





BPPSA: SCALING BACK-PROPAGATION BY PARALLEL SCAN ALGORITHM

Shang Wang 1 Yifan Bai 2 Gennady Pekhimenko 1

ABSTRACT

In an era when the performance of a single compute device plateaus, software must be designed to scale on massively parallel systems for better runtime performance. However, in the context of training deep learning models, the popular back-propagation (BP) algorithm imposes a *strong sequential dependency* in the process of gradient computation. Under model parallelism, BP takes $\Theta(n)$ steps to complete which hinders its scalability on parallel systems (n represents the number of compute devices into which a model is partitioned).

In this work, in order to improve the scalability of BP, we reformulate BP into a scan operation which is a primitive that performs an in-order aggregation on a sequence of values and returns the partial result at each step. We can then scale such reformulation of BP on parallel systems by our modified version of the Blelloch scan algorithm which theoretically takes $\Theta(\log n)$ steps. We evaluate our approach on a vanilla Recurrent Neural Network (RNN) training with synthetic datasets and a RNN with Gated Recurrent Units (GRU) training with the IRMAS dataset, and demonstrate up to $2.75\times$ speedup on the overall training time and $108\times$ speedup on the backward pass. We also demonstrate that the retraining of pruned networks can be a practical use case of our method.

1 Introduction

The training of deep learning models demands more and more compute resources (Amodei et al., 2018) as the models become more powerful and complex with an increasing number of layers in recent years (Krizhevsky et al., 2012; Szegedy et al., 2015; Simonyan & Zisserman, 2015; He et al., 2015; Huang et al., 2016). For example, ResNet can have more than a thousand layers (He et al., 2016), and ResNet-152 takes days to train on eight state-of-the-art GPUs (Coleman et al., 2017). Now that the performance of a single compute device plateaus (Esmaeilzadeh et al., 2011; Arunkumar et al., 2017), training has to be designed to scale on massively parallel systems.

Data parallelism (Shallue et al., 2018) is the most popular way to scale training by partitioning the training data among multiple devices, where each device contains a full replica of the model. As the number of devices increases, data parallelism faces the trade-off between synchronization cost in synchronous parameter updates and staleness in asynchronous parameter updates (Ben-Nun & Hoefler, 2018). Furthermore, recent studies demonstrate the scaling limit of data parallelism even when assuming perfect implementations and zero synchronization cost (Shallue et al., 2018).

Proceedings of the 3^{rd} MLSys Conference, Austin, TX, USA, 2020. Copyright 2020 by the author(s).

Lastly, data parallelism cannot be applied when a model does not fit into one device due to memory constraints (e.g., caused by deep network architecture, large batch size, or high input resolution (Rhu et al., 2016; Zhu et al., 2018)).

Model parallelism (Krizhevsky, 2014; Huang et al., 2018; Shazeer et al., 2018; Narayanan et al., 2019) is another approach to distributed training which partitions a model and distributes its parts among devices. It covers a wide spectrum of training deep learning models where data parallelism does not suffice. Naïve training under model parallelism does not scale well with the number of devices due to under-utilization of the hardware resources, since at most one device can be utilized at any given point in time (Narayanan et al., 2019). To address the aforementioned issue, prior works on pipeline parallelism, including PipeDream (Narayanan et al., 2019) and GPipe (Huang et al., 2018), propose pipelining across devices for better resource utilization; however, as the number of layers and devices increases, pipeline parallelism still faces the trade-off between resource utilization in synchronous parameter updates and staleness in asynchronous parameter updates (Narayanan et al., 2019). Moreover, to fully fill the pipeline with useful computation, each device needs to store the activations at the partition boundaries for all mini-batches that enter the pipeline. Therefore, the maximum number of devices that pipeline parallelism can support is limited by the memory capacity of a single device.

Algorithmically, the fundamental reason for this scalability limitation observed from prior works is that the backpropagation (BP) algorithm (Rumelhart et al., 1988) im-

¹Department of Computer Science, University of Toronto ²Department of Electrical Engineering and Computer Sciences, University of California, Berkeley; work done at the University of Toronto. Correspondence to: Shang Wang <wangsh46@cs.toronto.edu>.

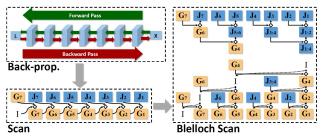


Figure 1: BP as a *scan* operation, scaled by our modified version of the *Blelloch scan* algorithm.

poses a *strong sequential dependency* between layers during the gradient computation. Since computing systems evolve to have more and more parallel nodes (Esmaeilzadeh et al., 2011; Arunkumar et al., 2017), in this work, we aim at exploring the following question: *How can BP scale efficiently when the number of layers and devices keeps increasing into the foreseeable future?*

To answer this question, we utilize a primitive operation called *Scan* (Blelloch, 1990) that performs an in-order aggregation on a sequence of values and returns the partial result at each step. Parallel algorithms (Hillis & Steele, 1986; Blelloch, 1990) have been developed to scale the scan operation on massively parallel systems. We observe that BP is mathematically similar to a scan operation on the transposed Jacobian matrix (Weisstein, 2019) of each layer and the gradient vector of the output from the last layer. Inspired by this key observation, we *restructure* the strong sequential dependency of BP, and present a new method to scale Back-propagation by Parallel Scan Algorithm (BPPSA). Our major contributions are summarized below.

- We reformulate BP as a scan operation and modify the $Blelloch \, scan$ algorithm (Blelloch, 1990) to efficiently scale BP in a parallel computing environment. Our method has a theoretical step complexity of $\Theta(\log n)$, where n represents the number of devices into which a model is partitioned, compared to $\Theta(n)$ of the naïve implementation of model parallelism. Moreover, our algorithm does not have the theoretical scalability limit by the memory capacity of a single device as pipeline parallelism does. As an example, Figure 1 shows how BP for training a network composed of 7 layers (blue cubes) can be reformulated into a scan operation on the transposed Jacobian matrices (blue squares) of this network and the final gradient vector (yellow squares), as well as how this scan operation can be scaled by BPPSA.
- Generating, storing and processing full Jacobian matrices are usually considered to be prohibitively expensive. However, we observe that the Jacobians of many types of layers (e.g., convolution, activation, max-pooling) can be extremely sparse where we can leverage sparse matrix format

(Saad, 2003) to reduce the runtime and storage costs; more importantly, the positions of input-independent zeros in this case are deterministic, which leads to potentially more optimized implementations of sparse matrix libraries. Based on these observations, we develop routines to efficiently generate sparse transposed Jacobians for various operators.

• As a proof of concept, we evaluate BPPSA on a vanilla Recurrent Neural Network (RNN) (Elman, 1990) training with synthetic datasets, as well as a RNN with Gated Recurrent Units (GRU) (Cho et al., 2014) training with the IRMAS dataset (Bosch et al., 2012). Our method achieves a maximum 2.75× speedup in terms of the overall (end-to-end) training time, and up to $108\times$ speedup on the backward pass, compared to the baseline BP approach which under-utilizes the GPU. Moreover, we demonstrate that the retraining of pruned networks (Han et al., 2015; See et al., 2016; He et al., 2017) (e.g., pruned VGG-11 (Simonyan & Zisserman, 2015)) can also be a practical use case of BPPSA.

2 BACKGROUND AND MOTIVATION

2.1 Problem Formulation

We conceptualize a deep learning model as a vector function f composed of sub-functions $\vec{x}_i = f_i(\vec{x}_{i-1}; \vec{\theta}_i)$:

$$f(.; \vec{\theta}_1, ..., \vec{\theta}_n) = f_1(.; \vec{\theta}_1) \circ ... \circ f_n(.; \vec{\theta}_n)$$
 (1)

where $\vec{\theta}_i, i \in \{1,...,n\}$ are the parameters of the model. The model is evaluated by an objective function $l(f(\vec{x}_0\,;\,\vec{\theta}_i,i\in\{1,...,n\}))$, where \vec{x}_0 is the initial input to the model. Figure 2 visualizes a convolutional neural network conceptualized in this formulation.

To train the model f, a first-order optimizer requires the gradients $\nabla_{\vec{\theta_i}} l$, which are derived from the gradients $\nabla_{\vec{x_i}} l$:

$$[\nabla_{\vec{\theta}_1} l, ..., \nabla_{\vec{\theta}_n} l] \leftarrow [(\frac{\partial \vec{x}_1}{\partial \vec{\theta}_1})^T \nabla_{\vec{x}_1} l, ..., (\frac{\partial \vec{x}_n}{\partial \vec{\theta}_n})^T \nabla_{\vec{x}_n} l] \quad (2)$$

where $\frac{\partial \vec{x}_i}{\partial \vec{\theta}_i}$ is the Jacobian matrix of the output \vec{x}_i of f_i to its parameters $\vec{\theta}_i$. To derive $\nabla_{\vec{x}_i} l$ given $\nabla_{\vec{x}_n} l$, BP (Rumelhart et al., 1988) solves the following recursive equation, from i=n-1 to i=1:

$$\nabla_{\vec{x}_i} l \leftarrow (\frac{\partial \vec{x}_{i+1}}{\partial \vec{x}_i})^T \nabla_{\vec{x}_{i+1}} l, \forall i \in \{n-1, ..., 1\}$$
 (3)

where $\frac{\partial \vec{x}_{i+1}}{\partial \vec{x}_i}$ is the Jacobian matrix of the output \vec{x}_{i+1} of f_{i+1} to its input \vec{x}_i . Equation 2 itself does not have dependency along i; therefore, the computation of $\nabla_{\vec{\theta}_i} l$ can be parallelized if $\nabla_{\vec{x}_i} l$ are available. However, Equation 3 imposes a strong sequential dependency along i where the computation of $\nabla_{\vec{x}_i+1} l$ finishes, and therefore, hinders the scalability when multiple workers (defined as instances of execution; e.g., a core in a multi-core CPU) are available.

¹Step complexity (detailed in Section 3.6) quantifies the runtime of a parallel algorithm.

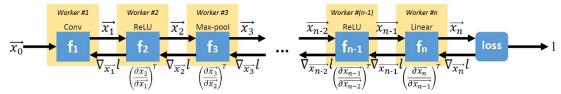


Figure 2: A visualization of the formulation in Section 2.1 on convolutional neural networks. Different parts of the model can be distributed to different devices (workers).

2.2 Prior Works

To increase the utilization of hardware resources in model parallelism, prior works, e.g., PipeDream (Narayanan et al., 2019) and GPipe (Huang et al., 2018), propose to pipeline the computation in the forward and backward passes across devices. However, these solutions are not "silver bullets" to scalability due to the following reasons.

First, both PipeDream (Narayanan et al., 2019) and GPipe (Huang et al., 2018) require storing activations and/or multiple versions of weights for all batches that enter the pipeline. Although GPipe's re-materialization (Chen et al., 2016) can mitigate memory usage, the theoretical per-device space complexity grows linearly with the length of the pipeline (i.e., the number of devices).² Thus, the maximum number of devices that pipeline parallelism can support is limited by the memory capacity of a single device (e.g., the GPU global memory), and such memory capacity is not expected to grow significantly in the foreseeable future (Mutlu, 2013).

Second, if the parameter updates are partially asynchronous as proposed in PipeDream (Narayanan et al., 2019), the resulting staleness may effect the convergence for adaptive optimizers such as Adam (Kingma & Ba, 2015) (Appendix C). If the gradient updates are fully synchronized as proposed in GPipe (Huang et al., 2018), the "bubble of idleness" between the forward and backward passes increases linearly with the length of the pipeline, which linearly reduces the hardware utilization and leads to diminishing returns.

