



# High-Performance Deep Learning via a Single Building Block

## Introduction

Deep learning (DL) is one of the most prominent branches of machine learning. Due to the immense computational cost of DL workloads, industry and academia have developed DL libraries with highly-specialized kernels for each workload/architecture, leading to numerous, complex code-bases that strive for performance, yet they are hard to maintain and do not generalize.

### Punchline of this work is the following:

• The introduction of the JIT-compiled batch-reduce GEMM kernel (this work leverages LIBXSMM's (<https://github.com/hfp/libxsmm>) infrastructure) and show how the most popular DL algorithms can be formulated with this kernel as the basic building-block

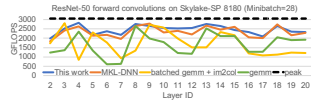
• The design and implementation of multi-threaded, high performance DL primitives covering Recurrent Neural Networks (RNN/LSTM), Convolution Neural Networks (CNN) and Multilayer Perceptron (MLP) inference and training algorithms with batch-reduce GEMM kernel being the basic building block. We need to optimize at low-level only this kernel for all DL primitives.

• significantly reducing the code size and code specialization and hence enhance code maintainability

• A detailed performance comparison of our DL primitives with state-of-the-art vendor-optimized libraries.

• CNN proof-of-concept results on integrated GPUs and CNN kernels within TVM that leverage the batch-reduce GEMM kernel demonstrating performance portability opportunities of our approach

A motivation example is given by this figure of our work:



The blue line shows the performance of the convolution primitive that leverages our new batch-reduce GEMM kernel achieving average efficiency of 83%, and outperforms even the ad hoc, vendor-optimized kernel. It is clearly superior over batched-gemm and classic im2col algorithms.

### Properties of our proposal:

• The batch-reduce GEMM kernel essentially multiplies a sequence of input sub-tensor blocks (which form a batch) and the partial multiplication results are reduced into a single accumulator/output sub-tensor block

• The new kernel is flexible enough to accommodate coarse-grained and fine-grained operations that arise in DL workloads

• The new kernel supports fusion of subsequent operators on the output sub-blocks

• The new kernel is optimized for the nuances of the architecture at hand, and serves as a robust building block

### These properties allow for transformative changes:

• by implementing and optimizing this single kernel for a given architecture, the development of DL primitives degenerates to mere loop tuning around this kernel

• addressing the issue of combinatorial explosion of low-level optimization work that is required for each pair <architecture, DL primitive>

• enabling efficient usage of tensor compiler frameworks such as TVM, GLOW, PlaidML, and/or MLIR (DL primitives are inherently tensor algebra programs)

## The Batch-Reduce GEMM Kernel

The batch-reduce GEMM kernel materializes the operation:

**Algorithm 1** The batch-reduce GEMM kernel

```
Inputs:  $A_j \in \mathbb{R}^{m \times k}$ ,  $B_i \in \mathbb{R}^{k \times n}$ ,  $i = 0, \dots, N-1$ ,  $C_j \in \mathbb{R}^{m \times n}$ ,  $\beta \in \mathbb{R}$ 
Output:  $C_j = \beta \cdot C_j + \alpha \sum_{i=0}^{N-1} A_i \cdot B_i$ 
1 for  $i_n = 0, \dots, n-1$  with step  $n_b$  do
2   for  $i_m = 0, \dots, m-1$  with step  $m_b$  do
3     acc_regs  $\leftarrow$  load  $m_b \times n_b$   $C_j$  subblock  $c_{i_m, i_n}$ 
4     for  $i = 0, \dots, N-1$  with step 1 do
5       for  $i_k = 0, \dots, k-1$  with step 1 do
6         Outer product GEMM microkernel
7         acc_regs  $\oplus= A_i$  subcolumn  $a_{i_m, i_k} \times B_i$  subrow  $b_{i_k, i_n}$ 
8      $C_j$  subblock  $c_{i_m, i_n} \leftarrow$  acc_regs
```

$$C_j = \beta \cdot C_j + \alpha \sum_{i=0}^{N-1} A_i \cdot B_i$$

- this kernel multiplies the specified blocks  $A_i \in \mathbb{R}^{m \times k}$  and  $B_i \in \mathbb{R}^{k \times n}$  and reduces the partial results to a block  $C_j \in \mathbb{R}^{m \times n}$
- $A$  and  $B$  can alias and also the blocks  $A_i$  and  $B_i$  can reside at any arbitrary position

