

Mixed Precision Training

计算平台事业部PAI团队





- What is mixed-precision & Why mixed-precision
- How mixed-precision
- Mixed-precision tools on PAI-tensorflow
- Experimental results

What is mixed-precision

- mixed-precision
 - FP32 and FP16
 - More precision format in the future
- TensorCore
 - Matrix-multiply and accumulate units
 - FP16 storage/inputs
 - FP32/Fp16 accumulator
 - Such as:
 - Conv
 - MatMul



Why mixed-precision



- Two key points which matter in training/inference:
 - Computation
 - Tensorcore 8X higher throughput in MP than FP32 (15Tflops v.s. 120Tflops)
 - Memory access
 - Inputs is FP16
 - Memory access is reduced by 2X





- What is mixed-precision & Why mixed-precision
- How mixed-precision
- Mixed-precision tools on PAI-tensorflow
- Experimental results

How mixed-precision



- Key Strategies in mixed-precision training
 - Issues using FP16 for training and the solutions
 - Less bits in fraction: \rightarrow Precision gap in sum
 - Less bits in exponent: → Gradients underflow
 - Arithmetic precision design
 - Considering both efficiency and performance

Issues using FP16 for training



- Less bits in fraction: Precision gap in sum
 - A+B, if A/B>2¹⁰, B will degrade to **zero**.
 - For FP32, the ratio can be up to 2²³
 - Common in weight update:

• *W*←*W*+*Ir*dW*



- Less bits in exponent: Gradients underflow
 - Gradients smaller that 2⁻²⁴ will become zero

Precision gap in sum



- variables v.s. gradients
 - weight update: $W \leftarrow W + lr^* dW$ (Ir normally in [10⁻¹, 10⁻⁴])



Variables: 2⁻¹⁶ to 2⁻⁴

gradients: 2⁻³⁰ to 2⁻⁵

Fig. Variables and gradients histogram in Faster RCNN

• Solution: Variables stored in FP32, and optimizer computation in FP32



Gradients of variables



Fig. Histogram for gradients of variables in Faster RCNN, respectively training in mixed-precision and FP32



Gradients of activations



Fig. Histogram for gradients of variables in Faster RCNN, respectively training in mixed-precision and FP32





Solution: gradients shift using loss scaling



- Constant loss scaling
 - Scale the loss by a factor S
 - Backprop to compute the dW
 - Unscale dW by 1/S
- Automatic loss scaling
 - Start with a large scaling factor S
 - For each training iteration:
 - Scale the loss by S
 - Backprop to compute the dW
 - Unscale dW by 1/S
 - If dW contains Inf/NaN, the decrease S by a step factor <u>S/step</u>
 - Otherwise, update dW to W
 - If there is no Inf/NaN for N updates, the increase S by a step factor <u>S*step</u>

How mixed-precision



- Key Strategies in mixed-precision training
 - Issues using FP16 for training and the solutions
 - Less bits in fraction: \rightarrow Precision gap in sum
 - Solution: Variables stored in FP32, and optimizer computation in FP32
 - Less bit in exponent: → Gradients underflow
 - Solution: loss scaling
 - Arithmetic precision design
 - Considering both efficiency and performance

How mixed-precision



- Key Strategies in mixed-precision training
 - Issues using FP16 for training and the solutions
 - Less bits in fraction: \rightarrow Precision gap in sum
 - Solution: Variables stored in FP32, and optimizer computation in FP32
 - Less bit in exponent: → Gradients underflow
 - Solution: loss scaling
 - Arithmetic precision design
 - Considering both efficiency and performance

Arithmetic precision design



- Arithmetic can be categorized into:
 - 1. Compute-bound
 - Convolution, Matmul
 - 2. Memory-bound
 - $\textcircled{1} \quad \text{Reductions} \quad$
 - Batch-norm/layer-norm/group-norm
 - Softmax / Average pooling
 - ② Element-wise operation
 - Add/mul, etc



Take advantage of Tensorcore:

- Inputs: FP16
- Accumulator: FP32
- Outputs: FP32

- Inputs/Ouputs in FP16
- Computation in FP32

- Inputs/Ouputs in FP16
- Computation in FP16
 - Computation in FP32

Arithmetic precision design



- Compute-bound operations:
 - Inputs in FP16
 - Computation using Tensorcore
 - Outputs in FP32
- Memory-bound operations:
 - Inputs/outputs in FP16
 - Computation in FP32

How mixed-precision training



FP32 training baseline



- Convert the computation part to MP
- Remain the optimizer part in FP32





tf.cast(var, tf.float16)



• Loss Scaling strategy (constant scaling)



• Auto Loss Scaling strategy





• Auto Loss Scaling strategy





• Auto Loss Scaling strategy





21

• Auto Loss Scaling strategy



计算平台事业部





- What is mixed-precision & Why mixed-precision
- How mixed-precision
- Mixed-precision tools on PAI-tensorflow
- Experimental results

MP training tools on PAI-TF



- Graph Optimization + Loss Scaling Training Strategy
 - Graph Optimization: AutoMixedPrecision Graph Optimization Pass

