



Natural Language Processing: Efficiency

HSE Faculty of Computer Science
Machine Learning and Data-Intensive Systems

Murat Khazhgeriev



Table of Content

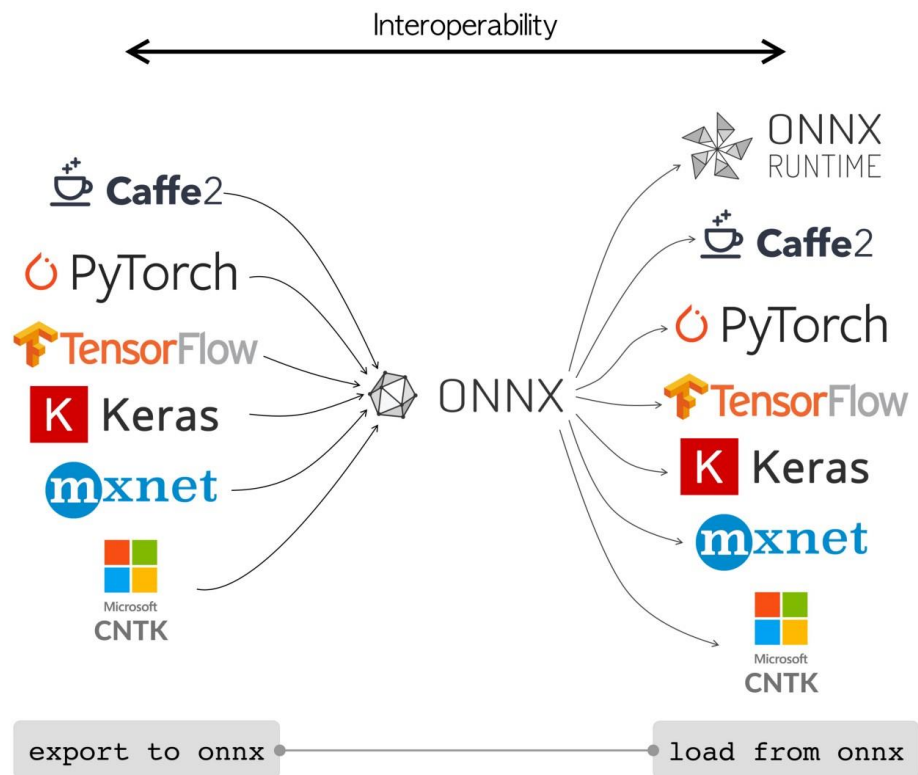
- Just finished training, now what?
- Engineering tricks
- Architectural tricks
- Going distributed



