Project Fiddle Fast & Efficient Infrastructure for Distributed Deep Learning Amar Phanishayee

with

Nikhil Devanur, *Aaron Harlap, Animesh Jain, Liang Luo, Deepak Narayanan*, Jacob Nelson, Gennady Pekhimenko, Vivek Seshadri, Jorgen Thelin, Shivaram Venkataraman, *Guanhua Wang*

Big Data Systems – Gen Al

• Motivating Scenarios:



Image Recognition, Classification Translation & Speech – Text - Speech

• Solutions:



Deep Learning for X

- Image recognition, classification
- Segmentation
- Captioning
- Stereo
- Depth estimation
- Video action recognition
- Sentiment Analysis
- Speech-Text, Text-Speech, OCR
- Translation
- •

Training is time and resource intensive



- Long running training jobs: Days to multiple weeks
- Train, Validate, Repeat Performance and accuracy specs validated only by running system
- Naïvely parallelizing work can be detrimental
- No luxury of limiting solution to one part of the stack

Need for speed

From EE Times – September 27, 2016

"Today the job of training machine learning models is limited by compute, if we had faster processors we'd run bigger models...in practice we train on a reasonable subset of data that can finish in a matter of months. We could use improvements of several orders of magnitude – 100x or greater."

- Greg Diamos, Senior Researcher, SVAIL, Baidu

Slide Credit: Joel Emer (MIT, NVIDIA)

Problem

Systematically speed-up distributed learning tasks while eking out the most from the resources used

Goal: 100x more efficient training

Day long training of model in 15-min coffee break

Approach

Cut through layers of the system stack and optimize all relevant parts

C	Single Machine Training	Multi-Machine Training		
utatio	Gist	Pi	ipeDream	
Memory, Compi	Train larger networks on a <i>single</i> GPU by reducing memory footprint.	Efficiently scale training New way to parallelize DNN computation.		
	Blink		Hub	
Interconnects	Blazingly fast transfe over PCIe + NVLink -	Fast parameter servers		

ently scale training. way to parallelize IN computation. Hub Fast parameter servers

Single Machine Training	Multi-Machine Training		
Gist Pi		peDream	
Train larger networks on a <i>single</i> GPU by reducing memory footprint. Up to 2x compression ratio	Efficiently scale training. New way to parallelize DNN computation. 3x faster training		
Blink	Hub		
Blazingly fast transfe over PCIe + NVLink +	Fast parameter servers		
Collective transfer prim 8x faster than NCC	Up to 3x speedup		

Memory, Computation

Interconnects

_	Single Machine Training	Multi-Machine Training		
utatio	Gist			
Memory, Comp	Train larger networks on a <i>single</i> GPU by reducing memory footprint. Up to 2x compression ratio			
Interconnects				







Example – AlexNet (Image Classification)



Slide Credit: http://vision03.csail.mit.edu/cnn_art/

Tasks, datasets, challenges

- 1989: MNIST
- 2001: Caltech 101
- 2007 2012: PASCAL Visual Object Classes
 - 20 classes
- 2014 ... : Imagenet 1K
 - Deng et al.
 - Classification and Detection

Profile of memory consumption



Feature Maps are a major consumer of GPU memory CNTK Profiling

Larger minibatch size → potential crash



Profile of memory consumption



Feature Maps are a major consumer of GPU memory CNTK Profiling

Larger minibatch size → potential crash

Heavy hitters: Relu->Pool, Relu/Pool-> Conv

Idea: Encode data structures when unused, Decode on use



Gist Architecture



- Statically construct usage schedule for all data structures
- Layer-specific encodings
 - LOSSIESS (Binarize, Sparse Storage/Dense Compute)
 - LOSSY (Delayed Precision Reduction)

Gist Lossless Encodings For Relu->Pool



Binarize 1 bit representation of 32bit values in Y

32x reduction

For Relu/Pool->Conv



dX = f(X, dY)

Naively applying Binarize for Relu/Pool followed by Conv can increase memory consumption!

Gist Lossless Encodings For Relu->Pool



Relu Backward Propagation

Binarize

1 bit representation of 32bit values in Y

For Relu/Pool->Conv





(sparse compute is computationally expensive)

Gist Lossy Encodings

- Used for all other feature maps
- Training with reduced precision (8/10/16 bits) done only when encoding/stashing values
- Forward pass uses full fidelity values



16 bits works out.

10 bits enough!

Gist – Putting it all together



Up to 2x compression ratio



Minimal runtime overhead (1 - 7%) for same batch size

With just lossless, we go faster at times (memory bandwidth bound)

	Single Machine Training	Multi-Machine Training
urario	Fiddle Gist	
INIEMOLY, COMPI	Train larger networks on a <i>single</i> GPU by reducing memory footprint. Up to 2x compression ratio	
Interconnects		

Distributed Deep Learning

- Data Parallelism
 - Replicas run in parallel on different machines
 - Occasionally exchange what they have learnt
 - Naïve setup can hurt convergence, accuracy
 - Total time = (Time per Epoch) * (#Epochs for given accuracy)



Hardware Efficiency

Statistical Efficiency



Manually tuned for performance of individual jobs

Pipeline Parallelism in Fiddle

- Idea: Pipeline computation across machines
 - For training tasks, optimize for throughput
 - Each machine runs a subset of layers
 - Compute one thing, data flows to compute task
 - Better FLOPs due to cache locality
 - Fits CPU, GPU, hybrid, heterogeneous clusters





- No bubbles in pipeline
 - How should we divvy work
 across machines?



