

A FRAMEWORK FOR MAPPING NEURAL NETWORKS ON PARALLEL ARCHITECTURES

A synopsis submitted in partial fulfillment of the requirements for the degree

of

Doctor of Philosophy

Submitted by:

SAURABH TEWARI
(2015CSZ8046)

under the guidance of

Prof. Anshul Kumar
Prof. Kolin Paul



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
INDIAN INSTITUTE OF TECHNOLOGY DELHI
NEW DELHI

List of Included Papers

This thesis is based on the following publications:

Published:

1. **S. Tewari**, A. Kumar and K. Paul, “*SACC: Split and Combine Approach to Reduce the Off-chip Memory Accesses of LSTM Accelerators*”, in DATE 2021.
2. **S. Tewari**, A. Kumar and K. Paul, “*Minimizing Off-Chip Memory Access for CNN Accelerators*”, in IEEE Consumer Electronics Magazine 2021.
3. **S. Tewari**, A. Kumar and K. Paul, “*Bus Width Aware Off-Chip Memory Access Minimization for CNN Accelerators*”, in ISVLSI 2020.
4. D. Stathis, Y. Yang, **S. Tewari**, A. Hemani, K. Paul, M. Grabherr and R. Ahmad, “*Approximate Computing Applied to Bacterial Genome Identification using Self-Organizing Maps*”, in ISVLSI 2019.

Contents

1	Introduction	1
1.1	Impact of off-chip memory access	1
2	About the Framework	2
2.1	Related Work	3
3	Off-Chip Memory Access Optimizations for CNN	4
3.1	Introduction	4
3.2	Related Work	5
3.3	Bus Width Aware Off-chip Memory Access Estimation	5
3.4	Optimal data partitioning for CNNs	5
3.5	Results	6
4	Optimizing data reuse of RNNs	7
4.1	Introduction	7
4.2	Previous Work	8
4.3	Proposed Approach	8
4.4	Result	9
5	Optimizing off-chip memory accesses using Quantization	10

1 Introduction

1.1 Impact of off-chip memory access

1.1.1 Performance

DNN accelerators have large number of processing elements (PEs) to exploit the parallelism of DNNs. While these PEs can perform several operations per cycle, their performance is limited by off-chip memory bandwidth. In the near foreseeable future, off-chip memory accesses limits the system performance [27]. Fig. 1 shows, two kernels with different compute intensity (FLOPS/byte), possibly due to different data reuse techniques. Kernel 2 can perform more computations per byte, compared to Kernel 1. Performance of Kernel 2 is limited by computational roof and it utilizes the compute resources fully (compute bound), while Kernel 1's performance is limited by memory bandwidth (memory bound). Memory bound kernels results

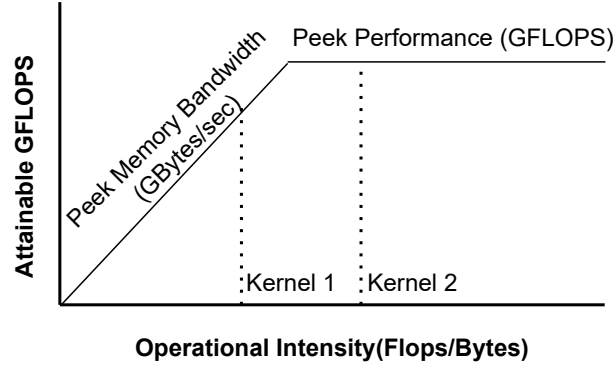


Figure 1: roofline model

in under-utilization of compute resources, which results in degraded performance and energy consumption. To improve the DNN accelerators' performance, maximizing the data reuse from on-chip memories is essential.

1.1.2 Energy efficiency

The energy of DNN accelerators is the sum of computations and data accesses energy.

$$E = E_{comp} + E_{data} \quad (1)$$

E_{comp} for a DNN layer depends on the number of operations in the layer, which is fixed as it depends on layer shape e.g. number of activations, filters, and filter sizes. However E_{data} can be optimized using data reuse techniques. Typically DNN accelerator system has multiple levels of memory hierarchy. Memory at each level has different size, access-energy and access-time. The off-chip memory (DDR), at the top of the hierarchy, has the largest size, and highest access-energy and access-time compared to memories at lower level of hierarchy. The off-chip memory access energy is upto two orders higher compared to on-chip memory access energy.

Reducing the off-chip memory accesses is the key to reduce the energy consumption of DNN accelerators. State of the art DNN accelerators aim to maximize the data reuse from on-chip memory to reduce the off-chip memory accesses. Figure 2 illustrates the impact of data reuse on the energy efficiency of a system with 2 levels of memory (DDR and L1). The energy efficiency can be expressed as below,

$$\begin{aligned} \text{energy-efficiency} &= \left(1 - \frac{E_{data}^2}{E_{data}^1}\right) \times 100 \\ &= \left(1 - \left(\frac{1}{n} + \frac{e_{L1}}{e_{ddr}}\right)\right) \times 100 \end{aligned}$$

As, $e_{L1} \ll e_{DDR}$, the ratio $\frac{e_{L1}}{e_{ddr}}$ is very small, the energy efficiency depends on number of reuses (n) from on-chip memory $L1$. It improves considerably with increasing the number of data reuse (n).

