Recurrent Neural Networks With Column-Wise Matrix–Vector Multiplication on FPGAs

Zhiqiang Que[®], Member, IEEE, Hiroki Nakahara[®], Member, IEEE, Eriko Nurvitadhi, Member, IEEE,

Andrew Boutros, *Member, IEEE*, Hongxiang Fan¹⁰, *Student Member, IEEE*, Chenglong Zeng,

Jiuxi Meng, Student Member, IEEE, Kuen Hung Tsoi, Xinyu Niu, Member, IEEE, and Wayne Luk, Fellow, IEEE

Abstract—This article presents a reconfigurable accelerator for REcurrent Neural networks with fine-grained cOlumn-Wise matrix-vector multiplicatioN (RENOWN). We propose a novel latency-hiding architecture for recurrent neural network (RNN) acceleration using column-wise matrix-vector multiplication (MVM) instead of the state-of-the-art row-wise operation. This hardware (HW) architecture can eliminate data dependencies to improve the throughput of RNN inference systems. Besides, we introduce a configurable checkerboard tiling strategy which allows large weight matrices, while incorporating various configurations of element-based parallelism (EP) and vector-based parallelism (VP). These optimizations improve the exploitation of parallelism to increase HW utilization and enhance system throughput. Evaluation results show that our design can achieve over 29.6 tera operations per second (TOPS) which would be among the highest for field-programmable gate array (FPGA)-based RNN designs. Compared to state-of-the-art accelerators on FPGAs, our design achieves 3.7-14.8 times better performance and has the highest HW utilization.

Index Terms—Hardware Accelerator, long short-term (LSTM), recurrent neural network (RNN).

I. INTRODUCTION

RECURRENT neural networks (RNNs) have shown remarkable successes in sequence-to-sequence processing applications such as natural language processing (NLP) [1], speech-to-text [2], and video analysis [3]. Since low latency is key for a seamless user experience in such applications, efficient and real-time RNN acceleration is required. There are many RNN variants. Long short-term memory (LSTM) and gated recurrent unit (GRU) are the two most popular ones. To speed up RNN inferences, fieldprogrammable gate arrays (FPGAs) have been utilized in

Manuscript received August 5, 2021; revised October 21, 2021; accepted November 20, 2021. Date of publication December 29, 2021; date of current version February 7, 2022. This work was supported in part by the EPSRC, U.K., under Grant EP/P010040/1, Grant EP/S030069/1, Grant EP/N031768/1, and Grant EP/L016796/1; in part by Intel; and in part by Corerain. (*Corresponding author: Zhiqiang Que.*)

Zhiqiang Que, Hongxiang Fan, Jiuxi Meng, and Wayne Luk are with the Imperial College London, London SW7 2AZ, U.K. (e-mail: z.que@imperial.ac.uk; h.fan17@imperial.ac.uk; jiuxi.meng16@imperial.ac. uk; w.luk@imperial.ac.uk).

Hiroki Nakahara is with the Tokyo Institute of Technology, Tokyo 152-8550, Japan (e-mail: nakahara.h.ad@m.titech.ac.jp).

Eriko Nurvitadhi and Andrew Boutros are with Intel PSG, Hillsboro, OR 97123 USA (e-mail: eriko.nurvitadhi@intel.com).

Chenglong Zeng, Kuen Hung Tsoi, and Xinyu Niu are with Corerain, Shenzhen 518048, China.

Color versions of one or more figures in this article are available at https://doi.org/10.1109/TVLSI.2021.3135353.

Digital Object Identifier 10.1109/TVLSI.2021.3135353

• Tesla V100 • Brainwave • NPU • Plasticine • This work (Medium) • This work (Large) 75% 50% 25% 25% 512 1024 1536 2048 LSTM hidden vector size

Fig. 1. HW utilization of various LSTM implementations.

various scenarios [4]–[7], achieving lower latency and power consumption compared to CPUs and GPUs.

However, the recurrent nature and data dependency in the RNN computation results in undesired system stall until the required hidden vectors return from the full pipeline to start the next time-step (TS) calculation [6]. Besides, while deep pipelining can be utilized to enhance operating frequency, it increases stall penalty due to longer drain time. Moreover, inefficient tiling can leave hardware (HW) resources idle, resulting in low utilization. For example, the Brainwave [6] has six matrix-vector multiplication (MVM) "tile engines," each processing 400×40 matrices, so they have a peak capability of processing a 400×240 matrix in parallel. Any MVM that does not map to this dimension will leave some resources idle. Fig. 1 shows that Brainwave's HW utilization ranges from less than 1% to only about 50% for various LSTM models. Another implementation, Google's TPU, also suffers from low HW utilization and achieves an average utilization of 3.5% for LSTMs [8]. The Brainwave-like neural processor unit (NPU) [7] with a fine-grained zero-padding scheme achieves around 75% for a large LSTM. However, it still suffers from low utilization, especially from medium- to small-sized LSTMs which are commonly used in many applications [3], [4], [9]. LSTM models with small-to-medium sizes that have a large number of TSs are the most tangible examples that require dealing with lots of dependencies as well as the parallel task of MVMs [10]. With the need for high-performance

1063-8210 © 2021 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.



Fig. 2. MVMs. (a) Row-wise MVM. (b) Column-wise MVM.



Fig. 3. EP and VP with tiles shaded in red or blue. Two sets of (EP, VP) are shown in this example.

systems, it is essential to maximize the HW utilization to achieve the highest possible effective performance and energy efficiency.

This article proposes a novel latency-hiding HW architecture and a configurable checkerboard tiling strategy for RNN/LSTM models to increase HW utilization and enhance the throughput of RNN inference. First, we propose columnwise MVM for RNN/LSTM gates, which is able to eliminate their data dependencies. The column-wise block-striped decomposition of a matrix in MVM, as shown in Fig. 2(b), is an effective outer-product-based parallel method for processing MVM in high-performance computing. Recently, there are also some outer-product-based matrix-matrix multiplication accelerators [11], [12]. However, most of the previous FPGA-based RNN implementations focus on row-wise MVM as shown in Fig. 2(a). The proposed architecture can start the calculation of the next TS without waiting for the system pipeline to be drained, which means that the system can be fully pipelined without stalling.

Moreover, a novel configurable checkerboard tiling strategy is proposed which incorporates element-based parallelism (EP) and vector-based parallelism (VP) to boost inference throughput, as shown in Fig. 3. To support EP and VP, a new HW architecture that supports hybrid kernels is proposed, which combines multiplier-adder-tree and multiply-accumulate architectures. The architecture deploying many parallel multipliers followed by a large balanced adder tree is commonly used in FPGA-based RNN/LSTM accelerators [4], [6], [7], [13]. These designs are based on row-wise MVM. To support MVM column-wise processing, our design deploys many parallel multipliers followed by accumulators, as shown in Fig. 2(b). Furthermore, unlike our previous work [14] which is based on a fixed-size tiling approach, this work proposes a configurable tiling technique that supports various configurations of EP and VP to further improve the performance since different sizes of RNN models prefer different configurations of (EP, VP), as shown in Fig. 3.

Our results indicate that the proposed acceleration architecture is not only the fastest compared to the state of the art for a large LSTM model, but also much more suitable for a wide range of RNN models in terms of complexity. Hence, it performs much better for RNN models with different sizes as shown in Fig. 1. We make the following contributions.

