

---

# Scaling Integer Arithmetic in Probabilistic Programs

---

William X. Cao<sup>1</sup> Poorva Garg<sup>\*1</sup> Ryan Tjoa<sup>\*2</sup> Steven Holtzen<sup>3</sup> Todd Millstein<sup>1</sup> Guy Van den Broeck<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of California, Los Angeles, California, USA

<sup>2</sup>Department of Computer Science, University of Washington, Seattle, Washington, USA

<sup>3</sup>Khoury College of Computer Sciences, Northeastern University, Boston, Massachusetts, USA

## Abstract

Distributions on integers are ubiquitous in probabilistic modeling but remain challenging for many of today’s probabilistic programming languages (PPLs). The core challenge comes from discrete structure: many of today’s PPL inference strategies rely on enumeration, sampling, or differentiation in order to scale, which fail for high-dimensional complex discrete distributions involving integers. Our insight is that there is structure in arithmetic that these approaches are not using. We present a binary encoding strategy for discrete distributions that exploits the rich logical structure of integer operations like summation and comparison. We leverage this structured encoding with knowledge compilation to perform exact probabilistic inference, and show that this approach scales to much larger integer distributions with arithmetic.

## 1 INTRODUCTION

Probabilistic programming languages (PPLs) are expressive languages for defining probability distributions. The core idea of a PPL is to enrich a programming language with the ability to define, observe, and compute with random variables: hence, the program itself defines a probabilistic model. This paper focuses on a particular programming feature: scaling inference for programs with random integers and integer arithmetic. Integers are very challenging for today’s approaches to probabilistic inference. The relationships between integer-valued random variables can be very complex: they can be added, multiplied, compared, etc. This rich structure is opaque to today’s inference strategies. Trace-based sampling like Markov-Chain Monte Carlo, importance sampling, and sequential Monte-Carlo all collapse integer distributions to a single sampled point [Gelman

et al., 2015, Bingham et al., 2019, Dillon et al., 2017, van de Meent et al., 2018, Lew et al., 2019]. These approximate inference strategies can scale well in many cases, but they struggle to find valid sampling regions in the presence of low-probability observations and non-differentiability (e.g., observing the sum of two large random integers to be a constant) [Gelman et al., 2015, Bingham et al., 2019, Dillon et al., 2017]. Exact inference strategies work by preserving the global structure of the distribution, but here there is a challenge: *what is the right strategy for efficiently representing and manipulating distributions on integers?* Today’s PPLs that support exact inference and integer manipulation – such as Dice [Holtzen et al., 2020], ProbLog [De Raedt et al., 2007], Psi [Gehr et al., 2016], and WebPPL [Goodman and Stuhlmüller, 2014] – model integer distributions using what is essentially a one-hot categorical encoding (i.e., an integer distribution  $[0 \mapsto 0.25, 1 \mapsto 0.25, 2 \mapsto 0.25, 3 \mapsto 0.25]$  is represented simply as a vector). This encoding style is not capable of exploiting the structure of addition: adding two random variables effectively requires full enumeration.

Our first contribution is a new representation of distributions on integers as distributions on *binary encodings*. For instance, in the above example, rather than representing the distribution as an exhaustive map from integer values to probabilities, we represent it as a joint distribution on binary bits  $[00 \mapsto 0.25, 01 \mapsto 0.25, 10 \mapsto 0.25, 11 \mapsto 0.25]$ . The upsides of this seemingly-equivalent representation are twofold. First, we can more efficiently represent the joint distribution itself when it has certain structure. In this case, because the distribution is uniform, we can represent it as a product of two independent Bernoulli distributions, one for each bit: we will show that the ability to factorize the distribution in this manner leads to significant performance improvements. Second, this binary representation reveals the structure of arithmetic: for instance, we can compare two integers by *independently* comparing each of their binary digits and aggregating the results.

Clearly a binary representation of integers reveals structure, but how can we automatically find and exploit this structure

<sup>\*</sup>These authors contributed equally to this work.

---

```

1 id = [discrete([0.72, 0.01, 0.01, 0.01, 0.01,
2             0.01, 0.2, 0.01, 0.01, 0.01]),...,
3       discrete([0.01, 0.01, 0.05, 0.01, 0.01,
4             0.63, 0.2, 0.01, 0.01, 0.05])]
5 check_digit = id[0]
6 remaining_id = id[1:] //tail of the array
7 check_val = luhn_checksum(remaining_id)
8 observe((check_digit + (check_val % 10)) == 10)
9 return id

```

---

Figure 1: A probabilistic program for the student ID probabilistic inference problem using integer random variables (discrete), integer arithmetic (the Luhn algorithm function), and Bayesian conditioning (observe)

---

```

1 def luhn_checksum(id)
2     sum = 0
3     for i in 0..length(id) - 1
4         if i % 2 == length(id) % 2:
5             if id[i] > 4:
6                 sum += 2 * id[i] - 9
7             else:
8                 sum += 2 * id[i]
9         else:
10            sum += id[i]
11    return sum

```

---

Figure 2: Luhn algorithm implementation

during inference in a PPL? As our second contribution, we show that two of today’s PPLs – Dice [Holtzen et al., 2020] and ProbLog [De Raedt et al., 2007, Fierens et al., 2015] – are already capable of exploiting this structure if it is properly encoded into the program, by virtue of their *knowledge compilation* approach to inference. We give a lightweight strategy for encoding integer distributions, and show empirically that when using our new binary-encoded distributions these two languages scale to significantly larger and more complex integer distributions without essential modifications to their existing inference strategies.

As our third contribution we show that scalable support for random integer arithmetic allows us to push the boundaries of discrete probabilistic programming systems in surprising ways. We demonstrate how to model a Beta distribution, a continuous distribution, using probabilistic integers. This modelling method exploits the conjugacy property of the Beta distribution, through which we can always characterize the distribution through its (integral) sufficient statistics. By doing so, we can use the Beta distribution as a prior for Bayesian learning.

The structure of this paper is as follows: Section 2 gives a motivating example for integer arithmetic. Section 3 explains our integer representation and explores how common integer operations on this representation have structure exploitable by knowledge compilation. Section 4 empirically evaluates our representation strategy against existing PPLs. Section 5 explores the representation of a continuous Beta prior with random integers. Sections 6 and 7 discuss related work and conclude respectively.

