whose last row is the allocation vector.

$$\Pi = \begin{bmatrix} \Lambda \\ \sigma \end{bmatrix} = \begin{bmatrix} \Lambda_1 \\ \Lambda_2 \\ \vdots \\ \Lambda_{n-1} \\ \sigma \end{bmatrix} = \begin{bmatrix} \Lambda_{1,1} & \Lambda_{1,2} & \cdots & \Lambda_{1,n} \\ \Lambda_{2,1} & \Lambda_{2,2} & \cdots & \Lambda_{2,n} \\ \vdots & \vdots & & \vdots \\ \Lambda_{n-1,1} & \Lambda_{n-1,2} & \cdots & \Lambda_{n-1,n} \\ \sigma_1 & \sigma_2 & \cdots & \sigma_n \end{bmatrix}$$

$$\bar{I} = \Pi I = (t_1, t_2, \dots, t_{n-1}, p)$$



Point  $I$  is computed at cell  $p$  at multi-dimensional time step  $(t_1, t_2, \dots, t_{n-1})$ . To realise the systolic array in hardware, we need an  $(n-1)$ -dimensional clock, whose most (least) significant hand represents  $t_1$  ( $t_{n-1}$ ). The scales for the  $n-1$  hands of the clock are represented by a vector  $(s_1, s_2, \dots, s_{n-1})$ . One unit of the  $i$ -th hand is equal to  $s_i$  units of the  $(i+1)$ -st hand.  $s_i$  is the difference between the maximum and minimum of the  $i$ -th time components of all space-time points. We shall later provide a procedure for deriving these scales.

**Theorem 4.4** *Let the validity of space-time mappings be defined in the absence of the delay rule of Def. 2.4. A space-time mapping  $\Pi$  is valid for the source UREs if*

- $(\forall V : V \in \mathcal{V} : \Lambda \vartheta_V > 0)$ . *(Precedence Constraint)*
- $\text{rank}(\Pi) = n$ . *(Computation and Communication Constraint)*

**Proof** The proof of the precedence constraint is similar to that of Thms. 4.1 and 4.2, except we need to consider  $(n-1)$ -dimensional time. When  $\Pi$  is non-singular, it is a bijection from  $\mathbb{Q}^n$  to  $\mathbb{Q}^n$ . The satisfaction of the computation and communication constraint is obvious. □

This theorem is due to Moldovan [50]. The mapping conditions depend on the data dependence vectors but not on the index space. The search for a space-time mapping  $\Pi$  that satisfies these conditions can be formulated as an integer programming problem [64]. The delay rule of Def. 2.4 is not part of mapping conditions of Thm. 4.4. It can be trivially satisfied by an appropriate scaling of the scheduling vector that is converted from a scheduling matrix.

Next, we convert the  $(n-1) \times n$  scheduling matrix to a scheduling vector by a procedure that depends on the index space or, more precisely, on the extended index space. This is accomplished by a calculation of scales  $s_1, s_2, \dots, s_{n-1}$ , from which we obtain a vector  $(g_1, g_2, \dots, g_{n-1})$  such that one unit of the  $i$ -th component is equal to  $g_i$  units of the  $(n-1)$ -st component. This vector is calculated by setting  $g_{n-1} = 1$  and  $g_i = g_{i+1} s_{i+1}$ . By setting  $g_{n-1}$  to 1, we assume that  $(0, 0, \dots, 0, 1)$  is the smallest unit of the  $(n-1)$ -dimensional clock. (Scale  $s_1$  does not contribute



to the calculation of the vector  $(g_1, g_2, \dots, g_{n-1})$ . Then  $(n-1)$ -dimensional time  $(t_1, t_2, \dots, t_{n-1})$  is converted to scalar time  $t = (\sum_{i: 0 < i < n} t_i g_i)$ .

**Procedure 4.1** (Conversion of  $(n-1)$ -dimensional to one-dimensional time)

INPUT: The extended index space  $\Psi$  and a scheduling matrix  $\Lambda$ .

OUTPUT: A scheduling vector  $\lambda$ .

1.  $(\forall i : 1 < i < n : s_i = (\max_{I, J : I, J \in \Psi} |\Lambda_i(I - J)| + 1))$ .
2.  $(\forall i : 0 < i < n : \Delta_i = (\min_{I, J : I \neq J \wedge I, J \in \Psi} |\Lambda_i(I - J)|))$ .
3.  $g_{n-1} = 1/\Delta_{n-1}$ ,  $(\forall i : 0 < i < n-1 : g_i = g_{i+1} s_{i+1} / \Delta_i)$ .
4.  $(\forall k : 0 < k \leq n : \lambda_k = (\sum_{i: 0 < i < n} g_i \Lambda_{i,k}))$ . □

Let us explain the rôle that quantities  $\Delta_1, \Delta_2, \dots, \Delta_{n-1}$  play in Proc. 4.1.  $\Delta_i$  represents the smallest difference among the  $i$ -th time components of all space-time points. If all the  $i$ -th time components are integers, then  $\Delta_i \in \mathbf{Z}^+$ .  $\Delta_i$  can be disregarded (i.e., treated as 1). But, retaining  $\Delta_i$  minimises the components  $g_1, g_2, \dots, g_i$  and thus the latency of the systolic array. If the  $i$ -th time components of some space-time points are not integers (because some points of  $\Psi$  are not integers),  $\Delta_i$  may not be integral.  $\Delta_i$  is essential in ensuring that different multi-dimensional time steps are converted to different scalar time steps.

**Lemma 4.5** *If  $\Pi$  satisfies Thm. 4.4, the transformation of  $\Lambda$  to  $\lambda$  by Proc. 4.1 preserves the following properties:*

1.  $(\forall I, J : I, J \in \Psi : \Lambda I = \Lambda J \iff \lambda I = \lambda J)$  (Bijectivity)
2.  $(\forall I, J : I, J \in \Psi : \Lambda I < \Lambda J \implies \lambda I < \lambda J)$  (Precedence)
3.  $(\forall V : V \in \mathcal{V} : (\forall I, J : I, J \in \Psi : \Lambda I = \Lambda J + \Lambda \vartheta_V \implies \lambda I = \lambda J + \lambda \vartheta_V))$  (Linearity)



**Proof** Algebraic manipulation. □

**Procedure 4.2** (Construction of the space-time mappings for the source UREs with respect to either the  $\pi$ -model or the  $\chi$ -model)

INPUT: The source UREs.

OUTPUT: Valid space-time mappings  $\pi$ .

1. Find space-time mappings  $\Pi$  that satisfy Thm. 4.4.
2. Obtain the extended index spaces  $\Psi$  by Def. 4.3 (with respect to either the  $\pi$ -model or the  $\chi$ -model).
3. Transform  $\Lambda$  to  $\lambda$  by Proc. 4.1.
4. Scale  $\lambda$  to satisfy the delay constraint for the  $\chi$ -model. □

**Theorem 4.5** *The space-time mapping  $\pi$  returned by Proc. 4.1 is valid.*

**Proof** If a space-time mapping  $\Pi$  satisfies Thm. 4.4,  $\pi$  satisfies the precedence and computation constraint by Lemma 4.5. Scaling  $\lambda$  ensures that  $\pi$  satisfies the delay constraint. We remain to prove that  $\pi$  satisfies the communication constraint.

$$\begin{aligned}
& \Pi \text{ is non-singular} \\
\Rightarrow & \{\text{Non-singularity}\} \\
& \Pi : \mathbb{Q}^n \mapsto \mathbb{Q}^n \\
\Rightarrow & \{\text{Restriction of } \Pi \text{ to a subset of } \mathbb{Q}^n\} \\
& \Pi : \text{fst}(\Psi_V, \delta_V) \mapsto \mathbb{Q}^n \\
\Rightarrow & \{\mathbb{Q}'_V \text{ of (4.6)}\} \\
& (\forall S : S \in \text{fst}(\Phi_V, \delta_V) / \mathbb{Q}'_V : \Pi : S \mapsto \mathbb{Q}^n) \\
\Rightarrow & \{\text{Lemma 4.5}\} \\
& (\forall S : S \in \text{fst}(\Phi_V, \delta_V) / \mathbb{Q}'_V : \pi : S \mapsto \mathbb{Q}^2) \\
\iff & \{(\forall I, J : I, J \in S : \text{place}(I) = \text{place}(J))\} \\
& (\forall S : S \in \text{fst}(\Phi_V, \delta_V) / \mathbb{Q}'_V : \text{step} : S \mapsto \mathbb{Q})\}
\end{aligned}$$



$\implies$  {Appropriate scaling of  $\lambda$ ; Lemma 4.3}

The communication constraint is satisfied □

The space-time mappings generated by Proc. 4.2 may be inefficient because Proc. 4.1 preserves the injectivity of the space-time mapping  $\Pi$  over the extended index space (Stat. 1 of Lemma 4.5). This follows from two observations. First, the parallel injection and ejection of the elements of different variables at the same cell is excluded by Proc. 4.1. But the communication constraint only requires the injectivity of  $\text{step}$  over every  $\mathbb{O}'_V$ -class. Second, Thm. 4.5 still holds if we allow two points scheduled at different multi-dimensional time steps to be converted to the same scalar time step, provided that they are mapped to different cells, i.e., if the scheduling vector obtained satisfies a weakened version of Stat. 1 of Lemma 4.5:

$$(\forall I, J : I, J \in \Psi : (\Lambda I = \Lambda J \implies \lambda I = \lambda J) \wedge (\Lambda I \neq \Lambda J \wedge \sigma I = \sigma J \implies \lambda I \neq \lambda J)).$$

## 4.7 Related Work

Depending on how the scheduling vector is derived, we distinguish two methods for the synthesis of data flow for one-dimensional arrays: the one-dimensional method searches the scheduling vector directly, and the multi-dimensional method does so via an  $(n-1)$ -dimensional scheduling matrix. In the one-dimensional method, one applies a  $2 \times n$  space-time matrix  $\pi$  to the source UREs. In the transformed UREs, one index represents time and the other space. In the multi-dimensional method, one applies an  $n \times n$  space-time matrix  $\Pi$  to the source UREs. In the transformed UREs,  $n-1$  indices represent time and the remaining index represents space. Then, the  $(n-1)$ -dimensional scheduling matrix is converted to a one-dimensional scheduling vector.

Ramakrishnan, Fussell and Silberschatz [61] were the first to introduce the  $\chi$ -model. Their synthesis of the data flow amounts to embedding two- or three-dimensional data dependence graphs in a space-time diagram. The correctness of an embedding is ensured by certain embedding rules, which are equivalent in



nature to the mapping conditions specified in Thm. 4.1. However, this graph embedding approach does not lead to a systematic embeddings of any UREs. This method has been recently improved by the one-dimensional method proposed by Lee and Kedem [41]. They adopt the  $\chi$ -model and formulate the previous embedding rules in terms of relations on the data dependence vectors (Thm. 4.2). This data dependence approach is superior to the previous graph embedding approach in that it applies for any  $n$ -dimensional UREs and that the search for valid space-time mappings can be formulated as linear and integer programming. Rao [64] presents a multi-dimensional method without providing a model; his concept of extended index space conforms to the  $\pi$ -model but not the  $\chi$ -model. Kumar and Tsai [33] present a method that converts two-dimensional arrays to one-dimensional arrays. They adopt the  $\pi$ -model, although this is not explicitly mentioned. Wong and Delosme [78] present a procedure for converting  $(n-2)$ -dimensional to one-dimensional time in the context of the synthesis of two-dimensional arrays. The index space rather than the extended index space is used in the conversion. Hence, input and output data may have to be handled at internal cells.

## 4.8 Conclusion

We have formally defined the two most frequently used one-dimensional array models: the  $\pi$ -model allows non-neighbouring connections and the  $\chi$ -model does not. We have extended previously known mapping conditions to both models. This allows synthesis methods formulated for one model to be applied to the other model. For example, Lee and Kedem's one-dimensional method [41] can be extended to the  $\pi$ -model. The only requirement is that we must choose the allocation vector before the scheduling vector, because the definition of  $\mathbb{I}_V$  (or  $\mathbb{I}'_V$ ) in the communication constraint depends on the allocation vector.

The synthesis procedure (Proc. 4.2) is built on previously known results, but it applies for both models. Valid space-time mappings can be obtained for either of the two models, depending on how the extended index space is defined.



We have represented the communication constraint in a variety of forms: each equivalently enforces the restriction that a channel cannot transfer more than one datum per step but conveys the properties of the space-time mapping from a distinct perspective. For example, Lemma 4.2 justifies the geometrical approach of [10] in which one-dimensional arrays are constructed by a manual embedding of  $\vartheta_V$ -paths in the space-time diagram.

We have investigated properties characteristic of one-dimensional arrays. First, we have characterised the regular distribution of pipelining points over the space-time diagram. This will be the basis of the construction of propagation control signals for one-dimensional arrays in Chap. 5. Second, we have shown how the non-injectivity of a space-time mapping can be understood by an extension of the index space. The non-injectivity of the space-time mapping makes the implementation of the control UREs inefficient. This is our reason to look for a different but more efficient specification of propagation control signals for one-dimensional arrays in Chap. 5.

The concept of the extended index space can be used in the optimisation of a schedule. The points that are mapped to  $t_{\text{fst}}$  and  $t_{\text{lst}}$  must be extreme points in  $(\bigcup V : V \in \mathcal{V} : \text{fst}(\Psi_V, \delta_V) \cup \text{lst}(\Psi_V, \delta_V))$ . Since a valid space-time mapping must satisfy the precedence constraint, some extreme points can be excluded as candidates. Assume that  $F$  ( $L$ ) is the set of extreme points that may be scheduled at the first (last) step  $t_{\text{fst}}$  ( $t_{\text{lst}}$ ). A space-time mapping that minimises the latency must minimise the objective function:

$$(\max I : I \in L : \lambda I) - (\min I : I \in F : \lambda I)$$

The computation constraint aims at preventing intra-cell parallelism. This is enforced by requiring too much: no two points that are scheduled simultaneously can be mapped to the same cell. The necessary and sufficient condition should state that no concurrent computations *for a fixed variable* can be mapped to the same cell. (So, two points can share the same image iff every variable is undefined at at least one of the two points). This is guaranteed by the communication constraint. Therefore, the computation constraint can be neglected.



The communication constraint defined in this chapter is overconstrained. Two  $\vartheta_V$ -paths may be mapped to the same  $\bar{\vartheta}_V$ -path if, for example, variable  $V$  is undefined at all the points in one of the two paths. In general, two  $\vartheta_V$ -paths can share the same image iff, for every pair of points that share the same image (one point per path), variable  $V$  is undefined at at least one of the two points, i.e., iff

$$(\forall x, y : x, y \in \text{paths}(\Psi_V, \delta_V) : \\ \bar{x} = \bar{y} \implies (\forall I, J : I \in x \wedge J \in y : \bar{I} = \bar{J} \implies V(\bar{I}) = \perp \vee V(\bar{J}) = \perp))$$

Unfortunately, the verification of this formula depends on both the space-time mapping and the index space. In addition, more complex control is required.



## Chapter 5

# Control Flow Synthesis for One-Dimensional Systolic Arrays

### 5.1 Introductory Remarks

In Chap. 3, we have provided a constructive method for the synthesis of control signals for  $(n-1)$ -dimensional arrays from  $n$ -dimensional UREs. The basic idea is to transform the source UREs to pipelined UREs. Then, finding  $(n-1)$ -dimensional arrays is just finding space-time mappings that are valid with respect to the pipelined UREs and the  $(n-1)$ -dimensional array model.

In Chap. 4, we have defined two one-dimensional array models and stated mapping conditions for the synthesis of data flow from  $n$ -dimensional UREs for these two models. The pipelined UREs constructed in Chap. 3 are a special form of UREs. Just like for the synthesis of  $(n-1)$ -dimensional arrays, finding one-dimensional arrays is just finding valid space-time mappings with respect to the pipelined UREs and one-dimensional array models.

This chapter is concerned with the construction of control UREs for one-dimensional arrays. To illustrate why the specification of the control UREs constructed previously may result in inefficient one-dimensional arrays, we state in a theorem the mapping conditions for the validity of space-time mappings with



respect to the pipelined UREs and the two one-dimensional array models; this theorem is a direct application of Thms. 4.1 and 4.2 to the pipelined UREs. From now on, we write  $\mathcal{V}_c$  ( $\mathcal{V}_p$ ) for the set of computation (pipelining) control variables, and  $\mathcal{V}_d$  for the set of data variables.  $\mathcal{V}$  is the union of  $\mathcal{V}_c$ ,  $\mathcal{V}_p$  and  $\mathcal{V}_d$ .

**Theorem 5.1** *A space-time mapping is valid for the pipelined UREs (in either the  $\pi$ -model or the  $\chi$ -model) iff*

$$\bullet (\forall V : V \in \mathcal{V} : \lambda \vartheta_V > 0) \quad (\text{Precedence Constraint})$$

$$\bullet \pi : \Phi \mapsto \mathbb{Z}^2 \quad (\text{Computation Constraint})$$

$$\bullet (\forall V : V \in \mathcal{V} : \sigma \vartheta_V \mid \lambda \vartheta_V) \quad (\text{Delay Constraint for the } \chi\text{-Model})$$

$$\bullet (\forall V : V \in \mathcal{V}_d : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \textcircled{1}_V : \text{input} : S \mapsto \mathbb{Z})) \quad (5.1)$$

$$(\forall V : V \in \mathcal{V}_c : (\forall S : S \in \text{fst}(\Psi_V, \vartheta_V) / \textcircled{1}'_V : \text{input} : S \mapsto \mathbb{Z})) \quad (5.2)$$

$$(\forall V : V \in \mathcal{V}_p : (\forall S : S \in \text{fst}(\Psi_V, \vartheta_V) / \textcircled{1}'_V : \text{input} : S \mapsto \mathbb{Z})) \quad (5.3)$$

(Communication Constraint)

**Proof** Thms. 4.1 and 4.2. □

Domain  $\Psi_V$  of a control variable  $V$  in (5.2) and (5.3) is defined in (3.24):

$$\Psi_V = \text{planes}(\Phi, \sigma) \cap \text{rays}(D_V, \pm \vartheta_V) \quad (5.4)$$

$D_V$  is the set of points  $I$  such that  $V(I)$  is referenced in the guards of the pipelined UREs. It has the following relationship with the index space  $\Phi$  and the extended index space  $\Psi$  (Sect. 3.4):

$$(\forall V : V \in \mathcal{V}_c : \text{rays}(D_V, \pm \vartheta_V) \subseteq \text{rays}(\Phi, \pm \vartheta_V)) \quad (5.5)$$

$$(\forall V : V \in \mathcal{V}_p : \text{rays}(D_V, \pm \vartheta_V) = \text{rays}(\Psi \cup H^c.V, \pm \vartheta_V) \cup \text{rays}(\Phi, \pm \vartheta_V)) \quad (5.6)$$

The communication constraint for a variable ensures that all input data of the variable are input at distinct time steps. In general, the higher the number of input data of a variable, the higher the latency of the systolic array. This makes it desirable to keep the number of inputs of a control variable as small as possible.



The CCUREs serve to distinguish different types of computation points. They need only be defined for the index space.  $\text{fst}(\Psi_V, \vartheta_V)$  of a computation control variable  $V$  contains no more elements than  $\text{fst}(\Phi, \vartheta_V)$ . The latency of the systolic array may be retained when it is possible to choose data dependence vectors as the control dependence vectors for computation control variables (Sect. 3.5). The PCUREs serve to distinguish pipelining points from computation points. They must be defined for the extended index space. However,  $\text{fst}(\Psi_V, \vartheta_V)$  of a propagation control variable  $V$  contains more elements than  $\text{fst}(\Phi, \vartheta_V)$ , leading to a higher latency of the systolic array. (This is not true for  $(n-1)$ -dimensional arrays because the corresponding mapping conditions depend only on data and control dependence vectors but not on the (extended) index space.) We avoid the higher latency by providing a different construction of the PCUREs. The domain of a propagation control variable  $V$  in the new PCUREs is defined for the extended index space  $\Psi$ . But, the construction of the new PCUREs will allow us to replace  $\Psi_V$  in the communication constraint (5.3) by a subset of  $\Psi$  called the *communication constraint domain* of  $V$  and denoted  $\Omega_V$ :

$$\Omega_V = \text{planes}(\Phi, \sigma) \cap \text{rays}(D_V, \pm \vartheta_V) \quad (5.7)$$

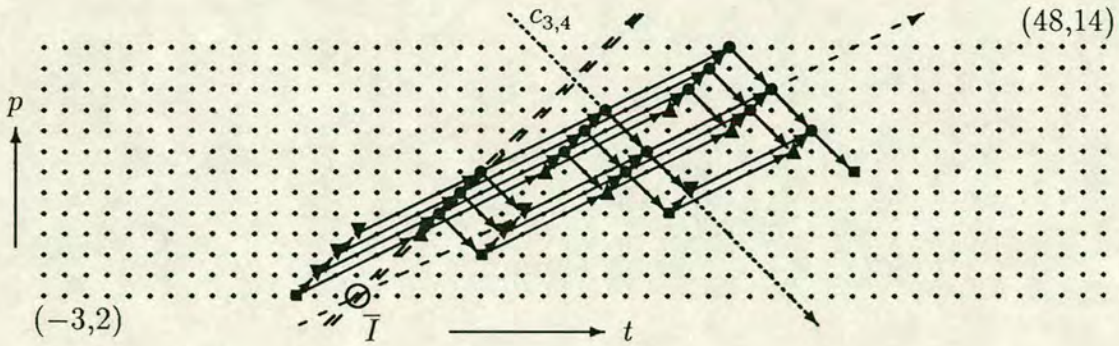
where  $D_V$  satisfies (5.5), eliminating the need for (5.6). Before doing so, we illustrate with LU-decomposition why the behaviour of the original PCUREs is not preserved if we directly substitute  $\Omega_V$  for  $\Psi_V$  of the communication constraint (5.3). This also indicates that the PCUREs play a more important rôle for one-dimensional arrays than for  $(n-1)$ -dimensional arrays.

### Example 5.1 $4 \times 4$ LU-Decomposition

Recall the pipelined UREs constructed in Sect. 3.7.2. The CCUREs consist of two computation control variables  $P$  and  $Q$ . The PCUREs were not presented. They follow from the set  $\mathfrak{B} = \{B_1, B_2, B_3\}$  given in Sect. 3.7.1.3 and (3.18). By our convention, we write  $C.B_1$ ,  $C.B_2$  and  $C.B_3$  for the three corresponding propagation control variables. By (3.19), we obtain

$$\begin{aligned} \mathcal{P}(\Phi) &= (\forall i : 0 < i \leq 3 : C.B_i(I) = 1) \\ \mathcal{P}(\mathcal{L}_\Psi \Phi) &= (\exists i : 0 < i \leq 3 : C.B_i(I) = 0) \end{aligned} \quad (5.8)$$





**Figure 5-1:** The space-time diagram for LU-decomposition ( $m = 4$ ). The two (identical) long-dashed paths depict  $p(\Upsilon, \bar{\vartheta}_{C.B_1}, \bar{I})$  and  $p(\Upsilon, \bar{\vartheta}_{C.B_3}, \bar{I})$ ; they are equal and should be viewed as overlapping each other. The short-dashed path depicts  $p(\Upsilon, \bar{\vartheta}_{C.B_2}, \bar{I})$ . The dotted path will be referred to in Sect. 5.3.

$\mathcal{P}(\Phi)$  holds at point  $I$  iff  $I \in \Phi$ . (That is,  $\mathcal{P}(\mathcal{L}_\Psi \Phi)$  holds at point  $I$  iff  $I \in \mathcal{L}_\Psi \Phi$ .) For the purpose of illustration, we choose  $\vartheta_P = \vartheta_{C.B_1} = \vartheta_{C.B_3} = (0, 1, 0)$  and  $\vartheta_Q = \vartheta_{C.B_2} = (1, 0, 0)$ . The following space-time mapping, which describes Ramakrishnan and Varman's one-dimensional array for matrix product, is valid by Thm. 5.1 if we replace  $\Psi_V$  of (5.3) by  $\Omega_V$  of (5.7) (Fig. 5-1):

$$\pi_{\text{even}} = \begin{bmatrix} 2m-2 & 1 & m/2 \\ m-1 & 1 & -m/2 \end{bmatrix} \quad \pi_{\text{odd}} = \begin{bmatrix} 2m & 1 & (m+1)/2 \\ m & 1 & -(m+1)/2 \end{bmatrix} \quad (5.9)$$

In the resulting systolic array, control variables  $P$ ,  $C.B_1$  and  $C.B_3$  each travel at the same velocity as data variable  $A$ , and control variables  $Q$  and  $C.B_2$  each travel at the same velocity as data variable  $B$ .

Two observations can be made on this example. First, the behaviour of the PCUREs is not preserved by the chosen space-time mapping. Consider the three dashed paths pointing towards the northeast shown in Fig. 5-1. If we input control signal **1** at each of these three paths, the pipelining point highlighted by a circle will be interpreted as a computation point, since  $\overline{\mathcal{P}(\Phi)}$  holds at this point. If we input control signal **0** at one of these three paths, then the computation points in the path will be interpreted as pipelining points, since  $\overline{\mathcal{P}(\mathcal{L}_\Psi \Phi)}$  holds at all points of the path. The problem is that two different  $\vartheta_{C.B_i}$ -paths for a fixed  $i$  are mapped to the same  $\bar{\vartheta}_{C.B_i}$ -path – a condition that is not permitted by the premises of Thm. 4.2.



Second, Lemma. 3.10 does not apply due to the non-injectivity of the space-time mapping. For example,  $\overline{\mathcal{P}(\Phi_\bullet)}$  holds not only at the computation points highlighted by fat dots, as intended, but also at the pipelining point highlighted by a circle in Fig. 5-1. This implies that  $\overline{\Phi_\bullet} \neq \overline{\Phi'_\bullet}$ , albeit  $\Phi_\bullet = \Phi'_\bullet$ . Thus, predicate  $\mathcal{P}(\Phi)$  in the pipelined equation of (3.21) is essential and cannot be disregarded as in the case of  $(n-1)$ -dimensional arrays. So Thm. 3.3.2 is not valid. This indicates that propagation control is more important for one-dimensional arrays than for  $(n-1)$ -dimensional arrays.  $\square$

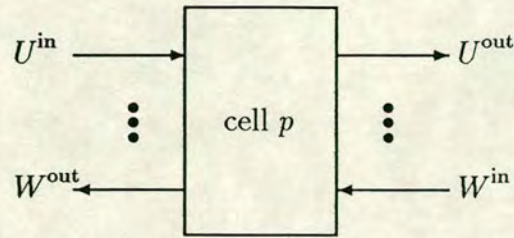
In this chapter, we focus on the construction of PCUREs for one-dimensional arrays. To do so, we need to lift one of the four restrictions imposed on the specification of control signals in Sect. 3.2, namely, Rst. 4. That is to say, we shall allow the control value of some propagation control variable to change during its propagation throughout the array. We call a control variable an *evolution control variable* if some value of the variable changes before it is ejected from the array.

The rest of this chapter is organised as follows. Sect. 5.2 describes some notation and the basic idea underlying the construction of the PCUREs. Sect. 5.3 presents the construction of the PCUREs for three-dimensional UREs. Sect. 5.4 generalises this construction to  $n$ -dimensional UREs. (We do not consider two-dimensional UREs here; they have already been covered in Chap. 3.) Sect. 5.5 comments on some related work. Sect. 5.6 illustrates our method with two examples. Sect. 5.7 contains the conclusion of the chapter.

## 5.2 The Synthesis of Propagation Control Flow

To avoid unduly complex notation, we present the specification of propagation control signals in the form of programs (rather than in the form of UREs). The program for the host, called the *host program*, corresponds to the input equations and injects propagation control signals to the array. The program for the cells of the array, called the *cell program*, corresponds to the computation equations and specifies the control signals of the output channels of a cell based on the control





**Figure 5-2:** The notation for specifying the cell program.

signals of the input channels of the cell. The cell program is identical for all cells and is executed by every cell at every time step.

How can UREs be expressed as programs? Since all cells are given the same cell program, we can abstract from the index vectors of propagation control variables. This results in one program that will be executed by every cell of the array.

Let us consider the situation where cell  $p$  is computing point  $\bar{I} = (t, p)$  at step  $t$  (Fig. 5-2). The cell receives its input signal  $V(\bar{I} - \bar{\vartheta}_V)$  of propagation control variable  $V$  at step  $t-1$  and sends its output signal  $V(\bar{I})$  at step  $t$ . We abstract from the index vectors of  $V$  by writing  $V^{\text{in}}$  for  $V(\bar{I} - \bar{\vartheta}_V)$  and  $V^{\text{out}}$  for  $V(\bar{I})$ . That is,  $V^{\text{in}}$  ( $V^{\text{out}}$ ) denotes the input (output) channel for  $V$  at this cell.

The domain of every propagation control variable is the extended index space  $\Psi$ . There are two types of propagation control variables  $V$ :

- *Pipelining Control Variables.* Every cell directly sends the control signal received at the input channel of  $V$  to its corresponding output channel. That is, the computation equation is given by

$$V^{\text{out}} = V^{\text{in}} \quad (5.10)$$

If the cell is an input cell,  $V^{\text{in}}$  represents one of the values in  $\text{sig}(V)$  (which is the set of control values of  $V$  in Sect. 3.3) to be injected by the host program.