The kernel takes the following arguments:

• two arrays of pointers to the corresponding blocks  $A_i$  and  $B_i$  to be multiplied

• a pointer to the output block  $C_j$

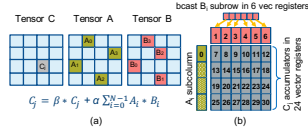
• the number  $N$  dictating the number of blocks

• the scaling parameters  $\alpha$  and  $\beta$

Algorithm 1 depicts the batch-reduce GEMM kernel:

- lines 1-2 block the computation of the result  $C_j$  in  $m_b \times n_b$  subblocks
- load accumulator (line 3) and loop over all pairs  $A_i, B_i$  and subblocks (line 4-7)
- Update  $C_j$  (line 8)

The following figures highlight the micro-kernel targeting an Intel AVX-512 enabled platform (32 SIMD registers, 16 elements each):



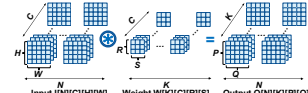
Most important properties of this kernel are:

- builds upon and extends the open source LIBXSMM library which leverages JIT techniques and generates small GEMMs achieving close to peak performance
- multiplies an  $m_b \times 1$  column of  $A_i$  with a  $1 \times n_b$  row of  $B_i$  (in this example  $m_b = 64, n_b = 6$ )
- accumulates into registers 7-30 that hold the partial  $C_j$  result
- broadcasts the row elements of  $B_i$  into registers 1-6
- load 16 elements of the  $A_i$  column in register 0

- performs 6 fused-multiply-add instructions (FMAs) with registers 0 and 1-6 we update the accumulators 7-12
- repeats the former for remaining 48 elements of the  $A_i$  column and we update all the accumulation registers 13-30

Note: our enhanced LIBXSMM leverages various strategies depending on the architecture at hand (i.e. vector length) and the  $m_b, n_b$  values and makes use of software prefetches for  $A_i$  and  $B_i$  elements.

## Convolution Neural Networks (CNN)



### Algorithm 2 CNN forward pass via batch-reduce GEMM

```
1  $C_b = C/b_h, K_b = K/b_h, Q_b = Q/b_h$ 
2 for  $n = 0, \dots, P-1$  do
3   for  $k_b = 0, \dots, K_b-1$  do
4     for  $o_b = 0, \dots, C_b-1$  do
5       for  $o_b = 0, \dots, Q_b-1$  do
6         for  $o_b = 0, \dots, S_b-1$  do
7            $o_i = o_b \cdot b_h, s_i = str \cdot o_i, i_j = str \cdot o_j, i = 0$ 
8           Prepare batch-reduce GEMM arguments
9           for  $r = 0, \dots, R-1$  do
10            for  $s = 0, \dots, S-1$  do
11              for  $c = 0, \dots, B_c-1$  do
12                 $A_{parr}[i] = \alpha W[b_h][c_b + c][r][s][o][0]$ 
13                 $B_{parr}[i+j] = \beta I[p][c_b + c][s][i+j] + [i] + [s][0]$ 
14                 $Out = \alpha O[p][k_b][o_b][0]$ 
15                batchreduce_gemm( $A_{parr}, B_{parr}, Out, R, S, B_c$ )
```

• The input activation tensors are convoluted with the weight tensors to yield the output activation tensors.

• Batch-reduce GEMM kernel obviates the need for a specialized convolution kernel.