## 2 MOTIVATION

We begin with a motivating example highlighting how integer distributions are used in probabilistic programs. Consider the following probabilistic model based on student ID numbers. Suppose that an optical character recognition system is attempting to parse a handwritten student ID number. For each digit of the ID, it produces a probability distri-

bution representing its beliefs about what the digit could be. Combining this output, we get a probability distribution over all possible student IDs.

The Luhn algorithm [Luhn, 1960] is a commonly used method of validating various ID numbers including student IDs. Given a starting ID such as 70733428, the algorithm provides for a way to compute a sum over the ID, giving us a check digit (4) which is then prepended to the original ID to get a final ID: 470733428. This ID is the one actually issued to a student; when provided with an ID, we can validate it by recomputing the sum and looking at the check digit.

We wish to use the fact that the student ID can be validated to additionally inform our single-digit distributions from the OCR system. We can implement this as a probabilistic program like the one in Figure 2. Figure 2 implements a function `luhn_checksum` that takes as input a list representing the digits of the student ID, excluding the check digit. It then does computation according to the Luhn algorithm to compute a sum over the digits, which is then returned. Figure 1 then uses this function: we create a list `id` which contains distributions over integers derived from the OCR system. The syntax `discrete(v)` for a vector  $v = [p_0, \dots, p_n]$  creates a distribution over the numbers  $0, \dots, n$ , in which the number  $i$  has the probability  $p_i$ , and is used in our program to represent said OCR distributions. We call the `luhn_checksum` function on these integer distributions to get a distribution over checksums and condition using the `observe` keyword in line 9 to get an updated distribution over IDs.

If we implement this program in today’s probabilistic programming languages, we will run into a problem. Even if we only wish to compute the marginal probability over a single digit of the ID, Figure 3 shows that the runtime will scale exponentially in the number of digits in the student ID. Each additional digit will contribute a multiplicative amount to the number of total possible ID instantiations, meaning that any approach involving enumeration is inherently exponential. In practice, this means that programs containing student IDs of a realistic length (9-10 digits) will not run.

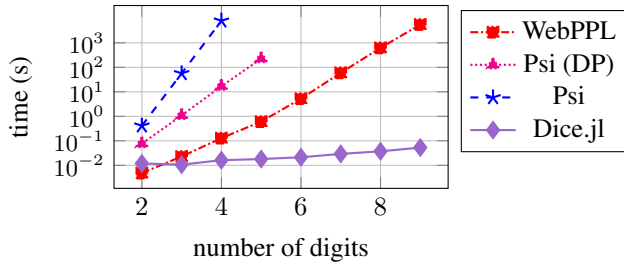


Figure 3: Single-marginal performance for ID example on increasing ID lengths. WebPPL and Psi scale exponentially due to having to enumerate all paths.

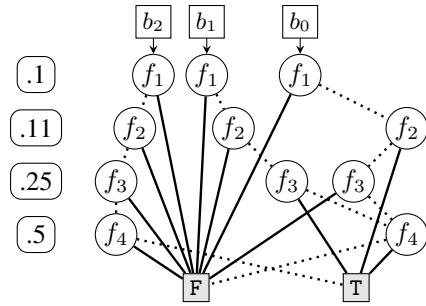
The fact that such straightforward programs fail to scale on existing probabilistic programming systems is the primary motivation behind our work. We have implemented our encoding of integer distributions in `Dice.jl`, a discrete PPL embedded in Julia that uses the same knowledge compilation approach as Dice [Holtzen et al., 2020]. Figure 3 shows that our technique allows inference for such programs to scale for a larger, more realistic number of digits.

### 3 REPRESENTING & MANIPULATING INTEGER DISTRIBUTIONS

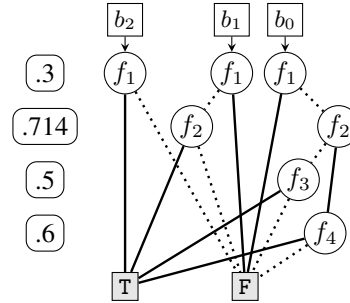
This section describes the key technical details behind a binary encoding approach and explains how such an encoding allows the knowledge compilation inference strategy used by Dice and ProbLog to automatically exploit arithmetic structure. We first provide a brief introduction to inference via knowledge compilation. We then demonstrate and analyze various approaches to constructing distributions over integers within probabilistic programs. Finally, we show how the binary encoding can be leveraged by knowledge compilation to identify and exploit conditional independencies for inference over distributions with integer arithmetic.

#### 3.1 INTEGER DISTRIBUTIONS VIA BDDS

Thus far we have seen how binary-represented distributions expose structure and can enable effective scaling in practice. In this section we explain exactly how this performance improvement is achieved during inference. In particular, we show how *knowledge compilation* is capable of automatically finding and exploiting the structure of integer distributions and operations. Inference via knowledge compilation is currently the state-of-the-art approach to exact discrete probabilistic inference in certain classes of probabilistic programs [Holtzen et al., 2020, De Raedt et al., 2007, Chavira et al., 2006, Chavira and Darwiche, 2008, Fierens et al., 2015]. The heart of inference via knowledge compilation is a reduction from inference to *weighted model counting* (WMC). Let  $\varphi$  be a Boolean formula and  $w$  be a



(a) Multi-rooted BDD for a CATEG\_INT encoded integer.



(b) Multi-rooted BDD for an BITWISE\_INT encoded integer.

Figure 4: BDDs representing the integer distribution  $[0 \mapsto 0.1, 1 \mapsto 0.1, 2 \mapsto 0.2, 3 \mapsto 0.3, 4 \mapsto 0.3]$  resulting from CATEG\_INT and BITWISE\_INT encoding methods. BITWISE\_INT achieves a smaller BDD by compactly representing higher order bits.

map from literals in  $\varphi$  to real-valued weights; the pair  $(\varphi, w)$  is called a *weighted Boolean formula*. Then, the *weighted model count*  $WMC(\varphi, w)$  is a weighted sum of models of  $\varphi$ :

$$WMC(\varphi, w) = \sum_{m \models \varphi} \prod_{\ell \in m} w(\ell). \quad (1)$$

This reduction to WMC is not useful on its own however: the WMC task is #P-hard for an arbitrary Boolean formula  $\varphi$ . This is where knowledge compilation comes into the picture:  $\varphi$  is compiled into a data structure that supports efficient weighted model counting in the size of the data structure. A common example of such a knowledge-compilation data structure is a binary decision diagram (BDD), which supports linear-time WMC, but there are many others [Darwiche and Marquis, 2002]. PPLs like Dice and ProbLog work by compiling a program into a BDD or related compilation target and thereby reducing probabilistic program inference to WMC on that target [Chavira and Darwiche, 2005, Sang et al., 2005b].

