EHzürich

T. HOEFLER

Demystifying Parallel and Distributed Deep Learning - An In-Depth Concurrency Analysis

Keynote at the 6th Accelerated Data Analytics and Computing Workshop (ADAC'19



https://www.arxiv.org/abs/1802.09941

Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis

TAL BEN-NUN* and TORSTEN HOEFLER, ETH Zurich

Deep Neural Networks (DNNs) are becoming an important tool in modern computing applications. Accelerating their training is a major challenge and techniques range from distributed algorithms to low-level circuit design. In this survey, we describe the problem from a theoretical perspective, followed by approaches for its parallelization. Specifically, we present trends in DNN architectures and the resulting implications on parallelization strategies. We discuss the different types of concurrency in DNNs; synchronous and asynchronous stochastic gradient descent; distributed system architectures; communication schemes; and performance modeling. Based on these approaches, we extrapolate potential directions for parallelism in deep learning.

 $\label{eq:ccs} Concepts: \bullet \mbox{General and reference} \rightarrow Surveys and overviews; \bullet \mbox{Computing methodologies} \rightarrow \mbox{Neural networks; Distributed computing methodologies; Parallel computing methodologies; Machine learning;}$

Additional Key Words and Phrases: Deep Learning, Distributed Computing, Parallel Algorithms

ACM Reference format:

Tal Ben-Nun and Torsten Hoefler. 2018. Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis 60 pages.

1 INTRODUCTION

Machine Learning, and in particular Deep Learning [LeCun et al. 2015], is a field that is rapidly taking over a variety of aspects in our daily lives. In the core of deep learning lies the Deep Neural Network (DNN), a construct inspired by the interconnected nature of the human brain. Trained properly, the expressiveness of DNNs provides accurate solutions for problems previously thought to be unsolvable, simply by observing large amounts of data. Deep learning has been successfully implemented for a plethora of subjects, ranging from image classification [Huang et al. 2017], through speech recognition [Amodei et al. 2016] and medical diagnosis [Cireşan et al. 2013], to autonomous driving [Bojarski et al. 2016] and defeating human players in complex games [Silver et al. 2017] (see Fig. 1 for more examples).



DINFK

spcl.inf.ethz.ch







^SPEL



spelintethich ETHZÜRICH

How does Deep Learning work?







Deep Learning is Supercomputing!



0.28 0.00 Dog Dog 0.07 0.00 Airplane Airplane 0.00 0.04 Horse Horse 0.33 0.00 Bicycle Bicycle 0.02 0.00 0.02 0.00 Truck Truck

- ImageNet (1k): 180 GB
- ImageNet (22k): A few TB
- Industry: Much larger

- 100-200 layers deep
- ~100M-2B parameters
- 0.1-8 GiB parameter storage

- 10-22k labels
- growing (e.g., face recognition)
- weeks to train



A brief theory of supervised deep learning



P. Land Street Frank



Stochastic Gradient Descent	W	$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}}[\ell(w, x)]$				
1:		$f_1(x)$ convolution	on 1			
2: 3: 4:		$f_2(f_1(x))$ convolution	on 2			
5: 6: 7:		pooling	5			
8: 9:		 convolutio	on 3			
10: 11: 12:		f(x) fully conne	cted			
• Layer storage = $ w_l + f_l(o_{l-1}) + \nabla w_l + \nabla o_l $	Learning Rate $w^{(t+1)} = w$ Adaptive Learning Rate $w^{(t+1)} = w$		<i>t</i>)			

A SALAR PARTY AND A SALAR SALAR



Learning Rate	$w^{(t+1)} = w^{(t)} - \eta \cdot \nabla \ell(w^{(t)}),$	$z) \qquad = w^{(t)} - \eta \cdot \nabla w^{(t)}$				
Adaptive Learning Rate	$w^{(t+1)} = w^{(t)} - \eta_t \cdot \nabla w^{(t)}$					
Momentum [Qian 1999]	$w^{(t+1)} = w^{(t)} + \mu \cdot (w^{(t)} - w^{(t-1)}) - \eta \cdot \nabla w^{(t)}$					
Nesterov Momentum [Nesterov 1983]	$w^{(t+1)} = w^{(t)} + v_t; v_{t+1} = v_t$	$= \mu \cdot \upsilon_t - \eta \cdot abla \ell(w^{(t)} - \mu \cdot \upsilon_t, z)$				
AdaGrad [Duchi et al. 2011]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A_{i,t} + \epsilon}};$	$A_{i,t} = \sum_{\tau=0}^{t} \left(\nabla w_i^{(t)} \right)^2$				
RMSProp [Hinton 2012]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A'_{i,t}} + \varepsilon};$	$A_{i,t}' = \beta \cdot A_{t-1}' + (1-\beta) \left(\nabla w_i^{(t)} \right)^2$				
Adam [Kingma and Ba 2015]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot M_{i,t}^{(1)}}{\sqrt{M_{i,t}^{(2)}} + \varepsilon};$	$M_{i,t}^{(m)} = \frac{\beta_m \cdot M_{i,t-1}^{(m)} + (1-\beta_m) \left(\nabla w_i^{(t)}\right)^m}{1-\beta_m^t}$				

T. Ben-Nun, T. Hoefler: Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis, arXiv Feb 2018



Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 227 papers in the area of parallel deep learning



Deep Learning is largely on distributed memory today!

