Low-precision training **EECE695D: Efficient ML Systems**

Spring 2025

Low-precision training

- Idea. Exploit low-precision for training
 - Faster computation
 - Reduced memory bandwidth



More throughput, less energy \rightarrow **Bigger batch** \rightarrow

mantissa = 0.395264 0 0 1 0 0 0 1 0 0 0 0 = 0.394531 1 1 1 1 0

1

= 0.40625

• The archenemy is the limited dynamic range

- Limited precision leads to over-/underflows
 - FP8 (E4M3) value covers maximum 448 and minimum 2^{-6}

- Can amplify the gradient noise
- Imprecise updates accumulate over time

Key challenge

• Further affects gradient, optimizer states, BN/LN statistics (quantize what?)



This week

- We take a very brief look at some notable examples:
 - Early IBM work (2015)
 - Binary networks (2015–2016)
 - BNN, XNOR, DoReFa, eventually to BitNets
 - FP16/FP8 training (2018–)

Gupta et al., (2015)

Background

- FP32 is a modern standard; old works used various precisions for DL:
 - Iwata et al. (1989) uses FP24
 - Hammerstrom (1990) uses 8–16 bits in fixed point
 - (...)
 - for their supercomputer

- This work. Train using 16bit fixed-point \approx FP32
 - + Hardware prototyping

• Chen et al. (2014) observes that at least 32bit fixed point is needed

(MNIST, CIFAR-10)

Iwata et al., "An artificial neural network accelerator using general purpose 24 bit floating point digital signal processors," IJCNN 1989 Hammerstrom, "A VLSI architecture for high-performance, low-cost, on-chip learning," IJCNN 1990 Chen et al., "DaDianNao: A Machine-Learning Supercomputer," MICRO 2014





Method

- To compute $\mathbf{a}^{\mathsf{T}}\mathbf{b}$ for two k-bit fixed-point vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$:
 - Step 1. Compute the MAC

- **z** requires bitwidth at most $2k + 1 + \log_2 d$
- Step 2. Convert the sum to k-bit
 - Round-to-nearest
 - Stochastic rounding
- Gradient in DL is matmul, thus no special consideration needed



Suyog et al., "Deep Learning with Limited Numerical Precision," ICML 2015



Method

- Quantized many things:
 - Weights
 - Biases
 - Activations
 - Back-propagated error
 - Weight update
 - Bias update



Observation

Stochastic rounding is quite essential (MNIST+MLP example)





Suyog et al., "Deep Learning with Limited Numerical Precision," ICML 2015



Observation

- Up to 16bit fixed-point, for CIFAR-10 training
 - Assign more bits for the integer bits than fraction
 - Late high-precision training helps



Suyog et al., "Deep Learning with Limited Numerical Precision," ICML 2015



Background

- with special focus on **binary networks**
 - ICLR workshop 2015
 - propagations" NeurIPS 2015
 - "Binarized Neural Networks" NeurIPS 2016
- Loses scaling
 - BatchNorm layers can handle these, presumably

Another team at Montréal worked on early shapings of low-precision training,

"Training deep neural networks with low precision multiplications"

"BinaryConnect: Training Deep Neural Networks with binary weights during



Method

- Binarize the weight during forward, and keep high-precision weights intact
 - Forward. Compute output with binarize(w)
 - Stochastic rounding; applies hard sigmoid
 - **Backward**. Compute gradients in full-precision
 - Update full-precision weight w
 - Straight-through estimator
 - Apply clipping to [-1,+1]; prevents w from growing endlessly

- Binarization noise is deemed independent and zero-mean
 - cancels out eventually

Courbariaux et al., "BinaryConnect: Training Deep Neural Networks with binary weights during propagations" NeurIPS 2015



Algorithm 1 SGD training with BinaryConnect. C is the cost function for minibatch and the functions binarize(w) and clip(w) specify how to binarize and clip weights. L is the number of layers.

and learning rate η .

Ensure: updated parameters w_t and b_t .