Our approach fundamentally differs from these key prior works (Narayanan et al., 2019; Huang et al., 2018) in the following ways. First, instead of following the dependency of BP, we reformulate BP so that scaling is achieved via the Blelloch scan algorithm (Blelloch, 1990) which is designed for parallelism. Second, the original BP is reconstructed exactly without introducing new sources of errors (e.g., staleness); therefore, our method is agnostic to the exact first-order optimizer being used. Third, our approach becomes more advantageous as the number of devices increases, instead of diminishing returns or hitting scalability limits due to linear per-device space complexity.

2.3 Definition of the Scan Operation

For a *binary* and *associative* operator \oplus with an identity value I, the *exclusive scan* (a.k.a., *prescan*) on an input array $[a_0, a_1, a_2, ..., a_{n-1}]$ produces an output array $[I, a_0, a_0 \oplus a_1, a_0 \oplus a_1 \oplus a_2, ..., a_0 \oplus ... \oplus a_{n-2}]$ (Blelloch, 1990). Parallel scan algorithms have been developed due to the importance of the scan operation and the need to leverage the computing power of emerging parallel hardware systems (Hillis & Steele, 1986; Blelloch, 1990).

3 PROPOSED METHOD: BPPSA

3.1 Back-propagation as a Scan Operation

We define a binary, associative, and non-commutative operator $A \diamond B = BA$, whose identity value is the identity matrix I, where A can be either a matrix or a vector and B is a matrix. Using operator \diamond , we can reformulate Equation 3 as calculation of the following array:

$$[\nabla_{\vec{x}_n} l, \nabla_{\vec{x}_n} l \diamond (\frac{\partial \vec{x}_n}{\partial \vec{x}_{n-1}})^T, \nabla_{\vec{x}_n} l \diamond (\frac{\partial \vec{x}_n}{\partial \vec{x}_{n-1}})^T \diamond (\frac{\partial \vec{x}_{n-1}}{\partial \vec{x}_{n-2}})^T, \\ ..., \nabla_{\vec{x}_n} l \diamond (\frac{\partial \vec{x}_n}{\partial \vec{x}_{n-1}})^T \diamond ... \diamond (\frac{\partial \vec{x}_2}{\partial \vec{x}_1})^T]$$

$$(4)$$

Equation 4 can be interpreted as an *exclusive scan* operation of \diamond on the following input array:

$$\left[\nabla_{\vec{x}_n}l, \left(\frac{\partial \vec{x}_n}{\partial \vec{x}_{n-1}}\right)^T, \left(\frac{\partial \vec{x}_{n-1}}{\partial \vec{x}_{n-2}}\right)^T, ..., \left(\frac{\partial \vec{x}_2}{\partial \vec{x}_1}\right)^T, \left(\frac{\partial \vec{x}_1}{\partial \vec{x}_0}\right)^T\right]$$
(5)

3.2 Scaling Back-propagation with the Blelloch Scan Algorithm

We parallelize the computation of Equation 4 on multiple workers with the Blelloch scan algorithm (Blelloch, 1990), formally described in Algorithm 1. The algorithm contains two phases: *up-sweep* and *down-sweep*. As an example, Figure 3 visualizes this algorithm applied on the convolutional layers of VGG-11 (Simonyan & Zisserman, 2015) with levels L0-L4 as the up-sweep and levels L5-L10 as the down-sweep. Only the up-sweep phase contains matrix-matrix multiplications. Due to the *non-commutative* property of the operator ⋄, we have to reverse the order of operands for ⋄ during the down-sweep phase. This modification is reflected on line 13 of Algorithm 1 and visualized in Figure 4b.

3.3 Jacobian Matrices in Sparse Format

A full Jacobian matrix $\frac{\partial \vec{x}_{i+1}}{\partial \vec{x}_i}$ of $f_{i+1}(.; \vec{\theta}_{i+1})$ can be too expensive to generate, store, and process. In fact, the Jacobian of the first convolution operator in VGG-11 (Simonyan

²Appendix B includes a detailed space complexity analysis.

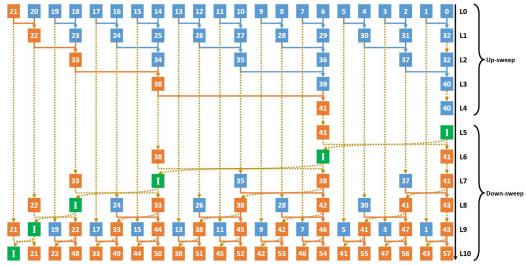


Figure 3: Applying our algorithm on the convolutional layers of VGG-11 (Simonyan & Zisserman, 2015). Blue, orange, and green squares represent transposed Jacobian matrices, gradient vectors, and *symbolic* identity matrices respectively. Blue solid lines, orange solid lines, and yellow dash lines represent matrix-matrix multiplications, matrix-vector multiplications, and *logical* data movements (that do not always have to be performed explicitly) respectively.

Algorithm 1 Modified Blelloch Scan Algorithm

```
Input: a = [\nabla_{\vec{x}_n} l, (\frac{\partial \vec{x}_n}{\partial \vec{x}_{n-1}})^T, ..., (\frac{\partial \vec{x}_1}{\partial \vec{x}_0})^T]
                                                                                  ⊳ Input array of Equation 5
Output: a = [I, \nabla_{\vec{x}_n} l, ..., \nabla_{\vec{x}_1} l] \quad \triangleright \nabla_{\vec{x}_i} l for Equation 2; computed in-place
                                                                                               ▷ Up-sweep Phase
 1: for d \leftarrow 0 to \lceil \log(n+1) \rceil - 2 do
           for all i \leftarrow 0 to (n-2^d) by 2^{d+1} do in parallel (l,r) \leftarrow (i+2^d-1, \min(i+2^{d+1}-1,n))
 3:
 4:
                 a[r] \leftarrow a[l] \diamond a[r]
 5.
           end for
 6: end for
 7: a[n] \leftarrow I
 8: for d \leftarrow \lceil \log(n+1) \rceil - 1 to 0 do
                                                                                           Down-sweep Phase
           for all i \leftarrow 0 to (n-2^d) by 2^{d+1} do in parallel
                 (l,r) \leftarrow (i+2^d-1, \min(i+2^{d+1}-1, n))
10:
                 T \leftarrow a[l]
11.
                 a[l] \leftarrow a[r]
12:
                 a[r] \leftarrow a[r] \diamond T

ightharpoonup Modification: Reverse the operands of \diamond.
13:
14:
           end for
15: end for
```

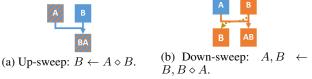


Figure 4: Visualizations of the primitive operations performed in the up-sweep and the down-sweep phases.

& Zisserman, 2015) processing a 32×32 image occupies 768 MB of memory if stored as a full matrix, which is prohibitively large. Fortunately, Jacobian matrices of major operators (such as convolution, ReLU, and max-pooling) are usually extremely sparse as shown in Figure 5. In comparison, representing the data contained in the same Jacobian of the aforementioned convolution operator in the Compressed Sparse Row (CSR) (Saad, 2003) format shrinks the memory consumption down to only 6.5 MB. We can observe that there are two reasons for zeros to appear in an operator's Ja-

cobian: guaranteed zeros that are input (\vec{x}_0) invariant (e.g., zeros that are not on the diagonal of ReLU's Jacobian) and related to the model's architecture; and possible zeros that depend on the input (e.g., zeros on the diagonal of ReLU's Jacobian). For any operator, the positions of guaranteed zeros (named as the sparsity pattern for brevity) in the Jacobian is deterministic with the model architecture and known ahead of training time. Thus, mapping non-zero elements in the input matrices to each non-zero element in the product matrix (e.g., calculating the number of non-zeros and index merging in CSR matrix-matrix multiplication (Kunchum et al., 2017)) can be performed prior to training and removed from a generic sparse matrix multiplication routine (e.g., cuSPARSE (NVIDIA, 2018)) to achieve significantly better performance during the training phase. As an example, the second last column of Table 1 shows the extremely high sparsity of guaranteed zeros (defined as the fraction over all elements in a matrix) for various operators in VGG-11 (Simonyan & Zisserman, 2015). In our implementation, the transposed Jacobian matrices are represented in the CSR format since it is the most straightforward and commonly used sparse matrix format; however, any other sparse matrix format can be used as an alternative, including a potentially more efficient customized sparse matrix format that utilizes the deterministic property of the current sparsity pattern, which we leave to investigate as part of our future work.

3.4 Generating Jacobian Matrix in CSR Analytically

To practically generate the Jacobian for an operator, instead of generating one column at a time either numerically (Mahaffy, 2019) or via automatic differentiation (Paszke et al., 2017; PyTorch Forums, 2019), we develop analyti-

Table 1: The sparsity expressions of guaranteed zeros for various operators.

Operator	Filter/Kernel Size	Input Size	Output Size	Sparsity	Examples †	Analytical Generation Speedup §
Convolution	$c_o \times c_i \times h_f \times w_f$	$c_i \times h_i \times w_i$	$c_o \times h_o \times w_o$	$1 - \frac{h_f w_f}{h_i w_i}^{\ddagger}$	0.99157	$8.3 \times 10^3 \times$
ReLU	N/A		$c \times h \times w$			$1.2 \times 10^6 \times$
Max-pooling	$h_f imes w_f$	$c_i \times h_i \times w_i$	$c_o \times h_o \times w_o$	$1 - \frac{h_f w_f}{c_i h_i w_i}$	0.99994	$1.5 \times 10^5 \times$

[†] The examples of sparsity for the first convolution, ReLU and max-pooling operators of VGG-11 (Simonyan & Zisserman, 2015) operating on 32 × 32 images are shown in the second last column of the table.

[‡] Approximation when h_i and w_i are much greater than the padding size.

Sover generating the transposed Jacobian through PyTorch's Autograd (Paszke et al., 2017) one column at a time; measured on a Ryzen Threadripper 1950X (AMD, 2019b) machine; averaged across 1000 trials.

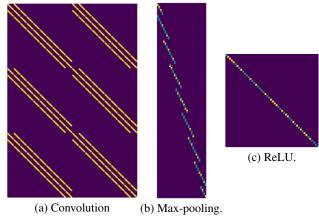
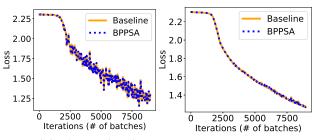


Figure 5: Transposed Jacobians for various operators. Yellow, cyan and purple dots represents locations of non-zero elements, possible zeros and guaranteed zeros in the matrix.

cal routines to generate the transposed Jacobian directly into the CSR format. Appendix D demonstrates such analytical routines in detail for the convolution, ReLU and max-pooling operators. As proofs of concept for the potential performance benefits, the last column of Table 1 shows the speedup on analytical generation of the transposed Jacobians for the aforementioned operators in VGG-11 (Simonyan & Zisserman, 2015). As part of our future work to build a mature framework with automatic differentiation capability that performs training via BPPSA, we aim to provide a library that implements a "sparse transposed Jacobian operator" (replacing the backward operator in the case of cuDNN (NVIDIA, 2019a)) for each forward operator.

3.5 Convergence

Theoretically, our algorithm is a reconstruction of BP instead of an approximation, and hence, expected to reproduce the exact same outputs. However, in practice, numerical differences could be introduced due to the change in the order of matrix multiplications. We apply our algorithm to train LeNet-5 (Lecun et al., 1998) on CIFAR-10 (Krizhevsky, 2009) to demonstrate that such numerical differences would not affect model convergence. We use a mini-batch size of 256 and the SGD (Qian, 1999) optimizer with a learning rate of 0.001 and a momentum of 0.9. We seed the experiments



(a) Training loss per iteration.