Trends in distributed deep learning: node count and communication

The field is moving fast – trying everything imaginable – survey results from 227 papers in the area of parallel deep learning



Deep Learning research is converging to MPI!



Minibatch Stochastic Gradient Descent (SGD)



A CONTRACTOR OF A CONTRACTOR O

T. Ben-Nun, T. Hoefler: Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis, arXiv Feb 2018



A primer of relevant parallelism



and communication theory

Parallel Reductions for Parameter Updates

 $y = x_1 \oplus x_2 \oplus x_3 \cdots \oplus x_{n-1} \oplus x_n$



Lower bound: $T \ge L \log_2 P + 2\gamma m G(P-1)/P$

E. Chan et al.: Collective communication: theory, practice, and experience. CCPE'07 9 TH, D. Moor: Energy, Memory, and Runtime Tradeoffs for Implementing Collective Communication Operations, JSFI'14



GoogLeNet in more detail







Parallelism in the different layer types

Layer Type	Eval.	Work (W)	Depth (D)

The second and the

W is linear and D logarithmic – large average parallelism

T. Ben-Nun, T. Hoefler: Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis, arXiv Feb 2018



Computing fully connected layers





a lar an and a second state

Computing convolutional layers

		Direct	Indirect				
-	4 1 9 8	1 -1 0 21.9 59.3 53.9 43.9	FFT Winograd				
Dire	Method	Work (W)	Depth (D)				
	Direct	$N \cdot C_{out} \cdot H' \cdot W' \cdot C_{in} \cdot K_y \cdot K_x$	$\left\lceil \log_2 C_{in} \right\rceil + \left\lceil \log_2 K_y \right\rceil + \left\lceil \log_2 K_x \right\rceil$				
	im2col	$N \cdot C_{out} \cdot H' \cdot W' \cdot C_{in} \cdot K_y \cdot K_x$	$\left\lceil \log_2 C_{in} \right\rceil + \left\lceil \log_2 K_y \right\rceil + \left\lceil \log_2 K_x \right\rceil$				
im2	FFT	$c \cdot HW \log_2(HW) \cdot (C_{out} \cdot C_{in} + N \cdot C_{in} + N \cdot C_{out}) + HWN \cdot C_{in} \cdot C_{out}$	$2\left\lceil \log_2 HW \right\rceil + \left\lceil \log_2 C_{in} \right\rceil$				
	Winograd $(m \times m \text{ tiles}, m \times m \text{ tiles})$	$\alpha(r^{2} + \alpha r + 2\alpha^{2} + \alpha m + m^{2}) + C_{out} \cdot C_{in} \cdot P$	$2\left\lceil \log_2 r \right\rceil + 4\left\lceil \log_2 \alpha \right\rceil + \left\lceil \log_2 C_{in} \right\rceil$				
	$r \times r$ kernels)	$(\alpha = m - r + 1, P = N \cdot H/m \cdot W/m)$	Activation				
		F_m O_m fficient Primitives for Deep Learning, arXiv 2014					

K. Chellapilla et al.: High Performance Convolutional Neural Networks for Document Processing, Int'l Workshop on Frontiers in Handwriting Recognition 2016 M. Mathieu et al.: Fast Training of Convolutional Networks through FFTs, ICLR'14

A. Lavin and S. Gray: Fast Algorithms for Convolutional Neural Networks, CVPR'16

Memory efficient

Microbatching (µ-cuDNN)

- Fast (up to 4.54x on DeepBench) tion implementations
- Performance depends on temporary memory (workspace) size 140 conv5 256 conv4 120 u Morkspace [MiB] 09 00 09 00 conv3 conv2 р 32 32 BF conv1 BF IMPLICIT_PRECOMP_GEMM FFT_TILING 48 48 60 а WINOGRAD_NONFUSED BD BD BF BF 0 5 20 Time [ms] 11 р 200 DP ILP etc. conv5 conv4 150 **Utilizes heterogeneous clusters** conv3 Time [ms] 100 conv2 conv1 400 750Ti 334.1 K20Xm 300 K80 50 Time [ms] 500 154.2 32 130.8 0 123 107.2 100 80.5 70.3 70.3 65.9 59.5 61.8 u p a u p а u p а 15 16 6 5 (8 MiB) (64 MiB) (512 MiB) 10 13 8 12 7 10 5 P100-SXM2 750Ti 750Ti 750Ti 750Ti K20Xm 750Ti K20Xm 750Ti K20Xm K80 K80×2 K20Xm K20Xm K80 K80×2 K80×2 K20Xm K80 K80 K80×2

