

# SIAM-IMA Etymo workshop - text extraction

Steven Elsworth

June 13, 2018

# Text extraction

Recently, significant accuracy improvement has been achieved for acoustic recognition systems by increasing the model size of Long Short-Term Memory (LSTM) networks. Unfortunately, the ever-increasing size of LSTM model leads to inefficient designs on FPGAs due to the limited on-chip resources. The previous work proposes to use a pruning based compression technique to reduce the model size and thus speeds up the inference on FPGAs. However, the random nature of the pruning technique transforms the dense matrices of the model to highly unstructured sparse ones, which leads to unbalanced computation and irregular memory accesses and thus hurts the overall performance and energy efficiency.

Recently, significant accuracy improvement has been achieved for acoustic recognition systems by increasing the model size of Long Short-Term Memory (LSTM) networks. Unfortunately, the ever-increasing size of LSTM model leads to inefficient designs on FPGAs due to the limited on-chip resources. The previous work proposes to use a pruning based compression technique to reduce the model size and thus speeds up the inference on FPGAs. However, the random nature of the pruning technique transforms the dense matrices of the model to highly unstructured sparse ones, which leads to unbalanced computation and irregular memory accesses and thus hurts the overall performance and energy efficiency.

# Methods

- ▶ pdftotxt
- ▶ Textract (Backend pdftotxt)
- ▶ PyPDF2
- ▶ pdfminer
- ▶ pyocr (backend tesseract)

See Jupyter Notebook.

# Problems faced

Give examples of :

- ▶ Equation conversion
- ▶ Figure conversion
- ▶ White space
- ▶ Merged words
- ▶ Page columns
- ▶ Time

dependencies between operators are complicated. So, it is difficult to evenly allocate computing resources under the FPGA resource constraints while guaranteeing the complex data dependencies.

In this work, we propose to compress the weight matrices in the LSTM inference model in a structured manner by using block-circulant matrices [22]. The circulant matrix is a square matrix, of which each row (column) vector is the circular reform of the row (column) vector. Any matrix could be transformed into a set of circulant submatrices via a block-circulant matrices. Therefore, by representing each block-circulant matrix with a vector, the storage requirement could be reduced from  $O(N^2)$  to  $O(N)$  if the block (vector) size is  $k$ . Since the compressed weight matrices are still dense, the block-circulant matrix based compression is infeasible to hardware acceleration on FPGAs. In order to further speed up the computation of LSTMs, we propose to accelerate the most computation-intensive circular correlation operator by applying Fast Fourier Transform (FFT) algorithms to reduce the computational complexity from  $O(N^2)$  to  $O(N \log N)$ .

After the model is compressed, we propose an automatic optimization and synthesis framework called C-LSTM to port efficient LSTM designs onto FPGAs. The framework is composed of model training and implementation flows. The former one is in charge of iteratively training the compressed LSTM model and exploring the trade-offs between compression ratio and prediction accuracy. As for the model implementation, it mainly consists of two parts which are (1) template generation and (2) automatic LSTM synthesis framework. For the former part, after analyzing a wide range of LSTM algorithms, we generate a suite of LSTM primitive operators which is general enough to accommodate even the most complicated LSTM variant [23]. Then, a suite of highly optimized C/C++ templates of the primitive operators are manually generated and by walking through a series of optimizations such as datapaths and activation quantization, FFT-FFT decoupling and etc. As for the latter part, the well-trained LSTM inference model is first analyzed and transformed into a directed acyclic dependency graph, where each node represents an operator and each edge indicates the associated data dependency between two operators. Secondly, we propose a specialized pipeline optimization algorithm considering both coarse-grained and fine-grained pipelining schemes to schedule the operators into appropriate stages. In the first step, we use an accurate performance and resource model to enable a fast design space exploration for optimal design parameters. Lastly, the scheduling results and optimization parameters are fed to code generator and hardware toolchains to implement the optimized LSTM acceleration design on FPGAs.

Overall, the contributions of this paper are listed as:

- We employ the block-circulant matrices based structured compression technique for LSTMs which largely reduces the computation complexity and memory footprint without incurring any computation and memory access irregularities. This method enables a high compression and acceleration of the LSTM models.
- We develop a general LSTM optimization and synthesis framework C-LSTM to enable automatic and efficient implementations of a wide range of LSTM variants on FPGAs.

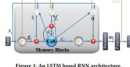


Figure 1: An LSTM based RNN architecture.

The framework mainly consists of a suite of highly optimized C/C++ based templates of primitive operators and an automatic LSTM synthesis flow.

- We present efficient implementations of LSTMs which achieve up to 35.8X and 13.5X gains in performance and energy efficiency, respectively, compared with the state-of-the-art. The proposed implementations incur very small accuracy degradation.