(b) Test loss per iteration.

Figure 6: Training and test loss per iteration for training LeNet-5 on CIFAR-10. *Baseline* represents training via the PyTorch Autograd, while *BPPSA* represents our method.

with the same constant. Figure 6 shows that the orange lines overlap with the blue lines for both training and test losses, which means our algorithm has negligible impact on the convergence compared to the original BP.

3.6 Complexity Analysis

Runtime Complexity We leverage the following definitions to quantify the complexity of a parallel algorithm: (1) step complexity (S) which evaluates the minimum number of steps to finish the execution on the critical path (endto-end) given the number of parallel workers; (2) *per-step* complexity(P) which evaluates the runtime of a single step; and (3) work complexity (W) which evaluates the number of total steps executed by all workers. For brevity, we refer to performing the scan operation serially as *linear scan*, which is essentially emulating BP by using the transposed Jacobian and multiplying it with the gradient (as shown in Equation 3) explicitly. Assuming the system can be conceptualized as a parallel random-access machine (PRAM) (Kruskal et al., 1990), the number of workers is p and the size of the input array in Equation 5 is n+1, the step and work complexity of our algorithm can be derived as:

$$S_{Blelloch}(n) = \begin{cases} \Theta(\log n) & p > n \\ \Theta(n/p + \log p) & \text{otherwise} \end{cases}$$
 (6)

$$W_{Blelloch}(n) = \Theta(n)$$
 (7)

compared to $S_{Linear}(n) = \Theta(n), W_{Linear}(n) = \Theta(n)$ of the linear scan (which has the same step and work complexity as BP). Therefore, in an ideal scenario where there is an

unbounded number of workers with unit per-step complexity, our algorithm reduces the runtime of BP from $\Theta(n)$ to $\Theta(\log n)$. If, however, we consider the difference in per-step complexity between our algorithm $(P_{Blelloch})$ and the baseline (P_{Linear}) due to runtime difference between matrix-matrix and matrix-vector multiplications, our algorithm has a runtime of $\Theta(\log n)P_{Blelloch}$ compared to $\Theta(n)P_{Linear}$ in the baseline. There are two approaches to make our algorithm achieve a lower runtime and better scaling than the baseline. First, we can reduce $P_{Blelloch}$, which is reflected in leveraging the sparsity in the transposed Jacobian as analyzed in Section 4.3 and Section 5.3. Second, without lowering $P_{Blelloch}$, our algorithm can still outperform the baseline if $P_{Blelloch}/P_{Linear} < \Theta(n/\log n)$. This can occur when $n/\log n$ grows larger than the dimension of \vec{x}_i . The performance benefit of such case is demonstrated in Section 4.1 and Section 5.1.

Space Complexity Assuming space of storing a transposed Jacobian matrix is bounded by M_{Jacob} and storing \vec{x}_i is bounded by $M_{\vec{x}}$ (note that $M_{Jacob} \ll O(M_{\vec{x}}^2)$ due to sparse matrix formats; both M_{Jacob} and $M_{\vec{x}}$ are not functions of p), in our method, each worker requires the space of $M_{Blelloch}(n) = \Theta(max(\frac{n}{p},1))M_{Jacob}$ which reduces as p increases until a constant M_{Jacob} , comparing to $M_{Pipeline} = \Theta(\frac{n}{p} + p)M_{\vec{x}}$ for pipeline parallelism which increases linearly as p increases. Therefore, our method does not have the limitation of scalability on p, as long as each worker has the memory capacity of at least M_{Jacob} .

4 METHODOLOGY

Although training deep learning models on thousands of devices has been proven feasible in the industry (MLPerf, 2019; Goyal et al., 2017), setting up an experiment for such a large number of devices would require a data center of GPUs and re-implementing/optimizing our entire experiment framework, which requires both monetary and engineering resources out of reach for a typical academic research group. Thus, we set up small-scale experiments that can reflect the large-scale workloads to demonstrate the potential performance benefits of our method.

Environment Setup Our experiments are performed on two platforms with RTX 2070 (NVIDIA, 2019b) and RTX 2080Ti (NVIDIA, 2019c) respectively (both are Turing architecture GPUs) whose specifications are listed in Table 2³.

Baselines We evaluate our method against *PyTorch Autograd* (Paszke et al., 2017) with cuDNN backend (NVIDIA, 2019a) which is a widely adopted and state-of-the-art implementation of BP.

Metrics We use three metrics to quantify the results from our evaluations: (1) *wall-clock time* which measures the

Table 2: Specifications of our experiment platforms.

GPU	RTX 2070	RTX 2080Ti
Number of Streaming Multiprocessors (SMs)	36	68
NVIDIA GPU Driver	430.50	440.33.01
CUDA (Nickolls et al., 2008)	10.0.130	10.0.130
cuDNN (Chetlur et al., 2014)	7.5.1	7.6.2
PyTorch (Paszke et al., 2017)	1.1.0	1.2.0
CPU	Ryzen Threadripper 1950X (AMD, 2019b)	EPYC 7601 (AMD, 2019a)
Host Memory	32GB, 2400MHz	128GB, 2133MHz
Linux Kernel (Torvalds, 2019)	4.15.0-76	4.19.49

system-wide actual time spent on a process, (2) *speedup* which is the ratio of the wall-clock time spent on the baseline over our method, and (3) *FLOP* which represents the number of floating-point operations executed.

We leverage three types of benchmarks to empirically evaluate BPPSA: (1) an end-to-end benchmark of a vanilla RNN training on synthesized datasets to demonstrate the scalability benefits of BPPSA on long sequential dependency; (2) an end-to-end benchmark of a GRU training on the IRMAS dataset (Bosch et al., 2012) to demonstrate the potential of BPPSA on a more realistic workload; and (3) a microbenchmark of a pruned VGG-11 (Simonyan & Zisserman, 2015) to evaluate the feasibility of using sparse matrix format to reduce the per-step complexity of BPPSA.

4.1 RNN End-to-end Benchmark

We set up experiments of training an RNN (Elman, 1990) on sequential data, which is a classical example of workloads where the runtime performance (in terms of the wall-clock time) is limited due to the *strong sequential dependency*. The large number of operators n is modeled through a large sequence length T. The large number of workers p is reflected in the total number of CUDA threads that can be executed concurrently in all SMs of a single GPU, which we model through the fraction of GPU per sample (derived as one over the mini-batch size B).

Datasets We synthesize the datasets $X = \{(x^{(k)}, c^{(k)})\}$ of 32000 training samples (i.e., $k \in \{0, 31999\}$) for the task of *bitstream classification*. Each sample consists of a class label $c^{(k)}$ where $c^{(k)} \in \{0, ..., 9\}$ and a bitstream $x^{(k)}$ where the value $x_t^{(k)}$ at each time step $t \in \{0, ..., T-1\}$ is sampled from the Bernoulli distribution (Bernoulli, 1713; Evans & Rosenthal, 2009):

$$x_t^{(k)} \sim Bernoulli(0.05 + c^{(k)} \times 0.1) \tag{8}$$

Equivalently, each bitstream $x^{(k)}$ can be viewed as a binomial experiment (Bernoulli, 1713; Evans & Rosenthal, 2009) of class $c^{(k)}$. The objective of this task is to classify each bitstream $x^{(k)}$ into its corresponding class $c^{(k)}$ correctly. We synthesize eight datasets with different T, where T increases up to 30000. In reality, long sequences of input can often be found in audio signals such as speech (Nagrani et al., 2017; Baumann et al., 2018; Barker et al., 2018) or

³Appendix G includes results on V100 (NVIDIA, 2019d).

music (Bertin-Mahieux et al., 2011; Benzi et al., 2016).

Model We leverage a vanilla RNN (Elman, 1990) (described in Equation 9) to solve the aforementioned task, since RNN is an intuitive, yet classical, deep learning model and often used to process sequential data:

$$\vec{h}_{t}^{(k)} = tanh(W_{ih}x_{t}^{(k)} + \vec{b}_{ih} + W_{hh}\vec{h}_{t-1}^{(k)} + \vec{b}_{hh})$$
 (9)

where $\vec{h}_t^{(k)}$, \vec{b}_{ih} , $\vec{b}_{hh} \in \mathbb{R}^{20}$. The output classes are predicted via the softmax function (Bridle, 1990) applied on a linear transformation to the last hidden states $\vec{h}_{T-1}^{(k)}$. The cross entropy (Goodfellow et al., 2016) is used as the loss function which is optimized in training via the Adam optimizer (Kingma & Ba, 2015) with the learning rate of 1×10^{-5} . The computation of $\nabla_{\vec{h}_t^{(k)}}l$ during the backward pass carries the *strong sequential dependency* which is the target for acceleration via BPPSA.

Implementation We implement our modified version of the Blelloch scan algorithm as two custom CUDA kernels for the up-sweep and down-sweep phases respectively, along with a few other CUDA kernels for the preparation of the input transposed Jacobian matrices. Each level during the up-/down-sweep is associated with a separate CUDA kernel launch (in the same CUDA stream); therefore, synchronization is ensured between two consecutive levels. Each thread block is responsible for the \diamond operation (i.e. multiplication in reverse) of two matrices as well as moving the intermediate results, and the shared memory is leveraged for caching input and output matrices. Our custom CUDA kernels are integrated into the Python front-end where the RNN and the training procedure are defined through PyTorch's Custom C++ and CUDA Extensions (Goldsborough, 2019). For the forward pass and the baseline of PyTorch Autograd (Paszke et al., 2017), we simply plug in the PyTorch's RNN module (PyTorch, 2019b) which calls into the cuDNN's RNN implementations (cudnnRNNForwardTraining and cudnnRNNBackwardData) (NVIDIA, 2019a); therefore, our baseline is already much faster than implementing RNN in Python using PyTorch's RNNCell module (Py-Torch, 2019b) due to GEMM streaming and kernel fusions (Appleyard et al., 2016).

4.2 GRU End-to-end Benchmark

To extend the aforementioned RNN end-to-end benchmark to a more realistic setting, we evaluate the runtime performance of training a GRU (Cho et al., 2014) on the IRMAS (Bosch et al., 2012) dataset for the task of *instrument classification* based on audio signals.

Datasets We preprocess the IRMAS dataset and compute the mel-frequency cepstral coefficients (MFCC) (Davis & Mermelstein, 1980) for each waveform audio sample via LibROSA's (Brian McFee et al., 2015) MFCC implementation. With different MFCC configurations as listed in

Table 3: MFCC configurations and the resulting feature sizes (represented as the number of frames F multiplied by the number of coefficients C) for the S, M and L sets.

Set Name	S	M	L
MFCC Coefficients FFT Window Length Hop Length	20 4096 512	13 2048 256	7 1024 128
Resulting Input Features $(F \times C)$	259×38	517×24	1034×12

Table 3, the preprocessing results in three sets (*S*, *M* and *L*), reflecting the trade-off between the temporal and frequency resolutions. For all samples, we normalize the values of each coefficient across the frames to have zero mean and unit variance. We remove the first coefficient because it only represents the average power of the audio signal.

