

Scaling Back-Propagation by Parallel Scan Algorithm

Shang Wang^{1,2}, Yifan Bai¹, Gennady Pekhimenko^{1,2}





The original PPTX file can be downloaded from <u>here</u>.

Executive Summary

The **back-propagation (BP)** algorithm is **popularly used** in training deep learning (DL) models and **implemented in many** DL frameworks (e.g., PyTorch and TensorFlow).

Problem: BP imposes a **strong sequential dependency** along layers during the gradient computations.



Executive Summary

The **back-propagation (BP)** algorithm is **popularly used** in training deep learning (DL) models and **implemented in many** DL frameworks (e.g., PyTorch and TensorFlow).

Problem: BP imposes a **strong sequential dependency** along layers during the gradient computations.

Key idea: We propose scaling BP by Parallel Scan Algorithm (BPPSA):

• Reformulate BP into a scan operation.



Executive Summary

The **back-propagation (BP)** algorithm is **popularly used** in training deep learning (DL) models and **implemented in many** DL frameworks (e.g., PyTorch and TensorFlow).

Problem: BP imposes a **strong sequential dependency** along layers during the gradient computations.

Key idea: We propose scaling BP by Parallel Scan Algorithm (BPPSA):

- Reformulate BP into a scan operation.
- Scaled by a customized parallel algorithm.

<u>Key Results</u>: O(log n) vs. O(n) steps on parallel systems.

Up to $108 \times$ backward pass speedup ($\rightarrow 2.17 \times$ overall speedup).

Back-propagation¹ (BP) Everywhere





¹Rumelhart et al. "Learning representations by back-propagating gerrors.", Nature (1986)

BP's Strong Sequential Dependency



BP's Strong Sequential Dependency



Strong Sequential Dependency along layers.

Data Parallel Training

Respects BP's strong sequential dependency.

Conceptually simple, widely used.

Effectively increases the batch size:

- Generalization gap¹
- Batch size scaling limit²

<u>Constraint:</u> The model **must** fit in one device.



5

¹Keskar, Nitish Shirish et al. "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima." ICLR (2017) ²Shallue, Christopher J. et al. "Measuring the Effects of Data Parallelism on Neural Network Training." Journal of Machine Learning Research 20 (2019)

Model Parallel Training

Used when the model cannot fit in one device.

BP's strong sequential dependency limits scalability.

Prior works on **pipeline parallel training**^{1,2} to mitigate such problem, but have their own limitations:

- Linear per-device space complexity.
- Trade-off between "bubble of idleness" vs. potential convergence affect.



¹Harlap, Aaron et al. "PipeDream: Fast and Efficient Pipeline Parallel DNN Training." SOSP (2019) ²Huang, Yanping et al. "GPipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism." NeurIPS (2019)

Rethinking BP from an Algorithm Perspective



Rethinking BP from an Algorithm Perspective

• Problems with strong sequential dependency were studied in the past (80'), but in a much simpler context.

- We propose scaling Back-Propagation by Parallel Scan Algorithm (BPPSA):
 - Reformulate BP as a scan operation.
 - Scale BP by a **customized Blelloch Scan** algorithm.
 - Leverage **sparsity** in the Jacobians.



What is a Scan¹ Operation?



Compute partial reductions at each step of the sequence.

¹Blelloch, Guy E. "Prefix sums and their applications". Technical Report (1990)

What is a Scan¹ Operation?



Compute partial reductions at each step of the sequence.

¹Blelloch, Guy E. "Prefix sums and their applications". Technical Report (1990)

Linear Scan

<u>Step:</u> executing the operator once.

Number of Elements (n)

<u>Worker (**p**):</u> an instance of execution; e.g., a core in a multi-core CPU

<u>On a single worker:</u> perform scan linearly; takes **n** steps.

With more workers: Can we achieve sublinear steps?



Blelloch Scan: 1) Up-sweep Phase



Compute partial sums via a **reduction tree**.

Blelloch Scan: 2 Down-sweep Phase



Blelloch Scan: Efficiency



Reformulate BP as a Scan Operation



Reformulate BP as a Scan Operation



Logarithmicsteps along thecritical path!



Logarithmicsteps along thecritical path!



Logarithmic steps along the critical path!

Down-sweep





Logarithmic steps along the critical path! Down-sweep





Reconstructs the Original BP Exactly

Our method produces gradients **mathematically equivalent** to BP. The Jacobians are multiplied in a different order \rightarrow numerical differences. Empirically show that such differences do not effect convergence.