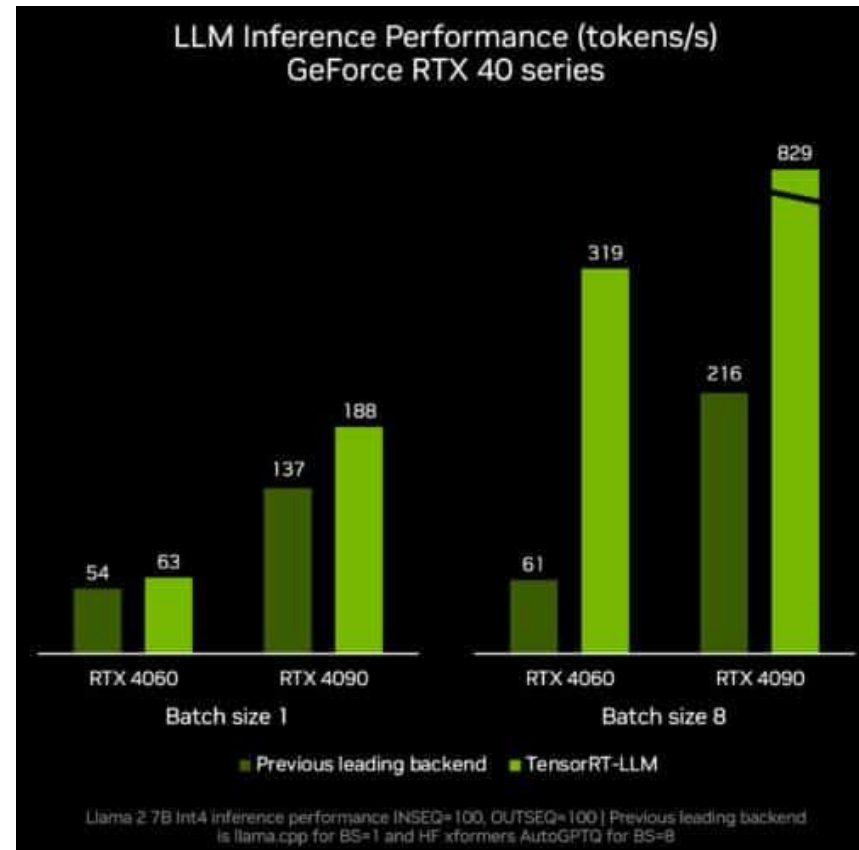
Table of Content

- **Just finished training, now what?**
- Engineering tricks
- Architectural tricks
- Going distributed

Move to engine-agnostic format ONNX



TensorRT if you have Nvidia GPUs



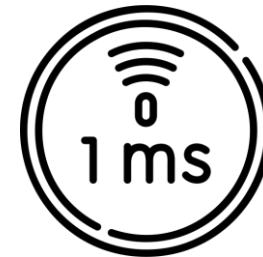
Know your enemy



Model size
(do we fit in RAM?)

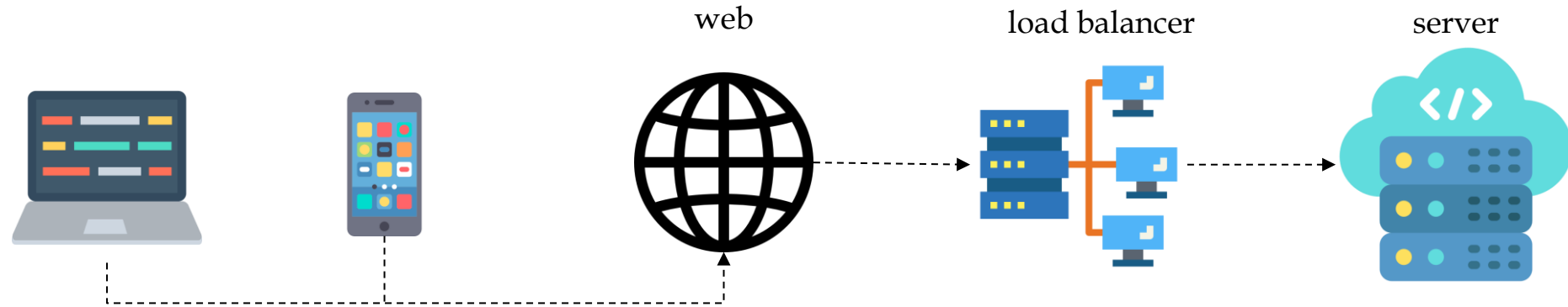


Throughput
(how many queries per second?)



Latency
(how much ms @ percentile till
first token?)

Inference server



Throughput is often critical



Monitor latency @ different percentiles



Key links:

Prototyping: [TorchServe](#)

If you don't like Nvidia: [vLLM](#)

Optimized for Nvidia GPUs: [Triton Inference Server](#)

Mobile

Memory is often a **huge** problem



Watch out for high latency



Key links:

In browsers: [TensorFlow.js](https://www.tensorflow.org/js)

Mobile Platform-agnostic [ExecuTorch](https://pytorch.org/mobile/)

Apple: [CoreML](https://developer.apple.com/machine-learning/)

Android: [NNAPI](https://developer.android.com/training/machine-learning/)

Workstation

Memory may be a problem, but often not



Watch out for latency



Key links:

1: [LLaMA.CPP](#)

2: [OLLAMA](#)