- Novel column-wise MVMs for RNNs to eliminate data dependencies, increasing the HW utilization and system throughput.
- 2) A flexible checkerboard tiling strategy supporting EP and VP to exploit the available parallelism while increasing HW utilization and scalability. Besides, the (EP, VP) parameter space is comprehensively explored.
- A latency-hiding HW architecture with novel hybrid kernels and configurable adder-tree tail (CAT) units to support the proposed optimizations.
- 4) Experimental evaluation of the proposed approach and HW architecture, showing performance improvement of 3.7–14.8 times over state-of-the-art FPGA-based LSTM designs [6], [7], [15], with the highest HW utilization among them.

Relationship to Prior Publication: This article extends our previous work in [14], which covers a medium-scale design with a fixed EP and VP configuration, limiting the exploitation of parallelism and diversity in RNN applications. Moreover, the requirements to achieve high HW utilization for a small LSTM workload and a large LSTM workload are different. The choice of the mapping is highly dependent on the detailed layer shape, and a fixed mapping is not ideal for different models or even different layers of a model. This issue becomes severe when the number of PEs is large. An additional contribution of this article is a runtime configurable tiling strategy which supports various sets of EP and VP, along with novel CAT units, addressing the limitation described above. Besides, an extensive design space exploration is performed to identify favorable tiling block sizes for RNN models. Our approach can benefit a large-scale design involving 65 536 PEs, which has four times more PEs than our previous design [14], while still achieving high run-time HW utilization with the same benchmarks. Furthermore, additional information about the HW architecture and related work is included to provide a more detailed comparison. These optimizations lead to significant throughput increase and latency reduction over the previous design [14].

II. BACKGROUND AND PRELIMINARIES

A. Recurrent Neural Networks

LSTMs, one of the variations of RNNs, were initially proposed in 1997 [16]. This work follows the standard LSTM cell [3], [6], [7], [15]. An LSTM cell is composed of four LSTM gates and an LSTM tail. Each gate has a back-to-back MVM and activation operations, while the tail mainly involves vector element-wise operations. The arithmetic of a single-cell iteration is listed in the following equations:

$$i_{t} = \sigma(W_{t}[x_{t}, h_{t-1}] + b_{t}), \quad f_{t} = \sigma(W_{f}[x_{t}, h_{t-1}] + b_{f})$$

$$g_{t} = \tanh(W_{g}[x_{t}, h_{t-1}] + b_{u}), \quad o_{t} = \sigma(W_{o}[x_{t}, h_{t-1}] + b_{o})$$

$$c_{t} = f_{t} \odot c_{t-1} + i_{t} \odot g_{t}, \quad h_{t} = o_{t} \odot \tanh(c_{t})$$

where symbols σ , tanh, and \odot are, respectively, the sigmoid function, the hyperbolic tangent function, and element-wise multiplication. i_t , f_t , g_t , and o_t stand for the input, forget, input modulation, and output gates at TS t, respectively. The input modulation gate is often considered as a sub-part of the input gate. W represents the weight matrix for both input and hidden units. b represents bias. c_t is the internal memory cell state and h_t is the output of the cell, also called the hidden state, which is passed to the next TS or next layer. Another variant of LSTMs is the GRU. The forget and input gates of the LSTM are combined into a single "update gate" in the GRU. It has fewer parameters since it has only three gates. This study focuses on the optimization techniques for the standard LSTMs, but these techniques can be generalized to other RNN variants.

B. Related Work

There has been much previous work on FPGA-based implementations of persistent LSTM whose weights are stored in on-chip memory [6], [7], [13], [17]–[19]. There are also some previous studies of LSTM implementations storing weights in off-chip memory on FPGA devices [4], [20]–[23], which had been identified as the performance bottleneck. Liu et al. [24] utilized stochastic computing to improve the energy efficiency of RNNs. Khalil et al. [25] proposed reversible logic gates for low-power circuit designs for LSTMs. BLINK [26] utilizes bit-sparse data representation for LSTMs. POLAR [27] and BRDS [28] presented FPGA-based pipelined and overlapped architecture for dense and sparse LSTM inferences, respectively. Sun and Amano [29] proposed a multi-FPGAbased approach for accelerating deep RNNs. Kwon et al. [30] explored RNN partitioning strategies to achieve scalable multi-FPGA acceleration for large RNNs. Some of the previous studies [2], [31]–[36] are focusing on weight pruning and model compression to achieve good performance and efficiency. Wang et al. [37] and [38], Li et al. [39], and Yalamarthy et al. [40] presented a block circulant matrix to reduce LSTM inference weight matrices. These studies are orthogonal to our proposed strategy and architecture and can be complementary to our approach to achieve even higher throughput of RNN inferences on FPGAs.

Zhao et al. [15] proposed the cross-kernel optimization within RNN cells targeting Plasticine, which achieves 30 times higher performance compared to a GPU. The Brainwave [6] is a single-threaded SIMD architecture for persistent RNNs. It achieves more than an order of magnitude improvement in latency and throughput over GPUs. Nurvitadhi et al. [7] introduced a Brainwave-like neural processing unit (NPU) for RNNs. A TensorRAM is also proposed for large persistent data-intensive RNN models. Related work [19], [41] proposed a novel HW architecture of 2D-LSTM and investigated the trade off between precision and performance. Boutros *et al.* [42] presented the first performance evaluation of Intel's new AI-optimized FPGA, the Stratix 10 NX, with AI tensor blocks. All these RNN designs are based on row-wise MVM and suffer from data dependencies. Deploying the proposed latency-hiding HW architecture involving column-wise

TABLE I Summary of Parameters Used in Our Study

W	Weights matrix
w_n	All the weights of row n in W
w'_n	All the weights of column n in W
Hw	Number of Rows of weights matrix
Lw	Number of columns of weights matrix
x_t	The input vector x at timestep t
h_t	The hidden vector h at timestep t
$x_t[j]$	The element j in the input vector x at timestep t
$h_t[j]$	The element j in the hidden vector h at timestep t
Lx	Number of elements in the input vector x
Lh	Number of elements in the hidden vector h
NPE	Number of processing elements
EP	Element-based Parallelism
VP	Vector-based Parallelism
TS	Timestep
MVM_x	The MVM involving input vector x
MVM_h	The MVM involving hidden vector h

MVM and the proposed tiling strategy, REcurrent Neural networks with fine-grained cOlumn-Wise matrix-vector multiplication (RENOWN) achieves high HW utilization and throughput.

III. DESIGN AND OPTIMIZATION

This section covers data dependency analysis and optimizations targeting RNN designs. Some system parameters are defined in Table I.

A. Weight Matrix of LSTM Gates

In this design, the four matrices of *i*, *f*, *o*, *u* gates in LSTMs are combined into one large matrix since they share the same size. Thus, in one TS calculation of the LSTM, we only need to focus on one large matrix multiplied by one vector for the whole LSTM cell instead of four small matrices multiplying one vector. This is a generic optimization that can be applied to any MVMs that share the same input vector. Since each gate matrix has the size of $Lh \times (Lx + Lh)$, the size of the combined matrix is $(4 \times Lh) \times (Lx + Lh)$. Then we have $Hw = 4 \times Lh$ and Lw = Lh + Lx. Besides, the weights of the four LSTM gates are also interleaved in the final weights matrix. Therefore, the related elements in the result vector from four gates are adjacent and can be reduced easily via the element-wise operations in the LSTM-tail units.

