

Efficient Sparse Data Computations for Deep Learning

Abstract

Deep learning has recently emerged as an important machine learning approach, with big deep neural networks (DNN) models trained on vast amounts of data demonstrating state-of-the-art accuracy on important yet challenging artificial intelligence tasks, such as image and speech recognition. However, training big DNN models using large training data is both compute and memory intensive, making distributed training on a cluster of server machines, leveraging the aggregate system resources, the standard approach.

This paper proposes hardware techniques for improving system performance and scalability for DNN training workloads by exploiting the sparse nature of computation to reduce the compute and memory requirements. Our techniques improve training efficiency by avoiding resource utilization on sparse data values (i.e., zeroes) which do not impact training quality. Our design is transparent to software, enabling existing codes to enjoy a performance boost without modifications.

Simulation-based evaluation using standard image recognition workloads shows that our techniques can improve DNN training performance significantly and outperform software approaches.

1. Introduction

Deep learning has recently attracted significant attention because of the state-of-the-art performance of deep neural networks (DNNs) on important but challenging artificial intelligence tasks, such as image recognition [26, 28, 11, 7], speech recognition [10, 18, 16], and text processing [8, 9, 31]. A key driver of these machine learning advancements is the ability to train big DNN models (billions of parameters) using large amounts of examples (TBs of data). However, the compute and memory resources required to train big models to reasonable accuracy in a practical amount of time (days instead of months) are significantly high, and surpass the capabilities of a single commodity server. Thus, big DNN models are, in practice, trained in a distributed fashion using a cluster of 100s/1000s of servers, leveraging parallel hardware resources [11, 7]. Addressing the high computational costs of DNN training is critical to sustain task accuracy improvements through model and data scaling.

This paper presents a hardware approach that exploits the computation pattern of DNN training to improve performance and scalability by reducing the compute and cache resource requirements. Our approach is based on the observation that the computation data of training are

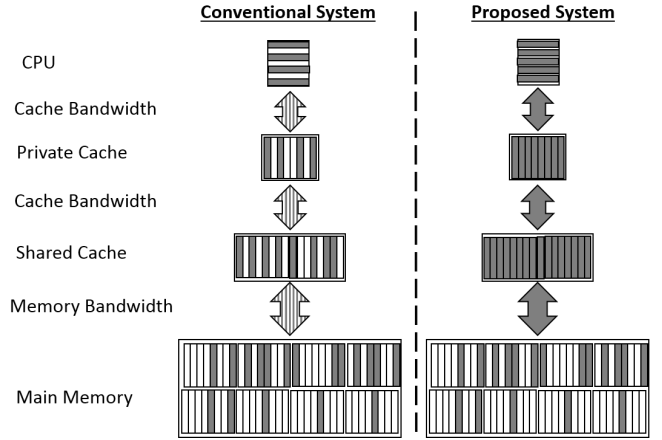


Figure 1: Processor and memory system utilization of sparse (white) and non-sparse (shaded) data.

significantly sparse, and since training kernels are dominated by multiply-accumulate operations, a significant portion of these computations are redundant to the training objective. We therefore improve training performance by avoiding the compute and memory system consumption of sparse data and the associated computations without harming model quality. The performance benefits should be larger for big DNN models where system resources (e.g., bandwidth) are typically oversubscribed.

Figure 1 illustrates a high-level comparison of processor and memory system utilization in a conventional (left) and our proposed (right) system. Resource utilization on zeroes (word granularity in the processor and cache line granularity in the memory system) are white, while those used on other values are shaded grey. Compared to a conventional system, our proposed system eliminates or significantly reduces the resource consumption for zero data computations to improve the performance of useful data computations.

Our proposed optimizations can improve DNN training performance and scalability in three ways. First, by eliminating computations on useless data while processing a training example, more iterations over the data set can be completed within a time budget (e.g., a week), which can improve model quality. Second, memory (and cache) bandwidth is a key bottleneck for *data-parallelism* within a machine (i.e., processing multiple examples at once using multiple CPU cores). By reducing the bandwidth utilization for each example, data-parallelism can scale training throughput more effectively. Finally, a standard approach for fitting big DNN models into the

last level cache is partitioning the model across multiple machines to exploit the aggregate cache capacity (a.k.a., *model-parallelism*). By reducing the cache consumption of a model partition, our techniques can increase the per-machine partition size and reduce the number of machines required to achieve a target training throughput, reducing model parallelism costs.

Prior software [14, 20] and hardware [6, 24, 15] approaches for sparse matrix-vector multiplications are less effective for DNN training for two reasons. First, those techniques assume that sparsity exists only in matrices but not in vectors, whereas in DNN training both matrices and vectors can be sparse. Second, those techniques assume that the sparse data structures are static, so that the cost of constructing a sparse representation is amortized over many uses. However, the matrices and vectors in DNN training change for each training example, thus representation cost is incurred repeatedly, hurting performance.