 FP32 graph_def
 Automatically conversion

 MP graph_def

- MP Training Strategy: MP optimizer wrapper
 - Wrap the standard optimizers to automatically adopt the constant/automatic loss scaling strategy
 - opt = tf.contrib.mixed_precision.MixedPrecisionOptimizer(opt)
 - Both constant/automatic loss scaling supported

Mixed-precision optimizer			
	Standard optimizer		



• ResNet50 on ImageNet





• Faster RCNN (VGG backbone) on PASCAL VOC 07

	mAP	mAP Diff
FP32(baseline)	69.48	-
FP16(no loss scaling)	67.83	-1.65
FP16(constant scaling=8)	69.71	+0.23
FP16(constant scaling=64)	69.96	+0.48
FP16(constant scaling=256)	69.78	+0.30
FP16(constant scaling=4096)	69.74	+0.26
FP16(auto scaling)	69.68	+0.20



• SSD (VGG backbone) on PASCAL VOC 07+12

训练精度	loss scaling策略	mAP
FP32	NO	77.46 (baseline)
Mixed precision	NO	diverge
	constant (=64)	77.63
	auto scaling	77.42



Small NMT on WMT German-English

- Encoder: 2 layers
- Decoder: 2 layers with attention

训练精度	Scaling策略	BLEU
FP32	NO	23.68
Mixed-precision	Auto scaling	23.79

PGAN



• PGAN (Progressive growth of GAN)



Karras, Tero, et al. "Progressive Growing of GANs for Improved Quality." Stability, and Variation.

PGAN

• G loss



Karras, Tero, et al. "Progressive Growing of GANs for Improved Quality." Stability, and Variation.

PGAN



• Generation results (cifar10 dataset)

fp32



mp-auto-scaling



mp-no-scaling



Exp.	fp32	mp-auto-scaling	mp-no-scaling
sliced_wasserstein	9.3764	9.1662	7.9601

Font Generation





Note:
 denotes element-wise subtraction

Font Generation



• G loss



Pyramid Embedded Generative Adversarial Network for Automated Font Generation

Font Generation



• Generation results (金陵刻经体)



Wide & Deep Learning



 Predict the probability that the individual has an annual income of over 50,000 dollars



Figure 1: The spectrum of Wide & Deep models.

Wide & Deep Learning



• Loss

loss



Ехр	fp32	mp-no-scaling
Accuracy	84.31%	84.27%

Wide & Deep Learning for Recommender Systems

More try: small inputs (for normalization layer (COMPUTING PLATFORM

- Underflow in FP16 gradients
 - Design the model to be more adaptive to FP16 representation
 - Move the gradient itself into the FP16 representable range
 - Especially the activation gradients
- Batch normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$; Parameters to be learned: γ, β **Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$ $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$ // mini-batch mean $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$ // mini-batch variance $\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ // normalize $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$ // scale and shift

- BN will stand for Batch Norm.
- f represents a layer upwards of the BN one.
- y is the linear transformation which scales x by γ and adds β .
- \hat{x} is the normalized inputs.
- μ is the batch mean.
- σ^2 is the batch variance.



• Derivatives of BN layer

$$\frac{\partial f}{\partial x_i} = \frac{m\frac{\partial f}{\partial \hat{x}_i} - \sum_{j=1}^m \frac{\partial f}{\partial \hat{x}_j} - \hat{x}_i \sum_{j=1}^m \frac{\partial f}{\partial \hat{x}_j} \cdot \hat{x}_j}{m\sqrt{\sigma^2 + \epsilon}}$$

$$rac{\partial f}{\partial x_i} \propto \mathcal{O}\left(1/\sigma
ight)$$

Smaller Inputs and Bigger derivatives

- Reduce the magnitude of the inputs
 - Reduce magnitude of the <u>forward</u> activations, so as to reduce the <u>overflow</u> in forward propagation when using FP16
 - <u>Improve</u> the magnitude of the <u>derivatives</u>
- Tips for Network with BN:
 - Normalize the layer to have std to be 1/S rather than 1.0

ResNet32+CIFAR10
Activations and the gradients







activations

activation gradients



- ResNet32+CIFAR10
 - All without loss scaling

训练精度	Input scale	Accuracy
FP32	1.0	93.58
Mixed-precision	1.0	93.41
Mixed-precision	1/8.0	93.43



• SSD on PASCAL VOC 07+12

Activations and the gradients





activations

gradients



• SSD on PASCAL VOC 07+12

• Activations and the gradients

训练精度	loss scaling策略	mAP (standard input)	mAP (small input, scaled 1/64)
FP32	NO	77.46 (baseline)	77.25
Mixed precision	NO	diverge	77.75
	constant (=64)	77.63	77.56
	auto scaling	77.42	77.17

Conclusion



- Mixed-precision tools have been supported on PAI-tensorflow
- More effort is still conducted to explore more in mixedprecision
 - More precision supported
 - More training strategy



Thank you