B. Row-Wise MVM for RNNs

Conventional designs of MVM for RNNs are row-wise, and they have a major problem of being stalled when their pipelines are not fast enough to bring data back to the input for the next time step. They involve the entire vector of (x_t, h_{t-1}) and one or several entire rows of the weight matrix at a time, as shown in Fig. 4(a). This approach imposes additional



Fig. 4. Row-wise and column-wise MVM with LSTM data dependencies analysis. (a) Row-wise MVM for LSTMs. (b) Pipeline analysis of (a). (c) Column-wise MVM for LSTMs. (d) Pipeline analysis of (c).

stalling since the system has to wait for a newly computed hidden vector before starting the calculation of next timestep.

Data hazard exists since the whole new hidden vector h_t is required to start the new computation of x_{t+1} in the conventional MVM design for RNN/LSTM. It is mainly due to the data dependencies between the output from the current timestep and the vector for the next timestep. It indicates that the whole system pipeline needs to be drained to get the new computed hidden vector h_t before the new matrix–vector operations can start. As [6] mentions, RNN programs have a critical loop-carry dependence on the h_t vector. If the full pipeline cannot return h_t to the vector register file in time to start the next timestep, then the MVM unit will stall, as shown in Fig. 4(b). Therefore, pipeline latency is important. On the other hand, deep pipelining is often required to achieve a high operating frequency for designs. This makes it difficult to achieve a design with the best trade off.

C. Proposed Column-Wise MVM for RNNs

This work proposes a new technique that can alleviate this problem by calculating the matrix-vector operations in a column-wise fashion. At the beginning, only a few elements from the x_t vector are used while h_{t-1} is not touched, but all the elements in the corresponding columns of the weight matrix are involved to perform the operations, as illustrated in Fig. 4(c). To illustrate the idea, the number of the pipeline stages of the example system is 4, as shown in Fig. 4(d). However, the pipeline stages of a real system can be much larger. In addition, the figure shows that only one element in the x_t vector is used to perform the calculation for simplicity, but the actual number of the involved elements in each cycle depends on the architecture's parallelism requirements. The number of elements employed in this work is EP which is explored and fine-tuned in Section IV. The partial result vector is generated from the small dot product of the partial x_t vector and the corresponding weights. Then it is accumulated over multiple cycles to generate the final result vector. This way, the calculation of the new inference of (x_{t+1}, h_t) can start without waiting for the system pipeline to be drained to get h_t since it only needs a partial input vector. It indicates that the system can be fully pipelined without stalls, as shown in Fig. 4(d). Each hidden vector can finish the calculation in the shadow region of processing x_t before it is required. The stall happens in the calculation of each timestep, and

the total potential stalling cycles equals the design pipeline stages. When the LSTM workload is large, for example, with a layer Lh as 2048, the number of processing cycles of such a workload will be much larger than the number of stall cycles and then the benefit of the column-based MVM will be small, because the ratio of the stall cycles to the processing cycles is small. However, when the workload is small, for example, with a layer Lh as 256, the number of processing cycles of such a workload is also small. In such a scenario, reducing stall cycles is vital because the ratio of stalls is large. For example, if the number of processing cycles is the same as to the design pipeline stages, the HW utilization can be increased from 50% to 100% if all the stalls can be hidden.

One disadvantage has been observed that although the column-wise MVM only needs a partial input vector, it produces the output vector later than the row-wise MVM, because it needs to wait for all the columns to be processed to get the final accumulated vector before producing the output [43]. It seems that the succeeding HW units that depend on the output vector (e.g., those that perform activation functions and element-wise operations in the RNNs) would need to wait longer. In contrast, the row-wise MVM completes one subset of the output vector before moving to the next subset. Therefore, a subset of the final output vector is completed sooner than in a column-wise case. However, in the columnwise case, the succeeding units can get an entire output vector but not a subset. Although the column-wise architecture starts the subsequent processing later than the one using row-wise MVM, the number of succeeding units can be increased to match the output bandwidth and finish the whole calculation sooner than the row-wise case. Practically, we do not need to significantly increase the number of these units since they can process the whole vector over multiple cycles, until the next vector is produced by the MVM engine. Besides, the row-based approach may also involve multiple succeeding units to increase parallelism. Moreover, these units are much smaller compared to the MVM kernel engine that has lots of multipliers and adders. Thus, the extra HW area is negligible compared to the whole accelerator.

When the size of x_t vector is small and/or the size of h_t vector is large, the design may still stall because the cycles of processing x_t vector cannot completely hide the whole pipeline latency to get the h_t ready before it is needed. However, with the column-wise MVM, we can still continue to process the MVM of x_t and its corresponding weights when we are

waiting for h_t to be computed. Besides, when the input vector is small, an LSTM model would rarely require a significantly larger hidden vector.

D. Tiling and Parallelism

To further exploit the available parallelism, we introduce EP and VP in our design, as shown in Fig. 3. The matrix of weights is split into small tiles with a size of (EP, VP). In each cycle, the HW engine is able to process a tile of the weights matrix and a sub-vector of $[x_t, h_{t-1}]$ with a size of EP. EP and VP need to be determined carefully so that the number of cycles to process the x_t vector, given by (Lx/EP), is larger than the system latency to ensure that the computation of hidden vectors can be fully hidden by processing x_t vector. This number is small when EP is large and it may still result in system stalls. To increase system parallelism, VP is chosen to be as large as possible. However, the largest number of VP is Hw, which equals $4 \times Lh$, since there are only four gates in LSTM. In summary, the HW utilization and system throughput can be improved via balancing EP and VP.

A fixed configuration of (EP, VP) can bring low HW utilization. Brainwave [6] has six MVM tile engines, each processing 400×40 matrices. The NPU [7] also has a fixed configuration of four tiles, 120-wide dot product engines, and 40 lanes. Any MVMs that do not map well to these dimensions will leave some resources idle. Since RNNs are used for various tasks, RNN accelerators should support diverse configurations. This work proposes a novel HW architecture to support various sets of (EP, VP) since models with different sizes may prefer different optimal EP and VP configurations. Fig. 3 shows the basic idea with two sets of (EP, VP).

One option is to adopt the row-wise MVM while cascading the computation of MVM_x and MVM_h, which are the MVMs involving input vector and hidden vector, respectively, instead of a unified MVM. Cascading them with a row-wise fashion may also help to eliminate the data dependencies between current and next timestep calculations. However, it brings many issues. First, it introduces new data dependencies. Since the partial results from the MVM x and MVM hneed to be added together to generate the final output, if we calculate the MVM_x first, we have to store the partial results somewhere locally waiting for the partial results from MVM_h operations. This introduces a memory overhead of size 4Lh (for LSTMs). When Lh is large, this overhead is large because the design not only needs to cache these partial results, but also has to fetch them to accumulate them with the corresponding partial results from the hidden vector part. In our design, only the VP partial results are needed which means that only VP registers are needed. Besides, these partial results are available in the next cycle without prefetching, so they can be accumulated easily. Second, it brings difficulty in enhancing parallelism. Usually, the length of x and h are different, resulting in different computation loads for MVM x and MVM h. The input vector x is usually application-dependent, while the value of h can be selected by designers to meet application requirements and to reduce HW utilization. Zero padding may be required to support both MVM_x and MVM_h , which

causes inefficiency. Actually the height of the MVM_x and MVM_h are both 4Lh. The column-wise version of this computation has more parallelism than the row-wise version, which improves design efficiency. Both the designs [6], [7] separate and cascade the MVM_x and MVM_h, but they still suffer from low HW utilization as explained above.