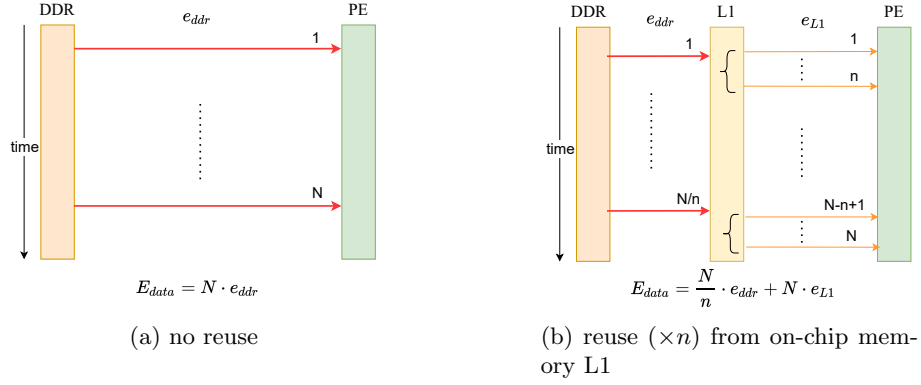


Figure 2: data access energy E_{data} with and without reuse.

In this work, we focus on the key issue to improve the performance and energy efficiency of DNN accelerators by reducing the off-chip memory accesses and maximizing the data reuse from on-chip memories.

2 About the Framework

It is arduous to create a generic framework for large variety of DNNs and accelerator architectures. However, most of the modern DNNs and accelerator’s architectures share similar architectural and algorithmic features. Therefore, we propose a mapping framework that addresses the most challenging problems faced by modern DNN accelerators and that are applicable to most of these DNNs accelerators.

1. Inferencing or usage of DNNs are frequently performed on edge devices. These edge devices are resource constraints and usually battery operated. They have tight energy budget which motivates to optimize the energy consumption of these DNN on edge devices. Secondly, the short response time is also expected to improve the user experience.
2. To optimize the performance of DNN based applications, DNN accelerators are commonly used. These accelerators have several compute units to perform large number of computations in parallel to improve the performance and user response time. However energy efficiency of these accelerators is far from desired, which restricts their usage in edge devices.
3. DNNs have large amount of parallelism, which can be exploited to speed up the processing, however limited memory bandwidth is the bottleneck.
4. It is observed, that more than 80% of overall energy consumption is due to off-chip memory accesses. Hence optimizing the off-chip memory accesses is the key to improve the energy efficiency of these DNN accelerators.

Our main contributions in this thesis dissertations are

1. We identify the key issue to be addressed for DNN accelerator’s wide usage in energy constraint devices is their off-chip memory accesses. We propose a mapping framework which takes input DNNs shape and sizes and analyzes the performance and energy consumption of DNN accelerators while considering accelerator’s architectural parameters.
2. We propose a model that precisely estimates the off-chip memory accesses and the energy consumption of the DNN’s layers, while taking into account the accelerator’s architectural parameters.
3. We propose an approach that determines the optimal partitoning of the DNN layers for reducing the off-chip memory accesses and energy consumptions. We analytically express the off-chip memory accesses of DNNs using the layer shapes and tiling parameters and express it as a constraint optimization problem.

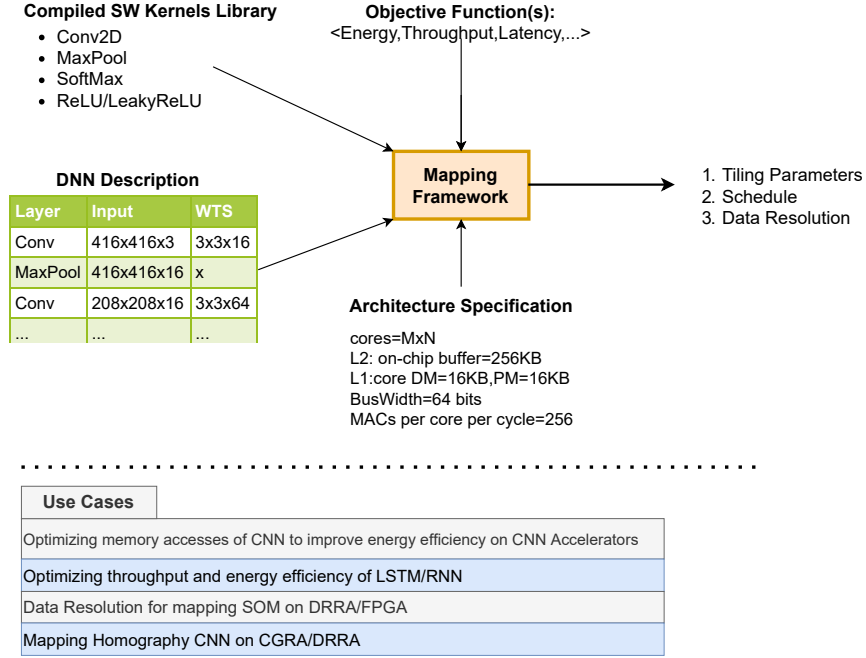


Figure 3: Mapping DNN on coarse grain parallel architectures

4. We propose a novel data reuse approach to improve the run time, energy efficiency of recurrent neural networks (RNN). The proposed approach schedules the computations in a way that reduces the memory accesses of large matrices by half and optimizes the performance and energy efficiency of RNN accelerators.
5. We analyzed the impact of data bit width on the performance, energy consumption and the accuracy of neural networks. We experimented with different data bit widths to analyze the trade-off between accuracy and energy efficiency of the NNs on reconfigurable parallel architectures.