Model Since instrument classification is a more complex task than the synthetic workloads in Section 4.1, a GRU (Cho et al., 2014) (described in Equations 10) is used in this set of experiments.

$$\vec{r}_{t} = \sigma(W_{ir}\vec{x}_{t} + \vec{b}_{ir} + W_{hr}\vec{h}_{t-1} + \vec{b}_{hr})$$

$$\vec{z}_{t} = \sigma(W_{iz}\vec{x}_{t} + \vec{b}_{iz} + W_{hz}\vec{h}_{t-1} + \vec{b}_{hz})$$

$$\vec{n}_{t} = \tanh(W_{in}\vec{x}_{t} + \vec{b}_{in} + \vec{r}_{t} \circ (W_{hn}\vec{h}_{t-1} + \vec{b}_{hn}))$$

$$\vec{h}_{t} = (1 - \vec{z}_{t}) \circ \vec{n}_{t} + \vec{z}_{t} \circ \vec{h}_{t-1}$$
(10)

where $\vec{h}_t \in \mathbb{R}^{20}$, $t \in \{0, ..., F-1\}$ and $\vec{x}_t \in \mathbb{R}^C$. Since cuDNN's GRU implementation (Appleyard et al., 2016) is closed source, we are unable to generate the transposed Jacobians efficiently (Appendix E), which leads to significant *overhead* in the forward pass. However, such overhead could potentially be reduced if cuDNN's source code becomes publicly available. Other settings are the same as Section 4.1 with the exception of a 3×10^{-4} learning rate.

Implementation We directly use PyTorch's GRU module (PyTorch, 2019b) which calls into the cuDNN's GRU implementations (cudnnRNNForwardTraining and cudnnRNNBackwardData with CUDNN_GRU) (NVIDIA, 2019a). We reuse the same CUDA implementation of the Blelloch scan algorithm as in Section 4.1.

4.3 Pruned VGG-11 Micro-benchmark

Despite the recent advances in network pruning algorithms (Han et al., 2015; See et al., 2016; He et al., 2017), there is no existing widely adopted software or hardware platform that can exploit performance benefits from pruning, as most techniques are evaluated through "masking simulation" which leads to the same (if not worse) runtime and memory usage. In contrast, in this work, we discover that the *retraining of pruned networks* could benefit from BPPSA due to the following reason: Since the values in the Jacobian of a convolution operator only depend on the filter weights (Appendix D.1), pruning the weights can lead to a higher sparsity in the Jacobian, which then reduces the per-step

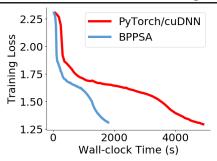


Figure 7: Training loss across wall-clock time when the RNN is trained via BPPSA (blue curve) and the PyTorch Autograd baseline with cuDNN's RNN backend (red curve).

complexity of sparse matrix-matrix multiplications.

To evaluate the feasibility of leveraging the sparsity in the transposed Jacobian of each operator, we set up a benchmark with VGG-11 (Simonyan & Zisserman, 2015): training on CIFAR-10 (Krizhevsky, 2009), pruning away 97% of the weights in all convolution and linear operators using the technique proposed by See et al. (See et al., 2016), and retraining the pruned network. We choose this pruning percentage so that a similar validation accuracy is reached (90.1% v.s. 88.9%) after retraining for the same number of epochs (100) as training. We then apply BPPSA on the convolutional layers of VGG-11 to compute Equation 3.

Since the sparsity pattern of the transposed Jacobian can be determined ahead of training time from the model architecture (as we show in Section 3.3), existing sparse matrix libraries which target generic cases are sub-optimal for our method. For example, cuSPARSE (NVIDIA, 2018) calculates the number of non-zeros in the product matrix and merges the indices of the input matrices before it can perform the multiplication. Such preparations do not need to repeat across iterations in BPPSA's case and could be performed ahead of time due to the deterministic nature of the sparsity pattern. This, in turn, saves considerable amount of execution time. Therefore, due to the lack of a fair implementation, we perform the evaluation by calculating the FLOPs needed for each step in our method and the baseline implementations through *static analysis*.

5 EVALUATION

In this section, we present the results from the RNN end-to-end benchmark (Section 4.1), the GRU end-to-end benchmark (Section 4.2) and the pruned VGG-11 microbenchmark (Section 4.3).

5.1 RNN End-to-end Benchmark

Figure 7 shows the training curves of loss values with respect to wall-clock time when we train the RNN for 80 epochs on the RTX 2070 GPU with the mini-batch size B=16 and the sequence length T=1000. This experiment can be viewed as the simplest mechanism to process

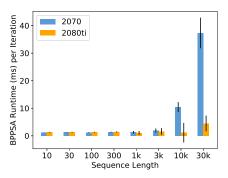
sequential data such as audio signals. We observe that the blue curve (BPPSA) is roughly equivalent to the red curve (PyTorch/cuDNN baseline) scaled down by 63% along the horizontal (time) axis. We conclude that, in this setting, training the RNN through BPPSA reconstructs the original BP algorithm while achieving a $2.73\times$ speedup on the overall training time and $16\times$ on the BP runtime.

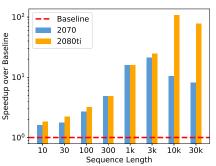
Sensitivity Analysis We measure the performance variation as the sequence length T and the fraction of GPU per sample (1/B) vary, since those two parameters represent the total number of operators n and the number of workers p respectively — key variables in the theoretical runtime of our method. To estimate the speedups, we measure the wall-clock time of training via BPPSA for a single epoch, and take the average of 20 measurements from different epochs. We then compare against training via the PyTorch/cuDNN baseline measured in the same way. We can also derive the backward pass runtime by measuring the wall-clock time of the training procedure without actually performing the backward pass, and subtracting from the total runtime (including the overhead of preparing the input transposed Jacobians).

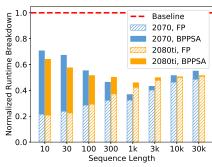
Figure 8a, Figure 8b and Figure 8c show how changing the sequence length T affects the backward pass and overall training time. We make three observations from these figures. First, our method scales as n increases when n is relatively in the same range as p. Second, when n increases to be much larger than p, the performance starts to be bounded by p. Third, even in the range of overly large n, our method still achieves better utilization on massively parallel hardware than the baseline.

Figure 8d, Figure 8e and Figure 8f show how changing the fraction of GPU per sample (1/B) affects the backward pass and overall training time. We can conclude that BPPSA scales as the "effective" number of workers p per sample increases (equivalently, as the batch size B decreases, since the total number of SMs in the GPU is constant). In reality, determining the appropriate mini-batch size can be nontrivial: training with large batch can lead to "generalization gap" (Keskar et al., 2016), while training with small batch would under-utilize the hardware resources and lead to longer training time. Here, BPPSA can be viewed as offering an alternative to train with smaller mini-batch while utilizing the hardware resources more efficiently than BP.

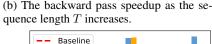
By comparing the speedup in Figure 8b and Figure 8e between RTX 2070 and RTX 2080Ti (RTX 2080Ti has a higher number of SMs than RTX 2070; 68 vs. 36 (NVIDIA, 2019c;b)), we can observe that: (1) BPPSA achieves its maximum speedup at a higher sequence length on RTX 2080Ti than RTX 2070; (2) as the batch size *B* increases, the speedup of BPPSA on RTX 2080Ti drops at a slower rate than RTX 2070. These two observations, together with Figure 8a and Figure 8d where the BPPSA latency per itera-

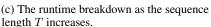


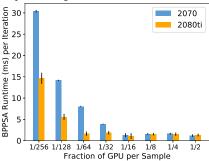


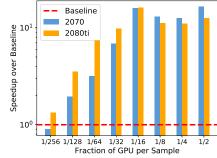


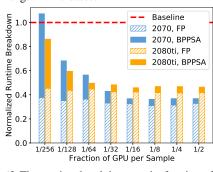
(a) The BPPSA runtime per iteration as the sequence length T increases.











(d) The BPPSA runtime per iteration as the fraction of GPU per sample (1/B) increases.

(e) The backward pass speedup as the fraction of GPU per sample (1/B) increases.

(f) The runtime breakdown as the fraction of GPU per sample (1/B) increases.

Figure 8: We report the BPPSA backward pass latency per iteration (Figure 8a and Figure 8d), the backward pass speedups (Figure 8b and Figure 8e) of BPPSA over the baseline, as well as the runtime (normalized by the baseline) breakdowns to demonstrate the overall speedups (Figure 8c and Figure 8f). The fraction of GPU per sample (which reflects the number of workers *p*) is computed as one over the batch size *B*. *FP* refers to the forward pass. The standard deviations of the BPPSA latency are reported as black lines in Figure 8a and Figure 8d.

tion on RTX 2080Ti is lower than RTX 2070, are consistent with the aforementioned conclusions regarding the performance variation with the number of workers p. We can observe a maximum of $108\times$ backward pass speedup on RTX 2080Ti and a maximum of $2.75\times$ overall speedup on RTX 2070 (the highest backward pass speedup might not lead to the highest overall speedup due to different forward pass runtime on which BPPSA has no impact).

5.2 GRU End-to-end Benchmark

We include the training curves of loss values with respect to the wall-clock time when we train the GRU with the preprocessed datasets in Appendix F. They leads to the same conclusions as the ones in Section 5.1.

Sensitivity Analysis To perform an analysis similar to Section 5.2, we only need to vary the batch size B since the preprocessed dataset type (S, M, L) already reflects the sequence length T. However, since the overhead of computing the transposed Jacobians during the forward pass cannot be neglected (as mentioned in Section 4.2), to achieve a deeper understanding of the performance variation, we demonstrate the runtime breakdowns among the forward pass, the backward pass and the overhead. We can derive the

overhead by taking the difference in the runtime of the training procedures without actually performing the backward pass between BPPSA and the PyTorch/cuDNN baseline. The measurements are averaged across 100 epochs.

Figure 9 shows how the sequence length T and batch size B affect the runtimes of the forward pass, the backward pass and the overhead. We make two observations from this figure. First, our method achieves a higher speedup on the backward pass as T increases (changing the preprocessed dataset from S to L), which reinforces the observation from Section 5.1 that our method scales well as the total number of operators n increases. Second, since the maximum sequence length (1034) is not as extreme as in Section 5.1, the backward pass runtime of BPPSA is less affected than the overhead by B and the GPU model, which means n is still within the same range as the number of workers p in this set of experiments. The maximum overall speedup and backward pass speedup (excluding the overhead) are $2.36 \times$ and $13.4 \times$ respectively.

5.3 Pruned VGG-11 Micro-benchmark

Since the sparsity of the product matrix might reduce after each multiplication, the per-step complexity might increase

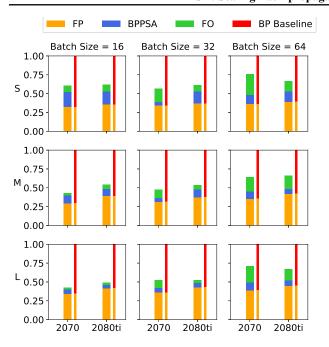


Figure 9: The runtime breakdowns in the GRU end-to-end benchmark as the dataset type (S, M, L) and batch sizes B vary. FP represents the forward pass; FO represents the forward pass overhead of computing the transposed Jacobians; BP represents the BP baseline; and BPPSA represents the backward pass via BPPSA. The measurements are normalized by the total runtime of the baseline (FP + BP).

as the up-sweep phase progresses into deeper levels. Fortunately, we can adopt BPPSA to balance the number of levels in the up-/down-sweep phases according to the sparsity of the products on each level to achieve an overall speedup. Specifically, in this experiment, BPPSA performs the up-sweep from L0 to L2 (consistent with the notations in Figure 3), calculates the partial results that are needed for the down-sweep phase through linear scan, and then performs the down-sweep from L7 to L10.