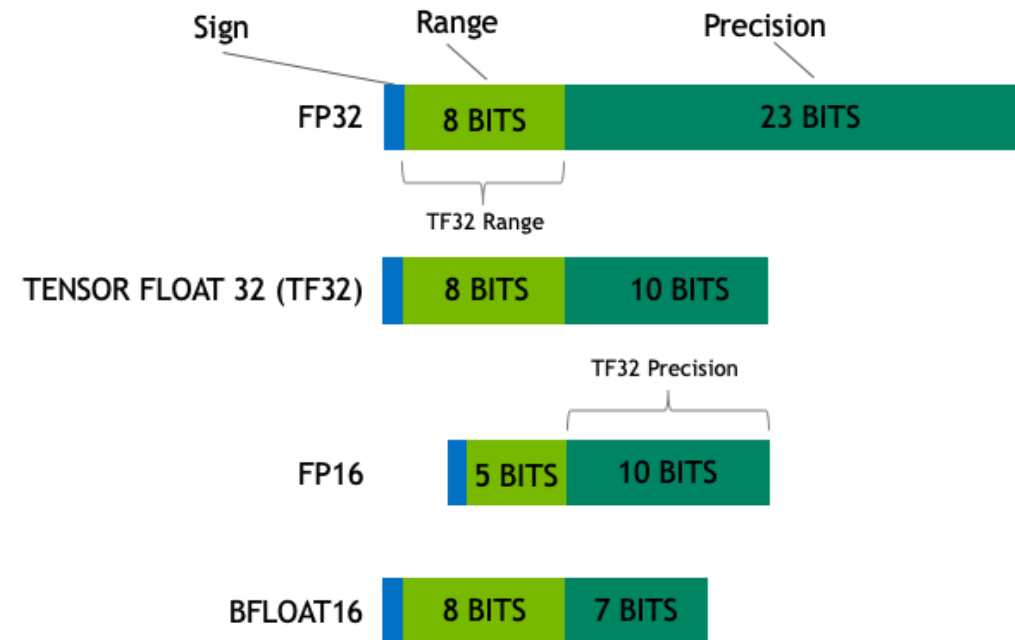
3: [ExLLaMA v2](#)



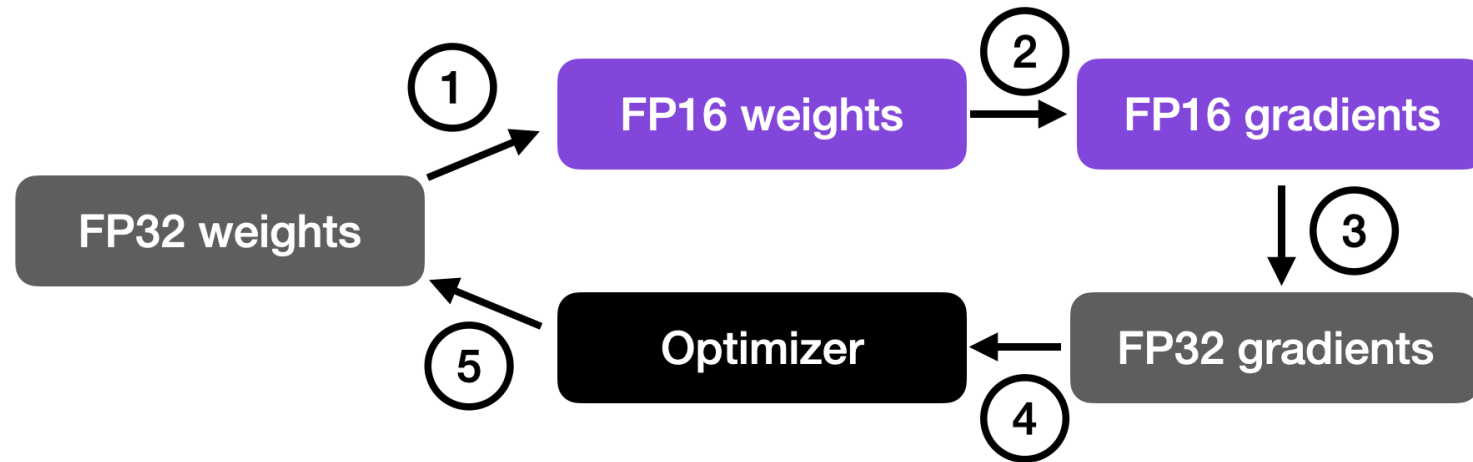
Table of Content

- Just finished training, now what?
- **Engineering tricks**
- Architectural tricks
- Going distributed

