

# Reducing the Memory Cost of Training Convolutional Neural Networks by CPU Offloading

|   |   |
|---|---|
| TRISTAN HASCOET* <sup>†</sup><br>Kobe Univeristy<br>tristan@people.kobe-u.ac.jp | WEIHAO ZHUANG*<br>Kobe Univeristy<br>zhuangweihao1996@gmail.com |
| QUENTI FEBVRE<br>Sicara<br>quentin.febvre@gmail.com                             | YASUO ARIKI<br>Kobe Univeristy<br>ariki@kobe-u.ac.jp            |
| TETSUYA TAKIGUCHI<br>Kobe Univeristy<br>takigu@kobe-u.ac.jp                     |   |

## Abstract

In recent years, Convolutional Neural Networks (CNN) have enabled unprecedented progress on a wide range of computer vision tasks. However, training large CNNs is a resource intensive task that requires specialized Graphical Processing Units (GPU) and highly optimized implementations to get optimal performance from the hardware. GPU memory is a major bottleneck of the CNN training procedure, limiting the size of both inputs and model architectures. In this paper, we propose to alleviate this memory bottleneck by leveraging an under-utilized resource of modern systems: the device to host bandwidth. Our method, termed CPU offloading, works by transferring hidden activations to the CPU upon computation, in order to free GPU memory for upstream layer computations during the forward pass. These activations are then transferred back to the GPU as needed by the gradient computations of the backward pass. The key challenge to our method is to efficiently overlap data transfers and computations in order to minimize wall time overheads induced by the additional data transfers. On a typical work station with a Nvidia Titan X GPU, we show that our method compares favorably to gradient checkpointing as we are able reduce the memory consumption of training a VGG19 model by 35% with a minimal additional wall time overhead of 21%. Further experiments detail the impact of the different optimization tricks we propose. Our method is orthogonal to other techniques for memory reduction such as quantization and sparsification so that they can easily be combined for further optimizations.

---

\*Equal contribution

<sup>†</sup>Corresponding author

# 1 Introduction

Over the last few years, Convolutional Neural Networks (CNN) have enabled unprecedented progress on a wide array of computer vision tasks. One disadvantage of these approaches is their resource consumption: Training deep models within a reasonable amount of time requires special Graphical Processing Units (GPU) with numerous cores and large memory capacity. Given the practical importance of these models, a lot of research effort has been directed towards algorithmic and hardware innovations to improve their resource efficiency such as low-precision arithmetic [1], network pruning [2], or efficient stochastic optimization algorithms [3].

In this paper, we focus on a particular aspect of resource efficiency: optimizing the GPU memory cost of training CNNs. Given the ubiquity of CNN for practical computer vision applications, optimizing the memory consumption of CNN training has the potential to impact a wide range of applications. Here, we only present a few of the most interesting potential impacts of such optimization:

**Low-memory GPUs:** Training large CNN requires special GPUs with large memory capacity. Typical desktop GPUs memory capacity is too small for training large CNNs. As a result, getting into deep learning research comes with the barrier cost of either buying specialized hardware or renting live instances from cloud service providers, while standard laptop GPUs remain idle untapped resources. Reducing the memory cost of deep model training allows training deep nets on standard graphic cards without the need for specialized hardware, effectively removing this barrier cost.

**Research in optimization:** Recent works on stochastic optimization algorithms have highlighted the benefits of large batch training [4, 5, 6]. For example, in Imagenet, linear speed-ups in training have been observed with increasing batch sizes up to tens of thousands of samples [5]. Optimizing the memory cost of CNN training may allow further research on the optimization trade-offs of large batch training. Very large batch training on small datasets like MNIST and CIFAR10 is computationally inefficient with current stochastic optimization algorithms [5]. However, for such small datasets, memory optimization would allow to process the full dataset in one pass through the networks. The ability to process the full dataset in one pass allows to easily train CNNs on the exact gradient of the loss function. Hence, memory optimization techniques opens the door for research on gradient descent optimization of neural networks outside the realm of Stochastic Gradient Descent.