- *Evolution Control Variables.* Every cell may change the control signal received at the input channel of  $V$  before sending the (changed) control signal



to the corresponding output channel. The computation equation is given by

$$V^{\text{out}} = \begin{cases} \text{if } B_1(W^{\text{in}}, \dots) \rightarrow f_1(V^{\text{in}}) \\ \square B_2(W^{\text{in}}, \dots) \rightarrow f_2(V^{\text{in}}) \\ \dots \\ \square B_b(W^{\text{in}}, \dots) \rightarrow f_b(V^{\text{in}}) \\ \mathbf{fi} \end{cases} \quad (5.11)$$

The guard  $B_i(W^{\text{in}}, \dots)$  can always be written in disjunctive normal form, where each disjunct consists of a conjunction of tests of an argument (i.e., a propagation control variable) for a control signal.  $f_i(V^{\text{in}})$  is a function that recursively defines the control value of  $V^{\text{out}}$  in terms of argument  $V^{\text{in}}$ .

Similarly, if the cell is an input cell,  $V^{\text{in}}$  represents one of the values in  $\text{sig}(V)$  to be injected by the host program.

The propagation control flow is *correct* if there exists a construction of a characteristic function  $\chi_{\bar{\Phi}}$  in the propagation control variables of  $\bar{\Phi}$  in  $\Upsilon$ :

$$\chi_{\bar{\Phi}} : \Upsilon \rightarrow \{\chi_c, \chi_p\}, \quad \chi_{\bar{\Phi}}(\bar{I}) = \text{if } \bar{I} \in \bar{\Phi} \rightarrow \chi_c \square \text{else} \rightarrow \chi_p \mathbf{fi} \quad (5.12)$$

That is,  $\bar{I}$  is a computation point if  $\chi_{\bar{\Phi}}(\bar{I}) = \chi_c$  and a pipelining point otherwise.

The notation  $\varphi(\bar{I})$  with  $\bar{I} = (t, p)$  stands for the tuple of propagation control signals associated with the input channels of cell  $p$  at step  $t-1$ , one component per input channel.

$$\varphi(\bar{I}) = (V^{\text{in}}, \dots) = (V(\bar{I} - \bar{\vartheta}_V), \dots) \quad (5.13)$$

The basic idea underlying the construction of the PCUREs is to exploit the regular distribution of pipelining points in the space-time diagram. Remember that the space-time points of  $\Upsilon$  are divided into five categories with respect to a fixed variable  $V$  (Sect. 4.4): the set of soaking points denoted by  $\Upsilon_V^s$ , the set of draining points denoted by  $\Upsilon_V^d$ , the set of relaying points denoted by  $\Upsilon_V^r$ , the set of undefined points denoted by  $\Upsilon_V^\perp$ , and the set of computation points  $\bar{\Phi}$ . By choosing one propagation control variable, say,  $V$  as a reference, we shall construct



the PCUREs in such a way that the five types of points with respect to  $V$ , i.e., the five sets  $\{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_V^s\}$ ,  $\{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_V^d\}$ ,  $\{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_V^r\}$ ,  $\{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_V^\perp\}$ , and  $\{\varphi(\bar{I}) \mid \bar{I} \in \bar{\Phi}\}$  are disjoint.

The PCUREs constructed this way apply for both one-dimensional array models. We shall conduct our presentation with respect to the  $\chi$ -model only. This is justified since a space-time mapping that is valid with respect to the  $\chi$ -model is also be valid with respect the  $\pi$ -model.

The construction of the PCUREs depends on the dimension of the index space. To provide a better understanding, we shall first present the PCUREs for three-dimensional UREs and then describe the generalisation to  $n$ -dimensional UREs.

We say a space-time mapping is *valid for a variable* if the computation constraint is satisfied and the precedence, delay and communication constraint for that variable are satisfied. Our construction of the PCUREs requires that all propagation control variables are moving. In Sect. 5.3.4, we shall state this constraint explicitly as part of the mapping conditions for the validity of space-time mappings. In the sections before, we shall assume this constraint implicitly.

## 5.3 The PCUREs for Three-Dimensional UREs

### 5.3.1 The Evolution Control Flow

We need to use only one evolution control variable; we name it  $E$ . Its specification exploits the regular distribution of pipelining points over the space-time diagram. In this section and subsequent sections, when we refer to soaking points, draining points, relaying points, undefined points, first computation points and last computation points with no explicit reference to a variable, we mean the points associated with the evolution control variable.

The PCUREs are supposed to ensure that the sets

$$\{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_E^s\}, \{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_E^d\}, \{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_E^r\}, \{\varphi(\bar{I}) \mid \bar{I} \in \Upsilon_E^\perp\}, \{\varphi(\bar{I}) \mid \bar{I} \in \bar{\Phi}\}$$



are disjoint. We classify the remaining propagation control variables, which are pipelining control variables, into two categories:

- *Initialisation Control Variables*  $F_{1,1}, F_{1,2}, F_{2,1}, F_{2,2}, \dots, F_{f,1}, F_{f,2}$ . Let  $B(F)$  be a predicate in the initialisation control variables. We shall construct initialisation control variables such that  $B(F)$  is a characteristic function of the first computation points of  $E$ :

$$B(F) = \text{if } \bar{I} \in \overline{\text{fst}(\Phi, \vartheta_E)} \rightarrow \text{true} \quad \square \quad \text{else} \rightarrow \text{false} \quad \text{fi} \quad (5.14)$$

The idea behind the construction of the initialisation control variables is to partition  $\text{fst}(\Phi, \vartheta_E)$  into a number of subsets,  $\text{fst}(\Phi, \vartheta_E)_1, \text{fst}(\Phi, \vartheta_E)_2, \dots, \text{fst}(\Phi, \vartheta_E)_f$ , such that  $\dim(\text{fst}(\Phi, \vartheta_E)_i) \leq n-1$ . We write

$$\textcircled{F} = \{\text{fst}(\Phi, \vartheta_E)_i \mid 0 < i \leq f\}$$

One possible solution is to first decompose  $\text{fst}(\Phi, \vartheta_E)$  into the union of facets of the index space  $\Phi$  and then build a partition from this decomposition. (These facets intersect at their boundaries.) In our construction, we associate two initialisation control variables,  $F_{i,1}$  and  $F_{i,2}$ , with  $\text{fst}(\Phi, \vartheta_E)_i$ .  $\text{sig}(F_{i,1}) = \{f_{i,1}, \bar{f}_{i,1}\}$ .  $\text{sig}(F_{i,2}) = \{f_{i,2}, \bar{f}_{i,2}\}$ . The reason for employing exactly two initialisation control variables will become clear later on in this section (Thm. 5.3).

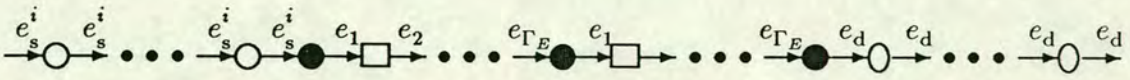
- *Termination Control Variables*  $L_1, L_2, \dots, L_\ell$ . Let  $B(L)$  be a predicate in the termination control variables. We shall construct termination control variables such that  $B(L)$  is a characteristic function of the last computation points of  $E$ :

$$B(L) = \text{if } \bar{I} \in \overline{\text{lst}(\Phi, \vartheta_E)} \rightarrow \text{true} \quad \square \quad \text{else} \rightarrow \text{false} \quad \text{fi} \quad (5.15)$$

For reasons of symmetry, the construction of the termination control variables proceeds in a similar way as that of the initialisation control variables. We decompose  $\text{lst}(\Phi, \vartheta_E)$  into a number of subsets,  $\text{lst}(\Phi, \vartheta_E)_1, \text{lst}(\Phi, \vartheta_E)_2, \dots, \text{lst}(\Phi, \vartheta_E)_\ell$ , such that  $\dim(\text{lst}(\Phi, \vartheta_E)_i) \leq n-1$  and the union of these subsets is  $\text{lst}(\Phi, \vartheta_E)$ . We write

$$\textcircled{L} = \{\text{lst}(\Phi, \vartheta_E)_i \mid 0 < i \leq \ell\}$$





**Figure 5-3:** The evolution of the evolution control variable. Arrows depict  $\bar{\vartheta}_E$ . Fat dots are computation points, circles are soaking points, boxes are relaying points and ovals are draining points.

(Our construction does not require  $\mathbb{L}$  to be a partition of  $\text{lst}(\Phi, \vartheta_E)$ .) One possible solution is to decompose  $\text{lst}(\Phi, \vartheta_E)$  into the union of facets of the index space  $\Phi$ . In our construction, we associate one termination control variable,  $L_i$ , with  $\text{lst}(\Phi, \vartheta_E)_i$ .  $\text{sig}(L_i) = \{\ell_i, \bar{\ell}_i\}$ . The reason for employing one termination control variable will be explained in Sect. 5.3.3.

Let us describe the specification of evolution control variable  $E$ . The communication constraint domain of  $E$  is given by

$$\Omega_E = \text{planes}(\Phi, \sigma) \cap \text{rays}(\Phi, \pm \vartheta_E) \quad (5.16)$$

This is to ensure that all  $\vartheta_E$ -paths that contain computation points are mapped to different  $\bar{\vartheta}_E$ -paths in the space-time diagram. Thus, each  $\bar{\vartheta}_E$ -path that contains computation points contains a unique first computation point of  $E$  (Lemma 4.2).

We consider two types of  $\bar{\vartheta}_E$ -paths in the space-time diagram:

- Path  $p(\Upsilon, \bar{\vartheta}_E, \bar{I})$  contains computation points (Fig. 5-3). It must contain one first computation point of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$  for some  $i$ . We input  $e_s^i$  at the first point of the path and adopt this value for all soaking points of  $E$ . (The superscript  $i$  in  $e_s^i$  identifies the soaking points in the  $\bar{\vartheta}_E$ -paths that contain first computation points of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ .)  $B(F)$  holds at the first computation point of the path. There,  $e_s^i$  is converted to  $e_1$ . Then, if a point receives element  $e_k$  of  $E$ , it sends  $e_{(k \bmod \Gamma_E)+1}$ ;  $k$  is the distance of a relaying point from its preceding computation point. Thus, the considered element of  $E$  periodically adopts the values  $e_1, e_2, \dots, e_{\Gamma_E}$ .  $B(L)$  holds at the last computation point of the path. There,  $e_{\Gamma_E}$  is converted to  $e_d$ , the value for all draining points of  $E$ . If path  $p(\Upsilon, \bar{\vartheta}_E, \bar{I})$  contains only one computation point, which is, therefore, both a first and last computation point,  $e_s^i$  changes to  $e_d$  at that point.



- Path  $p(\Upsilon, \bar{\vartheta}_E, \bar{I})$  contains no computation points; it contains only undefined points of  $E$ . We input  $e_u$  at the first point of  $p(\Upsilon, \bar{\vartheta}_E, \bar{I})$  and adopt this value for all the points of the path.

By the construction of  $E$ , we inject  $e_s^i$  at steps  $\text{input}(E(I))$  for  $I \in \text{fst}(\Phi, \vartheta_E)_i$ ; for every  $i$  and  $e_u$  at the remaining steps.  $\text{sig}(E) = \{e_u, e_s^1, e_s^2, \dots, e_s^f, e_1, e_2, \dots, e_{\Gamma_E}, e_d\}$ . For the host and cell program for three-dimensional UREs, see Tabs. 5-1 and 5-2. The notation  $\text{inject}(s, \text{pi}(V^{\text{in}}))$  stands for the injection of control value  $s \in \text{sig}(V)$  to input cell  $\text{pi}(V^{\text{in}})$ .

In Tab. 5-2, the first guard of  $E^{\text{out}}$  selects first computation points of  $E$  but excludes points that are both a first and last computation point of  $E$ . The second guard of  $E^{\text{out}}$  establishes whether a point is a last computation point of  $E$ .  $\mathcal{B}(F)$  holds at first computation points of  $E$  and  $E^{\text{in}} = e_{\Gamma_E}$  holds at computation points that are not first computation points of  $E$ . Thus, if  $\text{fst}(\Phi, \vartheta_E)$  and  $\text{lst}(\Phi, \vartheta_E)$  are disjoint, the first guard of  $E$  can be simplified to  $\mathcal{B}(F)$ . The third guard handles soaking points, draining points and undefined points of  $E$ . The **else** guard of  $E^{\text{out}}$  captures points that are relaying or computation points but neither first nor last computation points of  $E$ .

**Lemma 5.1** *Let the space-time mapping be valid for the evolution control variable.*

1.  $E^{\text{in}} \neq e_s^i$  for every  $i$  holds at the points of  $\Upsilon_E^\perp$ .
2. Let  $\bar{I}$  be a point in a  $\bar{\vartheta}_E$ -path that contains a first computation point of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ . If  $\bar{I}$  is the first point in the path at which  $\mathcal{B}(F)$  or  $\mathcal{B}(L)$  holds, then  $E^{\text{in}} = e_s^i$  holds at  $\bar{I}$  and all points that precede  $\bar{I}$  but not at any point that succeeds  $\bar{I}$  in the path.
3. If  $\bar{I}_1, \bar{I}_2, \dots, \bar{I}_p$  is a subpath of a  $\bar{\vartheta}_E$ -path such that (a)  $\mathcal{B}(F)$  holds at  $\bar{I}$  but not at the remaining points of the path and (b)  $\bar{I}_p$  is the first point of the path at which  $\mathcal{B}(L)$  holds, then  $E^{\text{in}} = e_{\Gamma_E}$  holds at all points  $\bar{I}_i$  of the subpath satisfying  $i \neq 1 \wedge (i-1) \bmod \Gamma_E = 0$  but not at the remaining points of the path.



**PROGRAM:** *HostProg*

$(\forall t, i : t_{\text{fst}} \leq t \leq t_{\text{lst}} \wedge 0 < i \leq f : f_{i,1}(t) \in \text{sig}(F_{i,1}))$

$(\forall t, i : t_{\text{fst}} \leq t \leq t_{\text{lst}} \wedge 0 < i \leq f : f_{i,2}(t) \in \text{sig}(F_{i,2}))$

$(\forall t, i : t_{\text{fst}} \leq t \leq t_{\text{lst}} \wedge 0 < i \leq \ell : \ell_i(t) \in \text{sig}(L_i))$

**for**  $t$  **from**  $t_{\text{fst}}$  **to**  $t_{\text{lst}}$  **do**

**for**  $i$  **from** 1 **to**  $f$  **do**

        inject( $f_{i,1}(t)$ , pi( $F_{i,1}^{\text{in}}$ ))

        inject( $f_{i,2}(t)$ , pi( $F_{i,2}^{\text{in}}$ ))

**for**  $i$  **from** 1 **to**  $\ell$  **do**

        inject( $\ell_i(t)$ , pi( $L_i^{\text{in}}$ ))

**if**  $t \notin \{\text{input}(E(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)\}$  **→** inject( $e_u$ ,  $E^{\text{in}}$ )

**□** **else** **→** **for**  $i$  **from** 1 **to**  $f$  **do**

**if**  $t \in \{\text{input}(E(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)_i\}$  **→** inject( $e_s^i$ , pi( $E^{\text{in}}$ ))

**□** **else** **→** skip

**fi**

**fi**

**Table 5–1:** The host program for three-dimensional UREs (to be refined). The specification of the initialisation and termination control variables are refined in Sects. 5.3.2 and 5.3.3.

4. If  $\bar{I}$  is the first point in a  $\bar{\vartheta}_E$ -path at which  $\mathcal{B}(L)$  holds and  $\mathcal{B}(F)$  does not hold at any point that succeeds  $\bar{I}$  of the path, then  $E^{\text{in}} = e_d$  does not hold at  $\bar{I}$  nor at any point that precedes  $\bar{I}$  but holds at all points that succeed  $\bar{I}$  in the path.

**Proof** Follows from the host and cell programs (Fig. 5–3). □

**Theorem 5.2** Let  $\mathcal{B}(F)$  and  $\mathcal{B}(L)$  be characteristic functions of first and last computation points, respectively. If the space-time mapping is valid for the evolution control variable, the PCUREs are correct.

**Proof** Proving the correctness of the PCUREs means proving that  $\chi_{\bar{\Phi}}$ , as defined in the cell program, is a characteristic function of  $\bar{\Phi}$ . By the hypothesis,  $\mathcal{B}(F)$



**PROGRAM:** *CellProg*

$$\begin{aligned}
(\forall i : 0 < i \leq \Gamma_E : e_i &= i) \\
(\forall i : 0 < i \leq f : F_{i,1}^{\text{out}} &= F_{i,1}^{\text{in}}) \\
(\forall i : 0 < i \leq f : F_{i,2}^{\text{out}} &= F_{i,2}^{\text{in}}) \\
(\forall i : 0 < i \leq \ell : L_i^{\text{out}} &= L_i^{\text{in}}) \\
E^{\text{out}} &= \text{if } \mathcal{B}(F) \wedge \neg \mathcal{B}(L) \rightarrow e_1 \\
&\quad \square \mathcal{B}(L) \rightarrow e_d \\
&\quad \square ((\exists i : 0 < i \leq f : E^{\text{in}} = e_s^i) \wedge \neg \mathcal{B}(F)) \vee \\
&\quad \quad E^{\text{in}} = e_d \vee E^{\text{in}} = e_u \rightarrow E^{\text{in}} \\
&\quad \square \text{else} \rightarrow (E^{\text{in}} \bmod \Gamma_E) + 1 \\
&\quad \text{fi} \\
\chi_{\overline{\Phi}} &= \text{if } \mathcal{B}(F) \vee E^{\text{in}} = e_{\Gamma_E} \rightarrow \chi_c \square \text{else} \rightarrow \chi_p \text{ fi}
\end{aligned}$$
**Table 5-2:** The cell program for three-dimensional UREs.

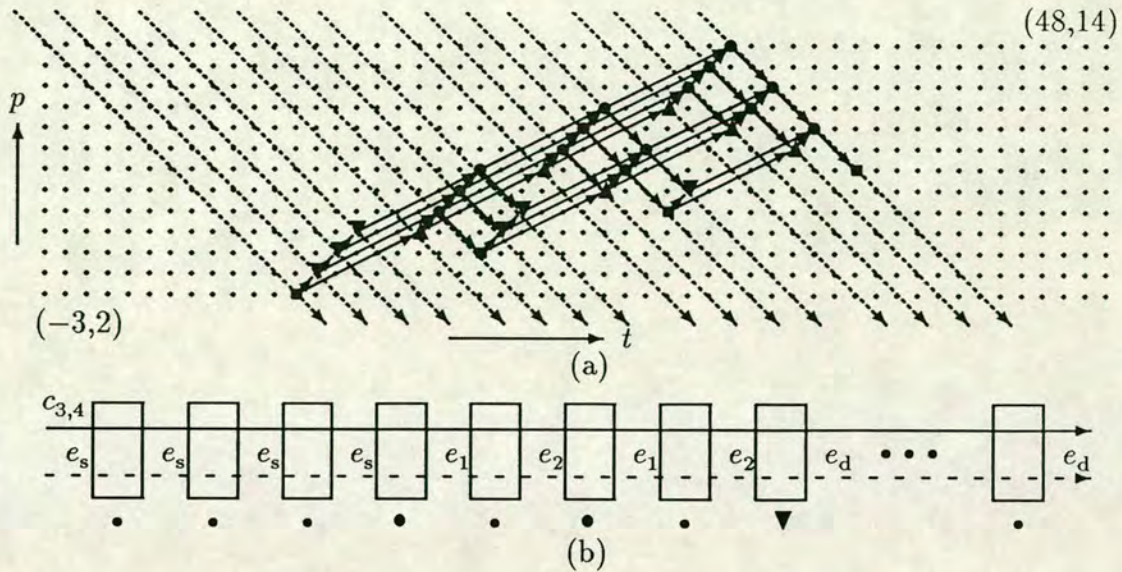
is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$ . It suffices to prove that  $E^{\text{in}} = e_{\Gamma_E}$  is a characteristic function of the set of the computation points that are not first computation points of  $E$ . This follows from (3) of Lemma 5.1, since  $\mathcal{B}(F)$  ( $\mathcal{B}(L)$ ) is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$ , ( $\overline{\text{lst}(\Phi, \vartheta_E)}$ ). (The validity of the space-time mapping for  $E$  ensures that Lemma 5.1 applies.)  $\square$

**Remark** If  $|\mathbb{F}| = 1$ , there are two initialisation control variables  $F_{1,1}$  and  $F_{1,2}$  and one control value  $e_s^1$  for the soaking points of  $E$ . The index  $i = 1$  will be omitted. Similarly, we write  $L$  instead of  $L_1$  for the termination control variable when  $|\mathbb{L}| = 1$ .  $\square$

**Example 5.2**  $4 \times 4$  LU-Decomposition

Without loss of generality, we choose  $\vartheta_E = \vartheta_C$ .  $\text{fst}(\Phi, \vartheta_E)$  by itself is a facet. We choose  $\mathbb{F} = \{\text{fst}(\Phi, \vartheta_E)\}$ . Because of the choice of  $\vartheta_E$ , the space-time mapping (5.9) is valid for evolution control variable  $E$  (Fig. 5-4(a)). The  $\vartheta_C$ -path depicted by the dashed line in the data dependence graph in Fig. 3-6 is mapped to the  $\overline{\vartheta}_C$ -path depicted by the dotted line in Fig. 5-1. This path is also depicted at the bottom of the array shown in Fig. 5-4(b). Let us consider the evaluation of the





**Figure 5-4:** The evolution control variable for LU-decomposition ( $m = 4$ ). (a)  $e_s$  is input at the  $\bar{\vartheta}_E$ -paths depicted by the dotted lines and  $e_u$  at the remaining  $\bar{\vartheta}_E$ -paths. (b) The evolution of the evolution control variable along the  $\bar{\vartheta}_E$ -path depicted by the dotted line in Fig. 5-1.

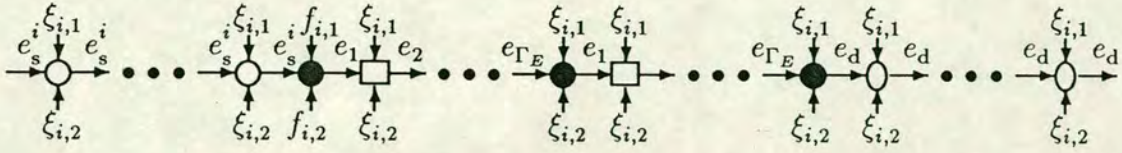
points in this path. Element  $c_{3,4}$  is input at this path. By the specification of the evolution control variable and the fact that  $\vartheta_E = \vartheta_C$ , we input  $e_s$  at the time step at which  $c_{3,4}$  is input.  $e_s$  and  $c_{3,4}$  travel with the same velocity.  $e_s$  changes to  $e_1$  at the first computation point. Then the control value changes alternatively from  $e_1$  to  $e_2$  ( $\text{sig}(E) = \{e_u, e_s, e_1, e_2, e_d\}$  since  $\Gamma_E = 2$ ).  $e_2$  changes to  $e_d$  at the last computation point; it is propagated to the output cell. The same reasoning applies for the other  $\bar{\vartheta}_E$ -paths. The correctness of the PCUREs is straightforward and is independent of the choice of  $\vartheta_E$ .  $\square$

Once evolution control variable  $E$  is constructed, or more precisely, once control dependence vector  $\vartheta_E$  is chosen, the specification of initialisation and termination control variables is completely determined. The following two sections present their respective constructions.

### 5.3.2 The Initialisation Control Flow

Initialisation control variables are pipelining control variables. Their specification must enable us to define  $\mathcal{B}(F)$  as a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$ . We shall





**Figure 5-5:** The interplay between the evolution control variable and the initialisation control variables. The arrows pointing down depict  $\vartheta_{F_{i,1}}$ . The arrows pointing up depict  $\vartheta_{F_{i,2}}$ .  $\xi_{i,1}$  and  $\xi_{i,2}$  satisfy  $\neg(\xi_{i,1} = f_{i,1} \wedge \xi_{i,2} = f_{i,2})$ .

construct  $F_{i,1}$  and  $F_{i,2}$  such that  $B(F_i)$ , defined below, is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ :

$$B(F_i) = E^{\text{in}} = e_s^i \wedge F_{i,1}^{\text{in}} = f_{i,1} \wedge F_{i,2}^{\text{in}} = f_{i,2} \quad (5.17)$$

Note that evolution control variable  $E$  also acts as an initialisation control variable. This is because a space-time point is a first computation point only if  $(\exists i : 0 < i \leq f : E^{\text{in}} = e_s^i)$  ((1) and (2) of Lemma 5.1). Clearly, if  $B(F_i)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ , for every  $i$ ,  $B(F)$  defined below is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$ :

$$B(F) = (\exists i : 0 < i \leq f : B(F_i)) \quad (5.18)$$

Let  $\text{soak}(E, i)$  be the set of soaking points in the  $\overline{\vartheta_E}$ -paths that contain first computation points of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ .  $\{\text{soak}(E, i) \mid 0 < i \leq f\}$  is a partition of  $\Upsilon_E^s$ . The next lemma is the basis of the construction of the initialisation control variables.

**Lemma 5.2** *Let the space-time mapping be valid for the evolution control variable.  $B(F_i)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$  iff  $F_{i,1}^{\text{in}} = f_{i,1} \wedge F_{i,2}^{\text{in}} = f_{i,2}$  holds at  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ ; but not at the soaking points of  $\text{soak}(E, i)$ .*

**Proof** The validity of the space-time mapping for evolution control variable  $E$  ensures that Lemma 5.1 applies.

“ $\implies$ ”: If  $F_{i,1}^{\text{in}} = f_{i,1} \wedge F_{i,2}^{\text{in}} = f_{i,2}$  does not hold at  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ , then  $B(F_i)$  is not a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ . If it holds at some soaking points of  $\text{soak}(E, i)$ , then  $B(F_i)$  must hold at some soaking points of  $\text{soak}(E, i)$ , by (2) of Lemma 5.1. Thus,  $B(F_i)$  is not a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ .



```

for  $t$  from  $t_{\text{fst}}$  to  $t_{\text{lst}}$  do
  for  $i$  from 1 to  $f$  do
    if  $t \in \{\text{input}(F_{i,1}(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)_i\}$  → inject( $f_{i,1}$ ,  $\text{pi}(F_{i,1}^{\text{in}})$ )
    [] else → inject( $\bar{f}_{i,1}$ ,  $\text{pi}(F_{i,1}^{\text{in}})$ )
    fi
    if  $t \in \{\text{input}(F_{i,2}(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)_i^{\text{in}}\}$  → inject( $f_{i,2}$ ,  $\text{pi}(F_{i,2}^{\text{in}})$ )
    [] else → inject( $\bar{f}_{i,2}$ ,  $\text{pi}(F_{i,2}^{\text{in}})$ )
    fi

```

**Table 5–3:** The specification of initialisation control variables.

“ $\Leftarrow$ ”: (1) of Lemma 5.1 asserts that  $\mathcal{B}(F_i)$  does not hold at the undefined points of  $E$ . A further application of (2) of Lemma 5.1 completes the proof.  $\square$

Let us describe the specification of initialisation control variables  $F_{i,1}$  and  $F_{i,2}$ . The communication constraint domains of  $F_{i,1}$  and  $F_{i,2}$  are given by

$$(\forall i : 0 < i \leq f : \Omega_{F_{i,1}} = \Omega_{F_{i,2}} = \emptyset) \quad (5.19)$$

Thus, the communication constraint is disregarded in the definition of the validity of the space-time mapping for the initialisation control variables. This is attributed to the construction of the initialisation control variables described below.

To make  $\mathcal{B}(F_i)$  a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ , we must establish  $\mathcal{B}(F_i)$  at  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ . To establish  $\mathcal{B}(F_i)$  at  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ , we must ensure that  $F_{i,1}^{\text{in}} = f_{i,1} \wedge F_{i,2}^{\text{in}} = f_{i,2}$  holds at  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$  (Lemma 5.2). Then, we must inject  $f_{i,k}$  ( $k=0,1$ ) at steps  $\text{fst}(F_{i,k}(I))$  for  $I \in \text{fst}(\Phi, \vartheta_E)_i$  and  $\bar{f}_{i,k}$  at the remaining steps (Fig. 5–5). For the refined specification of the initialisation control variables in the host program, see Tab. 5–3. However, this specification is insufficient to ensure that  $\mathcal{B}(F_i)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$  by Lemma 5.2, since  $F_{i,1}^{\text{in}} = f_{i,1} \wedge F_{i,2}^{\text{in}} = f_{i,2}$  and, consequently,  $\mathcal{B}(F_i)$  may also hold at some soaking points of  $\text{soak}(E, i)$  due to the non-injectivity of the space-time mapping. This happens, for example, if  $f_{i,1}$  which is injected at step  $\text{input}(F_{i,1}(I))$  meets  $f_{i,2}$  which is injected at step  $\text{input}(F_{i,2}(J))$  for different  $I, J \in \text{fst}(\Phi, \vartheta_E)_i$ . We avoid this situation by means of an appropriate choice of control dependence vectors for initialisation control variables.