IV. DESIGN SPACE EXPLORATION

The HW design space is characterized by the tiling block size of (EP, VP) and the number of processing elements (NPEs) after combining the configurations discussed previously. The effective performance varies with the tile size and the number of PEs. To figure out the optimal parameters of the system configuration for our in-depth analysis, we develop a cycle-accurate simulator to conduct the design space exploration. A greedy algorithm is proposed to explore design space. It starts with EP = 1 while the VP is given according to the system constraints, as shown in the following equation:

$$VP \le 4 \times Lh$$
 and $VP \le \frac{NPE}{EP}$. (1)

Practically, EP and VP should be as large as possible because this would maximize the potential parallelism and the system throughput. However, when EP increases, the number of cycles needed for processing the input vector (Lx) decreases so that the system may not have sufficient cycles to completely hide the processing of the hidden vector as discussed in Section III-C. In this exploration, the VP is set as large as possible, which is min $(4 \times Lh, (NPE/EP))$. Fig. 5 illustrates the exploration results for different sizes of LSTM models with Lh from 256 to 2048 with different colors and point shapes, using our HW design when the NPE is 4096, 16384, or 65536. Fig. 6 shows the feasible sets of (EP, VP) when NPE is 65536. The number of processing cycles determines the throughput of the system and the lower it is, the better the overall performance will be. As shown in Fig. 5, when EP is small, the number of processing cycles is high because the VP is constrained by (1) so that the number of effective PEs is less than NPE, which leads to severe underutilization. For instance, VP should be no larger than $4 \times Lh$ which is 1024 when targeting the LSTMs with the size of *Lh* being 256, illustrated by the blue line in Fig. 6. When EP increases, the number of processing cycles decreases until EP reaches these sweet spots. When EP is larger than the ones in sweet spots, processing cycles increase gradually. Note that EP = 1 is not shown in Fig. 5 (left) as its value is too large and it will make the sweet spot too small to be shown.

From our design space exploration result, the optimal configuration has an EP value between 4 and 16 when NPE = 16 382, and between 16 and 64 when NPE = 65 536. In these sweet spots, high parallelism can be achieved, which results in high system throughput. Note that for a given vector size, better performance and utilization can be obtained by adapting the (EP, VP) design parameters. As shown in Fig. 5 (left), the LSTM workload of 512 has lower latency with (EP, VP) of (32, 2048) than that of (16, 4096) as shown by the red line, while the workload of 1024 has lower latency with (16, 4096) than (32, 2048) as shown by the gray line. Different choices



Fig. 5. Number of processing cycles in our proposed column-wise approach for different values of EP, NPE, and model sizes. Different colors combining different point shapes represent different model sizes with *Lh* from 256 to 2048.



Fig. 6. Feasible sets of (VP, EP) when NPE equals 65 536. The sweet spot is set according to the sweet spot in Fig. 5 (left).



Fig. 7. Overview of the system.

of EP and VP impact the HW utilization and performance of the architecture when running RNN models of different sizes. There is a trade off between performance and design complexity with extra HW resources for supporting various values of (EP, VP). More details about this trade off are given in Section V.

V. HW IMPLEMENTATION AND OPTIMIZATION

This section presents our proposed HW architecture (see Fig. 7), based on the optimization techniques introduced above. It consists of the kernel units, an adapter unit, an activation function unit, and tail units.

A. Kernel Units

The architecture consists of VP kernel units, and each unit has EP processing elements (PEs), so the number of effective



Fig. 8. Various kernels. (a) Row-wise kernel. (b) Column-wise kernel. (c) Hybrid kernel.

PEs is $VP \times EP$. The VP and EP values are determined via the design space exploration described in detail in Section IV. In this design, each PE is one fully-pipelined multiplier. Fig. 8(c) shows the details of a computational kernel unit in our design. The architecture of kernel units, as shown in Fig. 8(a), which employs many parallel multipliers followed by an balanced adder-tree is commonplace in FPGA-based designs of RNNs [4], [6]. This architecture is for row-wise MVMs. The column-wise MVM is based on the architecture of many parallel multipliers followed by many parallel accumulators, as shown in Fig. 8(b), since the elements in the partial result vector are not related. To support EP, we propose a hybrid HW architecture that combines these two architectures. A small balanced adder tree is placed between the multipliers and the accumulators, as shown in Fig. 8(c). This small adder tree, which provides the summation of the products of EP multiplications, can help balance EP and VP for a proper shape of a tile. For example, when the VP is limited by $4 \times Lh$ as shown in Fig. 6, the design can increase the EP to enable more PEs to increase the throughput since the number of effective PEs is $VP \times EP$.

The hybrid kernel might look more complex than the row-wise or column-wise kernel but it does not consume more HW resources when targeting the same problem size. For example, if the number of the multipliers is N (for simplicity, N is a power of 2) for all three types of kernels and each hybrid kernel has a EP-to-1 small adder-tree, the design with row-wise kernels requires N multipliers and N - 1 adders (the design is supposed to be a fully pipelined one with a balanced tree), the design with column-wise kernels requires N multipliers and N accumulators, and the design with hybrid kernels needs N multipliers, $(N/\text{EP}) \times (\text{EP} - 1)$ adders, and (N/EP) accumulators. Because one EP-to-1 balanced adder-tree unit needs EP - 1 adders and there are (N/EP)



Fig. 9. Three modes of configurable four-input adder-tree tail with accumulators.

of such units, so there are $(N/\text{EP}) \times (\text{EP} - 1)$ adders in total. Generally, an accumulator is just an adder, so the hybrid kernel also needs $(N/\text{EP}) \times (\text{EP} - 1) + (N/\text{EP}) = N$ adders. Thus, the design with hybrid kernels has the same amount of HW resources as the one using pure column-wise kernels and just one more adder than the one using row-wise kernels. Note that some row-wise architectures also have accumulators after the large reduction tree since it can never guarantee to fully unroll the matrix dimension [7]. In such a case, the design using row-wise kernels also requires N adders.

B. Configurable Adder-Tree Tail

To support various versions of EP and VP, we design novel adder reduction based on a custom adder-tree with the CAT. With various EPs, the number of levels of the adder-tree needs to be changed correspondingly. If a fixed structure of the adder-tree is designed for a large EP (EP is also the number of the input elements for the adder-tree), the results from the last several levels of adder-tree can be used for small EPs to update the accumulators directly instead of entering next level adders, as shown in modes 2 and 3 in Fig. 9. For example, if EP is 64 and the number of input is also 64, with the proposed fourinput CAT, the number of the output elements of this adder-tree becomes 2 when mode 2 is enabled. Thus, instead of enabling a 64-to-1 adder-tree reduction, the design now has a 64-to-2 adder-tree which actually includes two 32-to-1 adder-trees. So, the configuration of (EP, VP) changes from (64, 1024) to (32, 2048). With mode 3, the same design can be enabled with (16, 4096). The detailed implementation can be achieved by additional accumulators while keeping the tree structure intact. However, the additional accumulators will result in resource overhead. This work proposes the CAT architecture to reuse the adders in the tail of the tree as the required accumulators with no extra adder components. The CAT with N-input (CAT-N) can be configured to update 1 to N accumulators when the data reach the last $\log_2(N)$ levels of the adder-tree. Fig. 9 shows the three modes using CAT-4. The results from the adder-tree can be used to update 1-4 accumulators. Fig. 10 shows the details of a CAT-4 unit. Different MUX settings configure CAT-4 to be one of the three modes shown in Fig. 9. For example, the red line in Fig. 10 shows the data flow when CAT-4 in mode 2. In our large-scale design, CAT-4 is sufficient since the sweet spot of EP is {16, 32, 64} according to design exploration for optimal system throughput.



