二次最適化を用いた巨大な言語モデルの学習 およびFRNNを用いたプラズマ挙動予測

東京工業大学学術国際情報センター

横田理央

ABCI グランドチャレンジ 2019 成果発表会 テレコムセンタービル東棟14階 2020年2月21日



スパコンでしかできない深層学習



Figures: M. Tan and Q. V. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks", ICML 2019



Google TPU v3 12.5PF/rack



AIST ABCI 17 PF/rack

数千GPU規模のImageNetの学習



数千GPU規模のImageNetの学習



二次最適化

一次最適化(確率的勾配降下法:SGD)

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)})$$

二次最適化

ニュートン法

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \left\{ \frac{H(\theta^{(t)})}{\sum_{n \neq t \in \mathcal{D}}} \right\}^{-1} \nabla \mathcal{L}(\theta^{(t)})$$

 自然勾配法 (S. Amari, 1998)
 $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \left\{ \frac{F(\theta^{(t)})}{\sum_{FIM}} \right\}^{-1} \nabla \mathcal{L}(\theta^{(t)})$

深層学習においては・・・

$$H(\theta^{(t)}), F(\theta^{(t)}) \in \mathbb{R}^{N \times N}$$
 × メモリに載らない
× 逆行列(O(N^3))計算は非現実的

クロネッカー因子分解(K-FAC)

Kronecker-Factored Approximate Curvature (J. Martens+, ICML 2015)



Step 2. Kronecker-factorization (for each layer)





much easier to invert $O(N_\ell^3) o O(N_\ell^{3/2})$

Distributed K-FAC (K. Osawa+, CVPR 2019)

 ∇

Local

Mini-batch

0000

GPU4

data-parallelism

Mini-batch

Local

Mini-batch

QOO(

600

GPU3

J

Local

Mini-batch

0000

GPU2

 $\overline{\mathbf{v}}$

Local

Mini-batch

GPU1





Training ResNet-50 (107 layers) on ImageNet classification

分散並列K-FACを用いたImageNetの学習

Our Distributed K-FAC results on ABCI supercomputer with extremely large-batches for training ResNet-50 on ImageNet (CVPR19)



The first empirical results showing <u>the advantage of 2nd-order optimization</u> (K-FAC) in Large-scale DL (faster convergence)

パフォーマンスの最適化



	Distrib. FIM inverse	Symm. comm. of FIM	float21 comm. of FIM	Overlap comm.	Hier. Comm.	Stale FIM	SGD Time/upd.	K-FAC Time/upd.	K-FAC/SGD ratio
	Sec. 3.3.1	Sec. 3.3.2	Sec. 3.3.3	Sec. 4.1	Sec. 4.1	Sec. 4.2			
Osawa <i>et al.</i> [19]	\checkmark	\checkmark				\checkmark^1	-	340 ms	-
Tsuji <i>et al.</i> [23]	\checkmark	\checkmark		\checkmark	\checkmark		85 ms	315 ms	3.71
Osawa <i>et al.</i> [18]	\checkmark	\checkmark		\checkmark	\checkmark		-	236 ms	-
	\checkmark	\checkmark		\checkmark	\checkmark	\checkmark	-	178 ms	-
This work	\checkmark							251 ms	4.40
	\checkmark	\checkmark						199 ms	3.49
	\checkmark	\checkmark	\checkmark				57 ms	182 ms	3.19
	\checkmark	\checkmark	\checkmark	\checkmark			37 1115	221 ms	3.88
	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark			192 ms	3.37
	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark		108 ms	1.89