We choose two linearly independent control dependence vectors  $\vartheta_{F_{i,1}}$  and  $\vartheta_{F_{i,2}}$  from  $\text{lin}(\text{fst}(\Phi, \vartheta_E)_i)$ . We write  $\mathfrak{F}_i$  for the intersection of  $\text{rays}(\text{fst}(\Phi, \vartheta_E)_i, \vartheta_{F_{i,1}})$  and  $\text{rays}(\text{fst}(\Phi, \vartheta_E)_i, \vartheta_{F_{i,2}})$ . We write  $\mathfrak{F}$  for the union of  $\mathfrak{F}_1, \mathfrak{F}_2, \dots, \mathfrak{F}_f$ .

$$\begin{aligned}\mathfrak{F}_i &= \text{rays}(\text{fst}(\Phi, \vartheta_E)_i, \vartheta_{F_{i,1}}) \cap \text{rays}(\text{fst}(\Phi, \vartheta_E)_i, \vartheta_{F_{i,2}}) \\ \mathfrak{F} &= (\cup i : 0 < i \leq f : \mathfrak{F}_i)\end{aligned}\tag{5.20}$$

It is easy to see that  $\mathfrak{F}_i \supseteq \text{fst}(\Phi, \vartheta_E)_i$  and thus  $\mathfrak{F} \supseteq \text{fst}(\Phi, \vartheta_E)$ . This choice of control dependence vectors ensures that  $\mathcal{B}(F_i)$  only holds at  $\overline{\mathfrak{F}_i}$  (by Lemma 5.3, which will be presented shortly). Still,  $\mathcal{B}(F_i)$  is not a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ ; if  $\overline{\mathfrak{F}_i} \setminus \overline{\text{fst}(\Phi, \vartheta_E)_i}$  contains soaking points of  $\text{soak}(E, i)$ . In order for  $\overline{\mathfrak{F}_i} \setminus \overline{\text{fst}(\Phi, \vartheta_E)_i}$  and  $\text{soak}(E, i)$  to be disjoint, we simply redefine (enlarge) the communication constraint domain of evolution control variable  $E$  to:

$$\Omega_E = \text{planes}(\Phi, \sigma) \cap \text{rays}(\mathfrak{F}, \pm \vartheta_E)\tag{5.21}$$

(Compare (5.16).) Under a valid space-time mapping for evolution control variable  $E$ , all  $\vartheta_E$ -paths of  $\text{paths}(\mathfrak{F}_i, \vartheta_E)$  are mapped to different  $\overline{\vartheta}_E$ -paths. Thus, no two points of  $\mathfrak{F}_i$  can have images in the same  $\overline{\vartheta}_E$ -path. That is, no point of  $\mathfrak{F}_i \setminus \text{fst}(\Phi, \vartheta_E)_i$  is mapped to a  $\overline{\vartheta}_E$ -path that contains a first computation point of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ , i.e.,  $(\overline{\mathfrak{F}_i} \setminus \overline{\text{fst}(\Phi, \vartheta_E)_i}) \cap \text{soak}(E, i) = \emptyset$ .

The advantage of choosing control dependence vectors this way is that we know precisely that  $\mathcal{B}(F_i)$  can hold only at the image of  $\mathfrak{F}_i$ .  $\mathfrak{F}_i$  depends on the choice of  $\vartheta_{F_{i,1}}$  and  $\vartheta_{F_{i,2}}$  but is independent of the space-time mapping. This enables us to enforce the disjointness of  $\overline{\mathfrak{F}_i} \setminus \overline{\text{fst}(\Phi, \vartheta_E)_i}$  and  $\text{soak}(E, i)$  by an appropriate choice of the communication constraint domain for the evolution control variable.

**Remark** By Lemma 5.2,  $\xi_{i,1}$  and  $\xi_{i,2}$ , as shown in Fig. 5-5 at all points succeeding the first computation point can take any value. However, this does not seem to be practically useful, because the verification of a space-time mapping that satisfies this constraint depends on the space-time mapping.  $\square$

**Theorem 5.3** *If the space-time mapping is valid for the evolution and initialisation control variable,  $\mathcal{B}(F)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$ .*



The proof of this theorem requires the following lemma, which states that two paths of variables moving with different velocities only share an image at the intersection point. This lemma is the reason why we use two initialisation control variables for every  $\text{fst}(\Phi, \vartheta_E)_i$  and why the communication constraint for the initialisation control variables can be disregarded.

**Lemma 5.3** *Assume  $0 \neq \text{flow}(V) \neq \text{flow}(W) \neq 0$ .*

$$(\forall I, J, K : I, J, K \in \mathbb{Z}^n \wedge I \in \text{ray}(K, \pm \vartheta_V) \wedge J \in \text{ray}(K, \pm \vartheta_W) : \bar{I} = \bar{J} \implies I = J = K)$$

**Proof**  $\bar{I} = \bar{J} \iff \text{step}(I) = \text{step}(J) \wedge \text{place}(I) = \text{place}(J) \iff \lambda I = \lambda J \wedge \sigma I = \sigma J$ .  
 $I \in \text{ray}(K, \pm \vartheta_V) \iff (\exists m : m \in \mathbb{Z} : K = m\vartheta_V + I)$  and  $J \in \text{ray}(K, \pm \vartheta_W) \iff$   
 $(\exists n : n \in \mathbb{Z} : K = n\vartheta_W + J)$ . A simple algebraic calculation establishes

$$m\lambda\vartheta_V = n\lambda\vartheta_W \tag{5.22}$$

$$m\sigma\vartheta_V = n\sigma\vartheta_W \tag{5.23}$$

We consider all four possible geometric relationships between  $I$  and  $J$  relative to  $K$  (the proof proceeds to show that only  $I = J = K$  can hold).

- $I = K \wedge J \neq K$ . This implies  $m = 0$  and consequently  $\lambda\vartheta_W = 0$ , contradicting the hypothesis  $\text{flow}(W) \neq 0$ .
- $I \neq K \wedge J = K$ . This implies  $n = 0$  and consequently  $\lambda\vartheta_V = 0$ , contradicting the hypothesis  $\text{flow}(V) \neq 0$ .
- $I \neq K \wedge J \neq K$ . By hypothesis we infer  $m\sigma\vartheta_V \neq 0 \neq n\sigma\vartheta_W$ . Dividing (5.23) by (5.22) yields  $\sigma\vartheta_V/\lambda\vartheta_V = \sigma\vartheta_W/\lambda\vartheta_W$ , i.e.,  $\text{flow}(V) = \text{flow}(W)$ , contradicting the hypothesis  $\text{flow}(V) \neq \text{flow}(W)$ .
- $I = K \wedge J = K$ . Trivially true. □

This proof does not require the containment of  $I$ ,  $J$  and  $K$  in the index space; it only relies on the geometric relationship of the three points. It is easy to see that, if  $\vartheta_V$  and  $\vartheta_W$  are not co-linear for two moving variables  $V$  and  $W$ , then they



must travel with distinct velocities, i.e.,  $\text{flow}(V) \neq \text{flow}(W)$ . This is why we require  $\vartheta_{F_{i,1}}$  and  $\vartheta_{F_{i,2}}$  to be linearly independent.

**Proof of Thm. 5.3** By the construction of the initialisation control variables,  $\vartheta_{F_{i,1}}$  and  $\vartheta_{F_{i,2}}$  are linearly independent vectors from the linear space  $\text{lin}(\text{fst}(\Phi, \vartheta_E)_i)$ . Thus, for every  $\vartheta_{F_{i,1}}$ -path that passes through some point of  $\text{fst}(\Phi, \vartheta_E)_i$ , there must exist a  $\vartheta_{F_{i,2}}$ -path that passes through some point of  $\text{fst}(\Phi, \vartheta_E)_i$  such that the two paths intersect, and vice versa. By definition,  $\mathfrak{F}_i$  denotes the set of these intersection points. By Lemma 5.3, the images of all these paths intersect only at the points of  $\overline{\mathfrak{F}}_i$  since the space-time mapping is valid for the initialisation control variables and by our assumption that all propagation control variables are moving. If the space-time mapping is valid for evolution control variable  $E$ , no point contained in  $\overline{\mathfrak{F}}_i \setminus \overline{\text{fst}(\Phi, \vartheta_E)_i}$  is in the same  $\overline{\vartheta_E}$ -path as a point contained in  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ . This means that  $(\overline{\mathfrak{F}}_i \setminus \overline{\text{fst}(\Phi, \vartheta_E)_i}) \cap \text{soak}(E, i) = \emptyset$ . That is,  $F_{i,1}^{\text{in}} = f_{i,1} \wedge F_{i,2}^{\text{in}} = f_{i,2}$  does not hold at  $\text{soak}(E, i)$ . But it does hold at  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$  by the construction of the initialisation control variables. An application of Lemma 5.2 establishes that  $\mathcal{B}(F_i)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$ . By the definition of  $\mathcal{B}(F)$ ,  $\mathcal{B}(F)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$ .  $\square$

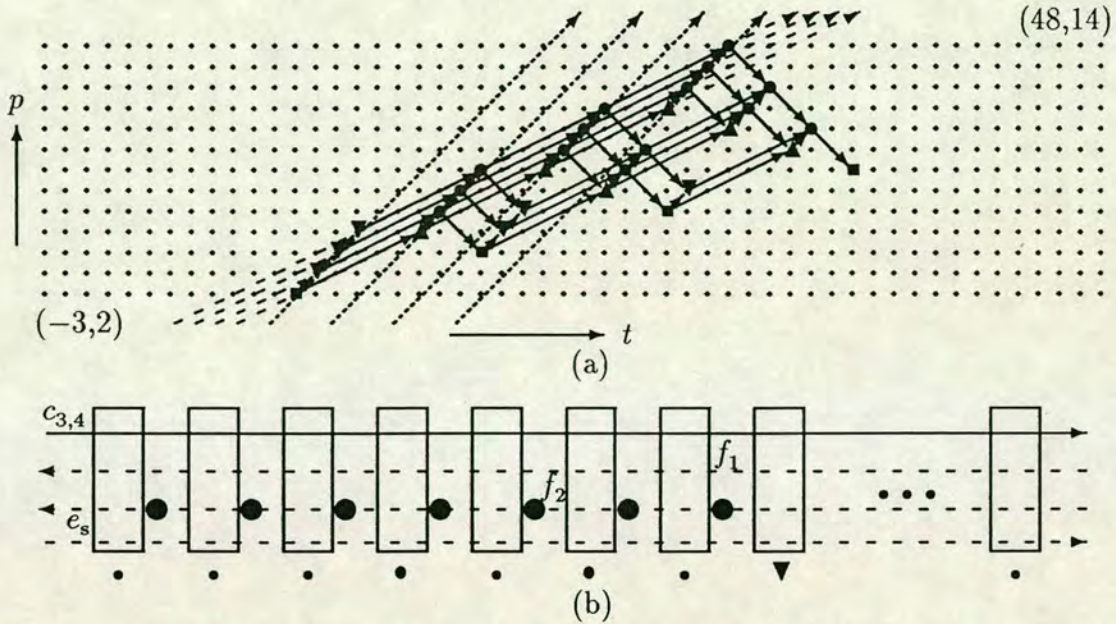
### Example 5.3 $4 \times 4$ LU-Decomposition

Let us continue to choose  $\vartheta_E = \vartheta_C$  and  $\textcircled{F} = \{\text{fst}(\Phi, \vartheta_E)\}$ . We need two initialisation control variables  $F_1$  and  $F_2$ . We choose  $\vartheta_{F_1} = \vartheta_A$  and  $\vartheta_{F_2} = \vartheta_B$ . This gives rise to  $\mathfrak{F} = \text{fst}(\Phi, \vartheta_E)$  (Fig. 5-6).  $\square$

### 5.3.3 The Termination Control Flow

Termination control variables are pipelining control variables. Their specification must enable us to define  $\mathcal{B}(L)$  as a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)}$ . However, the specification of the termination control variables can be greatly simplified for two reasons. First, we can construct them incrementally based on the specification of evolution and initialisation control variables. A space-time point is a last computation point only if it is a computation point, i.e., only when  $\mathcal{B}(F) \vee E^{\text{out}} = e_{\Gamma_E}$





**Figure 5-6:** The initialisation control variables for LU-decomposition ( $m = 4$ ). (a)  $f_1$  ( $f_2$ ) is input at the  $\bar{\vartheta}_{F_1}$ -paths ( $\bar{\vartheta}_{F_2}$ -paths) depicted by the dotted (dashed) lines and  $\bar{f}_1$  ( $\bar{f}_2$ ) at the remaining  $\bar{\vartheta}_{F_1}$ -paths ( $\bar{\vartheta}_{F_2}$ -paths). (b) The interplay between evolution control variable  $E$  and initialisation control variables  $F_1$  and  $F_2$  in the  $\vartheta_E$ -path depicted in Fig. 5-1. Control signals  $e_s$ ,  $f_1$  and  $f_2$  will meet at the fourth cell from the left three time steps after  $c_{3,4}$  is input. At this step,  $B(F)$  holds at that cell, indicating that the point to be computed by the cell is a first computation point.

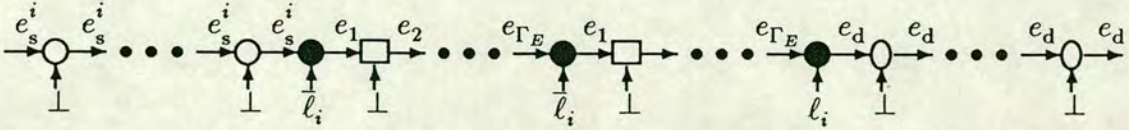
holds at the point. Second, we have dealt with the non-injectivity of the space-time mapping during the construction of initialisation control variables (Thm. 5.3). This enables us to associate one rather than two control variables with every  $\text{lst}(\Phi, \vartheta_E)_i$ . Besides,  $\mathbb{D}$  needs not be a partition of  $\text{lst}(\Phi, \vartheta_E)$ . ( $\mathbb{F}$  must be a partition of  $\text{fst}(\Phi, \vartheta_E)$ . If two  $\mathbb{F}$ -blocks  $\text{fst}(\Phi, \vartheta_E)_i$  and  $\text{fst}(\Phi, \vartheta_E)_j$  intersected, the construction of  $E$  would require the simultaneous injection of  $e_s^i$  and  $e_s^j$  for the  $\bar{\vartheta}_E$ -paths that intersect both  $\overline{\text{fst}(\Phi, \vartheta_E)_i}$  and  $\overline{\text{fst}(\Phi, \vartheta_E)_j}$ .)

We shall construct  $L_i$  such that  $B(L_i)$ , defined below, is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)_i}$ :

$$B(L_i) = (B(F) \vee E^{\text{in}} = e_{\Gamma_E}) \wedge L_i^{\text{in}} = \ell_i$$

Note that the evolution and initialisation control variables also act as termination





**Figure 5-7:** The interplay between the evolution control variable and the termination control variables. The arrows pointing up depict  $\vartheta_{L_i}$ .  $\perp$  denotes any value in  $\text{sig}(L_i)$ .

control variables. If  $\text{lst}(\Phi, \vartheta_E)_i$  does not contain first computation points of  $E$ , i.e., if  $\text{fst}(\Phi, \vartheta_E)$  and  $\text{lst}(\Phi, \vartheta_E)_i$  are disjoint,  $\mathcal{B}(F) \vee E^{\text{in}} = e_{\Gamma_E}$  can be simplified to  $E^{\text{in}} = e_{\Gamma_E}$ . Clearly, if  $\mathcal{B}(L_i)$  is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)_i}$  for every  $i$ ,  $\mathcal{B}(L)$ , defined below, is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)}$ :

$$\mathcal{B}(L) = (\exists i : 0 < i \leq \ell : \mathcal{B}(L_i))$$

The next lemma is the basis of the construction of the termination control variables.

**Lemma 5.4** *Let the space-time mapping be valid for the evolution and initialisation control variables.  $\mathcal{B}(L_i)$  is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)_i}$  iff  $L_i^{\text{in}} = \ell_i$  holds at  $\overline{\text{lst}(\Phi, \vartheta_E)_i}$  and  $L_i^{\text{in}} = \bar{\ell}_i$  holds at  $\overline{\Phi \setminus \text{lst}(\Phi, \vartheta_E)_i}$ .*

**Proof** Sufficiency follows from the definition of a characteristic function (Sect. 5.2). Let us prove necessity. Since the space-time mapping is valid for the evolution and initialisation control variables,  $\mathcal{B}(F)$  is a characteristic function of  $\overline{\text{fst}(\Phi, \vartheta_E)}$  by Thm. 5.3. An application of (3) and (4) of Lemma 5.1 completes the proof.  $\square$

Let us describe the construction of the termination control variables. The communication constraint domain of  $L_i$  is given by

$$(\forall i : 0 < i \leq \ell : \Omega_{L_i} = \text{planes}(\Phi, \sigma) \cap \text{rays}(\Phi, \pm \vartheta_{L_i})) \quad (5.24)$$

This is to ensure that all  $\vartheta_{L_i}$ -paths that contain computation points are mapped to different  $\bar{\vartheta}_{L_i}$ -paths. Following the same reasoning as in the specification of initialisation control variables, we inject  $\ell_i$  at time steps  $\text{input}(L_i(I))$  for  $I \in \text{lst}(\Phi, \vartheta_E)_i$ , and  $\bar{\ell}_i$  at the remaining steps (Fig. 5-7). For the final version of the host program for three-dimensional UREs, see Tab. 5-4.



```

PROGRAM: HostProg
for  $t$  from  $t_{\text{fst}}$  to  $t_{\text{lst}}$  do
  for  $i$  from 1 to  $f$  do
    if  $t \in \{\text{input}(F_{i,1}(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)_i\}$   $\rightarrow$  inject( $f_{i,1}$ ,  $\text{pi}(F_{i,1}^{\text{in}})$ )
    [] else  $\rightarrow$  inject( $\bar{f}_{i,1}$ ,  $\text{pi}(F_{i,1}^{\text{in}})$ )
    fi
    if  $t \in \{\text{input}(F_{i,2}(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)_i\}$   $\rightarrow$  inject( $f_{i,2}$ ,  $\text{pi}(F_{i,2}^{\text{in}})$ )
    [] else  $\rightarrow$  inject( $\bar{f}_{i,2}$ ,  $\text{pi}(F_{i,2}^{\text{in}})$ )
    fi
  for  $i$  from 1 to  $\ell$  do
    if  $t \in \{\text{input}(L_i(I)) \mid I \in \text{lst}(\Phi, \vartheta_E)_i\}$   $\rightarrow$  inject( $\ell_i$ ,  $\text{pi}(L_i^{\text{in}})$ )
    [] else  $\rightarrow$  inject( $\bar{\ell}_i$ ,  $\text{pi}(L_i^{\text{in}})$ )
    fi
  if  $t \notin \{\text{input}(E(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)\}$   $\rightarrow$  inject( $e_u$ ,  $E^{\text{in}}$ )
  [] else  $\rightarrow$  for  $i$  from 1 to  $f$  do
    if  $t \in \{\text{input}(E(I)) \mid I \in \text{fst}(\Phi, \vartheta_E)_i\}$   $\rightarrow$  inject( $e_s^i$ ,  $\text{pi}(E^{\text{in}})$ )
    [] else  $\rightarrow$  skip
    fi
  fi
fi

```

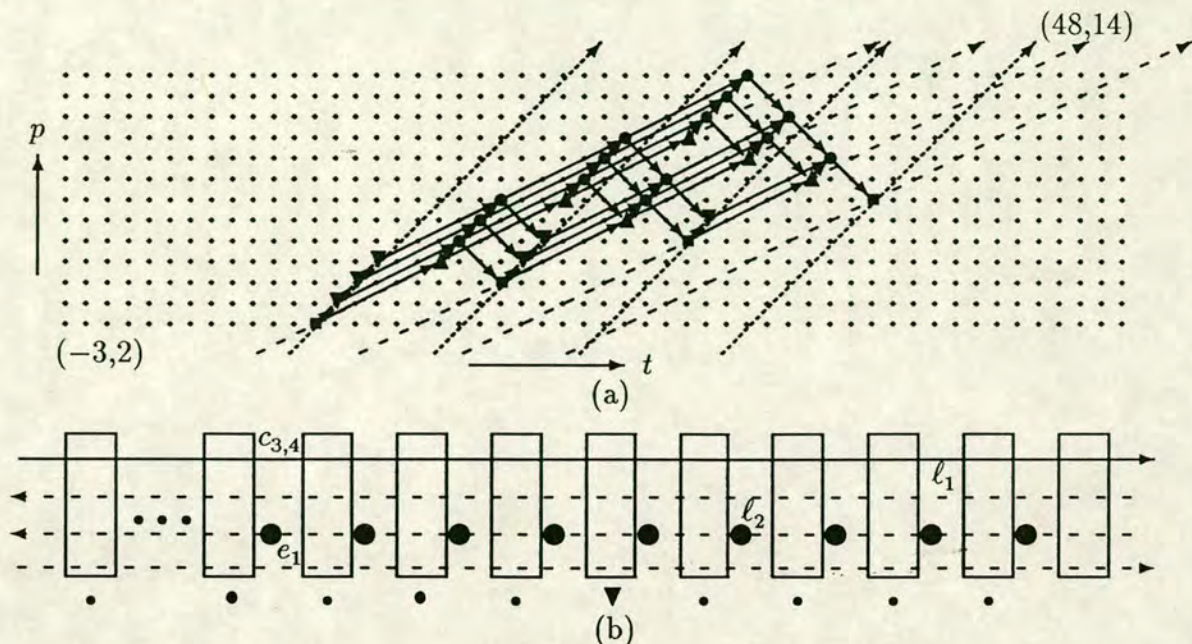
**Table 5–4:** The host program for three-dimensional UREs (the final version).

Similarly to the choice of control dependence vectors for initialisation control variables, we choose  $\vartheta_{L_i}$  from the linear space  $\text{lin}(\text{lst}(\Phi, \vartheta_E)_i)$ . We introduce domain  $\mathfrak{L}_i$  analogously to  $\mathfrak{F}_i$  of (5.20). Since each  $\text{lst}(\Phi, \vartheta_E)_i$  is associated with one termination control variable,  $\mathfrak{L}_i$  is defined as the intersection of the index space and the  $\vartheta_{L_i}$ -paths that pass through some point of  $\text{lst}(\Phi, \vartheta_E)_i$ . We write  $\mathfrak{L}$  for the union of  $\mathfrak{L}_1, \mathfrak{L}_2, \dots, \mathfrak{L}_\ell$ .

$$\begin{aligned}
 \mathfrak{L}_i &= \Phi \cap \text{rays}(\text{lst}(\Phi, \vartheta_E)_i, \pm \vartheta_{L_i}) \\
 \mathfrak{L} &= (\cup i : 0 < i \leq \ell : \mathfrak{L}_i)
 \end{aligned}
 \tag{5.25}$$

**Theorem 5.4** *If the space-time mapping is valid for all propagation control variables, then  $B(L)$  is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)}$ .*





**Figure 5-8:** The termination control variables for LU-decomposition ( $m = 4$ ). (a)  $\ell_1$  ( $\ell_2$ ) is input at the  $\bar{\vartheta}_{L_1}$ -paths ( $\bar{\vartheta}_{L_2}$ -paths) depicted by the dotted (dashed) lines and  $\bar{\ell}_1$  ( $\bar{\ell}_2$ ) at the remaining  $\bar{\vartheta}_{L_1}$ -paths ( $\bar{\vartheta}_{L_2}$ -paths). (b) The interplay between evolution control variable  $E$  and termination control variables  $L_1$  and  $L_2$  in the  $\vartheta_E$ -path depicted in Fig. 5-1.  $\ell_1$  and  $\ell_2$  and the control value of  $E$  will meet at the sixth cell from the right in three time steps. Then,  $B(L)$  evaluates to true, indicating that the point to be computed by the cell is a last computation point.

**Proof** Under a valid space-time mapping for the termination control variables,  $L_i^{\text{in}} = \ell_i$  holds at  $\overline{\text{lst}(\Phi, \vartheta_E)_i}$  and  $L_i^{\text{in}} = \bar{\ell}_i$  holds at  $\overline{\Phi \setminus \text{lst}(\Phi, \vartheta_E)_i}$ . By Lemma 5.4,  $B(L_i)$  is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)_i}$ . By the definition of  $B(L)$ ,  $B(L)$  is a characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)}$ .  $\square$

**Remark** It is not difficult to see that  $\mathfrak{L} = \text{lst}(\Phi, \vartheta_E)$  is necessary and sufficient for  $B(L)$  to be characteristic function of  $\overline{\text{lst}(\Phi, \vartheta_E)}$ . Decomposing  $\text{lst}(\Phi, \vartheta_E)$ , as described here, and choosing  $\vartheta_{L_i}$  from the linear space  $\text{lin}(\text{lst}(\Phi, \vartheta_E)_i)$  is unnecessary, but there is an advantage to doing so:  $\mathfrak{L} = \text{lst}(\Phi, \vartheta_E)$  is automatically satisfied. Otherwise, we may have to verify the equality of  $\mathfrak{L}$  and  $\text{lst}(\Phi, \vartheta_E)$  explicitly.  $\square$

**Example 5.4**  $4 \times 4$  LU-Decomposition

We continue to choose  $\vartheta_E = \vartheta_C$ ,  $\vartheta_{F_1} = \vartheta_A$  and  $\vartheta_{F_2} = \vartheta_B$ . We choose  $\mathbb{L} = \{\text{lst}(\Phi, \vartheta_E)_1, \text{lst}(\Phi, \vartheta_E)_2\}$ , where  $\text{lst}(\Phi, \vartheta_E)_1$  and  $\text{lst}(\Phi, \vartheta_E)_2$  are two facets of the



index space:

$$\begin{aligned}\text{lst}(\Phi, \vartheta_E)_1 &= \Phi \cap \{(i, j, k) \mid i - k = 0\} \\ \text{lst}(\Phi, \vartheta_E)_2 &= \Phi \cap \{(i, j, k) \mid j - k = 0\}\end{aligned}$$

We choose  $\vartheta_{L_1} = \vartheta_A \in \text{lin}(\text{fst}(\Phi, \vartheta_E)_1)$  and  $\vartheta_{L_2} = \vartheta_B \in \text{lin}(\text{fst}(\Phi, \vartheta_E)_2)$ . The space-time mapping (5.9) is valid for the termination control variables (Fig. 5-8).  $\square$

### 5.3.4 The Space-Time Mapping

Let us summarise the previous construction of the PCUREs in a procedure.

**Procedure 5.1** (Construction of the PCUREs for one-dimensional arrays from three-dimensional source UREs)

1. Choose control dependence vector  $\vartheta_E \in \mathbf{Z}^n$  for evolution control variable  $E$ .
2. Find a partition  $\mathbb{F} = \{\text{fst}(\Phi, \vartheta_E)_i \mid 0 < i \leq f\}$  of  $\text{fst}(\Phi, \vartheta_E)$  such that  $\dim(\text{fst}(\Phi, \vartheta_E)_i) \leq n - 1$ . Choose two linearly independent vectors  $\vartheta_{F_{i,1}}$  and  $\vartheta_{F_{i,2}}$  for initialisation control variables  $F_{i,1}$  and  $F_{i,2}$  from the linear space  $\text{lin}(\text{fst}(\Phi, \vartheta_E)_i)$  for every  $\text{fst}(\Phi, \vartheta_E)_i$ .
3. Find a set  $\mathbb{L} = \{\text{lst}(\Phi, \vartheta_E)_i \mid 0 < i \leq \ell\}$  such that  $\dim(\text{lst}(\Phi, \vartheta_E)_i) \leq n - 1$  and the union of its elements is  $\text{lst}(\Phi, \vartheta_E)$ . Choose control dependence vector  $\vartheta_{L_i}$  for termination control variable  $L_i$  from the linear space  $\text{lin}(\text{lst}(\Phi, \vartheta_E)_i)$  for every  $\text{lst}(\Phi, \vartheta_E)_i$ .
4. The cell program and host program are given in Tabs. 5-2 and 5-4.  $\square$

The communication constraint domains of the propagation control variables are given in (5.19), (5.21) and (5.24). Our objective for a different specification of PCUREs stated in Sect. 5.1 was to replace  $\Psi_V$  in the formulation of the communication constraint (5.3) by  $\Omega_V$ , which satisfies (5.7). This objective has been achieved.



The mapping conditions for the validity of space-time mappings are stated in Thm. 5.1. Remember that we require all propagation control variables to be moving (Sect. 5.1). We call this constraint the *velocity constraint*. For completeness, we state in a theorem the mapping conditions, including the velocity constraint, for the validity of space-time mappings. We rewrite the communication constraint for propagation control variables to dispense with the concept of the extended index space. We write  $\mathcal{V}_\ell$  for the set of termination control variables.

**Theorem 5.5** *A space-time mapping is valid for the pipelined UREs (in either the  $\pi$ -model or the  $\chi$ -model) iff*

- $(\forall V : V \in \mathcal{V} : \lambda \vartheta_V > 0)$  *(Precedence Constraint)*
- $\pi : \Phi \mapsto \mathbf{Z}^2$  *(Computation Constraint)*
- $(\forall V : V \in \mathcal{V} : \sigma \vartheta_V \mid \lambda \vartheta_V)$  *(Delay Constraint for the  $\chi$ -Model)*
- $(\forall V : V \in \mathcal{V}_d : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \mathbb{I}_V : \text{input} : S \mapsto \mathbf{Z}))$   
 $(\forall V : V \in \mathcal{V}_c : (\forall S : S \in \text{fst}(\Psi_V, \vartheta_V) / \mathbb{I}'_V : \text{input} : S \mapsto \mathbf{Z}))$   
 $(\forall S : S \in \text{fst}(\mathfrak{F}, \vartheta_E) / \mathbb{I}_E : \text{input} : S \mapsto \mathbf{Z})$   
 $(\forall V : V \in \mathcal{V}_\ell : (\forall S : S \in \text{fst}(\Phi, \vartheta_V) / \mathbb{I}_V : \text{input} : S \mapsto \mathbf{Z}))$   
*(Communication Constraint)*
- $(\forall V : V \in \mathcal{V}_p : \text{flow}(V) \neq 0)$  *(Velocity Constraint)*

**Proof** Thms. 5.1, 5.2, 5.3 and 5.4. □

Evolution control variable  $E$  may be stationary. Then, additional control is needed to direct its loading, access and recovery, just like for stationary data variables. In general, we should choose  $\vartheta_E$  and the space-time mapping such that  $E$  is moving. Initialisation and termination control variables must be moving for the following reasons. Their (control) values serve to change the control values of  $E$ ; they must be supplied from the outside environment such that they can arrive at the right cells at the right time. If they are stationary, we have to pre-load



their control signals – but we cannot, because the presence of these control signals ahead of time will validate some guards of  $E^{\text{out}}$  in the cell program at the wrong time (unless we add a second level of control variables, this one governing the access of stationary control variables, and so on).