Fig. 10. Details of CAT-4.

C. Other Units

The adapter converts the parallelism between kernels and tails. Then de-quantization (De-Quant) converts quantized values into fixed-point values to reduce HW resources. The σ /tanh unit performs the Sigmoid (σ) and hyperbolic tangent (tanh) functions. Both target programmable lookup tables of size 2048 [2], [7]. The LSTM-tail unit and GRU-tail unit mainly perform the element-wise operations. The output hidden vector (h_t) needs to be quantized before it can be used in the MVM kernels, so a Quant unit is deployed after the final output of Tail units as shown in Fig. 7.

D. Low-Precision Multiplications With DSP Block Sharing

Reducing the precision of operations in DNN inference accelerators can achieve high efficiency with little or no accuracy loss compared to floating point by fitting more multipliers per unit area. With careful retraining, low precision, even binarized RNNs can still have decent accuracy [19], [44], [45]. Rybalkin et al. [19] trained an LSTM model using 1-b weights and 2-b activations, which achieved a classification accuracy of 94% for OCR applications. Besides, narrow bit-width multiplications can be mapped efficiently onto lookup tables and DSPs. For example, Brainwave [6] deploys 96000 MACs on a Stratix 10 2800 FPGA by packing 2-b or 3-b multiplications into DSP blocks combined with cell-optimized soft logic multipliers and adders. Our fixed 8-b design has 16384 MACs. And our fixed 2-b design has 65 536 MACs, which deploys the same word length for multipliers as those in Brainwave targeting the same FPGA device for a fair comparison. The resource consumption of 65 535 fixed 8-b multipliers exceeds the available resources in the Stratix 10 2800 FPGA. Our column-wise MVM optimizations and fine-grained tiling strategy are applicable to multipliers of any numerical precision.

In current FPGAs, there are highly configurable DSP blocks which are often underutilized in implementing low-precision DNN designs. References [46] and [47] demonstrated methods to pack two 8-b multiplications into one Xilinx and Intel 18-b multiplier, respectively. Both methods require two multiplications to share one input operand. With the proposed columnwise MVM, one column of the weight matrix naturally shares the same element of the input vector, which helps us to pack four 8-b or ten 2-b multiplications into one DSP block on Intel FPGAs [47] to reduce the HW resources. Moreover, this would not be a restriction (and will come at lower cost) if we use a novel DSP similar to what was proposed in [48] and will be adopted in the next-generation Agilex devices [49].

VI. EVALUATION AND ANALYSIS

This section presents evaluation results and analysis of two generations of Intel FPGAs to show the scalability of the proposed RNN accelerator optimizations.

A. Experimental Setup

To evaluate our RNN design, we utilize the same LSTM and GRU workloads as the Brainwave design [6], the Brainwavelike NPU [7] and the Plasticine [15] for comparison. These workloads come from the DeepBench suite which is a set of micro-benchmarks containing representative layers from popular DNN models such as DeepSpeech [50]. These workloads are single layers with various sizes of hidden vectors and different timesteps (or sequence lengths) from 1 to 1500, as shown in Table III. This work takes the LSTMs with Lh =256 as small LSTM workloads, Lh = 512 or 1024 as medium ones, and Lh = 1536 or 2048 as large ones. Two generations of Intel FPGAs, an Arria 10 1150 (A10) and Stratix 10 2800 (S10), are evaluated and compared with previous work. Both run persistent LSTM/GRU of inference. Our proposed HW architecture is implemented in Verilog HW description language and implemented on the target A10 and S10 devices using Quartus Pro 18.1.

B. Resource Utilization

Table II shows the resource utilization of our designs with three configurations on FPGAs. We implement a small RENOWN with the configuration of (EP, VP) as (4, 1024) using an Arria 10 FPGA which has 4096 8-b multipliers in the MVM kernels. A medium RENOWN with the configuration of (EP, VP) as (16, 1024) is implemented using a Stratix 10 FPGA which includes 16384 8-b multipliers. A large RENOWN with 65536 2-b multipliers is also implemented on a Stratix 10 FPGA with the configurations of (EP, VP) as (16, 4096), (32, 2048), and (64, 1024). All our designs consume most of the FPGAs' available resources. Note that HW utilization is different from resource utilization and it reflects how often the HW computational units would be "not idle." Although we achieve a similar frequency to that reported in the Brainwave [6] and Intel-NPU [7] articles, we believe that further low-level optimizations can lead to higher frequencies for better performance. We leave that for future work since it has a limited impact on the conclusions in this article.

C. Performance and Efficiency Comparison

To illustrate the benefits of our proposed approach, some existing LSTM/GRU accelerator designs using the same benchmark are compared with ours in Table III. This table illustrates the latency, HW utilization, and throughput with various workloads under different numbers of hidden units (h) and TS. The HW utilization is the percentage of achieved tera operations per second (TOPS) to the peak performance for each layer. The DeepBench published results [15] on a modern NVIDIA Tesla V100 GPU with 16-b precision are also included. Some existing FPGA-based LSTM accelerator

TABLE II Resource Utilization

		ALMs	M20K	DSP	Freq.
Arria 10 1150 $\begin{pmatrix} Precision = 8-bit \\ NPE = 4096 \end{pmatrix}$	Avail. Used R. Util.	427,200 186,534 44%	2713 1,178 43%	1518 1,176 77%	259Mhz
Stratix 10 2800 $\begin{pmatrix} Precision = 8-bit \\ NPE = 16,384 \end{pmatrix}$	Avail. Used R. Util.	933,120 487,232 52%	11,721 10,061 86%	5760 4,368 76%	260Mhz
Stratix 10 2800 $\begin{pmatrix} Precision = 2-bit \\ NPE = 65,536 \end{pmatrix}$	Avail. Used R. Util.	933,120 768,406 82%	11,721 6,803 58%	5760 5,368 93%	250Mhz

