

# Efficient Deep Learning Systems

## Optimizing training pipelines

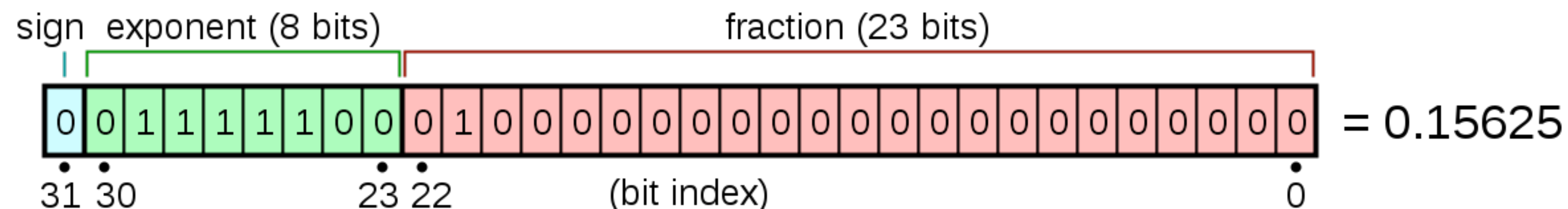
Max Ryabinin

# Plan for today

- Mixed precision training
  - When and why to use it
  - How to enable it and utilize it to the fullest
  - Dealing with stability in training
- Training pipeline optimization
  - Hardware considerations
  - Storing and loading data efficiently
- Profiling DL code

# Floating point numbers

- Neural networks require real numbers...
- ...which need to be represented in finite memory
- Single precision (FP32) is the default format with 4 bytes of storage



$$\text{value} = (-1)^{\text{sign}} \times 2^{\text{E}-127} \times \left( 1 + \sum_{i=1}^{23} b_{23-i} 2^{-i} \right)$$

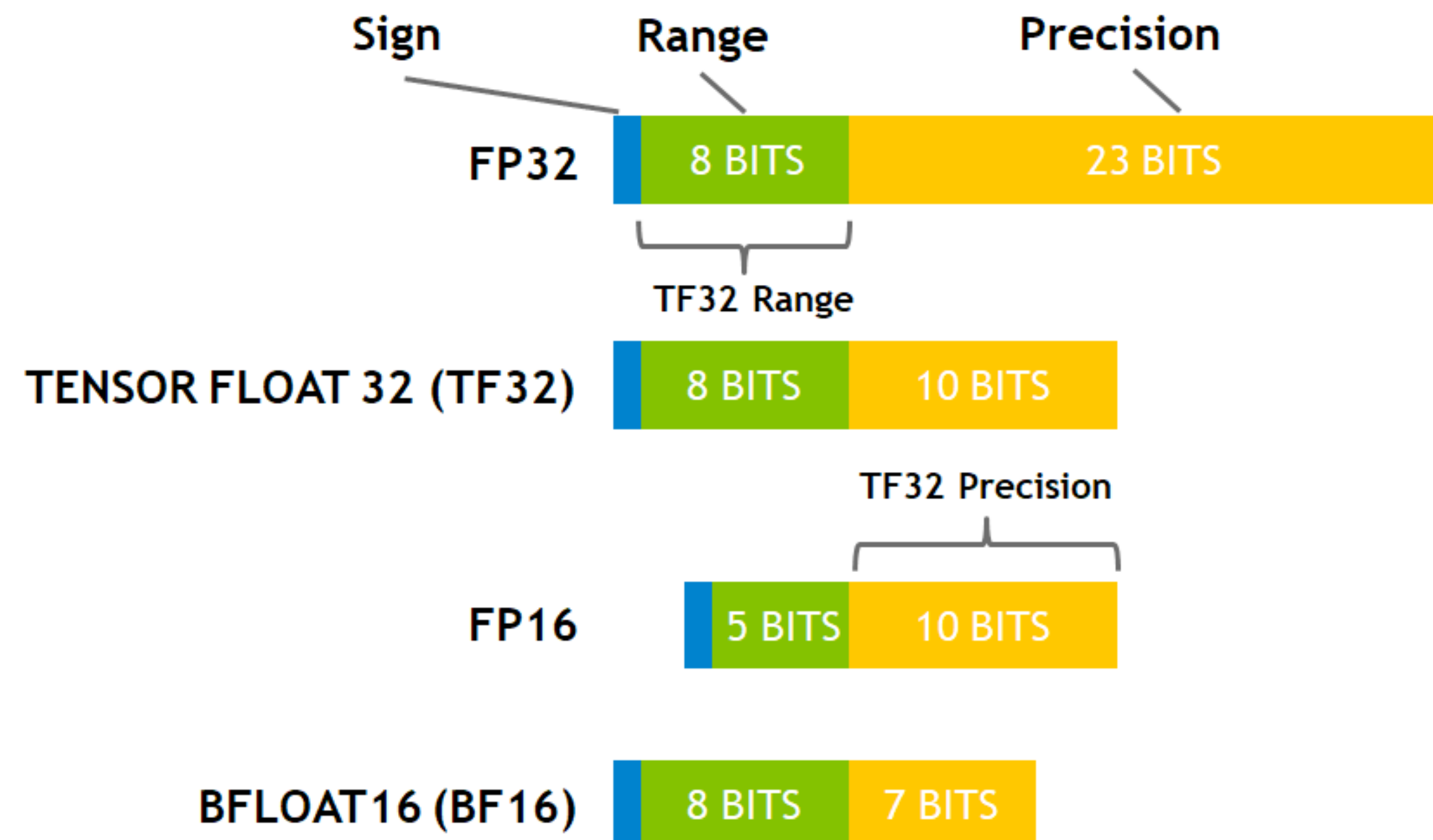
- Special values (0, NaN,  $\pm\text{inf}$ ) are encoded by exponent values

# Why use low precision?

- Can we go smaller than 32 bits? Should we?
- Key benefits:
  - Reduced memory usage (duh)
  - Faster performance (due to higher arithmetic intensity or smaller communication footprint)
  - Can use specialized hardware for even faster computation
- Makes your code prone to spectacular explosions :)

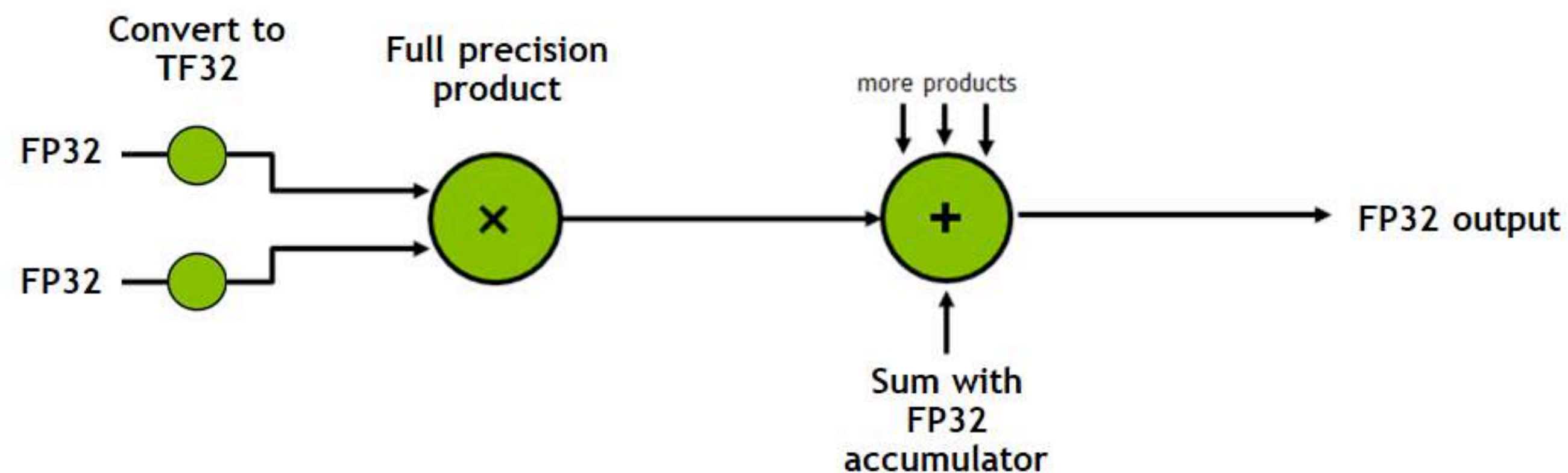
# Floating point formats

- Naive FP16 is not the only option!
- Specialized formats preserve dynamic range for computations



# Switching to lower precision

- FP16 exists since CUDA 8, just allocate the tensor/cast it to `half`
- BF16 is available on CPUs and TPUs [1], `Tensor.bfloat16()` in PyTorch
- TF32 can be enabled for you on Ampere GPUs  
(was enabled in PyTorch by default until 1.12)
  - Never exposed as a data type, only as a type for specific operations [2]