Assuming the sparse transposed Jacobian matrices are encoded in the CSR format, Figure 10 shows the calculated FLOP of each step in BPPSA and each "gradient operator" in the baseline (BP) for re-training pruned VGG-11 on CIFAR-10. We observe that the green circles (baseline) have similar expected performance as the other circles (BPPSA). Thus, we can conclude that exploiting the sparsity in the transposed Jacobian is an efficient strategy that reduces the per-step complexity of our method $P_{Belloch}$ to a level similar with the baseline P_{Linear} . This strategy makes the overall scalability to be "ensured" algorithmically.

6 CONCLUSION

In this work, we explore a novel direction to scale BP by challenging its fundamental limitation: the *strong sequen*-

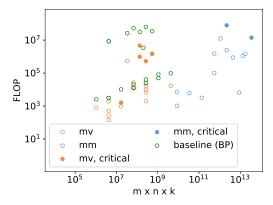


Figure 10: Measuring FLOP for each step when retraining pruned VGG-11 on CIFAR-10. *mv* and *mm* represent matrix-vector and matrix-matrix multiplications in BPPSA respectively. *critical* indicates that the step is on the critical path. The x-axis represents the theoretical runtime complexity of the step **if** the transposed Jacobian were not encoded in a sparse format. The green circles represent the FLOP estimated for each "gradient operator" in the BP baseline.

tial dependency. We reformulate BP into a scan operation which is scaled by our modified version of the Blelloch scan algorithm. Our proposed algorithm, BPPSA, achieves a logarithmic, rather than linear, step complexity. In addition, BPPSA has a constant per-device space complexity; hence, its scalability is not limited by the memory capacity of each device. In our detailed evaluations, we demonstrate that performance benefits can be achieved in two important use cases. First, for the case where there is a long dependency in BP, we evaluate BPPSA by training a RNN with synthetic datasets and training a GRU with the IRMAS dataset (Bosch et al., 2012), where our method achieves up to $2.75 \times$ speedup on the overall (end-to-end) training time and 108× speedup on the backward pass runtime. Second, we can reduce the per-step complexity by leveraging the sparsity in the Jacobian itself. To this end, we develop efficient routines to generate the transposed Jacobian in the CSR format, and demonstrate that the retraining of pruned networks can potentially benefit from BPPSA (as we show for a pruned VGG-11 benchmark when re-training on the CIFAR-10 dataset). We hope that our work will inspire radically new ideas and designs to improve distributed DNN training beyond the existing theoretical frameworks.

ACKNOWLEDGEMENTS

We want to thank Xiaodan (Serina) Tan, James Gleeson, Geoffrey Yu, Roger Grosse, Jimmy Ba, Andrew Pelegris, Bojian Zheng, Kazem Cheshmi and Maryam Mehri Dehnavi for their constructive feedback during the development of this work. This work was supported in part by the NSERC Discovery grant, the Canada Foundation for Innovation JELF grant, the Connaught Fund, and Huawei grants.

REFERENCES

- Amazon. Amazon ec2 p3 instances. https://aws.amazon.com/ec2/instance-types/p3/, 2019. Accessed: 2019-08-28.
- AMD. Epyc 7601 dual amd server processors. https://www.amd.com/en/products/cpu/amd-epyc-7601, 2019a. Accessed: 2019-08-28.
- AMD. Ryzen threadripper 1950x processor. https://www.amd.com/en/products/cpu/amd-ryzen-threadripper-1950x, 2019b. Accessed: 2019-08-28.
- Amodei, D., Hernandez, D., Sastry, G., Clark, J., Brockman, G., and Sutskever, I. Ai and compute. https://openai.com/blog/ai-and-compute/, 2018. Accessed: 2019-08-28.
- Appleyard, J., Kociský, T., and Blunsom, P. Optimizing performance of recurrent neural networks on gpus. *CoRR*, abs/1604.01946, 2016. URL http://arxiv.org/abs/1604.01946.
- Arunkumar, A., Bolotin, E., Cho, B., Milic, U., Ebrahimi, E., Villa, O., Jaleel, A., Wu, C.-J., and Nellans, D. Mcmgpu: Multi-chip-module gpus for continued performance scalability. *ACM SIGARCH Computer Architecture News*, 45(2):320–332, 2017.
- Barker, J., Watanabe, S., Vincent, E., and Trmal, J. The fifth 'chime' speech separation and recognition challenge: Dataset, task and baselines. *CoRR*, abs/1803.10609, 2018. URL http://arxiv.org/abs/1803.10609.
- Baumann, T., Köhn, A., and Hennig, F. The spoken wikipedia corpus collection: Harvesting, alignment and an application to hyperlistening. *Language Resources and Evaluation*, Jan 2018. doi: 10.1007/s10579-017-9410-y. URL https://doi.org/10.1007/s10579-017-9410-y.
- Ben-Nun, T. and Hoefler, T. Demystifying parallel and distributed deep learning: An in-depth concurrency analysis. *CoRR*, abs/1802.09941, 2018. URL http://arxiv.org/abs/1802.09941.
- Benzi, K., Defferrard, M., Vandergheynst, P., and Bresson, X. FMA: A dataset for music analysis. *CoRR*, abs/1612.01840, 2016. URL http://arxiv.org/abs/1612.01840.
- Bernoulli, J. Jacobi Bernoulli, ... Ars conjectandi, opus posthumum. Accedit Tractatus de seriebus infinitis, et epistola Gallice scripta De ludo pilae reticularis. impensis Thurnisiorum, fratrum, 1713.

- Bertin-Mahieux, T., Ellis, D. P., Whitman, B., and Lamere, P. The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.
- Blelloch, G. E. Prefix sums and their applications. Technical Report CMU-CS-90-190, School of Computer Science, Carnegie Mellon University, Nov. 1990.
- Bosch, J., Janer, J., Fuhrmann, F., and Herrera, P. A comparison of sound segregation techniques for predominant instrument recognition in musical audio signals. *Proceedings of the 13th International Society for Music Information Retrieval Conference, ISMIR 2012*, pp. 559–564, 01 2012.
- Brian McFee, Colin Raffel, Dawen Liang, Daniel P.W. Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. librosa: Audio and Music Signal Analysis in Python. In Kathryn Huff and James Bergstra (eds.), *Proceedings of the 14th Python in Science Conference*, pp. 18 24, 2015. doi: 10.25080/Majora-7b98e3ed-003.
- Bridle, J. S. Probabilistic interpretation of feedforward classification network outputs, with relationships to statistical pattern recognition. In Soulié, F. F. and Hérault, J. (eds.), *Neurocomputing*, pp. 227–236, Berlin, Heidelberg, 1990. Springer Berlin Heidelberg.
- Chen, T., Xu, B., Zhang, C., and Guestrin, C. Training deep nets with sublinear memory cost. *CoRR*, abs/1604.06174, 2016. URL http://arxiv.org/abs/1604.06174.
- Chetlur, S., Woolley, C., Vandermersch, P., Cohen, J., Tran, J., Catanzaro, B., and Shelhamer, E. cudnn: Efficient primitives for deep learning. *CoRR*, abs/1410.0759, 2014. URL http://arxiv.org/abs/1410.0759.
- Cho, K., van Merrienboer, B., Gülçehre, Ç., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014. URL http://arxiv.org/abs/1406.1078.
- Coleman, C. A., Narayanan, D., Kang, D., Zhao, T. J., Zhang, J., Nardi, L., Bailis, P., Olukotun, K., Re, C. B., and Zaharia, M. A. Dawnbench: An end-to-end deep learning benchmark and competition. 2017.
- Davis, S. and Mermelstein, P. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(4):357–366, August 1980. ISSN 0096-3518.

- Deng, J., Dong, W., Socher, R., Li, L., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255, June 2009. doi: 10.1109/CVPR.2009.5206848.
- Elman, J. L. Finding structure in time. *Cognitive Science*, 14 (2):179–211, 1990. doi: 10.1207/s15516709cog1402_1.
- Esmaeilzadeh, H., Blem, E., St. Amant, R., Sankaralingam, K., and Burger, D. Dark silicon and the end of multicore scaling. In *Proceedings of the 38th Annual International Symposium on Computer Architecture*, ISCA '11, pp. 365–376. ACM, 2011.
- Evans, M. and Rosenthal, J. *Probability and Statistics: The Science of Uncertainty*. W. H. Freeman, 2009. ISBN 9781429281270.
- Goldsborough, P. Custom c++ and cuda extensions. https://pytorch.org/tutorials/advanced/cpp_extension.html, 2019. Accessed: 2019-08-28.
- Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org.
- Goyal, P., Dollár, P., Girshick, R. B., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017. URL http://arxiv.org/abs/1706.02677.
- Han, S., Pool, J., Tran, J., and Dally, W. Learning both weights and connections for efficient neural network. In Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems* 28, pp. 1135–1143. 2015.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. URL http://arxiv.org/abs/1512.03385.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. *CoRR*, abs/1603.05027, 2016. URL http://arxiv.org/abs/1603.05027.
- He, Y., Zhang, X., and Sun, J. Channel pruning for accelerating very deep neural networks. *CoRR*, abs/1707.06168, 2017. URL http://arxiv.org/abs/1707.06168.
- Hillis, W. D. and Steele, Jr., G. L. Data parallel algorithms. *Commun. ACM*, 29(12):1170–1183, December 1986.
- Huang, G., Liu, Z., and Weinberger, K. Q. Densely connected convolutional networks. *CoRR*, abs/1608.06993,

- 2016. URL http://arxiv.org/abs/1608.06993.
- Huang, Y., Cheng, Y., Chen, D., Lee, H., Ngiam, J., Le, Q. V., and Chen, Z. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *CoRR*, abs/1811.06965, 2018. URL http://arxiv.org/abs/1811.06965.
- Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. *CoRR*, abs/1609.04836, 2016. URL http://arxiv.org/abs/1609.04836.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, 2009.
- Krizhevsky, A. One weird trick for parallelizing convolutional neural networks. *CoRR*, abs/1404.5997, 2014. URL http://arxiv.org/abs/1404.5997.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems Volume 1*, NIPS'12, pp. 1097–1105, 2012.
- Kruskal, C. P., Rudolph, L., and Snir, M. A complexity theory of efficient parallel algorithms. *Theoretical Computer Science*, 71(1):95 132, 1990. doi: https://doi.org/10.1016/0304-3975(90)90192-K.
- Kunchum, R., Chaudhry, A., Sukumaran-Rajam, A., Niu, Q., Nisa, I., and Sadayappan, P. On improving performance of sparse matrix-matrix multiplication on gpus. In *Proceedings of the International Conference on Supercomputing*, ICS '17, pp. 14:1–14:11, 2017.
- Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pp. 2278–2324, 1998.
- Mahaffy, J. Numerical evaluation of jacobians personal.psu.edu. http://www.personal.psu.edu/jhm/ME540/lectures/NumJacobian.html, 2019. Accessed: 2019-08-28.
- MLPerf. Mlperf training v0.6 results. https://mlperf.org/training-results-0-6/, 2019. Accessed: 2019-08-28.
- Mutlu, O. Memory scaling: A systems architecture perspective. pp. 21–25, 05 2013. doi: 10.1109/IMW.2013. 6582088.