TABLE III

PERFORMANCE COMPARISON OF DEEPBENCH INFERENCE FOR THE PREVIOUS WORK AND OUR DESIGNS

Benchmark		GPU	[6]	[7]	[15]	This Work (Medium)	This Work (Large)
NPE		-	96000	19200	12288	16384	65536
GRU h=512	latn. (ms) HW Util. Perf.	0.39	0.013	0.0015 21.7%	0.0004 30.9%	0.0006 64.1%	0.0004 24.8%
TS=1	(TOPS)	0.01	0.25	2.17	7.6	5.46	8.13
GRU h=1024 TS=1500	latn. (ms)	33.77	3.792	3.139	1.4430	2.59	0.879
	HW Util. Perf	1.8%	10.4%	60.2%	53.3%	85.5%	65.5%
	(TOPS)	0.56	4.98	6.01	13.1	7.28	21.5
GPU	latn. (ms)	13.12	0.951	1.454	0.7463	1.36	0.428
h=1536	HW Util.	2.6%	23.3%	73.2%	57.8%	91.4%	75.8%
TS=375	(TOPS)	0.81	11.17	7.30	14.2	7.79	24.8
CDU	latn. (ms)	17.70	0.954	-	1.283	-	0.695
h=2048	HW Util.	3.4%	41.2%	-	59.81%	-	82.8%
TS=375	Perf. (TOPS)	1.07	19.79	-	14.7	-	27.1
CDU	latn. (ms)	23.57	0.993	-	1.973	-	1.076
h=2560	HW Util.	4.0%	61.8%	-	61.0%	-	83.6%
TS=375	(TOPS)	1.25	29.69	-	15.0	-	27.4
LSTM h=256 TS=150	latn. (ms)	1.69	0.425	0.110	0.0419	0.033	0.029
	HW Util.	0.3%	0.8%	14.3%	15.5%	56.1%	16.7%
	Perf. (TOPS)	0.09	0.37	1.43	3.8	4.79	5.49
LOTM	latn. (ms)	0.60	0.077	0.027	0.0139	0.014	0.0061
h=512	HW Util.	0.6%	2.8%	38.8%	30.9%	85.9%	52.6%
TS=25	(TOPS)	0.18	1.37	3.89	7.6	7.33	17.2
ISTM	latn. (ms)	0.71	0.074	0.064	0.0292	0.054	0.015
h=1024	HW Util.	1.9%	2.8%	65.7%	58.6%	90.7%	86.6%
TS=25	(TOPS)	0.59	5.68	6.56	14.4	7.73	28.4
I STM	latn. (ms)	4.38	0.145	0.246	0.1224	0.236	0.066
h=1536 TS=50	HW Util.	1.4%	27.1%	76.9%	63.7%	94.1%	87.4%
	(TOPS)	0.43	13.01	7.67	15.4	8.02	28.7
LOTM	latn. (ms)	1.55	0.074	-	0.106	-	0.113
h=2048	HW Util.	3.4%	47.1%	-	64.3%	-	90.2%
TS=25	Pert. (TOPS)	1.08	22.62	-	15.8	-	29.6

designs are listed in Table IV. For a fair comparison, we only show the previous work with a detailed implementation of the LSTM system in this table. We show the FPGA chips, model storage, precision, NPEs, run-time frequency, throughput, power efficiency, and HW utilization.

The GPU is significantly underutilized even when cuDNN library API calls, since it is designed for throughput-oriented workloads, and it prefers BLAS level-3 (matrix–matrix) operations which are not common in RNN workloads [15]. Our design can provide promising latency under 3 ms for all Deepbench RNN layers at batch size of 1, reaching up to 29.6 effective TOPS for a large LSTM workload (h = 2048) which is the largest reported performance in all these LSTM



Fig. 11. Performance comparison. (a) Performance speedup for the various LSTM dimensions. (b) HW utilization of various GRUs compared to prior work.

designs as shown in Tables III and IV. Our work achieves 27.4–95.8 times higher performance than the Tesla V100, as shown in Fig. 11(a). The performance of the specific case (h = 512) is an approximate two orders of magnitude advantage over the Tesla V100.

Our experiments show that the utilization is low with small RNN applications that are composed of sequences of small MVMs due to small hidden unit sizes and large number of TSs. However, with our proposed optimizations, we can get higher throughput and HW utilization than the counterparts using a similar number of PEs. With a similar number of PEs to [7], our RENOWN (Medium) achieves up to 94.1% HW utilization which is the highest with respect to state-of-the-art implementations on FPGAs, as shown in Figs. 1 and 11(b). Achieving high utilization using a small number of PEs is easier than using a large number of PEs. In our RENOWN (Large) design, we use around 31.7% fewer multipliers than [6], while achieving 30.7% higher performance than [6] when targeting large LSTM workloads. When targeting small LSTM models, our design achieves 14.8 times higher performance. This is due to better HW utilization which comes from our optimizations incorporating novel column-wise MVM and fine-grained tiling strategy. With a large number of PEs, our RENOWN can still achieve up to 90.2% HW utilization which is higher than the other work. When targeting small LSTMs and GRUs, RENOWN (large) has a similar utilization as [7] and [15]. However, the number of PEs is, respectively, 3.41 times and 5.3 times more than those in [7] and [15]. The configurable tiling of RENOWN (Large) also results in an additional up to



Fig. 12. Performance speedup due to configurable tiling.

27% higher system throughput as shown in Fig. 12. The figure shows the speedup of the system throughput compared to the lowest one among them with a given EP.

The previous designs of [19] and [41] in Table IV explore various low-precision designs, showing their scalability to different bitwidth designs. Most of the other designs in Table IV only support one bitwidth, while our work demonstrates both 8-b and 2-b designs which show the scalability of our architecture to cover different bit-width designs. Our implementations are produced using the Verilog templates with configurable parameters for each HW module instance. Generally, Rybalkin et al. [19] and Rybalkin and Wehn [41] are more scalable in bitwidth since they do not use the DSP hard blocks on FPGAs for the computational kernels. The FPGA DSP block has a fixed bit-width and it needs a tailor-made wrapper for packaging several small multipliers into one DSP HW block. Besides, Rybalkin et al. [19] and Rybalkin and Wehn [41] target a particular model which has four parallel and independent LSTM layers. These four independent LSTM layers can be scheduled in an interleave manner for an HW engine to alleviate the issues of recurrent dependencies.

Some of the designs target smaller FPGAs [19], [41] than Stratix 10. They have fewer logic resources and DSPs but they also have a smaller TDP (thermal design power) power consumption. We follow the convention in all related articles and use GOP/S for performance comparison. In addition, we also use GOPS/W and power efficiency, so that we are not taking advantage of the FPGA chip size when compared to the designs on small chips, such as the designs in [19] and [41].

RENOWN Overall, our (Medium) provides over 1.05–3.35 times higher performance and 1.22–3.92 times higher HW utilization than the state-of-the-art design [7], as shown in Table III. In addition, the RENOWN (Large) achieves 3.7-14.8 times better performance than state-of-theart FPGA-based LSTM designs [6], [7], [15]. This article focuses on minimizing latency and maximizing throughput by increasing the HW utilization. The results show flexible customizability of our architecture for different scenarios. The column-wise approach exposes the most parallelism while minimizing stalls due to data dependencies.

To minimize latency, our design places the model weights onto the on-chip memory, achieving high memory read

	FPGA17 ESE [2]	FPL18- [13]	ISCA18- BW [6]	FCCM19- NPU [7]	FCCM19- NPU [7]	JSPS20- [19]	FPGA20- [41]	TVLSI20- [27]	This work (Medium)	This work (Large)
FPGA chips	Kintex KU060	Zynq ZU7EV	Stratix10 GX2800	Stratix10 GX2800	Stratix10 GX2800	Zynq ZU9EG	Zynq ZU9EG	Zynq ZU6EG	Stratix10 GX2800	Stratix10 GX2800
Model Storage	-	on-chip	on-chip	on-chip	on-chip	on-chip	on-chip	on-chip	on-chip	on-chip
Precision (bits)	12 fixed	1-8 fixed	BFP-1s5e2m	4 fixed	8 fixed	2 fixed	4/8 fixed	16 fixed	8 fixed	2 fixed
DSP Avail.	2,760	1,728	5,760	5,760	5,760	2,520	2,520	1,973	5,760	5,760
DSP Used	1,504	-	5,245	4,880	5,600	6	64	-	4,368	5,368
NPE	1024	-	19,200	96,000	38,400	-	-	-	16,384	65,536
Freq. (MHz)	200	240	250	255	260	240	240	385	260	250
Power (W)	41	-	125	65.3	67.6	15.5	15.5	0.95	62	98
Perf. (GOPS)	282	1,833	370 ^a 22,620	<u>-</u> 14,112	1,431 ^a 7,980	3,618	3,072	18	4,790 ^a 8,015	5,489 ^a 29,574
Power Effi. (GOPS/W)	6.87	-	180	216	118	234	199	19	129	302
LSTM HW Utilization	-	-	$0.8\%^{a}$ 47.1%	-	14.3% ^a 76.9%	-	-	-	56.1% ^a 94.1%	16.7% ^a 90.2%

TABLE IV Comparison With Previous Implementations of LSTM on FPGAs

^a When targeting a small LSTM model (*Lh*=256).