Machine numbers refresher



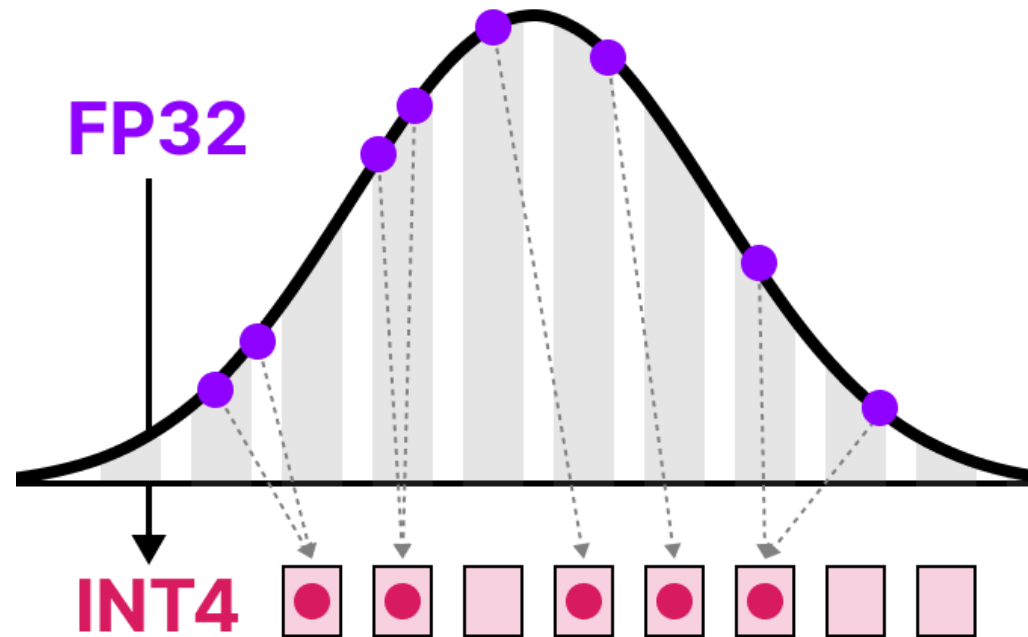
Mixed Precision



Key links:

- 1: [Automatic Mixed Precision in Torch](#)
- 2: [What Every User Should Know About Mixed Precision Training in PyTorch](#)

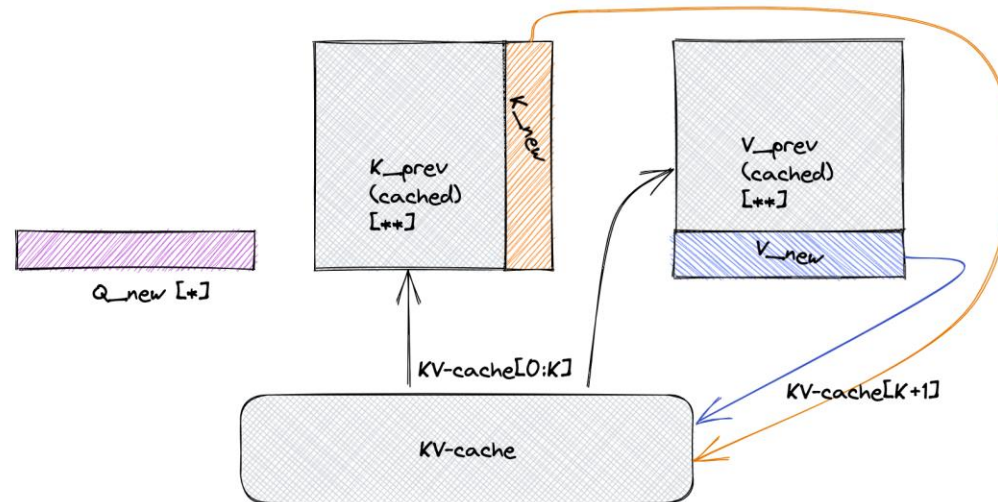
Quantization



Key links:

- 1: [Quantization Recipe](#)
- 2: [HF Optimum](#)
- 3: [Practical Quantization](#)
- 4: [HF bitsandbytes integration](#)

KV-cache



Notes:

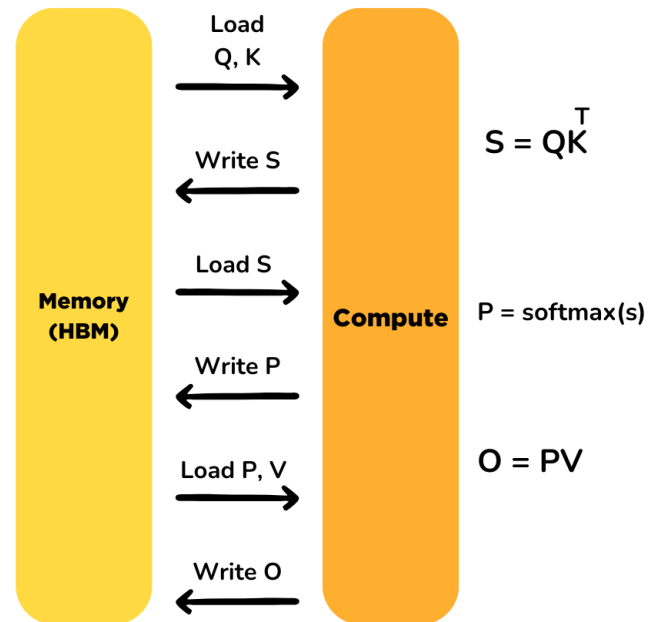
- * When processing token[K], we only need the K'th row of Q
- ** When processing token[K], we require the full K & V tensors, but we can mostly reuse the cached values (This enables skipping the computation of K & V)