The cost of the knowledge compilation approach to inference is almost entirely determined by the structure of the program; the more structure that exists, the more compact the resulting BDD or related data structure can be. This leads to our core contribution: a new logical representation of integer distributions that is amenable to efficient

compilation into BDDs. To demonstrate how BDDs can encode integer distributions, Figure 4 shows two different multi-rooted BDDs that represent *the same distribution* on integers. In both cases the roots in each BDD represent random variables for each binary digit:  $b_0$  is the 0-order digit,  $b_1$  is the 1-order digit. Positive weights of each literal are shown on the left (with the negative weight being 1 minus the positive weight); dotted edges represent a false assignment and solid edges represent true assignment. Intuitively, a binary representation has the potential to be more compact than a naive categorical representation due to the reduction in the number of roots: for instance, Dice [Holtzen et al., 2020] requires one root for each possible integer value.

As an example of how to use these data structures, consider computing the marginal probability of the high-order bit  $b_2$  being true. This is  $\text{WMC}(b_2, w)$ , which is 0.3 – that is clear in Figure 4b since the sole path from  $b_2$  to the true node has weight 0.3, but it is also true for the sole path from  $b_2$  to the true node in Figure 4a, which has weight  $(0.9 * 0.89 * 0.75 * 0.5) \approx 0.3$ . In general, to compute the probability of an arbitrary integer, we convert it into binary and conjoin the appropriate roots: for instance, to compute the probability of the integer 0, we compute  $\text{WMC}(\overline{b_0} \wedge \overline{b_1} \wedge \overline{b_2}, w)$ .

### 3.2 INTEGER ENCODINGS

The previous section demonstrated the potential for binary encodings in knowledge compilation, but how do we connect this to probabilistic programs? In this section we give lightweight encoding strategies for translating `discrete(...)` syntax for an arbitrary distribution over integers into distributions on Booleans, which knowledge-compilation-based languages like Dice and ProbLog are already capable of representing. How might such distributions be represented in a probabilistic programming language in practice? To make the problem concrete, we define the integer representation problem as follows: given an input vector  $[p_0, \dots, p_w]$ , we want a method which returns a distribution over integers taking on value  $i$  with probability  $p_i$ . While this is relatively restricted by demanding that our distribution is contiguous with lowest value 0, we can convert this to other distributions (for example) by adding an offset or multiplying by a constant.

#### 3.2.1 A First Approach

One natural way of constructing such a categorical distribution, is as a set of if-else statements, with each branch corresponding to a different value. For example, the following probabilistic program snippet would correspond to the integer distribution with probability vector  $[0.1, 0.2, 0.3, 0.4]$ . The syntax `flip( $\theta$ )` used in the program is commonly used in discrete PPLs to represent a Bernoulli random variable with bias  $\theta$ .

---

```

1  if flip(0.1) // Bernoulli(0.1)
2    return 0
3  elseif flip(0.2/0.9)
4    return 1
5  elseif flip(0.3/0.7)
6    return 2
7  else
8    return 3

```

---

We use a sequence of these random flips as arguments to the if-else statements to generate the mixture of numbers; note that we renormalize the flip probability at each step to get the correct distribution. This approach is generalized in Algorithm 1. This and future algorithms should be interpreted as a general method to represent a distribution over integers in any probabilistic programming language supporting Bernoulli random variables and (non-probabilistic) integers. Note that representing a categorical variable in this way is a probabilistic program framing of the SBK encoding presented by Sang et al. [2005a].

---

**Algorithm 1:** CATEG\_INT ( $v \in [0, 1]^w$ )

---

**Input:** Vector  $v$  such that  $v[i] \propto \text{pr}(i)$

**Output:** Distribution over the integers  $0, \dots, w - 1$

```

      matching  $v$ 
if  $w == 1$  or  $\text{flip}\left(\frac{v[0]}{\sum v}\right)$  then
  | return 0
else
  | // Recurse on the remainder of  $v$ 
  | return  $1 + \text{CATEG\_INT}(v[1:])$ 

```

---

What would occur if we use Algorithm 1 to represent a distribution over binary-encoded integers in a language such as Dice? The BDD for one distribution represented using this approach is given in Figure 4a. Note that for each bit, the decision diagram is essentially a linear chain; intuitively, this corresponds to checking each if-else guard in sequence. In addition, notice that there is almost no node reuse occurring in this BDD; each root has its own linear chain.

#### 3.2.2 A More Compact Encoding

We propose an alternative method of representing integers from a probability vector that produces provably more compact BDDs. Rather than constructing our mixture by a linear pass through the probability vector, we can instead divide the vector into two parts, using a divide-and-conquer approach. Consider the same example as above, where we are again given as input a probability vector  $[0.1, 0.2, 0.3, 0.4]$ . To get the wanted distribution, we can conditionally add the value 2 with probability  $\frac{0.3+0.4}{0.1+0.2+0.3+0.4}$ , corresponding to the latter half of the vector. Depending on if 2 is added, we then conditionally add the value 1, with probability derived from the subvectors  $[0.1, 0.2]$  and  $[0.3, 0.4]$ . An example

program implementing this is given below:

---

```

1 num = 0
2 if flip(0.7) // 0.3 + 0.4
3   num += 2
4   if flip(0.4/0.7)
5     num += 1
6 else
7   if flip(0.2/0.3)
8     num += 1
9 return num

```

---

This approach is formalized in Algorithm 2. For the sake of simplicity, we assume the input vector is always of length  $2^b$  for some number  $b$ ; this means that we always divide the vector into its two halves. For an arbitrary input vector, we can simply pad 0 probability values to fulfill this condition; in practice, this algorithm can easily be adapted to work without this explicit padding.