Let us have a look at the requirements that the PCUREs impose on hardware. Recall our discussions about the injection of input data and the ejection of output data in the systolic array in Sect. 4.4. Some  $\bar{\vartheta}_V$ -paths may start at internal cells (see Figs. 5-4, 5-6 and 5-8). If the corresponding input values at these paths are required to be input, the input can be implemented by a system reset or by pipelining before the first step  $t_{\text{fst}}$ . In the construction of the PCUREs, the specification of the evolution control variable requires that all corresponding input channels of  $E$  be reset to  $e_u$ . Similarly, the specification of initialisation control variables requires that all corresponding input channels of  $F_{i,1}$  be reset to  $\bar{f}_{i,1}$  and of  $F_{i,2}$  to  $\bar{f}_{i,2}$ , and the specification of termination control variables requires that all corresponding input channels be reset to  $\bar{\ell}_i$ .

The bit width required to communicate values of  $V$  is  $\lceil \log_2 |\text{sig}(V)| \rceil$ . For example, evolution control variable  $E$  needs  $\lceil \log_2(\Gamma_E + |\textcircled{F}| + 2) \rceil$  bits –  $\Gamma_E - 1$  for relaying points, one for computation points,  $|\textcircled{F}|$  for soaking points, one for draining points and one for undefined points. We should choose  $\vartheta_E$  and the space-time mapping such that  $\Gamma_E$  becomes a constant in order to satisfy Rst. 2 of Sect. 3.2. The modulo operation in the cell program can be implemented by a circular shift (or rotate) operation if  $\Gamma_E$  is a power of 2. It is superfluous if  $\Gamma_E = 1$ .

## 5.4 The PCUREs for $n$ -Dimensional UREs

Both Thm. 5.2 and Thm. 5.4 work for  $n$ -dimensional UREs. But, unfortunately, Thm. 5.3 cannot be generalised to  $n$ -dimensional UREs. Recall that, in the specification of the initialisation control variables for three-dimensional UREs,  $\text{fst}(\Phi, \vartheta_E)$  is partitioned into a union of subsets each of which is contained in a hyperplane. We associate two initialisation control variables with each subset. The proof of



Thm. 5.3 relies on Lemma 5.3, which assumes that each subset has no more than two dimensions. This is because the two linearly independent vectors, which are associated with the two initialisation control variables of a subset forms a basis of two-dimensional space. We failed to find an analogue of Lemma 5.3 that applies for more than two control variables.

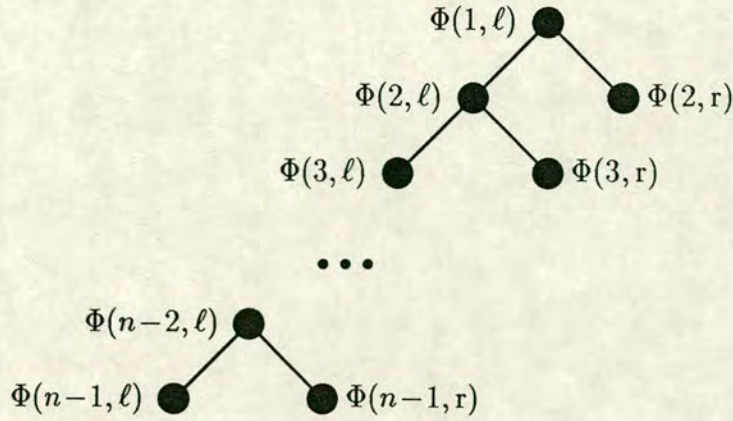
There is a brute-force application of Thm. 5.3 to  $n$ -dimensional UREs. All we need to do is to partition  $\text{fst}(\Phi, \vartheta_E)$  into a union of subsets each of which has no more than two dimensions. The number of initialisation control variables needed will be twice the number of these subsets. Clearly, this number will depend on the size of the problem. In the extreme, each subset contains only one first computation point. Such a simple-minded scheme is rather impractical.

We choose a construction scheme for initialisation control variables for  $n$ -dimensional UREs that is based on the observation that our previous construction of the PCUREs supports a hierarchical specification of control signals. To identify the first computation points in  $\text{fst}(\Phi, \vartheta_E)$ , we treat  $\text{fst}(\Phi, \vartheta_E)$  as a new index space and subject it to the previous construction of the PCUREs. We repeat this procedure until we obtain an index space of three dimensions. We then apply Thm. 5.3 to the three-dimensional index space. This decomposition produces a binary tree structure. The left (right) node at a level denotes the set of first (last) computation points of the set denoted by the father node. Since the initial index space is bounded and the new index space resulting from each decomposition always contains less points than the preceding one, this decomposition must eventually lead to a three-dimensional domain.

To speed up the convergence and thus reduce the number of control variables, we make the new index space at one level one dimension less than that at the level above, by an appropriate choice of control dependence vector for the evolution control variable at every level in the decomposition tree. This choice is always possible if we extend the initial index space appropriately prior to the decomposition, as will be discussed in Sect. 5.7.

Let us assume that the new index space at one level is one dimension less than the index space at the level above. Let us denote the evolution control variable





**Figure 5-9:** A hierarchical construction of the PCUREs for  $n$ -dimensional UREs.

associated with the resulting index space at level  $i$  by  $E(i)$ . The decomposition tree has  $n-2$  levels; the first level is labelled 1. The following recursive equation assigns a unique identifier to each node in the decomposition tree (Fig. 5-9):

$$\begin{aligned} \Phi(1, \ell) &= \Phi \\ (\forall i : 1 < i < n : \Phi(i, \ell) &= \text{fst}(\Phi(i-1, \ell), \vartheta_{E(i-1)})) \\ (\forall i : 1 < i < n : \Phi(i, r) &= \text{lst}(\Phi(i-1, \ell), \vartheta_{E(i-1)})) \end{aligned} \quad (5.26)$$

$\Phi(i, \ell)$  is called a *left* node and  $\Phi(i, r)$  a *right* node. Note that the range of  $i$  is extended by 1 so that the specification of the PCUREs for three-dimensional UREs, as presented previously, becomes a special case.

Next, we present the specification of the PCUREs as constructed previously, for an index space that is  $n$ -parallelepiped. Let  $\vartheta_{E(1)}, \vartheta_{E(2)}, \dots, \vartheta_{E(n)}$  be the  $n$  linearly independent vectors that generate the parallelepiped.  $\Phi(i+1, \ell)$  and  $\Phi(i+1, r)$  are the sets of first and last computation points of  $\Phi(i, \ell)$ . We choose  $\mathbb{F}_i = \{\Phi(i+1, \ell)\}$  and  $\mathbb{L}_i = \{\Phi(i+1, r)\}$  for  $\Phi(i, \ell)$ ;  $\mathbb{F}_i$  plays the rôle of  $\mathbb{F}$  and  $\mathbb{L}_i$  the rôle of  $\mathbb{L}$  in the case of a three-dimensional index space.

- *Evolution Control Variables.* Every left non-leaf node  $\Phi(i, \ell)$  is a facet of its father node  $\Phi(i-1, \ell)$ . It is associated with one evolution control variable  $E(i)$ .  $\text{sig}(E(i)) = \{e(i)_u, e(i)_s, e(i)_1, e(i)_2, \dots, e(i)_{\Gamma_{E(i)}}, e(i)_d\}$ . The communication constraint domains of the evolution control variables are given by:

$$(\forall i : 0 < i \leq n-2 : \Omega_{E(i)} = \text{planes}(\Phi, \sigma) \cap \text{rays}(\Phi(i, \ell), \pm \vartheta_{E(i)}))$$



- *Initialisation Control Variables.* The left leaf node  $\Phi(n-1, \ell)$  is a (two-dimensional) facet of  $\Phi(n-2, \ell)$ . It is associated with two initialisation control variables,  $F(n-1)_1$  and  $F(n-1)_2$ .  $\text{sig}(F(n-1)_1) = \{f(n-1)_1, \bar{f}(n-1)_1\}$ .  $\text{sig}(F(n-1)_2) = \{f(n-1)_2, \bar{f}(n-1)_2\}$ . We choose  $\vartheta_{F(n-1)_1}$  and  $\vartheta_{F(n-1)_2}$  such that they are parallel to two adjacent edges of  $\Phi(n-1, \ell)$ , which are parallel to  $\vartheta_{E(n-1)}$  and  $\vartheta_{E(n)}$ , respectively. Clearly, this ensures that  $\vartheta_{F(n-1)_1}, \vartheta_{F(n-1)_2} \in \text{lin}(\Phi(n-1, \ell))$ . The communication constraint domains of the initialisation control variables are given by

$$\Omega_{F(n-1)_1} = \Omega_{F(n-1)_2} = \emptyset$$

- *Termination Control Variables.* Every right node  $\Phi(i, r)$  is a facet of its father node  $\Phi(i-1, \ell)$ . It is associated with one termination control variable,  $L(i)$ .  $\text{sig}(L(i)) = \{\ell(i), \bar{\ell}(i)\}$ . We choose  $\vartheta_{L(i)}$  from the linear space  $\text{span}(\vartheta_{E(i+1)}, \vartheta_{E(i+2)}, \dots, \vartheta_{E(n)})$ . Clearly,  $\vartheta_{L(i)} \in \text{lin}(\Phi(i, r))$ . The communication constraint domains of the termination control variables are given by

$$(\forall i : 1 < i < n : \Omega_{L(i)} = \text{planes}(\Phi, \sigma) \cap \text{rays}(\Phi(i-1, \ell), \pm \vartheta_{L(i)}))$$

Because of the recursive nature of the specification of the PCUREs, the control variables associated with all the nodes in the subtree rooted at node  $\Phi(i, \ell)$  act as initialisation control variables for the evolution control variable associated with the father node  $\Phi(i-1, \ell)$ .

Next, we define the characteristic functions of all nodes in the decomposition tree. Let  $\mathcal{B}(F(i))$  be the characteristic function of every left node  $\Phi(i, \ell)$ , except the root  $\Phi(1, \ell)$ , and  $\mathcal{B}(L(i))$  the characteristic function of every right node  $\Phi(i, r)$ . These characteristic functions are recursively defined as follows:

$$\begin{aligned} \mathcal{B}(F(n-1)) &= E(n-2)^{\text{in}} = e(n-2)_s \wedge F(n-1)_1^{\text{in}} = f(n-1)_1 \wedge F(n-1)_2^{\text{in}} = f(n-1)_2 \\ (\forall i : 1 < i < n-1 : \mathcal{B}(F(i)) &= E(i-1)^{\text{in}} = e(i-1)_s \wedge \mathcal{B}(F(i+1)) \wedge E(i)^{\text{in}} = e(i)_{\Gamma_{E(i)}}) \\ (\forall i : 1 < i < n : \mathcal{B}(L(i)) &= (\mathcal{B}(F(i)) \vee E(i-1)^{\text{in}} = e(i-1)_{\Gamma_{E(i-1)}}) \wedge L(i)^{\text{in}} = \ell(i)) \end{aligned}$$

The characteristic function  $\chi_{\bar{\Phi}}$  of  $\bar{\Phi}$  in  $\Upsilon$  is given by

$$\chi_{\bar{\Phi}} = \text{if } \mathcal{B}(F(2)) \vee E(1)^{\text{in}} = e(1)_{\Gamma_{E(1)}} \rightarrow \chi_c \text{ \texttt{[]}} \text{ else } \rightarrow \chi_p \text{ \texttt{fi}} \quad (5.27)$$



**PROGRAM:** *HostProg*

```

for  $t$  from  $t_{\text{fst}}$  to  $t_{\text{lst}}$  do
  if  $t \in \{\text{input}(F(n-1)_1(I)) \mid I \in \Phi(n-1, \ell)\} \rightarrow \text{inject}(f(n-1)_1, \text{pi}(F(n-1)_1^{\text{in}}))$ 
   $\square$  else  $\rightarrow \text{inject}(\bar{f}(n-1)_1, \text{pi}(F(n-1)_1^{\text{in}}))$ 
  fi
  if  $t \in \{\text{input}(F(n-1)_2(I)) \mid I \in \Phi(n-1, \ell)\} \rightarrow \text{inject}(f(n-1)_2, \text{pi}(F(n-1)_2^{\text{in}}))$ 
   $\square$  else  $\rightarrow \text{inject}(\bar{f}(n-1)_2, \text{pi}(F(n-1)_2^{\text{in}}))$ 
  fi
  for  $i$  from 2 to  $n-1$  do
    if  $t \in \{\text{input}(L(i)(I)) \mid I \in \Phi(i, r)\} \rightarrow \text{inject}(\ell(i), \text{pi}(L(i)^{\text{in}}))$ 
     $\square$  else  $\rightarrow \text{inject}(\bar{\ell}(i), \text{pi}(L(i)^{\text{in}}))$ 
    fi
  for  $i$  from 1 to  $n-2$  do
    if  $t \in \{\text{input}(E(i)(I)) \mid I \in \Phi(i, \ell)\} \rightarrow \text{inject}(e(i)_s, \text{pi}(E(i)^{\text{in}}))$ 
     $\square$  else  $\rightarrow \text{inject}(e(i)_u, \text{pi}(E(i)^{\text{in}}))$ 
    fi

```

**Table 5-5:** The host program for  $n$ -dimensional UREs ( $n > 2$ ).

$\mathcal{B}(F(n-1))$  is the previous  $\mathcal{B}(F)$  and  $\mathcal{B}(L(n-1))$  is the previous  $\mathcal{B}(L)$ . These characteristic functions can be simplified by removing redundant predicates.

For the host and cell program for an  $n$ -dimensional parallelepiped index space, see Tabs. 5-5 and 5-6. Thm. 5.1 provides the mapping conditions for the validity of space-time mappings for the pipelined UREs whose index space is an  $n$ -dimensional parallelepiped index space. For completeness, we state in a theorem the mapping conditions, including the velocity constraint, for the validity of space-time mappings. We rephrase the communication constraint for propagation control variables to dispense with the concept of the extended index space.

**Theorem 5.6** *A space-time mapping is valid for the pipelined UREs (in either the  $\pi$ -model or the  $\chi$ -model) iff*

- $(\forall V : V \in \mathcal{V} : \lambda \vartheta_V > 0)$  (Precedence Constraint)



**PROGRAM:** *CellProg*

$(\forall i, k : 0 < i \leq n-2 \wedge 0 < k \leq \Gamma_{E(i)} : e(i)_k = k)$

$F(n-1)_1^{\text{out}} = F(n-1)_1^{\text{in}}$

$F(n-1)_2^{\text{out}} = F(n-1)_2^{\text{in}}$

$(\forall i : 1 < i < n : L(i)^{\text{out}} = L(i)^{\text{in}})$

$E(n-2)^{\text{out}} = \text{if } \mathcal{B}(F(n-1)) \wedge \neg \mathcal{B}(L(n-1)) \rightarrow e(n-2)_1$

$\square \mathcal{B}(L(n-1)) \rightarrow e(n-2)_d$

$\square (E(n-2)^{\text{in}} = e(n-2)_s \wedge \neg \mathcal{B}(F(n-1))) \vee$

$E(n-2)^{\text{in}} = e(n-2)_d \vee E(n-2)^{\text{in}} = e(n-2)_u \rightarrow E(n-2)^{\text{in}}$

$\square \text{else} \rightarrow (E(n-2)^{\text{in}} \bmod \Gamma_{E(n-2)}) + 1$

**fi**

$(\forall i : 0 < i < n-2 :$

$E(i)^{\text{out}} = \text{if } \mathcal{B}(F(i-1)) \wedge \neg \mathcal{B}(L(i-1)) \rightarrow e(i)_1$

$\square \mathcal{B}(L(i-1)) \rightarrow e(i)_d$

$\square (E(i)^{\text{in}} = e(i)_s \wedge \neg \mathcal{B}(F(i-1)))$

$\vee E(i)^{\text{in}} = e(i)_d \vee E(i)^{\text{in}} = e(i)_u \rightarrow E(i)^{\text{in}}$

$\square \text{else} \rightarrow (E(i)^{\text{in}} \bmod \Gamma_{E(i)}) + 1$

**fi**

)

$\chi_{\Phi} = \text{if } \mathcal{B}(F(2)) \vee E(1)^{\text{in}} = e(1)_{\Gamma_{E(1)}} \rightarrow \chi_c \square \text{else} \rightarrow \chi_p \text{ fi}$

**Table 5-6:** The cell program for  $n$ -dimensional UREs ( $n > 2$ ).

- $\pi : \Phi \mapsto \mathbf{Z}^2$  (Computation Constraint)
- $(\forall V : V \in \mathcal{V} : \sigma \vartheta_V \mid \lambda \vartheta_V)$  (Delay Constraint for the  $\chi$ -Model)
- $(\forall V : V \in \mathcal{V}_d : (\forall S : S \in \text{fst}(\Phi_V, \vartheta_V) / \textcircled{1}_V : \text{input} : S \mapsto \mathbf{Z}))$   
 $(\forall V : V \in \mathcal{V}_c : (\forall S : S \in \text{fst}(\Psi_V, \vartheta_V) / \textcircled{1}'_V : \text{input} : S \mapsto \mathbf{Z}))$   
 $(\forall i : 0 < i \leq n-2 : (\forall S : S \in \text{fst}(\Phi(i, \ell), \vartheta_{E(i)}) : \text{input} : S \mapsto \mathbf{Z}))$   
 $(\forall i : 1 < i < n : (\forall S : S \in \text{fst}(\Phi(i-1, \ell), \vartheta_{L(i)}) : \text{input} : S \mapsto \mathbf{Z}))$   
(Communication Constraint)
- $(\forall V : V \in \mathcal{V}_p : \text{flow}(V) \neq 0)$  (Velocity Constraint)



**Proof** The proof is an induction on the reverse order of the decomposition of the index space (Fig. 5–9). To prove that  $\mathcal{B}(F(i))$  is a characteristic function of  $\Phi(i, \ell)$ , it suffices to show that  $\mathcal{B}(F(i+1))$  is a characteristic function of  $\Phi(i+1, \ell)$  and  $\mathcal{B}(L(i+1))$  is a characteristic function of  $\Phi(i+1, r)$  by Thm. 5.2. To prove that  $\mathcal{B}(L(i+1))$  is a characteristic function of  $\Phi(i+1, r)$ , it suffices to show that  $\mathcal{B}(F(i+1))$  is a characteristic function of  $\Phi(i+1, \ell)$  by Thm. 5.4. Thus, to prove the correctness of the PCUREs, it suffices to show that  $\mathcal{B}(F(n-1))$  is a characteristic function of  $\Phi(n-1, \ell)$ . This follows from Thm. 5.3.  $\square$

## 5.5 Related Work

Ramakrishnan and Varman [9] present a one-dimensional array for matrix product. They supply the flow of data and control and use linear algebra to prove the correctness of the array. Kumar and Tsai [10] propose a method that transforms two-dimensional to one-dimensional arrays. Their method requires an explicit choice of the communication topology and the sequencing of input data. The flow of data and control is then derived by solving a set of constraint equations on timing. Lang [11] shows how control signals can be pipelined to solve problems such as matrix product and merging two matrices. This method only applies for iterative algorithms defined over three-dimensional index spaces. All these methods do not include procedures that allow the systematic derivation of control signals. Rather, the synthesis of control signals is ad hoc and problem-specific.

## 5.6 Examples

### 5.6.1 Matrix Product

We recommend to choose a data dependence vector as the control dependence vector  $\vartheta_E$ . For reasons of symmetry, it makes no difference which data dependence vector we choose:  $\vartheta_A$ ,  $\vartheta_B$  or  $\vartheta_C$ . Let us choose  $\vartheta_E = \vartheta_C$ . Both  $\text{fst}(\Phi, \vartheta_E)$



and  $\text{lst}(\Phi, \vartheta_E)$  are facets of the index space  $\Phi$ . We choose  $\textcircled{F} = \{\text{fst}(\Phi, \vartheta_E)\}$  and  $\textcircled{L} = \{\text{lst}(\Phi, \vartheta_E)\}$ . We need two initialisation control variables,  $F_1$  and  $F_2$ , and one termination control variable,  $L$ . Hence,  $\vartheta_{F_1}, \vartheta_{F_2} \in \text{lin}(\text{fst}(\Phi, \vartheta_E))$  and  $\vartheta_L \in \text{lin}(\text{lst}(\Phi, \vartheta_E))$ . We should choose two data dependence vectors  $\vartheta_A$  and  $\vartheta_B$  as  $\vartheta_{F_1}$  and  $\vartheta_{F_2}$ . This makes  $\mathfrak{F} = \text{fst}(\Phi, \vartheta_E)$ . Without loss of generality, we choose  $\vartheta_{F_1} = \vartheta_A$  and  $\vartheta_{F_2} = \vartheta_B$ . We should choose either  $\vartheta_A$  or  $\vartheta_B$  as  $\vartheta_L$ . Let us take  $\vartheta_L = \vartheta_A$ . Because  $\Phi/\textcircled{L} = \{\Phi\}$ , the CCUREs are unnecessary. The PCUREs specify the control signals for all one-dimensional systolic arrays that realise matrix product whose three data variables  $A$ ,  $B$  and  $C$  are moving.

Choose the space-time mapping that describes Ramakrishnan and Varman's array (Fig. 4-3). Our array runs in  $(9m^2 - 9m + 2)/2$  steps. The array of Ramakrishnan and Varman runs in  $(12m^2 + 7m + 2)/2$  steps. Both arrays require the same hardware for data flow. But our array is superior in terms of control requirement: it has a total of bit-width of four bits for control signals. Ramakrishnan and Varman's array requires five bits and an encoder and decoder for every cell.

As discussed in Sect. 5.3.4, the evolution control variable can be stationary. If all  $\vartheta_E$ -paths of  $\text{paths}(\Phi, \vartheta_E)$  are mapped to distinct cells, we only need one buffer in each cell to store the corresponding element of  $E$  mapped to the cell. Here is a space-time mapping that has this property:

$$\pi = \begin{bmatrix} m+1 & m & 1 \\ & 1 & m & 0 \end{bmatrix}$$

Both data variable  $C$  and control variable  $E$  are stationary. Each cell needs a buffer for accumulating one element of  $C$ . Similarly, each cell needs a buffer to store one element of  $E$ . The buffers for  $C$  must be initialised to zero and the buffers for  $E$  to  $e_s$  before the first step.

The systolic array consists of  $m^2$  cells and runs in  $m^3 + m^2 - 1$  steps. There is an array proposed by Kumar and Tsai [33] that has the same number of cells but runs in  $3m^2 - m + 3$  steps. The lower latency is obtained at the expense of one additional link for data variable  $B$ . Elements of  $B$  must be frequently routed from one link to the other, complicating the control circuitry.



To compare the hardware required for the control flow: both arrays use a two-bit buffer for  $E$  in each cell. In Kumar and Tsai's array, there are three one-bit moving control variables; two control variables need two buffers per channel; one needs one buffer. In our array, there are also three one-bit moving control variables:  $F_1$ ,  $F_2$ , and  $L$ ;  $F_1$  and  $F_2$  do not need buffers;  $L$  needs  $m$  buffers per channel.

### 5.6.2 Dynamic Programming

Recall that the CCUREs consist of two computation control variables  $P$  and  $Q$  (Sect. 3.7.1.2). We choose  $\vartheta_P = (-1, 0, 0)$  and  $\vartheta_Q = (0, 2, 1)$ .

Since data variable  $E$  has the same variable name as the evolution control variable, we rename it to  $E'$  to avoid confusion. It is desirable to choose one of the three data dependence vectors  $\vartheta_A$ ,  $\vartheta_C$  and  $\vartheta_{E'}$  as the evolution control dependence vector  $\vartheta_E$ . We choose  $\vartheta_E = \vartheta_{E'}$ .

Let us first consider the specification of the initialisation control variables. We choose  $\mathbb{F} = \{\text{fst}(\Phi, \vartheta_E)_1, \text{fst}(\Phi, \vartheta_E)_2\}$ :

$$\begin{aligned}\text{fst}(\Phi, \vartheta_E)_1 &= \Phi \cap \{(i, j, k) \mid j - i = 2k - 1\} \\ \text{fst}(\Phi, \vartheta_E)_2 &= \Phi \cap \{(i, j, k) \mid j - i = 2k - 2\}\end{aligned}$$

We need to choose  $\vartheta_{F_{1,1}}, \vartheta_{F_{1,2}} \in \text{lin}(\text{fst}(\Phi, \vartheta_E)_1)$  and  $\vartheta_{F_{2,1}}, \vartheta_{F_{2,2}} \in \text{lin}(\text{fst}(\Phi, \vartheta_E)_2)$ . Let us choose  $\vartheta_{F_{1,1}} = \vartheta_{F_{2,1}} = (0, 2, 1)$  and  $\vartheta_{F_{1,2}} = \vartheta_{F_{2,2}} = (-2, 0, 1)$ .

The specification of termination control is simple.  $\text{lst}(\Phi, \vartheta_E)$  is a facet of the index space. We choose  $\mathbb{L} = \{\text{lst}(\Phi, \vartheta_E)\}$ . We name the associated termination control variable  $L$  and choose  $\vartheta_L = (0, 1, 0) \in \text{lin}(\text{lst}(\Phi, \vartheta_E))$ .

The following space-time mapping describes the one-dimensional array in [62]:

$$\pi = \begin{bmatrix} -4 & 2(n+2) & -2 \\ -1 & & 1 & 0 \end{bmatrix}$$

Data variable  $C$  is stationary. Elements of  $C$  must be loaded before the computation starts and recovered after the computation completes. Elements of  $A$ ,  $B$ ,  $D$ ,



$E'$ ,  $E$ ,  $F_{1,1}$ ,  $F_{1,2}$ ,  $F_{2,1}$ ,  $F_{2,2}$  and  $L$  encounter a delay of  $2n+3$ ,  $1$ ,  $2n+1$ ,  $3$ ,  $3$ ,  $2n+2$ ,  $2n+2$ ,  $2$ ,  $2$  and  $2n+3$  at the respective channel (remember that our analysis is conducted with respect to the  $\chi$ -model).

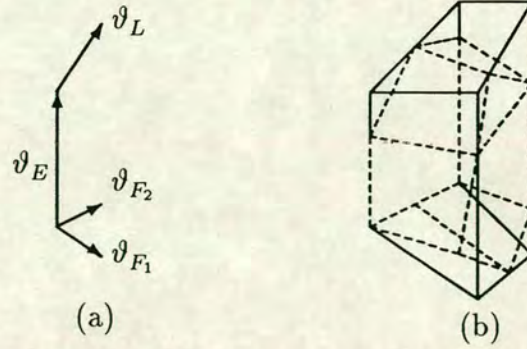
## 5.7 Conclusion

We have provided a construction of propagation control flow for one-dimensional arrays from  $n$ -dimensional UREs. For notational convenience, we have presented the PCUREs in terms of two programs: a host program that is executed by the host and that is responsible for the injection of the propagation control signals, and a cell program that is executed by all cells of the array and that specifies the output control signals of a cell based on the input control signals received by that cell. Similarly to the control UREs constructed in Chap. 3, the control dependence vectors of the propagation control variables in the PCUREs remain to be determined. We have summarised in a procedure how to choose these dependence vectors. We have stated the mapping conditions for the pipelined UREs constructed for one-dimensional arrays. Once the control dependence vectors have been chosen, as demonstrated by two examples in Sect. 5.6, finding one-dimensional arrays is just finding valid space-time mappings.

The idea underlying the construction of the PCUREs may also shed light on solving other issues in systolic design. As has been proved in [64], a systolic array is a system of space-time UREs. Once we have formulated a systolic array model and the corresponding mapping conditions for the validity of space-time mappings, we can conduct our analysis at the source level, i.e., independently of the space-time mapping. This has the benefit that we do not have to be concerned with temporal and spatial issues of the systolic array. Rather, the synthesis of control signals can be considered as a process of program construction and transformation.