Our approach consists of independent techniques for tackling sparse data computation overheads in the processor and memory system. By separating these concerns, we achieve a flexible and modular design which makes it easy to evaluate the benefits of each technique. Our processor extensions are based on “zero-optimizable” instructions, which are arithmetic instructions (e.g., multiplication) whose results and side effects are pre-determined when an input operand is a zero. Our optimizations exploit zero-optimizable instructions to reduce execution cycles and pressure on critical processor resources. Our memory system extensions efficiently track zero data at cache granularity in the caches and main memory to avoid the bandwidth costs associated with moving zero cache lines. Our approach does not require software modifications and so existing binaries can benefit from our optimizations.

We have evaluated our proposed hardware extensions in a simulation environment using real-world DNN training workloads for image recognition. The results show that our approach can significantly improve DNN training performance in single threaded, multi-threaded, and model-parallelism scenarios.

This paper makes the following contributions.

- We propose hardware techniques for improving the performance and scalability of DNN training by reducing computational requirements of sparse data computations, without requiring software changes.
- We study the impact of sparse data on computation in a real-world DNN training for an image recognition task.
- We present a detailed design of our proposed hardware extensions, which add negligible logic on the critical path of processor execution and memory accesses.
- We quantitatively evaluate how our optimizations improve DNN training performance and scalability using standard image recognition workloads.

The rest of the paper is organized as follows. Section 2 provides background on DNN and DNN training. Sec-

tion 3 studies sparse data computations in DNN training. Our processor optimizations are described in Section 4, while our cache optimizations are described in Section 5. We present our evaluation results in Section 6, review related work in Section 7, and conclude in Section 8.

2. Background

2.1. Deep Neural Networks

DNNs consist of large numbers of neurons with multiple inputs and a single output called an activation. Neurons are connected hierarchically, layer by layer, with the activations of neurons in layer $l - 1$ serving as inputs to neurons in layer l . This deep hierarchical structure enables DNNs to learn complex AI tasks, such as image recognition, speech recognition and text processing [3]. DNNs comprise *convolutional* layers (possibly interleaved with *pooling* layers) at the bottom of the hierarchy followed by *fully connected* layers. Convolutional layers, which are inspired by the visual cortex [19], extract features from input samples, and consist of neurons that are only connected to spatially local neurons in the lower layer [29]. Pooling layers summarize the features learned by convolutional layers (e.g., identify the maximum intensity in a cluster of image pixels, reduce spectral variance in speech samples). Fully connected layers classify the learned features into a number of categories (e.g., handwritten digits) and consist of neurons that are connected to all neurons in the lower layer.

2.2. DNN Training

A common approach for training DNNs is using learning algorithms, such as stochastic gradient descent (SGD) [5], and labeled training data to tune the neural network parameters for a specific task. The parameters are the *bias* of each neuron and the *weight* of each neural connection. Each training input is processed in three steps: *feed-forward evaluation*, *back-propagation*, and *weight updates*.

Feed-forward evaluation: Define a_i as the activation of neuron i in layer l . It is computed as a function of its J inputs from neurons in the preceding layer $l - 1$:

$$a_i = f \left(\left(\sum_{j=1}^J w_{ij} \times a_j \right) + b_i \right), \quad (1)$$

where w_{ij} is the weight associated with the connection between neurons i at layer l and neuron j at layer $l - 1$, and b_i is a bias term associated with neuron i . The activation function, f , associated with all neurons in the network is a pre-defined non-linear function, typically sigmoid or hyperbolic tangent.

Back-propagation: Error terms δ are computed for each neuron i in the output layer L :

$$\delta_i = (\text{true}_i - a_i) \times f'(a_i), \quad (2)$$

where $\text{true}(x)$ is the true value of the output and $f'(x)$ is the derivative of $f(x)$. These error terms are back-propagated to each neuron i in the layer l from its S connected neurons in layer $l + 1$:

$$\delta_i = \left(\sum_{s=1}^S \delta_s \times w_{si} \right) \times f'(a_i). \quad (3)$$

Weight updates: These error terms are used to compute the weight deltas, Δw_{ij} , for updating the weights:

$$\Delta w_{ij} = \alpha \times \delta_i \times a_j \text{ for } j = 1 \dots J, \quad (4)$$

where α is the learning rate and J is the number of neurons of the layer.

This process is repeated for each input until the entire training data has been processed, which constitutes a training *epoch*. Typically, training continues for multiple epochs, reprocessing the training data set each time, until the error converges to a desired (low) value.

3. Sparsity in DNN Training

In this section we motivate our work and project the benefits of our optimizations for DNN training. First, we provide some intuition of why sparsity exists in DNN training. Next, using an image recognition workload, we empirically demonstrate the amount of sparsity that exists in practice. Finally, we identify the portion of available sparsity that can be exploited by our techniques to improve training efficiency.