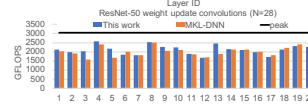
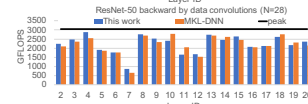
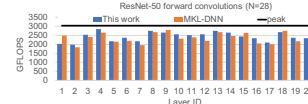
• The batch-reduce GEMM kernel inherently optimizes load/store of the output  $O$  in case of  $R, S > 1$  and in case the input feature map loop is reordered as the innermost loop.

– By properly selecting the sub-tensors of weights/inputs to be multiplied and reduced onto an  $O$  sub-tensor, the accumulation takes place entirely in registers

• The data reuse optimizations/parallelization tasks then simply translates to tuning the surrounding loops

• Only slight modifications to the forward kernel are required in order to implement the back-propagation kernels, as they can be mapped through linear index transformations into the forward convolution loop nest (“duality”).

Performance of ResNet-50 convolutions on Xeon 8180:



• The weighted efficiency of the FWD convolutions is 83% of peak, while MKL-DNN exhibits weighted efficiency of 81% of peak (2.5% slower than our work).

• Our kernels have weighted efficiency 80% and 73.6% for the BWD pass and the UPD pass respectively, whereas MKL-DNN BWD and UPD kernels illustrate weighted efficiencies of 78.9% and 68.9% respectively, and are 1% and 7% slower than our kernels.

## Long-Short Term Memory (LSTM)

Parameters:  $N$  batch-size,  $T$  sequence length,  $C$  state size and  $K$  hidden state size. The input/weight tensors are combined based on the Equations 1-6 and yield the output sequence  $h \in \mathbb{R}^{T \times K \times N}$  and tensor  $s \in \mathbb{R}^{T \times K \times N}$ :

$$\begin{aligned} i_t &= \sigma(W_i \cdot x_t + R_i \cdot h_{t-1} + b_i) & (1) \\ c_t &= \tanh(W_c \cdot x_t + R_c \cdot h_{t-1} + b_c) & (2) \\ f_t &= \sigma(W_f \cdot x_t + R_f \cdot h_{t-1} + b_f) & (3) \\ o_t &= \sigma(W_o \cdot x_t + R_o \cdot h_{t-1} + b_o) & (4) \\ s_t &= i_t \circ s_{t-1} + i_t \circ c_t & (5) \\ h_t &= o_t \circ \tanh(s_t) & (6) \end{aligned}$$

### Algorithm 3 Forward propagation pass of LSTM cell

```
Inputs: Weight tensors  $W_i[K_h][C_h][b_i][b_c]$ ,  $R_i[K_h][K_h][b_i][b_c]$ 
Input sequence  $x[T][N][C]$ , Bias  $b_i[K]$ , blocking factors  $b_h, b_c, b_i$ 
Outputs: Output sequence  $h[T][N][K]$  and  $s[T][N][K]$ 
1  $N_b \leftarrow N/b_h$ 
2 Based on  $thread\_id$  calculate  $K_{h\_start}, K_{h\_end}, N_{h\_start}$  and  $N_{h\_end}$  to assign output work items
3 for  $t = 0, \dots, T-1$  do
4   for  $o_b = N_{h\_start}, \dots, N_{h\_end}$  do
5     for  $i_b = N_{h\_start}, \dots, N_{h\_end}$  do
6       Compute a block of  $x = \sigma(W_i \cdot x_t + R_i \cdot h_{t-1} + b_i)$ 
7        $i_k \leftarrow b_h, i_n \leftarrow i_b \cdot b_h$ 
8        $i[f][i_k] \leftarrow b_i[b_h]$ 
9       for  $i_b = 0, \dots, C_h-1$  do
10         $A_{parr}[i_k] = \alpha W_i[b_h][i_b][i_k][0]$ 
11         $B_{parr}[i_k] = \beta R_i[b_h][i_b][i_k][0]$ 
12        batchreduce_gemm( $A_{parr}, B_{parr}, s[t][i][i_k], C_h$ )
13       for  $i_b = 0, \dots, K_h-1$  do
14         $A_{parr}[i_k] = \alpha R_i[b_h][i_b][i_k][0]$ 
15         $B_{parr}[i_k] = \beta h[t-1][i][i_b - b_h]$ 
16        batchreduce_gemm( $A_{parr}, B_{parr}, s[t][i][i_k], K_h$ )
17        $i[f][i_k] \leftarrow \sigma(i[f][i_k])$ 
18       Dirac for blocks of  $c_t, o_t$  via Equations 2-4
19        $s[t][i][i_k] \leftarrow f[t][i][i_k] \circ s[t-1][i][i_k] +$ 
20        $i[f][i_k] \circ c[t][i][i_k]$ 
21        $h[t][i][i_k] \leftarrow o[t][i][i_k] \circ \tanh(s[t][i][i_k])$ 
```

