

Sony's deep learning development environment for massively large scale training

Takuya Narihira & Hisahiro Suganuma 8th ADAC Workshop, Oct 30, 2019

Sony Corporation R&D center

Copyright 2019 Sony Corporation

Agenda

- Deep Learning at Sony
- Large Scale Training
 - Large model training
 - Distributed training



Deep Learning at Sony

Sony's Product



Deep Learning is being utilized in many application domains.

Software 2.0

Software 1.0



- Write in code (C++, Python …)
- Requires domain expertise
- Hard to maintain
- Does not scale

Software 2.0



- Written in the weight of NN by optimization
- Requires data
- Easier to maintain (Data tells)
- Scale empowered by data

"Gradient Descent can write code better than you. I'm sorry." Andrej Karpathy

We need techniques, infrastructures and echo-systems for software 2.0 era

Sony's deep learning software



Neural Network Console

Neural Network Libraries nnabla.org



GUI based deep learning IDE

Windows Desktop (free) & Cloud (paid)

Deep learning framework with Python API

Open source (Apache 2.0 license)

Intuitive, fast and easy to deploy

	4			
4	(2	2	
	I	1	Г	ŝ
	1	÷	2	

PROJECT

DATASET

2

, and the set				Literation
01_logistic_regression.sdcproj	Inputs : x[1,28,28] Outputs : y[1] Training Dataset : small_mnist_4or9_training.csv (1500 datas, 2 columns) Validation Dataset : small_mnist_4or9_test.csv (491 datas, 2 C:\sdeep_console_free\samples\sample_project\tutorial\basics\	2015/06/29 15:58:00	▷ € □ ▷ □ Overview	
02_binary_cnn.sdcproj	Inputs : x[1,28,28] Outputs : y[1] Training Dataset : small_mnist_4or9_training.csv (1500 datas, 2 columns) Validation Dataset : small_mnist_4or9_test.csv (491 datas, 2 C:\sdeep_console_free\samples\sample_project\tutorial\basics\	2015/ <mark>06/19</mark> 10:21:59		
06_auto_encoder.sdcproj	Inputs : x[1,28,28] Outputs : x[1,28,28] Training Dataset : small_mnist_4or9_training.csv (1500 datas, 2 columns) Validation Dataset : small_mnist_4or9_test.csv (491 datas, 2	2015/06/18 14:15:40		
	C:\sdeep_console_free\samples\sample_project\tutorial\basics\		Statistics	
10 deep mlp.sdcproj	Inputs : x[1,28,28] Outputs : y[1]	2015/06/18 14:15:35	Output	0
	Training Dataset : small_mnist_4or9_training.csv (1500 datas, 2 columns)		CostParameter	0
	Validation Dataset : small_mnist_4or9_test.csv (491 datas, 2		CostAdd	0
	C:\soeep_console_free\samples\sample_project\tutonai\basics\		CostMultiply	0
12 residual learning edenroi	Inputs v[1 28 28] Outputs v[1]	2016/03/16 15:24:37	CostMultiplyAdd	0
	Training Dataset : mnist_training.csv (60000 datas, 2 columns)	constast to totenst	CostDivision	0
	Validation Dataset : mnist_test.csv (10000 datas, 2 columns)		CostExp	0
	C:\sdeep_console_free\samples\sample_project\tutorial\basics\		CostIf	0
LeNet.sdcproj	Dataset "Training": mnist_training.csv (60000 datas, 2 columns) Dataset "Validation": mnist_test.csv (10000 datas, 2 columns)	2017/08/07 14:46:41	Tasks	
	C\sdeep console free\samples\sample project\image recognition\	MNIST	Training:	

Ē

⇔

Easy model building, experiment, and deployment



Training can be distributed over GPUs by just clicking a couple of times in GUI

AI Bridging Cloud Infrastructure (ABCI) <u>https://github.com/aistairc/abci-docs/</u>

The world's first large-scale Open AI Computing Infrastructure, constructed and operated by <u>National Institute of Advanced</u> <u>Industrial Science and Technology (AIST)</u>.



Neural Network Libraries



Easy network definition

Both static & dynamic graph paradigms are supported

nnabla-examples: A bunch of SoTA training scripts

https://github.com/sony/nnabla-examples

Generative models for content creation

StarGAN

Black \rightarrow Blonde

CycleGAN

Pix2PixHD



Mask → Texture

MUNIT



Drawing \rightarrow Texture



Horse \rightarrow Zebra

InstaGAN



Skirt \rightarrow Pants

Self-Attention GAN

"dog"



Label → Image

Please use research baselines & application prototyping etc.