Training LeNet-5 on CIFAR-10 (baseline: PyTorch Autograd)



A full Jacobian can be prohibitively expensive to handle.

• e.g., 1st convolution in VGG-11 on CIFAR-10 images occupy **768 MB** of memory.



- e.g., 1st convolution in VGG-11 on CIFAR-10 images occupy **768 MB** of memory.
- Generated one row at a time by passing basis vectors into Op_Grad() (the VJP function).



- e.g., 1st convolution in VGG-11 on CIFAR-10 images occupy 768 MB of memory.
- Generated one row at a time by passing basis vectors into Op_Grad() (the VJP function).



- e.g., 1st convolution in VGG-11 on CIFAR-10 images occupy **768 MB** of memory.
- Generated one row at a time by passing basis vectors into Op_Grad() (the VJP function).



- e.g., 1st convolution in VGG-11 on CIFAR-10 images occupy **768 MB** of memory.
- Generated one row at a time by passing basis vectors into Op_Grad() (the VJP function).



A full Jacobian can be prohibitively expensive to handle.

- e.g., 1st convolution in VGG-11 on CIFAR-10 images occupy **768 MB** of memory.
- Generated one row at a time by passing basis vectors into Op_Grad() (the VJP function).

Conventional ML algorithms avoid using Jacobians directly (including BP).



The Jacobians of Many Operators are Sparse



Fast Sparse Jacobians Generation

Therefore, instead of calculating the Jacobians row-wise, generate **directly** into **Compressed Sparse Row (CSR)**:



| First three ops of VGG-11 on CIFAR-10 | Convolution | ReLU | Max Pooling |
|---------------------------------------|-----------------------|-----------------------|-----------------------|
| Jacobian Calculation Speedup | 8.3×10 ³ x | 1.2×10 ⁶ x | 1.5×10 ⁵ x |

| F | Runtime | • |
|-----------------------------|---------|----------------------|
| BPPSA | | BP |
| С _{вррза} Ө(log n) | VS. | C _{BP} Θ(n) |

Per-step Complexity (C): runtime of each step.



<u>Per-step Complexity (C)</u>: runtime of each step.

Performance benefits:

1. Large **n**: deep network, long sequential dependency.



Performance benefits:

- 1. Large **n**: deep network, long sequential dependency.
- 2. Reducing per-step complexity: SpGEMM.

Per-step Complexity (C): runtime of each step.



Performance benefits:

- 1. Large **n**: deep network, long sequential dependency.
- 2. Reducing per-step complexity: SpGEMM.

Constant per-device space complexity!

Per-step Complexity (C): runtime of each step.



Methodology: Benchmark

Model: RNN <u>Task:</u> Bitstream Classification

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



Methodology: Environment



Implementation: custom CUDA 10 kernels.

End-to-end Training Speedup

Training curve of BPPSA v.s. the baseline when batch size **B**=16, sequence length **T**=1000:



Numerical differences do **not** effect convergence.

End-to-end Training Speedup

Training curve of BPPSA v.s. the baseline when batch size **B**=16, sequence length **T**=1000:



Sensitivity Analysis: Model Length



Sensitivity Analysis: Model Length



Sensitivity Analysis: Model Length



Sequence length (T) reflects the model length **n**.

BPPSA **scales** with the model length (**n**);

until being bounded by the number of workers (**p**).









Sensitivity Analysis: 2070 v.s. 2080Ti

#SMs(2070) < #SMs(2080Ti) → Latency(2070) > Latency(2080Ti)

■ 2070 ■ 2080Ti

<u>SM</u>: Streaming Multiprocessor; i.e., "Parallel Cores".



More Results in the Paper

- End-to-end benchmarks of GRU training on IRMAS.
 - A more realistic version of the RNN results.
- Pruned VGG-11 retraining on CIFAR-10.
 - Microbenchmark via FLOP measurements.
 - Evaluate the effectiveness of leveraging the Jacobians' sparsity in CNNs.

Conclusion

BP imposes a **strong sequential dependency** among layers during the gradient computations, limiting its scalability on parallel systems.

We propose scaling Back-Propagation by Parallel Scan Algorithm (BPPSA):

- Reformulate BP as a scan operation.
- Scale by a **customized Blelloch scan** algorithm.
- Leverage **sparsity** in the Jacobians.

<u>Key Results:</u> Θ(log n) vs. Θ(n) steps on parallel systems.

Up to $108 \times$ speedup on the backward pass ($\rightarrow 2.17 \times$ overall speedup).