2.1 Related Work

We compare our approaches and show that our work significantly improves the performance and energy efficiency of DNN accelerators compared to the state of the art. To address the computational and energy efficiency of DNNs, several ASIC [3, 7, 26] and FPGA based accelerators [4, 8, 11, 12, 16] are proposed. The energy efficiency of DNN accelerators is critical for their widespread usage, and off-chip memory access is the key to improving energy consumption. Most of these works focused on improving energy efficiency by reducing off-chip memory accesses. Figure 4 shows broad categories of state of the art approaches that aim to improve energy efficiency of DNN accelerators by reducing the off-chip memory accesses.

Some approaches [8, 16, 22] used on-chip memory to store all the weights. Sizes of weights in recent multi-layer LSTM models can be several MB's, and using large on-chip memory is expensive. These approaches are not scalable and effective only for small LSTM models. The proposed approach is independent of model size and effective for large LSTM models.

Several approaches used the fact that neural networks are error-tolerant and have lots of redundancy. They used the quantization and pruning techniques to compress the models' size. Approaches [8, 25] used 18-bit, Chang et al. [4] used 16-bit, Han et al. [12] used 12-bits precision for storing the inputs and weights, Lee et al. [16] used 8-bit inputs and 6-bits for weights to reduce the model size. The proposed approach is orthogonal to the quantization techniques and can be integrated with different quantization techniques to reduce the memory accesses further.

Han et al. [12] used pruning to compress the model. However, pruning results in irregular network structure, and the sparse matrix require additional computational and storage resources and causes unbalanced load distribution. To overcome this Wang et al. [25] used block-circulant matrices representations to com-

$$\text{Access Energy} = e_l \times (\text{Num. Elements} \times \text{Trips Count} \times \text{data width})$$

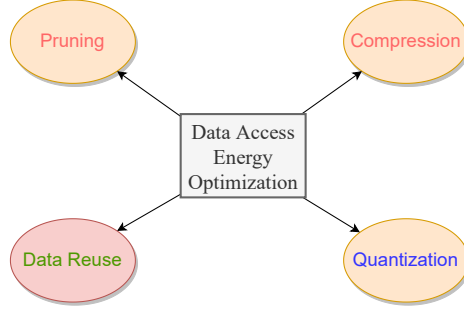


Figure 4: Energy efficiency optimization approaches

press the LSTM/RNN model and to eliminate irregularities resulted from compression. Some approaches ([12, 18, 19]) used load balance aware pruning techniques to overcome the unbalanced load distribution problem. Quantization and pruning approaches impacts the accuracy of the networks and may not be suitable where accuracy is at priority.

The other line of works reduced the memory accesses without compromising the accuracy of the network by applying the data-reuse techniques. Que et al. [21] proposed a blocking-batching scheme to reuse the LSTM weights fetched from external memory. Their approach reuses the weights of W matrix on a group of input vectors, by processing those input vectors as a batch. The input vectors in the same batch share the same weight matrices (W). However, it is difficult to collect required number of input vectors. As the LSTM cell states computations depend on cell states of the previous time-step, benefit of their batching schemes is limited to W matrix. Reusing weights of R across different time-steps has not been successful because of the state dependency.

Park et al. ([20]) proposed a time step interleaved weight reuse scheme (TSI-WR) which reuses the weights of R matrix between two adjacent time steps by performing computations in a time-interleaved manner. Their approach logically partitions the R matrix into blocks. A block is accessed from off-chip memory to compute the two consecutive hidden state vectors partially. However, their approach do not fully exploit the data reuse, and several weights are accessed repeatedly from the off-chip memory. In addition, the data reuse in TSI-WR approach depends on the on-chip storage size which limits the benefits of their approach.

3 Off-Chip Memory Access Optimizations for CNN

The mapping framework optimizes the performance of feed-forward and recurrent neural networks. Each of these networks has specific challenges that are addressed differently.

3.1 Introduction

CNNs have large number operations in each layer which can be speed up by exploiting the parallel compute capabilities of the CNN accelerators. However, CNN accelerators' performance are limited by memory bandwidth. In addition, the DNN accelerators have limited on-chip memory, which constrains the data size that can be stored in on-chip memory and limits the data reuse from on-chip memory. This results in increase in off-chip memory accesses.