3.1. Sources of Sparsity

Machine learning experts have long observed that DNN training using back-propagation and gradient descent involves a considerable amount of computations on sparse matrices and vectors [33, 32, 26, 4, 37]. The performance-critical data of training are neuron activations and errors (both implemented as vectors) and synaptic weights and corresponding deltas (both implemented as matrices) can be sparse. Some of the sparsity arise naturally from the training process and its matrix-vector multiplication kernel. For example, correct predictions of a neuron’s output activation, during feed-forward evaluation, result in zero-valued neuron error terms, during back-propagation, which can introduce sparsity in the rest of the network. Beyond this, standard techniques for boosting training quality often introduce additional sparsity in the network. These include techniques such as Rectified Linear Units (ReLUs) [32, 26] for faster convergence, and L_1 [33, 4] and Dropout [37] regularization methods for reducing overfitting. [Trishul to help with this content](#)

3.2. Sparsity in real-world image recognition task

For a better insight into the amount of sparsity in real-world DNN training workloads, we profile training on *CIFAR-10* [25], a standard image recognition task (described in 6.1). In our study, we reason about sparsity from 2 perspectives: (i) data sparsity and (ii) computation sparsity. Specifically, data sparsity measures the amount of zeroes in performance-critical data (e.g., activation vectors), while computation sparsity measures the amount of multiply-accumulate operations performed on zero values in the main phases of training (e.g., feed-forward evaluation). Both perspectives are useful because they capture different effects of sparsity on system performance. Memory capacity bandwidth impact is captured by data sparsity, processing cycles impact is captured by computation sparsity, and bandwidth impact is captured by both data and computation sparsity. We measure both sparsity metrics over 10 training epochs using a data set of 60000 images.

Data Sparsity Figure 2(a) illustrates the sparsity of activation and error vectors, and weight delta and weight matrices in CIFAR-10 training. We see that the sparsity amount and rate of change is quite different among the data structures. While the weight matrix is dense, the activation and error vectors and the weight delta matrix are noticeably sparse. We see that sparsity generally increases with training epochs, albeit at varying rates. The error vector has the greatest amount of sparsity (83%—85%), followed by the weight delta matrix (66%—83%), and finally the activation vector (26%—33%). The results show the memory/cache consumption of activations, errors, and weight deltas for this workload can be reduced significantly.

Computation Sparsity Figure 2(b) reports the computation sparsity of the different training steps. Compared to data sparsity, the results illustrate how the different vectors and matrices are combined through multiplication and addition operations. For example, the sparsity in feed-forward evaluation is the result of multiplying the dense weight matrix and sparse activation vector. We see that considerable sparsity exists in each training step (from 29% for feed-forward evaluation to 92% for computing weight deltas). We also see that the amount of sparsity generally grows with training epochs (e.g., 29%—42% for feed-forward evaluation). In summary, these results shows the potential for significant saving in processing cycles by eliminating cycle consumption for generating zero values that do not impact training quality.

3.3. Sparsity at cache line granularity

Our proposed hardware mechanisms track data sparsity at cache line granularity, so it is unlikely that we can fully exploit the amounts of sparsity presented above because

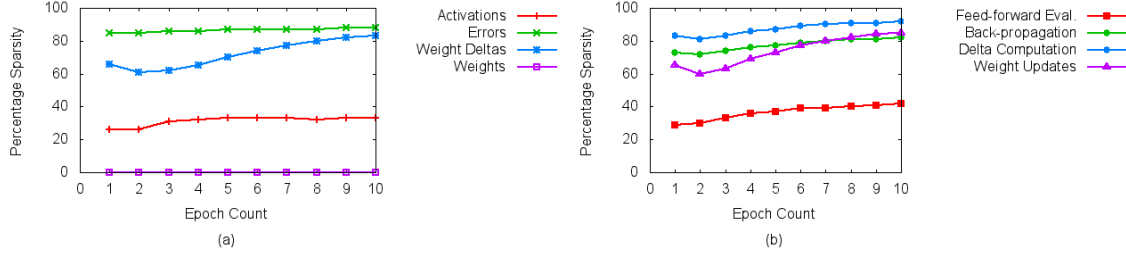


Figure 2: (a) Data and (b) computation sparsity in CIFAR-10 training.