Fig. 13. Power consumption of the accelerator with the decomposition for each of the major underlying blocks.

bandwidth suitable for real-time services. However, some large-scale RNNs recently emerge with large on-chip memory requirements. Our design adopts a scale-out network of utilizing multiple accelerators like [6], [30] which partition the design to multiple FPGAs to address the challenges of fast-growing RNN models where weights exceed on-chip memory capacity on a single FPGA. Some cloud-based services are able to tolerate a slightly longer latency of response. It means a small amount of batching can be employed if necessary. This benefits the GPU-based RNN inferences [6]. In [13], [42], [44], and [51], the batching technique is used to improve the HW throughput and utilization for LSTM inferences. Since our design executes a single input at a time, increasing batch size does not affect the utilization. Thus, our architecture's utilization is not affected with or without batching. Furthermore, some designs use binarized datapaths [13], [19], [41] for LSTMs with negligible or no effect on accuracy. Utilizing very low precision, for example, binary, is orthogonal to our proposed approach which transforms computation to eliminate data dependencies. Reducing precision can be combined with our approach to achieve even higher performance and efficiency.

Fig. 13 illustrates the power consumption of the proposed RENOWN (Medium) design. The power is estimated by the Intel early power estimator (EPE) tool and verified by the Intel Quartus Power Analyzer. We note that this work considers only the chip power, for a fair comparison. The bar chart in Fig. 13(left) shows the decomposition for each of the major underlying FPGA components. The static power consumption of the device is 7.16 W. The dynamic power consumption of the accelerator engine unit is the largest which is over 50%. The pie chart in Fig. 13(right) shows the dynamic power consumption of each major unit of the accelerator engine unit. The Buffer Units which store the weights and input-output data have the largest power consumption which reaches nearly half of the power consumed by the whole engine unit. The Kernel Units also consume nearly half of the power while the Tail units only consume less than 1% power.

VII. CONCLUSION AND FUTURE WORK

This article proposes fine-grained tile-based column-wise MVMs with a novel latency-hiding HW architecture for RNN/LSTM to improve both HW utilization and design throughput. The proposed accelerator has been implemented using Arria 10 and Stratix 10 FPGAs, achieving superior performance and power efficiency compared to prior state-ofthe-art implementations. Further research includes combining our approach with binarized RNNs, pruning techniques and in-memory computing, and automating it to support rapid development of efficient FPGA-based RNN/LSTM designs.

REFERENCES

- Y. Goldberg, "A primer on neural network models for natural language processing," J. Artif. Intell. Res., vol. 57, pp. 345–420, Nov. 2016.
- [2] S. Han et al., "ESE: Efficient speech recognition engine with sparse LSTM on FPGA," in Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays, Feb. 2017, pp. 75–84.

- [3] J. Donahue et al., "Long-term recurrent convolutional networks for visual recognition and description," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2015, pp. 2625–2634.
- [4] Y. Guan, Z. Yuan, G. Sun, and J. Cong, "FPGA-based accelerator for long short-term memory recurrent neural networks," in *Proc. 22nd Asia South Pacific Design Autom. Conf. (ASP-DAC)*, Jan. 2017, pp. 625–634.
- [5] Z. Sun *et al.*, "FPGA acceleration of LSTM based on data for test flight," in *Proc. IEEE Int. Conf. Smart Cloud (SmartCloud)*, Sep. 2018, pp. 1–6.
 [6] J. Fowers *et al.*, "A configurable cloud-scale DNN processor for real-
- [6] J. Fowers et al., A cominguiable clouescale Divide processor for real-time AI," in *Proc. ACM/IEEE 45th Annu. Int. Symp. Comput. Archit. (ISCA)*, Jun. 2018, pp. 1–14.
 [7] E. Nurvitadhi et al., "Why compete when you can work together:
- [7] E. Nurvitadhi et al., "Why compete when you can work together: FPGA-ASIC integration for persistent RNNs," in Proc. IEEE 27th Annu. Int. Symp. Field-Program. Custom Comput. Mach. (FCCM), Apr. 2019, pp. 199–207.
- [8] N. P. Jouppi *et al.*, "In-datacenter performance analysis of a tensor processing unit," in *Proc. 44th Annu. Int. Symp. Comput. Archit.*, 2017, pp. 1–12.
- [9] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," 2014, arXiv:1409.2329.
- [10] R. Yazdani, O. Ruwase, M. Zhang, Y. He, J.-M. Arnau, and A. Gonzalez, "LSTM-sharp: An adaptable, energy-efficient hardware accelerator for long short-term memory," 2019, arXiv:1911.01258.
- [11] S. Pal et al., "OuterSPACE: An outer product based sparse matrix multiplication accelerator," in Proc. IEEE Int. Symp. High Perform. Comput. Archit. (HPCA), Feb. 2018, pp. 724–736.
- [12] Z. Zhang, H. Wang, S. Han, and W. J. Dally, "SpArch: Efficient architecture for sparse matrix multiplication," in *Proc. IEEE Int. Symp. High Perform. Comput. Archit. (HPCA)*, Feb. 2020, pp. 261–274.
- [13] V. Rybalkin, A. Pappalardo, M. M. Ghaffar, G. Gambardella, N. Wehn, and M. Blott, "FINN-L: Library extensions and design trade-off analysis for variable precision LSTM networks on FPGAs," in *Proc. 28th Int. Conf. Field Program. Log. Appl. (FPL)*, Aug. 2018, pp. 89–897.
- [14] Z. Que et al., "Optimizing reconfigurable recurrent neural networks," in Proc. IEEE 28th Annu. Int. Symp. Field-Program. Custom Comput. Mach. (FCCM), May 2020, pp. 10–18.
- [15] T. Zhao, Y. Zhang, and K. Olukotun, "Serving recurrent neural networks efficiently with a spatial accelerator," 2019, *arXiv:1909.13654*.
- [16] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural Comput., vol. 9, no. 8, pp. 1735–1780, 1997.
- [17] J. C. Ferreira and J. Fonseca, "An FPGA implementation of a long short-term memory neural network," in *Proc. Int. Conf. ReConFigurable Comput. FPGAs (ReConFig)*, Nov. 2016, pp. 1–6.
- [18] V. Rybalkin, N. Wehn, M. R. Yousefi, and D. Stricker, "Hardware architecture of bidirectional long short-term memory neural network for optical character recognition," in *Proc. Design, Autom. Test Eur. Conf. Exhib. (DATE)*, Mar. 2017, pp. 1394–1399.
- [19] V. Rybalkin, C. Sudarshan, C. Weis, J. Lappas, N. Wehn, and L. Cheng, "Efficient hardware architectures for 1D-and MD-LSTM networks," *J. Signal Process. Syst.*, vol. 92, no. 11, pp. 1219–1245, 2020.
- [20] A. X. M. Chang, B. Martini, and E. Culurciello, "Recurrent neural networks hardware implementation on FPGA," 2015, arXiv:1511.05552.
- [21] Y. Guan et al., "FP-DNN: An automated framework for mapping deep neural networks onto FPGAs with RTL-HLS hybrid templates," in Proc. IEEE 25th Annu. Int. Symp. Field-Program. Custom Comput. Mach. (FCCM), Apr. 2017, pp. 152–159.
- [22] Z. Que, Y. Zhu, H. Fan, J. Meng, X. Niu, and W. Luk, "Mapping large LSTMs to FPGAs with weight reuse," J. Signal Process. Syst., vol. 92, no. 9, pp. 965–979, Sep. 2020.
- [23] N. Park, Y. Kim, D. Ahn, T. Kim, and J.-J. Kim, "Time-step interleaved weight reuse for LSTM neural network computing," in *Proc. ACM/IEEE Int. Symp. Low Power Electron. Design*, Aug. 2020, pp. 13–18.
- [24] Y. Liu, L. Liu, F. Lombardi, and J. Han, "An energy-efficient and noise-tolerant recurrent neural network using stochastic computing," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 27, no. 9, pp. 2213–2221, Sep. 2019.
- [25] K. Khalil, B. Dey, A. Kumar, and M. Bayoumi, "A reversiblelogic based architecture for long short-term memory (LSTM) network," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, May 2021, pp. 1–5.
- [26] Z. Chen, G. J. Blair, H. T. Blair, and J. Cong, "BLINK: Bit-sparse LSTM inference kernel enabling efficient calcium trace extraction for neurofeedback devices," in *Proc. ACM/IEEE Int. Symp. Low Power Electron. Design*, Aug. 2020, pp. 217–222.