The off-chip memory accesses of a layer depends significantly on the tile dimensions [17, 29]. Determining optimal tile dimensions requires analyzing off-chip memory accesses for different tile dimensions. Memory accesses also depend on the data bit width and architectural parameters of the accelerator like bus width. Towards this end, we have proposed an approach to compute the off-chip memory accesses accurately and used it to determine the optimal tile dimensions and data reuse scheme for CNN layers.

3.2 Related Work

To meet the computation and energy demands of CNNs, researchers have proposed efficient accelerators [2, 6, 10, 28]. CNN accelerators apply loop tiling to partition the layer data to fit into on-chip memory. Loop tiling is a compiler technique [1] that partitions the loop iteration space and large arrays into smaller tiles to increase the data locality and ensures that data fits into smaller memories.

Zhang et al. [29] and Li et al. [17] used loop tiling to optimize the off-chip memory accesses. These approaches expressed off-chip memory access as a function of tile dimensions and layer shape. Zhang et al. [29] determined optimal tile dimensions by enumerating all the legal tile dimensions. To reduce the hardware design complexity, they determined a global optimal tile dimension and used a common data reuse scheme for all the layers. Li et al. [17] proposed a layer-wise adaptive data partitioning and scheduling scheme. However, previous approaches ignored the architectural parameters and address alignment, and assumed all tiles of the same dimensions have the same off-chip memory accesses. With this assumption, the tile dimensions determined by previous approaches lead to a suboptimal solution.

3.3 Bus Width Aware Off-chip Memory Access Estimation

The CNN accelerators use a wide data bus to access off-chip memory to meet the high memory bandwidth requirement [5, 6]. If the number of bytes accessed from an off-chip memory address is not a multiple of bus width or the address is not aligned to the word boundary, it results in unused bytes lanes of the data bus. Figure 5 shows examples of memory accesses on a 64-bit data bus. Fig. 5a shows a read transaction of 8 bytes from an aligned address and uses the full bus width. However, if only 5 bytes are read from an aligned address, as shown in Figure 5b, 8 bytes are still accessed. If 5 bytes are read from an unaligned address, it results in 16 bytes of data access, as shown in Fig. 5d. The unused byte lanes do not carry any useful data, but they contribute to overall energy consumption. The length of the data read should be chosen such that bus utilization is high, and off-chip memory accesses and energy consumption are minimized. To estimate the performance and energy efficiency, we have developed an offline model that estimates the

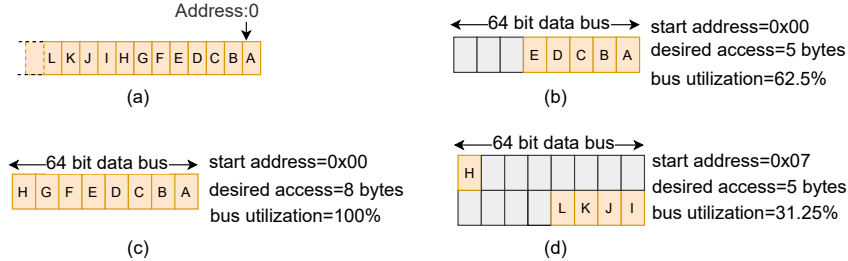


Figure 5: Off-chip memory accesses on 64-bit wide data bus

off-chip memory accesses while taking into account the data alignment and memory bus width.

3.4 Optimal data partitioning for CNNs

Fig. 6 shows a CL data stored in off-chip memory and its tiles in the accelerator’s on-chip buffer. The off-chip memory access of 3D data can be computed as following

$$\mathbb{B}_{3D} = r \times \sum_{t=1}^{N_{tiles}} \mathbb{B}_t \quad (2)$$

where N_{tiles} is the number of tiles. r and \mathbb{B}_t are the trips count and the number of bytes accessed from off-chip memory of the t^{th} tile, respectively. The trips count r depends on the data reuse scheme, and \mathbb{B}_t is computed using the model described in section 3.3

Tiles of ifm , ofm and wts reside in on-chip memory. If the on-chip memory buffer size is $buffSize$, and V_{ifm} , V_{ofm} , and V_{wts} are the size of ifm , ofm , and wts tiles, respectively, then constraints on tile dimensions are

$$(V_{ifm} + V_{wts} + V_{ofm}) \cdot DW \leq buffSize \quad (3)$$

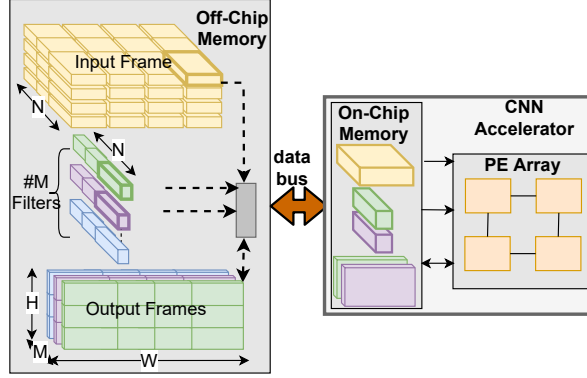


Figure 6: CNN layer tiles in off-chip and on-chip memory

Determining the tile dimensions which minimize the off-chip memory accesses, is a constraint optimization problem. The number of bytes accessed from off-chip memory and the constraints (3) are non-linear functions of tile dimensions and solving it is non-trivial.