Key links:

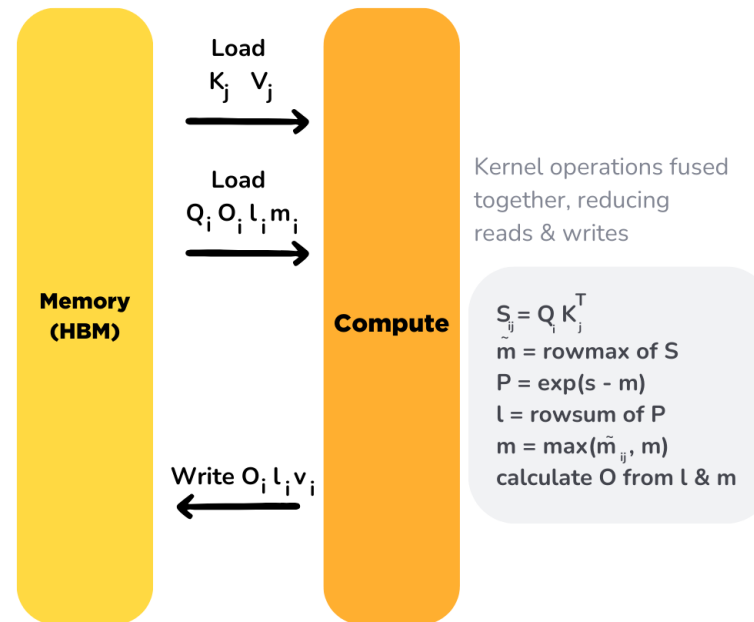
- 1: [Key-Value Cache Quantization](#)