There is an inherent trade-off between the memory consumption and

computation wall time of the CNN training procedure: Existing approaches to optimize the memory consumption of CNN training, (gradient checkpointing, reversible network architectures), trade off memory consumption for additional computations by recomputing all or a subset of the hidden layers activations during the backward pass.

Instead, our approach reduces the GPU memory consumption without introducing any additional computation by leveraging an under-utilized resource: host-device communication. We propose to temporarily offload GPU memory buffers to the CPU during the forward pass of the computation, and transferring these memory buffers back into GPU memory as needed by the gradient computations during the backward of the backpropagation algorithm.

The key challenge in our approach is to efficiently overlap the GPU computations with the data transfers between CPU and GPU in order to minimize the overhead in wall time introduced by these data transfer. In this paper, we describe an efficient implementation of this approach that allows us to reduce by up to 35% the memory cost of training a VGG network with a minimal wall time overhead of 21%. We compare the memory vs. wall time trade off of our approach to gradient checkpointing to illustrate the efficiency of our approach.

The remainder of this paper is organized as follows: In Section 2, we briefly review the literature for related work. Section 3 introduces the preliminary notions necessary to understand the root of the GPU memory bottleneck. Section 4 presents our approach and details the different tricks needed for efficient implementation. Finally, Section 5 presents the results of our evaluation.

## 2 Related Work

Research into resource optimization of CNNs covers a wide array of techniques, most of which are orthogonal to our work. We briefly present some of these works:

On the architectural side, Squeezenet [7] was first proposed as an efficient neural architecture reducing the number of model parameters while maintaining high classification accuracy. MobileNet [8] uses depth-wise separable convolutions to further reduce the computational cost of inference for embedded device applications.

Network pruning [2] is a set of techniques developed to decrease the model weight size and computational complexity. Network pruning works

by removing the network weights that contribute the least to the model output. Pruning deep models has been shown to efficiently reduce the memory and computational cost of inference without significantly hurting model accuracy. Although pruning methods focus on the optimization of inference, the recently proposed lottery ticket hypothesis [9] has shown that specifically pruned networks could be trained from scratch to high accuracy. This may be an interesting and complementary line of work to investigate in the future to reduce training memory costs.

Low precision arithmetic has been proposed as a mean to reduce both memory consumption and computation time of deep learning models. Mixed precision training [10] combines float16 with float32 operations to avoid numerical instabilities due to either overflow or underflow. For inference, integer quantization [1, 11] has been shown to drastically improve the computation and memory efficiency and has been successfully deployed on both edge devices and data centers. Integrating mixed-precision training to our proposed architecture would allow us to further reduce training memory costs.

Most related to our work, gradient checkpointing was introduced as a mean to reduce the memory cost of deep neural network training. Gradient checkpointing, first introduced in [12], trades off memory for computational complexity by storing only a subset of the activations during the forward pass. During the backward pass, missing activations are recomputed from the stored activations as needed by the backpropagation algorithm. Follow-up work [13] has since built on the original gradient checkpointing algorithm to improve this memory/computation trade-off.

In contrast, our approach does not introduce any additional computation: Instead of computing a set of missing hidden activations during the backward pass, we propose to offload the hidden activations to the CPU during the forward pass, and to transfer these activations back to GPU memory during the backward pass.

Reversible models [14, 15] constrain the CNN architecture to feature invertible transformations. This allows the activation values of lower layers to be reconstructed from those of higher layers during the backward pass. Reversible networks have been shown to offer a better memory/computation trade-off than gradient checkpointing at the cost of constraining the CNN architecture.

Our approach combines revertible operations with CPU offloading: we use the invertible BN-Leaky ReLu block design proposed in [16] to efficiently deal with normalization and non-linearity layers, and only offload to CPU the activations of the pooling and convolution layers.