Oyama et al.: µ-cuDNN: Accelerating Deep Learning Frameworks with Micro-Batching, arXiv 2018



Model parallelism



- Parameters can be distributed across processors
- Mini-batch has to be copied to all processors
- Backpropagation requires all-to-all communication every layer



Pipeline parallelism



- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors



Data parallelism



The second second

- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors

X. Zhang et al.: An Efficient Implementation of the Back-propagation Algorithm on the Connection Machine CM-2, NIPS'89



- Layers/parameters can be distributed across processors
- Can distribute minibatch
- Often specific to layer-types (e.g., distribute fc layers but handle conv layers data-parallel)
 - Enables arbitrary combinations of data, model, and pipeline parallelism very powerful!
- A. Krizhevsky: One weird trick for parallelizing convolutional neural networks, arXiv 2014
- J. Dean et al.: Large scale distributed deep networks, NIPS'12.
- T. Ben-Nun, T. Hoefler: Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis, arXiv Feb 2018



Updating parameters in distributed data parallelism



and the second second



Parameter (and Model) consistency - centralized

Parameter exchange frequency can be controlled, while still attaining convergence:







- Started with Hogwild! [Niu et al. 2011] shared memory, by chance
- DistBelief [Dean et al. 2012] moved the idea to distributed
- Trades off "statistical performance" for "hardware performance"

J. Dean et al.: Large scale distributed deep networks, NIPS'12.

F. Niu et al.: Hogwild: A lock-free approach to parallelizing stochastic gradient descent, *NIPS*'11.



Parameter (and Model) consistency - decentralized

 Parameter exchange frequency can be controlled, while still attaining convergence:



Training Agent Training Agent Training Agent Training Agent



May also consider limited/slower distribution – gossip [Jin et al. 2016]



Parameter consistency in deep learning



Carlan - - - - - - - - -

S. Zhang et al.: Deep learning with Elastic Averaging SGD, NIPS'15



Parameter consistency in deep learning





T. G. Dietterich: Ensemble Methods in Machine Learning, MCS 2000

***SPEL

Communication optimizations

- Different options how to optimize updates
 - Send ∇w , receive w
 - Send FC factors (o_{l-1}, o_l), compute ∇w on parameter server Broadcast factors to not receive full w
 - Use lossy compression when sending, accumulate error locally!
- Quantization
 - Quantize weight updates and potentially weights
 - Main trick is stochastic rounding [1] expectation is more accurate Enables low precision (half, quarter) to become standard
 - TernGrad ternary weights [2], 1-bit SGD [3], ...
- Sparsification
 - Do not send small weight updates or only send top-k [4]
 Accumulate them locally







[1] S. Gupta et al. Deep Learning with Limited Numerical Precision, ICML'15

^[2] F. Li and B. Liu. Ternary Weight Networks, arXiv 2016

^[3] F. Seide et al. 1-Bit Stochastic Gradient Descent and Application to Data-Parallel Distributed Training of Speech DNNs, In Interspeech 2014

^[4] C. Renggli et al. SparCML: High-Performance Sparse Communication for Machine Learning, arXiv 2018



SparCML – Quantified sparse allreduce for decentral updates



System	Dataset	Model	# of nodes Algorithm	Speedup
Piz Daint	ImageNet	VGG19	8 Q4	1.55 (3.31)
Piz Daint	ImageNet	AlexNet	16 Q4	1.30 (1.36)
Piz Daint EC2	MNIST	MLP	8 Top16_Q4 Top16_Q4	3.65 (4.53) 19.12 (22.97)



MNIST test accuracy

Evolutionary Algorithms [4]

Hyperparameter and Architecture search

- Meta-optimization of hyper-parameters (momentum) and DNN architecture
 - Using Reinforcement Learning [1] (explore/exploit different configurations)
 - Genetic Algorithms with modified (specialized) mutations [2]
 - Particle Swarm Optimization [3] and other meta-heuristics



Reinforcement Learning [1]

- [1] M. Jaderberg et al.: Population Based Training of Neural Networks, arXiv 2017
- [2] E. Real et al.: Regularized Evolution for Image Classifier Architecture Search, arXiv 2018
- [3] P. R. Lorenzo et al.: Hyper-parameter Selection in Deep Neural Networks Using Parallel Particle Swarm Optimization, GECCO'17
- [4] H. Liu et al.: Hierarchical Representations for Efficient Architecture Search, ICLR'18