We propose an approach to determine the optimal tile dimensions by computing the number of bytes accessed from off-chip memory (\mathbb{B}) at all the feasible points using section 3.3. Our approach determines the optimal tile dimensions for each layer for different data reuse schemes and it can be configured for different on-chip memory sizes, bus widths, and data bit width.

3.5 Results

We experimented with three popular CNN networks, AlexNet [15], VGG16 [23], and ResNet [13] having 8, 16, and 50 layers, respectively, with varying layer shapes and using filters of dimensions 1×1 , 3×3 , 5×5 , 7×7 , and 11×11 . To compare the results with other approaches, we have used the on-chip buffer size of 108 KB, batch size of 3 for VGG16, and 4 for ResNet and AlexNet. Figure 7 shows the comparison of the proposed approach with SmartShuttle (SS) [17] for VGG16. The SS approach determines the tile dimensions and data reuse scheme adaptively for each layer to minimize the data size (volume) accessed from off-chip memory.

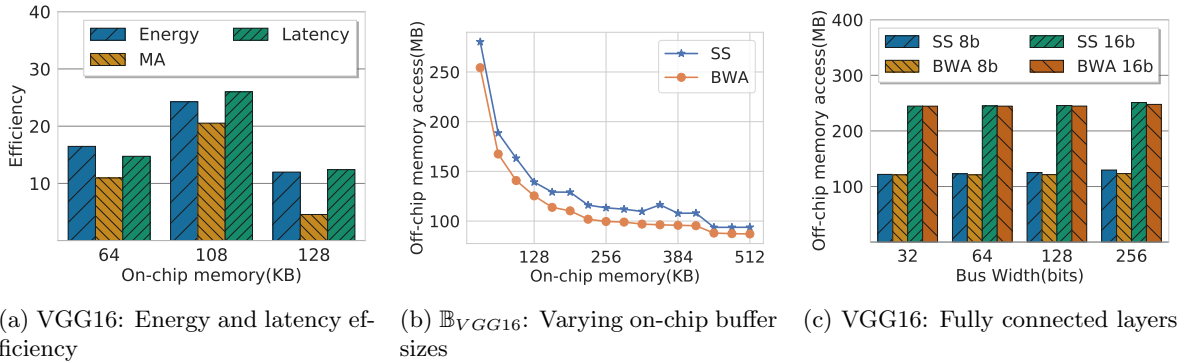


Figure 7: Latency and off-chip memory access. BWA: Bus Width Aware, SS: SmartShuttle

3.5.1 Latency And Energy Analysis

Figure 7a show the energy, off-chip memory accesses, and latency efficiency achieved using the BWA compared to the SS approach for VGG16 for 8 bits data width and 64 bits bus width. Using our FPGA implementation, we measured the memory access latencies and execution time for CLs to estimate the energy efficiency of the BWA approach compared to the SS approach. We observed that the changes in energy and latency are proportional to the changes in memory access. This observation confirms that off-chip memory access dominates the energy consumption of the CNN accelerators.

3.5.2 On-chip Memory Size

The size of the on-chip memory constrains the tile dimensions. Larger on-chip memory can accommodate bigger tiles which reduces the trips and number of off-chip memory access transactions. Figure 7b shows off-chip memory access for VGG16. Off-chip memory access reduces with increasing on-chip memory size for both approaches. However, the BWA approach performs better than SS for all the on-chip memory sizes.

3.5.3 Impact of Bus Width on \mathbb{B} of FCLs

In FCLs, the height and width of tiles and data are the same. Tiles in FCL can be accessed from off-chip memory using a single transaction which results in fewer transactions in FCLs compared to CLs. Whereas in CLs, multiple transactions are required to access different rows of the tiles. Thus the impact of bus width is less on \mathbb{B} of FCLs as compared to CLs. Figure 7c shows the off-chip memory accesses of FCLs of VGG16 for 8 bits data width. BWA reduces \mathbb{B}_{VGG16}^{FC} by 1%, 2%, 3%, and 4% on 32, 64, 128, and 256 bits wide buses, respectively.

4 Optimizing data reuse of RNNs

4.1 Introduction

Many applications involve sequential data processing and time-series predictions, e.g., natural language processing, speech recognition, video activity recognition, sentiment classification. Processing sequential data requires remembering the contextual information from previous data. Recurrent neural networks (RNNs) are deep learning algorithms specialized in handling such problems by maintaining an internal state based on previously seen data. LSTMs [14] are variants of RNNs designed to handle long-range dependencies by storing useful information about previous inputs for a long duration.

Customized accelerators are proposed to speed up the computations of LSTM. These accelerators have limited on-chip memory, specifically the accelerators targeted for embedded devices. LSTM computations involve multiple matrix-vector multiplications, and these matrix-vector multiplications are performed several times. The size of these matrices can be significant in several MB's and often exceed the size of the accelerator's on-chip memory. These matrices are partitioned into blocks and accessed from off-chip memory repeatedly by the accelerator, which results in a large volume of off-chip memory accesses and energy consumption. The high energy consumption limits the usage of these accelerators on embedded devices.