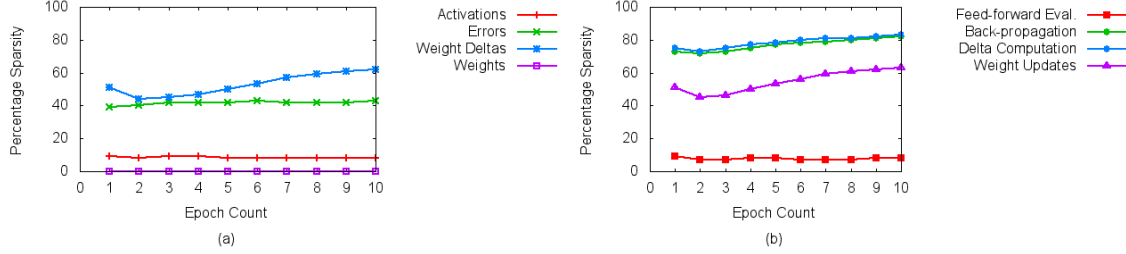


Figure 3: (a) Data and (b) computation sparsity in CIFAR-10 training at cacheline granularity.

the profiling was conducted at a finer granularity (e.g., individual activation values). To get a more accurate view of the effectiveness of our optimizations we repeat the profiling study at a cache line granularity. Since data values are represented as 4-byte floats (or word), each cache line contains up to 16 data values. We measure sparsity at cache line granularity in the following manner. Data sparsity represents the percentage of cache lines of a data structure that contain only zero values. For computation, we adopt a coarse-grained view of computations, i.e., a unit of computation operates on a pair of cache lines (e.g., an activation cache line and a weight cache line in feed-forward evaluation). Thus, computation sparsity represents the percentage of such computations that operate on a sparse cache line. The results are presented in Figures 3(a) and 3(b) for data and computation sparsity respectively.

Figure 3(a) shows that data sparsity at cache line granularity is generally lower compared to word granularity (Figure 2(a), e.g., error sparsity is about half). This indicates that sparse data values are not clustered, which limits the capacity savings achievable by our approach. However, the computation sparsity results in Figure 3(b) presents a more promising picture. We see that computation sparsity at cache line granularity is not much lower than sparsity at word granularity for 3 phases of training: back-propagation, delta computation, and weight updates. This indicates that even though the non-sparse cache line ratio is relatively high, non-sparse cache lines are more likely to be combined with sparse cache lines leading to sparse computations. These results suggest that reducing the cycles and bandwidth consumption of sparse data computations can yield significant performance benefits.

4. Processor Optimizations

Our processor optimizations are based on the observation that certain arithmetic operations, such as addition and multiplication, which are performance critical in training computations have predetermined results when one of the input operands is a zero. We refer to machine instructions that perform such arithmetic operations as “zero-optimizable” instructions. Exploiting zero-optimizable instructions to improve training performance is promising because, as shown in our profiling studies, a significant portion of training computations involve zeroes.

4.1. Opportunities

Zero-optimizable instructions present a number of opportunities to improve the ILP and resource pressure of training computations on modern out-of-order processors. These opportunities arise because of the predetermined results of a zero-optimizable instruction with a zero input operand which can make some data dependencies and pipeline stages redundant for the instruction and dependent instructions. Thus, zero-optimizable and dependent instructions can be issued or committed earlier than normal or skipped completely in program execution, as discussed below.

4.1.1. Training Code Example. We use the training code snippets presented in Figure 4 to describe our optimizations. Figure 4(a) illustrates a simplified version of the code for computing the gradients to update the weights of a layer during back-propagation. Gradients are computed as an inner product of the activation and error vectors. From Figure 2, we can see that a promising optimization is to skip multiply and addition operations in the inner

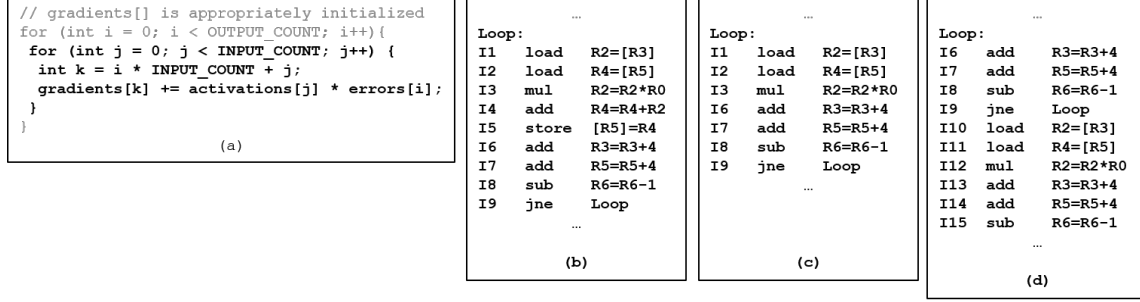


Figure 4: (a) Source code for computing gradients, (b) machine code of inner loop, (c) optimized code after basic instruction quashing, and (d) optimized code after advanced instruction quashing.

loop if *activations[j]* or *errors[i]* is zero. Moreover, if *errors[i]* is zero then the inner loop execution can be skipped.

Although these optimizations could be implemented in software, such an approach has a couple of practical limitations. First, it requires software modifications which might not be possible for existing binaries. Second, it incurs the runtime overheads of software checks for zeroes, which could be significant. For example, checking *errors[i]* is likely beneficial because it can be done outside the inner loop and skips a large amount of computation. In contrast, checking *activations[j]* is likely harmful as it occurs in the inner loop and skips a small amount of computation. We compare the performance benefits of software and hardware approaches in our evaluation.