- No bubbles in pipeline
 - How should we divvy work
 across machines?

Running at speed of the slowest layer

Idle Cycles



- No bubbles in pipeline
 - How should we divvy work across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent



- No bubbles in pipeline
 - How should we divvy work across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?



- No bubbles in pipeline
 - How should we divvy work across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?





- No bubbles in pipeline
 - How should we divvy work
 across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?
- In steady state:

Each new item learns from a previous (new) mini-batch



- No bubbles in pipeline
 - How should we divvy work
 across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?
- In steady state:

Each new item learns from a previous (new) mini-batch



- No bubbles in pipeline
 - How should we divvy work across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?
- In steady state:
- Each new item learns from a previous (new) mini-batch
 - Stage alternates between F/B



- No bubbles in pipeline
 - How should we divvy work
 across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?
- In steady state:

Each new item learns from a previous (new) mini-batch



- No bubbles in pipeline
 - How should we divvy work
 across machines?
 - Replicate stages (if needed)
 - Static load balancing
 - mini-batch-id % stage-replica-id
 - F/B paths consistent
 - How many items to admit?
 - Wild, S, less?
- In steady state:

Each new item learns from a previous (new) mini-batch



- No bubbles in pipeline
 - How should we divvy work
 across machines?
 - Replicate stages (if needed)
 - Static load balancing

mini-batch-id % stage-replica-id

- F/B paths consistent
- How many items to admit?
 - Wild, S, less?
- In **steady** state:

Each new item learns from a previous (new) mini-batch

- Statistical Efficiency
 - Affected by feature map and model parameter versions
 - Each incomplete mini-batch needs to stash feature maps between F and B passes
 - Should admitted items see the updates of the same set of minibatches across levels (not newer ones)?



Valid Gradients Or Bust

$$w^{(t+1)} = w^{(t)} - \mu * \nabla f(w_1^{(t)}, w_2^{(t)}, \dots, w_n^{(t)})$$

PipeDream with wt stashing is critical for valid gradients (use same wts across F and B *within* every stage)



Valid Gradients Or Bust

$$w^{(t+1)} = w^{(t)} - \mu * \nabla f(w_1^{(t)}, w_2^{(t)}, \dots, w_n^{(t)})$$

PipeDream with **wt stashing** is critical for valid gradients (use same wts across F and B *within* every stage)

$$w^{(t+1)} = w^{(t)} - \mu * \nabla f(w_1^{(t-n+1)}, w_2^{(t-n+2)}, \dots, w_n^{(t)})$$

PipeDream with **vertical sync** maps to bounded staleness (use same wts across all stages) $w^{(t+1)} = w^{(t)} - \mu * \nabla f(w_1^{(t-n+1)}, w_2^{(t-n+1)}, \dots, w_n^{(t-n+1)})$

- No bubbles in pipeline
 - How should we divvy work
 across machines?
 - Replicate stages (if needed)
 - Static load balancing

mini-batch-id % stage-replica-id

- F/B paths consistent
- How many items to admit?
 - Wild, S, less?
- In **steady** state:

Each new item learns from a previous (new) mini-batch

Stage alternates between F/B

- Statistical Efficiency
 - Affected by feature map (FM) and model parameter versions
 - Each incomplete mini-batch needs to stash feature maps between F and B passes
 - Should admitted items see the updates of the same set of minibatches across levels (not newer ones)?

Memory Manager

- Stashing: Static memory buffer pool (FMs, wts, wt-updates)
- Use same wts across F/B within every stage (valid gradient)
- Vertical Sync

Stage Runtime Architecture



Visualization / Debugging

Record Save L	.oad eg_trace (dp).html					View Options +	← → » ?
		-500 µs	0 µs	1 1 I	500 µs	1,000 µs	1,500 µs
 Process 0 							· · · · · · · · · · · · · · · · · · ·
• main			FowardC	omput		PSRead	k. +
 Process 1 							
▪ main			Fox	Fow Fo Fo		PSRead	
 Process 2 							
 main 			Fow.	Fow Fo Fow Fow Fow		PSRead	
1 item selected.	Slice (1)						
Title User Friendly Catego Start Wall Duration Self Time	FowardCompute- conv4_4 batch 0 ory other 0.539 ms 0.032 ms 0.031 ms						

Encouraging results



Up to 3x faster, more efficient

Effect of smaller minibatches on hardware efficiency and stat efficiency in data parallel runs

Restructure computation to push a lot less data on the network



Single Machine Training	Multi-Machine Training		
Gist Pi		peDream	
Train larger networks on a <i>single</i> GPU by reducing memory footprint. Up to 2x compression ratio	Efficiently scale training. New way to parallelize DNN computation. 3x faster training		
Blink		Hub	
Blazingly fast transfe over PCIe + NVLink -	Fast parameter servers		
Collective transfer prim 8x faster than NCC	Up to 3x speedup		

Interconnects

Memory, Computation



Fin