Flash, Grouped-query, Sparse, Paged Attention

Standard Attention Implementation



Flash Attention

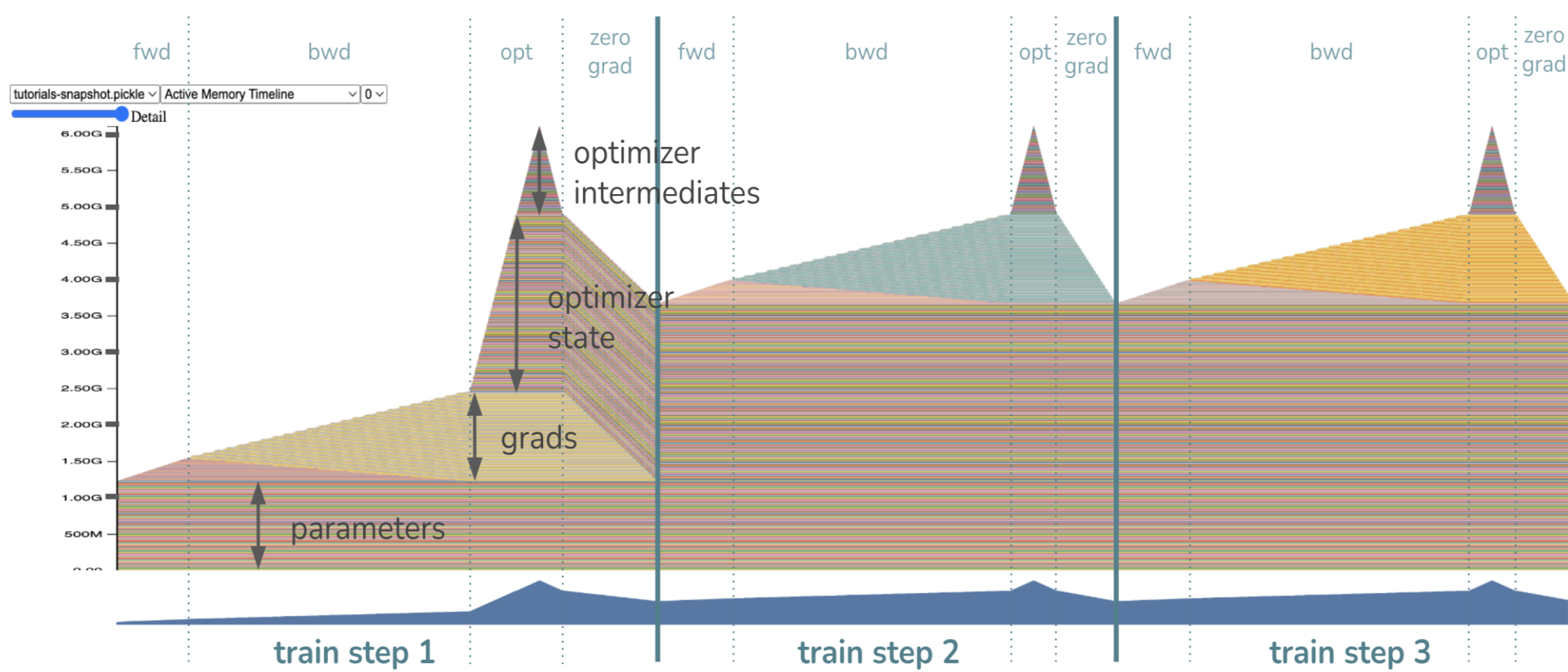


Initialize O, l and m matrices with zeroes. m and l are used to calculate cumulative softmax. Divide Q, K, V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

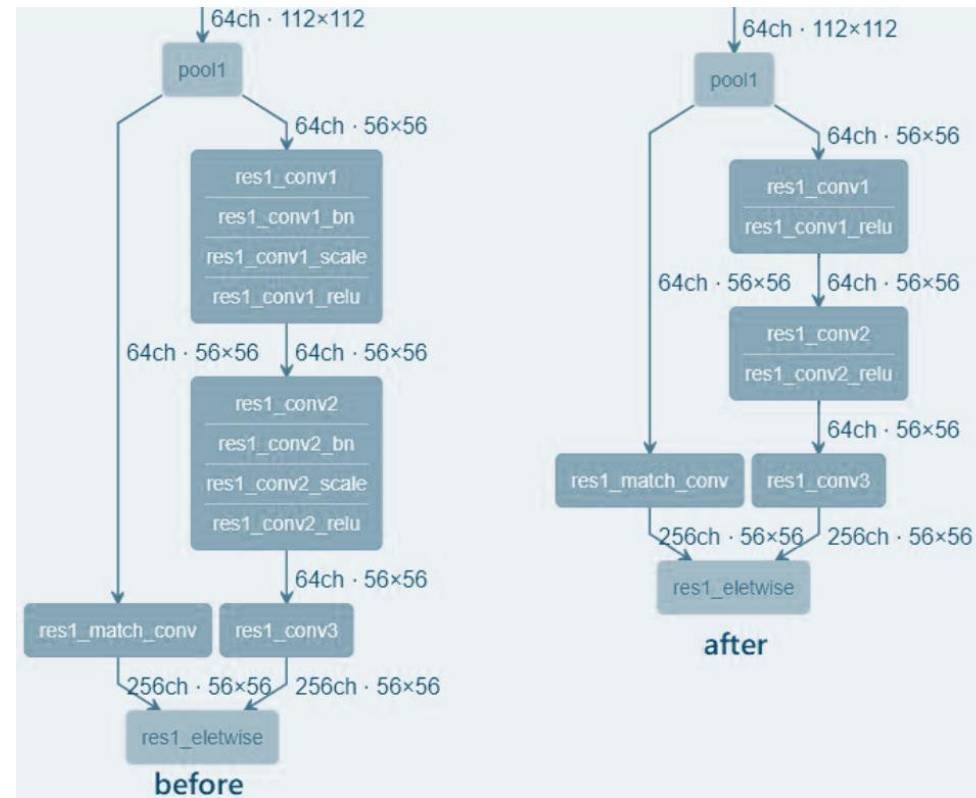
Key links:

- 1: [Grouped-query Attention](#)
- 2: [Flash Attention](#)
- 3: [Sparse Attention](#)
- 4: [Paged Attention](#)