---

**Algorithm 2:** BITWISE\_INT ( $v \in [0, 1]^{2^b}$ )

---

**Input:** Vector  $v$  such that  $v[i] \propto \text{pr}(i)$

**Output:** Distribution over the integers  $\{0, \dots, 2^b - 1\}$  matching  $p$

$$p \leftarrow \frac{\sum_{i=2^{b-1}}^{2^b-1} v[i]}{\sum_{i=0}^{2^b-1} v[i]}$$

**if**  $\text{length}(v) == 1$  **then**

  | **return** 0

**else**

  | **if**  $\text{flip}(p)$  **then**

    | // Recurse on second half  $\geq 2^{b-1}$

    | **return** BITWISE\_INT( $v[2^{b-1} : 2^b] + 2^{b-1}$ )

  | **else**

    | // Recurse on first half  $< 2^{b-1}$

    | **return** BITWISE\_INT( $v[0 : 2^{b-1}]$ )

---

Note that while Algorithm 2 uses arithmetic to produce the distribution, it only ever adds or return powers of two, which directly correspond to the bits of the integer. Therefore, when implementing the algorithm as a distribution on a tuple of bits, we encode each such addition by simply setting the appropriate bit. For the same example as above, our implementation constructs a tuple of bits  $(b_1, b_0)$  corresponding to a binary number such that  $b_1 = \text{flip}(0.7)$  and  $b_0 = \text{if } b_1 \text{ then } \text{flip}(\frac{0.4}{0.7}) \text{ else } \text{flip}(\frac{0.2}{0.3})$ .

How does this method of representing integer distributions differ than the one given before? To see this, we look at the BDD for a distribution written in this manner given in Figure 4b. We can see a clear difference between this BDD and that for the approach given in Algorithm 1. The most significant bit corresponds to a BDD depending on only one flip, as this corresponds to the largest power of two: only one flip is used to determine its value. For the less significant bits, we add an additional layer of variables for each one, with the number of layers in total corresponding to the number

of bits needed to represent the input distribution. This is in contrast to the CATEG\_INT encoding, which requires checking a linear chain of variables for each bit, and so achieves a much more compact BDD representation.

We formalize this difference in BDD size in Proposition 1. Note that variable order can greatly influence the size of a BDD, and finding the optimal variable order is an NP-hard problem [Meinel and Theobald, 1998]; we follow the Dice convention of ordering logical variables using (strict, left-to-right) evaluation order. For example, in Figure 1, the Boolean variables encoding the discrete distribution on Line 1 occur before the variables in the distribution on Line 2 in the order.

**Proposition 1** *A discrete distribution over the integers  $\{0, 1, \dots, 2^b - 1\}$  compiles to a BDD of size  $\Theta(b2^b)$  when represented using CATEG\_INT (Algorithm 1) and a BDD of size  $\Theta(2^b)$  when represented using BITWISE\_INT (Algorithm 2), with variables in flip evaluation order.*

It is BITWISE\_INT that we have implemented in `Dice.jl` and experimentally evaluate in the next section.

### 3.2.3 Uniform Integers

The previous encoding strategy works for arbitrary distributions on integers, but in practice one often encounters common highly-structured distributions such as the uniform. One advantage of our approach is that we can exploit the structure of such distributions in order to scale significantly better than the general approach presented in Algorithm 2. In particular, the structure of the uniform distribution allows for a special encoding with fully independent flips.

Since the probability of every integer is equal, we can encode a uniform distribution over integers  $\{0, 1, \dots, 2^b - 1\}$  by adding the values  $2^0, 2^1, \dots, 2^{b-1}$  independently with probability 0.5. From a bitwise perspective, this is same as independently setting each bit of the number to be true with probability 0.5. As an example, consider the uniform distribution over integers  $\{0, 1 \dots 15\}$ . Clearly, each possible instantiation of  $(\text{flip}(0.5), \text{flip}(0.5), \text{flip}(0.5), \text{flip}(0.5))$  is equally likely, and thus equivalently the probability of each integer in the range.

The method described above works for uniform distributions whose range is  $2^n$  for some  $n$ ; for ranges that are not a power of 2, we use the fact we can any decompose natural number into a sum of powers of 2. This enables a uniform distribution over any range to be represented as a mixture of multiple uniform distributions over smaller power-of-two ranges. We formalize this idea in Algorithm 3, which gives a method for representing uniform distributions starting at 0; the correctness of this approach is shown in the supplementary material. We can then use this approach

---

**Algorithm 3:** UNIFORM( $n$ )

---

**Input:** Positive integer  $n$ **Output:** Integer uniformly distributed over $0, \dots, n - 1$  $b \leftarrow \lfloor \log_2(n) \rfloor$ **if** flip( $\frac{2^b}{n}$ ) **then**     $sum \leftarrow 0$     **for**  $i \leftarrow 0$  **to**  $b - 1$  **do**        **if** flip( $\frac{1}{2}$ ) **then**             $sum \leftarrow sum + 2^i$     **return**  $sum$ **else**    **return** UNIFORM( $n - 2^b$ ) +  $2^b$ 

---

to achieve any uniform distribution by adding an offset. Just like the previous algorithms, this algorithm is implemented by constructing sequences of bits in a manner equivalent to arithmetic.

We note that unlike the BITWISE\_INT algorithm, where less significant bits have a dependence on more significant bits, our uniform algorithm leverages independence between the bits. Therefore, the BDD obtained when using UNIFORM is more compact than for our other algorithms, and fewer variables are needed to represent such a distribution.

### 3.3 EFFICIENT INTEGER OPERATIONS

While the binary representations of discrete and uniform distributions over integers are interesting, they do not by themselves necessarily give much advantage. If adding two such distributions still requires an explicit enumeration of all possible sums, then we have not gained much over the existing inference approaches. However, the binary encoding enables us to leverage the structure of integers to do much better than this for common operations. In this section, we demonstrate this for integer comparisons and addition.

#### 3.3.1 Integer Comparisons

The comparison operator on binary tuples can be implemented using logic circuits like those in computer hardware. Suppose we compute  $a < b$  for two binary numbers  $a = 001$  and  $b = 100$ . The circuit first compares the most significant bits (MSBs) of these numbers, which are 0 and 1 respectively - enough to know that  $a < b$  is true. If the two numbers instead had the same MSB, we would need to start this comparison over on remaining bits. This process of computing  $a < b$  highlights its key conditional independencies. First, given the MSBs of the operands are different, the result of  $a < b$  does not depend on the remaining bits. Second, given the MSBs of the operands are same, the computation on the remaining bits does not depend on the value

of the MSBs. This structure gets automatically exploited when we use this standard logic circuit to compare integer distributions, where the inputs are now weighted Boolean formulas represented as BDDs, rather than bits.