### 3 Preliminaries

Let us consider a model  $F$  of  $N$  sequential layers trained to minimize an error  $e$  defined by a loss function  $\mathcal{L}$  for an input  $x$  and ground-truth label  $\bar{y}$ :

$$F : x \rightarrow y \quad (1a)$$

$$y = f_N \circ \dots \circ f_2 \circ f_1(x) \quad (1b)$$

$$e = \mathcal{L}(y, \bar{y}) \quad (1c)$$

During the forward pass, each layer  $f_i$  takes as input the activations  $z_{i-1}$  from the previous layer and outputs activation features  $z_i = f_i(z_{i-1})$ , with  $z_0 = x$  and  $z_N = y$  being the input and output of the network respectively. During the backward pass, the gradient of the loss with respect to the hidden activations are propagated backward through the layers of the networks using the chain rule as:

$$\frac{\delta \mathcal{L}}{\delta z_{i-1}} = \frac{\delta \mathcal{L}}{\delta z_i} \times \frac{\delta z_i}{\delta z_{i-1}} \quad (2)$$

Before propagating the loss gradient with respect to its input to the previous layer, each parameterized layer computes the gradient of the loss with respect to its parameters. In vanilla SGD, for a given learning rate  $\eta$ , the weight gradients are subsequently used to update the weight values as:

$$\frac{\delta \mathcal{L}}{\delta \theta_i} = \frac{\delta \mathcal{L}}{\delta z_i} \times \frac{\delta z_i}{\delta \theta_i} \quad (3a)$$

$$\theta_i \leftarrow \theta_i - \eta \times \frac{\delta \mathcal{L}}{\delta \theta_i} \quad (3b)$$

For most layers, the computation of either gradients are functions of the layer's input activations  $z_{i-1}$ : For example, convolution layers need the values of input activations to compute the weight gradients:

$$\frac{\delta \mathcal{L}}{\delta \theta_i} = z_{i-1}^F \star \frac{\delta \mathcal{L}}{\delta z_i} \quad (4)$$

while Rectified Linear Unit layers need the input activations values to compute the gradients of the loss with respect to its inputs:

$$\frac{\delta \mathcal{L}}{\delta z_{i-1}^j} = \begin{cases} \frac{\delta \mathcal{L}}{\delta z_i^j}, & \text{if } z_{i-1}^j \geq 0 \\ 0 & \text{if } z_{i-1}^j < 0 \end{cases} \quad (5)$$

Hence, backpropagation implementations in deep learning frameworks store hidden layers activations in GPU memory upon computation during the forward pass. Activations accumulate in live memory buffers throughout the full forward pass until used for gradients computations of the backward pass. Once the gradients computed during the backward pass, their associated hidden activation buffers can be freed from live memory. However, the accumulation of activation values stored within each layer along the forward pass creates a major bottleneck in GPU memory. In the next section, we detail our approach to alleviate this memory bottleneck.

## 4 Propose Method

### 4.1 Framework

The input activations of each layer are kept in GPU memory only to be used for the computation of the layer weight gradients during the backward pass. Hence, the activations of lower layers are kept idle in GPU memory during the forward and backward computations through higher layers. We propose to offload these activations to the CPU during this idle time in order to free up some GPU memory space for the computation of higher layers activations.

Figure 1 illustrates our approach. During the forward pass (top), activations are computed forward through the network layers. Instead of keeping these activations idle in GPU memory, activation values are transferred to the CPU memory immediately after their computation. In the backward pass (bottom), gradients are backpropagated backward through the network layers following equations (3). Our implementation synchronizes the transfer of the layers input activations back to GPU right before they are needed for their layer’s gradient computation. Hence the key challenge in our implementation consists in synchronizing the data transfers with the computations so that only the minimal amount of activation values is loaded in GPU memory at any given time, while the least amount of time is spent waiting for the data transfer.

To achieve this goal, we propose a set of optimization tricks organized along two axes: The first consists in optimizing the data transfer speed between CPU and GPU memory, using efficient memory accesses and data compression schemes. The second consists in efficient parallelization to maximally overlap the computations with the data transfer. The following subsections details optimizations along these two axes.