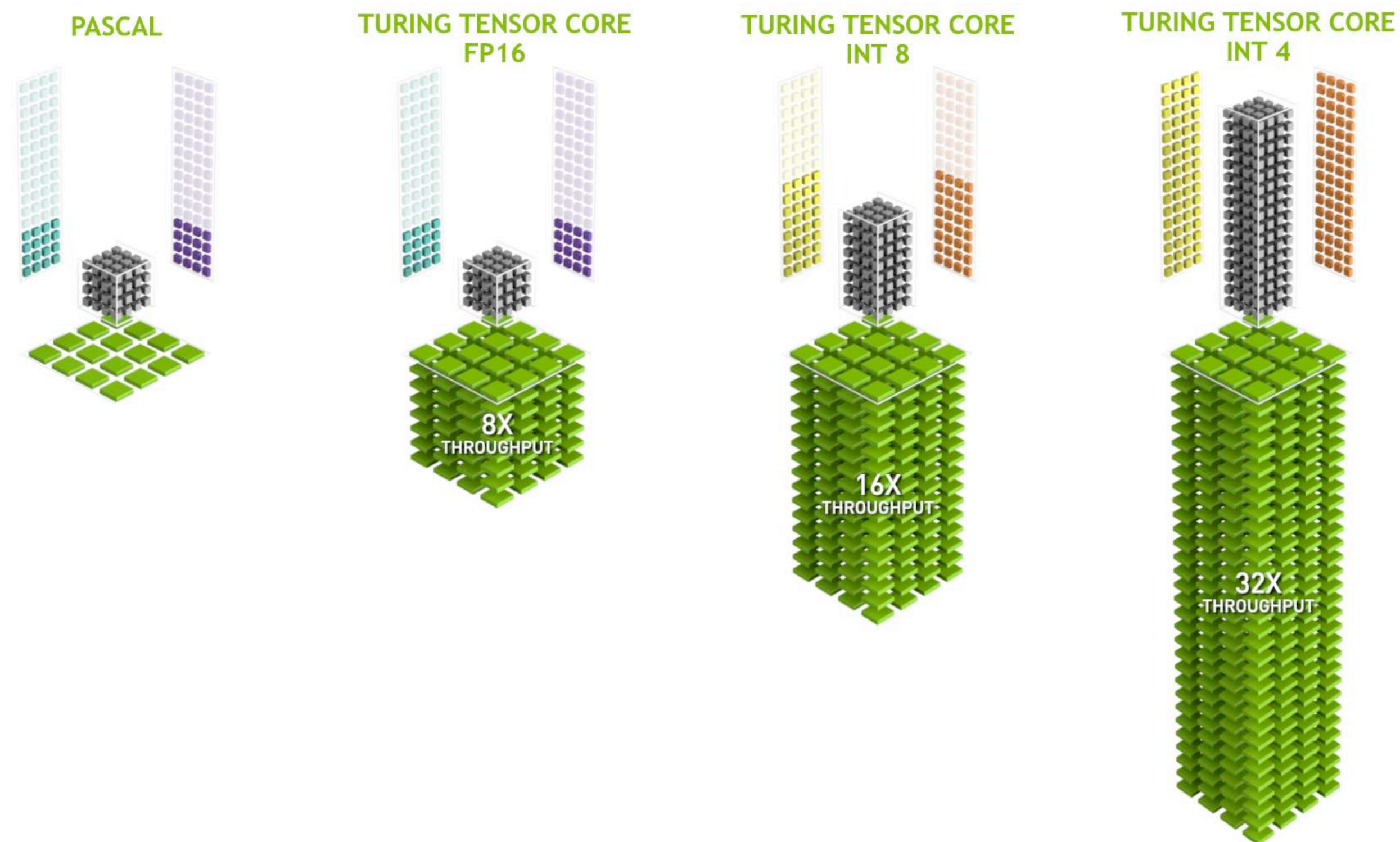
[1] [pytorch.org/xla/release/1.9/index.html#xla-tensors-and-bfloat16](https://pytorch.org/xla/release/1.9/index.html#xla-tensors-and-bfloat16)

[2] [developer.nvidia.com/blog/accelerating-ai-training-with-tf32-tensor-cores](https://developer.nvidia.com/blog/accelerating-ai-training-with-tf32-tensor-cores)



# Tensor Cores

- Specialized computation units available in latest generations of NVIDIA GPUs (since Volta)
- Allow the user to speed up  $D = A \times B + C$  by up to 8-16x (claimed)



# Tensor Cores

- Specialized computation units available in latest generations of NVIDIA GPUs (since Volta)
- Allow the user to speed up  $D = A \times B + C$  by up to 8-16x (claimed)
- Enabled not only for TF32/FP16/BF16 (Ampere), but even for INT8/INT4
- You do not specify their usage manually!



# Utilizing Tensor Cores

- To enable them, you either need recent CUDA or specific size constraints:

Table 1. Tensor Core requirements by cuBLAS or cuDNN version for some common data precisions. These requirements apply to matrix dimensions M, N, and K.

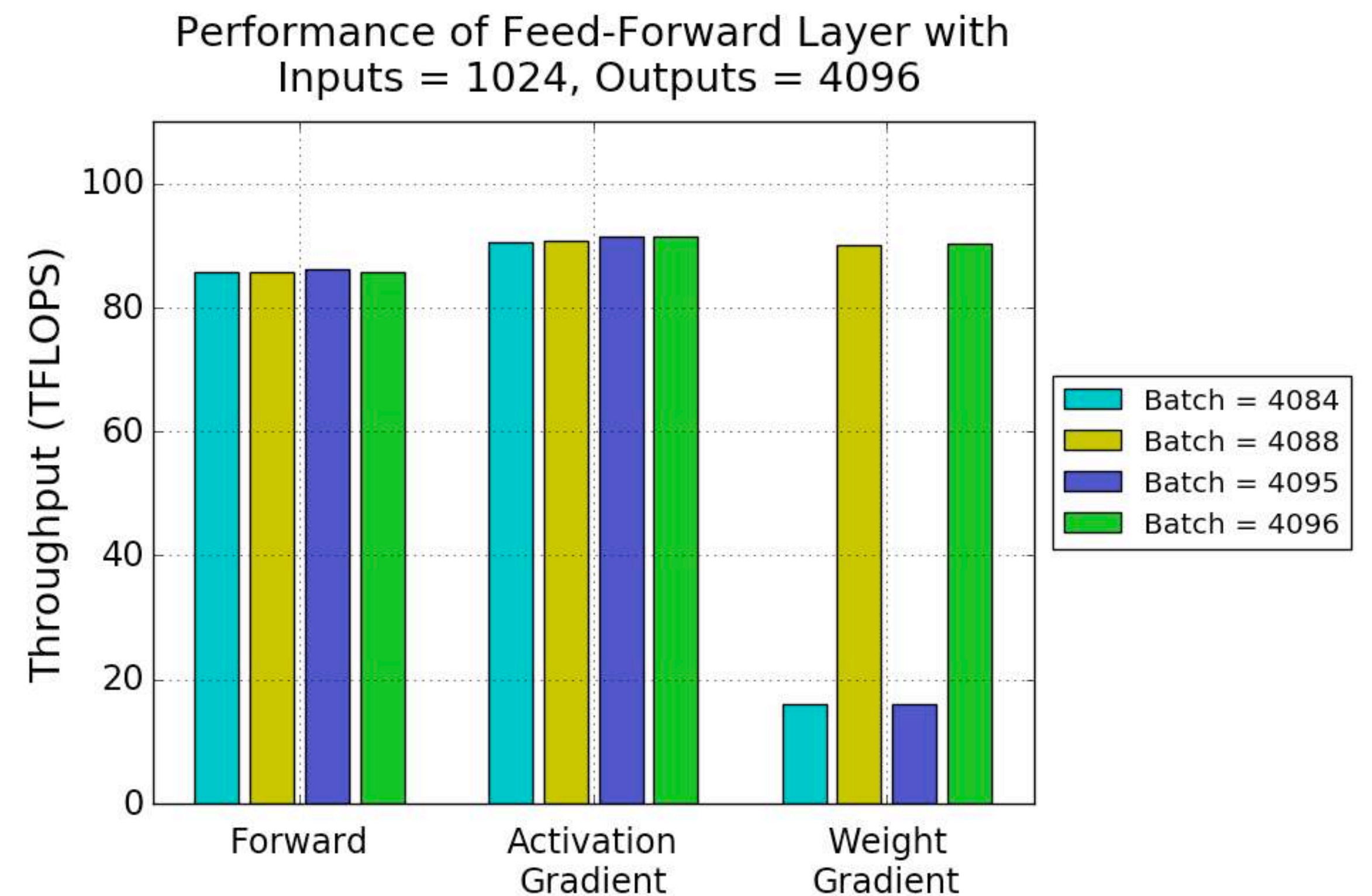
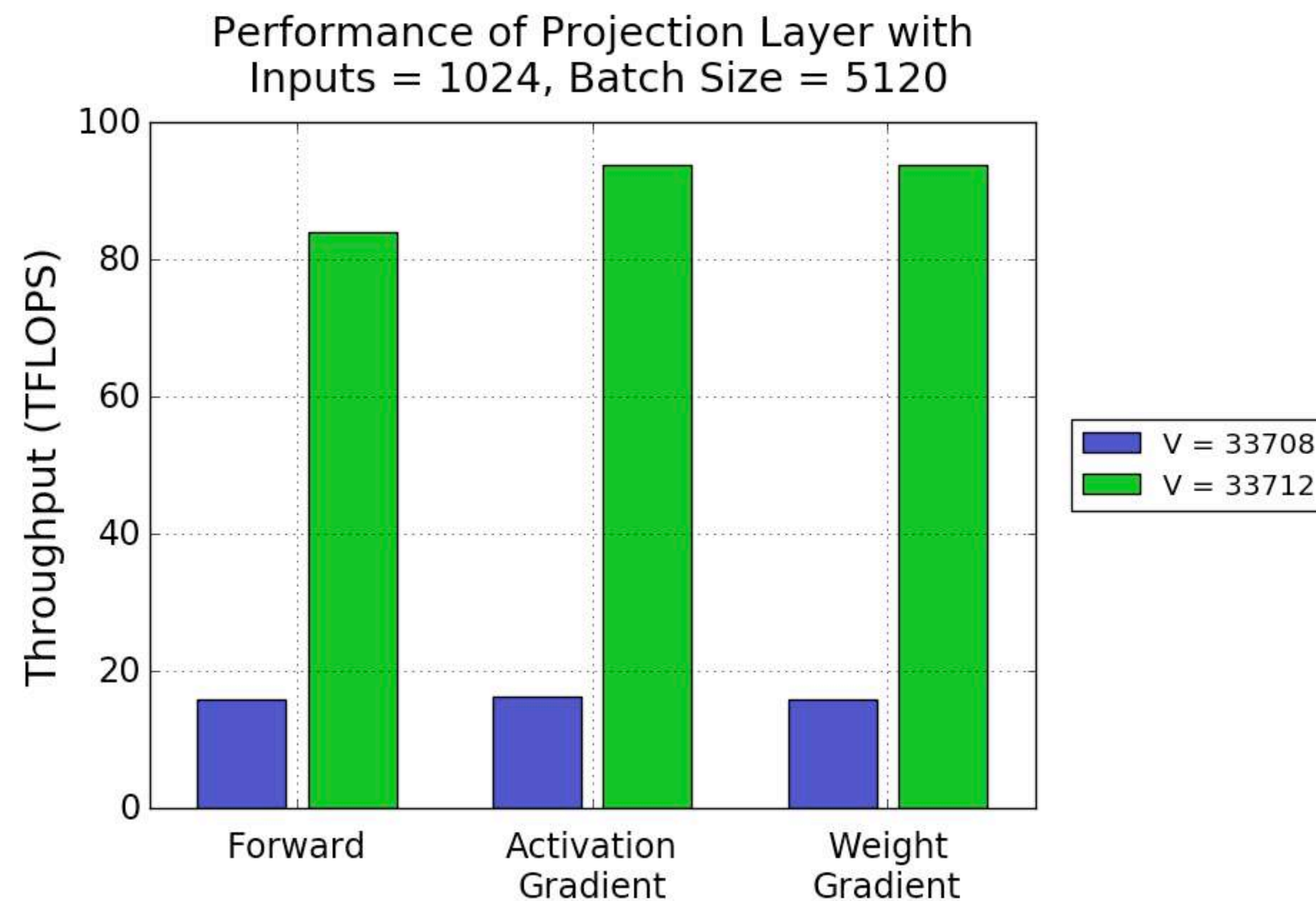
Tensor Cores can be used for...	cuBLAS version < 11.0 cuDNN version < 7.6.3	cuBLAS version ≥ 11.0 cuDNN version ≥ 7.6.3
INT8	Multiples of 16	Always but most efficient with multiples of 16; on A100, multiples of 128.
FP16	Multiples of 8	Always but most efficient with multiples of 8; on A100, multiples of 64.
TF32	N/A	Always but most efficient with multiples of 4; on A100, multiples of 32.
FP64	N/A	Always but most efficient with multiples of 2; on A100, multiples of 16.