- Nagrani, A., Chung, J. S., and Zisserman, A. Voxceleb: a large-scale speaker identification dataset. *CoRR*, abs/1706.08612, 2017. URL http://arxiv.org/abs/1706.08612.
- Narayanan, D., Harlap, A., Phanishayee, A., Seshadri, V., Devanur, N., Granger, G., Gibbons, P., and Zaharia, M. Pipedream: Generalized pipeline parallelism for dnn training. In SOSP 2019, October 2019.
- Nickolls, J., Buck, I., Garland, M., and Skadron, K. Scalable parallel programming with cuda. *Queue*, 6(2):40–53, March 2008. ISSN 1542-7730.
- NVIDIA. cusparse :: Cuda toolkit documentation, 2018. URL https://docs.nvidia.com/cuda/cusparse/index.html. [Accessed: 2018-11-06].
- NVIDIA. cudnn developer guide :: Deep learning sdk documentation. https://docs.nvidia.com/deeplearning/sdk/cudnn-developer-guide/index.html, 2019a. Accessed: 2019-08-28.
- NVIDIA. Geforce rtx 2070 graphics card | nvidia. https://www.nvidia.com/en-us/geforce/graphics-cards/rtx-2070/, 2019b. Accessed: 2019-08-28.
- NVIDIA. Geforce rtx 2080 ti graphics card | nvidia. https://www.nvidia.com/en-us/geforce/graphics-cards/rtx-2080-ti/, 2019c. Accessed: 2019-08-28.
- NVIDIA. Nvidia v100 tensor core gpu. https://www.nvidia.com/en-us/data-center/v100/, 2019d. Accessed: 2019-08-28.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in PyTorch. In NIPS Autodiff Workshop, 2017.
- PyTorch. Torch.nn. https://pytorch.org/docs/stable/nn.html#torch.nn.MaxPool2d, 2019a. Accessed: 2019-05-23.
- PyTorch. torch.nn pytorch master documentation. https://pytorch.org/docs/stable/nn.html, 2019b. Accessed: 2019-08-28.
- Qian, N. On the momentum term in gradient descent learning algorithms. *Neural Netw.*, 12(1):145–151, January 1999.
- Rhu, M., Gimelshein, N., Clemons, J., Zulfiqar, A., and Keckler, S. W. Virtualizing deep neural networks for memory-efficient neural network design. *CoRR*,

- abs/1602.08124, 2016. URL http://arxiv.org/abs/1602.08124.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. Neuro-computing: Foundations of research. chapter Learning Representations by Back-propagating Errors, pp. 696–699. MIT Press, 1988.
- Saad, Y. *Iterative Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2nd edition, 2003.
- See, A., Luong, M., and Manning, C. D. Compression of neural machine translation models via pruning. *CoRR*, abs/1606.09274, 2016. URL http://arxiv.org/abs/1606.09274.
- Shallue, C. J., Lee, J., Antognini, J. M., Sohl-Dickstein, J., Frostig, R., and Dahl, G. E. Measuring the effects of data parallelism on neural network training. *CoRR*, abs/1811.03600, 2018. URL http://arxiv.org/abs/1811.03600.
- Shazeer, N., Cheng, Y., Parmar, N., Tran, D., Vaswani, A., Koanantakool, P., Hawkins, P., Lee, H., Hong, M., Young, C., Sepassi, R., and Hechtman, B. Mesh-tensorflow: Deep learning for supercomputers. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 10414–10423. Curran Associates, Inc., 2018.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR* 2015, 2015.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S.,
 Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich,
 A. Going deeper with convolutions. In *Computer Vision*and Pattern Recognition (CVPR), 2015.
- PyTorch Forums. How to compute jacobian matrix in pytorch? https://discuss.pytorch.org/t/how-to-compute-jacobian-matrix-in-pytorch/14968, 2019. Accessed: 2019-08-28.
- Torvalds, L. Linux kernel source tree. https://github.com/torvalds/linux, 2019. Accessed: 2019-08-28.
- Varoquaux, G., Gouillart, E., and Vahtras, O. Compressed sparse row format (csr), 2019. URL https://scipy-lectures.org/advanced/scipy_sparse/csr_matrix.html. Accessed: 2019-05-23.

Weisstein, E. W. "jacobian." from mathworld-a wolfram web resource. http://mathworld.wolfram.com/Jacobian.html, 2019. Accessed: 2019-08-28.

Zhu, H., Akrout, M., Zheng, B., Pelegris, A., Jayarajan, A., Phanishayee, A., Schroeder, B., and Pekhimenko, G. Benchmarking and analyzing deep neural network training. In 2018 IEEE International Symposium on Workload Characterization, IISWC 2018, pp. 88–100, 2018. doi: 10.1109/IISWC.2018.8573476.

SUMMARY OF APPENDICES

Due to space constraints, we are unable to include several important details in the main text of this paper. We provide these details in the appendix below that includes the following content:

- Appendix A describes our open-sourced artifact and explains how to reproduce all major experiments in this work.
- Appendix B performs an analysis on the space complexity for one of the key prior works, GPipe (Huang et al., 2018). Appendix C describes our initial attempts to analyze PipeDream's (Narayanan et al., 2019) behavior on VGG-16 with the Adam optimizer (instead of a vanilla SGD). We use these two appendices to support our arguments in Section 2.2.
- Appendix D lists the routines that we developed to generate the transposed Jacobians for various operators directly into the CSR format. This is the complementary material for Section 3.4.
- Appendix E shows how to derive the transposed Jacobians for GRU, and demonstrates the source of the overhead in the forward pass for our GRU end-to-end benchmark (Section 4.2).
- Appendix F includes the training curves for our GRU end-to-end benchmark, which serves as a complementary material to Section 5.2.
- Appendix G reports the additional hardware sensitivity results on the Volta-based V100 GPU (NVIDIA, 2019d) to validate BPPSA potential across different GPU generations.

A ARTIFACT APPENDIX

A.1 Abstract

We provide the source code, scripts and data that corresponds to Section 4 as our artifact to reproduce the results

in Section 5 and Table 1. We require an x86-64 based machine with at least one NVIDIA GPU to evaluate the artifact, and NVIDIA Container Toolkit is the only software dependency to prepare. After the installation, the entire workflow (from building the needed Docker image to plotting the final results) is automated by a single workflow. sh script. Although the the exact numerical results produced by the artifact might vary across hardware platforms, the general trends and conclusions should be similar to the results reported in this paper.

A.2 Artifact check-list (meta-information)

- Algorithm: Back-propagation by Parallel Scan Algorithm (BPPSA)
- **Program:** RNN and GRU end-to-end benchmarks (Section 4.1, 4.2); a VGG-11 micro-benchmark Section 4.3. All benchmarks are public, included, and automated.
- Compilation: GCC 7.4.0 and CUDA 10.0 are recommended, included, and tested, although other versions of GCC and CUDA might work as well.
- **Binary:** Scripts included to build binaries from the source code
- **Data set:** The synthetic datasets (Section 4.1) and IRMAS are included; approximately 4.5 GB in total.
- Run-time environment: The main software dependency is NVIDIA Container Toolkit (https://github.com/NVIDIA/nvidia-docker) which dictates the OS requirements. We recommend and tested on Ubuntu 18.04.
- Hardware: An x86-64 based machine with at least one NVIDIA GPU and internet access. No SUDO access needed.
- Run-time state: No contentions on hardware resources (CPU, GPU, RAM, PCIe) with other processes.
- Execution: Around 57 hours in total on the RTX 2080Ti platform listed in Table 2.
- Metrics: Wall-clock time, speedup, and FLOP (Section 4).
- Output: Figures that are similar to Figure 7, 8, 9, 10 and 13; Text that contains speedups similar to the last column of Table 1.
- **Experiments:** A single script is provided that automates the entire workflow.
- How much disk space required (approximately)?: Approximately a total of 19.7 GB is needed.
- How much time is needed to prepare workflow (approximately)?: Around one hour to install NVIDIA Container Toolkit with its dependencies.
- How much time is needed to complete experiments (approximately)?: Refer to the Execution part above.
- Publicly available?: Yes.
- Code licenses (if publicly available)?: MIT.

- Data licenses (if publicly available)?: MIT.
- Workflow framework used?: No.
- Archived (provide DOI)?: 10.5281/zenodo.3605368

A.3 Description

The source code is publicly available on GitHub (https://github.com/UofT-EcoSystem/BPPSA-open) and Zenodo (https://doi.org/10.5281/zenodo.3605368). The source code and scripts only require 37.9 kB of disk space. However, the workflow.sh script builds a 7.7 GB Docker image, then downloads and unzips 12 GB of data.

A.3.1 Hardware dependencies

The hardware specifications used are listed in Table 2. In general, an x86-64 based machine with at least one NVIDIA GPU and internet access is required.

A.3.2 Software dependencies

Although it is possible to run the experiments natively on the host machine (and, in fact, this is how our RTX 2070 platform was set up), we do not recommend this approach since installing the dependencies can be tedious, non-portable, and unsafe (due to the SUDO access requirements). Instead, we package all of the original dependencies into a Docker image which can be built natively by the workflow.sh script. Therefore, our artifact only requires NVIDIA Container Toolkit (https://github.com/NVIDIA/nvidia-docker). We recommend and tested on Ubuntu 18.04, however, it is possible to evaluate the artifact on other Linux distributions that NVIDIA Container Toolkit supports as well.

A.3.3 Data sets

The workflow.sh script downloads all the required datasets automatically.

A.3.4 Models

The RNN (Section 4.1) and GRU (Section 4.2) are included. The transposed Jacobians of VGG-11 are downloaded by the workflow.sh script.

A.4 Installation

Assuming the hardware listed in Section A.3.1 is available, the following steps are needed to perform the installation:

1. Clone the project by git clone https:
 //github.com/UofT-EcoSystem/
 BPPSA-open.git.

- 2. Install a NVIDIA GPU driver that is compatible with the GPU, the CUDA version (10.0 recommended) and the NVIDIA Container Toolkit.
- 3. Install Docker Engine Community (https://docs.docker.com/install/), then configure the docker group to use Docker as a non-root user. https://docs.docker.com/install/linux/linux-postinstall/).
- 4. Install NVIDIA Container Toolkit (https://github.com/NVIDIA/nvidia-docker).

We provide the install.sh script as a reference to the above steps 2 to 4.

A.5 Experiment workflow

We provide the workflow.sh script that automates the entire workflow consisting of the following stages:

- 1. Build the Docker image used across experiments.
- 2. Download and unzip the synthetic datasets (Section 4.1) and IRMAS.
- 3. Execute the RNN (Section 4.1) and GRU (Section 4.2) end-to-end benchmarks.
- 4. Plot the results for the RNN and GRU end-to-end benchmarks.
- 5. Evaluate the speedups for sparse transposed Jacobian generation (Section 3.4).
- 6. Download the sparse transposed Jacobians of a regular and pruned VGG-11.
- 7. Execute the VGG-11 micro-benchmark (Section 4.3) and plot the results.

After the installation in Section A.4, the user only need to run the command "./workflow.sh" in the project root directory, which takes around 57 hours on our reference plarform with the RTX 2080Ti GPU (Table 2).

A.6 Evaluation and expected result

After ./workflow.sh finishes, a results/ directory is created to contain the following results:

- fig_7.png corresponding to Figure 7.
- fig_8_a.png, fig_8_b.png, fig_8_c.png, fig_8_d.png, fig_8_e.png and fig_8_f.png corresponding to Figure 8.
- fig_9.png corresponding to Figure 9.

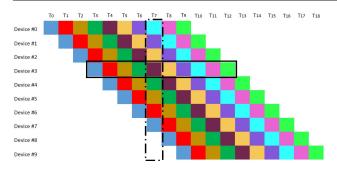


Figure 11: Timing diagram of the forward pass when distributing a model via pipeline parallelism. Each color represents an individual batch.

- fig_10.png corresponding to Figure 10.
- fig_13.png corresponding to Figure 13.
- table_1_last_column.txt corresponding to the last column of Table 1.

The exact numerical results might vary across hardware platforms, but the general trends should be similar to the results presented in this paper where we conducted the experiments on platforms with the RTX 2070 and 2080Ti GPUs (Table 2). In addition, the speedups of BPPSA over BP should be easily observable in the RNN and GRU end-to-end benchmarks.