Let us now have a look at how an adaptation of the index space can simplify the specification of the initialisation and termination control variables. We describe the basic idea with respect to an artificial example shown in Fig. 5–10.  $\text{fst}(\Phi, \vartheta_E)$





**Figure 5–10:** Adaptation of the index space for the construction of the PCUREs  
 (a) The control dependence vectors  $\vartheta_E$ ,  $\vartheta_{F_1}$ ,  $\vartheta_{F_2}$ ,  $\vartheta_L$  are shown. (b) The dashed polytope depicts the index space  $\Phi$ . Both  $\text{fst}(\Phi, \vartheta_E)$  and  $\text{lst}(\Phi, \vartheta_E)$  are not facets of  $\Phi$ . The polytope that encloses  $\Phi$  depicts the adapted index space  $\Phi_A$ .

and  $\text{lst}(\Phi, \vartheta_E)$  each contain two facets of the index space. Therefore,  $|\mathbb{F}| \geq 2$  and  $|\mathbb{L}| \geq 2$ . Assume that  $|\mathbb{F}| = |\mathbb{L}| = 2$ . We then need four initialisation control variables and two termination control variables. A simple adaptation of the index space can reduce the number of each type of control variable by half. Let us denote the two initialisation control variables by  $F_1$  and  $F_2$  and the termination control variable by  $L$ . Let their control dependence vectors be as shown in Fig. 5–10(a). Let  $H(F_1, F_2)$  be the unique hyperplane that is spanned by  $\vartheta_{F_1}$  and  $\vartheta_{F_2}$  and that bounds the index space from below:

$$H(F_1, F_2) = \{I \mid (\vartheta_{F_1} \times \vartheta_{F_2})I = \alpha\}$$

where “ $\times$ ” denotes vector product. We extend the index space along  $-\vartheta_E$  until the points intersect  $H(F_1, F_2)$ . Similarly, let  $H(L)$  be a hyperplane whose normal is orthogonal to  $\vartheta_L$  and that bounds the index space from above (the choice of this hyperplane is not unique):

$$H(L) = \{I \mid \pi_L I = \beta \wedge \pi_L \vartheta_L = 0\}$$

We then extend the index space along  $\vartheta_E$  until the extension intersects  $H(L)$ . The *adapted* index space, denoted by  $\Phi_A$ , is given by

$$\begin{aligned} \Phi_A = & (\text{planes}(\text{fst}(\Phi, \vartheta_E), \vartheta_{F_1} \times \vartheta_{F_2}) \cap \text{rays}(\text{fst}(\Phi, \vartheta_E), -\vartheta_E)) \cup \Phi \cup \\ & (\text{planes}(\text{lst}(\Phi, \vartheta_E), \pi_L) \cap \text{rays}(\text{lst}(\Phi, \vartheta_E), \vartheta_E)) \end{aligned}$$

Both  $\text{fst}(\Phi_A, \vartheta_E)$  and  $\text{lst}(\Phi_A, \vartheta_E)$  are facets of the adapted index space  $\Phi_A$ .



$\Phi_A$  is a superset of  $\Phi$ . A fixed space-time mapping may yield the same latency and processor space in both cases. This happens when the extension does not change the longest dependence path and the size of the projection space. If such an adaptation of the index space degrades the performance of the systolic array too much, it had better be avoided.

Finally, we point out that the PCUREs presented in this chapter work for systolic arrays of  $r$ -dimensions ( $1 < r < n$ ) that conform to the systolic array model defined by Def. 2.2 and that are described by the space-time mappings defined in Def. 2.4, whose validity is given by Def. 2.4. The reason is that Lemma 5.3, on which the construction of the PCUREs (or more precisely, the construction of the initialisation control variables) depends, is independent of the dimensionality of the systolic array. Lemma 5.3 is still valid if  $\sigma$  is an  $r \times n$  allocation matrix rather than an allocation vector.

However, the PCUREs constructed here can be greatly simplified in the case of  $(n-1)$ -dimensional systolic arrays, since valid space-time mappings are bijections from  $\mathbb{Q}^n$  of  $\mathbb{Q}^n$ . This has the following two consequences. First, the hierarchical decomposition of the index space, as in Sect. 5.4, in order to construct the PCUREs for  $n$ -dimensional UREs is not necessary. The PCUREs that are presented with respect to three-dimensional UREs apply for  $n$ -dimensional UREs. Second, the construction of the initialisation control variables can be conducted in exactly the same way as that of the termination control variables. This halves the number of initialisation control variables required.

Now, we have two different constructions of the PCUREs for  $(n-1)$ -dimensional arrays. In general, the PCUREs that are constructed in Chap. 3 are superior to the PCUREs constructed here, partly because they are simpler and partly because they do not degrade the latency of the systolic array as they do in the case of one-dimensional arrays (Sect. 5.1).



## Chapter 6

# The Elimination of Propagation Control Flow

### 6.1 Introductory Remarks

Why would we want to eliminate propagation control signals for systolic arrays? Consider source UREs that do not require any CCUREs, i.e., that satisfy  $\Phi/\textcircled{\text{T}} = \{\Phi\}$ . That is, all cells in the systolic array are identical. At any given step, a cell is in one of two possible states: the *computation state* when it is performing a common computation as specified by the source UREs and the *propagation state* when it is propagating input/output data. The PCUREs specify when a cell is in which of these two states.

Recall that *pattern* specifies the distribution of the input data of a variable for  $(n-1)$ -dimensional arrays and *input* and *pi* together play the same rôle as *pattern* for one-dimensional arrays. If a variable is not given a value at a given step, the undefined value  $\perp$  is injected. If we substitute any value for  $\perp$ , we still obtain the output data specified by the source UREs from the systolic array. We use this fact to eliminate the propagation control signals in the systolic array. Rather than replacing  $\perp$  arbitrarily, we choose the substitute value carefully. We then let all cells of the array adopt the computation state at every step. If we still obtain the output data specified by the source UREs, we can do without propagation control signals.



Let us examine this scheme formally. Because of our assumption that  $\Phi/\textcircled{\oplus} = \{\Phi\}$ , i.e., that  $(\forall V : V \in \mathcal{V} : \Phi_V = \Phi)$ , a variable  $V$  in the extended source UREs is specified by the two equations:

$$I \in \Phi \rightarrow V(I) = f(W(I - \vartheta_W), \dots) \quad (6.1)$$

$$I \in \Psi_V \setminus \Phi_V \rightarrow V(I) = V(I - \vartheta_V) \quad (6.2)$$

The first equation is from the source UREs. The second equation results from the extension of the index space (Def. 2.5). Propagation control is required iff the two equations for some data variable differ.

Since the space-time mappings for one-dimensional arrays are non-injective, we discuss the elimination of propagation control with respect to space-time UREs (Sect. 2.4) rather than the source UREs. Assume that  $W(\bar{I} - \bar{\vartheta}_W) = \perp$  at pipelining point  $\bar{I} = (t, p)$ . Instead of letting cell  $p$  evaluate equation  $V(\bar{I}) = V(\bar{I} - \bar{\vartheta}_V)$  at step  $t$ , we let it evaluate the equation  $V(\bar{I}) = f(W(\bar{I} - \bar{\vartheta}_W), \dots)$  by substituting  $W(\bar{I} - \bar{\vartheta}_W) = w_{\bar{I}}$  for  $W(\bar{I} - \bar{\vartheta}_W) = \perp$ . The choice of value  $w_{\bar{I}}$  will be explained later. We still obtain the output data specified by the source UREs from the systolic array iff the following two conditions are satisfied.

- Cond. 1.  $W$  is undefined at all points of the  $\bar{\vartheta}_W$ -path that passes through the point  $\bar{I}$ , i.e.,  $(\forall \bar{J} : \bar{J} \in p(\bar{\Psi}, \bar{\vartheta}_W, \bar{I}) : W(\bar{J}) = \perp)$ . (This ensures that value  $w_{\bar{I}}$  can be injected from input cells, as is required in systolic array models.)
- Cond. 2.  $V(\bar{I} - \bar{\vartheta}_V) = f(W(\bar{I} - \bar{\vartheta}_W), \dots)$ . (This preserves the behaviour of the (extended) source UREs.)

To satisfy Cond. 1, we choose the space-time mapping carefully. To satisfy Cond. 2, we exploit the algebraic properties of the operators in the equation (6.1).

There cannot be a general method for all systolic algorithms, since the algebraic properties of operators are problem-specific. We shall present our method with respect to a restricted class of UREs called *inner-product UREs*. The inner-product UREs are a subclass of UREs as defined in [55]. The basic idea can be generalised to other systolic algorithms.



Inner-Product UREs:

$$\begin{aligned}
 A(I) &= \begin{cases} I \in \text{in}(\Phi_A, \vartheta_A) & \rightarrow A(I) = a_{\text{in}_A(I)} \\ I \in \Phi_A & \rightarrow A(I) = A(I - \vartheta_A) \end{cases} \\
 B(I) &= \begin{cases} I \in \text{in}(\Phi_B, \vartheta_B) & \rightarrow B(I) = b_{\text{in}_B(I)} \\ I \in \Phi_B & \rightarrow B(I) = B(I - \vartheta_B) \end{cases} \\
 C(I) &= \begin{cases} I \in \text{in}(\Phi_C, \vartheta_C) & \rightarrow C(I) = c_{\text{in}_C(I)} \\ I \in \Phi_C & \rightarrow C(I) = C(I - \vartheta_C) \oplus (A(I - \vartheta_A) \otimes B(I - \vartheta_B)) \end{cases}
 \end{aligned}$$

where  $\Phi_A = \Phi_B = \Phi_C = \Phi$ , and the functions  $\text{in}_A$ ,  $\text{in}_B$  and  $\text{in}_C$  are defined in Sect. 2.2. The output equations are irrelevant for our discussion and are thus omitted. We assume that the values represented by variables  $A$ ,  $B$  and  $C$  belong to a set  $M$ .  $\otimes$  and  $\oplus$  are binary operations on the set  $M$ . To satisfy Cond. 2, we postulate the existence of elements  $u, z_1, z_2$  (which may be equal) in  $M$  satisfying:

$$\begin{aligned}
 (\forall s : s \in M : s \oplus u = s) \\
 (\forall s : s \in M : z_1 \otimes s = u) \\
 (\forall s : s \in M : s \otimes z_2 = u)
 \end{aligned} \tag{6.3}$$

The element  $u$  is a right unit of  $\oplus$ . The idea is to replace  $A(\bar{I}) = \perp$  ( $B(\bar{I}) = \perp$ ) by  $A(\bar{I}) = z_1$  ( $B(\bar{I}) = z_2$ ). Examples that are instances of inner-product UREs are 1-D convolution and matrix product. In this case,  $\oplus$  becomes addition,  $\otimes$  becomes multiplication,  $M$  becomes  $\mathbf{R}$ , and  $z_1 = z_2 = u = 0$ . 0 is both the unit element for addition and the zero element for multiplication.

The rest of this chapter is organised as follows. Sect. 6.2 motivates the elimination of propagation control signals with the example of 1-D convolution. Sect. 6.3 discusses the elimination of propagation control signals for  $(n - 1)$ -dimensional systolic arrays. Sect. 6.4 continues our presentation for one-dimensional arrays.

## 6.2 Example: 1-D Convolution

Given two input sequences  $x_1, x_2 \dots, x_m$  and  $w_1, w_2 \dots, w_m$ , 1-D convolution computes an output sequence  $y_1, y_2 \dots, y_{2m-1}$  where  $y_i$  is given by:



Specification:  $(\forall i : 0 < i < 2m : y_i = (\sum k : 0 < k \leq m \wedge 0 \leq i - k < m : x_{i-k+1} w_k))$

Source UREs:

$$\begin{aligned}
 X(i, k) &= \begin{cases} 0 \leq i < m \wedge k = 0 & \rightarrow x_{i+1} \\ m \leq i \leq 2m - 2 \wedge k = 0 & \rightarrow 0 \\ 0 = i \wedge 0 < k < m & \rightarrow 0 \\ 0 < i < 2m \wedge 0 < k \leq m & \rightarrow X(i-1, k-1) \bullet \end{cases} \\
 W(i, k) &= \begin{cases} 0 = i \wedge 0 < k \leq m & \rightarrow w_k \\ 0 < i < 2m \wedge 0 < k \leq m & \rightarrow W(i-1, k) \bullet \end{cases} \\
 Y(i, k) &= \begin{cases} 0 < i < 2m \wedge k = 0 & \rightarrow 0 \\ 0 < i < 2m \wedge 0 < k \leq m & \rightarrow Y(i, k-1) + X(i-1, k-1)W(i-1, k) \bullet \end{cases} \\
 y_i &= \begin{cases} 0 < i < 2m \wedge k = m & \rightarrow Y(i, k) \end{cases}
 \end{aligned}$$

Index Space:  $\Phi = \{(i, k) \mid 0 < i < 2m \wedge 0 < k \leq m\}$

Data Dependence Matrix:  $\mathcal{D} = [\vartheta_X, \vartheta_Y, \vartheta_W] = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$

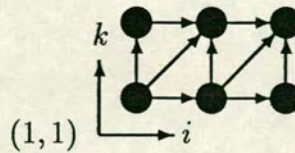
First Computation Points:  $\text{fst}(\Phi, \vartheta_X) = \{(i, k) \mid (0 < i < 2m \wedge 1 = k) \vee (1 = i \wedge 1 < k \leq m)\}$

$\text{fst}(\Phi, \vartheta_Y) = \{(i, k) \mid 0 < i < 2m \wedge 1 = k\}$

$\text{fst}(\Phi, \vartheta_W) = \{(i, k) \mid 1 = i \wedge 0 < k \leq m\}$

Last Computation Points:  $\text{lst}(\Phi, \vartheta_Y) = \{(i, k) \mid 0 < i < 2m \wedge k = m\}$

Data Dependence Graph ( $m=2$ ):



The UREs for 1-D convolution satisfy  $\Phi/\oplus = \{\Phi\}$ . Thus, CCUREs are not required. Only PCUREs are necessary to tell whether a cell is in a state of computation or propagation. Let us consider the specification of the PCUREs.  $\mathfrak{B}$  contains the following supporting half-spaces to the index space:



$$\begin{aligned}
B_1 &= \{(i, j, k) \mid i \geq 1\} \\
B_2 &= \{(i, j, k) \mid i \leq 2m-1\} \\
B_3 &= \{(i, j, k) \mid k \geq 1\} \\
B_4 &= \{(i, j, k) \mid k \leq m\}
\end{aligned}$$

The PCUREs follow from (3.18). If we observe that the two hyperplanes corresponding to  $B_1$  and  $B_2$  are parallel, we can combine the two associated control variables into one, denoted by  $P_1$ . This reduces the number of control values from four to three. Variable  $P_1$  defines three control values at the domains  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_1$ ,  $\Psi \cap B_1 \cap B_2$  and  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_2$ . It suffices to use the same control value for the two domains  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_1$  and  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_2$ , because both contain pipelining points. We denote the control value defined at  $\Psi \cap B_1 \cap B_2$  by  $p_1$  and the control value at  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_1$  and  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_2$  by  $\bar{p}_1$ . For reasons of symmetry, we use one control variable,  $P_2$ , for  $B_3$  and  $B_4$ . We denote the control value defined at  $\Psi \cap B_3 \cap B_4$  by  $p_2$  and the control value defined at  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_3$  and  $\Psi \cap \mathcal{L}_{\mathbf{Z}^n} B_4$  by  $\bar{p}_2$ . By (3.19), cell  $p$  is in the computation state at step  $t$  if  $P_1(\bar{I})=p_1 \wedge P_2(\bar{I})=p_2$  holds at point  $\bar{I}=(t, p)$  and in the propagation state otherwise.

We choose the space-time mapping (Figs. 6-1(a), (c) and (d)).

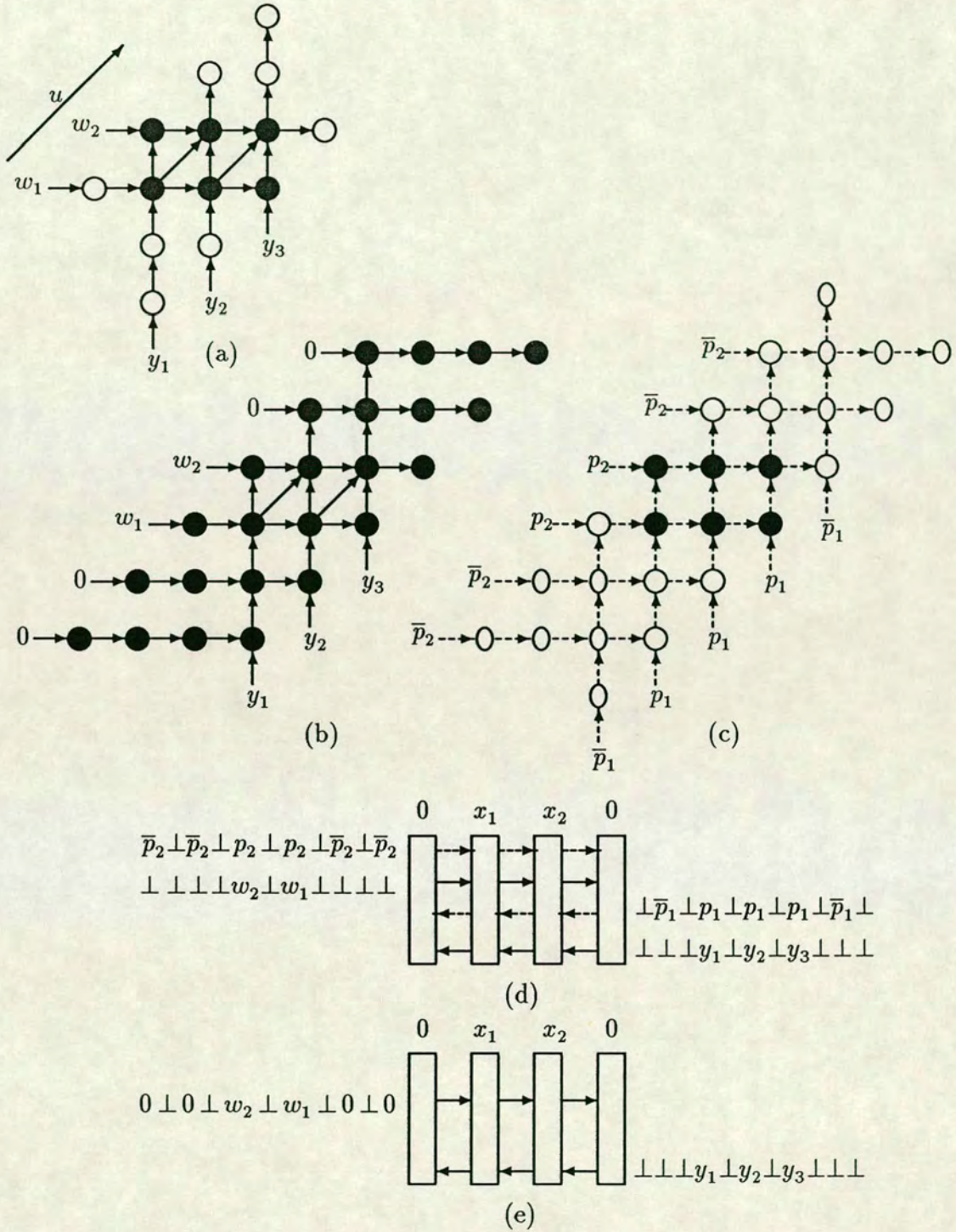
$$\pi = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

The projection vector is  $u=(1, 1)$ .

Both  $X$  and  $W$  are specified by pipelining equations over the extended index space. The only reason why we need propagation control is that  $Y$  is defined by the equation  $Y(I)=Y(I-\vartheta_Y)$  at the pipelining points and by the equation  $Y(I)=Y(I-\vartheta_Y)+X(I-\vartheta_X)W(I-\vartheta_W)$  at the computation points. We can eliminate propagation control signals in the systolic array if we do the following:

- Inject 0 instead of  $\perp$  for variable  $W$  at steps  $t_{\text{fst}}, t_{\text{fst}}+2, t_{\text{fst}}+4, \dots$  excluding those at which  $w_1, w_2, \dots, w_m$  are injected.
- Evaluate  $Y$  at the pipelining points  $\bar{I}$  using equation  $Y(\bar{I})=Y(\bar{I}-\bar{\vartheta}_Y)+X(\bar{I}-\bar{\vartheta}_X)W(\bar{I}-\bar{\vartheta}_W)$  rather than  $Y(\bar{I})=Y(\bar{I}-\bar{\vartheta}_Y)$ .





**Figure 6-1:** (a) The extended dependence graph with respect to projection vector  $u=(1,1)$  ( $m=2$ ). (b) The extension of the extended data dependence graph of (a) for the purpose of substituting  $W(\bar{I}-\bar{\vartheta}_W)=0$  for  $W(\bar{I}-\bar{\vartheta}_W)=\perp$ . (c) The extension of the extended data dependence graph of (a) by Def. 2.5 for propagation control variables  $P_1$  and  $P_2$ . The ovals depict the new points generated in the extension. (d) The systolic array with the distribution of both data and control signals at the first step. (e) The systolic array without the propagation control signals.



To see why Conds. 1 and 2 are satisfied, we classify the pipelining points of  $\Psi^P$  (Def. 2.5) into two categories:

- $\Psi_Y^\perp$  contains the points at which  $Y$  is undefined. At these points  $I$ , we can replace the equation  $Y(I) = Y(I - \vartheta_Y)$  by any equation.
- $\Psi_Y^s \cup \Psi_Y^d$  contains the points at which input/output data of  $Y$  are propagated. Because  $\vartheta_Y$  and  $\vartheta_W$  are orthogonal,  $W$  is undefined at these points  $I$ :

$$(\forall I : I \in \Psi_Y^s \cup \Psi_Y^d : (\forall J : J \in p(\Psi, \vartheta_W, I) : W(J) = \perp))$$

The space-time mapping (6.4) is a bijection from  $\mathbb{Q}^n$  to  $\mathbb{Q}^n$  ( $n=2$ ). So

$$(\forall \bar{I} : \bar{I} \in \bar{\Psi}_Y^s \cup \bar{\Psi}_Y^d : (\forall \bar{J} : \bar{J} \in p(\bar{\Psi}, \bar{\vartheta}_W, \bar{I}) : W(\bar{J}) = \perp))$$

Thus, Cond. 1 is satisfied. 0 is both the unit element for addition and the zero element for multiplication. Substituting  $W(\bar{I} - \bar{\vartheta}_W) = 0$  for  $W(\bar{I} - \bar{\vartheta}_W) = \perp$  at the points of  $\bar{\Psi}_Y^s \cup \bar{\Psi}_Y^d$  gives rise to

$$Y(\bar{I} - \bar{\vartheta}_Y) = Y(\bar{I} - \bar{\vartheta}_Y) + X(\bar{I} - \bar{\vartheta}_X)W(\bar{I} - \bar{\vartheta}_W)$$

Thus, Cond. 2 is satisfied.

The previous illustration can be understood by an extension of the extended index space (Fig. 6-1(b)); the extension is carried out by first replicating  $\vartheta_W$  throughout the extended index space and then extending  $\vartheta_W$  and  $-\vartheta_W$  until the points so created are at the boundary of the processor space. Fig. 6-1(e) displays the systolic array without propagation control signals.

### 6.3 Systolic Arrays of $n-1$ Dimensions

We present necessary and sufficient conditions for the elimination of propagation control signals for  $(n-1)$ -dimensional arrays synthesised from inner-product UREs. These conditions depend on the data dependence patterns of the source UREs and on the projection vector.



**Theorem 6.1** *PCUREs are unnecessary for inner-product UREs if*

$$\begin{aligned} (\forall I : I \in \Psi_C^s \cup \Phi_C^d : C(I) \neq \perp \implies \\ (\forall J : J \in p(\Psi, \vartheta_A, I) : A(J) = \perp) \vee (\forall J : J \in p(\Psi, \vartheta_B, I) : B(J) = \perp)) \end{aligned} \quad (6.4)$$

**Proof** We show that the inner-product UREs can be refined to a system of unconditional UREs defined over the extended index space. Thus, propagation control is not needed. We explain at the end of the proof why Conds. 1 and 2 are satisfied.

The extended source UREs of the inner-product UREs follow from Def. 2.5:

$$\begin{aligned} A(I) &= \begin{cases} I \in \text{in}(\Psi_A, \vartheta_A) & \rightarrow A(I) = a_{\text{in}_A(\mathcal{I}(I))} \\ I \in \Psi_A & \rightarrow A(I) = A(I - \vartheta_A) \end{cases} \\ B(I) &= \begin{cases} I \in \text{in}(\Psi_B, \vartheta_B) & \rightarrow B(I) = b_{\text{in}_B(\mathcal{I}(I))} \\ I \in \Psi_B & \rightarrow B(I) = B(I - \vartheta_B) \end{cases} \\ C(I) &= \begin{cases} I \in \text{in}(\Psi_C, \vartheta_C) & \rightarrow C(I) = c_{\text{in}_C(\mathcal{I}(I))} \\ I \in \Phi_C & \rightarrow C(I) = C(I - \vartheta_C) \oplus (B(I - \vartheta_A) \otimes B(I - \vartheta_B)) \\ I \in \Psi_C^\perp \cup \Psi_C^s \cup \Psi_C^d & \rightarrow C(I) = C(I - \vartheta_C) \end{cases} \end{aligned} \quad (6.5)$$

The computation equations of both  $A$  and  $B$  over the extended index space  $\Psi$  are pipelining equations. The only reason why we need propagation control signals is that  $C$  is defined by different equations at pipelining and computation points. To eliminate propagation control signals, as discussed informally in Sect. 6.2, we replace the defining equation at the pipelining points by the defining equation at the computation points and add some new input equations for  $A$  or  $B$  or both.

Let  $A_\perp$  ( $B_\perp$ ) be the union of the  $\vartheta_A$ -paths ( $\vartheta_B$ -paths) that intersect  $\Psi_C^s \cup \Psi_C^d$  and whose points are undefined points of  $A$  ( $B$ ). The extension of the index space ensures

$$\begin{aligned} \text{in}(A_\perp, \vartheta_A) \cap \text{in}(\Psi_A, \vartheta_A) &= \emptyset \\ \text{in}(B_\perp, \vartheta_B) \cap \text{in}(\Psi_B, \vartheta_B) &= \emptyset \\ I \in \text{in}(A_\perp, \vartheta_A) &\rightarrow A(I) = \perp \\ I \in \text{in}(B_\perp, \vartheta_B) &\rightarrow B(I) = \perp \end{aligned}$$



Hence, the following defining equations of  $A$  and  $B$  are a refinement of their counterparts of (6.5):

$$\begin{aligned}
 A(I) &= \begin{cases} I \in \text{in}(A_{\perp}, \vartheta_A) & \rightarrow A(I) = z_1 \\ I \in \text{in}(\Psi_A, \vartheta_A) & \rightarrow A(I) = a_{\text{in}_A(I)} \\ I \in \Psi_A & \rightarrow A(I) = A(I - \vartheta_A) \end{cases} \\
 B(I) &= \begin{cases} I \in \text{in}(B_{\perp}, \vartheta_B) & \rightarrow B(I) = z_2 \\ I \in \text{in}(\Psi_B, \vartheta_B) & \rightarrow B(I) = b_{\text{in}_B(I)} \\ I \in \Psi_B & \rightarrow B(I) = B(I - \vartheta_B) \end{cases}
 \end{aligned} \tag{6.6}$$

To eliminate propagation control, it suffices to show that the following defining equations of  $C$  are a refinement of their counterparts of (6.5):

$$C(I) = \begin{cases} I \in \text{in}(\Psi_C, \vartheta_C) & \rightarrow C(I) = c_{\text{in}_C(I)} \\ I \in \Psi_C & \rightarrow C(I) = C(I - \vartheta_C) \oplus (A(I - \vartheta_A) \otimes B(I - \vartheta_B)) \end{cases} \tag{6.7}$$

Variable  $C$  is undefined for the domain  $\Psi_C^{\perp}$ . For the points of  $\Psi_C^{\perp}$  and the points  $I$  of  $\Psi_C^s \cup \Psi_C^d$  at which  $C(I) = \perp$ , the defining equation of  $C$  is free of interpretation.

For point  $I \in \Psi_C^s \cup \Psi_C^d$  at which  $C(I) \neq \perp$ , hypothesis (6.4) states:

$$\begin{aligned}
 &(\forall J : J \in p(\Psi, \vartheta_A, I) : A(J) = \perp) \text{ or} \\
 &(\forall J : J \in p(\Psi, \vartheta_B, I) : B(J) = \perp)
 \end{aligned}$$

So, the introduction of the two new equations in (6.6) gives rises to:

$$\begin{aligned}
 &(\forall J : J \in p(\Psi, \vartheta_A, I) : A(J) = z_1) \text{ or} \\
 &(\forall J : J \in p(\Psi, \vartheta_B, I) : B(J) = z_2)
 \end{aligned}$$

In either case, assumption (6.3) ensures the equality of  $C(I - \vartheta_C)$  and  $C(I - \vartheta_C) \oplus (A(I - \vartheta_A) \otimes B(I - \vartheta_B))$ .

(The satisfaction of Cond. 1 is enforced by hypothesis (6.4) and the bijectivity of space-time mappings. The satisfaction of Cond. 2 is enforced by assumption (6.3).)  $\square$

The PCUREs are not needed for a special instance of inner-product UREs, independently of the projection vector, namely, when

$$\text{rays}(\Phi_A, \vartheta_A) \cap \text{rays}(\Phi_C, \vartheta_C) = \text{rays}(\Phi_B, \vartheta_B) \cap \text{rays}(\Phi_C, \vartheta_C) = \Phi$$

In this case, hypothesis (6.4) is always satisfied, independently of the projection vector. Examples include matrix product and 1-D convolution.



## 6.4 Systolic Arrays of One Dimension

Thm. 6.1 is not valid for one-dimensional arrays due to the non-injectivity of the space-time mapping. A counterexample is Ramakrishnan and Varman's one-dimensional array for matrix product (Fig. 4-3). Hypothesis (6.4) is only necessary but no longer sufficient for the elimination of propagation control in one-dimensional arrays. A sufficient condition must depend on the scheduling vector.