Fuse optimizer's kernels into one



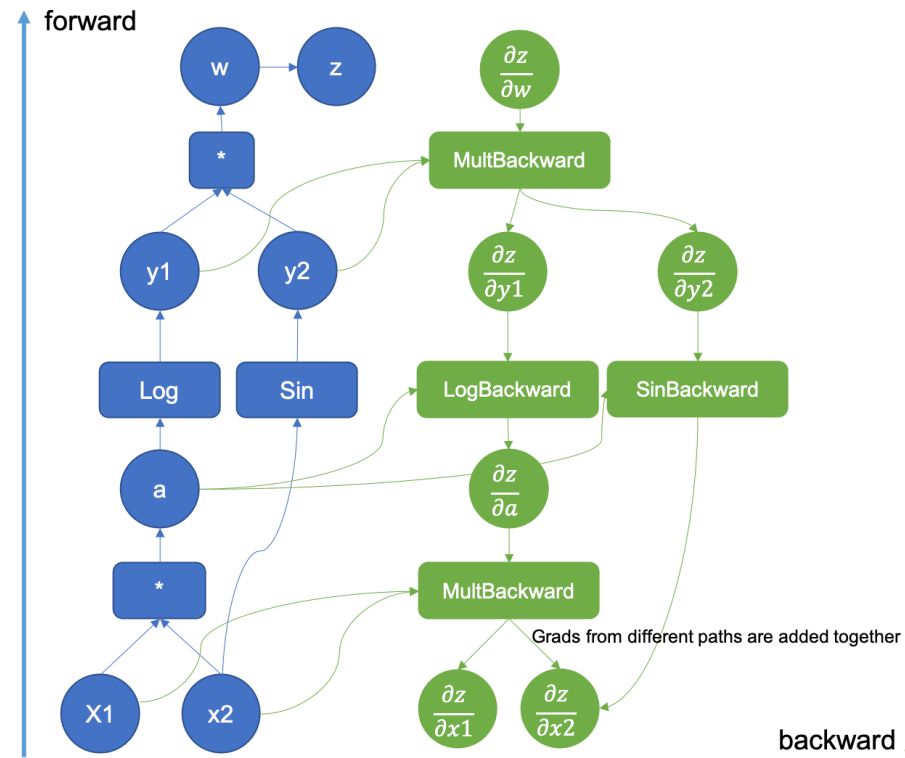
Fuse Torch modules



Key links:

1: [Fuse Modules Recipe](#)

Make Torch's Dynamic Computation Graph Static



Key links:

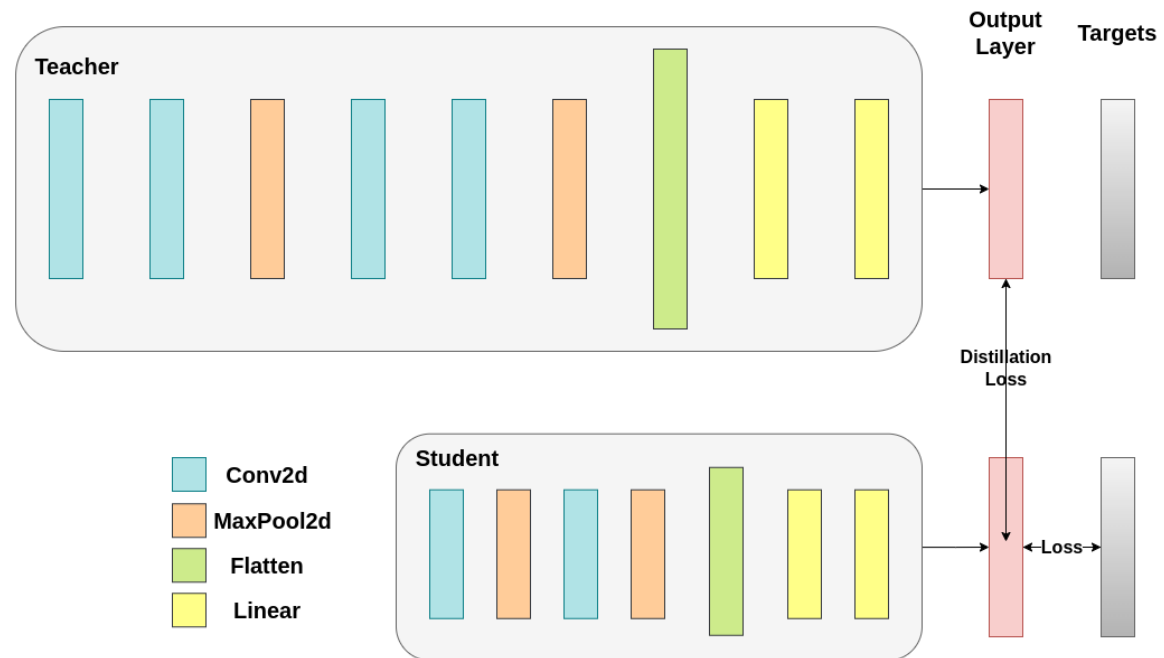
1: [torch.compile](https://pytorch.org/docs/stable/torch_compile.html)



