# An Implementation of List Successive Cancellation Decoder with Large List Size for Polar Codes

ChenYang Xia*, YouZhe Fan*, Ji Chen*, Chi-ying Tsui◇, ChongYang Zeng†, Jie Jin†, and Bin Li†

*◇Department of Electronic and Computer Engineering, the HKUST, Hong Kong

†Communications Technology Research Lab., Huawei Technologies, P. R. China

*{cxia, jasonfan, jchenbh}@connect.ust.hk, ◇eestui@ust.hk

†{zengchongyang, steven.jinjie, binli.binli}@huawei.com

*Abstract*—**Polar codes are the first class of forward error correction (FEC) codes with a provably capacity-achieving capability. Using list successive cancellation decoding (LSCD) with a large list size, the error correction performance of polar codes exceeds other well-known FEC codes. However, the hardware complexity of LSCD rapidly increases with the list size, which incurs high usage of the resources on the field programmable gate array (FPGA) and significantly impedes the practical deployment of polar codes. To alleviate the high complexity, in this paper, two low-complexity decoding schemes and the corresponding architectures for LSCD targeting FPGA implementation are proposed. The architecture is implemented in an Altera Stratix V FPGA. Measurement results show that, even with a list size of 32, the architecture is able to decode a codeword of 4096-bit polar code within 150 μs, achieving a throughput of 27Mbps.**

*Index Terms*—**polar codes, list successive cancellation decoding, FPGA implementation, low-complexity design.**

## I. INTRODUCTION

As an emerging class of *forward error correction* (FEC) codes with a provably capacity-achieving capability, *polar codes* [1] attract a lot of research interests recently. To decode the polar codes, *list successive cancellation decoding* (LSCD) [2], [3] was proposed, which outputs $\mathcal{L}$ (called *list size*) decoding paths by using $\mathcal{L}$ parallel *successive cancellation decodings* (SCDs) [4], [5]. By concatenating the polar codes with *cyclic redundancy check* (CRC) codes [3], [6] and using the checksums to choose the most reliable path from the list, LSCD with a large list size ($\mathcal{L} \geq 16$) achieves a similar or even better performance than other well-known FEC codes [3], such as low-density parity-check codes and turbo codes. However, this comes at a high hardware cost as the complexity scales with the list size $\mathcal{L}$. Thus, a low-complexity implementation of the corresponding LSCD is very desirable.

The existing LSCD architectures [7], [8], which were designed for a small or medium list size ($\mathcal{L} \leq 8$), are not suitable for a large list size due to their high complexity that is mainly due to two computational blocks. Firstly, several crossbars are required for executing the *list management* (LM) operation [7] and they have complexity of $O(\mathcal{L}^2)$. Secondly, a sorter with $2\mathcal{L}$ inputs is needed to compare and select the $\mathcal{L}$ best out of $2\mathcal{L}$ decoding paths to keep the list size to $\mathcal{L}$ during the decoding process. To reduce the logic delay, usually, a parallel sorter is used [7]. However, this parallel sorter has

$O(\mathcal{L}^2)$ comparators and hence dictates the clock frequency and incurs high hardware complexity.

Recently, two *field programmable gate array* (FPGA) implementations of LSCD architectures were presented in [9], [10], which can be used as the emulation platforms for evaluating the performance of polar codes. Due to the high complexity of LSCD, these platforms cannot support an LSCD of $\mathcal{L} > 4$. Moreover, to the best of our knowledge, hardware implementation for LSCD with $\mathcal{L} = 32$ has not been investigated in the literatures yet. In this work, we first propose two low-complexity decoding schemes for the LSCD with a large list size based on the analysis of the design constraints. Then, an LSCD architecture using these schemes is developed and implemented in an Altera FPGA. Measurement results show that our LSCD of $\mathcal{L} = 32$ decodes a 4096-bit polar code within 150 μs to achieve a 27Mbps throughput.