Large Scale Training

Why large scale training is important? 1/2

Bigger models are better in performance (accuracy)

Recognition Cite: GPipe



Generation Cite: BigGAN

Class conditional image generator



Image generation models has also become bigger (requires 512 cores TPUv3 pod to train)

Batch	Ch.	Param (M)	Shared	Skip-z	Ortho.	Itr $\times 10^3$	FID	IS
256	64	81.5	SA-0	GAN Base	eline	1000	18.65	52.52
512	64	016	×	×	×	1000	15.00	(± 1.18)
1024	64	Ridde	r X	×	×	1000	14. Ro [.]	$tt \rho r(\pm 1.42)$
2048	64	Digge	×	×	×	732	12.	(± 3.83)
2048	96	173.5	×	×	×	$295(\pm 18)$	$9.54(\pm 0.62)$	$92.98(\pm 4.27)$
2048	96	160.6	 Image: A set of the set of the	×	×	$185(\pm 11)$	$9.18(\pm 0.13)$	$94.94(\pm 1.32)$
2048	96	158.3	1	-	×	$152(\pm 7)$	$8.73(\pm 0.45)$	$98.76(\pm 2.84)$
2048	96	158.3	/	 Image: A start of the start of	 Image: A start of the start of	$165(\pm 13)$	$8.51(\pm 0.32)$	$99.31(\pm 2.10)$
2048	64	71.3	1	1	 Image: A set of the set of the	$371(\pm 7)$	$10.48(\pm 0.10)$	$86.90(\pm 0.61)$

Why large scale training is important? 2/2

Large dataset gives better performance



Target task: ImageNet

Data size becomes bigger



https://research.fb.com/publications/exploring-the-limits-of-

Bigger model, dataset, and data size impose large memory and longer training time

Large scale training: Large model execution in nnabla

Presentation only

SONY 16 2019.10.30 Sony Corporation R&D Center



Distributed Training



Introduction of techniques for scaling up distributed training

- toward over 100K mini-batch training -

Hisahiro SUGANUMA

Sony Corporation R&D center

Copyright 2019 Sony Corporation

Agenda

- Overview of data parallel distributed training
- Problems regarding scaling up disrtibuted training
- Techniques for large scale distributed training
 - Brief introduction of our new optimizer "STiLL"
- Experimental Result : ImageNet/ResNet-50 training

Brief introduction of our recent activities

- We are studying distributed training to train
- And to benchmark our framework, we have ImageNet/ResNet-50 that contain <u>the time</u> techniques for large batch training.

ImageNet/ResNet-50 Training in 224 Seconds

Hiroaki Mikami, Hisahiro Suganuma, Pongsakorn U-chupala, Yoshiki Tanaka and Yuichi Kageyama

Sony Corporation {Hiroaki.Mikami, Hisahiro.Suganuma, Pongsakorn.Uchupala, Yoshiki.Tanaka, Yuichi.Kageyama}@sony.com

Abstract

Scaling the distributed deep learning to a massive GPU cluster level is challenging due to the instability of the large mini-batch training and the overhead of the gradient synchronization. We address the instability of the large mini-batch training with batch size control. We address the overhead of the gradient synchronization with 2D-Torus all-reduce. Specifically, 2D-Torus all-reduce arranges GPUs in a logical 2D grid and performs a series of collective operation in different orientations. These two techniques are implemented with Neural Network Libraries (NNL)¹. We have successfully trained ImageNet/ResNet-50 in 224 seconds without significant accuracy loss on ABCl² cluster.



Brief overview of Distributed Training

Paradigms of distributed training : Data and Model Parallelism



DNN Training on single GPU.

- DNN Training basically employs "Mini-batch training"
 - Training consist of <u>fetch data</u>, forward, backward and update weights.
 - Iterative training with each split dataset "mini-batch" to optimize DNN parameters(weight) gradually.



Data-parallel distributed DNN Training

- Data parallel training needs to synchronize the gradients among the all workers before it updates the weights since it must maintain the consistency of the model.
- Synchronization is the obviously overhead to the time of training.



Problems regarding large scale training, and how do we solve them ?

Problems of distributed large-batch training.

- 1. <u>To achieve convergence with very small number of update steps</u>
 - The larger the mini-batch size is, the more number of GPUs we can use.
 - However, larger mini-batch size also makes it more difficult for convergence.



Mini-batch size will be N-times size, if using N GPUs. The point is developing the way to train even with large mini-batch size.

2. <u>To reduce gradient synchronization time</u>

 Ring topology becomes too slow, if using more than hundreds of GPUs . Out of scope in this talk

The point is developing a topology (for collective communication) which allow for an efficient synchronization even with the thousands of GPUs.

Convergence degradation in large-scale distributed DNN Training

Mini-batch size is also increased by a multiple of the number of GPUs.
 With SGD, it is hard to train DNN with large batch-size because ...



Example: optimization in 2-dimension.

% N. S. Keskar et al. "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima" ICLR 2017

Real world approaches to scale batch-size.