1. Forward propagation:

 $w_b \leftarrow \text{binarize}(w_{t-1})$

For k = 1 to L, compute a_k knowing a_{k-1} , w_b and b_{t-1} 2. Backward propagation:

Initialize output layer's activations gradient $\frac{\partial C}{\partial a_T}$ For k = L to 2, compute $\frac{\partial C}{\partial a_{k-1}}$ knowing $\frac{\partial C}{\partial a_k}$ and w_b **3. Parameter update:** Compute $\frac{\partial C}{\partial w_b}$ and $\frac{\partial C}{db_{t-1}}$ knowing $\frac{\partial C}{\partial a_k}$ and a_{k-1} $w_t \leftarrow \operatorname{clip}(w_{t-1} - \eta \frac{\partial C}{\partial w_b})$ $b_t \leftarrow b_{t-1} - n \frac{\partial C}{\partial t}$

$$b_{t-1} - \eta \overline{\partial b_{t-1}}$$

Require: a minibatch of (inputs, targets), previous parameters w_{t-1} (weights) and b_{t-1} (biases),

Courbariaux et al., "BinaryConnect: Training Deep Neural Networks with binary weights during propagations" NeurIPS 2015



Observations

- Requires much longer training than full-precision
- Interestingly, advantages in generalization



Courbariaux et al., "BinaryConnect: Training Deep Neural Networks with binary weights during propagations" NeurIPS 2015



2016 work

- Binarizes both W&A
 - Keeps both W&A in full-precision as well
- Uses hard-sigmoid as activation function, with additional binarization
 - Gradients for activations outside [-1,+1] is set to zero
 - Empirically works well, rather than being well-justified (in fact, later works remove this; e.g., XNOR-Net, DoReFA-Net)



Hubara et al., "Binarized neural networks" NeurIPS 2016



Limitation

- Keeps track of full-precision W&A
 - Much memory cost
- Requires full-precision gradient
 - Weak computation advantage during backward

XNOR-Net

- Introduces scaling factor to the forward binary operation
 - This change force us to update weights greater than 1, as well.
 - **BNN**.

Forward: r_{o} **Backward:** $\frac{\partial}{\partial t}$

- XNOR-Net. Forward: r
 - **Backward:** $\frac{\partial}{\partial t}$

(expectation denotes averaging w.r.t. output channel)

$$b_{o} = \operatorname{sign}(r_{i})$$

 $\partial c_{o} = \frac{\partial c}{\partial r_{o}} \mathbb{I}_{|r_{i}| \leq 1}.$

$$b_{o} = \operatorname{sign}(r_{i}) \times \mathbf{E}_{F}(|r_{i}|)$$

 $\frac{\partial c}{\partial r_{i}} = \frac{\partial c}{\partial r_{o}}.$



DoReFa-Net

Generalizes the idea to binarizing gradients as well

- Unlike activations, gradients are **unbounded** (typically larger than activations) Thus perform quantization with normalization factors

$$f_{\gamma}^{k}(\mathrm{d}r) = 2\max_{0}(|\mathrm{d}r|) \left[\operatorname{quantize}_{k} \left[\frac{\mathrm{d}r}{2\max_{0}(|\mathrm{d}r|)} + \frac{1}{2} + N(k) \right] - \frac{1}{2} \right) \right]$$

gradient quantization

• Here, N(k) is a uniform noise which compensates for any bias due to

Zhou et al., "DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients" arXiv 2016



Further readings

- This line of work has been followed-up by BitNet, in the context of LLMs
- <u>https://arxiv.org/abs/2310.11453</u>

FP16/FP8 training

Background

- Concerns regarding "will fixed-point scale up?"
 - ImageNet-scale experiments
 - Applications to language modeling

• **Solution.** Use floating point?

FP16 Training (NVIDIA ver.)

- NVIDIA & Baidu worked on extending the idea to FP16
 - Weights are also kept in FP32, but not others
 - <u>Reason</u>. Weight updates will be too small, when multiplied with LR



Micikevicius et al., "Mixed Precision Training" ICLR 2018



FP16 Training (NVIDIA ver.)