[1] [docs.nvidia.com/deeplearning/performance/dl-performance-matrix-multiplication/index.html#requirements-tc](https://docs.nvidia.com/deeplearning/performance/dl-performance-matrix-multiplication/index.html#requirements-tc)

[2] <https://developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/s9926-tensor-core-performance-the-ultimate-guide.pdf>

# Utilizing Tensor Cores

- To enable them, you either need recent CUDA or specific size constraints:



[1] [docs.nvidia.com/deeplearning/performance/dl-performance-matrix-multiplication/index.html#requirements-tc](https://docs.nvidia.com/deeplearning/performance/dl-performance-matrix-multiplication/index.html#requirements-tc)

[2] <https://developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/s9926-tensor-core-performance-the-ultimate-guide.pdf>

# Utilizing Tensor Cores

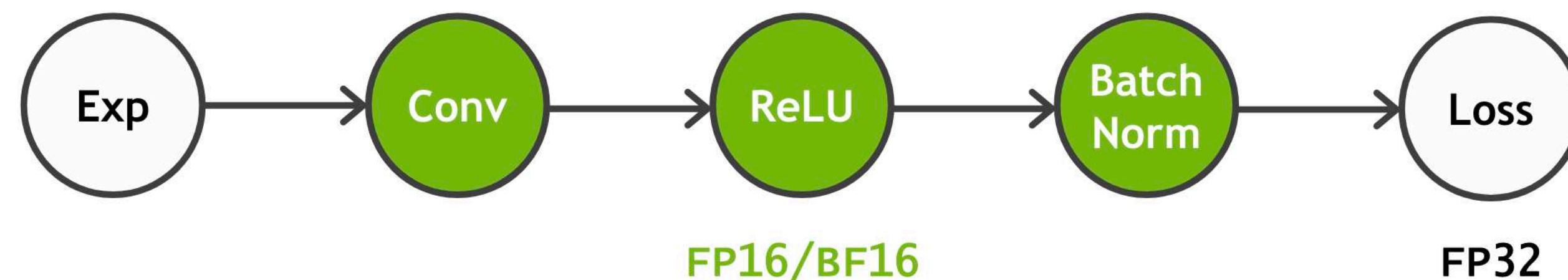
- To enable them, you either need recent CUDA or specific size constraints:
- Run GPU profiler to check if they are used ([i|s|h](\d)+ in kernel names)
- Also, DL profilers can indicate Tensor Core eligibility and usage

[1] [docs.nvidia.com/deeplearning/performance/dl-performance-matrix-multiplication/index.html#requirements-tc](https://docs.nvidia.com/deeplearning/performance/dl-performance-matrix-multiplication/index.html#requirements-tc)

[2] <https://developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/s9926-tensor-core-performance-the-ultimate-guide.pdf>

# Mixed precision training

- Training in pure FP16 hardly works
- Some operations (e.g. matrix multiplication) can work, others (softmax, batch normalization) need higher precision
- Mixed precision training casts layer activations to appropriate data types
- Supported in popular DL frameworks (e.g. torch.cuda.amp)

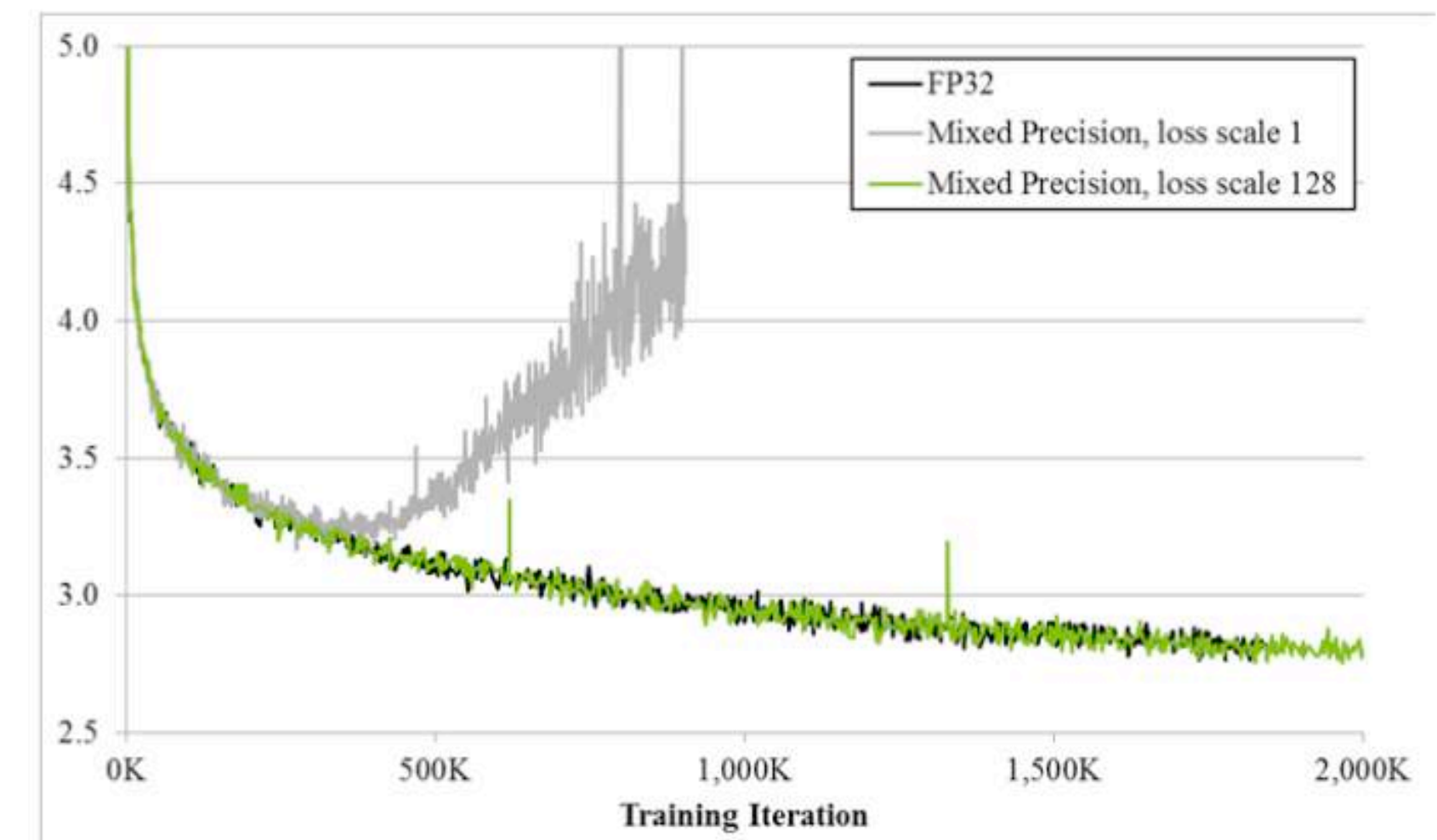
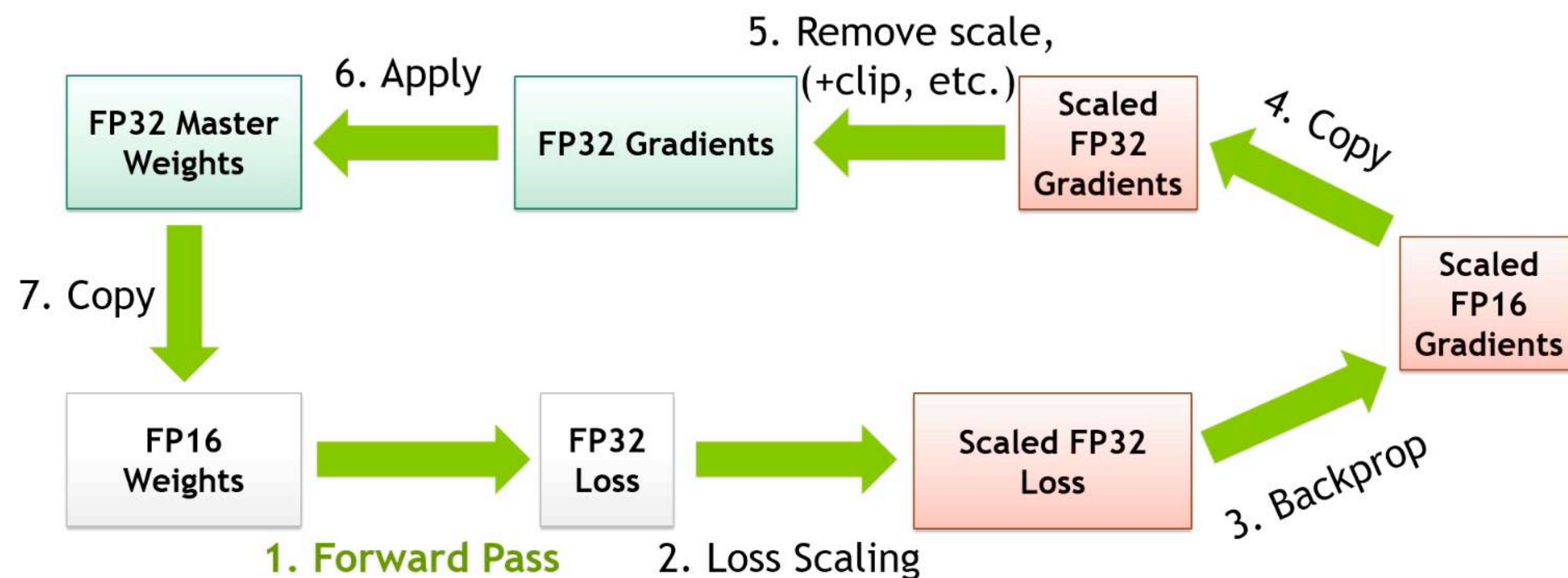


- Increases the training throughput due to the use of Tensor Cores
- Decreases the memory usage by half... or not?