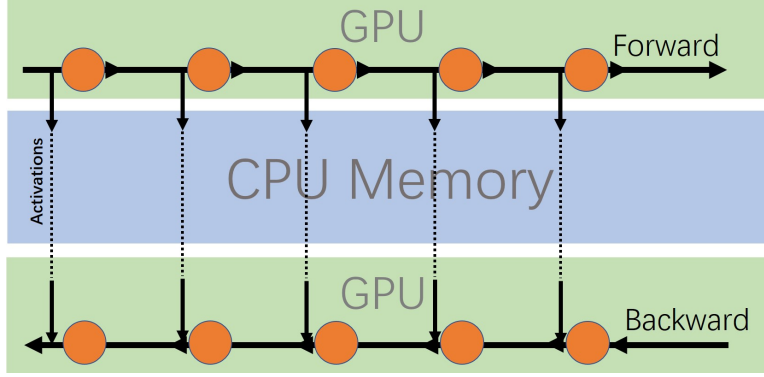


Figure 1: Illustration of CPU Offloading. During the forward pass (top), hidden activation buffers are transferred to CPU memory upon computation. During the backward pass (bottom), hidden activation buffers are transferred back to GPU memory just in time to compute the weight gradients.

## 4.2 Parallelization

Figure 2 illustrates the execution through time of a forward and backward pass through a toy network with and without parallelization of the data transfer. Without parallelization, computation and data transfers are performed sequentially so that the total wall time is given by the sum of the computation and data transfer time.  $\mathcal{T}_{total} = \mathcal{T}_{comp} + \mathcal{T}_{data}$ . Parallelization aims to overlap the computation and data transfer so that the total wall time is given by  $\mathcal{T}_{total} = \mathcal{T}_{comp} + \mathcal{T}_{idle}$ , where  $\mathcal{T}_{idle}$  represents synchronization delays in cases where the computation is stopped to await for the required data transfer to complete.

The key challenge in this parallelization scheme consists in efficiently managing memory allocation and data transfer so as to minimize the time  $\mathcal{T}_{idle}$  spent awaiting for data transfer. In this paper, we adopt a simple parallelization strategy: During the forward pass, activations  $z_i$  are transferred to CPU upon computation by layer  $i$  as  $z_i = f_i(z_{i-1})$ . The GPU memory buffer of  $z_i$  must then await for both the data transfer to CPU and the next layer computation  $z_{i+1} = f_{i+1}(z_i)$  to complete before de-allocation in order to avoid unfortunate overwriting by concurrent operations.

One important exception to this rule concerns skip connections, as illustrated in Figure 3. Skip connections induce a delay in the GPU buffer de-allocation as the input of residual blocks must be kept in memory until

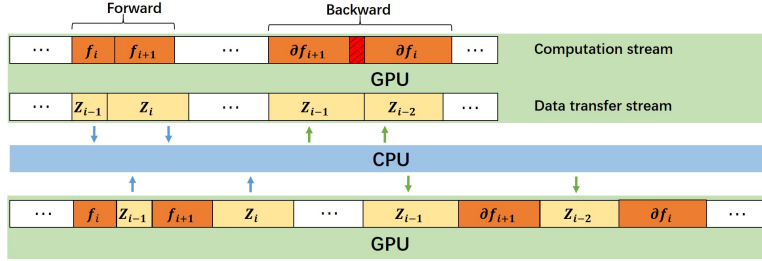


Figure 2: Illustration of the parallel CPU Offloading execution through time. (bottom): Without parallelism, computations and data transfers are sequentially executed within the same stream so that the overhead in wall time corresponds to the total time of data transfers  $\mathcal{T}_{data}$ . (top): Computations and data transfer are executed in parallel in their dedicated stream. The overhead in wall time  $\mathcal{T}_{idle}$  is spent as the computation stream awaits for data transfers to complete (illustrated in red).

the end of the residual block computation to be added to the output. In ResNet architectures, this delay is short enough to have little impact on the memory/wall time trade off. However, this means that our method is poorly suited to densely connected architectures such as DenseNet [17] or UNet [18] as their long-range skip connections introduce large delays in the GPU buffer de-allocation.

During the backward pass, the input activation  $z_{i-1}$  of layer  $i$  must be transferred back into GPU memory for the backward gradient computations  $\delta f_i$  to proceed. Hence, we overlap the transfer of  $z_{i-1}$  to GPU with the backward computations  $\delta f_{i+1}$  of the upper layer  $i + 1$  to avoid GPU idle time. As formalized in Algorithm 2, we thus synchronize the data transfer of  $z_{i-1}$  with the beginning of the backpropagation through  $f_{i+1}$ .