**Theorem 6.2** *PCUREs are unnecessary for inner-product UREs if*

$$\begin{aligned} & (\forall \bar{I} : \bar{I} \in \Upsilon_C^s \cup \Upsilon_C^r \cup \Upsilon_C^d : C(\bar{I}) \neq \perp \implies \\ & (\forall \bar{J} : \bar{J} \in p(\Upsilon, \bar{\vartheta}_A, \bar{I}) : A(\bar{J}) = \perp) \vee (\forall \bar{J} : \bar{J} \in p(\Upsilon, \bar{\vartheta}_B, \bar{I}) : B(\bar{J}) = \perp)) \end{aligned} \quad (6.8)$$

**Proof** Similar to that of Thm. 6.1, except that the proof is conducted with respect to the space-time diagram  $\Upsilon$  instead of the extended index space  $\Psi$ .  $\square$

The following space-time mapping satisfies hypothesis (6.8):

$$\pi = \begin{bmatrix} 6m-1 & 1 & 1 \\ & 1 & -1 \end{bmatrix}$$

The elimination of propagation control may increase the latency of the array in order to satisfy the hypothesis (6.8). In the presence of propagation control signals, this hypothesis need not be satisfied. One valid improvement of the previous space-time mapping that violates the hypothesis is:

$$\pi = \begin{bmatrix} 2m-2 & 1 & 1 \\ & 1 & -1 \end{bmatrix}$$

The latency of the former space-time mapping is  $18m^2 - 18m + 1$  steps, the latency of the latter is only  $6m^2 - 9m + 4$ . According to the PCUREs for matrix product in Sect. 5.6.1, the latter array needs three one-bit links for propagation control signals.



## 6.5 Conclusion

We have made an attempt of a formal treatment of the elimination of propagation control for systolic arrays that do not require computation control. We first described informally two prerequisites, Conds. 1 and 2, for the elimination of the propagation control. We then presented necessary and sufficient conditions for the elimination of propagation control with respect to a restricted class of UREs: inner-product UREs. We distinguish two cases: the synthesis of  $(n-1)$ -dimensional arrays and of one-dimensional arrays. We have not given a general method for  $n$ -dimensional UREs. But the underlying idea is general enough to be applied to other systolic algorithms.

We have stated before that the synthesis of systolic arrays from programs is a process of program transformation. Most transformations that are reported in the literature are syntactic adaptations of the source program; they amount to a rearrangement of the data dependences in the source program. Transformations such as uniformisation and space-time mapping are typical syntactic adaptations. The transformations that we have applied here to inner-product UREs for the purpose of eliminating propagation control are semantic adaptations; they take semantic aspects of the source programs into account. The relevance of semantic adaptations to systolic design cannot be overemphasised. By realising that some variables always denote the same value, wrap-around channels can be eliminated from systolic arrays for transitive closure [39]. The latency of two-dimensional arrays for band matrix product can be reduced by counting the accumulating index ( $k$  in Ex. 2.1) down rather than up [27]. We succeeded in eliminating diagonal channels in systolic arrays for pyramidal algorithms by means of removing redundant computations in the source program [45]. A more formal treatment of the exploitation of commutativity and associativity of operators in systolic design is reported in [59]. We expect to see more work along this line.



## Chapter 7

# The Loading, Recovery and Access of Stationary Data

### 7.1 Introductory Remarks

The control UREs provide a complete specification of control signals for systolic arrays if all the data variables are moving. In the presence of stationary data variables, local processor memory is required. The space-time mapping does not provide any help for handling stationary data. Thus, a scheme for loading, recovering and accessing stationary data and the according control signals required remain to be developed. In practice, it is sometimes advantageous to selectively make certain data variables stationary in order to satisfy certain performance requirements, e.g., to reduce the latency and/or the size of the systolic array. Thus, the systematic derivation of a scheme for handling stationary data is of practical importance.

We describe a systematic method for handling stationary data; the method exhibits a continuation of the idea underlying the construction of the control UREs, especially the idea underlying the construction of the PCUREs for one-dimensional arrays. The basic idea is to rewrite the source UREs so as to include the specification for the loading, recovery and access of stationary data. Essentially, this new system of UREs contains two additional types of variables:



- *Load-recovery (data) variables* are responsible for the loading and recovery of the elements of stationary data variables of the source UREs.
- *Address (control) variables* provide the addresses for the access of the elements of stationary data variables of the source UREs.

They are pipelining variables; each is associated with one data dependence vector.

To specify the control signals for systolic arrays synthesised from this new system of UREs, we need only construct the control UREs according to the methods described in Chaps. 3 and 5.

The place function determines which data variables are stationary. Therefore, the construction of load-recovery and address variables depends on the choice of place function. We describe the construction of load-recovery and address variables for both  $(n-1)$ -dimensional and one-dimensional arrays synthesised from  $n$ -dimensional UREs. We distinguish the dimensionality of the systolic array under consideration by the rank of the allocation matrix  $\sigma$ . It describes an  $(n-1)$ -dimensional array if  $\text{rank}(\sigma) = n-1$  and a one-dimensional array if  $\text{rank}(\sigma) = 1$ .

The rest of this chapter is organised as follows. Sect. 7.2 describes the construction of load-recovery variables. Sect. 7.3 describes the construction of address variables. Sect. 7.4 contains the conclusion of this chapter.

## 7.2 The Loading and Recovery of Stationary Data

We associate a distinct load-recovery variable with every stationary data variable in the source UREs. We shall conduct our presentation with respect to a fixed but arbitrary stationary data variable  $V$ . We write  $P_V$  for the load-recovery variable associated with  $V$ . Variable  $P_V$  serves to load the input elements of  $V$  into and recover the output elements of  $V$  from the respective cells. Thus, we must choose  $\vartheta_{P_V}$  such that  $P_V$  is moving.

Recall that the input data of  $V$  are initialised at the points of  $\text{in}(\Phi_V, \vartheta_V)$  and used for the first time at the points of  $\text{fst}(\Phi_V, \vartheta_V)$ , and the output data of  $V$  are



available at the points of  $\text{lst}(\Phi_V, \vartheta_V)$ . The construction of load-recovery variable  $P_V$  proceeds in two independent and symmetric steps: one step deals with the loading of the input data of  $V$  and the other step deals with the recovery of the output data of  $V$ .

Let us consider the input data first.  $\Phi_{P_V}^\ell$  contains the points across which the input data of  $V$  are loaded:

$$\Phi_{P_V}^\ell = \begin{cases} \text{if } \text{rank}(\sigma) = n-1 & \rightarrow \text{rays}(\Phi, \pm u) \cap \text{rays}(\text{fst}(\Phi_V, \vartheta_V) - \vartheta_V, -\vartheta_{P_V}) \\ \square \text{ rank}(\sigma) = 1 & \rightarrow \text{planes}(\Phi, \sigma) \cap \text{rays}(\text{fst}(\Phi_V, \vartheta_V) - \vartheta_V, -\vartheta_{P_V}) \\ \text{fi} \end{cases}$$

To determine the initialisation of variable  $P_V$  with the input data of  $V$ , we partition  $\text{fst}(\Phi_V, \vartheta_V)$  into two subsets: the set  $\Phi_V^i$  contains the points that are mapped to border cells and the set  $\text{fst}(\Phi_V, \vartheta_V) \setminus \Phi_V^i$  contains the points that are mapped to internal cells, where

$$\Phi_V^i = \begin{cases} \text{if } \text{rank}(\sigma) = n-1 & \rightarrow \{I \mid I \in \text{fst}(\Phi_V, \vartheta_V) \wedge I - \vartheta_{P_V} \notin \text{rays}(\Phi, \pm u)\} \\ \square \text{ rank}(\sigma) = 1 & \rightarrow \{I \mid I \in \text{fst}(\Phi_V, \vartheta_V) \wedge I - \vartheta_{P_V} \notin \text{planes}(\Phi, \sigma)\} \\ \text{fi} \end{cases}$$

The loading of the input data  $V(I - \vartheta_V)$  for  $I \in \Phi_V^i$  is unnecessary;  $P_V(I - \vartheta_{P_V})$  is initialised with the input value that  $V(I - \vartheta_V)$  is initialised with. The initialisation of  $P_V$  with the remaining input data is defined as follows. The two sets  $\text{in}(\Phi_{P_V}^\ell, \vartheta_{P_V})$  and  $\text{fst}(\Phi_V, \vartheta_V) \setminus \Phi_V^i$  are isomorphic by the bijection:

$$\ell_V : \text{in}(\Phi_{P_V}^\ell, \vartheta_{P_V}) \rightarrow \text{fst}(\Phi_V, \vartheta_V) \setminus \Phi_V^i, \quad \ell_V(I) = J, \quad \text{where } I \xrightarrow{\vartheta_{P_V}} J \quad (7.1)$$

$P_V(I)$  is initialised with the input value that  $V(J - \vartheta_V)$  is initialised with. It is then pipelined along  $\vartheta_{P_V}$  to the corresponding point of  $\text{lst}(\Phi_{P_V}^\ell, \vartheta_{P_V})$ .

Finally, we initialise variable  $V$  with its input data loaded via variable  $P_V$  by substituting  $P_V(I - \vartheta_{P_V})$  for all arguments  $V(I - \vartheta_V)$  for every  $I \in \text{fst}(\Phi_V, \vartheta_V)$ .

Next, we consider the recovery of the output data.  $\Phi_{P_V}^r$  contains the points across which the output data of  $V$  are recovered:

$$\Phi_{P_V}^r = \begin{cases} \text{if } \text{rank}(\sigma) = n-1 & \rightarrow \text{rays}(\Phi, \pm u) \cap \text{rays}(\text{lst}(\Phi_V, \vartheta_V) + \vartheta_V, \vartheta_{P_V}) \\ \square \text{ rank}(\sigma) = 1 & \rightarrow \text{planes}(\Phi, \sigma) \cap \text{rays}(\text{lst}(\Phi_V, \vartheta_V) + \vartheta_V, \vartheta_{P_V}) \\ \text{fi} \end{cases}$$



The recovery of the output data of  $V$  is symmetric to the initialisation of  $P_V$ . We partition  $\text{lst}(\Phi_V, \vartheta_V)$  into two subsets: the set  $\Phi_V^\circ$  contains the points that are mapped to border cells and the set  $\text{lst}(\Phi_V, \vartheta_V) \setminus \Phi_V^\circ$  contains the points that are mapped to internal cells, where

$$\Phi_V^\circ = \begin{cases} \text{if } \text{rank}(\sigma) = n-1 & \rightarrow \{I \mid I \in \text{lst}(\Phi_V, \vartheta_V) \wedge I + \vartheta_{P_V} \notin \text{rays}(\Phi, \pm u)\} \\ \text{fi} \\ \text{if } \text{rank}(\sigma) = 1 & \rightarrow \{I \mid I \in \text{lst}(\Phi_V, \vartheta_V) \wedge I + \vartheta_{P_V} \notin \text{planes}(\Phi, \sigma)\} \\ \text{fi} \end{cases}$$

The recovery of the output data  $V(I)$  for  $I \in \Phi_V^\circ$  is unnecessary. The recovery of the output data  $V(I)$  for  $I \in \text{lst}(\Phi_V, \vartheta_V) \setminus \Phi_V^\circ$  is defined as follows. The two sets  $\text{lst}(\Phi_{P_V}^r, \vartheta_{P_V})$  and  $\text{lst}(\Phi_V, \vartheta_V) \setminus \Phi_V^\circ$  are isomorphic by the bijection

$$r_V : \text{lst}(\Phi_{P_V}^r, \vartheta_{P_V}) \rightarrow \text{lst}(\Phi_V, \vartheta_V) \setminus \Phi_V^\circ, \quad r_V(I) = J, \quad \text{where } J \xrightarrow{\vartheta_{P_V}} I$$

$V(J)$  is first copied to  $P_V(J + \vartheta_{P_V})$ , which is then pipelined along  $\vartheta_{P_V}$  and finally copied to  $P_V(I)$ , which is mapped to a border cell.

The specification of stationary data variable  $V$  must be updated to reflect the change caused by the introduction of load-recovery variable  $P_V$  in the source UREs. This update and the specification of  $P_V$  are summarised in a procedure.

**Procedure 7.1** (Construction of the load-recovery variable  $P_V$  for stationary data variable  $V$  of the source UREs)

- The specification of load-recovery variable  $P_V$  is as follows:

$$\begin{aligned} I \in \text{in}(\Phi_V^i, \vartheta_{P_V}) & \rightarrow P_V(I) = v_{\text{in}_V(I + \vartheta_{P_V} - \vartheta_V)} \\ I \in \text{in}(\Phi_{P_V}^\ell, \vartheta_{P_V}) & \rightarrow P_V(I) = v_{\text{in}_V(\ell_V(I) - \vartheta_V)} \\ I \in \Phi_{P_V}^\ell \cup (\Phi_{P_V}^r \setminus \Phi_{P_V}^c) & \rightarrow P_V(I) = P_V(I - \vartheta_{P_V}) \\ I \in \Phi_{P_V}^c & \rightarrow P_V(I) = V(I - \vartheta_{P_V}) \\ I \in \text{lst}(\Phi_V^\circ, \vartheta_{P_V}) & \rightarrow v_{\text{out}_V(I)} = V(I) \\ I \in \text{lst}(\Phi_{P_V}^r, \vartheta_{P_V}) & \rightarrow v_{\text{out}_V(r_V(I))} = P_V(I) \end{aligned}$$

where  $\Phi_{P_V}^c = (\text{lst}(\Phi_V, \vartheta_V) \setminus \Phi_V^\circ) + \vartheta_{P_V}$ .

- Delete the input and output equations of variable  $V$ .
- Substitute  $P_V(I - \vartheta_{P_V})$  for all arguments  $V(I - \vartheta_V)$ , where  $I \in \text{fst}(\Phi_V, \vartheta_V)$ .  $\square$



In the rest of this section, we illustrate the construction of load-recovery variables with the example of matrix product. Let us consider S. Y. Kung's two-dimensional array for matrix product of Fig. 2-3. Its space-time mapping is:

$$\Pi = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad (7.2)$$

The projection vector is  $u=(0,0,1)$ . Variable  $C$  is stationary.

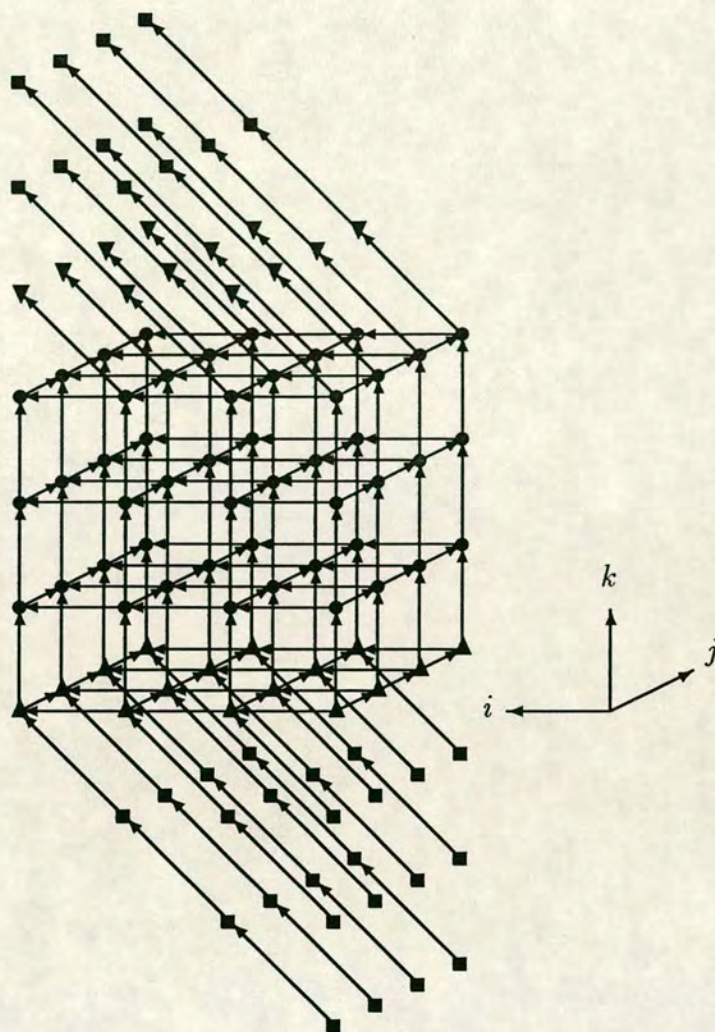
$P_C$  is the load-recovery variable associated with  $C$ . We choose  $\vartheta_{P_C}=(1,0,1)$ . An application of Proc. 7.1 rewrites the original source UREs for matrix product as described in Sect. 2.2 to the following set of equations:

*UREs with load-recovery variable  $P_C$ :*

$$\begin{aligned} A(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 = j \wedge 0 < k \leq m & \rightarrow a_{i,k} \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow A(i, j-1, k) \quad \blacktriangle \bullet \end{cases} \\ B(i, j, k) &= \begin{cases} 0 = i \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow b_{k,j} \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 0 < k \leq m & \rightarrow B(i-1, j, k) \quad \blacktriangle \bullet \end{cases} \\ P_C(i, j, k) &= \begin{cases} i = 0 \wedge 0 < j \leq m \wedge -m < k \leq 0 & \rightarrow 0 \\ (0 < i \wedge 0 < j \leq m \wedge i-m < k \leq 0) \vee \\ i \leq m \wedge 0 < j \leq m \wedge m+1 < k < i+m & \rightarrow P_C(i-1, j, k-1) \quad \blacksquare \\ 1 < i \leq m \wedge 0 < j \leq m \wedge k = m+1 & \rightarrow C(i-1, j, k-1) \quad \blacktriangledown \end{cases} \\ C(i, j, k) &= \begin{cases} 0 < i \leq m \wedge 0 < j \leq m \wedge k = 1 & \rightarrow P_C(i-1, j, k-1) \\ & + A(i, j-1, k) B(i-1, j, k) \quad \blacktriangle \\ 0 < i \leq m \wedge 0 < j \leq m \wedge 1 < k \leq m & \rightarrow C(i, j, k-1) \\ & + A(i, j-1, k) B(i-1, j, k) \quad \bullet \end{cases} \\ c_{2m-k,j} &= \begin{cases} i = m \wedge 0 < j \leq m \wedge k = m & \rightarrow C(i, j, k) \\ i = m \wedge 0 < j \leq m \wedge m < k < 2m & \rightarrow P_C(i, j, k) \end{cases} \end{aligned}$$

Fig. 7-1 depicts the corresponding data dependence graph. The cube depicts the original index space. The input elements of  $C$  are now initialised at the right facet of the lower triangle translated by  $-\vartheta_{P_C}$  and propagated towards the bottom facet of the cube. This implements the loading of  $C$ . Similarly, the output elements





**Figure 7-1:** The data dependence graph of the UREs with load-recovery variable  $P_C$  for matrix product ( $m=4$ ).

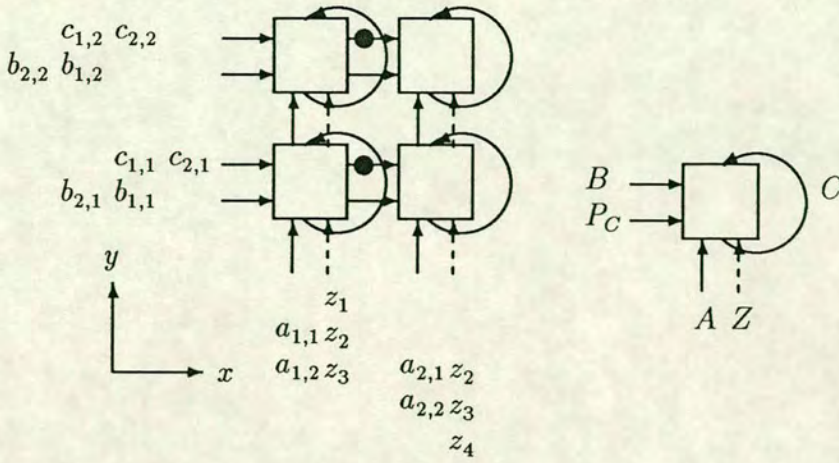
of  $C$  that are available at the top facet of the cube are now propagated towards the left facet of the top triangle. This implements the recovery of  $C$ .

Once the loading and recovering scheme is devised, the derivation of the corresponding control signals is a matter of construction of the control UREs. There are four  $\oplus$ -classes:

$$\Phi/\oplus = \{\Phi_{\blacksquare}, \Phi_{\bullet}, \Phi_{\blacktriangle}, \Phi_{\blacktriangledown}\}$$

PCUREs are unnecessary because of the particular choice of the projection vector:  $\Phi = \Psi$ . We construct the CCUREs by identifying a separation set  $\mathfrak{S} =$





**Figure 7-2:** The systolic array for matrix product that handles the loading and recovery of stationary variable  $C$  ( $m=2$ ). The distribution of the data and control signals is the snapshot at the first step.

$\{S_1, S_2, S_3, S_4\}$ :

$$\begin{aligned} S_1 &= \{(i, j, k) \mid k \leq 0\} \\ S_2 &= \{(i, j, k) \mid k \leq 1\} \\ S_3 &= \{(i, j, k) \mid k \leq m\} \\ S_4 &= \{(i, j, k) \mid k \leq m+1\} \end{aligned}$$

Since the corresponding hyperplanes of these three half-spaces are parallel, we can combine the three control variables associated with these half-spaces into one, namely  $Z$ .  $\text{sig}(Z) = \{z_1, z_2, z_3, z_4\}$ . The CCUREs follow from (3.8) and (3.23):

$$(0, 0, 1)\vartheta_Z = 0$$

$$Z(I) = \begin{cases} I \in \text{in}(\Psi, \vartheta_Z) \cap (\cap S_1 \cup \complement_{\mathbf{Z}^n} S_4) & \rightarrow z_1 \\ I \in \text{in}(\Psi, \vartheta_Z) \cap \complement_{\mathbf{Z}^n} S_1 \cap S_2 & \rightarrow z_2 \\ I \in \text{in}(\Psi, \vartheta_Z) \cap \complement_{\mathbf{Z}^n} S_2 \cap S_3 & \rightarrow z_3 \\ I \in \text{in}(\Psi, \vartheta_Z) \cap \complement_{\mathbf{Z}^n} S_3 \cap S_4 & \rightarrow z_4 \\ I \in \Psi & \rightarrow Z(I - \vartheta_Z) \end{cases}$$

The replacement of the domain predicates by predicates in control variable  $Z$  is straightforward and omitted. Let us choose  $\vartheta_Z = (0, 1, 0)$ . The previous space-time mapping (7.2) describes the systolic array shown in Fig. 7-2. This array runs



in  $5m - 2$  steps;  $2m + 1$  extra steps are used to load and recover the elements of stationary variable  $C$ . Note that the loading (recovery) of some elements of  $C$  overlaps with the computations of other elements of  $C$ .

**Remark** The search of loading and recovery schemes reduces to a search for dependence vectors of load-recovery variables. Consider the previous example. If we choose  $\vartheta_{P_C} = (2, 0, -1)$  instead, the latency of the original array (which is  $3m + 2$  execution steps) is retained. That is, the loading and recovery of  $C$  overlaps completely with the computation of  $C$ . The price paid is that the loading and recovery of  $C$  requires non-neighbouring channel connections (albeit without delay buffers, since  $\lambda\vartheta_{P_C} = 1$ ).  $\square$

### 7.3 The Access of Stationary Data

Address variables provide the addresses for accessing the stationary elements stored in the local memory of a cell. They can be dispensed with in  $(n - 1)$ -dimensional arrays, because at most one element per variable can be mapped to a cell. (Recall that the systolic arrays we consider are parameterised systolic arrays (Sect. 2.7).) The size of the local processor memory associated with a stationary variable is at most 1. The access of the corresponding stationary element – if any – can be predesigned in hardware. However, address variables are essential for one-dimensional arrays, because more than one element per variable can be mapped to the local memory of a cell. We need to supply the right address at the right cell at the right step so that the right stationary element can be accessed. We shall present the construction of address variables for one-dimensional arrays only, i.e., we only consider the allocation vector  $\sigma$ .

Similarly to the construction of load-recovery variables, we associate a distinct address variable with every stationary data variable of the source UREs. Again, pick a fixed but arbitrary stationary data variable  $V$ . We write  $X_V$  for the address variable associated with  $V$ . We must choose  $\vartheta_{X_V}$  such that  $X_V$  is moving. Essen-



tially, variable  $X_V$  serves to supply the addresses for the access of the elements of  $V$  stored in the local processor memories of the array.

Assume that  $m$  elements of  $V$  are mapped to a cell. Thus,  $m$  different  $\vartheta_V$ -paths are mapped to that cell. A stationary element is successively computed at the points of the corresponding  $\vartheta_V$ -path. Therefore, all the points in one  $\vartheta_V$ -path are mapped to the same memory location. (Of course, different stationary elements mapped to a cell must be stored at different memory locations.)  $X_V(I)$  must be the same for all points  $I$  of a  $\vartheta_V$ -path. The construction of address variable  $X_V$  is as follows. We choose an  $n$ -vector  $\pi_{X_V}$  and, thus, obtain a sequence of parallel hyperplanes  $\text{addr}_{\ell_{\min}^V}, \text{addr}_{\ell_{\min}^V+1}, \dots, \text{addr}_{\ell_{\max}^V}$  ( $\ell_{\min}^V \leq \ell_{\max}^V$ ) that intersect the domain of stationary data variable  $V$ . We choose the control dependence vector  $\vartheta_{X_V}$  of address variable  $X_V$  such that  $\vartheta_{X_V} \pi_{X_V} = 0$  and construct  $X_V$  such that  $X_V(I) = \ell$  for all points  $I$  in hyperplane  $\text{addr}_{\ell}$ . Therefore, we must choose  $\pi_{X_V}$  such that

- $\pi_{X_V}$  and  $\sigma$  are not co-linear, and
- $\pi_{X_V} \vartheta_V = 0$ .

The first condition ensures that  $X_V$  is a moving variable. The second condition ensures that  $X_V(I)$  is the same address for all points of a  $\vartheta_V$ -path.

The specification of address variable  $X$  is as follows:

$$X_V(I) = \begin{cases} I \in \text{in}(\Phi_V, \vartheta_{X_V}) \cap \text{addr}_{\ell_{\min}^V} & \rightarrow \ell_{\min}^V \\ I \in \text{in}(\Phi_V, \vartheta_{X_V}) \cap \text{addr}_{\ell_{\min}^V+1} & \rightarrow \ell_{\min}^V + 1 \\ \dots & \dots \\ I \in \text{in}(\Phi_V, \vartheta_{X_V}) \cap \text{addr}_{\ell_{\max}^V} & \rightarrow \ell_{\max}^V \\ I \in \Phi_V & \rightarrow X_V(I - \vartheta_{X_V}) \end{cases} \quad (7.3)$$

According to the specification of address variable  $X_V$ , different elements of  $V$  that are contained in the same hyperplane  $\text{addr}_{\ell}$  are necessarily projected to different cells. They are stored at the same location  $\ell$  at the respective cells. This ensures that the address variable  $X_V$  can be pipelined across these cells.



The specification of stationary data variable  $V$  must be updated to reflect the change caused by the introduction of address variable  $X_V$  in the source UREs. This update and the specification of  $X_V$  are summarised in a procedure.

**Procedure 7.2** (Construction of the address variable  $X_V$  for stationary data variable  $V$  of the source UREs)

- The specification of address variable  $X_V$  is given by (7.3).
- Replace all  $V(I)$  in the source UREs by  $V(X_V(I))$ . □

In the rest of this section, we use dynamic programming as an example. Let us consider the one-dimensional array of Sect. 5.6.2. Its space-time mapping is:

$$\pi = \begin{bmatrix} -4 & 2(n+2) & -2 \\ -1 & 1 & 0 \end{bmatrix} \quad (7.4)$$

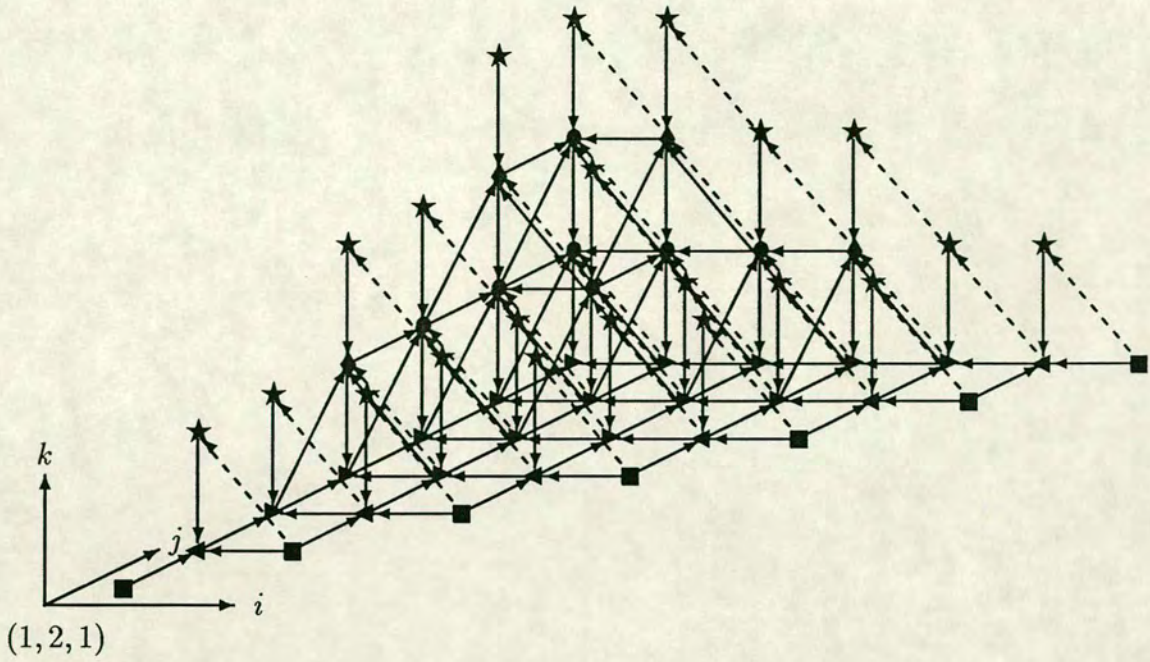
Variable  $C$  is stationary. The specification of the loading and recovery of the stationary elements of  $C$  follows from Proc. 7.1.