## 2 LSTM BACKGROUND

LSTM is a key component of the accurate model in modern large-scale automatic speech recognition (ASR) systems [9, 21], and also the most computation and memory-intensive part. Due to the complicated and flexible data dependencies among gates, cells, and outputs, a lot of LSTM variants have been proposed. In this paper, we use a widely deployed variant called Google LSTM [23] as an example throughout this paper without loss of generality. The architecture details of the Google LSTM is shown in Figure 1. The LSTM accepts an input sequence  $x_1, x_2, x_3, \dots, x_T$  (each of  $x_i$  is a vector corresponding to time  $i$ ) with the output sequence from last step  $y_1^{(t-1)}, y_2^{(t-1)}, \dots, y_{T-1}^{(t-1)}$  (each of  $y_i$  is a vector). The input of Google LSTM at time  $t$  depends on the output at  $t-1$ . The LSTM contains a special memory cell storing the temporal state of the network. It also contains three special multiplicative units which are input, output and forget gates. The output sequence  $Y = (y_1, y_2, y_3, \dots, y_T)$  is computed by using the following equations iteratively from  $t = 1$  to  $T$ :

$$i_t = \sigma(W_{ix}x_t + W_{iy}y_{t-1} + W_{ic}c_{t-1} + b_{i_t}) \quad (1)$$

$$f_t = \sigma(W_{fx}x_t + W_{fy}y_{t-1} + W_{fc}c_{t-1} + b_{f_t}) \quad (2)$$

$$g_t = \sigma(W_{gx}x_t + W_{gy}y_{t-1} + b_{g_t}) \quad (3)$$

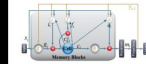
$$c_t = f_t \odot c_{t-1} + g_t \odot i_t \quad (4)$$

$$o_t = \sigma(W_{ox}x_t + W_{oy}y_{t-1} + W_{oc}c_{t-1} + b_{o_t}) \quad (5)$$

$$y_t = o_t \odot h_t \odot h_{t-1} \quad (6)$$

$$y_t = W_{y0}h_{t-1} \quad (7)$$

where symbols  $f, g, c, o, m$ , and  $y$  are respectively the input gate, forget gate, output gate, cell state, cell output, and a projected output; the  $\odot$  operator denotes the element-wise multiplication, and the  $+$  operator denotes the element-wise addition. The  $W$  terms denote weight matrices (e.g.  $W_{ix}$  is the matrix of weights from the input vector  $x_t$  to the input gate), and the  $b$  terms denote bias vectors. Please note  $W_{ix}, W_{iy}, W_{ic}$  and  $W_{ox}$  are diagonal matrices for preclude connections, thus they are essentially a vector, and the



---

<sup>1</sup> This estimation considers both weights and indices (there is at least one index per weight after compression in ESE). However, this is a pessimistic estimation for ESE because indices can use fewer bits for representation than weights;

the ADM-7v3 platform. Compared with ESE, we achieve 10.2X and 18.8X performance speedups and 19.1X and 33.5X energy efficiency gains using FFT8 and FFT16, respectively. Since the power consumption of C-LSTM is only half of the ESE, the energy efficiency gain is higher than performance. It is necessary to note that as shown in Table 2, the manufacturing process of XCKU060 FPGA is 28nm while the process of Virtex-7 is 28nm, which means the energy efficiency gain reported here is pessimistic.

Although the promising performance and energy gains are achieved by C-LSTM, the resource utilization for LUT, FF, and BRAM are less than ESE, and more important, the relative PER degradation is very small, which are 0.32% and 1.23% using FFT8 and FFT16, respectively. After detailed analysis, we summarize the fundamental reasons for the high performance and power gains in three aspects. First, the structured compression used in this work eliminates the irregular computation and memory accesses which not only makes the design more regular but also exposes more parallelism. This could be verified in that the DSP resource consumption of the proposed method is much more than ESE. Secondly, the whole model (weights matrices and the projection matrix) could be stored on-chip without fetching data from off-chip DRAM, making the LSTM not bounded by memory. Lastly, the more efficient implementation of LSTM on FPGAs contributes to the high efficiency. For example, we use the 22-segment piece-wise linear function to approximate the activation functions while ESE employs look-up tables which break the activation down into 2048 segments and consume more resources. Moreover, we propose to employ FFT based block-circulant

matrix multiplication while ESE uses sparse matrix multiplication which needs to store extra indices for sparse matrices and thus prevents from storing the whole model on-chip.

### 6.3 Experimental Results of Small LSTM

In order to validate that proposed C-LSTM is not only appropriate for Google LSTM model, we also implement a Small LSTM [20] model on both FPGA platforms.

In KU060 platform, the FFT8 and FFT16 designs could achieve 19.3X and 35.9X performance speedup compared with ESE, respectively. In the ADM-7v3 platform, the performance speedups are 17.5X and 31.9X and the energy efficiency gains are 34.2X and 59.4X compared with ESE, respectively. For both platforms, the PER degradation is 0.29% and 1.16% for FFT8 and FFT16, respectively.

## 7 RELATED WORK

Recently, FPGA has emerged as a promising hardware acceleration platform for DNNs as it provides high performance, low power and reconfigurability. A lot of FPGA based accelerators have been proposed for convolutional neural networks (CNNs) to overcome the computing and energy efficiency challenges. [28] proposes to utilize systolic array based convolution architecture to achieve better frequency and thus performance for CNNs on FPGAs. [18] employs the Winograd algorithm to reduce the multiplication operations as to save DSP resources and accelerate matrix multiplication in CNNs. [30] proposes to take advantage of the heterogeneous algorithms to maximize the resource utilization for convolutional layers on

This estimation considers both weights and indices (there is at least one index per weight after compression in ESE). However, this is a pessimistic estimation for ESE because indices can use fewer bits for representation than weights; matrix multiplication while ESE uses sparse matrix multiplication which needs to store extra indices for sparse matrices and thus prevents from storing the whole model on-chip. the ADM-7v3 platform. Compared with ESE, we achieve 10.2X and 18.8X performance speedups and 19.1X and 33.5X energy efficiency gains using FFT8 and FFT16, respectively. Since the power