A.7 Experiment customization

Each stage of the workflow can be turned off independently by commenting out the corresponding lines in workflow.sh. The software environment can be customized by modifying docker/Dockerfile and rebuilding the Docker image. The parameters of the RNN and GRU end-to-end benchmarks can be customized by modifying the code/rnn_grid_run.sh and code/gru_grid_run.sh scripts which are launched by workflow.sh through Docker containers.

A.8 Methodology

Submission, reviewing and badging methodology:

- http://cTuning.org/ae/ submission-20200102.html
- http://cTuning.org/ae/ reviewing-20200102.html
- https://www.acm.org/publications/ policies/artifact-review-badging

B SPACE COMPLEXITY OF GPIPE

Using the notations consistent with GPipe (Huang et al., 2018), with re-materialization enabled, each device reserves

 $\Theta(L/K)$ space for re-computing the intermediate activations of each sample in a "micro-batch", where L and K are the length of the network and the number of devices in the pipeline correspondingly. As we show in Figure 11, to fully fill the pipeline with useful computation, the number of "micro-batches" entering the pipeline (the solid black box) should be equal to the length of the pipeline (the dashed black box); thus, each device needs to store at least $\Theta(K)$ activations at the partition boundary for each sample, resulting in a $\Theta(L/K+K)$ per-device space complexity.

C AFFECT OF PIPEDREAM'S STALENESS ON ADAM

Using the source code from https://github.com/msr-fiddle/pipedream, we reproduce PipeDream's results on VGG-16 with the same settings except the following:

- 4 RTX 2080Ti GPUs (instead of 16 V100 GPUs).
- Mini-batch size of 32 (instead of 64).
- Adam optimizer with the learning rate of 0.00003 and zero weight decay (instead of SGD with the learning rate of 0.01 and the weight decay of 0.0005).
- 90 epochs in total (instead of 60).
- Step decay learning rate schedule (instead of polynomial decay).

For the baseline, we use the source code from https://github.com/pytorch/examples/tree/master/imagenet (which is a plain VGG implementation used as one of PyTorch's official examples for ImageNet (Deng et al., 2009)) and the same aforementioned settings except using one GPU (instead of four). We choose these settings for the following purpose: (1) to fit in the hardware resources available to us; (2) to match with a widely adopted baseline; and (3) to use the Adam optimizer instead of SGD. We run both experiments three times and record the Top-1 and Top-5 validation accuracy across epochs. We present our results in Figure 12.

We observe a 2.6% top-1 and 1.9% top-5 accuracy loss from PipeDream compared to the baseline. Combining with the observation that the error bars are negligible (i.e., negligible variance across runs), we conclude that at least in this case, PipeDream is not fully equivalent to the baseline for some adaptive optimizers (e.g., Adam), which differs from the optimizer-oblivious property of our work. It is important to emphasize that we **do not** imply that PipeDream will always have a negative impact on the convergence. A deeper analysis is needed on a much greater space of hyper-parameters to understand how general is such an effect on convergence with PipeDream that is beyond the scope of our work.

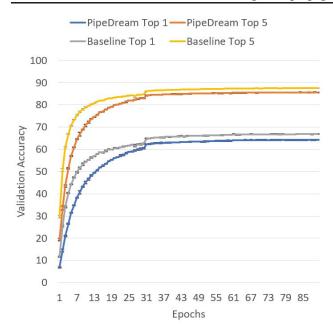


Figure 12: Top-1 and Top-5 validation accuracy across epochs for both PipeDream and the baseline. We report both the mean (as curves) and the range (as error bars).

D SPARSE JACOBIAN GENERATION ROUTINES

D.1 Convolution

Algorithm 2, Algorithm 3 and Algorithm 4 show how to generate the CSR indptr, indices and data arrays (Varoquaux et al., 2019) respectively for the transposed Jacobian of a convolution operator that has a 3×3 filter and padding size of 1.4

Algorithm 2 Compute the CSR indptr array for the transposed Jacobian of a 3×3 convolution.

```
Input: input channels c_i, output channels c_o, input height h_i, input width w_i
Output: indptr \leftarrow malloc(c_i h_i w_i + 1)
 1: for all i \leftarrow 0 to (c_i h_i w_i) do in parallel
        a \leftarrow |i/(h_i w_i)|
        b \leftarrow i \bmod (h_i w_i)
        if b \leqslant w_i then
            indptr[i] \leftarrow ac_o(3w_i(3h_i-2)) + 6c_ob
        else if b \leqslant w_i(h_i - 1) then
 6:
            indptr[i] \leftarrow ac_o(3w_i(3h_i-2)) + 6c_ow_i + 9c_o(b-w_i)
 7:
 8:
            indptr[i] \leftarrow ac_o(3w_i(3h_i-2)) + 6c_ow_i + 9c_o(w_i(h_i-2))
                             +6c_o(b-w_i(h_i-1))
        end if
11: end for
```

D.2 ReLU

Our methods of generating the CSR indptr, indices and data arrays (Varoquaux et al., 2019) for the trans-

Algorithm 3 Compute the CSR indices array for the transposed Jacobian of a 3×3 convolution.

```
Input: input channels c_i, output channels c_o, input height h_i, input width w_i,
    indptr computed from Algorithm 2
Output: indices \leftarrow malloc(3w_i(3h_i)
 1: for all i \leftarrow 0 to (c_i h_i w_i - 1) do in parallel
        r \leftarrow i \mod (h_i w_i)
        base \leftarrow \text{malloc}(9c_o)
 3:
 4:
        for all j \leftarrow 0 to (c_o - 1) do in parallel
            for all k \leftarrow 0 to 2 do in parallel
 5:
                base[9j + 3k : 9j + 3(k+1)] \leftarrow (
 6:
                      [-1,0,1] + (jh_i + k - 1)w_i + r) \bmod (c_o h_i w_i)
 7:
            end for
 8:
        end for
 9:
        if r < w_i or r \geqslant w_i(h_i - 1) then
10:
            row \leftarrow malloc(6c_o)
            (left, right) \leftarrow (3,9) if r < w_i; (0,6) otherwise
11:
            for all j \leftarrow 0 to (c_o - 1) do in parallel
12:
                row[6j:6j+6] \leftarrow base[9j+left:9j+right]
13:
14:
            end for
15:
16:
            row ← base
17:
        end if
        indices[indptr[i]:indptr[i+1]] \leftarrow sorted(row)
18:
19: end for
```

Algorithm 4 Compute the CSR data array for the transposed Jacobian of a 3×3 convolution.

```
\begin{aligned} & \textbf{Input:} & \text{input channels } c_i, \text{ output channels } c_o, \text{ input height } h_i, \text{ input width } w_i, \text{ filter weights, indptr computed from Algorithm 2} \\ & \textbf{Output:} & \text{data} \leftarrow \text{malloc}(3w_i(3h_i-2)c_ic_o) \\ & \text{I:} & \textbf{ for all } i \leftarrow 0 \text{ to } (c_ih_iw_i-1) \text{ do in parallel}} \\ & \text{2:} & r \leftarrow i \text{ mod } (h_iw_i) \\ & \text{3:} & m \leftarrow \lfloor i/(h_iw_i) \rfloor \\ & \text{range} \leftarrow (1::-1) \text{ if } (r < w_i); \\ & \text{4:} & (2:0:-1) \text{ otherwise} \\ & \text{data}[\text{indptr}[i]: \text{indptr}[i+1]] \leftarrow \text{flatten}(\\ & \text{5:} & \text{weights}[:, m, \text{range}, ::-1]) \\ & \text{6:} & \text{Fix comer cases when } (i \text{ mod } w_i) = 0 \text{ or } (i \text{ mod } w_i) = (w_i-1). \\ & \text{7:} & \text{end for} \end{aligned}
```

posed Jacobian of a ReLU operator are formally described in Algorithm 5, Algorithm 6 and Algorithm 7 respectively.

Algorithm 5 Compute the CSR indptr array for the transposed Jacobian of ReLU.

```
Input: size d of the (flattened) input tensor x
Output: indptr \leftarrow malloc(d+1)
1: for all i \leftarrow 0 to d do in parallel
2: indptr[i] \leftarrow i
3: end for
```

D.3 Max-pooling

Assuming the stride size and the window size are the same, and we can access a tensor (named as pool_indices for brevity) which specifies the indices of the elements in the input tensor that are "pooled" for the output tensor (documented in (PyTorch, 2019a)), our methods of generating the CSR indptr, indices and data arrays (Varoquaux et al., 2019) are formally described in Algorithm 8, Algorithm 9 and Algorithm 10 respectively.

⁴Although we assume a specific configuration of the convolution operator here, deriving a generic routine is doable.

Algorithm 6 Compute the CSR indices array for the transposed Jacobian of ReLU.

```
Input: size d of the (flattened) input tensor x

Output: indices \leftarrow malloc(d)

1: for all i \leftarrow 0 to (d-1) do in parallel

2: indices \leftarrow i

3: end for
```

Algorithm 7 Compute the CSR data array for the transposed Jacobian of ReLU.

```
Input: the (flattened) input tensor x, and its size d
Output: \mathtt{data} \leftarrow \mathtt{malloc}(d)
1: for all i \leftarrow 0 to (d-1) do in parallel
2: if x[i] > 0 then
3: \mathtt{data}[i] \leftarrow 1
4: else
5: \mathtt{data}[i] \leftarrow 0
6: end if
7: end for
```

Algorithm 8 Compute the CSR indptr array for the transposed Jacobian of max-pooling.

```
Input: pool_indices, input height h_i, input width w_i, output height h_o, output width w_o, output channels c_o
```

```
Output: indptr \leftarrow malloc(c_oh_iw_i+1), mapping \leftarrow malloc(c_oh_iw_i)
 1: for all i \leftarrow 0 to c_o h_i w_i - 1 do in parallel
        mapping[i] \leftarrow -1
 3: end for
 4: for all c \leftarrow 0 to c_o - 1 do in parallel
         for all h \leftarrow 0 to h_o - 1 do in parallel
              for all w \leftarrow 0 to w_o - 1 do in parallel
                  i \leftarrow ch_i w_i + \texttt{pool\_indices}[c, h, w]
                  j \leftarrow (ch_o + h)w_o + w
 g.
                  \texttt{mapping}[i] \leftarrow j
10:
              end for
11:
         end for
12: end for
13: ptr \leftarrow 0
14: for i \leftarrow 0 to (c_o h_i w_i - 1) do
         indptr[i] \leftarrow ptr
if mapping[i] \neq -1 then
             ptr \leftarrow ptr + 1
         end if
19: end for
20: indptr[-1] \leftarrow ptr
```

Algorithm 9 Compute the CSR indices array for the transposed Jacobian of max-pooling.