We use the machine code sequences in Figures 4(b), 4(c), and 4(d), which correspond to the inner loop, to illustrate the impact of our optimizations on instructions in the instruction queue. Although there are six zero-optimizable instructions in the loop (I3, I4, I6, I7, and I8), the optimizations discussed involve only I3 and I4. We assume that R0, which corresponds to *errors[i]*, is zero.

4.1.2. Early Instruction Issue/Commit. First, a zero-optimizable instruction can be issued once the zero operand is available if it makes other operands redundant. For example, I3 can be issued early because it is a multiplication and the zero value of R0 makes R2 redundant. Second, a zero-optimizable instruction could be committed early if the zero input determines its results and side effects. This is also the case for I3. Early issue and commit of zero-optimizable instructions can reduce pressure on processor resources and wait times of data dependent instructions, such as I4, since the dependencies are satisfied sooner.

4.1.3. Instruction Squashing. A zero-optimizable instruction can be squashed in the instruction queue if a zero input operand makes it an identity function and thus redundant. For this reason, I4 can be squashed since it is an addition and R2 is zero. Squashing an instruction can make the instructions that it depends on (producers) and those that depend on it (consumers) redundant, leading to more instruction squashing. For example, I5 becomes

redundant (a silent store) and can be squashed, if I5 is squashed. Figure 4(c) shows the impact of squashing I4 and I5. We further observe that I1, I2, and I3 are now redundant in all but the last loop iteration, since their results (R2 and R4) are not used. We can squash these three instructions in all but the last iteration as shown in Figure 4(d). Compared to the original machine code sequence, the optimized code sequence will run much faster because of the squashed instructions, especially loads which often have high latency. Thus, by exploiting zero-optimizable instructions we can improve the performance of the inner loop of the gradient computation code.

4.2. Mechanisms

Our optimizations can be realized with minor extensions to the front-end processing of a modern out-of-order processor. The extensions are lightweight, and despite being on the critical path should not introduce noticeable execution delays. Specifically, we propose processor extensions to do the following operations: (i) identify zero-optimizable instructions, (ii) detect when zero-optimizable input operand is a zero, (iii) modify producer and consumer data dependencies, and (iv) squash instructions. These steps can be done in parallel with existing pipeline front-end stages, as we discuss below.

4.2.1. Identify Zero-Optimizable Instructions. We can detect zero-optimizable instructions during instruction decoding by matching the opcode against a predefined set of opcodes. Since only a small set of arithmetic instructions qualify as zero-optimizable, the storage requirements of the opcode set is modest, and opcode matching can be done in parallel to avoid extra delays. Zero-optimizable instructions that are identified in the decode stage are marked for easy identification in later pipeline stages.

4.2.2. Detect Zero Operands. We can detect zero input operands while a zero-optimizable instruction is waiting in the instruction queue for data dependencies. Current mechanisms for signaling operand availability can be extended to also indicate whether or not the value is zero.

4.2.3. Modify Data Dependencies. We can extend current mechanisms for tracking data dependencies among

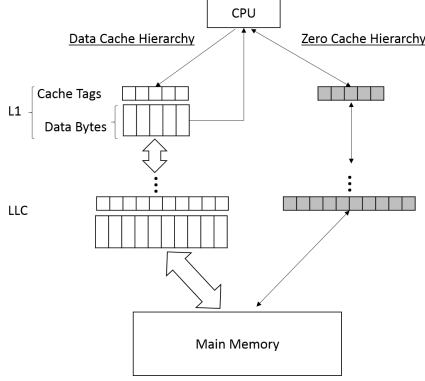


Figure 5: A Memory System with Zero Cache Hierarchy.

instructions to clear dependencies of zero-optimizable instructions that become redundant due to a zero input operand becoming available. Furthermore, dependencies from instructions that consume the results of a zero-optimizable instruction should be cleared when a zero input makes the instruction an identity function, and thus redundant.

5. Cache Optimizations

Our cache optimizations for DNN training are based on the sparse nature of the performance critical data (e.g., activations, errors, etc.). Our approach improves cache performance through a compact representation of cache lines containing only zeroes (a.k.a. *zero cache lines*) in the caches, which helps to avoid the normal bandwidth and storage costs of zero cache lines. These optimizations enable efficient scaling of model size and training threads.

Managing zero data at cache line granularity enables implementation of our optimizations through simple and efficient extensions of existing memory systems. Our current design comprises of mechanisms for achieving the following: (i) compact representation of zero cache lines, (ii) a decoupled cache hierarchy for zero cache lines, and (iii) tracking zero cache lines in the memory system. We describe these mechanisms in the rest of this section.