Typically the computations of LSTM cell is described by the following equations

$$\begin{aligned}
 i &= \sigma(W^i \cdot x_t + R^i \cdot h_{t-1} + b^i) \\
 f &= \sigma(W^f \cdot x_t + R^f \cdot h_{t-1} + b^f) \\
 g &= \tanh(W^g \cdot x_t + R^g \cdot h_{t-1} + b^g) \\
 o &= \sigma(W^o \cdot x_t + R^o \cdot h_{t-1} + b^o) \\
 c_t &= f \odot c_{t-1} + i \odot g \\
 h_t &= o \odot \tanh(c_t)
 \end{aligned} \tag{4}$$

where x_t is the input, h_t is the hidden state, and c_t is the cell state at time t . i, f, g, o are the computed gate values. \odot denotes the element-wise multiplications. W^j and R^j are the input and hidden state weight matrices, and b^j is the bias vector, learned during the training process, where $j \in \{i, f, g, o\}$. The dimension of h_t is referred to as the number of hidden states of the LSTM (N). At every time step, x_t is taken as input, and cell state (c_t) and hidden state (h_t) are computed using (4). The dependency of h_t on h_{t-1} and c_{t-1} prevents the parallel processing of multiple time steps and limits the data reuse.

Figure 8 shows the LSTM cell computations for two consecutive time steps using conventional and the proposed approach. To compute the hidden state vector h_t , conventional approaches access R matrix at each time step t , as shown in Fig. 8a. Accessing the weights at each time step results in a large volume of off-chip memory accesses.

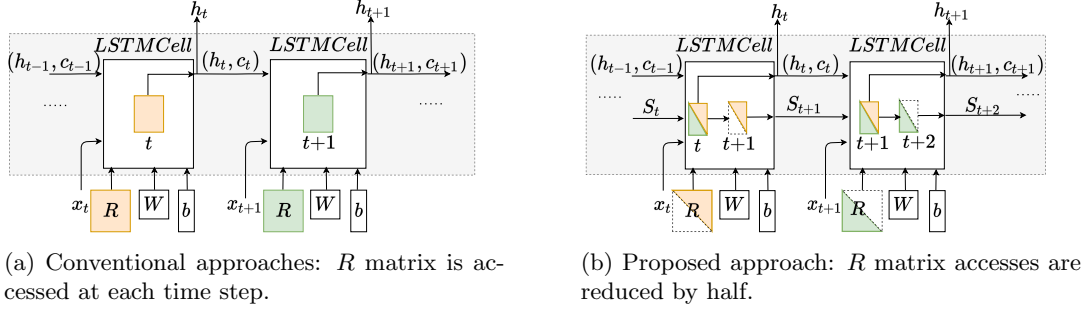


Figure 8: LSTM cell computations for consecutive time-steps showing the weight accesses.

4.2 Previous Work

Que et al. [21] proposed a blocking-batching scheme to reuse the LSTM weights fetched from external memory. Their approach reuses the weights of W matrix on a group of input vectors, by processing those input vectors as a batch. The input vectors in the same batch share the same weight matrices (W). However, it is difficult to collect required number of input vectors. As the LSTM cell states computations depend on cell states of the previous time-step, benefit of their batching schemes is limited to W matrix. Reusing weights of R across different time-steps has not been successful because of the state dependency.

Park et al. ([20]) proposed a time step interleaved weight reuse scheme (TSI-WR) which reuses the weights of R matrix between two adjacent time steps by performing computations in a time-interleaved manner. Their approach logically partitions the R matrix into blocks. A block is accessed from off-chip memory to compute the two consecutive hidden state vectors partially. However, their approach do not fully exploit the data reuse, and several weights are accessed repeatedly from the off-chip memory. In addition, the data reuse in TSI-WR approach depends on the on-chip storage size which limits the benefits of their approach.

4.3 Proposed Approach

We propose an approach that reduces off-chip memory accesses by splitting the computation in two. At each time step, while the computations of a hidden state vector of one time step h_t completes, partial computation of next time step (S_{t+1}) is also performed by reusing the weights. Our approach schedules the computations in a way that reuses all the weights of R between two adjacent time steps. The data reuse in our approach is independent of on-chip buffer sizes which makes it suitable for accelerators with very small on-chip memory.

The proposed data reuse approach splits and combines the LSTM cell computations in a way that reduces the off-chip memory accesses of hidden state matrices by 50%. The computation of the h_t can be expressed as shown below

$$h_t[k] = F(S_t[k] + q_t[k]) \quad (5)$$

where F is a non-linear function. q_t is computed as $W \cdot x_t + b$ and its computations are independent of previous step cell states. $S_t[k]$ is the sum of N product terms as shown below,

$$S_t[k] = \sum_{n=0}^{N-1} R[k][n] \cdot h_{t-1}[n] \quad (6)$$

$S_t[k]$ can be computed as a sum of the following two partial sums $S_t^L[k]$ and $S_t^U[k]$