# Loss scaling

- To prevent underflows, we need to scale the FP16 gradients by a small number
- Ironically, this can lead to overflows when unscaling
- Dynamic loss scaling detects such overflows and repeatedly halves the scale



# Memory savings of AMP

- Let's count the number of bytes per parameter for standard training with Adam:

FP32:

- Parameters — 4 bytes
- Gradients — 4 bytes
- Optimizer statistics — 8 bytes

**16 bytes per parameter** in total

AMP:

- Parameters — 2 bytes
- Master parameters — 4 bytes
- Gradients — 2 bytes (sometimes 4)
- Optimizer statistics — 8 bytes

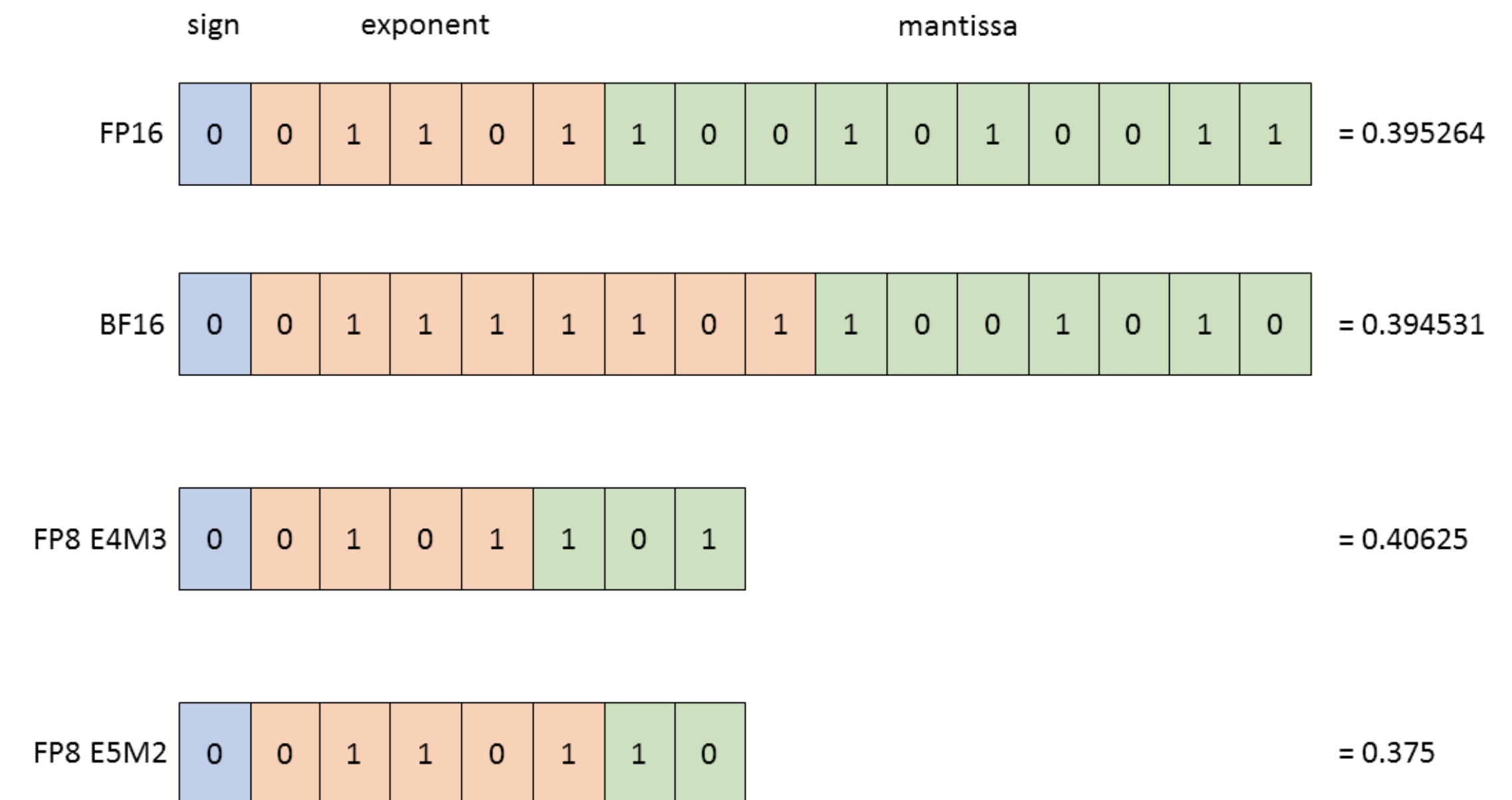
**Also 16 bytes per parameter!**

- The only major savings come from reduced activation memory



# FP8 training

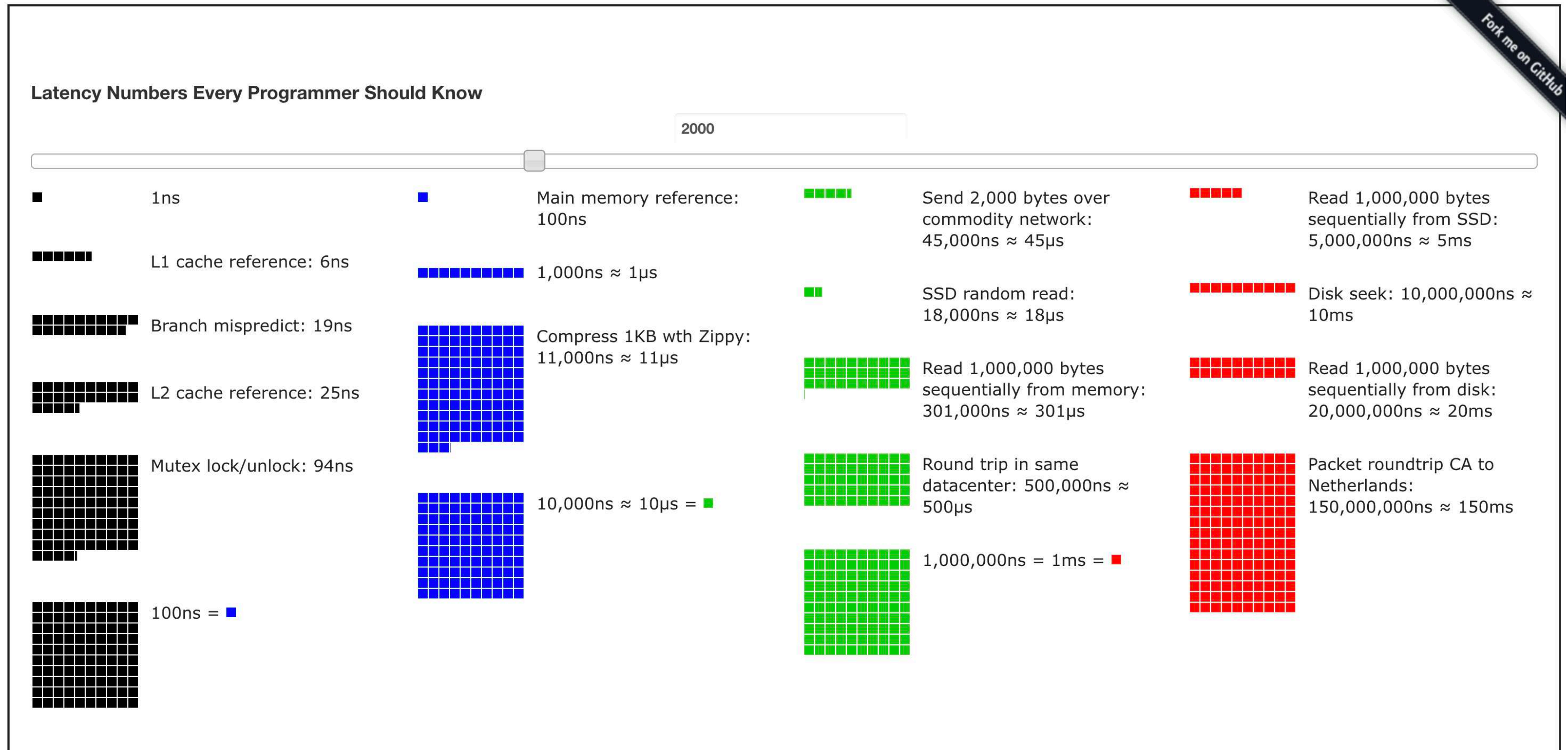
- On latest hardware (e.g., H100), we have even lower precision formats
- E4M3 is used for weights and activations, E5M2 is best for gradients
- Extra tricks (per-tensor scaling) required to maintain accuracy
- Use [github.com/NVIDIA/TransformerEngine](https://github.com/NVIDIA/TransformerEngine) to leverage this
- In PyTorch (alpha): [github.com/pytorch-labs/float8\\_experimental](https://github.com/pytorch-labs/float8_experimental)



# AMP: takeaways

- Use more efficient data types when available
- Mind the sizes/operation types
- Don't expect significant memory savings for large models
- In many cases, this is easy to integrate through standard tools