$X_C$  is the address variable associated with stationary variable  $C$ . We choose  $\vartheta_{X_C} = (-1, 0, 1)$ . ( $\pi_{X_C} = (0, 1, 0)$ ). A simple calculation shows that  $\ell_{\min}^C = 2$  and  $\ell_{\max}^C = m+1$ . An application of Proc. 7.2 rewrites the source UREs for dynamic programming, as described in Sect. 3.7.1, to the following set of equations:

*UREs with address variable  $X_C$ :*

$$\begin{aligned}
 X_C(i, j, k) &= \begin{cases} 1 < i \leq j \wedge 2 \leq j \leq m+1 \wedge k=0 & \rightarrow j \\ 0 < i < j \leq m+1 \wedge 0 < k \leq (j-i)/2+1 & \rightarrow X_C(i+1, j, k-1) \end{cases} \\
 & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \blacktriangleleft \blacktriangleright \blacklozenge \bullet \blacksquare \star \\
 A(i, j, k) &= \begin{cases} j-i=2k \rightarrow D(i, j, k) & \blacktriangleleft \blacklozenge \\ j-i \neq 2k \rightarrow A(i, j-1, k) & \blacktriangleright \bullet \blacksquare \star \end{cases} \\
 B(i, j, k) &= \begin{cases} k=1 \rightarrow C(X_C(i+1, j, k)) & \blacktriangleleft \blacktriangleright \blacksquare \\ k \neq 1 \rightarrow B(i+1, j, k-1) & \blacklozenge \bullet \star \end{cases}
 \end{aligned}$$





**Figure 7-3:** The data dependence graph of the UREs with address variable  $X_C$  for dynamic programming ( $m=7$ ). The dashed vectors depict  $\vartheta_{X_C}$ .

$$\begin{aligned}
 D(i, j, k) &= \begin{cases} k=1 \rightarrow C(X_C(i, j-1, k)) & \blacktriangleleft \blacktriangleright \blacksquare \\ k \neq 1 \rightarrow D(i, j-1, k-1) & \blacklozenge \bullet \star \end{cases} \\
 E(i, j, k) &= \begin{cases} j-i=2k \rightarrow B(i, j, k) & \blacktriangleleft \blacklozenge \\ j-i \neq 2k \rightarrow E(i+1, j, k) & \blacktriangleright \bullet \blacksquare \star \end{cases} \\
 c_{1, m+1} &= C(X_C(1, m+1, 1)) \\
 C(X_C(i, j, k)) &= \begin{cases} k=1 \wedge j-i=2k-1 \rightarrow w_{i,j} & \blacksquare \\ k=1 \wedge j-i \geq 2k \rightarrow \min \begin{pmatrix} A(i, j, k) + B(i, j, k) \\ C(X_C(i, j, k+1)) \\ D(i, j, k) + E(i, j, k) \end{pmatrix} + w_{i,j} & \blacktriangleleft \blacktriangleright \\ k \neq 1 \wedge j-i \geq 2k \rightarrow \min \begin{pmatrix} A(i, j, k) + B(i, j, k) \\ C(X_C(i, j, k+1)) \\ D(i, j, k) + E(i, j, k) \end{pmatrix} & \blacklozenge \bullet \\ k \neq 1 \wedge j-i < 2k \wedge j-i \geq 2k-2 \rightarrow \infty & \star \end{cases}
 \end{aligned}$$



Note that  $\Phi_C = \Phi$ . The data dependence graph for this new system of UREs is the data dependence graph for the original source UREs, as shown in Sect. 3.7.1, with data dependence vector  $\vartheta_{X_C}$  replicated uniformly across the index space (Fig. 7-3, compare Fig. 3-4). The CCUREs are specified in Sect. 3.7.1.2. The PCUREs are specified in Sect. 5.6.2. Under the previous space-time mapping (7.4), address variable  $X_C$  encounters a delay of one unit along a neighbouring channel (in the  $\chi$ -model). This control variable was derived informally in [62].

**Remark** Similarly to the specification of load-recovery variables, the search of schemes for the access of stationary data variables reduces to a search of dependence vectors for address variables.  $\square$

## 7.4 Conclusion

We have presented a systematic method for the construction of the specification of loading, recovery and access of stationary variables. The method was presented for  $(n - 1)$ -dimensional and one-dimensional arrays synthesised from UREs. It generalises directly to systolic arrays of  $r$  dimensions ( $0 < r < n$ ) synthesised from systolic algorithms that are defined over a domain that can be expressed by a convex set of integer coordinates.

The advantage of our methods on handling stationary variables is that they are only dependent on the place function but independent of the step function. That leaves plenty of freedom of choice for control dependence vectors for load-recovery and address variables and the scheduling vector.



## Chapter 8

### Conclusion

We have shown how control signals for systolic arrays can be systematically derived by program transformation and construction. The method we propose applies for systolic algorithms defined over a convex set of integer coordinates. The basic idea is to transform the domain predicates in the initial program to UREs that we call control UREs. Then, the search of systolic arrays with a description of both data and control signals reduces to a search of space-time mappings.

In Chap. 2, we have reviewed the basic techniques for the synthesis of data flow for  $(n - 1)$ -dimensional arrays. In Chap. 3, we have presented our method for the construction of control UREs from  $n$ -dimensional UREs. We have provided necessary and sufficient conditions for the correctness of the control UREs and a mechanisable procedure that constructs the control UREs from the source UREs. We have provided some heuristics for guiding the user of interactive systolic design systems in the optimisation of the domain predicates (optimised domain predicates lead to a more efficient control flow). We have also shown how redundant control hardware can be completely removed once the space-time mapping has been selected. In Chap. 4, we have reviewed the synthesis of data flow for one-dimensional arrays. We have improved and extended previous results. We have provided unified mapping conditions for the validity of space-time mappings for two one-dimensional array models. We have also presented a variety of equivalent mapping conditions that highlight the properties of one-dimensional arrays from different perspectives. We have concluded that the control UREs constructed for



$(n-1)$ -dimensional arrays can be inefficient for one-dimensional arrays. In Chap. 5, we have provided a more appropriate construction. In Chap. 6, we have shown how to eliminate the control signals for a class of systolic arrays by exploiting the algebraic properties of some operators in the source UREs. Finally, in Chap. 7, we have described a systematic method for the construction of UREs that specify loading, recovery and access of stationary elements in the local processor memories of a systolic array.

There are several areas that require further research. In Chap. 3, we have presented a procedure that delivers the control UREs from the source UREs. The control UREs returned by this procedure are completely determined by the specification of the domain predicates of the source UREs. We have shown by example that different specifications of the domain predicates of the source UREs lead to different specifications of the control UREs. It is worthwhile searching for a method for finding the specifications of the domain predicates that lead to (in some sense) optimal control signals.

Once the control UREs are constructed, there are many possible control dependence vectors. The choice of control dependence vector affects the quality of the systolic array. This thesis does not provide a systematic way of making a good choice.

Finally, the control UREs are targeted directly at parameterised systolic arrays. In reality, we need to map both the source and the control UREs to fixed-size arrays (i.e., partitioned arrays). This can be done in three successive steps. First, we obtain a parameterised systolic array by the space-time mapping technique. Second, we partition the array [11,17,52,82]. There is no need for distinguishing which recurrence equation specifies data or control signals. The partitioning demands additional control, which ensures that the right values (data and control signals) as specified by the source and the control UREs arrive at the right cells in the partitioned array at the right steps. This is the task of the third step. The regularity exhibited in the partitioning can be exploited in the derivation of this additional control. This regularity has been recently captured by a more general form of **step** and **place** that contains **mod** or **div** operators [17,82].



## Appendix A

### The Normalisation of Domain Predicates

Proc. 3.1 for the construction of separation set  $\mathfrak{S}_{\text{src}}$  expects that domain predicates are in normal form. This appendix describes a procedure for the normalisation of certain more general domain predicates. Tab. A-1 gives a syntactic definition of these domain predicates in terms of a set-theoretic abstract syntax [68]. The domain predicates defined in Tab. A-1 are distinguished from the domain predicates in normal form in that they may also contain the six operators in the following set:

$$O = \{ |, \max, \min, [ ], [ ], \text{div} \} \quad (\text{A.1})$$

These six operators are the most frequently used in practice. For example, the URE specification for dynamic programming of [14] contains operator  $[ ]$ , the URE specification for transitive closure of [72] contains operators  $\max$  and  $\min$ , and the URE specification for LU-decomposition of [44] contains operator  $\min$ .

We specify a language's syntax by listing its syntax domains and its BNF rules. In Tab. A-1, the phrase  $D \in \text{Domain-Predicate}$  indicates that  $D$  is the non-terminal that represents an arbitrary member of the syntax domain Domain-Predicate. The structure of domain predicates is given by the BNF rule  $D ::= C; (D); D_1 \wedge D_2; D_1 \vee D_2$  which says that any domain predicate must be either a conditional, or a conditional enclosed in a pair of parentheses, or a conjunction of two conditionals or a disjunction of two conditionals.

An expression  $E$  is called a *basic* expression if  $E$  is of the form  $\pi I + \delta$ , i.e., an affine expression of index vector  $I$  in  $\mathbf{Z}^n$ . A conditional  $E @ 0$  is called a *basic* conditional



$D \in \text{Domain-Predicate}$
$C \in \text{Conditional}$
$E \in \text{Expression}$
$A \in \{\pi I + \delta \mid \pi \in \mathbf{Z}^n \wedge I \in \mathbf{Z}^n \wedge \delta \in \mathbf{Z}\}$
$n \in \mathbf{Z}^+$
$D ::= C; (D); D_1 \wedge D_2; D_1 \vee D_2$
$C ::= E \textcircled{R} 0$
$E ::= L;  E ;$
$E \textcircled{\otimes} L; \max(E, L); \min(E, L); \lceil (E \textcircled{\otimes} L) / n \rceil; \lfloor (E \textcircled{\otimes} L) / n \rfloor; (E \textcircled{\otimes} L) \text{ div } n$
$L \textcircled{\otimes} E; \max(L, E); \min(L, E); \lceil (L \textcircled{\otimes} E) / n \rceil; \lfloor (L \textcircled{\otimes} E) / n \rfloor; (E \textcircled{\otimes} L) \text{ div } n$
$L ::= A; nL;  L ; \max(L_1, L_2); \min(L_1, L_2)$
$\textcircled{R} ::= <; \leq; >; \geq; =; \neq$
$\textcircled{\otimes} ::= +; -$

**Table A-1:** The abstract syntax of the domain predicates.

if  $E$  is a basic expression. Conditionals of the form  $E=0$  are called *equalities*, those of the form  $E \neq 0$  *inequalities*. A domain predicate is a *basic* domain predicate if all its constituent conditionals are basic conditionals. A domain predicate is said to be in normal form if it is a basic domain predicate. For convenience, we refer to anything that is not a domain predicate as a *non-domain predicate*.

Domain predicates exhibit a nested parentheses property. To see this, let us make a syntactic substitution in rule E: we replace  $(E \textcircled{\otimes} L) \text{ div } n$  by  $\uparrow E \textcircled{\otimes} L, n \uparrow$  and  $(L \textcircled{\otimes} E) \text{ div } n$  by  $\uparrow L \textcircled{\otimes} E, n \uparrow$ . If we delete all symbols except those in  $\{[, \lfloor, \lceil, \rfloor, \rceil, \uparrow\}$  from a domain predicate and let  $n$  be the length of the resulting string, then the  $i$ -th symbol ( $0 < i \leq n/2$ ) from the left forms an operator (either  $\lceil \ ]$  or  $\lfloor \ ]$  or  $\uparrow \uparrow$ ) with the  $i$ -th symbol from the right.

**Example A.1** These examples illustrate the structure of domain predicates.

*Domain predicates:*

1.  $\lfloor i/2 \rfloor - j \leq 0$
2.  $|i+2| - k \neq 0$



3.  $(i+j) \text{ div } 2 - k > 0$
4.  $\max(i, j) + k = 0$
5.  $\min(\lfloor i/2 \rfloor, k) - j \geq 0$

*Non-domain predicates:*

1.  $2\lfloor i/2 \rfloor - j = 0$
2.  $2\lceil (i \text{ div } 2 + j/3) \rceil + 1 \geq 0$
3.  $i \text{ div } 2 + \lceil (j+k)/2 \rceil = 0$
4.  $\lfloor i/2 \rfloor - \lceil (j+k)/3 \rceil < 0$
5.  $\max(\lceil i/2 \rceil, \lceil i/2 \rceil + \lfloor j/2 \rfloor) \neq 0$  □

The normalisation of domain predicates relies on the notion of rewriting. In general, it is achieved according to a number of pre-defined rewriting rules. A rewriting rule rewrites a domain predicate of a certain form to a domain predicate of a different form. The following lemma states algebraic laws of the six operators in the set  $O$  of (A.1). It is the basis of the formulation of the rewriting rules in Lemma A.2.

**Lemma A.1** *Let  $u, v, w \in \mathbb{R}$ ,  $x, y \in \mathbb{Z}$ ,  $c \in \mathbb{Z}^+$ .*

$$\begin{aligned}
 \max: \quad & w \textcircled{\mathbb{R}} \max(u, v) \iff (w \textcircled{\mathbb{R}} u \wedge u \geq v) \vee (v \geq u \wedge w \textcircled{\mathbb{R}} v) \\
 \min: \quad & w \textcircled{\mathbb{R}} \min(u, v) \iff (w \textcircled{\mathbb{R}} u \wedge u \leq v) \vee (v \leq u \wedge w \textcircled{\mathbb{R}} v) \\
 | |: \quad & |u| \textcircled{\mathbb{R}} v \iff (u \textcircled{\mathbb{R}} v \wedge u \geq 0) \vee (-u \textcircled{\mathbb{R}} v \wedge u \leq 0) \\
 \lceil \rceil: \quad & \lceil y/c \rceil \textcircled{\mathbb{R}} x \iff (y \textcircled{\mathbb{R}} cx \vee y \textcircled{\mathbb{R}} cx - 1 \vee \dots \vee y \textcircled{\mathbb{R}} cx - c + 1) \\
 \lfloor \rfloor: \quad & \lfloor y/c \rfloor \textcircled{\mathbb{R}} x \iff (y \textcircled{\mathbb{R}} cx \vee y \textcircled{\mathbb{R}} cx + 1 \vee \dots \vee y \textcircled{\mathbb{R}} cx + c - 1) \\
 \text{div}: \quad & (y \text{ div } c) \textcircled{\mathbb{R}} x \iff (x \geq 0 \wedge (y \textcircled{\mathbb{R}} cx \vee y \textcircled{\mathbb{R}} cx + 1 \vee \dots \vee y \textcircled{\mathbb{R}} cx + c - 1) \vee \\
 & (x \leq 0 \wedge (y \textcircled{\mathbb{R}} cx \vee y \textcircled{\mathbb{R}} cx - 1 \vee \dots \vee y \textcircled{\mathbb{R}} cx - c + 1))
 \end{aligned}$$

**Proof** Algebraic manipulation. □

The following lemma states the rewriting rules that are used in the normalisation of domain predicates. The notation  $\overline{\textcircled{\mathbb{R}}}$  in rule  $| |$ -elim1 has the following meaning: the overbar denotes the following mapping (whose argument is  $\textcircled{\mathbb{R}}$ ) from



$\{<, \leq, >, \geq, =, \neq\}$  onto itself:

$$\overline{\mathbb{R}} = \left\{ \begin{array}{l} \text{if } \mathbb{R} = < \rightarrow > \\ \square \quad \mathbb{R} = \leq \rightarrow \geq \\ \square \quad \mathbb{R} = > \rightarrow < \\ \square \quad \mathbb{R} = \geq \rightarrow \leq \\ \square \quad \mathbb{R} = = \rightarrow = \\ \square \quad \mathbb{R} = \neq \rightarrow \neq \\ \text{fi} \end{array} \right.$$

### Lemma A.2

$$\text{0-law: } E\mathbb{R}0 \iff E+0\mathbb{R}0$$

$$\text{max-dist: } n \max(L_1, L_2)\mathbb{R}0 \iff \max(nL_1, nL_2)\mathbb{R}0$$

$$\text{min-dist: } n \min(L_1, L_2) \iff \min(nL_1, nL_2)\mathbb{R}0$$

$$| \text{-law: } n|L|\mathbb{R}0 \iff |nL|\mathbb{R}0$$

$$\text{max-elim1: } \max(E, L_1) \pm L_2 \mathbb{R}0 \iff (E \pm L_2 \mathbb{R}0 \wedge E - L_1 \geq 0) \vee (L_1 - E \geq 0 \wedge L_1 \pm L_2 \mathbb{R}0)$$

$$\text{max-elim2: } L_2 \pm \max(E, L_1) \mathbb{R}0 \iff (L_2 \pm E \mathbb{R}0 \wedge E - L_1 \geq 0) \vee (L_1 - E \geq 0 \wedge L_2 \pm L_1 \mathbb{R}0)$$

$$\text{min-elim1: } \min(E, L_1) \pm L_2 \mathbb{R}0 \iff (E \pm L_2 \mathbb{R}0 \wedge E - L_1 \leq 0) \vee (L_1 - E \leq 0 \wedge L_1 \pm L_2 \mathbb{R}0)$$

$$\text{min-elim2: } L_2 \pm \min(E, L_1) \mathbb{R}0 \iff (L_2 \pm E \mathbb{R}0 \wedge E - L_1 \leq 0) \vee (L_1 - E \leq 0 \wedge L_2 \pm L_1 \mathbb{R}0)$$

$$| \text{-elim1: } |E| \pm L \mathbb{R}0 \iff (E \pm L \mathbb{R}0 \wedge E \geq 0 \vee E \mp L \overline{\mathbb{R}}0 \wedge E \leq 0)$$

$$| \text{-elim2: } L \pm |E| \mathbb{R}0 \iff (L \pm E \mathbb{R}0 \wedge E \geq 0 \vee L \mp E \mathbb{R}0 \wedge E \leq 0)$$

$$[ \ ]\text{-elim1: } [E/n] \pm L \mathbb{R}0 \iff E \pm nL \mathbb{R}0 \vee E \pm nL + 1 \mathbb{R}0 \vee \dots \vee E \pm nL + n - 1 \mathbb{R}0$$

$$[ \ ]\text{-elim2: } L \pm [E/n] \mathbb{R}0 \iff nL \pm E \mathbb{R}0 \vee nL + 1 \pm E \mathbb{R}0 \vee \dots \vee nL + n - 1 \pm E \mathbb{R}0$$

$$[ \ ]\text{-elim1: } [E/n] \pm L \mathbb{R}0 \iff E \pm nL \mathbb{R}0 \vee E \pm nL - 1 \mathbb{R}0 \vee \dots \vee E \pm nL - n + 1 \mathbb{R}0$$

$$[ \ ]\text{-elim2: } L \pm [E/n] \mathbb{R}0 \iff nL \pm E \mathbb{R}0 \vee nL - 1 \pm E \mathbb{R}0 \vee \dots \vee nL - n + 1 \pm E \mathbb{R}0$$

$$\text{div-elim1: } E \text{ div } n \pm L \mathbb{R}0 \iff$$

$$(\pm L \geq 0 \wedge (E \pm nL \mathbb{R}0 \vee E \pm nL + 1 \mathbb{R}0 \vee \dots \vee E \pm nL + n - 1 \mathbb{R}0)) \vee$$

$$(\pm L \leq 0 \wedge (E \pm nL \mathbb{R}0 \vee E \pm nL - 1 \mathbb{R}0 \vee \dots \vee E \pm nL - n + 1 \mathbb{R}0))$$

$$\text{div-elim2: } L \pm E \text{ div } n \mathbb{R}0 \iff$$

$$(L \geq 0 \wedge (nL \pm E \mathbb{R}0 \vee nL + 1 \pm E \mathbb{R}0 \vee \dots \vee nL + n - 1 \pm E \mathbb{R}0)) \vee$$

$$(L \leq 0 \wedge (nL \pm E \mathbb{R}0 \vee nL - 1 \pm E \mathbb{R}0 \vee \dots \vee nL - n + 1 \pm E \mathbb{R}0))$$



**Proof** The first rule follows from the fact that 0 is the zero element in  $\mathbb{Z}$ . The next three rewriting rules follow from the conventional semantics of max, min, and  $\lfloor \cdot \rfloor$ . The remaining rewriting rules follow from Lemma A.1.  $\square$

To illustrate the rewriting rules in Lemma A.2, we normalise the domain predicates of Ex. A.1.

1.  $\lfloor i/2 \rfloor - j \leq 0$   
 $\iff \{ \lfloor \cdot \rfloor\text{-elim1} \}$   
 $i - 2j \leq 0 \vee i - 2j + 1 \leq 0$
2.  $|i+2| - k \neq 0$   
 $\iff \{ | \cdot | \text{-elim1} \}$   
 $(i+2-k \neq 0 \wedge i+2 \geq 0) \vee (-i-2-k \neq 0 \wedge i+2 \leq 0)$
3.  $(i+j) \text{ div } 2 - k > 0$   
 $\iff \{ \text{div-elim1} \}$   
 $(-k \geq 0 \wedge (i+j-2k > 0 \vee i+j-2k+1 > 0)) \vee$   
 $(-k \leq 0 \wedge (i+j-2k > 0 \vee i+j-2k-1 > 0))$
4.  $\max(i, j) + k = 0$   
 $\iff \{ \text{max-elim1} \}$   
 $(i+k=0 \wedge i-j \geq 0) \vee (j-i \geq 0 \wedge j+k=0)$
5.  $\min(\lfloor i/2 \rfloor, k) - j \geq 0$   
 $\iff \{ \text{min-elim1} \}$   
 $(\lfloor i/2 \rfloor - j \geq 0 \wedge \lfloor i/2 \rfloor - k \leq 0) \vee (k - \lfloor i/2 \rfloor \leq 0 \wedge k - j \geq 0)$   
 $\{ \text{Apply } \lfloor \cdot \rfloor \text{-elim1 for the first two conditionals and } \lfloor \cdot \rfloor \text{-elim2 for the third} \}$   
 $\iff ((i-2j \geq 0 \vee i-2j-1 \geq 0) \wedge (i-2k \leq 0 \vee i-2k-1 \leq 0)) \vee$   
 $(2k-i \leq 0 \vee 2k-1-i \leq 0) \wedge k-j \geq 0$

**Theorem A.1** *Any domain predicate can be put into normal form.*

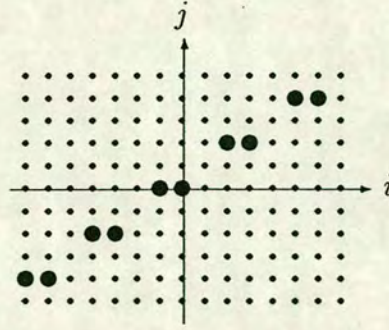
**Proof** It suffices to show that the theorem holds for conditionals of the form  $E \textcircled{R} 0$ . For reasons of symmetry, we shall not consider expressions that are of the last seven forms of rule E. The proof is conducted by a structural induction on the



number of the operators in the set  $O$  of (A.1) in a conditional. Assume that the theorem holds for conditionals that contain less than  $m$  such operators. Let  $E \textcircled{R} 0$  be a conditional that contains  $m$  such operators.  $E$  can be put into the following forms by rule E:

- $\max(E, L) \textcircled{R} 0$ . Apply 0-law and max-elim1.
- $\min(E, L) \textcircled{R} 0$ . Apply 0-law and min-elim1.
- $|E| \textcircled{R} 0$ . Apply 0-law and  $|$ -elim1.
- $(E \odot L) \text{div } n \textcircled{R} 0$ . Apply 0-law and div-elim1.
- $[(E \odot L)/n] \textcircled{R} 0$ . Apply 0-law and  $[$ ]-elim1.
- $[(E \odot L)/n] \textcircled{R} 0$ . Apply 0-law and  $]$ -elim1.
- $L \textcircled{R} 0$ .  $L$  can be put into the following forms by rule L:
  - $A \textcircled{R} 0$ . It is already in normal form.
  - $nL \textcircled{R} 0$ . Each domain predicate is finite. We apply rule L recursively until we are in one of the other four cases.
  - $n \max(L_1, L_2) \textcircled{R} 0$ . Apply max-dist, 0-law and max-elim1.
  - $n \min(L_1, L_2) \textcircled{R} 0$ . Apply min-dist, 0-law and min-elim1.
  - $n|L| \textcircled{R} 0$ . Apply  $|$ -law, 0-law and  $|$ -elim1.
- $E \odot L \textcircled{R} 0$ . Again,  $E$  can be put into the following forms by rule E.
  - $\max(E, L) \odot L \textcircled{R} 0$ . Apply max-elim1.
  - $\min(E, L) \odot L \textcircled{R} 0$ . Apply min-elim1.
  - $|E| \odot L \textcircled{R} 0$ . Apply  $|$ -elim1.
  - $(E \odot L_1) \text{div } n \odot L \textcircled{R} 0$ . Apply div-elim1.
  - $[(E \odot L_1)/n] \odot L \textcircled{R} 0$ . Apply  $[$ ]-elim1.





**Figure A-1:** The points  $(i, j, 1)$  satisfying predicate  $2\lfloor i/2 \rfloor - j = 0$ .

- $\lfloor (E \odot L_1) / n \rfloor \odot L \mathbb{R} 0$ . Apply  $\lfloor \rfloor$ -elim1.
- $L_1 \odot L \mathbb{R} 0$ .  $L_1$  can be put into the following forms by rule L:
  - \*  $A \odot L \mathbb{R} 0$ . L can be put into the following forms by rule L:
    - $A \odot A_1 \mathbb{R} 0$ . It is already in normal form.
    - $A \odot nL_1 \mathbb{R} 0$ . Consider rule L for  $L_1$  recursively.
    - $A \odot \max(L_2, L_3) \mathbb{R} 0$ . Apply max-elim2.
    - $A \odot \min(L_2, L_3) \mathbb{R} 0$ . Apply min-elim2.
    - $A \odot |L| \mathbb{R} 0$ . Apply  $| \rfloor$ -elim2.
  - \*  $nL_2 \odot L \mathbb{R} 0$ . Consider rule L for  $L_2$  recursively.
  - \*  $\max(L_2, L_3) \odot L \mathbb{R} 0$ . Apply max-elim1.
  - \*  $\min(L_2, L_3) \odot L \mathbb{R} 0$ . Apply min-elim1.
  - \*  $|L_1| \odot L \mathbb{R} 0$ . Apply  $| \rfloor$ -elim1.

In each of the above cases, the resulting predicate is still a domain predicate, each of whose constituent conditionals contains less than  $m$  operators in the set  $O$  of (A.1). The proof is completed by an application of the induction hypothesis.  $\square$

Finally, we discuss informally how non-domain predicates can be put in normal form. Let us consider non-domain predicate  $2\lfloor i/2 \rfloor - j = 0$  given in Ex. A.1. Rule  $\lfloor \rfloor$ -elim1 does not apply. Fig. A-1 exposes geometrically the difficulty in normalising this predicate. The points in any of the two lines are not consecutive. They cannot be specified by a finite number of basic conditionals in  $i$  and  $j$ . Algebraically, predicate  $2\lfloor i/2 \rfloor - j = 0$  implies the evenness of  $j$ .



A non-domain predicate can be put in normal form if it can be rewritten as a domain predicate. The motivation for this transformation is to reuse the indices that do not appear in the non-domain predicate or to introduce extra indices. For example, letting  $\lfloor i/2 \rfloor = k$  rewrites the previous example predicate to domain predicate  $\lfloor i/2 \rfloor - k = 0 \wedge j - 2k = 0$ , where  $j - 2k = 0$  specifies the evenness of  $j$ . By the same token, we rewrite the other four non-domain predicates in Ex. A.1 to the following domain predicates:

$$\begin{aligned} 2. \quad & 2\lfloor (i \operatorname{div} 2 + j)/3 \rfloor + 1 \geq 0 \\ \iff & \{\text{Let } k = \lfloor (i \operatorname{div} 2 + j)/3 \rfloor\} \\ & 2k + 1 = 0 \wedge k - \lfloor (i \operatorname{div} 2 + j)/3 \rfloor = 0 \end{aligned}$$

$$\begin{aligned} 3. \quad & i \operatorname{div} 2 + \lfloor (j + k)/2 \rfloor = 0 \\ \iff & \{\text{Let } \ell = \lfloor (j + k)/2 \rfloor\} \\ & i \operatorname{div} 2 + \ell = 0 \wedge \ell - \lfloor (j + k)/2 \rfloor = 0 \end{aligned}$$

$$\begin{aligned} 4. \quad & \lfloor i/2 \rfloor - \lfloor (j + k)/3 \rfloor < 0 \\ \iff & \{\text{Let } \ell = \lfloor i/2 \rfloor\} \\ & \ell - \lfloor (j + k)/3 \rfloor < 0 \wedge \ell - \lfloor i/2 \rfloor = 0 \end{aligned}$$

$$\begin{aligned} 5. \quad & \max(\lfloor i/2 \rfloor, \lfloor i/2 \rfloor + \lfloor j/2 \rfloor) \neq 0 \\ \iff & \{\text{Let } k = \lfloor i/2 \rfloor + \lfloor j/2 \rfloor\} \\ & \max(\lfloor i/2 \rfloor, k) \neq 0 \wedge k - \lfloor i/2 \rfloor - \lfloor j/2 \rfloor = 0 \\ \iff & \{\text{Let } \ell = \lfloor i/2 \rfloor\} \\ & \max(\ell, k) \neq 0 \wedge k - \ell - \lfloor j/2 \rfloor = 0 \wedge \ell - \lfloor i/2 \rfloor = 0 \end{aligned}$$

The five non-domain predicates in Ex. A.1 become domain predicates if we replace rule E (this implies the deletion of rule L) by:

$$E ::= A; |E|;$$

$$E_1 \odot E_2; \max(E_1, E_2); \min(E_1, E_2); \lceil (E_1 \odot E_2)/n \rceil; \lfloor (E_1 \odot E_2)/n \rfloor; (E_1 \odot E_2) \operatorname{div} n$$

It is not difficult to see that all domain predicates defined this way can be put in normal form if extra indices are introduced as discussed previously. However, the introduction of extra indices increases the dimensionality of the index space and consequently degrades the efficiency of systolic arrays.