$$S_t^L[k] = \sum_{n=0}^k R[k][n] \cdot h_{t-1}[n] \quad (7)$$

$$S_t^U[k] = \sum_{n=k+1}^{N-1} R[k][n] \cdot h_{t-1}[n] \quad (8)$$

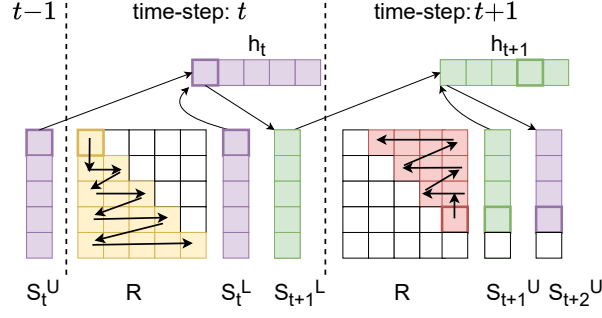


Figure 9: Splitting the hidden state vector computations into partial sums

Equation (7) uses the lower-diagonal and diagonal elements of R (R^L), and (8) uses the upper diagonal elements of R (R^U). As shown in Figure 9, R^L and R^U are accessed in consecutive time steps and reused in the partial sum computations of two steps. At time step t , S_t^U and h_{t-1} are the inputs from the previous time step, and R^L is reused to compute the partial sums S_t^L and S_{t+1}^L . Input S_t^U is added to S_t^L to compute h_t , and S_{t+1}^L is passed to $(t+1)^{th}$ step computations. In the same way, at time step $t+1$, R^U is reused to compute S_{t+1}^U and S_{t+2}^U . Elements of R^L are accessed from top to bottom, left to right, while elements of R^U are accessed in the reverse order to satisfy the dependencies. As shown in Figure 9, the proposed approach accesses the weight matrix R once, to compute h_t and h_{t+1} .

4.4 Result

We have compared our approach with conventional approaches and TSI-WR approach [20]. Conventional approaches access the hidden state weight matrices R at each step from the off-chip memory. We have used the same on-chip buffer size to store the weight matrices ($4 \times B^2$) to perform a fair comparison. The proposed approach requires additional on-chip memory ($4N + 4B$) to store the four partial sum and temporary vectors. We have experimented with LSTM models used in speech recognition (for TIMIT [9]) and character level Language Modelling (LM) [24].

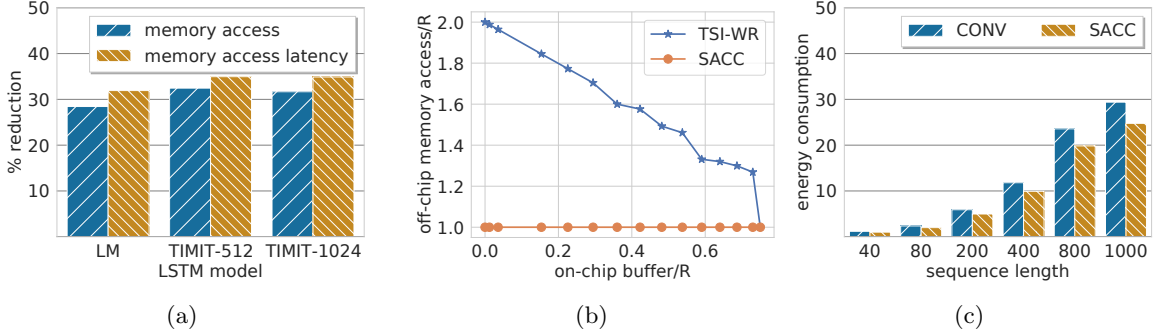


Figure 10: (a) off-chip memory access and latency (b) Off-chip memory access comparison between TSI-WR and SACC approach for two consecutive time steps. (c) normalized energy. CONV: conventional

4.4.1 Memory Accesses

Figure 10a shows the proposed approach's reduction in off-chip memory accesses compared to conventional approaches. The proposed approach reduces the memory accesses of the R matrices by half. Total memory accesses, including R , W , and b , for LM and TIMIT models reduce by 28% and 32%, and memory access latencies reduce by 29% and 31%, respectively. The proposed approach reduces the run-time by 12.8% and 16% for LM and TIMIT models.

4.4.2 Comparison with TSI-WR

Figure 10b compares the off-chip memory accesses of the SACC and the TSI-WR approaches using the simulation results. The performance of the TSI-WR approach depends on on-chip buffer sizes. TSI-WR reduces 50% off-chip memory accesses when the on-chip buffer size is 70% of the R matrix. The proposed approach reduces R 's memory access by 50%, irrespective of the on-chip buffer size.

4.4.3 Energy Efficiency

We computed the energy consumption using the design power reported by the Vivado synthesis tool, execution time, and off-chip memory accesses. Figure 10c shows the normalized energy consumption of the conventional and the proposed approach for TIMIT-1024. Due to the reduction in the run-time and off-chip memory accesses, the SACC approach reduces the energy consumption on average by 13% and 16% for LM and TIMIT models.