We use threading and locks to handle the parallelism and synchronization on the CPU side, and CUDA streams to handle the parallelization on the GPU device. CUDA events provides synchronization primitives to precisely track the completion of kernel executions. Algorithm 1 and 2 provide pseudo-code for the forward and backward pass respectively.

With the above parallelization scheme, the overhead  $\mathcal{T}_{idle}$  in wall time is given by the sum of the difference in the computation and data transfer time at each layer. If data transfer is faster than the computation in each layer, CPU offloading will reduce memory consumption without any overhead in wall time. Hence, CPU offloading is best suited for computation intensive operations like convolutions. In contrast, batch normalization layers come



---

**Algorithm 1:** Forward procedure through layer  $i$  with parallel CPU offloading.

---

**Data:** Layer  $f_i$ , input activation  $z_{i-1}$   
CPU pinned memory buffer  $P_{i-1}$   
CPU thread  $T_{data}$   
CUDA events  $E_{data}^i, E_{comp}^i$   
CUDA Streams  $S_{data}, S_{comp}$

**Result:**  $z_i$

$Allocate(z_i);$   
 $S_{comp} \Leftarrow z_i \Leftarrow f_i(z_{i-1});$   
 $S_{comp} \Leftarrow E_{comp}^i;$   
**In Thread  $T_{data}$ :**  
 $S_{data} \Leftarrow P_{i-1} \Leftarrow z_{i-1};$   
 $S_{data} \Leftarrow E_{data}^i;$   
 $Wait(E_{data}^i, E_{comp}^i);$   
 $Free(z_{i-1});$

---

Algorithm 1: Double arrows indicate the asynchronous execution of a CUDA directive within a stream. Data transfers are executed within dedicated CUDA stream and CPU thread to synchronize the memory de-allocation without blocking the execution of upward layers

with negligible computational cost, but significant data transfer time so that offloading the input activations of such layer comes with a significant overhead. Instead of offloading these layers, we reconstruct their input activations from their output using their inverse operation, as proposed in [16], and only offload the input activations of the convolution and pooling layers to CPU.

### 4.3 Data Transfer Optimization

The time needed to transfer data between devices is given by the product of the transfer speed bandwidth  $\mathcal{B}$ , in bytes per second, by the volume  $\mathcal{V}$  of the data buffer in bytes:

$$\mathcal{T}_{data} = \frac{\mathcal{V}}{\mathcal{B}} \quad (6)$$

Hence, minimizing data transfer time can be achieved by either maximizing the transfer speed  $\mathcal{B}$  or minimizing the activation volume  $\mathcal{V}$  through

---

**Algorithm 2:** Backward procedure through layer  $i$  with parallel CPU offloading.

---

**Data:** Layer  $f_i$ , output gradients  $\frac{\delta \mathcal{L}}{\delta z_i}$   
CPU pinned memory buffer  $P_{i-1}$   
CPU thread  $T_{comp}$   
CUDA events  $E_{data}^i, E_{data}^{i+1}$   
 $, E_{comp}^i$   
CUDA Streams  $S_{data}, S_{comp}$

**Result:**  $\frac{\delta \mathcal{L}}{\delta z_{i-1}}, \frac{\delta \mathcal{L}}{\delta \theta_i}$

$Allocate(z_{i-1});$   
 $S_{data} \leftarrow z_{i-1} \leftarrow P_{i-1};$   
 $S_{data} \leftarrow E_{data}^i;$   
 $Wait(E_{data}^{i+1});$   
 $Allocate(\frac{\delta \mathcal{L}}{\delta z_{i-1}}, \frac{\delta \mathcal{L}}{\delta \theta_i});$   
 $S_{comp} \leftarrow \frac{\delta \mathcal{L}}{\delta z_{i-1}} \leftarrow \frac{\delta \mathcal{L}}{\delta z_i} \times \frac{\delta z_i}{\delta z_{i-1}};$   
 $S_{comp} \leftarrow \frac{\delta \mathcal{L}}{\delta \theta_i} \leftarrow \frac{\delta \mathcal{L}}{\delta z_i} \times \frac{\delta z_i}{\delta \theta_i};$

---

data compression. In this paper, we focus on maximizing data transfer speed using pinned CPU memory buffers.