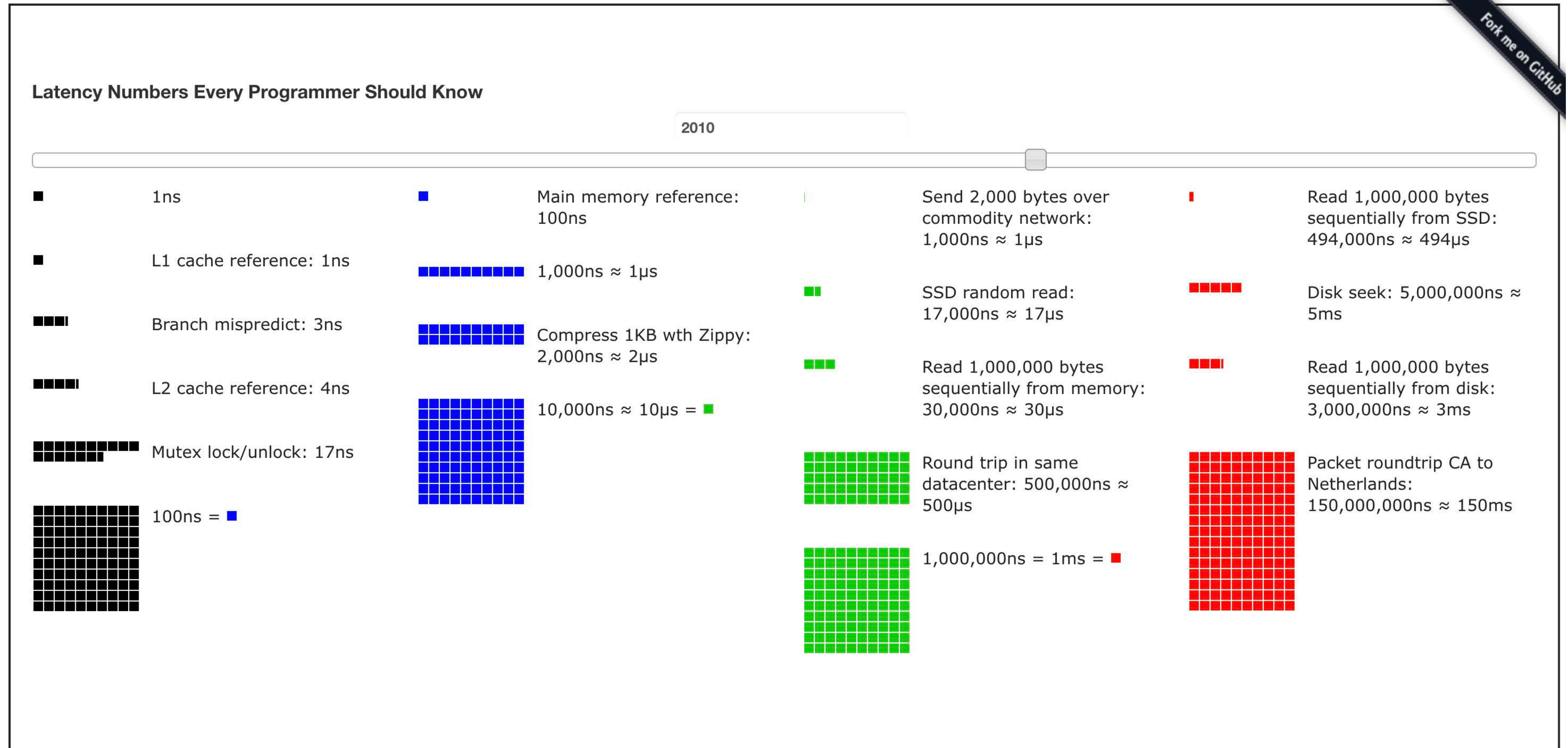
# Bottlenecks in data loading



Fork me on GitHub



# Bottlenecks in data loading



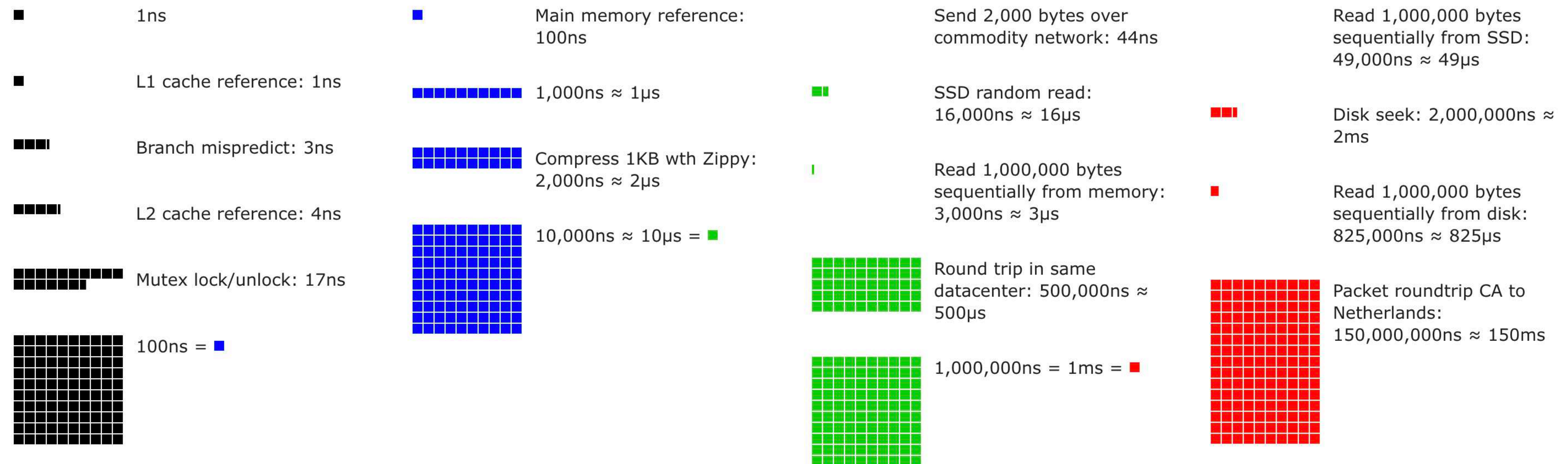
Fork me on GitHub

# Bottlenecks in data loading

Fork me on GitHub

## Latency Numbers Every Programmer Should Know

2020



# Bottlenecks in data loading

- Sometimes the models aren't so compute-intensive...
- We still want to process the data efficiently!
- Need to be mindful of hardware/network performance and the CPU code
- Two components: what to read and how to read
- Obvious part: read data in parallel  
(several processes, asynchronously with computation)



# Storage formats

- Raw files are often easy to visualize, but storage-inefficient (especially when accessing external storage)
- In some cases, you might benefit from better formats:
  - For structured data, Apache Arrow/Protobuf/msgpack etc.
  - For images, apply non-random “heavy” processing before training
  - For language data, tokenize the texts and store integer indices only

# Minimizing preprocessing time

- Reading the data and feeding it into the model can also be slow
  - For large images, you can be bound by CPU operations
  - For sequence data, you can waste time on padding tokens

# Performance of image loading

- When reading images, consider the code that reads them :)
- Default PIL.Image.Open can be highly inefficient!  
Use at least Pillow-SIMD
- Use better decoders (e.g. jpegturbo, nvJPEG from DALI)



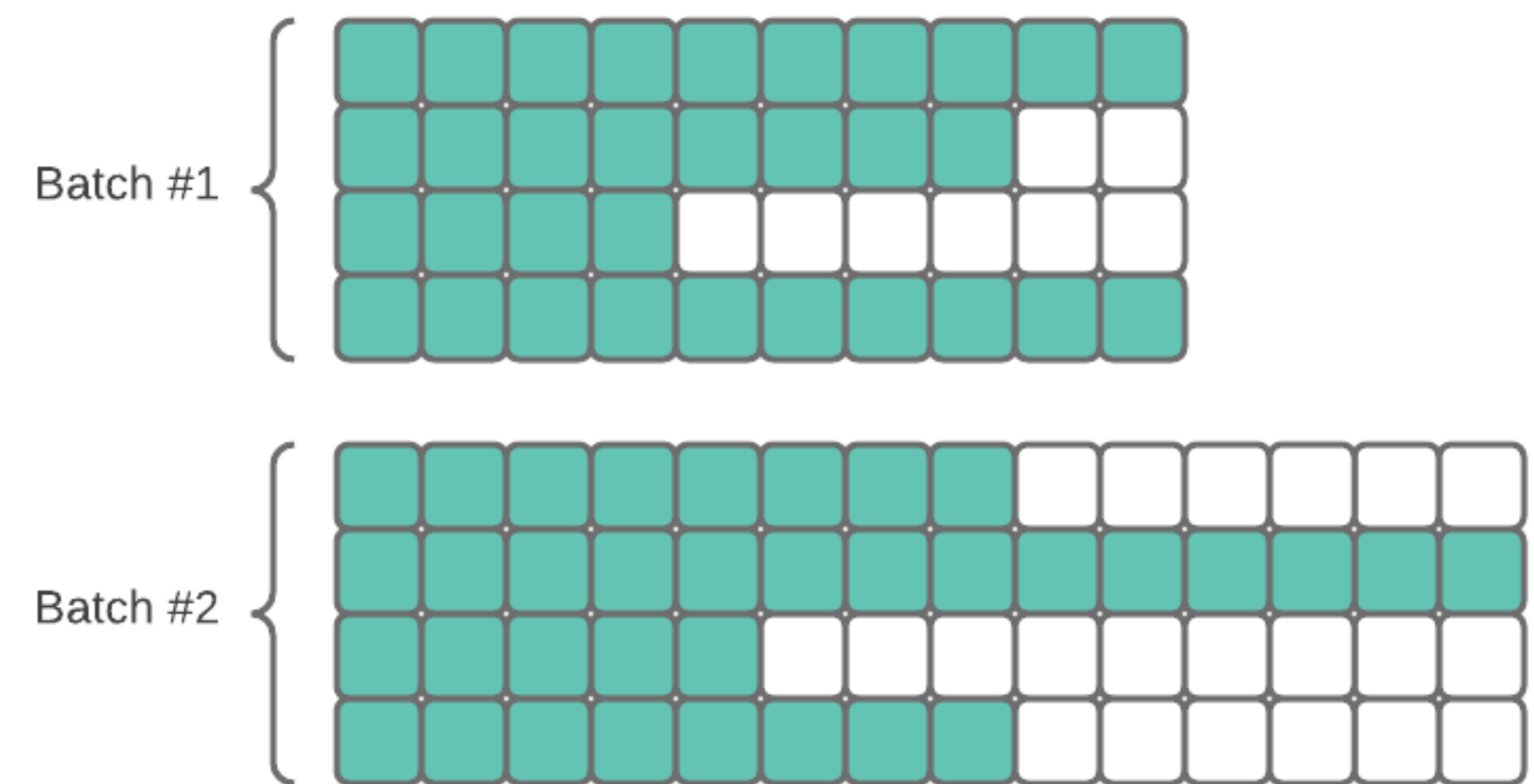
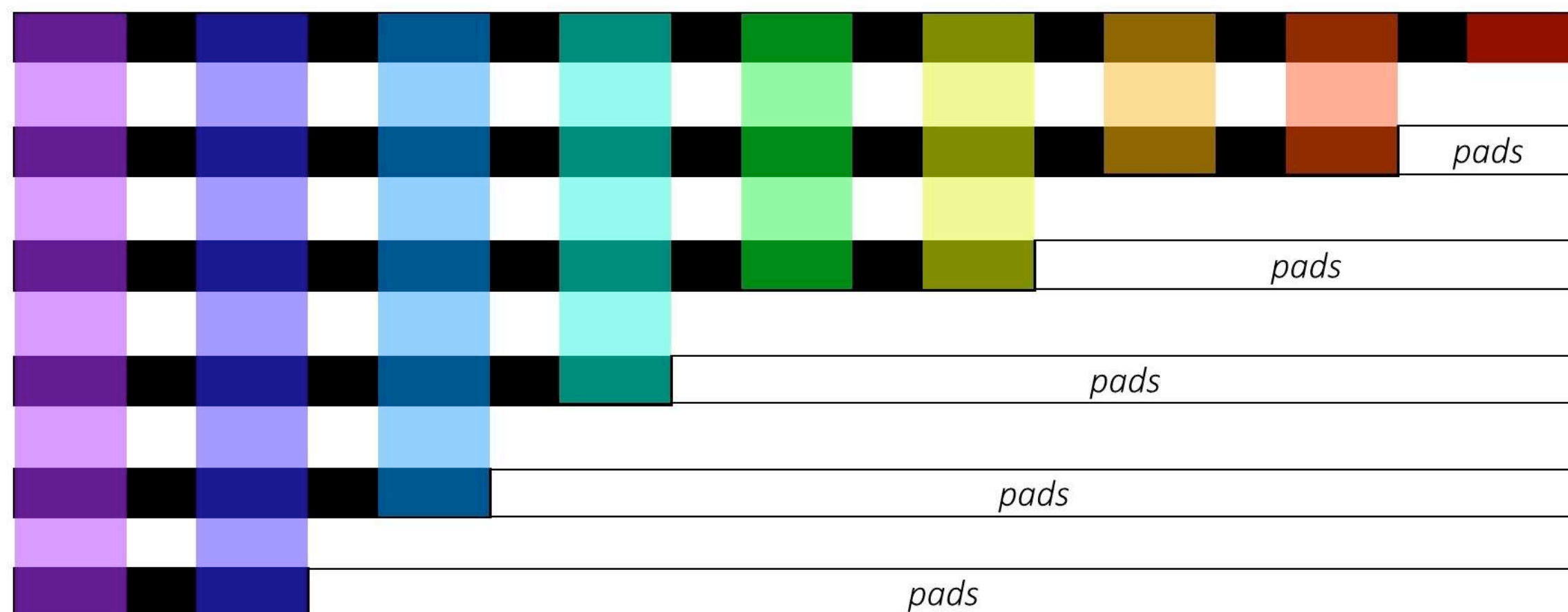
# Performance of image loading