1. Tuning hyperparameters on SGD

Our method

- Many techniques were proposed.
 - E.g. "Learning rate warmup", "LARS", "LAMB", ...etc.
- We employ this way for several reasons. Mainly we consider that SGD to train DNN is wellstudied and we can use more wisdom of the past researches
- 2. Using second or higher order optimization like "Newton's method"
 - They could optimize with smaller steps without sharp minima. But most of them need enormous computing resource due to using hessian.
 - Recently proposed approximation have made us optimize DNN in real.

Overview of common techniques for large batch training

Learning rate warmup*

- Linearly scaling up LR in several epochs.



LARS optimizer**

- <u>Layer-wise Adaptive LR Scaling with the ratio of</u> weights and grads.

$$\lambda^{l} = \eta \times \frac{||w^{l}||}{||\nabla L(w^{l})|| + \beta * ||w^{l}|}$$

* T.He et al. Bag of Tricks for Image Classification with Convolutional Neural Networks (<u>arxiv.org/abs/1812.01187</u>) ** Y. you et al. Large Batch Training of Convolutional Networks (arxiv.org/abs/1708.03888)

Up to 64K mini-batch these are enough to train.

However we need another technique to increase batch-size more and more !

Label Smoothing***

Smoothing i-th label of Softmax output **q**(0 < ε < 1, K is the number of classes.)

$$q_i = \begin{cases} 1 - \varepsilon & \text{if } i = y, \\ \varepsilon/(K - 1) & \text{otherwise,} \end{cases}$$
(4)

*** C. Szegedy et al Rethinking the Inception Architecture for Computer Vision (arxiv.org/abs/1512.00567)

Zero-γ-init

- Initilizing γ of batchnorm in residual path to 0.



STILL : Smoothly Transition from LAMB* to LARS.

- We proposed the new optimizer "STiLL" that uses two optimizer, LAMB and LARS at the same time.
- STILL starts from LAMB:LARS=100:0 to update the weights, then it linearly changes the ratio of them during warmup duration.
- After the warmup duration, it uses only LARS to update weights.



*Y.You et al, Large Batch Optimization for Deep Learning: Training BERT in 76 minutes https://arxiv.org/abs/1904.00962

Why we mixed them?

• We expected that LAMB accelerates convergence in the beginning of training and LARS helps to get better generalization performance by combining them.



Other Technique – Smoothout, that strengthens regularization.

- Smoothout* is one of the regularization technique proposed by W.Wen.
 - 1. Add noise from Uniform dist. to the weights before forward computation.
 - 2. Do forward and backward computation.
 - 3. Sub. noise from the weights.
 - 4. Update weights.
- "Smoothout" can avoid sharp-minima, that is more important in large-batch training. We use this over 64K batch training.



*W.Wen et al. SmoothOut: Smoothing Out Sharp Minima to Improve Generalization in Deep Learning https://arxiv.org/abs/1805.07898



Experiments

Experiments(1): Performance comparison with LARS, LAMB and STILL

- Firstly, we conducted ImageNet/ResNet-50 training while changing optimizers, "LAMB", "LARS" and "STILL", and increasing batch-size from 72K to 96K.
- STiLL outperforms "LAMB" and "LARS" in the final validation accuracy.



Training curve of LAMB, LARS and STiLL

Top-1 accuracy of ImageNet/ResNet-50 model after training for 90 epochs

Mini-Batch Size	LAMB	LARS	STiLL
72К	70.9%	75.5%	76.1%
80K	N/A	75.1%	75.9%
96К	N/A	N/A	75.5%

Experiments(2): Performance of NNL in ABCI.

• We evaluated wallclock time to train for 90 epochs and convergence accuracy with STiLL in ABCI.

- Settings

Batch size	64K, 80K, 96K
STiLL's Initial/MAX LR	Init: 1.36(LARS)/0.0012(LAMB), Max: 26.1
LARS Momentum	0.955
Warmup/LR Decay	Warmup:25epoch, and then poly decay.
Label-smoothing α	0.1
Weight decay	8e-5

<u>- Results</u>

BS	epoch	#GPU	Wallclock	Accuracy
64K	90	2048	97 sec	76.13 %
80K	90	2560	79.7 sec	75.93 %
96K	90	3072	67.65 sec	75.46 %

- Training curve



Experiments (3): Training longer epochs could give better performance

- Then we made extra experiments to see if the model get better generalization performance by training longer epochs.
- In the result, the model could reach good solution even over 100K batch training.

- Settings
- Use the same settings as the previous exp.

- Results

BS	epoch	#GPU	Acc
96K	110	3072	76.19%
128K	150	4096	76.12 %



- Training curve

- We proposed the new optimizer "STiLL" for large-scale training.
- Then empirically showed that "STiLL" outperforms "LARS" and "LAMB" in Imageclassification task.
- Also we succeeded to train with 128K batch, however even longer epoch we need now.

Question?

SONY

SONY is a registered trademark of Sony Corporation.

Names of Sony products and services are the registered trademarks and/or trademarks of Sony Corporation or its Group companies.

Other company names and product names are registered trademarks and/or trademarks of the respective companies.