- [27] E. Bank-Tavakoli, S. A. Ghasemzadeh, M. Kamal, A. Afzali-Kusha, and M. Pedram, "POLAR: A pipelined/overlapped FPGA-based LSTM accelerator," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 28, no. 3, pp. 838–842, Mar. 2020.
- [28] S. A. Ghasemzadeh, E. B. Tavakoli, M. Kamal, A. Afzali-Kusha, and M. Pedram, "BRDS: An FPGA-based LSTM accelerator with rowbalanced dual-ratio sparsification," 2021, arXiv:2101.02667.
- [29] Y. Sun and H. Amano, "FiC-RNN: A multi-FPGA acceleration framework for deep recurrent neural networks," *IEICE Trans. Inf. Syst.*, vol. E103.D, no. 12, pp. 2457–2462, 2020.
- [30] D. Kwon, S. Hur, H. Jang, E. Nurvitadhi, and J. Kim, "Scalable multi-FPGA acceleration for large RNNs with full parallelism levels," in *Proc.* 57th ACM/IEEE Design Autom. Conf. (DAC), Jul. 2020, pp. 1–6.
- [31] Z. Chen, A. Howe, H. T. Blair, and J. Cong, "CLINK: Compact LSTM inference kernel for energy efficient neurofeedback devices," in *Proc. Int. Symp. Low Power Electron. Design*, Jul. 2018, pp. 1–6.
- [32] C. Gao, D. Neil, E. Ceolini, S.-C. Liu, and T. Delbruck, "DeltaRNN: A power-efficient recurrent neural network accelerator," in *Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays*, Feb. 2018, pp. 21–30.
- [33] J. Wu, F. Li, Z. Chen, and X. Xiang, "A 3.89-GOPS/mW scalable recurrent neural network processor with improved efficiency on memory and computation," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 27, no. 12, pp. 2939–2943, Dec. 2019.
- [34] S. Cao et al., "Efficient and effective sparse LSTM on FPGA with bankbalanced sparsity," in Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays, Feb. 2019, pp. 63–72.
- [35] R. Shi, J. Liu, H. K.-H. So, S. Wang, and Y. Liang, "E-LSTM: Efficient inference of sparse LSTM on embedded heterogeneous system," in *Proc.* 56th Annu. Design Autom. Conf., Jun. 2019, pp. 1–6.
- [36] G. Nan, C. Wang, W. Liu, and F. Lombardi, "DC-LSTM: Deep compressed LSTM with low bit-width and structured matrices," in *Proc. IEEE Int. Symp. Circuits Syst. (ISCAS)*, Oct. 2020, pp. 1–5.
- [37] Z. Wang, J. Lin, and Z. Wang, "Accelerating recurrent neural networks: A memory-efficient approach," *IEEE Trans. Very Large Scale Integr.* (VLSI) Syst., vol. 25, no. 10, pp. 2763–2775, Oct. 2017.
- [38] S. Wang et al., "C-LSTM: Enabling efficient LSTM using structured compression techniques on FPGAs," in Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays, Feb. 2018, pp. 11–20.
- [39] Z. Li *et al.*, "E-RNN: Design optimization for efficient recurrent neural networks in FPGAs," in *Proc. IEEE Int. Symp. High Perform. Comput. Archit. (HPCA)*, Feb. 2019, pp. 69–80.
- [40] K. P. Yalamarthy, S. Dhall, M. T. Khan, and R. A. Shaik, "Lowcomplexity distributed-arithmetic-based pipelined architecture for an LSTM network," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 28, no. 2, pp. 329–338, Feb. 2020.
- [41] V. Rybalkin and N. Wehn, "When massive GPU parallelism ain't enough: A novel hardware architecture of 2D-LSTM neural network," in *Proc. ACM/SIGDA Int. Symp. Field-Program. Gate Arrays*, Feb. 2021, pp. 1–35.
- [42] A. Boutros et al., "Beyond peak performance: Comparing the real performance of AI-optimized FPGAs and GPUs," in Proc. Int. Conf. Field-Program. Technol. (ICFPT), Dec. 2020, pp. 10–19.
- [43] E. Nurvitadhi, A. Mishra, Y. Wang, G. Venkatesh, and D. Marr, "Hardware accelerator for analytics of sparse data," in *Proc. Design*, *Autom. Test Eur. Conf. Exhib. (DATE)*, 2016, pp. 1616–1621.
- [44] A. Ardakani, Z. Ji, S. C. Smithson, B. H. Meyer, and W. J. Gross, "Learning recurrent binary/ternary weights," 2018, arXiv:1809.11086.
- [45] W. Zheng *et al.*, "Binarized neural networks for language modeling," Stanford Univ., Stanford, CA, USA, Tech. Rep. cs224d, 2016.
- [46] Deep Learning With INT8 Optimization on Xilinx Devices. Accessed: 2017. [Online]. Available: https://www.xilinx.com/support/ documentation/white_papers/wp486-deep-learning-int8.pdf
- [47] M. Langhammer et al., "Extracting INT8 multipliers from INT18 multipliers," in Proc. Field Program. Logic Appl. (FPL), 2019, pp. 114–120.
- [48] A. Boutros, S. Yazdanshenas, and V. Betz, "Embracing diversity: Enhanced DSP blocks for low-precision deep learning on FPGAs," in *Proc. 28th Int. Conf. Field Program. Log. Appl. (FPL)*, Aug. 2018, pp. 35–357.
- [49] Intel Agilex Variable Precision DSP Blocks User Guide, Intel, Santa Clara, CA, USA, 2020.
- [50] A. Hannun et al., "Deep speech: Scaling up end-to-end speech recognition," 2014, arXiv:1412.5567.
- [51] F. Silfa, J. M. Arnau, and A. Gonzalez, "E-BATCH: Energy-efficient and high-throughput RNN batching," 2020, arXiv:2009.10656.