## II. PRELIMINARIES

### A. Code Construction

Considering a polar code with length $N = 2^n$. Its generator matrix, $\mathbf{F}^{\otimes n}$, is the $n^{th}$ Kronecker power of $\mathbf{F} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$. Source word $\mathbf{u}$ and code word $\mathbf{x}$ are two $N$-bit binary vectors related by $\mathbf{x} = \mathbf{u} \cdot \mathbf{F}^{\otimes n}$. The bits in $\mathbf{u}$ have different reliabilities. The indices of the $K$ most reliable bits compose the information set $\mathcal{A}$ while its complement $\mathcal{A}^c$ is called the frozen set. Accordingly, $u_i$s $(i \in \mathcal{A})$ are called information bits and are used to deliver message; while the rest are called frozen bits and fixed to 0. The code rate is thus defined as $R = K/N$. When an $r$-bit CRC code is concatenated, the last $r$ information bits are used to deliver the CRC checksums of the other $K - r$ information bits.

### B. List Successive Cancellation Decoding

As shown in Figure 1(a), the LSCD is made up of $\mathcal{L}$ copies of SCD operations (each described by the full binary tree) and the LM operations (represented by the squares).

The SCD operation is a depth-first traversal of the full binary tree with $n + 1$ stages which is also called a scheduling tree. The channel *log-likelihood ratios* (LLRs), $L_i = \log(\Pr(\mathbf{y}|0)) - \log(\Pr(\mathbf{y}|1)), i \in [0, N-1]$, are the inputs at the root node of the scheduling tree, where $\mathbf{y}$ is the channel
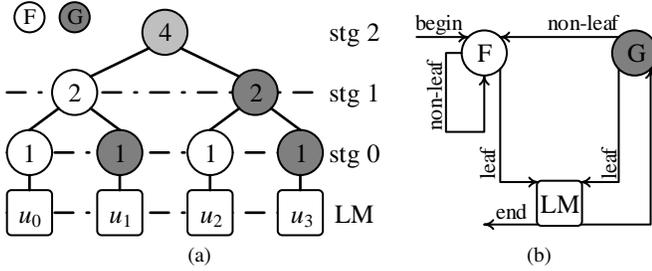
Figure 1: (a) scheduling tree of polar codes of $N = 4$ and (b) the corresponding state transfer diagram.

Table I: Crossbar complexity for 4096-bit polar codes on Altera 5SGXEA7N2F45C2 (available ALMs: 234,720)

| List size | 2 | 4 | 8 | 16 | 32 |
|---|---|---|---|---|---|
| Req. ALMs | 10,240 | 15,360 | 87,040 | 414,720 | 1,479,680 |

output on **x**. The left and right children of a node are called F- and G-node whose functions are[1]

$$L_F(L_a, L_b) = (\text{sgn}(L_a) \oplus \text{sgn}(L_b)) \cdot \min(|L_a|, |L_b|), \quad (1)$$

$$L_G(\hat{s}, L_a, L_b) = (-1)^{\hat{s}} L_a + L_b, \quad (2)$$

respectively, where $L_a$ and $L_b$ are the inputs of the both functions from the previous stage. $\hat{s}$ in (2) is a binary input called partial-sum, which is calculated from the bits already decoded up to the corresponding G-node. In the LSCD, all the $\mathcal{L}$ copies of SCDs are executed in parallel.

An LM is executed after a leaf node is reached by the SCDs. Assuming that after $\hat{u}_{i-1}$ is decoded, the list is full of $\mathcal{L}$ paths and each path has a different decoded sub-vector $[\hat{u}_0, ..., \hat{u}_{i-1}] \in \{0,1\}^i$. A *path metric* (PM), $\gamma_{i-1}^l$, is associated with each path to represent its reliability. When $\hat{u}_i$ is decoded, the LM of $\hat{u}_i$ is executed in two steps. First, each path is expanded into two with $\hat{u}_i$ instantiated to 0 and 1, respectively. For a path $l$, its *path metric update* (PMU) is[1]