Input: mapping computed from Algorithm 8, input height h_i , input width w_i , output height h_o , output width w_o , output channels c_o

```
Output: indices \leftarrow malloc(c_oh_ow_o)

1: indices_ptr \leftarrow 0

2: for i \leftarrow 0 to (c_oh_iw_i-1) do

3: if mapping[i]=-1 then

4: continue

5: end if

6: indices[indices_ptr] \leftarrow mapping[i]

7: indices_ptr \leftarrow indices_ptr + 1

8: end for
```

Algorithm 10 Compute the CSR data array for the transposed Jacobian of max-pooling.

```
\begin{array}{ll} \textbf{Input:} \  \, \text{output channels} \  \, c_o, \, \text{output height} \  \, h_o, \, \text{output width} \  \, w_o \\ \textbf{Output:} \  \, \text{data} \leftarrow \text{malloc}(c_oh_ow_o) \\ \text{1:} \  \, \textbf{for all} \  \, i \leftarrow 0 \  \, \text{to} \  \, (c_oh_ow_o-1) \  \, \textbf{do in parallel} \\ \text{2:} \quad \quad \text{data} \leftarrow 1 \\ \text{3:} \  \, \textbf{end for} \end{array}
```

E OVERHEAD ANALYSIS OF THE GRU END-TO-END BENCHMARK

We can rewrite the GPU in Equation 10 into the following form:

$$\begin{split} \vec{R}_t &= W_{ir} \vec{x}_t + \vec{b}_{ir} + W_{hr} \vec{h}_{t-1} + \vec{b}_{hr} \\ \vec{Z}_t &= W_{iz} \vec{x}_t + \vec{b}_{iz} + W_{hz} \vec{h}_{t-1} + \vec{b}_{hz} \\ \vec{M}_t &= W_{hn} \vec{h}_{t-1} + \vec{b}_{hn}, \quad \vec{N}_t = W_{in} \vec{x}_t + \vec{b}_{in} + \vec{r}_t \circ \vec{M}_t \\ \vec{r}_t &= \sigma(\vec{R}_t), \quad \vec{z}_t = \sigma(\vec{Z}_t), \quad \vec{n}_t = tanh(\vec{N}_t) \\ \vec{h}_t &= (1 - \vec{z}_t) \circ \vec{n}_t + \vec{z}_t \circ \vec{h}_{t-1} \end{split}$$

Given the GRU expressed in the above form, the transposed Jacobian between consecutive hidden states $\frac{\partial \vec{h}_t}{\partial \vec{h}_{t-1}}$ can be computed analytically:

$$J_{1} = \left(\frac{\partial \vec{R}_{t}}{\partial \vec{h}_{t-1}}\right)^{T} = W_{hr}^{T}$$

$$\vec{j}_{2} = Diag\left(\left(\frac{\partial \vec{r}_{t}}{\partial \vec{R}_{t}}\right)^{T}\right) = \vec{r}_{t} \circ (1 - \vec{r}_{t})$$

$$\vec{j}_{3} = Diag\left(\left(\frac{\partial \vec{N}_{t}}{\partial \vec{r}_{t}}\right)^{T}\right) = \vec{M}_{t}$$

$$J_{4} = \left(\frac{\partial \vec{M}_{t}}{\partial \vec{h}_{t-1}}\right)^{T} = W_{hn}^{T}$$

$$\vec{j}_{5} = Diag\left(\left(\frac{\partial \vec{N}_{t}}{\partial \vec{M}_{t}}\right)^{T}\right) = \vec{r}_{t}$$

$$\vec{j}_{6} = Diag\left(\left(\frac{\partial \vec{N}_{t}}{\partial \vec{N}_{t}}\right)^{T}\right) = 1 - \vec{n}_{t} \circ \vec{n}_{t}$$

$$\vec{j}_{7} = Diag\left(\left(\frac{\partial \vec{h}_{t}}{\partial \vec{n}_{t}}\right)^{T}\right) = 1 - \vec{z}_{t}$$

$$J_{8} = \left(\frac{\partial \vec{Z}_{t}}{\partial \vec{h}_{t-1}}\right)^{T} = W_{hz}^{T}$$

$$\vec{j}_{9} = Diag\left(\left(\frac{\partial \vec{z}_{t}}{\partial \vec{Z}_{t}}\right)^{T}\right) = \vec{z}_{t} \circ (1 - \vec{z}_{t})$$

$$\vec{j}_{10} = Diag\left(\left(\frac{\partial \vec{h}_{t}}{\partial \vec{Z}_{t}}\right)^{T}\right) = \vec{h}_{t-1} - \vec{n}_{t}$$

$$J_{11} = \left(\frac{\partial \vec{h}_{t}}{\partial \vec{h}_{t-1}}\right)^{T}_{\text{direct}} = I \circ \vec{z}$$

$$\frac{\partial \vec{h}_{t}}{\partial \vec{h}_{t-1}} = (J_{1} \circ (\vec{j}_{2} \circ \vec{j}_{3})^{T} + J_{4} \circ \vec{j}_{5}^{T}) \circ (\vec{j}_{6} \circ \vec{j}_{7})^{T} + J_{8} \circ (\vec{j}_{9} \circ \vec{j}_{10})^{T} + J_{11}$$

where Diag(.) represents taking the diagonal of a square matrix, and \circ represents the broadcasting element-wise (Hadamard) product. Since cuDNN's GRU implementation (Appleyard et al., 2016) is closed source, we cannot access the values of the gates $(\vec{r_t}, \vec{z_t}, \vec{n_t})$. Therefore, we

have to recompute the gates (but in a more parallelized way) during the forward pass in order to compute $\frac{\partial \vec{h}_{t-1}}{\partial \vec{h}_{t-1}}$ as shown in Equations 11. This engineering challenge results in significant *overhead* during the forward pass in our experiments, however, can potentially be resolved if cuDNN's source code were publicly available and modifiable.

F GRU TRAINING CURVE

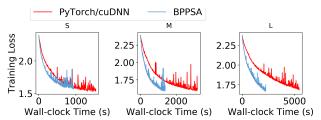


Figure 13: Training loss across wall-clock time when the GRU is trained via BPPSA (blue curve) and the PyTorch Autograd baseline with cuDNN's RNN backend (red curve).

Figure 13 shows the training curves of loss values with respect to wall-clock time when we train the GRU with the (S, M, L) preprocessed datasets for 400 epochs on the RTX 2070 GPU when the mini-batch size B is 16. We observe that the blue curve (BPPSA), if horizontally-scaled, maintains a similar shape as the red curve (PyTorch/cuDNN baseline), which reinforces our observation in Section 5.1 that BPPSA reconstructs the original back-propagation algorithm but achieves a shorter training time.

G ADDITIONAL HARDWARE SENSITIVITY RESULTS

To further validate BPPSA on a different GPU architecture, we repeat the experiments in Section 4.1 and Section 4.2 on an NVIDIA V100 (Volta architecture) (NVIDIA, 2019d) through an AWS p3.2xlarge instance (Amazon, 2019) with the same software stack as our RTX 2080Ti platform (in Table 2). The results are summarized in:

- Table 4 and Table 5 for the RNN end-to-end benchmark (Section 4.1). We can derive the backward pass runtime of the *baseline* by subtracting the "Forward Pass Only" column from the "Baseline" column for each row (T or 1/B); while we can also derive the backward pass runtime of *our* method by subtracting the "Forward Pass Only" column from the "BPPSA" column for each row (T or 1/B).
- Table 6 for the GRU end-to-end benchmark (Section 4.2). We can derive the backward pass runtime of the *baseline* by subtracting the "FP" row from the "Baseline" row for each column (batch size); while we can also derive the backward pass runtime of *our*

method by subtracting the "FP + Overhead" row from the "BPPSA" row for each column (batch size).

From the aforementioned results, we can observe trends in both benchmarks similar to the ones in the results from RTX2070 and RTX2080Ti (Section 5.1 and Section 5.2).

Table 4: The wall-clock time (s) for running a single epoch of the RNN end-to-end benchmark (Section 4.1) as the sequence length T increases.

$\begin{array}{c} \text{Sequence} \\ \text{Length} \ (T) \end{array}$	Forward Pass Only	Baseline	BPPSA
10 30 100 300 1000 3000 10000 30000	$\begin{array}{c} 1.57 \pm 0.01 \\ 2.13 \pm 0.01 \\ 3.82 \pm 0.02 \\ 8.79 \pm 0.03 \\ 25.86 \pm 0.12 \\ 75.33 \pm 0.36 \\ 250.28 \pm 0.32 \\ 748.99 \pm 0.71 \end{array}$	$\begin{array}{c} 4.49 \pm 0.04 \\ 5.5 \pm 0.06 \\ 8.87 \pm 0.05 \\ 18.49 \pm 0.1 \\ 53.33 \pm 0.39 \\ 157.11 \pm 0.58 \\ 510.13 \pm 1.11 \\ 1557.94 \pm 6.53 \end{array}$	$\begin{array}{c} 4.24 \pm 0.04 \\ 4.91 \pm 0.06 \\ 6.64 \pm 0.08 \\ 11.71 \pm 0.07 \\ 29.29 \pm 0.18 \\ 79.48 \pm 0.45 \\ 265.02 \pm 0.64 \\ 764.31 \pm 0.63 \end{array}$

Table 5: The wall-clock time (s) for running a single epoch of the RNN end-to-end benchmark (Section 4.1) as the fraction of GPU per sample (1/B) increases.

Fraction of GPU per Sample $(1/B)$	Forward Pass Only	Baseline	BPPSA
1/256	2.13 ± 0.02	4.02 ± 0.06	3.72 ± 0.02
1/128	3.87 ± 0.02	7.96 ± 0.11	5.49 ± 0.01
1/64	7.51 ± 0.06	14.77 ± 0.2	9.27 ± 0.02
1/32	13.62 ± 0.11	29.51 ± 0.34	15.0 ± 0.12
1/16	25.86 ± 0.12	53.33 ± 0.39	29.29 ± 0.18
1/8	51.29 ± 0.18	102.77 ± 0.61	56.37 ± 0.35
1/4	100.55 ± 0.25	209.67 ± 0.66	112.23 ± 0.47
1/2	200.65 ± 0.24	409.33 ± 0.67	227.13 ± 0.74

Table 6: The wall-clock time (s) of running one epoch in the GRU end-to-end benchmark (Section 4.2) as the dataset type (S, M, L) and the batch size B varies. "FP" represents running the forward pass only; "FP + Overhead" represents running the forward pass with GRU's Jacobian generation overhead; "Baseline" represents training normally via the BP baseline; "BPPSA" represents training via our method.

		Batch Size $B = 16$	Batch Size $B = 32$	Batch Size $B = 64$
S	FP Baseline FP + Overhead BPPSA	$1.66 \pm 0.01 3.71 \pm 0.02 2.1 \pm 0.01 2.94 \pm 0.02$	$ \begin{vmatrix} 0.91 \pm 0.01 \\ 1.9 \pm 0.01 \\ 1.09 \pm 0.01 \\ 1.48 \pm 0.01 \end{vmatrix} $	$ \begin{vmatrix} 0.52 \pm 0.0 \\ 1.0 \pm 0.0 \\ 0.62 \pm 0.0 \\ 0.82 \pm 0.01 \end{vmatrix} $
М	FP Baseline FP + Overhead BPPSA	$\begin{array}{c} 3.21 \pm 0.01 \\ 6.19 \pm 0.04 \\ 3.52 \pm 0.02 \\ 4.2 \pm 0.03 \end{array}$	$ \begin{vmatrix} 1.73 \pm 0.01 \\ 3.3 \pm 0.02 \\ 1.83 \pm 0.01 \\ 2.22 \pm 0.01 \end{vmatrix} $	$ \begin{vmatrix} 0.86 \pm 0.01 \\ 1.82 \pm 0.02 \\ 1.18 \pm 0.01 \\ 1.2 \pm 0.01 \end{vmatrix} $
L	FP Baseline FP + Overhead BPPSA		$ \begin{vmatrix} 3.04 \pm 0.02 \\ 6.09 \pm 0.07 \\ 3.23 \pm 0.02 \\ 3.57 \pm 0.04 \end{vmatrix} $	$ \begin{vmatrix} 1.64 \pm 0.01 \\ 3.19 \pm 0.04 \\ 2.13 \pm 0.01 \\ 2.33 \pm 0.01 \end{vmatrix} $