More concretely, consider the following probabilistic program which defines two random variables having a uniform distribution over the integers  $\{0, 1, \dots, 7\}$  and then outputs the probability of one integer being less than the other.

---

```
1 a = uniform(0, 8)
2 b = uniform(0, 8)
3 return (a < b)
```

---

Enumerating all the values that  $a$  and  $b$  can take in the above program would lead to 64 combinations. In contrast, the BDD for the comparison operation has size linear in the number of bits, as it exploits conditional independences. We later present empirical results demonstrating that this leads directly to better scalability for discrete inference.

#### 3.3.2 Integer Addition

Consider two binary numbers  $a = 001$  and  $b = 100$  that we wish to add. The least significant bit (LSB) of  $a + b$  is computed as the *xor* of the LSBs of  $a$ ,  $b$  and 0 (the initial carry bit). The carry, computed as the *and* of the LSBs of  $a$  and  $b$ , is passed on to the next bit and the same process will be repeated for the remaining bits. The process described above shows that given the carry bit, each bit of the result is independent of the lesser significant bits of the operands. Similar to the comparison operation, encoding addition on integer distributions as a logical circuit directly exploits such conditional independences to produce a compact BDD, which in turn leads to significant performance gains. The manner in which addition corresponds to a compact BDD has been explored before; Wegener [2004] show that for an optimal variable ordering, there is a linear bound on the BDD for addition.

In this section, we described the structure of two arithmetic operations, comparison and addition, independently. When composing these operations together, the compilation of weighted Boolean formulas will naturally compose as described in previous work [Holtzen et al., 2020]. The sizes of the resulting BDDs depend highly on the variable ordering — but even when the variable ordering is not optimal, conditional independences can still be identified and exploited, as shown by the experiments in the next section.

## 4 EMPIRICAL EVALUATION

In this section, we empirically evaluate our integer compilation strategy. While we have demonstrated that a binary encoding exposes structure that knowledge compilation can exploit, it remains to be seen if this can improve the per-

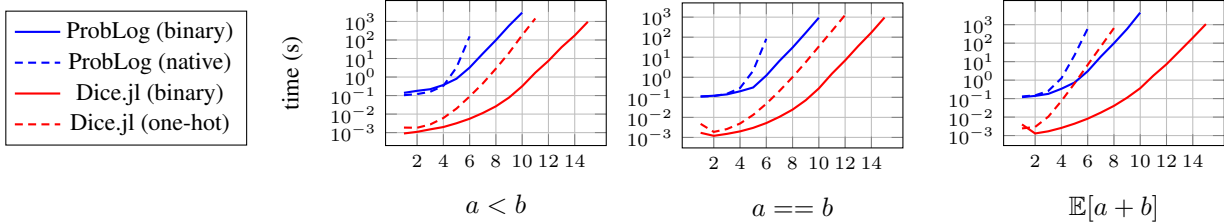


Figure 5: Time needed to compute the given operation on two random integers with varying bitwidth (x-axis).

formance of probabilistic programs. In addition, we have yet to show how our approach compares to other inference methods on larger, more complex arithmetic models. We seek to answer the following questions.

- 1) Does a binary encoding benefit existing knowledge compilation based languages?
- 2) Does our approach outperform those of existing PPLs that support exact discrete inference?

To this end, we have implemented our integer representation in `Dice.jl`, a PPL embedded in Julia that uses the same knowledge compilation approach as `Dice` [Holtzen et al., 2020]. In `Dice.jl` (binary), unsigned random integers are implemented using the strategy described in the previous section; signed random integers are additionally implemented as a natural extension. The language provides the syntax `discrete(..)` for arbitrary integer distributions and `uniform(..)` for uniform distributions, implemented using the algorithms in Section 3, as well as the operators `=`, `<`, `+`, `-`, `*`, `/`, and `%`. Simple operators are implemented as logical circuits, while more complex operators are implemented by composing simpler ones.<sup>1</sup>

Reported runtimes are a median over at least 5 runs; all experiments were run with a 1 hour timeout. A best-effort attempt was made for each language to implement benchmarks in a maximally performant manner. More experimental details are available in the supplementary material.

#### 4.1 IMPROVING KNOWLEDGE COMPILATION LANGUAGES WITH A BINARY ENCODING

In this section, we demonstrate the benefit a binary encoding brings to knowledge compilation based languages. We compare our integer representation method to the methods of the PPL `Dice` [Holtzen et al., 2020] and the probabilistic logic programming language `ProbLog` [Fierens et al., 2015], both of which use knowledge compilation as their approach to inference. We do this by comparing the time needed to compute the simple arithmetic operations  $a + b$ ,  $a < b$ , and

$a == b$  on random integers of width  $2^n$  for varying  $n$  from 1 to 15. Here, the runtime of each method corresponds to the time needed to compile and run inference on each representation, effectively measuring how well the knowledge compilation based inference can exploit the structure of each simple function. For addition, we use the expectation of the sum as our target computation to avoid an output distribution with an exponentially increasing support.

To allow for a fair comparison between `ProbLog`’s native integer representation and our binary representation, we implement equivalent `ProbLog` programs computing the arithmetic operations, one using native `ProbLog` encodings and one using a binary representation. Both programs can then be run using `ProbLog`, controlling for the specific knowledge compilation system. To compare with `Dice`’s native one-hot integer encoding, we implement both the one-hot encoding and our binary encoding in `Dice.jl`. We then run the simple arithmetic programs using both encodings.

The results of these experiments are presented in Figure 5. We can clearly see that our binary encoding outperforms the existing integer representation strategy used in each language; while at small distribution widths (on the order of  $2^4$ ), they are roughly comparable, our approach scales much better to larger integer distributions.

#### 4.2 COMPLEX ARITHMETIC MODELS

We also evaluated `Dice.jl` on more complex models involving distributions over integers. The models were taken from a variety of sources. These include examples involving integers from the existing PPL literature, various examples using continuous distributions adapted to a discrete space, natural modelling tasks using integers such as ranking and text manipulation, and traditional algorithms in a probabilistic setting. A short description of our baselines and their sources is given in the supplementary material.