$$\begin{cases} \gamma_i^{2l} = \gamma_{i-1}^l, & \text{if } \hat{u}_i = \Theta\left(\Lambda_i^l\right), \\ \gamma_i^{2l+1} = \gamma_{i-1}^l + \left|\Lambda_i^l\right|, & \text{if } \hat{u}_i = 1 - \Theta\left(\Lambda_i^l\right), \end{cases} \quad (3)$$

where $\gamma_i^{2l}$ and $\gamma_i^{2l+1}$ are the PMs of the two expanded paths and $\Lambda_i^l$ is the output LLR of the $i^{th}$ leaf node. The hard decision is made by $\Theta(x) = (x < 0)$. If the number of paths exceeds $\mathcal{L}$ after the path expansion, the *list pruning operation* (LPO) is executed to find the $\mathcal{L}$ smallest PMs and keep them as the survival paths. Note that if $i \in \mathcal{A}^c$, only one of the equations is executed and the LPO is not needed.

### C. Problems in the Existing LSCD Architectures

Based on the algorithms presented above, several LSCD architectures were proposed [7]–[11]. One common feature of them is that some $\mathcal{L} \times \mathcal{L}$ crossbars are needed to align the data in the $\mathcal{L}$ blocks of SCD hardware according to the LM results. Table I shows the synthesis results of the crossbars used in the architecture of [11]. Here, an 8-bit quantization is used for the LLRs. It can be seen the complexity scales far beyond $O(\mathcal{L})$

[1]The exact forms of these functions are non-linear. To have an efficient hardware implementation, approximate forms, (1) and (3), are used [4], [7].
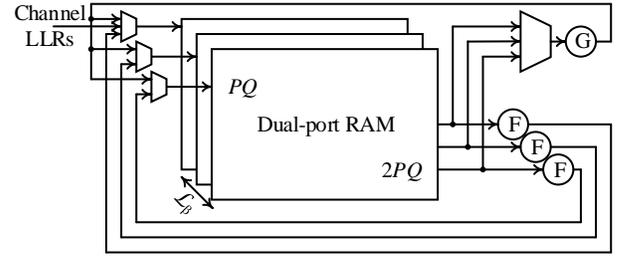


Figure 2: The structure of parallel-F serial-G computation.

for a given polar code and the required resources far exceed the logic resources, i.e. *adaptive logic modules* (ALMs), available in the FPGA for the LSCD with large list sizes.

Another issue is the implementation of the sorter which is required for finding the smallest $\mathcal{L}$ PMs after each path expansion. According to [7], the delay and complexity of radix-$2\mathcal{L}$ sorter are significantly increased with $\mathcal{L}$. The complexity of this sorter is further increased if a low-latency LM scheme, such as *multi-bit decoding* (MBD) [8], is used.

From the above discussion, implementing the architecture of LSCD with a large list size on hardware, especially in a resource-limited device such as an FPGA, is a very challenging task. In the following sections, we will present some schemes to reduce the complexity of LSCD.

### III. LOW-COMPLEXITY DECODING SCHEMES FOR LSCD

#### A. Parallel-F Serial-G Computation

From Section II-C, it is beneficial to avoid using crossbars in the architecture of LSCD with a large list size. A straightforward method is to integrate $\mathcal{L}$ blocks of LLR memories into a single memory and evaluating the SCD functions of each path serially. By doing so, the required operands are obtained by accessing the memory in the right locations. However, since the $\mathcal{L}$ SCDs are executed serially, the decoding latency is $\mathcal{L}$ times that of the traditional SCD. To reduce this latency, the following proposition related to the LSCD is used.