CPU memory can be accessed through either pageable memory or pinned memory address spaces. Pinned memory allows for Direct Memory Access (DMA), which can significantly speedup data transfer between CPU and GPU devices through PCIe lanes. However, allocating pinned memory buffers is a time consuming operation, as illustrated in Table 1. Hence, naively allocating pinned memory buffers within each forward pass would slow down training due to the overhead of pinned memory allocation. Instead, we propose to allocate the pinned memory buffers beforehand, during the initialization of the CNN. Our implementation scans the network architecture upon instantiation and allocates dedicated pinned memory buffers for each target layer activations (convolutions and max-pooling). During the forward pass, we transfer the GPU data directly into the pre-allocated pinned memory buffer, and, during the backward pass, we transfer the activations back to GPU from these pinned buffers. Table 1 quantifies the gains in transfer speed brought by the use of pinned memory buffers and, in Section 5, we evaluate the impact of using pinned memory buffers on the overall efficiency of our method.

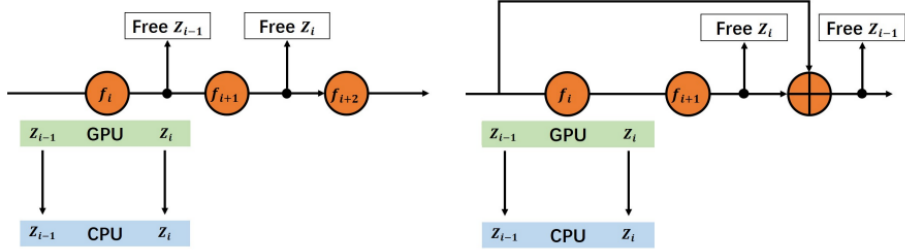


Figure 3: Illustration of the delay in memory de-allocation induced by residual connections. (left) without residual connection, the input  $z_{i-1}$  of layer  $f_i$  can be freed from GPU memory as  $f_i$  computation terminates. (right) Residual connections induce a delay in memory de-allocation as  $z_{i-1}$  must be kept in GPU memory to be added to the output of the residual block.

| Operation                        | Speed (GB/s) |
|----------------------------------|--------------|
| Pinned Memory Allocation         | 1.95         |
| GPU $\rightarrow$ CPU (Pageable) | 1.75         |
| GPU $\rightarrow$ CPU (Pinned)   | 6.5          |
| CPU $\rightarrow$ GPU (Pageable) | 4.1          |
| CPU $\rightarrow$ GPU (Pinned)   | 6.1          |

Table 1: Measurement of data transfer and pinned memory allocation speed on our workstation. Pinned memory buffers bring a substantial speed up in both directions of the data transfer.

Another way to reduce the time of data transfer between CPU and GPU is to minimize  $\mathcal{V}$  by compressing the hidden layer activations before transferring them between devices. Compression can be achieved by diverse means including sparsification, quantization or low rank factorization techniques. Low-rank factorization would induce non-negligible additional computational costs, which would slow down the computation. Either sparsifying or quantizing the layers activations are interesting solutions. However, compression would incur minor losses of precision in the activation values, which might affect the training dynamics. Although previous works suggest that such minor losses in precision does not negatively impact training, investigation of compression techniques would require a detailed analysis of the impact of compression on accuracy, which is out of the scope of this paper. Hence we leave this question open for future work as an interesting

direction for further optimization.

## 5 Experiments

We conduct our experiments on a workstation with an NVIDIA Titan X GPU connected to the host through PCIe1 16x lanes, which represents a peak bandwidth of 6.5 GB/s. We believe this set-up to be representative of a typical workstation for individual practitioners, if not slightly unfavorable, as our device to host bandwidth is rather limited. Unless specified otherwise, we conduct our experiments on a batch of 32 RGB images with a resolution of 128 by 128 pixels. For each experiment, we are interested in the total memory usage and Wall time of the computation for a full training iteration including the forward and backward pass through the network. It is important to note that our method does not modify the actual computations performed by the network so that the training dynamics and accuracy are unaffected by CPU Offloading. Hence, we focus on reporting the memory vs. wall time trade-off, as the accuracy reached by the model is strictly unaffected.