## Bibliography

- [1] M. Annaratone, E. Arnould, T. Gross, H. T. Kung, M. Lam, O. Menzilcioglu, and J. A. Webb. The Warp computer: Architecture, implementation, and performance. *IEEE Trans. on Computers*, C-36(12):1523–1538, Dec. 1987.
- [2] J. Annevelink and P. Dewilde. HIFI: A functional design system for VLSI processing arrays. In K. Bromley, S.-Y. Kung, and E. E. Swartzlander, editors, *Proc. Int. Conf. on Systolic Arrays*, pages 413–452. IEEE Computer Society, 1988.
- [3] J.-P. Banâtre, A. Coutant, and D. Le Métayer. A parallel machine for multi-set transformation and its programming style. *Future Generation Computer Systems*, 4(2):133–144, Sept. 1988.
- [4] J.-P. Banâtre and D. Le Métayer. The GAMMA model and its discipline of programming. *Science of Computer Programming*, 15(1):55–77, Nov. 1990.
- [5] U. Banerjee. *Dependence Analysis for Supercomputing*. Kluwer Academic Publishers, 1988.
- [6] M. Barnett and C. Lengauer. The synthesis of systolic programs. In *Proc. Seminar on Research Directions in High-Level Parallel Programming Languages*. Springer-Verlag, 1991. To appear.
- [7] M. Barnett and C. Lengauer. A systolizing compilation scheme. Technical Report TR-91-03, Department of Computer Sciences, The University of Texas at Austin, Jan. 1991. Also: Technical Report ECS-LFCS-91-134, Department



- of Computer Science, University of Edinburgh. Abstract: *Proc. 1991 Int. Conf. on Parallel Processing, Vol. II*, Pennsylvania State University Press, 1991, 305–306.
- [8] A. Benaini, P. Quinton, Y. Robert, and B. Tourancheau. Synthesis of a new systolic architecture for the algebraic path problem. *Science of Computer Programming*, 15(2-3):135–158, Dec. 1990.
- [9] A. Benaini and Y. Robert. Spacetime-minimal systolic architectures for Gaussian elimination and the algebraic path problem. *Parallel Computing*, 15(1):211–226, 1990.
- [10] A. Benaini and M. Tchuenté. Matrix product on modular linear systolic arrays. In M. Cosnard, Y. Robert, P. Quinton, and M. Raynal, editors, *Parallel & Distributed Algorithms*, pages 79–88. North-Holland, 1989.
- [11] J. Bu, E. F. Deprettere, and P. Dewilde. A design methodology for fixed-size systolic arrays. In S. Y. Kung and E. E. Swartzlander, editors, *Application Specific Array Processors*, pages 591–602. IEEE Computer Society Press, 1990.
- [12] P. R. Cappello. A spacetime-minimal systolic array for matrix product. In J. V. McCanny, J. McWhirter, and E. E. Swartzlander, editors, *Systolic Array Processors*, pages 347–356. Prentice-Hall, 1989. To appear in *IEEE Trans. on Parallel and Distributed Systems* titled: A processor-time-minimal systolic array for cubical mesh algorithms.
- [13] K. M. Chandy and J. Misra. *Parallel Program Design*. Addison-Wesley, 1988.
- [14] M. C. Chen. A design methodology for synthesizing parallel algorithms and architectures. *J. Parallel and Distributed Computing*, 3(4):461–491, 1986.
- [15] C. Choffrut and K. Culik. On real-time cellular automata and trellis automata. *Acta Informatica*, 21:393–407, 1984.



- [16] P. Clauss, C. Mongenet, and G. R. Perrin. Calculus of space-optimal mappings of systolic algorithms on processor arrays. In S. Y. Kung and E. E. Swartzlander, editors, *Application Specific Array Processors*, pages 5–18. IEEE Computer Society Press, 1990.
- [17] A. Darte. Regular partitioning for synthesizing fixed-size systolic arrays. Technical Report 91-10, Laboratoire LIP-IMAG, Ecole, Normale Supérieure de Lyon, Apr. 1991.
- [18] J.-M. Delosme. A parallel algorithm for the algebraic path problem. In M. Cosnard, Y. Robert, P. Quinton, and M. Raynal, editors, *Parallel & Distributed Algorithms*, pages 67–78. North-Holland, 1989.
- [19] J. M. Delosme and I. C. F. Ipsen. Systolic array synthesis: Computability and time cones. In M. Cosnard, P. Quinton, Y. Robert, and M. Tchente, editors, *Parallel Algorithms & Architectures*, pages 295–312. North-Holland, 1986.
- [20] E. W. Dijkstra and C. S. Scholten. *Predicate Calculus and Program Semantics*. Texts and Monographs in Computer Science. Springer-Verlag, 1990.
- [21] B. R. Engstrom and P. R. Cappello. The SDEF systolic programming system. In S. K. Tewksbury B. W. Dickinson and S. C. Schwartz, editors, *Concurrent Computations*, chapter 15. Plenum Press, 1987.
- [22] H. A. Fencl and C.-H. Huang. On the synthesis of programs for various parallel architectures. In *Proc. 1991 Int. Conf. on Parallel Processing, Vol. II*, pages 202–206. Pennsylvania State University Press, 1991.
- [23] J. A. B. Fortes and D. I. Moldovan. Parallelism detection and algorithm transformation techniques useful for VLSI architecture design. *J. Parallel and Distributed Computing*, 2(3):277–301, Aug. 1985.



- [24] P. Gachet, B. Joinnault, and P. Quinton. Synthesizing systolic arrays using DIASTOL. In W. Moore, A. McCabe, and R. Urquart, editors, *Systolic Arrays*, pages 25–36. Adam Hilger, 1987.
- [25] I. Graham and T. King. *The Transputer Handbook*. Prentice-Hall, 1990.
- [26] L. J. Guibas, H. T. Kung, and C. D. Thompson. Direct VLSI implementation of combinatorial algorithms. In *Proc. Caltech Conf. on VLSI*, pages 509–525, 1979.
- [27] C.-H. Huang and C. Lengauer. The derivation of systolic implementations of programs. *Acta Informatica*, 24(6):595–632, Nov. 1987.
- [28] R. P. Hughey. *Programmable Systolic Arrays*. PhD thesis, Department of Computer Science, Brown University, May 1991. Technical Report CS-91-34.
- [29] K. Hwang and Y. H. Chung. Partitioned algorithms and VLSI structures for large-scale matrix computations. In *Proc. Symp. on Computing Arithmetic*, pages 222–232, 1981.
- [30] O. H. Ibarra, S. M. Kim, and M. A. Palis. Designing systolic algorithms using sequential machines. *IEEE Trans. on Computers*, C-35:531–542, June 1991.
- [31] F. Irigoien and R. Triolet. Dependence approximation and global parallel code generation for nested loops. In M. Cosnard, Y. Robert, P. Quinton, and M. Raynal, editors, *Parallel & Distributed Algorithms*, pages 297–308. North-Holland, 1989.
- [32] R. M. Karp, R. E. Miller, and S. Winograd. The organization of computations for uniform recurrence equations. *J. ACM*, 14(3):563–590, July 1967.
- [33] V. K. Prasanna Kumar and Y.-C. Tsai. Designing linear systolic arrays. *J. Parallel and Distributed Computing*, 7(3):441–463, Nov. 1989.
- [34] H. T. Kung. Let's design algorithms for VLSI systems. In *Proc. Caltech Conf. on VLSI*, pages 65–87, 1979.



- [35] H. T. Kung. Why systolic architectures? *Computer*, 15(1):37–46, Jan. 1982.
- [36] H. T. Kung and M. S. Lam. Wafer-scale integration and two-level pipelined implementations. *J. Parallel and Distributed Computing*, 1(1):33–63, Aug. 1984.
- [37] H. T. Kung and C. E. Leiserson. Algorithms for VLSI processor arrays. In C. Mead and L. Conway, editors, *Introduction to VLSI Systems*, chapter 8.3. Addison-Wesley, 1980.
- [38] S.-Y. Kung. *VLSI Processor Arrays*. Prentice-Hall Int., 1988.
- [39] S.-Y. Kung, S.-C. Lo, and P. S. Lewis. Optimal systolic design for the transitive closure and shortest path problems. *IEEE Trans. on Computers*, C-36(5):603–614, May 1987.
- [40] L. Lamport. The parallel execution of DO loops. *Comm. ACM*, 17(2):83–93, Feb. 1974.
- [41] P. Lee and Z. Kedem. Synthesizing linear-array algorithms from nested for loop algorithms. *IEEE Trans. on Computers*, C-37(12):1578–1598, Dec. 1988.
- [42] P. Lee and Z. M. Kedem. On high-speed computing with a programmable linear array. *J. Supercomputing*, 4:223–249, 1990.
- [43] C. E. Leiserson and J. B. Saxe. Optimising synchronous circuitry by retiming. In R. Bryant, editor, *Proc. Caltech Conf. on VLSI*, pages 87–116, 1983.
- [44] C. Lengauer. Code generation for a systolic computer. *SOFTWARE—Practice and Experience*, 20(3):261–282, Mar. 1990.
- [45] C. Lengauer and J. Xue. A systolic array for pyramidal algorithms. *J. VLSI Signal Processing*, 3(3):239–259, 1991.
- [46] G. J. Li and B. W. Wah. The design of optimal systolic algorithms. *IEEE Trans. on Computers*, C-34(10):66–77, Jan. 1985.



- [47] B. Lisper. *Synthesizing Synchronous Systems by Static Scheduling in Space-Time*. Lecture Notes in Computer Science 362. Springer-Verlag, 1989.
- [48] B. Lisper. Synthesis of time-optimal systolic arrays with cells with inner structure. *J. Parallel and Distributed Computing*, 10(2):182–187, 1990.
- [49] W. Luk, G. Jones, and M. Sheeran. Computer-based tools for regular array design. In J. McCanny, J. McWirth, and E. E. Swartzlander, editors, *Systolic Array Processors*, pages 589–598. Prentice Hall Int., 1989.
- [50] D. I. Moldovan. On the design of algorithms for VLSI systolic arrays. *Proc. IEEE*, 71(1):113–120, Jan. 1983.
- [51] D. I. Moldovan. ADVIS: A software package for the design of systolic arrays. *IEEE Trans. on Computer-Aided Design*, CAD-6(1):33–40, Jan. 1987.
- [52] D. I. Moldovan and J. A. B. Fortes. Partitioning and mapping algorithms into fixed-size systolic arrays. *IEEE Trans. on Computers*, C-35(1):1–12, Jan. 1986.
- [53] F. J. Núñez and M. Valero. A block algorithm and optimal fixed-size systolic array processor for the algebraic path problem. *J. VLSI Signal Processing*, 1(2):153–162, October 1991.
- [54] E. T. L. Omtzigt. Systars: A CAD tool for the synthesis and analysis of VLSI systolic/wavefront arrays. In K. Bromley, S.-Y. Kung, and E. E. Swartzlander, editors, *Proc. Int. Conf. on Systolic Arrays*, pages 383–391. IEEE Computer Society, 1988.
- [55] P. Quinton. Automatic synthesis of systolic arrays from uniform recurrent equations. In *Proc. 11th Ann. Int. Symp. on Computer Architecture*, pages 208–214. IEEE Computer Society Press, 1984.
- [56] P. Quinton. Derivation of regular parallel algorithms with the ALPHA language. In *Proc. Seminar on Research Directions in High-Level Parallel Programming Languages*, June 1991.



- [57] P. Quinton and V. van Dongen. The mapping of linear recurrence equations on regular arrays. *J. VLSI Signal Processing*, 1(2):95–113, Oct. 1989.
- [58] S. V. Rajopadhye. Synthesizing systolic arrays with control signals from recurrence equations. *Distributed Computing*, 3:88–105, 1989.
- [59] S. V. Rajopadhye. Algebraic transformations in systolic array synthesis: A case study. In L. J. M. Claesen, editor, *Formal VLSI Specification and Synthesis (VLSI Design Methods-I)*, pages 361–370. North-Holland, 1990.
- [60] S. V. Rajopadhye and R. M. Fujimoto. Synthesizing systolic arrays from recurrence equations. *Parallel Computing*, 14(2):163–189, June 1990.
- [61] I. Ramakrishnan, D. Fussell, and A. Silberschatz. Mapping homogeneous graphs on linear arrays. *IEEE Trans. on Computers*, C-35(3):189–209, Mar. 1986.
- [62] I. V. Ramakrishnan and P. J. Varman. Dynamic programming and transitive closure on linear pipelines. In *Proc. Int. Conf. on Parallel Processing*, pages 359–364. IEEE Press, 1984.
- [63] J. Ramanujam and P. Sadayappan. A methodology for parallelizing programs for multicomputers and complex memory multiprocessors. In *Supercomputing '89*, pages 637–646. ACM Press, 1989.
- [64] S. K. Rao. *Regular Iterative Algorithms and their Implementations on Processor Arrays*. PhD thesis, Department of Electrical Engineering, Stanford University, Oct. 1985.
- [65] H. B. Ribas. *Automatic Generation of Systolic Programs from Nested Loops*. PhD thesis, Department of Computer Science, Carnegie-Mellon University, June 1990. Technical Report CMU-CS-90-143.
- [66] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.



- [67] G. Rote. A systolic array algorithm for the algebraic path problem (shortest paths; matrix inversion). *Computing*, 34(3):191–219, 1985.
- [68] D. A. Schmidt. *Denotational Semantics*. Wm. C Brown, 1988.
- [69] W. Shang and J. A. B. Fortes. On the optimality of linear schedules. *J. VLSI Signal Processing*, 1(2):209–220, Oct. 1989.
- [70] W. Shang and W. A. B. Fortes. Time optimal linear schedules for algorithms with uniform dependences. *IEEE Trans. on Computers*, C-40:723–742, June 1991.
- [71] J. Teich and L. Thiele. Control generation in the design of processor arrays. *J. VLSI Signal Processing*, 3(1/2):77–92, June 1991.
- [72] J. D. Ullman. *Computational Aspects of VLSI*. Computer Science Press, 1984.
- [73] M. Valero-García, J. J. Navarro, J. M. Llabería, and M. Valeo. Implementation of systolic algorithms using pipelined functional units. In S. Y. Kung and E. E. Swartzlander, editors, *Application Specific Array Processors*, pages 272–283. IEEE Computer Society Press, 1990.
- [74] V. van Dongen. Quasi-regular arrays: Definition and design methodology. In *Systolic Array Processors*, pages 126–135. Prentice-Hall Int., 1989.
- [75] V. van Dongen and M. Petit. PRESAGE: A tool for the parallelization of nested loop programs. In L. J. M. Claesen, editor, *Formal VLSI Specification and Synthesis (VLSI Design Methods-I)*, pages 341–359. North-Holland, 1990.
- [76] M. E. Wolf and S. lam. A loop transformation theory and an algorithm to maximize parallelism. *IEEE Trans. on Computers*, C-2(2):452–471, Oct. 1991.
- [77] M. Wolfe. *Optimizing Supercompilers for Supercomputers*. MIT Press, 1989.



- [78] Y. Wong and J. M. Delosme. Optimal systolic implementations of  $n$ -dimensional recurrences. In *Proc. IEEE Int. Conf. on Computer Design (ICCD 85)*, pages 618–621. IEEE Press, 1985. Also: Technical Report 8810, Department of Computer Science, Yale University, Apr. 1988.
- [79] Y. Wong and J.-M. Delosme. Broadcast removal in systolic algorithms. In K. Bromley, S.-Y. Kung, and E. E. Swartzlander, editors, *Proc. Int. Conf. on Systolic Arrays*, pages 403–412. IEEE Computer Society, 1988.
- [80] Y. Wong and J. M. Delosme. Optimization of computation time for systolic arrays. Technical Report YALEU/DCS/RR-651, Department of Computer Science, Yale University, May 1989.
- [81] Y. Wong and J. M. Delosme. Optimization of processor count for systolic arrays. Technical Report YALEU/DCS/RR-697, Department of Computer Science, Yale University, May 1989.
- [82] X. Zhong and S. Rajopadhye. Deriving fully efficient systolic arrays by quasi-linear allocation functions. In E. H. L. Aarts, J. van Leeuwen, and M. Rem, editors, *Parallel Architectures and Languages Europe (PARLE '91), Vol. I: Parallel Architectures and Algorithms*, Lecture Notes in Computer Science 505, pages 219–235. Springer-Verlag, 1991.
- [83] X. Zhong, S. Rajopadhye, and I. Wong. Systematic generation of linear allocation functions in systolic array design. *J. VLSI Signal Processing*. To appear.



# Index

- $(p, \vartheta, I)$ , 92  
 $I \xrightarrow{\vartheta} J$ , 20  
 $I \xrightarrow{\vartheta} J$ , 20  
 $S'$ , 49  
 $S^=$ , 48  
 $V^{\text{in}}$ , 113  
 $V^{\text{out}}$ , 113  
 $[\pi : \delta]$ , 12  
 $[\pi \sqcap \delta]$ , 12  
 $\Omega_V$ , 110  
 $\mathfrak{m}$ , 10  
 $\Gamma_V$ , 95  
 $\Phi$ , 16  
 $\Phi_V$ , 16  
 $\Psi$ , 29, 97  
 $\Psi^{\text{P}}$ , 29, 98  
 $\Psi_V$ , 29, 97  
 $\Psi_V^\perp$ , 30, 98  
 $\Psi_V^{\text{d}}$ , 29, 97  
 $\Psi_V^{\text{r}}$ , 97  
 $\Psi_V^{\text{s}}$ , 29, 97  
 $\text{in}_V$ , 18  
 $\text{out}_V$ , 18  
 $\Upsilon$ , 92  
 $\Upsilon_V^\perp$ , 114  
 $\Upsilon_V^{\text{d}}$ , 114  
 $\Upsilon_V^{\text{r}}$ , 114  
 $\Upsilon_V^{\text{s}}$ , 114  
 $\text{aff}(S)$ , 12  
 $\mathfrak{B}$ , 59  
 $\mathbb{C}_S A$ , 10  
 $\mathfrak{S}$ , 48  
 $\mathfrak{S}_{\text{src}}$ , 55  
 $\square$ , 12  
 $\mathcal{H}(X, \mathfrak{S})$ , 51  
 $\mathcal{I}_V$ , 30  
 $\mathcal{O}_V$ , 30  
 $\mathcal{P}$ , 24  
 $\mathcal{P}(S)$ , 49  
 $\mathcal{V}$ , 17  
 $\chi$ -model, 87  
 $\chi_T$ , 10  
 $\mathbb{F}$ , 116  
 $\mathbb{1}_V$ , 86  
 $\mathbb{1}'_V$ , 98  
 $\mathbb{L}$ , 117  
 $\mathbb{O}_V$ , 86  
 $\mathbb{P}$ -partition, 46  
 $\mathbb{P}_{\text{src}}$ , 47  
 $\mathbb{T}$ , 42  
 $\text{cone}(x_1, x_2, \dots, x_m)$ , 13  
 $\text{conv}(S)$ , 12



- $\delta_V$ , 96
- $\text{deq}(V)$ , 16
- $\text{deq}(\Phi)$ , 17
- $\text{dim}(S)$ , 12
- $\text{facets}(S)$ , 13
- $\text{fst}(\Phi_V, \vartheta_V)$ , 20
- flow, 25
- $\text{in}(\Phi_V, \vartheta_V)$ , 20
- inject, 119
- input, 86
- $\text{lst}(\Phi_V, \vartheta_V)$ , 20
- $\text{lin}(S)$ , 12
- $\text{span}(x_1, x_2, \dots, x_m)$ , 10
- output, 87
- $\text{paths}(S, \vartheta)$ , 92
- pattern, 25
- $\pi$ -model, 84
- place, 24
- $\text{planes}(S, \pi)$ , 97
- pi, 86
- po, 86
- $\text{rank}(A)$ , 11
- $\text{ray}(I, \vartheta)$ , 20
- $\text{rays}(D, \vartheta)$ , 20
- $\text{sig}(V)$ , 45
- step, 24
- $\text{sup}(S, F)$ , 13
- $\mathfrak{S}^c$ , 48
- $\bar{x}$ , 25
- $\varphi(\bar{I})$ , 114
- i.x.*, 11
- $k$ -parallelepiped, 14
- $t_{\text{fst}}$ , 25, 87
- $t_{\text{lst}}$ , 87
- adapted index space, 145
- address variable, 159
- allocation
  - function, 25
  - matrix, 24
  - vector, 83
- AREs, 16
- border cell, 23, 85, 87
- broadcast dependence, 19
- CCUREs, 5, 44
- cell program, 112
- communication constraint domain, 110
- computation
  - control flow, 4
  - equation, 18
  - point, 19, 92
- conditional, 17, 173
- control
  - channel, 14
  - UREs, 5, 40
  - variable, 40
- convex
  - class partition, 46
  - domain partition, 47
- data
  - channel, 14
  - dependence graph, 19
  - dependence matrix, 19



- dependence vector, 17
- pipelining, 33
- variable, 40
- domain
  - of a variable, 16
  - of an equation, 16
  - predicate, 17, 173
- draining point, 30, 95
- evolution control variable, 112
- extended
  - data dependence graph, 30
  - domain of a variable, 30
  - index space, 29, 97
  - source UREs, 30, 98
- first
  - computation point, 19
  - step number, 25, 87
- host program, 112
- image under a space-time mapping (over-  
bar), 25
- index
  - mapping, 17
  - space, 16
  - vector, 17
- initialisation control variable, 116
- inner-product UREs, 148
- input
  - cell, 23, 85, 87
  - data, 18
  - equation, 18
- intermediate data, 18
- internal cell, 23, 85, 87
- last
  - computation point, 19
  - step number, 87
- latency, 25
- link, 85
- load-recover variable, 159
- non-injectivity of a space-time map-  
ping, 99
- normal form of domain predicates, 17,  
173
- output
  - cell, 23, 85, 87
  - data, 18
  - equation, 18
- parameterised systolic array, 35
- partitioning, 35
- PCUREs, 5, 44
- pipelined UREs, 45
- pipelining
  - dependence, 19
  - equation, 20
  - point, 30, 92
  - variable, 20
  - vector, 33
- place function, 25
- point, 17
- processor space, 24
- projection, 36



- vector, 26
- propagation control flow, 4
- relaying point, 95
- scheduling
  - matrix, 100
  - vector, 24
- separation set, 48
- set of defining half-spaces, 48
- soaking point, 30, 95
- source UREs, 4
- space-time
  - diagram, 94
  - mapping, 24
  - matrix, 25, 83, 100
  - point, 92
  - UREs, 28
- step function, 25
- systolic algorithm, 3
- systolic array model, 23
- temporal hyperplane, 25
- termination control variable, 116
- timing function, 25
- type of a point, 42
- undefined point, 30, 95
- UREs, 19
- validity of a space-time mapping, 26,
  - 88, 89, 101, 133, 139
  - for a variable, 115



## ERRATA (4 March, 1992)

- p21, ¶1 and p141, last paragraph and p165, (7.4):** Replace  $n$  by  $m$ .
- p26, Precedence Rule and p34, last sentence:** Replace “smaller” by “larger”.
- p28:** Replace  $0 = i \wedge 0 < j \leq m \wedge 0 < k \leq m$  by  $0 = x \wedge 0 < y \leq m \wedge 0 < t - x - y \leq m$ .
- p43, ¶1:** Replace the sentence starting with “Therefore” by “Therefore, the defining equations for a variable agree at points of the same type and the defining equations for some variable do not agree at points of different types.”.
- p43, ¶2, below (3.3):** Replace  $|\Phi/\oplus|$  by  $|\mathcal{V}|$ .
- p49, (3.8):** Delete  $C_i(I) =$  at the right hand side of each  $\rightarrow$ .
- p52, Ex. 3.2:** Replace Figs. 3-1 by Figs. 3-2.
- p57, Step 3 of Proc. 3.1:** Replace the second  $\pi I < \delta$  ( $\pi I \leq \delta$ ) by  $\pi I \leq \delta$   $\pi I < \delta$ .
- p57, below Lemma 3.9 and p65, above and below Thm. 3.9:** Replace  $[\pi I \sqcap \delta]$  by  $[\pi \sqcap \delta]$ .
- p59, (3.18):** Delete  $P_i(I) =$  at the right hand side of each  $\rightarrow$ .
- p62, (3.23):** Add  $I \in \Psi \rightarrow C(I - \vartheta_C)$ .
- p62:** Replace  $\text{ray}(\Phi, \pm \vartheta_C)$  ( $\text{ray}(D_C, \pm \vartheta_C)$ ) by  $\text{rays}(\Phi, \pm \vartheta_C)$  ( $\text{rays}(D_C, \pm \vartheta_C)$ ).
- p62, below (3.24) and p109, (5.6):** Replace  $\Psi \cup H^c.C$  by  $\Psi \cap H^c.C$ .
- p63, Fig. 3-3:** Add “extended” before “index space”.
- p65, last paragraph:** Replace “=” by “ $\geq$ ”, and the succeeding sentence by “This implies that some control signals can be broadcast across the array.”.
- p80, ¶1, last line:** Replace “control” by “pipelined”.
- p80, last paragraph:** Replace the definition of  $\oplus_{\perp}$  by
- $$I \oplus_{\perp} J \iff (\forall V : V \in \mathcal{V} : (\forall D : D \in \text{deq}(V) : I \in D \implies J \in D \vee V(J) = \perp))$$
- p83, the definition of place:** Replace the first  $Z$  by  $Z^r$ .
- p95, above (4.3):** Replace “oaking” by “soaking”.
- p101, ¶ 1:** Replace the first  $s_i$  by  $s_{i+1}$ .
- p102:, Step. 1 of Proc. 4.1:** Replace 1 by  $\Delta_i$ :
- p102:, Step. 2 of Proc. 4.1:** Change it to:
- $$\left( \forall i : 0 < i < n : \Delta_i = \begin{cases} \text{if } (\forall I, J : I, J \in \Psi : \Lambda_i I = \Lambda_i J) \rightarrow 1 \\ \square \text{ else } \rightarrow (\min I, J : I, J \in \Psi \wedge \Lambda_i I \neq \Lambda_i J : |\Lambda_i(I - J)|) \\ \text{fi} \end{cases} \right)$$
- p118, Stat. 3 of Lemma 5.1, after (a):** Replace  $\bar{I}$  by  $\bar{I}_1$ .
- p124, last paragraph:** Replace  $\text{fst}(F_{i,k}(I))$  by  $\text{input}(F_{i,k}(I))$ .
- p126:** Replace  $m\sigma\vartheta_V \neq 0 \neq n\sigma\vartheta_V$  by  $m\lambda\vartheta_V \neq 0 \neq n\lambda\vartheta_V$ .
- p152, (6.4):** Replace  $\Phi_C^d$  by  $\Psi_C^d$ .
- p154, last sentence:** Replace “three” by “four”.
- p158, ¶2:** Replace  $\text{rays}(\text{fst}(\Phi_V, \vartheta_V) - \vartheta_V, -\vartheta_{P_V})$  by  $\text{rays}(\text{fst}(\Phi_V, \vartheta_V), -\vartheta_{P_V}) \setminus \text{fst}(\Phi_V, \vartheta_V)$ .
- p158, last paragraph:** Replace  $\text{rays}(\text{lst}(\Phi_V, \vartheta_V) + \vartheta_V, \vartheta_{P_V})$  by  $\text{rays}(\text{lst}(\Phi_V, \vartheta_V), \vartheta_{P_V}) \setminus \text{lst}(\Phi_V, \vartheta_V)$ .
- p167, the end of § 7.4:** Add the following:

The construction of addressing control presented in Sect. 7.3 works only for three-dimensional UREs. For  $n$ -dimensional UREs, we simply associate  $n - 2$  address variables,  $X_V^1, X_V^2, \dots, X_V^{n-2}$  with every stationary data variable  $V$ . Element  $V(I)$  is associated with memory location  $(X_V^1(I), X_V^2(I), \dots, X_V^{n-2}(I))$  at cell  $\text{place}(I)$ . The specification of these address variables is as in (7.3). We must choose  $\pi_{X_V^i}$  and  $\vartheta_{X_V^i}$  with  $\pi_{X_V^i} \vartheta_{X_V^i} = 0$  such that (1)  $\sigma\vartheta_{X_V^i} \neq 0$  (this ensures that  $X_V^i$  is a moving variable), (2)  $\pi_{X_V^i} \vartheta_V = 0$  (this ensures that  $X_V^i(I)$  is the same address for all points of a  $\vartheta_V$ -path), and (3)  $\pi_{X_V^1}, \pi_{X_V^2}, \dots, \pi_{X_V^{n-2}}, \sigma$  are linearly independent (this ensures that different elements of stationary variable  $V$  that are mapped to the same cell are given different memory locations).