- When reading images, consider the code that reads them :)
- Default PIL.Image.Open can be highly inefficient!  
Use at least Pillow-SIMD
- Use better decoders (e.g. jpegturbo, nvJPEG from DALI)
- Heavy groups of augmentations can also slow you down
  - Consider moving them to GPU (e.g. kornia, DALI)
  - In most cases, you can switch to efficient implementations

# Optimal sequence processing

- For sequential data, padding in batches is necessary
- However, padding the ENTIRE dataset can lead to redundant timesteps
- It's usually better to store samples without padding and use `collate_fn`
- Also, bucket examples by length to further minimize padding

*Padded sequences sorted by decreasing lengths*



# Data pipelines: takeaways

- Consider the performance/size of your storage when loading the data
- Use better deserialization primitives when available
- Try to avoid obvious inefficiencies when building task-specific pipelines

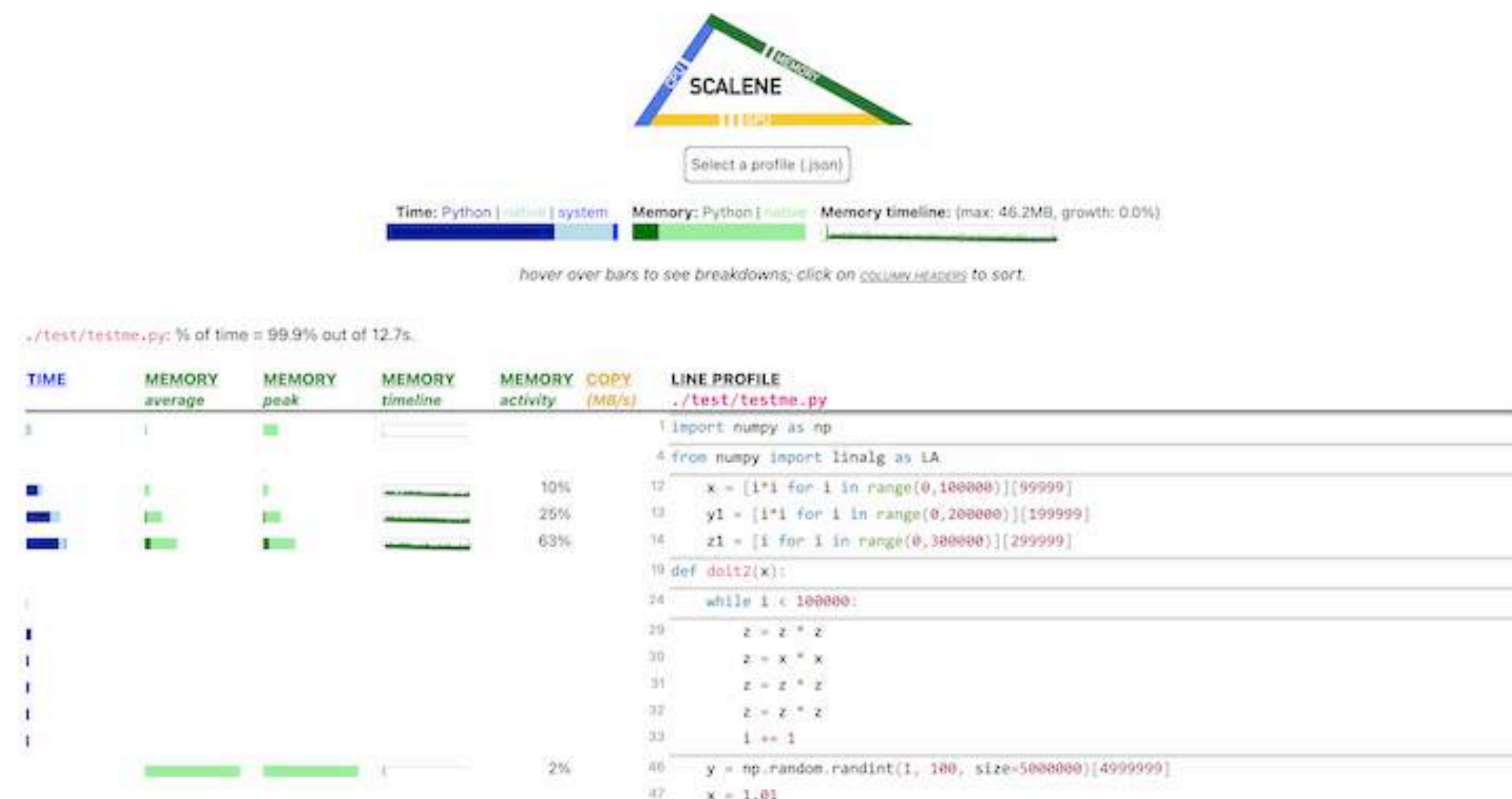


# Profiling: what and why

- In benchmarking, we measure the speed of our program as a black box
- Profiling is a process of determining the runtime of parts of your program
- More of a “white box” approach

# How to profile Python code?

- cProfile as a standard tool built into Python
- Sampling-based profilers (scalene etc.)
- Some of them (e.g. py-spy) even allow to attach to running code!



Collecting samples from 'python examples/lastfm.py' (python v3.6.3)  
Total Samples 300  
GIL: 24.33%, Active: 100.00%, Threads: 2

%Own	%Total	OwnTime	TotalTime	Function (filename:line)
36.00%	36.00%	0.360s	0.360s	fit (implicit/als.py:139)
26.00%	26.00%	0.260s	0.260s	fit (implicit/als.py:141)
23.00%	23.00%	0.230s	0.230s	fit (implicit/als.py:160)
13.00%	13.00%	0.580s	0.580s	tocsr (scipy/sparse/csc.py:151)
1.00%	1.00%	0.120s	0.230s	_call_with_frames_removed (<frozen importlib._bootstrap>:219)
1.00%	1.00%	0.010s	0.010s	__subclasscheck__ (abc.py:231)
0.00%	0.00%	0.040s	0.040s	__init__ (scipy/sparse/coo.py:159)
0.00%	0.00%	0.000s	0.040s	<module> (numpy/lib/__init__.py:8)
0.00%	0.00%	0.000s	0.010s	<module> (h5py/tests/common.py:37)
0.00%	0.00%	0.000s	0.010s	<module> (email/parser.py:12)
0.00%	0.00%	0.130s	0.130s	bm25_weight (implicit/nearest_neighbours.py:149)
0.00%	0.00%	0.000s	0.010s	<module> (tqdm/_tqdm.py:24)
0.00%	0.00%	0.000s	0.010s	<module> (h5py/tests/old/__init__.py:4)
0.00%	0.00%	0.030s	0.030s	__init__ (scipy/sparse/coo.py:158)
0.00%	0.00%	0.000s	0.170s	__init__ (scipy/sparse/coo.py:170)
0.00%	0.00%	0.000s	0.020s	_check (scipy/sparse/coo.py:277)
0.00%	0.00%	0.050s	0.050s	_mul_vector (scipy/sparse/coo.py:572)
0.00%	0.00%	0.000s	0.210s	bm25_weight (implicit/nearest_neighbours.py:146)
0.00%	0.00%	0.000s	0.010s	_find_spec (<frozen importlib._bootstrap>:914)
0.00%	0.00%	0.010s	0.010s	_get_default_tempdir (tempfile.py:216)
0.00%	0.00%	0.000s	0.160s	_handle_fromlist (<frozen importlib._bootstrap>:1017)
0.00%	1.00%	0.000s	0.210s	_find_and_load_unlocked (<frozen importlib._bootstrap>:955)
0.00%	0.00%	0.000s	0.010s	<module> (numpy/compat/py3k.py:14)

Press **Control-C** to quit, or **?** for help.