### 5.1 Memory vs. Wall Time Trade-Off

Our implementation allows to control for the memory foot-print vs. computation wall time trade-off by either using more aggressive parallelization schemes, or by offloading different subsets of hidden layers input activations. In this section, we use the parallelization scheme described in Section 4.2 and explore the memory vs. time trade-off by off-loading different subsets of the model’s layers.

Figure 4 compares the memory/Wall time trade-off of CPU off-loading to the trade-off provided by gradient checkpointing, as implemented in the Pytorch [19] library on the VGG19 architecture. In gradient checkpointing, this trade-off is defined by the number of checkpointed layers through the full architecture.

Figure 4 shows the pareto frontier for an extensive search in both checkpointing and CPU offloading layer subset. CPU offloading leads to a better trade-off between memory consumption and wall time. This observed performance is due to the fact that the data transfer between host and device can be efficiently parallelized with the computation, hence incurring a minor overhead, while gradient checkpointing can not parallelize the gradient computations and activation reconstructions. Hence, the additional cost of

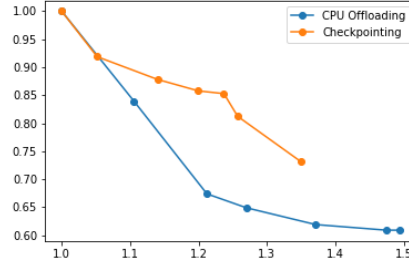


Figure 4: Comparison of the memory vs. Wall time trade-off achieved by CPU Offloading and sequential gradient checkpointing.

sequentially reconstructing the hidden activations during the backward pass induces a higher overhead.

## 5.2 Ablation study

Table 2 illustrates the speed up brought by both parallelization and pinned memory buffers. We record the timings of a training iteration needed to achieve maximal memory reduction on the VGG19 architecture for an input batch of 32 images with 218 by 128 pixel resolution. We start by recording the time for a baseline CPU offloading approach without parallelization nor data transfer optimization. We then successively integrate the pinned memory buffers and parallelization to this baseline and observe the improvement.

| Operation       | Time (ms) | Rel. Improvement |
|-----------------|-----------|------------------|
| Baseline        | 560       | 0%               |
| Pinned Memory   | 355       | −36.6%           |
| Parallelization | 236       | −21.2%           |

Table 2: Speed-ups achieved by the pinned memory buffers and parallelization techniques

Both parallelization and pinned memory buffers bring substantial improvements in speed and have proven critical for performance. Compression of the hidden states to either sparse representations or lower precision encodings may bring similar speed-up and would be very interesting to investigate in future work. We expect that simply converting the hidden activations to

FP16 before sending them to the host device should further reduce the data transfer time by half.

### 5.3 Architecture comparison

In this section, we compare the memory gains and time overhead achieved through CPU offloading on different popular architectures.

| Architecture | Wall Time Overhead | Memory Reduction |
|--------------|--------------------|------------------|
| VGG16        | 19.5%              | 32.5%            |
| VGG19        | 21.0%              | 36.9%            |
| Resnet18     | 17.6%              | 54.4%            |
| Resnet50     | 15.3%              | 56.8%            |
| MobileNet    | 63.2%              | 67.8%            |

Table 3: Memory vs.Time trade-off achieved by different architectures. The Wall Time Overhead column shows the additional time needed to achieve the given memory reduction (the lower the better). Memory Reduction represents the reduction of peak memory usage achieved with CPU offloading (the higher the better).