5.1. Zero Cache Line Representation

Our compact representation exploits the fact that the data bytes of a zero cache line are not required to represent the line in cache, the cache tag is sufficient for this purpose. Also, it is not necessary to transfer the data bytes of a zero cache line across the caches since they can be synthesized in the processor (read) or main memory (on a writeback) as appropriate. However, in event of a cache hit, we must quickly determine whether it is a zero cache line that is referenced so that the appropriate data transfer is done promptly. We consider two alternatives for handling this: (i) an extra bit in the cache tags to identify zero cache lines, or (ii) a decoupled hierarchy of cache tags for zero cache lines. Although the first option avoids the extra

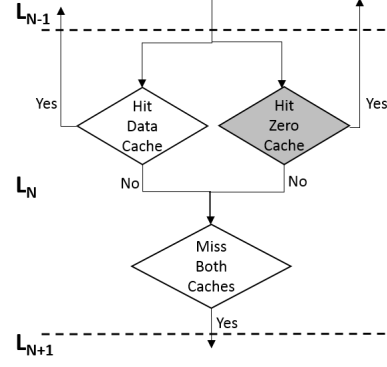


Figure 6: Handling read requests.

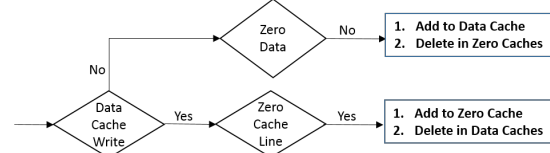


Figure 7: Handling processor writes.

cost of zero cache line tags, the data bytes space of zero cache lines are unused. To avoid this waste, we adopt the second option in our current work.

5.2. Hierarchy of Zero Cache Lines

Figure 5 illustrates a memory system that is augmented with a cache hierarchy for zero cache lines, which we call the *zero cache* hierarchy. The zero cache hierarchy is a multi-level structure with caches (a.k.a., *zero caches*) containing tags but no data bytes. Since zero cache lines are not maintained in the conventional data caches, both cache hierarchies are mutually exclusive. The zero cache hierarchy and the data cache hierarchy have the same number of levels, and can additionally share other properties, such as number of entries, ways, associativity, replacement policies, etc. The coherence of zero caches is maintained across cores using the same protocol as the data caches.

Data access requests from the processor are satisfied by accessing the two cache hierarchies in parallel to avoid introducing extra latency. Figure 6 shows the processing of a read request by the N th level caches. The request is processed in parallel by the data and zero caches, and forwarded to the next level if it is a miss in both. If the request is a hit in either cache, then the appropriate response is sent to the processor or lower levels of the cache hierarchy. The data cache responds, as normal, with the requested data bytes (or cache line), while the zero cache responds by signaling a zero cache line hit.

5.3. Tracking Zero Cache Lines

Our optimizations are based on the invariant that cache lines reside in the appropriate cache hierarchy: zero cache

lines in zero caches and other cache lines in the data caches. To maintain this invariant, we track the zero status of cache lines to ensure that a cache line is placed in the right hierarchy in the following events: (i) update by processor writes, (ii) cache fill from main memory, and (iii) writebacks from lower level caches (e.g., due to evictions). Our tracking operations do not increase cache access latencies as they execute off the critical path of cache accesses. We leverage zero detector hardware [13] to detect that an entire cache line (i.e., 32/64 bytes) is zero.

5.3.1. Processor Writes. The zero-status of a cache line can be changed by a processor write depending on the current status and write data. The following four situations could arise: (i) write zeroes to a zero cache line, (ii) write non-zeroes to a non-zero cache line, (iii) write non-zeroes to a zero cache line, and (iv) write zeroes to a non-zero cache line. The first two situations are irrelevant since the non-zero status of the cache line is unchanged. Writing a non-zero value to a zero cache line moves the cache line to the data cache in the the same level, and removes the cache line from the zero cache hierarchy. This may require data cache evictions to accomodate the new cache line. Writing zeroes to a non-zero cache line moves the cache line from the data caches into the zero-cache hierarchy if the cache line contains only zeroes after the update. Naturally, the cache tag is moved as well. Figure 7 illustrates how cache updates by processor writes are handled to ensure that cache lines reside in the right hierarchy.

5.3.2. Cache Fills from Main Memory. Since our optimizations are focused on the cache capacity and bandwidth, the data bytes of zero cache lines are stored in main memory, similar to other cache lines. We extend the memory controller to avoid sending data bytes when handling a cache fill request for a zero cache line. Requests for non-zero cache lines are handled normally.

5.3.3. Writebacks from Lower Level Caches. Our decoupled cache hierarchies approach implicitly handles writebacks from lower level caches because the zero-status of a cache line is unchanged. Thus, the data caches are not involved by zero cache writebacks, and vice versa.