As a point of comparison, we use two other PPLs supporting exact discrete inference. We identify two major classes of exact inference approaches used for discrete probabilistic programs: enumerative methods, which work by enumerating all paths through the program, and symbolic methods, which represent and compute the probability distribution through

<sup>1</sup>Our implementation and code for all experiments are available at <https://github.com/Juice-jl/Dice.jl/tree/arithmetic>.

Table 1: Runtimes in seconds for probabilistic models using integers in various PPLs.  $\times$  indicates a timeout (over 1 hour).

Benchmarks	Dice.jl (binary)	Dice.jl (one-hot)	WebPPL	Psi (DP)	Psi
book	<b>5.297</b>	$\times$	$\times$	$\times$	$\times$
tugofwar	<b>0.106</b>	2660.373	21.012	$\times$	$\times$
caesar-small	<b>0.041</b>	4.968	0.074	2.022	402.196
caesar-medium	0.239	39.518	<b>0.135</b>	12.505	$\times$
caesar-large	0.556	122.109	<b>0.227</b>	30.387	$\times$
ranking-small	<b>0.007</b>	0.025	0.83	103.572	$\times$
ranking-medium	<b>0.022</b>	0.077	$\times$	318.658	$\times$
ranking-large	<b>0.048</b>	0.150	$\times$	330.51	$\times$
radar1	<b>0.034</b>	0.664	118.002	394.525	2.517
floydwarshall-small	<b>0.009</b>	0.152	<b>0.009</b>	0.115	113.467
floydwarshall-medium	<b>0.515</b>	624.220	9.51	2792.14	$\times$
floydwarshall-large	<b>3.406</b>	$\times$	$\times$	$\times$	$\times$
linear extensions-small	<b>0.003</b>	0.004	0.016	0.351	5.153
linear extensions-medium	<b>0.007</b>	0.013	0.465	111.38	$\times$
linear extensions-large	<b>0.072</b>	0.164	162.009	$\times$	$\times$
triangle-small	<b>0.086</b>	102.544	3.693	616.746	482.14
triangle-medium	<b>0.455</b>	1123.171	28.354	$\times$	$\times$
triangle-large	<b>17.365</b>	$\times$	$\times$	$\times$	$\times$
gcd-small	2.876	$\times$	<b>0.189</b>	24.33	$\times$
gcd-medium	103.614	$\times$	<b>2.501</b>	467.581	$\times$
gcd-large	$\times$	$\times$	<b>46.626</b>	$\times$	$\times$
disease-small	7.91	$\times$	<b>1.093</b>	109.242	1009.848
disease-medium	764.212	$\times$	<b>327.545</b>	$\times$	$\times$
luhn-small	<b>0.039</b>	0.594	0.428	44.164	$\times$
luhn-medium	<b>4.575</b>	23.933	42.372	$\times$	$\times$

symbolic expressions. We compare against WebPPL [Goodman and Stuhlmüller, 2014] from the former category and Psi [Gehr et al., 2016] from the latter.

We also compare against a version of `Dice.jl` that uses a one-hot encoding of integer distributions as a proxy for existing knowledge compilation approaches; this comparison avoids language-specific differences in performance. Compiled BDD sizes for the programs are provided in the supplementary material as an additional metric.

Our results are summarized in Table 1. Many of our benchmark models can naturally be scaled to different sizes; they are implemented in small, medium, and large (corresponding to model size) variants to display the scaling behavior. Psi supports two exact inference algorithms: a default symbolic exact inference algorithm ("Psi") and its specialized dynamic programming inference algorithm ("Psi (DP)").

For the majority of benchmarks, our approach outperforms the current exact inference approaches, often achieving an orders-of-magnitude speedup. This result empirically validates the ability of our compilation strategy to exploit arithmetic structure in order to improve inference performance. We observe that the binary encoding outperforms the one-hot encoding with the same underlying knowledge compilation approach, demonstrating its superiority in exposing arithmetic structure.

We note that `Dice.jl` does not always outperform WebPPL, the enumeration based approach. We make special note of these examples. The caesar example introduces many

random integers but immediately makes an observation on their value, thereby reducing the enumeration task and making this tractable for all approaches. The disease example contains parametric distributions on integers; for example, a binomial distribution with parameter  $n$  distributed by a uniform distribution. These distributions have much less structure to exploit, and so our approach becomes essentially enumerative, but with additional overhead in compilation. The GCD example, which makes repeated use of the mod (%) operator, similarly has a harder to exploit structure. However, we can see in general that our approach scales well to larger examples and outperforms existing PPLs that support exact discrete inference.

## 5 ENABLING CONTINUOUS PRIORS WITH DISCRETE DISTRIBUTIONS

Previous sections presented an inference strategy for integer distributions allowing for the scaling of integer arithmetic. We now demonstrate an interesting application of these integers: representing the continuous Beta distribution in a discrete space. Continuous priors are an essential part of Bayesian reasoning. One particularly useful prior for discrete PPLs like `Dice.jl` is the *Beta prior*  $\text{Beta}(\alpha, \beta)$ , which is conjugate to the Bernoulli. The Beta distribution is continuous and thus not amenable to direct representation in `Dice.jl`. However, we observe that if a Beta prior for a Bernoulli random variable has integral parameters  $\alpha$  and  $\beta$ , then the posterior distribution is also a Beta with



integral parameters. This section explains how we use this observation to represent Beta distribution in `Dice.jl`.

Assume we have the following program where `Dice.jl` is extended to permit restricted classes of Beta priors, where the parameters  $\alpha$  and  $\beta$  must be constant integers:

---

```

1  $\theta$  = Beta(1, 2)
2  $x$  = flip( $\theta$ )
3 observe( $x$ )
4 return  $\theta$ 

```

---

We can perform inference for the posterior by exploiting well-known conjugacy results between Beta priors and Bernoullis. In particular, observing that  $x$  is true, as is done on Line 3 above, increases the *pseudocount*  $\alpha$  by 1, making the posterior for  $\theta$  become `Beta(2, 2)`. Similarly, observing that  $x$  is false increases the pseudocount  $\beta$  by 1.

To automate this approach in `Dice.jl`, we introduce program variables `A` and `B` to represent the pseudocounts  $\alpha$  and  $\beta$ , respectively. We then conditionally update these pseudocounts after each `flip`: increment  $\alpha$  if the `flip` returns true; otherwise increment  $\beta$ . Doing so ensures that later observations will have the desired effect on the pseudocounts.