• We adopt a “dataflow” based approach for optimizations

• Blocked layout to exploit locality & avoid conflict misses

• Perform computation with fused time steps

• Also, allow blocked inputs / weights to be passed directly from framework

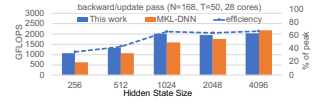
• Use JIT batch-reduce GEMM kernels

– Implement optimized blocked GEMM

– Implement fused kernel for elementwise operations ( $i_t, f_t, o_t, c_t, s_t, h_t$ )

– Once a block of batch-reduce GEMM is computed, apply element-wise operations on it while hot in cache

• Same optimizations applied to backward/update passes



• MKL-DNN approach: large GEMMs followed by elementwise ops – potentially poor locality for elementwise part

• For small/medium sized problems, our implementation of LSTM forward pass is up to 1.4x faster than the MKL-DNN cell, while for backward/weight update it is up to 1.3x faster on Xeon 8180.

• Large weight matrices yield similar performance

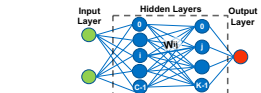
– Cubic GEMM scaling VS quadratic elementwise scaling

Evangelos Georganas<sup>1</sup>, Kunal Banerjee<sup>2</sup>, Dhiraj Kalamkar<sup>2</sup>, Sasikanth Avancha<sup>2</sup>, Anand Venkat<sup>1</sup>, Michael Anderson<sup>1</sup>, Greg Henry<sup>3</sup>, Hans Pabst<sup>4</sup>, Alexander Heinecke<sup>4</sup>

<sup>1</sup>Intel Corporation, Intel Labs, Mission College Boulevard 2000, Santa Clara 95054, CA, USA  
<sup>2</sup>Intel Technology India Private Limited, Intel Labs, #2350F, Dewandehra Hall, Outer Ring Road, Bangalore-560113, India  
<sup>3</sup>Intel Corporation, Software and Services Group, 2111 NE 25<sup>th</sup> Avenue, Hillsboro 97124, OR, USA  
<sup>4</sup>Intel Semiconductor AG, Software and Services Group, Badenstrasse 549, 8048 Zurich, Switzerland



## Multilayer Perceptron (MLP)



• An input layer  $x \in \mathbb{R}^C$  is mapped to an output layer  $y \in \mathbb{R}^N$  via:  $y = W \cdot x$ , where  $W \in \mathbb{R}^{K \times C}$  is the weight tensor of the connections between the neurons.

• During the training process,  $N$  multiple inputs ( $N$  being mini-batch size) are grouped together yielding:  $Y = W \cdot X$  with  $W \in \mathbb{R}^{K \times C}$ ,  $X \in \mathbb{R}^{C \times N}$  and  $Y \in \mathbb{R}^{K \times N}$ .

• After the output tensor  $Y$  is computed, a non-linear activation function  $g(\cdot)$  is applied on it.