Both ResNet architectures and MobileNet [8] feature batch normalization layers, for which we used the invertible operation of [16]. Invertible batch normalization layers further reduce the memory consumption with little time overhead, which explains the slightly better trade-off of ResNet architectures. On the other hand, the MobileNet architecture uses depth-wise convolutions instead of the regular convolution layers used in both VGG and ResNet architectures. Depth-wise convolutions have lower arithmetic intensity [20] than convolution layers, i.e.; for input activations of similar size, they perform less, hence faster, operations. As computations are performed faster while the data transfer time remains constant, CPU offloading of depth-wise convolutions incur a more important overhead. This result highlights the fact that CPU offloading works best for computationally intensive layers.

### 5.4 Discussion

The memory/time trade-off achievable by CPU offloading is a function of the ratio between computation and data transfer speed. This ratio varies with

the specific hardware being run on, the algorithm used for the convolutions, as well the efficiency of their implementation.

On the software side, the experiments presented in this paper were run with state-of-the art cudnn implementations of the winograd convolution algorithm through the PyTorch framework [19]. On the hardware side, we ran experiments on a single Nvidia Titan X GPU connected to the host via a 6GB/s bandwidth connection, which we believe to be representative of a typical deep learning student/researcher’s workstation.

The exact trade-offs achieved will differ depending on the specific hardware and software used in different systems, as well as the parallelization strategy used for CPU offloading. However, the fact that data transfers between CPU and GPU can be efficiently parallelized with the GPU computations means that memory gains can be achieved with minimal time overhead on any system given an appropriate parallelization strategy.

Searching for an optimal parallelization strategy given a fixed computational graph and hardware system is a very interesting problem left open to investigate. In this paper, we provide a proof of concept by showing that even a rather naive parallelization strategies with minimum optimization leads to a competitive memory/time trade-off relatively to the well established method of gradient checkpointing.

## 6 Conclusion

Convolutional Neural Networks have become the backbone of computer vision systems. Despite their great success, one major drawback of these models is their intense resource consumption: Training CNNs needs highly optimized implementations leveraging all possible hardware resources. In this paper, we proposed to leverage GPU-CPU communication, an under-utilized resource in typical CNN training pipelines, to alleviate the GPU bottleneck in training CNNs. Our experiment on a standard single-GPU work station shows a more effective memory/wall time trade-off than gradient checkpointing. Our approach is orthogonal to other resource optimization approaches such as pruning and quantization so these approaches can easily be combined in future work.

## References

- [1] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko.

- Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2704–2713, 2018.
- [2] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *arXiv preprint arXiv:1611.06440*, 2016.
  - [3] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
  - [4] Christopher J Shallue, Jaehoon Lee, Joe Antognini, Jascha Sohl-Dickstein, Roy Frostig, and George E Dahl. Measuring the effects of data parallelism on neural network training. *arXiv preprint arXiv:1811.03600*, 2018.
  - [5] Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. An empirical model of large-batch training. *arXiv preprint arXiv:1812.06162*, 2018.
  - [6] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.
  - [7] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
  - [8] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
  - [9] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
  - [10] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.



- [11] Shuang Wu, Guoqi Li, Feng Chen, and Luping Shi. Training and inference with integers in deep neural networks. *arXiv preprint arXiv:1802.04680*, 2018.
- [12] James Martens and Ilya Sutskever. Training deep and recurrent networks with hessian-free optimization. In *Neural networks: Tricks of the trade*, pages 479–535. Springer, 2012.
- [13] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- [14] Aidan N Gomez, Mengye Ren, Raquel Urtasun, and Roger B Grosse. The reversible residual network: Backpropagation without storing activations. In *Advances in Neural Information Processing Systems*, pages 2214–2224, 2017.
- [15] Jörn-Henrik Jacobsen, Arnold Smeulders, and Edouard Oyallon. i-revnet: Deep invertible networks. *arXiv preprint arXiv:1802.07088*, 2018.
- [16] Samuel Rota Bulò, Lorenzo Porzi, and Peter Kotschieder. In-place activated batchnorm for memory-optimized training of dnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5639–5647, 2018.
- [17] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [18] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [19] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- [20] Bichen Wu, Alvin Wan, Xiangyu Yue, Peter Jin, Sicheng Zhao, Noah Golmant, Amir Gholaminejad, Joseph Gonzalez, and Kurt Keutzer.

Shift: A zero flop, zero parameter alternative to spatial convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9127–9135, 2018.