The only remaining challenge is that discrete PPLs that employ knowledge compilation, like `Dice.jl`, only support `flips` whose parameters are constants, so the `flip( $\theta$ )` on Line 2 above is not supported. To encode it, we use the fact that  $\mathbb{E}[\text{Beta}(\alpha, \beta)] = \frac{\alpha}{\alpha + \beta}$ , so the `flip` on line 2 can be simplified to `flip( $\frac{A}{A+B}$ )`. Unfortunately, `Dice.jl` still does not support this construct, since  $\frac{A}{A+B}$  is not a constant. However,  $A + B$  is always a deterministic integer, since each observation increases it by exactly 1. Therefore, we can introduce a variable `T` representing  $A+B$  and encode `flip( $\frac{A}{A+B}$ )` as `uniform(0, T) < A` (where `uniform(0, n)` is the uniform distribution over the integers  $0, \dots, n - 1$ ). Finally, we observe that it is not necessary to maintain both variables `A` and `B`, since `B` is derivable from `A` and `T`. The final transformed version of the program is as follows:

---

```

1 A = 1, T = 3
2  $x$  = uniform(0, T) < A
3 A = if  $x$  then A+1 else A
4 T = T + 1
5 observe( $x$ )
6 return (A, T-A)

```

---

If we wish to draw another `flip` on the same Beta prior, we can simply repeat the code in lines 2-4 above. In this way, what we have actually implemented is a Beta-Bernoulli process via a Polya urn model - something analyzed in detail in probabilistic programs by Staton et al. [2018]. We also note that this representation strategy — an implementation of an urn model — can also be used for many other distributions in addition to the Beta [Mahmoud, 2008]. An example application of this Beta prior — learning Bayesian network

parameters — is given in the supplementary material.

## 6 RELATED WORK

*PPL Inference.* Knowledge-compilation-based PPLs are most closely-related to this work [Holtzen et al., 2020, De Raedt et al., 2007, Fierens et al., 2015, Saad et al., 2021, Pfanschilling et al., 2022]. All these languages stand to benefit from our new binary-encoding. Other PPLs perform exact discrete inference by eliminating discrete variables via enumeration or variable elimination; these approaches lose global structure and hence cannot exploit arithmetic structure as in our approach [Goodman and Stuhlmüller, 2014, Bingham et al., 2019]. Symbolic methods support integers by representing them as a symbolic formula or program [Gehr et al., 2016, Narayanan et al., 2016]; we believe that in principle it may be possible to adapt these symbolic representations to use a binary representation, but currently these systems do not. Recent work uses probability generating functions (PGFs) to represent (potentially unbounded) discrete distributions [Klinkenberg et al., 2023]. PGFs represent the distribution symbolically, but do not appear to be compatible with our strategy for binary encodings. Sampling based inference algorithms work very well for probabilistic program with continuous distributions, but do not exploit the global structure of integer arithmetic [Kantas et al., 2009, Hoffman and Gelman, 2014, Arouna, 2004, Wu et al., 2016]. Finally, there are algorithms that seek to efficiently model integer distributions with recursion [Knuth and Yao, 1976, Saad et al., 2020]; these approaches are orthogonal to ours, as we do not use recursion.

*Graphical models.* Probabilistic graphical model (PGM) based inference methods [Pfeffer, 2009, McCallum et al., 2009, Sanner and Abbasnejad, 2012, Koller and Friedman, 2009] support integers by treating them as categorical distributions. PGMs struggle to represent arithmetic: for instance, a CPT for adding two  $n$ -bit numbers requires  $O(2^n)$  entries.

## 7 CONCLUSION

We presented a strategy for encoding random integers in probabilistic programs via a binary representation, which allows arithmetic operations to be performed through standard Boolean circuits. When combined with the knowledge compilation approach to probabilistic inference, this strategy naturally exploits structure in arithmetic that current approaches do not account for. We showed empirically that this allows existing discrete PPLs to scale to significantly more complex probabilistic models. One interesting consequence is that we can now leverage conjugacy to represent the Beta distribution, a continuous distribution, in purely discrete programs.

## Acknowledgements

This work was funded in part by the DARPA PTG Program under contract number HR00112220005, NSF grants #IIS-1943641, #IIS-1956441, #CCF-1837129, #CCF-2220408, and a gift from RelationalAI. GVdB discloses a financial interest in RelationalAI.