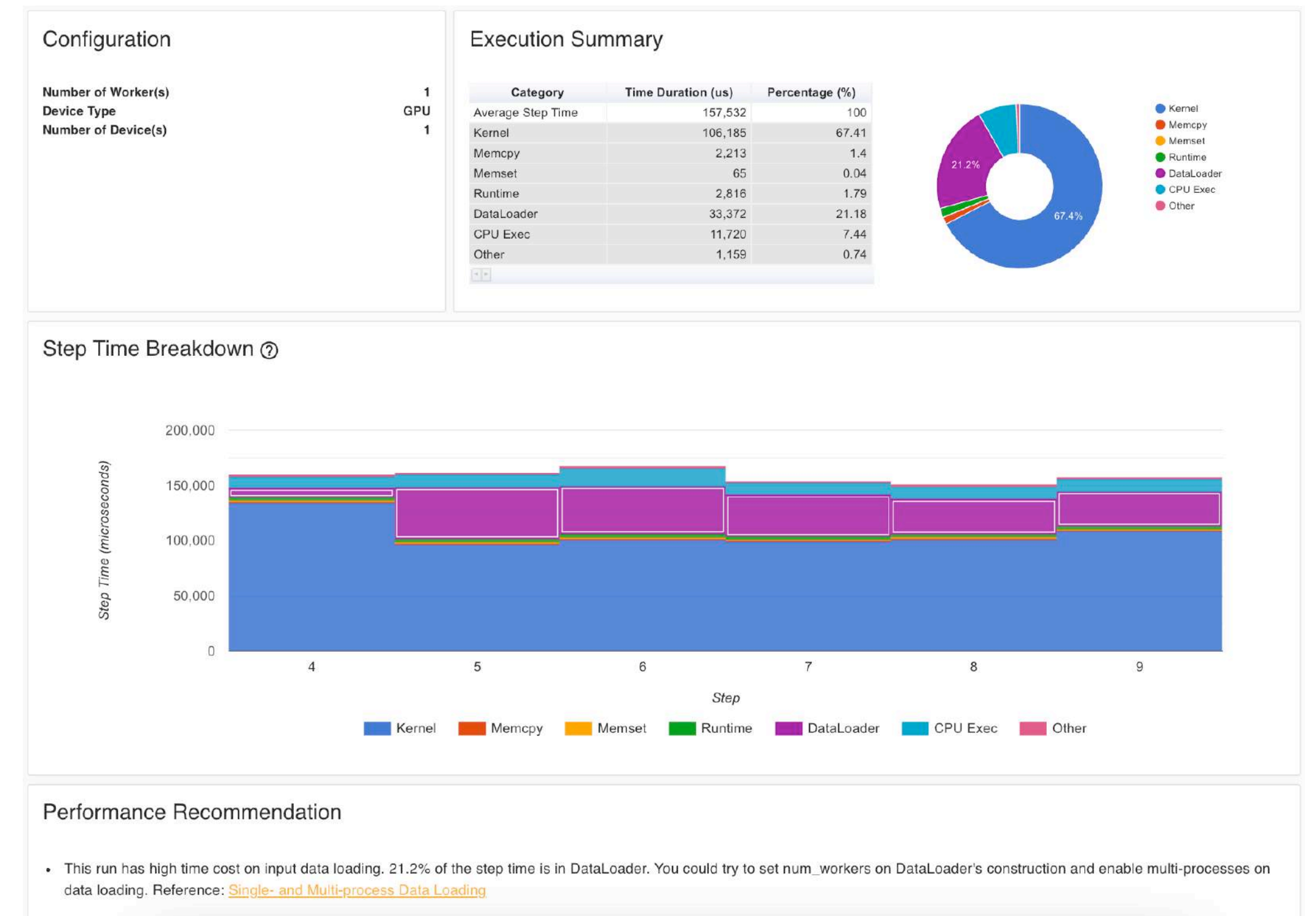
# How to profile GPU code?

- nvprof is the low-level profiling tool
- Gives you the performance of low-level kernel launches and copies

```
==9261== Profiling application: ./tHogbomCleanHemi
==9261== Profiling result:
Time(%)      Time      Calls      Avg      Min      Max      Name
58.73%    737.97ms    1000    737.97us  424.77us  1.1405ms subtractPSFLoop_kernel(float co
38.39%    482.31ms    1001    481.83us  475.74us  492.16us findPeakLoop_kernel(MaxCandidat
 1.87%     23.450ms         2    11.725ms  11.721ms  11.728ms [CUDA memcpy HtoD]
 1.01%     12.715ms    1002    12.689us  2.1760us  10.502ms [CUDA memcpy DtoH]
```

# How to profile PyTorch code?

- High-level: [torch.utils.bottleneck](#)
- Older API: [torch.autograd.profiler](#)
- Newer one: [torch.profiler](#)





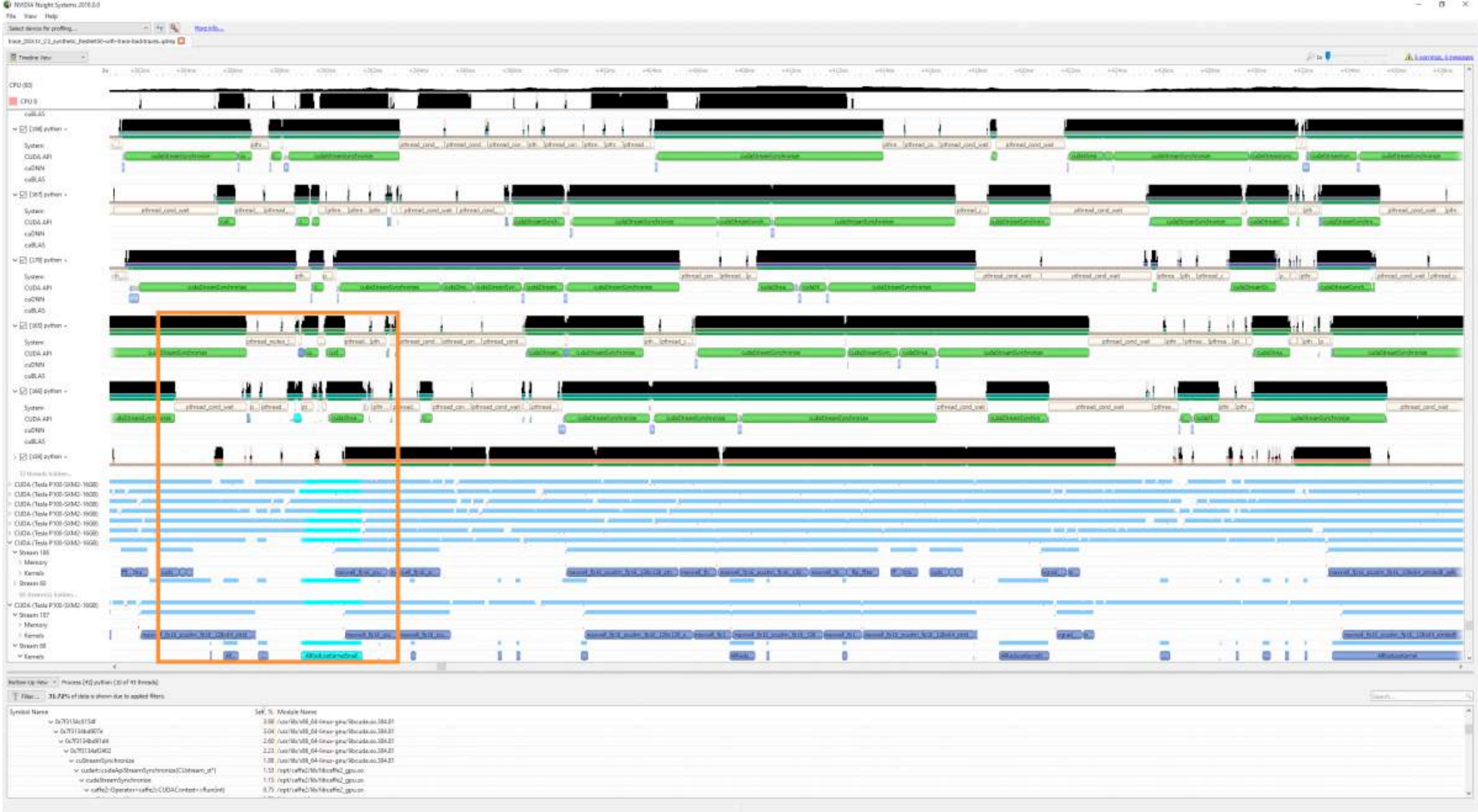
# PyTorch Profiler + TensorBoard



```
with torch.profiler.profile(  
    schedule=torch.profiler.schedule(wait=1, warmup=1, active=3, repeat=2),  
    on_trace_ready=torch.profiler.tensorboard_trace_handler('./log/resnet18'),  
    record_shapes=True,  
    profile_memory=True,  
    with_stack=True  
) as prof:  
    for step, batch_data in enumerate(train_loader):  
        if step >= (1 + 1 + 3) * 2:  
            break  
        train(batch_data)  
        prof.step() # Need to call this at the end of each step to notify profiler of steps'  
boundary.
```



# Nsight Systems/Nsight Compute



**[developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/s9339-profiling-deep-learning-networks.pdf](https://developer.download.nvidia.com/video/gputechconf/gtc/2019/presentation/s9339-profiling-deep-learning-networks.pdf)**