6. Evaluation

We now evaluate the effectiveness of our processor and memory system optimizations for DNN training. We conduct our evaluations along 3 dimensions: (i) the impact on single thread performance (6.2), (ii) the impact on multi-threading scalability in a server (6.3), and (iii) the impact on model parallelism performance (6.4).

6.1. Methodology

Image Recognition Task: Although we expect our optimizations to be effective in general for training with gradient descent methods, we focus on image recog-

```

...
for (i = 0; i < OutputNeuronCount; i++) {
    if (Error[i] != 0) {
        for (j = 0; j < InputNeuronCount; j++) {
            ... += Error[i] * weights[i,j]
        }
    }
}
...

```

Figure 8: Zero error signal optimization in back-propagation of a linear layer.

nition because it represents an important class of AI problems for which significant accuracy improvements have been achieved through gradient descent training [26, 28, 11, 7, 17]. Specifically, we measure the impact of our optimizations on the training of high quality DNN models on 3 common image recognition workloads: (i) *MNIST* [30], (ii) *CIFAR-10* [25], and (iii) *ImageNet* [12]. We describe each benchmark and corresponding DNN in more details below.

- **MNIST:** The task is to classify 28x28 grayscale images of handwritten digits into 10 categories. The DNN is relatively small, containing about 2.5 million connections in 5 layers: 2 convolutional layers with pooling, 2 fully connected layers, and a 10-way output layer [7].
- **CIFAR-10:** The task is to classify 32x32 color images into 10 categories. The DNN is moderately-sized, containing about 28.5 million connections in 5 layers: 2 convolutional layers with pooling, 2 fully connected layers, and a 10-way output layer [26].
- **ImageNet:** The task is to classify 256x256 color images from a dataset of about 15 million images into a number of categories. There are 2 standard versions of this benchmark: (i) classifying 1.2 million images into 1000 categories (a.k.a., *ImageNet-1K*), and (ii) classifying the entire data set into 22000 categories (a.k.a. *ImageNet-22K*). We used the largest ImageNet task (i.e., *ImageNet-22K*), which is to classify 256x256 color images into 22,000 categories. This DNN is extremely large, containing over 2 billion connections in 8 layers: 5 convolutional layers with pooling, 2 linear layers, and a 22,000-way output layer [7].

Comparison to Software Approach: We compare our technique to a software approach that avoids zero-value computations without compact representations of sparse matrices or vectors, unlike CSR [20]. As discussed earlier, the dynamic nature of sparsity in training makes CSR less effective because the construction cost of the representation is incurred for each example. Rather, we test for zero values and skip the corresponding multiply-add operations. To derive the most benefit, this optimization is only applied when multiple multiply-add operations can be skipped for each zero value, such as during back-propagation or weight updates computation. The code

snippet in Figure 8 illustrates this optimization for back-propagation of a linear layer, where computations involving an entire row of the weight matrix can be skipped for a zero error signal.

Simulation-based Approach: We conduct our performance experiments in a simulation environment, which allows us to easily prototype our proposed memory system extensions (Section ??). We use a modified version of Memsim [1, 34], a multi-core simulator that models out-of-order cores coupled with a DDR3-1066 [22] DRAM simulator. All systems use a three-level cache hierarchy with a uniform 64B cache line size. We do not enforce inclusion in any level of the hierarchy. We use the state-of-the-art DRRIP cache replacement policy [?] for the last-level cache. All our evaluated systems use an aggressive multi-stream prefetcher [36] similar to the one implemented in IBM Power 6 [27].

6.2. Single Thread Performance

Here we show that our techniques improve single thread training performance by reducing the number of computations performed per training example.

6.3. Multi-threaded Performance

Here we show improved scalability by reducing the cache capacity and bandwidth pressure.

6.4. Model-parallelism Performance

Here we show improved scalability of model-parallelism because the reduced cache capacity and bandwidth pressure allows bigger models to fit on a machine.

7. Related Work

Our work is related to prior work on efficient sparse matrix-vector computations [14, 20, 6, 24, 15, 35], and reducing the cache [38, 39, 13, 21] and physical register file [23, 2] overheads of zero values.