5 Optimizing off-chip memory accesses using Quantization

References

- [1] A. V. Aho, M. S. Lam, R. Sethi, and J. D. Ullman. Compilers: Principles, techniques, and tools, 2006.
- [2] M. Alwani, H. Chen, M. Ferdman, and P. Milder. Fused-layer CNN accelerators. In *MICRO*, 2016.
- [3] E. Azari and S. Vrudhula. Elsa: A throughput-optimized design of an lstm accelerator for energy-constrained devices. *ACM Transactions on Embedded Computing Systems (TECS)*, 19(1):1–21, 2020.
- [4] A. X. M. Chang, B. Martini, and E. Culurciello. Recurrent neural networks hardware implementation on fpga. *arXiv preprint arXiv:1511.05552*, 2015.
- [5] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam. DianNao: A small-footprint high-throughput accelerator for ubiquitous machine-learning. *ASPLOS*, 2014.
- [6] Y.-H. Chen, J. S. Emer, and V. Sze. Eyeriss: A spatial architecture for energy-efficient dataflow for convolutional neural networks. *ACM/IEEE ISCA*, 2016.
- [7] F. Conti, L. Cavigelli, G. Paulin, I. Susmelj, and L. Benini. Chipmunk: A systolically scalable 0.9 mm², 3.08 gop/s/mw@ 1.2 mw accelerator for near-sensor recurrent neural network inference. In *2018 IEEE Custom Integrated Circuits Conference (CICC)*, pages 1–4. IEEE, 2018.
- [8] J. C. Ferreira and J. Fonseca. An fpga implementation of a long short-term memory neural network. In *2016 International Conference on ReConFigurable Computing and FPGAs (ReConFig)*, pages 1–8. IEEE, 2016.
- [9] J. S. Garofolo. Timit acoustic phonetic continuous speech corpus. *Linguistic Data Consortium*, 1993, 1993.
- [10] V. Gokhale, J. Jin, A. Dundar, B. Martini, and E. Culurciello. A 240 G-ops/s mobile coprocessor for deep neural networks. In *CVPR Workshops*, 2014.
- [11] Y. Guan, Z. Yuan, G. Sun, and J. Cong. Fpga-based accelerator for long short-term memory recurrent neural networks. In *2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC)*, pages 629–634. IEEE, 2017.
- [12] S. Han, J. Kang, H. Mao, Y. Hu, X. Li, Y. Li, D. Xie, H. Luo, S. Yao, Y. Wang, et al. Ese: Efficient speech recognition engine with sparse lstm on fpga. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 75–84, 2017.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

- [14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [16] M. Lee, K. Hwang, J. Park, S. Choi, S. Shin, and W. Sung. Fpga-based low-power speech recognition with recurrent neural networks. In *2016 IEEE International Workshop on Signal Processing Systems (SiPS)*, pages 230–235. IEEE, 2016.
- [17] J. Li, G. Yan, W. Lu, S. Jiang, S. Gong, J. Wu, and X. Li. SmartShuttle: Optimizing off-chip memory accesses for deep learning accelerators. *DATE*, 2018.
- [18] J. Park, J. Kung, W. Yi, and J.-J. Kim. Maximizing system performance by balancing computation loads in lstm accelerators. In *2018 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 7–12. IEEE, 2018.
- [19] J. Park, W. Yi, D. Ahn, J. Kung, and J.-J. Kim. Balancing computation loads and optimizing input vector loading in lstm accelerators. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39(9):1889–1901, 2019.
- [20] N. Park, Y. Kim, D. Ahn, T. Kim, and J.-J. Kim. Time-step interleaved weight reuse for lstm neural network computing. In *Proceedings of the ACM/IEEE International Symposium on Low Power Electronics and Design*, pages 13–18, 2020.
- [21] Z. Que, T. Nugent, S. Liu, L. Tian, X. Niu, Y. Zhu, and W. Luk. Efficient weight reuse for large lstms. In *2019 IEEE 30th International Conference on Application-specific Systems, Architectures and Processors (ASAP)*, volume 2160, pages 17–24. IEEE, 2019.
- [22] 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 *2018 28th international conference on field programmable logic and applications (FPL)*, pages 89–897. IEEE, 2018.
- [23] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [24] M. Sundermeyer, H. Ney, and R. Schlüter. From feedforward to recurrent lstm neural networks for language modeling. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(3):517–529, 2015.
- [25] S. Wang, Z. Li, C. Ding, B. Yuan, Q. Qiu, Y. Wang, and Y. Liang. C-lstm: Enabling efficient lstm using structured compression techniques on fpgas. In *Proceedings of the 2018 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 11–20, 2018.
- [26] Z. Wang, J. Lin, and Z. Wang. Accelerating recurrent neural networks: A memory-efficient approach. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 25(10):2763–2775, 2017.
- [27] S. Williams, A. Waterman, and D. Patterson. Roofline: an insightful visual performance model for multicore architectures. *Communications of the ACM*, 52(4):65–76, 2009.
- [28] L. Xie, X. Fan, W. Cao, and L. Wang. High throughput CNN accelerator design based on FPGA. In *FPT*, 2018.
- [29] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong. Optimizing FPGA-based accelerator design for deep convolutional neural networks. In *FPGA*, 2015.