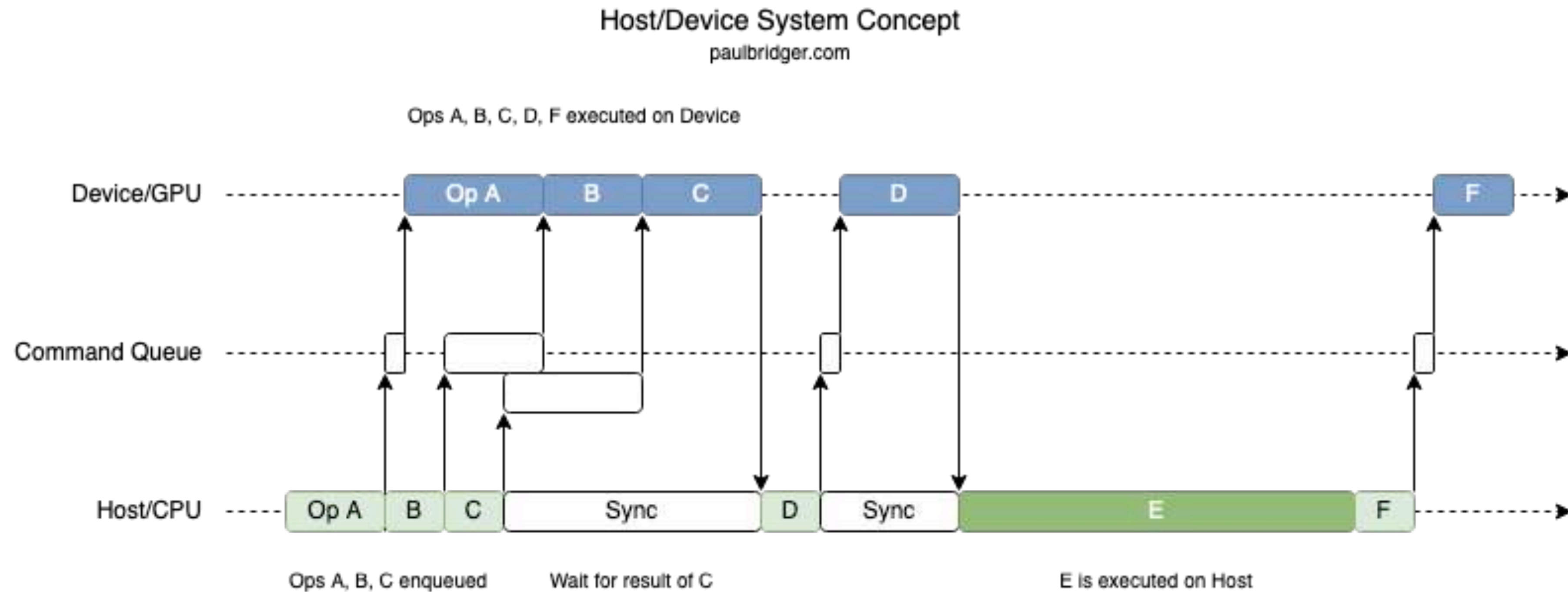


# Profiling: typical patterns





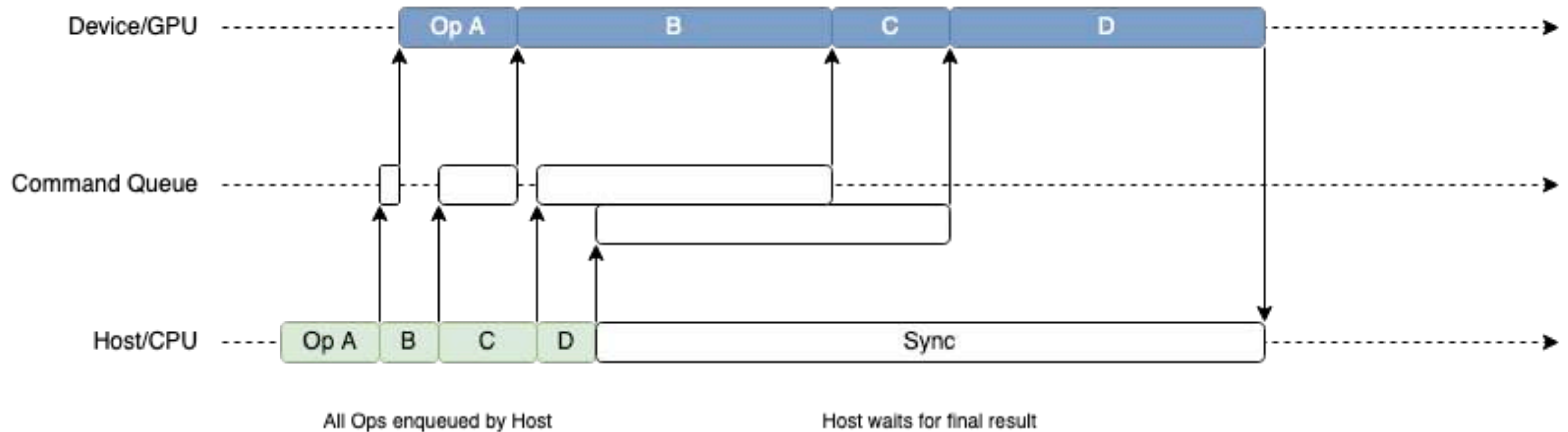
# Profiling: typical patterns



## Pattern: GPU Compute Bound

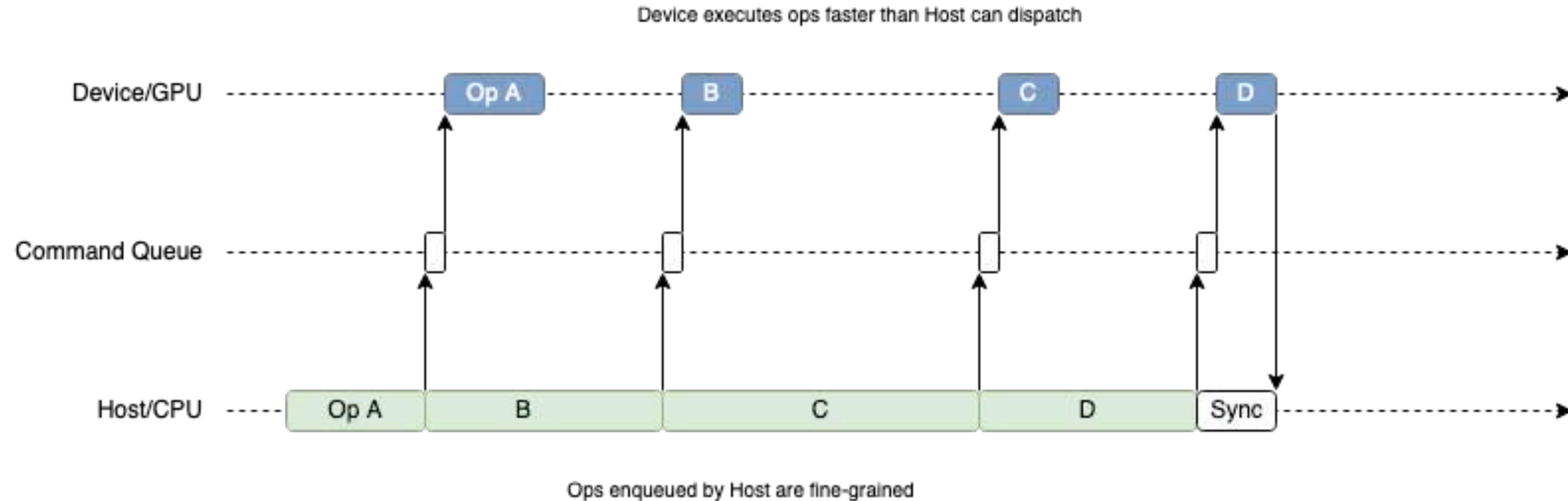
paulbridger.com

All Ops executed by Device



- This is the best case!
- You need to optimize the model itself (lower precision, faster kernels etc)

Pattern: CPU CUDA API Bound  
paulbridger.com

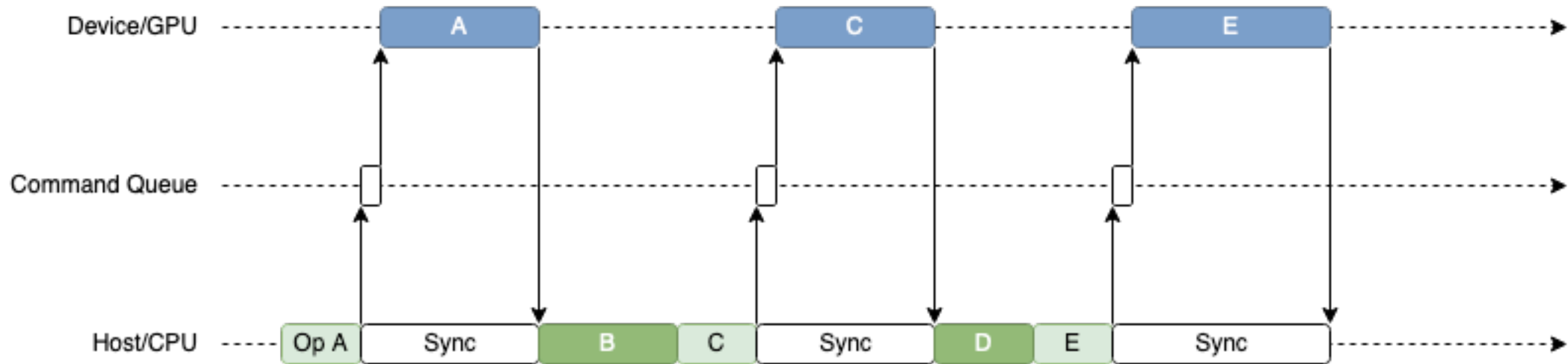


- Operations run faster than kernels are scheduled
- Also happens during inference
- You need to optimize the CUDA API calls (torch.compile, TorchScript) or have more compute-intensive operations (e.g. larger batches)

## Pattern: Synchronization Bound

paulbridger.com

Op computation is interleaved between Device (A, C, E) and Host (B, D)



Because ops are sequentially dependent constant synchronization is required

- CPU and GPU processing are too heavily interleaved
- Remove unnecessary synchronization points, execute as much work on the GPU as possible

# Profiling: takeaways

- A very useful tool for understanding the performance of your pipeline
- Can be applied to both CPU and GPU code
- Depending on the required granularity of measurements, you can use different approaches