## References

- Bouhari Arouna. Adaptive monte carlo method, a variance reduction technique. *10(1):1–24*, 2004. doi: 10.1515/156939604323091180.
- Eli Bingham, Jonathan P Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D Goodman. Pyro: Deep universal probabilistic programming. *The Journal of Machine Learning Research*, 20(1):973–978, 2019.
- Mark Chavira and Adnan Darwiche. Compiling Bayesian networks with local structure. In *IJCAI*, pages 1306–1312, 2005.
- Mark Chavira and Adnan Darwiche. On probabilistic inference by weighted model counting. *J. Artificial Intelligence*, 172(6-7):772–799, April 2008. ISSN 0004-3702. doi: 10.1016/j.artint.2007.11.002.
- Mark Chavira, Adnan Darwiche, and Manfred Jaeger. Compiling relational Bayesian networks for exact inference. *International Journal of Approximate Reasoning*, 42(1): 4–20, 2006.
- A. Darwiche and P. Marquis. A knowledge compilation map. *Journal of Artificial Intelligence Research*, 17:229–264, sep 2002. doi: 10.1613/jair.989. URL <https://doi.org/10.1613%2Fjair.989>.
- Luc De Raedt, Angelika Kimmig, and Hannu Toivonen. ProbLog : A Probabilistic Prolog and Its Applications to Link. *Proc. of IJCAI*, pages 2468–2473, 2007.
- Luc De Raedt, Angelika Kimmig, and Hannu Toivonen. Problog: A probabilistic prolog and its application in link discovery. In *Proceedings of IJCAI*, volume 7, pages 2462–2467, 2007.
- Joshua V Dillon, Ian Langmore, Dustin Tran, Eugene Brevdo, Srinivas Vasudevan, Dave Moore, Brian Patton, Alex Alemi, Matt Hoffman, and Rif A Saurous. Tensorflow distributions. *arXiv preprint arXiv:1711.10604*, 2017.
- Daan Fierens, Guy Van den Broeck, Joris Renkens, Dimitar Shterionov, Bernd Gutmann, Ingo Thon, Gerda Janssens, and Luc De Raedt. Inference and learning in probabilistic logic programs using weighted boolean formulas. *J. Theory and Practice of Logic Programming*, 15(3):358 – 401, 2015. doi: 10.1017/S1471068414000076.
- Timon Gehr, Sasa Misailovic, and Martin Vechev. Psi: Exact symbolic inference for probabilistic programs. In *International Conference on Computer Aided Verification*, pages 62–83. Springer, 2016.
- Andrew Gelman, Daniel Lee, and Jiqiang Guo. Stan: A probabilistic programming language for bayesian inference and optimization. *Journal of Educational and Behavioral Statistics*, 40(5):530–543, 2015.
- Noah D Goodman and Andreas Stuhlmüller. The Design and Implementation of Probabilistic Programming Languages. <http://dippl.org>, 2014. Accessed: 2022-10-26.
- Matthew D Hoffman and Andrew Gelman. The no-u-turn sampler: adaptively setting path lengths in hamiltonian monte carlo. *Journal of Machine Learning Research*, 15 (1):1593–1623, 2014.
- Steven Holtzen, Guy Van den Broeck, and Todd Millstein. Scaling exact inference for discrete probabilistic programs. In *Proc. ACM Program. Lang.*, OOPSLA 2020, pages 140:1–140:31. Association for Computing Machinery, 2020. doi: 10.1145/3428208.
- N. Kantas, A. Doucet, S.S. Singh, and J.M. Maciejowski. An overview of sequential monte carlo methods for parameter estimation in general state-space models. *IFAC Proceedings Volumes*, 42(10):774–785, 2009. ISSN 1474-6670. doi: <https://doi.org/10.3182/20090706-3-FR-2004.00129>. 15th IFAC Symposium on System Identification.
- Lutz Klöckner, Tobias Winkler, Mingshuai Chen, and Joost-Pieter Katoen. Exact probabilistic inference using generating functions. 2023.
- D. Knuth and A. Yao. *Algorithms and Complexity: New Directions and Recent Results*, chapter The complexity of nonuniform random number generation. Academic Press, 1976.
- D. Koller and N. Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Alexander K Lew, Marco F Cusumano-Towner, Benjamin Sherman, Michael Carbin, and Vikash K Mansinghka. Trace types and denotational semantics for sound programmable inference in probabilistic languages. *Proceedings of the ACM on Programming Languages*, 4(POPL): 1–32, 2019.
- H.P. Luhn. Computer for verifying numbers. U.S. Patent US2950048A, 1960.
- Hosam Mahmoud. *Pólya urn models*. Chapman and Hall/CRC, 2008.

- A McCallum, K Schultz, and S Singh. Factorie: Probabilistic programming via imperatively defined factor graphs. *Proc. of NIPS*, 22:1249–1257, 2009. ISSN 03643417.
- Christoph Meinel and Thorsten Theobald. *Algorithms and Data Structures in VLSI Design: OBDD-foundations and applications*. Springer Verlag, 1998. doi: 10.1007/978-3-642-58940-9.
- Praveen Narayanan, Jacques Carette, Wren Romano, Chungchieh Shan, and Robert Zinkov. Probabilistic inference by program transformation in hakaru (system description). In *International Symposium on Functional and Logic Programming - 13th International Symposium, FLOPS 2016, Kochi, Japan, March 4-6, 2016, Proceedings*, pages 62–79. Springer, 2016. doi: 10.1007/978-3-319-29604-3\_5. URL [http://dx.doi.org/10.1007/978-3-319-29604-3\\_5](http://dx.doi.org/10.1007/978-3-319-29604-3_5).
- Viktor Pfanschilling, Hikaru Shindo, Devendra Singh Dhami, and Kristian Kersting. Sum-product loop programming: From probabilistic circuits to loop programming. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning*, volume 19, pages 453–462, 2022.
- Avi Pfeffer. Figaro: An object-oriented probabilistic programming language. *Charles River Analytics Technical Report*, 137, 2009.
- Feras Saad, Cameron Freer, Martin Rinard, and Vikash Mansinghka. The fast loaded dice roller: A near-optimal exact sampler for discrete probability distributions. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 1036–1046. PMLR, 26–28 Aug 2020. URL <https://proceedings.mlr.press/v108/saad20a.html>.
- Feras A. Saad, Martin C. Rinard, and Vikash K. Mansinghka. SPPL: probabilistic programming with fast exact symbolic inference. In *PLDI 2021: Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Design and Implementation*, pages 804–819, New York, NY, USA, 2021. ACM. doi: 10.1145/3453483.3454078.
- Tian Sang, Paul Beame, and Henry Kautz. Performing bayesian inference by weighted model counting. In *Proceedings of the 20th National Conference on Artificial Intelligence - Volume 1, AAAI’05*, page 475–481. AAAI Press, 2005a. ISBN 157735236x.
- Tian Sang, Paul Beame, and Henry A Kautz. Performing bayesian inference by weighted model counting. In *AAAI*, volume 5, pages 475–481, 2005b.
- Scott Sanner and Ehsan Abbasnejad. Symbolic variable elimination for discrete and continuous graphical models. In *AAAI*, 2012.
- Sam Staton, Dario Stein, Hongseok Yang, Nathanael L. Ackerman, Cameron E. Freer, and Daniel M. Roy. The beta-bernoulli process and algebraic effects. 2018. doi: 10.4230/LIPICS.ICALP.2018.141.
- Jan-Willem van de Meent, Brooks Paige, Hongseok Yang, and Frank Wood. An introduction to probabilistic programming. *arXiv preprint arXiv:1809.10756*, 2018.
- Ingo Wegener. Bdds—design, analysis, complexity, and applications. *Discrete Applied Mathematics*, 138(1):229–251, 2004. ISSN 0166-218X. doi: [https://doi.org/10.1016/S0166-218X\(03\)00297-X](https://doi.org/10.1016/S0166-218X(03)00297-X). Optimal Discrete Structures and Algorithms.
- Yi Wu, Lei Li, Stuart Russell, and Rastislav Bodik. Swift: Compiled inference for probabilistic programming languages. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI’16*, page 3637–3645. AAAI Press, 2016. ISBN 9781577357704.