Application: Neural Code Comprehension

- In 2017, GitHub reports 1 billion git commits in 337 languages!
- Can DNNs understand code?
- Previous approaches read the code directly → suboptimal (loops, functions)



Ben-Nun et al.: Neural Code Comprehension: A Learnable Representation of Code Semantics, arXiv 2018

Application: Neural Code Comprehension

Embedding space (using the Skip-gram model)



Application: Neural Code Comprehension

				Ta	able 3: A	Algorit	hm classif	fication test	t accurac	y				
		Metric Si (RBF			urface Features [46] SVM + Bag-of-Trees)		RNN [46]	TBCNN	[46] inst2v	vec				
			Test A	Accuracy [%		88	.2	84.8	94.0	94.8	33			
× Pre	dicts whic	h device is	s faste	r (CPU c	or GPU)					Opt	imal tiling	ous Code E		
	Table 4: Het	terogeneous de	evice ma	pping result	ts			Table 5: Speedups achieved by coarsening threads						
Architecture	Predic	ction Accuracy [%]			Speedup			Computin	g Platform	Magni et al. [43]	DeepTune [17]	DeepTune-TL [17]	inst2vec	
	Grewe et al. [27]	DeepTune [17]	inst2vec	Grewe et al.	DeepTune	inst2vec	STM [AMD Rad	leon HD 5900	1.21	1.10	1.17	1.25	
AMD Tahiti 7970 NVIDIA GTX 970	73.38 72.94	83.68 80.29	82.79 81.76	2.91 1.26	3.34 1.41	3.42 1.39		NVIDIA (NVIDIA (GTX 480 Tesla K20c	0.86 0.94	1.05 1.10 0.99	1.23 1.14 0.93	1.07 1.02 1.03	
	tmuL 'y	~		Λŏ.		7	Shits L			///20	Code	Optimizatio	n	
	%AF 100							Table	2: Analogy	and test sco	ores for inst	2vec		
	phi O						Context	Syntactic Analogies		Seman	Semantic Analogies		Semantic Distance Test	
	50	199		Ni O			Size	Types	Options	Conversions	Data Structure	es		
	8 0						1 2 3	101 (18.04%) 226 (40.36%) 125 (22.32%)	13 (24.53%) 45 (84.91%) 24 (45.28%)	100 (6.63%) 134 (8.89%) 48 (3.18%)	3 (37.50% 7 (87.50% 7 (87.50%	(6) (6) (6)	60.98% 79.12% 62.56%	
	-50	-50 0	50											

2 Carlos and

Ben-Nun et al.: Neural Code Comprehension: A Learnable Representation of Code Semantics, arXiv 2018



Outlook

- Full details in the survey (60 pages)
 - Parallelism, distribution, synchronization
- Additional content:
 - Unsupervised (GAN/autoencoders)
 - Recurrent (RNN/LSTM)
- Call to action to the HPC and ML/DL communities to join forces!
 - It's already happening on the tool basis
 - Need more joint events!

Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis

https://www.arxiv.org/abs/1802.09941

TAL BEN-NUN* and TORSTEN HOEFLER, ETH Zurich

Deep Neural Networks (DNNs) are becoming an important tool in modern computing applications. Accelerating their training is a major challenge and techniques range from distributed algorithms to low-level circuit design. In this survey, we describe the problem from a theoretical perspective, followed by approaches for its parallelization. Specifically, we present trends in DNN architectures and the resulting implications on parallelization strategies. We discuss the different types of concurrency in DNNs; synchronous and asynchronous stochastic gradient descent; distributed system architectures; communication schemes; and performance modeling. Based on these approaches, we extrapolate potential directions for parallelism in deep learning.

CCS Concepts: • General and reference \rightarrow Surveys and overviews; • Computing methodologies \rightarrow Neural networks; Distributed computing methodologies; Parallel computing methodologies; Machine learning;

Additional Key Words and Phrases: Deep Learning, Distributed Computing, Parallel Algorithms

ACM Reference format:

Tal Ben-Nun and Torsten Hoefler. 2018. Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis. 60 pages.

1 INTRODUCTION

Machine Learning, and in particular Deep Learning [LeCun et al. 2015], is a field that is rapidly taking over a variety of aspects in our daily lives. In the core of deep learning lies the Deep Neural Network (DNN), a construct inspired by the interconnected nature of the human brain. Trained properly, the expressiveness of DNNs provides accurate solutions for problems previously thought to be unsolvable, simply by observing large amounts of data. Deep learning has been successfully implemented for a plethora of subjects, ranging from image classification [Huang et al. 2017], through speech recognition [Amodei et al. 2016] and medical diagnosis [Cireşan et al. 2013], to autonomous driving [Bojarski et al. 2016] and defeating human players in complex games [Silver et al. 2017] (see Fig. 1 for more examples).

1