

Accelerating Transformers with FP8 Griffin Lacey | 2023 May 04



- More energy efficient
- First introduced in Volta GPU microarchitecture (2017)

	Ada	Hopper	Ampere	Turing	Volta
Supported Tensor Core precisions	FP64, TF32, bfloat16, FP16, FP8, INT8	FP64, TF32, bfloat16, FP16, FP8, INT8	FP64, TF32, bfloat16, FP16, INT8, INT4, INT1	FP16, INT8, INT4, INT1	FP16
Supported CUDA' Core precisions	FP64, FP32, FP16, bfloat16, INT8	FP64, FP32, FP16, bfloat16, INT8	FP64, FP32, FP16, bfloat16, INT8	FP64, FP32, FP16, INT8	FP64, FP32, FP16, INT8

Some History on Tensor Cores

 Tensor Cores are specialized hardware execution units that perform matrix multiplications 8x - 64x faster than scalar instructions in "CUDA Cores"

• Higher densities of Tensor Cores in datacenter GPUs (Tesla)





Allocate 1 bit to either range or precision

Inside FP8 Tensor Cores



Support for multiple accumulator and output types



- Accelerate math-intensive operations
- - FP8 requires 25% of the memory of FP32
- Facilitates preparation for inference
 - Many models are deployed in low precision

Why use FP8?

• On H100 using FP8 with Tensor Cores has 60x more theoretical throughput than FP32 with CUDA cores

Accelerates memory-intensive operations and reduces storage requirements



Partition the DL graph into safe and unsafe ops

- Safe ops benefit from reduced precision and have dynamic ranges are similar from input to output
- During forward pass compute multiplications in FP16 and accumulations in FP32 for safe ops
- During backward pass use loss scaling factor to avoid gradient under/overflows

Mixed Precision Recipe The FP16 Way





Partition the DL graph into safe and unsafe ops

- During forward pass use E4M3 and in backward pass use E5M2
 - E4M3 dynamic range: 18 powers of 2
 - E5M2 dynamic range: 32 powers of 2

- During backward pass use per-tensor scaling factors
 - Pick scaling factor based on window of past N instances of the tensor

Mixed Precision Recipe The FP8 Way

• Unsafe ops don't necessarily need to be FP32, can be FP16/BF16

Dynamic range (in powers-of-two) needed for different tensors in GEMMs.				
Network	Activations	Weights	Activation Gradients	Weight Gradients
GPT 126M	10	10	17	11
Transformer-XL	8	8	11	8
ResNet50	7	5	16	9
ResNet18	5	7	14	9





Operations that are safe in FP8

- Convolutions
- Linear layers in Transformer blocks
- BatchNorm has been tested for specific cases (e.g. ResNet50)
- Operations that should remain in higher precision
 - Most nonlinear functions and normalizations (GeLU, Softmax, ...)
 - Residual connections in LLMs
 - Loss functions

What is Considered Safe?

Context dependent – some networks or tasks may tolerate more layers in FP8



import torch import transformer engine.pytorch as te from transformer engine.common import recipe

Set dimensions. in features = out features = hidden size =

Initialize model and inputs. model = te.Linear(in features, out features, bias=True) inp = torch.randn(hidden size, in features, device="cuda")

Create FP8 recipe. fp8 recipe = recipe.DelayedScaling()

Enable autocasting to FP8. with te.fp8 autocast(enabled=True,

out = model(inp)

loss = out.sum()loss.backward()

Transformer Engine

fp8 recipe=fp8 recipe):

- An open-source library implementing the FP8 recipe for Transformer building blocks
- Optimized for FP8 and other datatypes
- Supports PyTorch, with JAX and TF support coming soon
- operators
- Supports model and sequence parallelism
- https://github.com/NVIDIA/TransformerEngine
- Composable with the native framework



Convolutional Networks

Model	16-bits	FP8
ResNet18	70.58	70.12
ResNet34	73.84	73.72
ResNet50	76.65	76.61
ResNet101	77.51	77.48 ¹
ResNext50	77.68	77.62
Wide ResNet50	78.13	77.901
InceptionV3	77.23	77.071
Xception	79.46	79.24 ¹
DenseNet121	75.59	75.39 ¹
DenseNet169	76.97	76.94 ¹
Dilated ResNet C-26	75.22	75.041
Dilated ResNet C-42	76.80	76.52 ¹
Dilated ResNet A-50	78.16	78.12 ¹
MobileNet V2	71.65	71.04 ¹

¹Experiments on different FP8 recipe choices

Image Classification

Model
DeiT Tiny
DeiT Small
ViT Tiny
ViT Small
ViT Base
ViT Large
Swin-V1 Tiny
Swin-V1 Small
Swin-V1 Base
Swin-V2 Tiny
Swin-V2 Small
Swin-V2 Base
GC-ViT Tiny
GC-ViT Small
GC-ViT Base
FAN Tiny
FAN Small
FAN Base

FP8 Formats for Deep Learning, Micikevicius et. al.

Vision Transformers

16-bits	FP8	
72.69	72.52	
80.08	79.75 ¹	
71.61	71.24	
77.19	77.24	
80.27	80.20	
78.38	78.39	
81.13	81.16	
83.02	83.07	
83.50	83.42	
82.79	82.76	
84.30	84.12	
84.50	84.60	
82.32	82.88	
84.03	83.98	
84.42	84.46	
78.65	78.75	
82.57	82.56	
83.41	83.14	



Network	Metric	16-bits	FP8
GNMT	BLEU	24.83	24.651
Vaswani Base	BLEU	26.87	26.83 ¹
Vaswani Large	BLEU	28.43	28.35 ¹
Transformer-XL Base	PPL*	22.71	22.76
Transformer-XL Large	PPL*	17.90	17.85
Megatron BERT Base	Loss*	1.352	1.3571
Megatron BERT Large	Loss*	1.163	1.1671
JoC BERT Large	F1	89.89	90.23
T5 Base	F1	91.68	91.88
T5 Large	F1	93.41	93.66
T5 Base	Rouge	42.88	42.88
T5 Large	Rouge	43.84	43.64
GPT 126M	PPL*	19.36	19.50
GPT 357M	PPL*	14.07	14.17
GPT 1.3B	PPL*	10.77	10.78
GPT 5B	PPL*	8.95	8.98
GPT 22B	PPL*	7.21	7.24
GPT 175B	PPL*	6.65	6.68

¹Experiments on different FP8 recipe choices

*Lower means better

Natural Language Processing

5B BF16 Perplexity

Results based on running on A100 (16-bit Tensor Cores) with FP8 IO and on H100 FP8 Tensor Cores



% of Training



• LLMs

Up to 3x faster with H100 FP8 than A100 BF16/FP16



DL TRAINING PERFORMANCE



GPT 40B





• What's needed for accuracy?

- Transformer Engine Recipes
 - Per-tensor scaling
 - Both FP8 types (E5M2 and E4M3)
 - Mixed precision
- What's needed for performance?
 - H100 GPU
 - FP8 Tensor Cores for faster math
 - FP8 I/O to reduce bandwidth pressure
- You should consider using if you:

 - Already use 16-bit mixed-precision on Ampere and other GPUs

Conclusions

• Are training large DL models (e.g. LLMs) and want further speedup on H100