巨大な言語モデルのパラメータ数



巨大な言語モデルの学習時間

	BERT	RoBERTa	DistilBERT	XLNet
Size (millions)	Base: 110 Large: 340	Base: 110 Large: 340	Base: 66	Base : ~110 Large : ~340
Training Time	Base: 8 x V100 x 12 days* Large: 64 TPU Chips x 4 days (or 280 x V100 x 1 days*)	Large: 1024 x V100 x 1 day; 4-5 times more than BERT.	Base : 8 x V100 x 3.5 days; 4 times less than BERT.	Large: 512 TPU Chips x 2.5 days; 5 times more than BERT.
Performance	Outperforms state-of- the-art in Oct 2018	2-20% improvement over BERT	3% degradation from BERT	2-15% improvement over BERT
Data	16 GB BERT data (Books Corpus + Wikipedia). 3.3 Billion words.	<mark>160 GB</mark> (16 GB BERT data + 144 GB additional)	16 GB BERT data. 3.3 Billion words.	Base : 16 GB BERT data Large : 113 GB (16 GB BERT data + 97 GB additional). 33 Billion words.
Method	BERT (Bidirectional Transformer with MLM and NSP)	BERT without NSP**	BERT Distillation	Bidirectional Transformer with Permutation based modeling







K-FACはTransformer用のチューニングが行われていないため途中で 時間切れになったが、早く収束する傾向は観測された。

これはパイパーパラメータチューニングをほとんどしていないK-FAC の性能であり、今後さらなる収束性の向上が見込める。

プラズマの挙動を予測する方法

Simulation

Deep Learning





Machine Learning Workflow



Background/Approach for DL/Al

• <u>Deep Learning Method</u>: distributed data-parallel approach to train deep neural networks \rightarrow <u>Python Framework using high-level Keras library with</u> <u>Google Tensorflow backend</u> <u>Reference</u>: Deep Learning with Python, François Chollet (Nov. 2017, 384 pages)

*** Major contrast with "Shallow Learning" approaches including SVM's, Random Forests, Single Layer Neural Nets, & modern Stochastic Gradient Boosting ("XG-BOOST") methods by enabling moving from ML software deployment on clusters to supercomputers:
→ Summit (ORNL); Tsubame-3 (TiTech); ABCI (U. Tokyo); ... also other architectures, e.g. – Intel Systems: beyond KNL to new designs for Aurora-21 @ANL

- -- <u>stochastic gradient descent (SGD)</u> used for large-scale (i.e., optimization on supercomputers) with parallelization via mini-batch training to reduce <u>communication costs</u>
- -- <u>DL Supercomputer Challenge</u>: need large-scale scaling studies to examine if convergence rate saturates with increasing mini-batch size (to thousands of GPU's)

プラズマ挙動の1ms未来の予測性能





FRNNID: 線計測データ



-		Single mad	chine	Cross-machin	Cross-machine with 'glimpse'	
_	Training set	DIII-D	JET (CW)	JET (CW)	DIII-D	$DIII\text{-}D+\delta$
	Testing set	DIII-D	JET (ILW)	DIII-D	JET (ILW)	JET (ILW) $-\delta$
XG-Boost	Best classical model	0.937	0.893	0.636	0.616	0.851
	FRNN OD	0.890	0.952	0.761	0.817	0.879
	FRNN 1D	0.922	_	_	0.836	0.911

Kates-Harbeck et al. Nature, 568, pp. 526-531 (2019)





パフォーマンスモデルよりも理想的な並列化効率に近いスケーリング ABCIのほぼ全体である4096 GPUまで理想的な並列化効率が得られた

関連する論文発表

二次最適化を用いた巨大な言語モデルの学習

Distributed Training of Large Language Models Using Second Order Optimization, KDD, submitted

Osawa, K., Swaroop, S., Jain, A., Eschenhagen, R., Turner, R. E., Yokota, R., Khan, M. E., ``Practical Deep Learning with Bayesian Principles'', The 33rd Conference on Neural Information Processing Systems (NeurIPS 2019).

Ueno, Y. and Yokota, R. ``Hierarchical Topology-aware Communication for Scaling Deep Learning to Thousands of GPUs'', The 19th Annual IEEE/ACM International Symposium in Cluster, Cloud, and Grid Computing (CCGrid 2019).

Osawa, K., Tsuji, Y., Ueno, Y., Naruse, A., Yokota, R., and Matsuoka, S., ``Large-scale Distributed Second-order Optimization Using Kronecker-factored Approximate Curvature for Deep Convolutional Neural Networks', Conference on Computer Vision and Pattern Recognition (CVPR 2019).

FRNNを用いたプラズマ挙動予測

Svyatkovskiy, A., Kates-Harbeck, J., and Tang, W., planning to submit to SC'20.