**Proposition 1.** *When the F-nodes are computed, the memories and PE arrays are one to one corresponding and the crossbars do not need to permute any data; only when the G-nodes are visited, crossbars need to permute the data from the memories.*

*Proof.* This can be easily proved from the state transfer diagram, as shown in Figure 1(b), which shows the execution order of the F-nodes, G-nodes and LMs. □

Based on Proposition 1, a *parallel-F serial-G* (PFSG) computation scheme is proposed. All the F-functions are calculated in parallel for all the paths as the crossbar is not needed in this situation and a direct connection between the corresponding memory and PE array can already support the calculations. In contrast, the G-functions of each path are serially evaluated to avoid using crossbars. As the latency for evaluating these two kinds of functions are the same in the SCD operation, the latency of LSCD using PFSG computation is $\frac{\mathcal{L}+1}{2}$ times that of the traditional SCD, which is reduced by almost one half comparing with that of straightforward mapping for large $\mathcal{L}$.

The corresponding PFSG structure is shown in Figure 2. Each block of RAM is implemented with a dual-port RAM with a $2PQ$-bit read port and a $PQ$-bit write port, where $Q$ is the number of quantization bits for the LLRs. $\mathcal{L} + 1$ groups of $P$ processing elements are used in this structure. One group is for the G-nodes, whose inputs are selected by an $\mathcal{L}$-to-1 multiplexer. The others are for the F-nodes, which can calculate the F-functions for $\mathcal{L}$ paths simultaneously.

It is noted when $\mathcal{L}$ is large, the utilization of RAMs is temporarily low as only the data from one of the $\mathcal{L}$ blocks of RAMs are valid in each cycle. So, in the real implementation, the number of blocks of RAMs can be reduced from $\mathcal{L}$ to $\mathcal{L}_\beta$, which is a power of 2, and each block of RAM stores the LLRs of $\mathcal{L}/\mathcal{L}_\beta$ paths. By choosing a proper $\mathcal{L}_\beta$, the balance between the complexity and the latency can be achieved.

### B. Low-Complexity List Management

In this section, a simplified LM operation of LSCD is proposed to reduce the computational complexity. To avoid the long latency brought by the G-nodes in the PFSG computation, the proposed method is based on the MBD. Specifically, the MBD simultaneously decodes all the $M$ bits of a sub-tree rooted at stage $m$, where $M = 2^m$. Let $\gamma_{\text{in}}$ be the PM of one survival path and $\gamma_{\text{out}}^{\text{MBD}}$ be the PM of one of its expanded paths, then the PMU of MBD is

$$\gamma_{\text{out}}^{\text{MBD}} = \gamma_{\text{in}} + \sum_{i=0}^{M-1} (v_i \oplus \Theta(L_i)) \cdot |L_i|, \qquad (4)$$

where $[L_0, ..., L_{M-1}]$ are the output LLRs at the root node of the sub-tree and $[v_0, ..., v_{M-1}] = [\hat{u}_0, ..., \hat{u}_{M-1}] \cdot F^{\otimes m}$. There are at most $2^M$ combinations of $v_i$s and hence at most $2^M$ paths are expanded from each survival path, which incurs a high complexity to the LPO even when $M$ is small.

To reduce the complexity, we combine one of our previously proposed algorithms, *selective expansion* (SE) [11], with the MBD. The SE efficiently reduces the number of the expanded paths by partitioning the information set $\mathcal{A}$ into an unreliable set $\mathcal{A}_u$ and a reliable set $\mathcal{A}_r$ based on the reliability of each information bit. The path expansions corresponding to the bits belonging to $\mathcal{A}_r$ do not need to be executed. We call the combined method *low-complexity list management* (LCLM). Supposing there are $M_u$ unreliable bits and $M_r$ reliable bits in a $M$-bit sub-tree. For a given set of values of the unreliable bits, the PMU of one of the expanded paths is calculated as

$$\gamma_{\text{out}}^{\text{LCLM}} = \min_{u_j \in \{0,1\}, \, j \in \mathcal{A}_r} (\gamma_{\text{out}}^{\text{MBD}}), \qquad (5)$$

where $\gamma_{\text{out}}^{\text{MBD}}$s are obtained from (4). The minimum in (5) is selected over the $2^{M_r}$ $\gamma_{\text{out}}^{\text{MBD}}$s. To expand each survival path, (5) needs to be calculated $2^{M_u}$ times as $2^{M_u}$ paths will be generated from the path expansion. Finally, LPO is used to select the $\mathcal{L}$ best paths from the $2^{M_u} \cdot \mathcal{L}$ expanded paths.