8. Conclusion

References

- [1] Memsim. <http://safari.ece.cmu.edu/tools.html>, 2012.
- [2] Saisanthosh Balakrishnan and Gurindar S. Sohi. Exploiting value locality in physical register files. In *MICRO*, 2003.
- [3] Yoshua Bengio. Learning deep architectures for AI. *Foundations and Trends in Machine Learning*, 2009.
- [4] Yoshua Bengio, Nicolas Boulanger-Lewandowski, and Razvan Pascanu. Advances in optimizing recurrent networks. In *Proc. ICASSP 38*, 2013.
- [5] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *COMPSTAT*, 2010.
- [6] John B. Carter, Wilson C. Hsieh, Leigh Stoller, Mark R. Swanson, Lixin Zhang, Erik Brunvand, Al Davis, Chen-Chi Kuo, Ravindra Kuramkote, Michael A. Parker, Lambert Schaelicke, and Terry Tateyama. Impulse: Building a smarter memory controller. In *HPCA*, 1999.
- [7] Trishul Chilimbi, Johnson Apacible, Karthik Kalyanaraman, and Yutaka Suzue. Project adam: Building an efficient and scalable deep learning training system. In *OSDI*, 2014.
- [8] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *ICML*, 2008.
- [9] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel Kuksa. Natural language processing (almost) from scratch. *The Journal of Machine Learning Research*, 2011.
- [10] G. E. Dahl, Dong Yu, Li Deng, and A. Acero. Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *Trans. Audio, Speech and Lang. Proc.*, 2012.
- [11] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc’Aurelio Ranzato, Andrew W. Senior, Paul A. Tucker, Ke Yang, and Andrew Y. Ng. Large scale distributed deep networks. In *NIPS*, 2012.
- [12] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [13] Julien Dusser, Thomas Piquet, and André Seznec. Zero-content augmented caches. In *ICS*, 2009.
- [14] C. Stanley Eisenstat, MC Gursky, H. Martin Schultz, and H. Andrew Sherman. Yale sparse matrix package i: The symmetric codes. *International Journal for Numerical Methods in Engineering*, 1982.
- [15] Jeremy Fowers, Kalin Ovtcharov, Karin Strauss, Eric Chung, and Greg Stitt. A high memory bandwidth fpga accelerator for sparse matrix-vector multiplication. In *International Symposium on Field-Programmable Custom Computing Machines*, 2014.
- [16] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, et al. Deepspeech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*, 2014.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *CoRR*, 2015.
- [18] Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, and Brian Kingsbury. Deep neural networks for acoustic modeling in speech recognition. *Signal Processing Magazine*, 2012.
- [19] D.H. Hubel and Wiesel T.N. Receptive fields of single neurons in the cat’s striate cortex. *Journal of Physiology*, 1959.
- [20] Intel. Sparse Matrix Storage Formats, Intel Math Kernel Library. <https://software.intel.com/en-us/node/471374>.
- [21] Mafijul Md. Islam and Per Stenstrom. Zero-value caches: Cancelling loads that return zero. In *PACT*, 2009.
- [22] JEDEC. DDR3 SDRAM, JESD79-3F, 2012.
- [23] Stephen Jourdan, Ronny Ronen, Michael Bekerman, Bishara Shomar, and Adi Yoaz. A novel renaming scheme to exploit value temporal locality through physical register reuse and unification.
- [24] Srinidhi Kestur, John D. Davis, and Eric S. Chung. Towards a universal fpga matrix-vector multiplication architecture. In *FCCM*, 2012.
- [25] Alex Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, Computer Science Department, University of Toronto, 2009.
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [27] H. Q. Le, W. J. Starke, J. S. Fields, D. Q. O’Connell, F. P. and Nguyen, B. J. Ronchetti, W. M. Sauer, E. M. Schwarz, and M. T. Väden. Ibm power6 microarchitecture. *IBM JRD*, 51(6), 2007.
- [28] Quoc Le, Marc’Aurelio Ranzato, Rajat Monga, Matthieu Devin, Kai Chen, Greg Corrado, Jeff Dean, and Andrew Ng. Building high-level features using large scale unsupervised learning. In *ICML*, 2012.
- [29] Yann LeCun and Yoshua Bengio. The handbook of brain theory and neural networks. chapter Convolutional Networks for Images, Speech, and Time Series. 1998.

- [30] Yann Lecun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324, 1998.
- [31] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [32] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- [33] Andrew Y. Ng. Feature selection, l1 vs. l2 regularization, and rotational invariance. In *ICML*, 2004.
- [34] V. Seshadri, O. Mutlu, M. A. Kozuch, and T. C. Mowry. The Evicted-Address Filter: A Unified Mechanism to Address Both Cache Pollution and Thrashing. In *PACT*, 2012.
- [35] Vivek Seshadri, Gennady Pekhimenko, Olatunji Ruwase, Onur Mutlu, Phillip B. Gibbons, Michael A. Kozuch, Todd C. Mowry, and Trishul Chilimbi. Page overlays: An enhanced virtual memory framework to enable fine-grained memory management. In *ISCA*, 2015.
- [36] S. Srinath, O. Mutlu, H. Kim, and Y. N. Patt. Feedback directed prefetching: Improving the performance and bandwidth efficiency of hardware prefetchers. In *HPCA*, 2007.
- [37] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 2014.
- [38] Luis Villa, Michael Zhang, and Krste Asanović. Dynamic zero compression for cache energy reduction. In *MICRO*, 2000.
- [39] Youtao Zhang, Jun Yang, and Rajiv Gupta. Frequent value locality and value-centric data cache design. In *ASPLOS*, 2000.