- Also, we cannot use the same dynamic range for weights and gradients
 - Many gradients have very small values
 - Solution. Scale up the loss (x8), so that the gradient become larger





Figure 5: bigLSTM training perplexity

Micikevicius et al., "Mixed Precision Training" ICLR 2018

Further readings

- Intel concurrently studied INT16 training (ICLR 2018):
 - <u>https://arxiv.org/abs/1802.00930</u>

- IBM introduced FP8 format (E5M2)
 - Accumulation in FP16 (E6M9); still a "master copy" of weights
 - AXPY (vector addition; $y = \alpha x + y$) done in FP16



Figure 2: A diagram showing the precision settings for (a) three GEMM functions during forward and backward passes, and (b) three AXPY operations during a standard SGD weight update process.



- Observation. "Swamping" happens a lot in DL
 - Losing information when adding large + small FP values
 - Especially problematic if:

 - Adding a lot of non-zero mean values (as the sum gradually grows) • Some of the elements are extremely large



- Solution. Chunk-based accumulation
 - Generate partial sums to avoid scale mismatch
 - DON'T:
 - a1 = a1 + a2
 - a1 = a1 + a3
 - a1 = a1 + a4
 - DO:
 - a1 = a1 + a2
 - a3 = a3 + a4
 - a1 = a1 + a3



- Solution. Stochastic rounding
 - Rounded-off values (least significant bits) give chance of rounding up, instead of simply being truncated

• For floating-point, rounding the value $x = s \cdot 2^{e} \cdot (1 + m)$ is done by: $\operatorname{Round}(x) = \begin{cases} s \cdot 2^{e} \cdot (1 + \lfloor m \rfloor + \epsilon) & \text{w.p.} & \frac{m - \lfloor m \rfloor}{\epsilon} \\ s \cdot 2^{e} \cdot (1 + \lfloor m \rfloor) & \text{w.p.} & 1 - \frac{m - \lfloor m \rfloor}{\epsilon} \end{cases}$



- Use Chunking for GEMM, and SR for AXPY



AXPY is done over data batches; no hierarchical aggregation doable



• Other details. No quantization for 1st and the last layer





- Later work from IBM uses E4M3 for forward and E5M2 for backward
 - Works better on MobileNetV2 and Transformers https://papers.nips.cc/paper_files/paper/2019/hash/ 65fc9fb4897a89789352e211ca2d398f-Abstract.html
 - Also much discussions on BatchNorm
- Similar findings by Graphcore
 - <u>https://arxiv.org/abs/2206.02915</u>
 - - More discussions on loss scaling per-tensor scaling

See also 2023 work: <u>https://openreview.net/forum?id=nErbvDkucY</u>

Sun et al., "Hybrid 8-bit Floating Point (HFP8) Training and Inference for Deep Neural Networks" NeurIPS 2019 Noune et al, "8-bit Numerical Formats for Deep Neural Networks," arXiv 2022 Perez et al., "Training and inference of large language models using 8-bit floating point," WANT workshop @ NeurIPS 2023

- - <u>https://arxiv.org/abs/2209.05433</u>



NVIDIA, ARM, and Intel proposes an interchange format of E4M3 and E5M2:

Conducts extensive empirical study, showing that FP8 is stable enough

Figure 1: Training loss (perplexity) curves for various GPT-3 models. x-axis is normalized number of iterations.

Micikevicius et al, "FP8 formats for deep learning," arXiv 2022



- Microsoft applies FP8 on all pipeline: including parallelism & optimizer
 - <u>https://arxiv.org/abs/2310.18313</u>
 - Precision decoupling: Reduced precision for insensitive weights





Per-tensor scaling of loss

Peng et al., "FP8-LM: Training FP8 Large Language Models," arXiv 2023



- DeepSeek–V3 applies E4M3 for all ops
 - <u>https://arxiv.org/abs/2412.19437</u>



Combine with fine-grained quantization (and many other details)

Takeaways

- Both families are seeing much progress nowadays
 - **PTQ.** Any lessons to learn from SOTA techniques?
 - QAT. Similar ideas applicable for QAT?

• We'll see more in the student presentations...