The LCLM expands fewer paths and hence achieves a lower complexity than the MBD. Also, Proposition 2 guarantees the decoding performance of LCLM is not worse than that of SE.

**Proposition 2.** *For a given $\gamma_{in}$ and $u_i$s $(i \in \mathcal{A}_u)$ in an $M$-bit tree, the updated PMs of LCLM and SE satisfy $\gamma_{out}^{LCLM} \leq \gamma_{out}^{SE}$.*
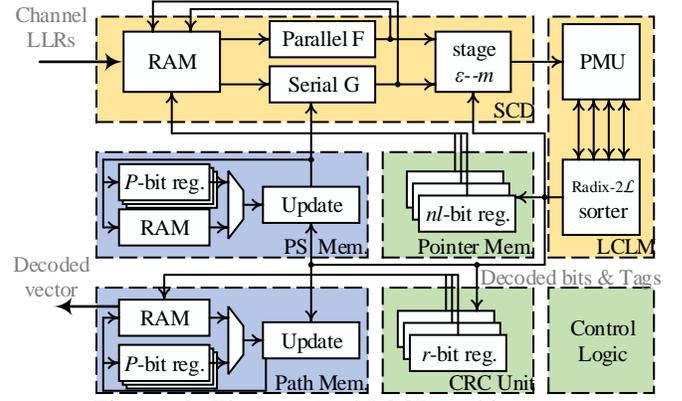


Figure 3: The proposed LSCD architecture.

*Proof.* For the given $u_i$s $(i \in \mathcal{A}_u)$, the corresponding $\gamma_{\text{out}}^{\text{SE}}$ equals to one of the $\gamma_{\text{out}}^{\text{MBD}}$s calculated by (4). So (5) ensures the validity of Proposition 2. □

An LSCD tries to find the best $\mathcal{L}$ paths with the locally smallest PMs. Proposition 2 ensures that the paths generated by the LCLM is not worse than those by the SE. Hence, the error performance of LCLM is at least as good as that of SE.

### IV. IMPLEMENTATION RESULTS

#### A. The Implementation of the Proposed LSCD Architecture

The implementation of the proposed LSCD architecture is shown in Figure 3, which mainly includes seven blocks.

The SCD module is used to compute the F- and G-nodes to obtain the LLR outputs of the stages higher than stage $m-1$. We further divide these stages into high stages (higher than a pre-determined stage $\epsilon$) and low stages (the rest). The high stages are calculated with the PFSG structure. The low stages are calculated in a parallel fashion as the PFSG brings a large latency overhead for these stages. Specifically, one memory is used to store the LLRs of all the paths and only one SCD hardware for the low stages is connected with it. Such structure is duplicated $\mathcal{L}$ times and the computations of all the paths can be executed simultaneously without a crossbar. The LCLM module receives the LLR outputs at stage $m$ from the SCD module. Here, a radix-$2\mathcal{L}$ parallel sorter is used. If $M_u > 1$ in a sub-tree, the $2\mathcal{L}$-to-$\mathcal{L}$ sorting is executed multiple times in serial to find the best $\mathcal{L}$ paths. The LCLM greatly reduces the number of expanded paths, so the latency for sorting is moderate. The outputs of the LCLM module include, for each path, $M$ decoded bits and a tag, indicating which survival path the expanded path is extended from.