### Algorithm 2 Forward pass of Fully Connected Layer

```
Inputs: Weight  $W[K_h][C_h][b_h]$ , Input  $X[N_h][C_h][b_h]$ 
Outputs: Output  $Y[N_h][K_h][b_h]$ 
1 Based on  $thread\_id$  calculate  $K_{h\_start}, K_{h\_end}, N_{h\_start}$  and  $N_{h\_end}$  to assign output work items
2 for  $i_b = N_{h\_start}, \dots, N_{h\_end}$  do
3   for  $i_b = N_{h\_start}, \dots, N_{h\_end}$  do
4     Prepare batch-reduce GEMM arguments
5     for  $i_b = 0, \dots, C_h-1$  do
6        $A_{parr}[i_k] = \alpha W[b_h][i_b][i_k][0]$ 
7        $B_{parr}[i_k] = \beta X[i_b][i_k][0]$ 
8        $Out = \alpha Y[i_b][i_k][0]$ 
9       batchreduce_gemm( $A_{parr}, B_{parr}, Out, C_h$ )
10       $Y[i_b][i_k][0] \leftarrow g(Y[i_b][i_k][0])$ 
```

• We implement the multiplication  $Y = W \cdot X$  via batch-reduce GEMM calls

• Apply the non-linear activation function  $g(\cdot)$  on the batch-reduce GEMM outputs once they are still hot in cache.

• Blocked layout to exploit locality & avoid conflict misses



• On Xeon 8180, the primitive based on batch-reduce GEMM outperforms MKL-DNN by 1.1x to 1.3x.

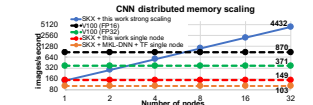
## Distributed memory results

• 32-node cluster (Intel Omnipath interconnect), each node having two Skylake-SP (SKX) 8180 processors

• ResNet-50 training: integrated our CNN kernels into the GxM framework which scales efficiently on CPU clusters

• GNMT training: We replaced the Tensorflow LSTM cell implementation with our optimized LSTM cell that leverages the batch-reduce GEMM kernel.

• Exceeded the performance of vendor-optimized implementations by up to 2.3x

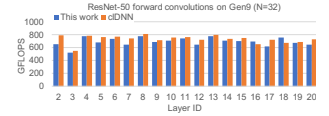


• LSTM distributed memory scaling



## Performance portability on integrated GPU

- We developed the batch-reduce GEMM kernel in OpenCL.
- We implemented PoC forward propagation CNN kernels targeting Intel's integrated GPU Gen9 (Core i7 6770HQ) by leveraging the batch-reduce GEMM kernel.
- Our approach is within 3% of the vendor-optimized, ad hoc implementation (cIDNN).

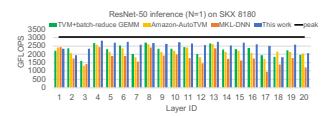


## PoC integration within TVM

• We provide to TVM the forward propagation loop recipes in high-level Python code, and at the innermost loop nest level we invoke our batch-reduce GEMM kernel.

• Benchmarked inference use-case, ResNet-50 forward kernels (mini-batch  $N=1$ ) on the SKX 8180 platform

• Our TVM implementation is 2% faster than the Amazon-AutoTVM *auto-tuned* code and 1.24x faster than the vendor optimized MKL-DNN library.



## Summary and Conclusions

On this poster we presented the following work:

• we showed how the most popular DL algorithms (RNN/LSTM, CNN and MLP) can be formulated with batch-reduce GEMM as basic building block

• we demonstrated that our methodology outperforms vendor-optimized, low-level DL primitives by factors up to 1.4x

• we integrated our DL kernels into distributed frameworks, and optimized end-to-end workflows for GNMT and ResNet-50 training: we exceeded the performance of vendor-optimized implementations by up to 2.3x

• we highlighted the architectural-agnostic aspect of our methodology by matching the CNN kernel performance of a vendor-provided library on integrated GPUs

• we prototyped CNN kernels in a tensor compiler framework by harnessing our batch-reduce GEMM kernel, and matched the performance of *auto-tuned* inference TVM primitives

• we upstreamed all changes to LIBXSMM and they can be found here: <https://github.com/hfp/libxsmm>. As future work, we plan to extend our DL primitives for a wider set of architectures/workloads, and also we intend to experiment with Tensor compilers' automatic polyhedral optimization, e.g. PlaidML.

For references of this work, please refer to the extended poster abstract and/or our extensive technical report: <https://arxiv.org/abs/1906.06440>.