The partial-sum memory and the path memory are used to store and update the partial-sums and the decoded vectors of the $\mathcal{L}$ paths, respectively. These memories are only activated when a G-node is calculated. Therefore the crossbars originally required in these two blocks in the existing architectures are not needed as the PFSG computation is used. A two-staged memory structure similar to the folded partial-sum network in [5] is used. The other parts, including the pointer memory, the CRC unit and the control logic, are similar to their counterparts in the existing architectures [7], [11].

Table II: The LSCD parameters for FPGA implementation.

| $N$ | $K$ | $r$[a] | CRC generator polynomial | $\mathcal{L}$ |
|---|---|---|---|---|
| 4096 | 2048 | 24 | 0x864cfb | 32 |

| $\mathcal{L}_\beta$ | $P$ | $Q$ | $Q_{\text{PM}}$ | $\eta$@SNR=2dB[b] | $m$ | $\epsilon$ |
|---|---|---|---|---|---|---|
| 4 | 128 | 8 | 9 | 0.3 | 2 | 3 |

[a] The effective code rate is $R = \frac{K-r}{N} = 0.494$.

[b] Following the method and notation of [11], $\eta$ determines $\mathcal{A}_u$ for SE.

Table III: Hardware usage of the LSCD architecture in FPGA.

| | ALMs | Registers | RAM blocks |
|---|---|---|---|
| LSCD usage | 67,211 | 31,247 | 1,122 |
| FPGA capacity | 234,720 | 939,000 | 2,560 |
| Utilization | 28.63% | 3.33% | 43.82% |

## B. Implementation and Measurement Results in the FPGA

To demonstrate the performance of the FPGA implementation of the above LSCD architecture, we implement it in an Altera Stratix V 5SGXEA7N2F45C2 FPGA. Table II summarizes the parameters of the target polar codes and the implemented decoder. The LSCD architecture is mapped on the FPGA with a clock frequency of 107MHz. The decoding latency of the LSCD is 16019 cycles for one codeword, translating into 149.71 $\mu$s under the target clock frequency. The hardware usage of our LSCD under the specified constraint is shown in Table III. Among all the resources, the RAM blocks (each with 20 kbits) have the highest usage, 22.44 Mbits, which is much higher than the theoretical value of about 2 Mbits. This is because the port width of one RAM block is limited. To guarantee the calculation parallelism, relatively wide port widths are used and multiple RAM blocks are then needed, leading to the high usage of RAM blocks.

Table IV compares our LSCD with other FPGA-based LSCD architectures in the literatures [9], [10]. Our architecture can support a longer code length and a much larger list size with even lower utilization of logic resources. The memory resources used per path are less than those of [9]. Though the memory usage of the architecture in [10] is lower, without any reported timing results, it is not easy to determine which architecture makes a better tradeoff between the complexity and the latency. The comparison results indicate the proposed low-complexity schemes are very efficient. At the same time, though the latency of our LSCD is supposed to scale linearly with the list size, the throughput degradation is less than linear. Also, for the other two architectures, it is not feasible to use them to implement LSCD with a large list size in an FPGA.

Finally, the measured *block error rate* (BLER) of the implemented LSCD is shown in Figure 4. For this measurement, an encoder and an additive white Gaussian noise channel are also implemented on-chip. As reference, the simulated BLER of the traditional LSCD with floating-point is also shown. It can be seen that our LSCD functions well and the performance degradation is less than 0.05dB at a BLER of $10^{-3}$, which is the target BLER of a typical cellular communication system. Also, comparing with the testing results presented in [9], a performance gain of about 0.8dB is achieved at this BLER.

## V. Conclusion

In this work, two low-complexity decoding schemes, namely PFSG and LCLM schemes, are proposed for the

Table IV: Comparison of the implementation results of several FPGA-based LSCD architectures.

| | Proposed | [9] | [10] |
|---|---|---|---|
| FPGA Device[c] | Altera Stratix V | Xilinx Kintex 7 | Altera Stratix V |
| $(N,\mathcal{L})$ | (4096,32) | (1024,4) | (1024,4) |
| ALMs(A)/LUTs(X)[d] | 67,211 | 142,961 | 101,160 |
| Registers | 31,247 | 19,795 | 13,544 |
| RAM (Mbits) | 22.440 | 4.404 | 0 |
| Clock rates (MHz) | 107 | 42.66 | N/A |
| Throughput (Mbps) | 27.35 | 115 | N/A |

[c] All the FPGAs are manufactured on 28nm process technology.

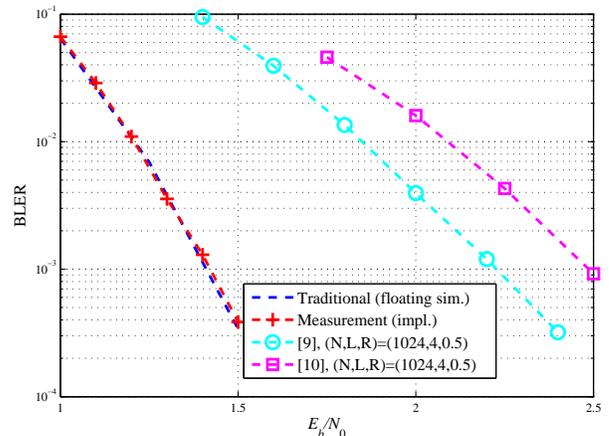[d] An ALM on Altera FPGA can be used as a 6-input LUT.



Figure 4: The BLERs of different LSCDs.

LSCD with a large list size, and the corresponding architecture for FPGA is developed and implemented. The measurement results show that the proposed LSCD ($\mathcal{L}$=32) has low hardware usage with negligible error performance degradation.

## References

[1] E. Arıkan. Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels. *IEEE Trans. Inf. Theory*, 55(7):3051–3073, 2009.

[2] K. Chen, K. Niu, and J. Lin. List successive cancellation decoding of polar codes. *IET Electron. Lett.*, 48(9):500–501, Apr 2012.

[3] I. Tal and A. Vardy. List decoding of polar codes. *IEEE Trans. Inf. Theory*, 61(5):2213–2226, May 2015.

[4] C. Leroux, A. J. Raymond, G. Sarkis, and W.J. Gross. A semi-parallel successive-cancellation decoder for polar codes. *IEEE Trans. Signal Process.*, 61(2):289–299, Jan 2013.

[5] Y. Fan and C.y. Tsui. An efficient partial-sum network architecture for semi-parallel polar codes decoder implementation. *IEEE Trans. Signal Process.*, 62(12):3165–3179, Jun 2014.

[6] B. Li, H. Shen, and D. Tse. An adaptive successive cancellation list decoder for polar codes with cyclic redundancy check. *IEEE Commun. Lett.*, 16(12):2044–2047, Dec 2012.

[7] A. Balatsoukas-Stimming, M. Bastani Parizi, and A. Burg. LLR-based successive cancellation list decoding of polar codes. *IEEE Trans. Signal Process.*, 63(19):5165–5179, Oct 2015.

[8] B. Yuan and K. K. Parhi. Low-latency successive-cancellation list decoders for polar codes with multibit decision. *IEEE Trans. VLSI Syst.*, 23(10):2268–2280, Oct 2015.

[9] C. Xiong, Y. Zhong, C. Zhang, and Z. Yan. An FPGA emulation platform for polar codes. In *Proc. IEEE Workshop on Signal Process. Syst. (SiPS)*, pages 148–153, 2016.

[10] X. Liang, J. Yang, C. Zhang, W. Song, and X. You. Hardware efficient and low-latency CA-SCL decoder based on distributed sorting. In *IEEE Conf. and Exhi. Glob. Telecom. (GLOBECOM)*, 2016.

[11] Y. Fan, C. Xia, J. Chen, C.y. Tsui, J. Jin, H. Shen, and B. Li. A low-latency list successive-cancellation decoding implementation for polar codes. *IEEE J. Sel. Areas Commun.*, 34(2):303–317, Feb 2016.