Table of Content

- Just finished training, now what?
- Engineering tricks
- **Architectural tricks**
- Going distributed

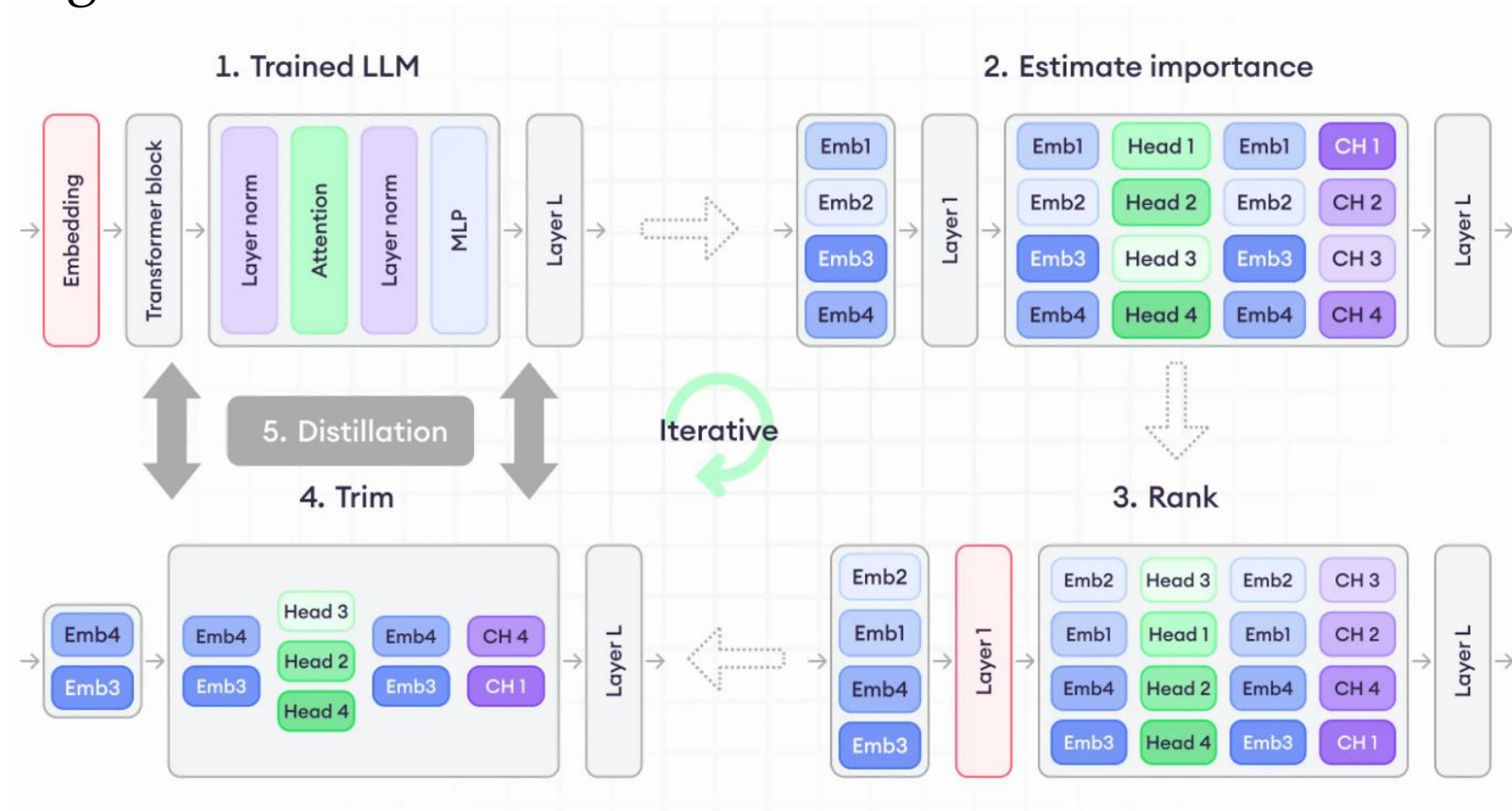
Knowledge Distillation



Key links:

1: [Distillation Tutorial Torch](#)

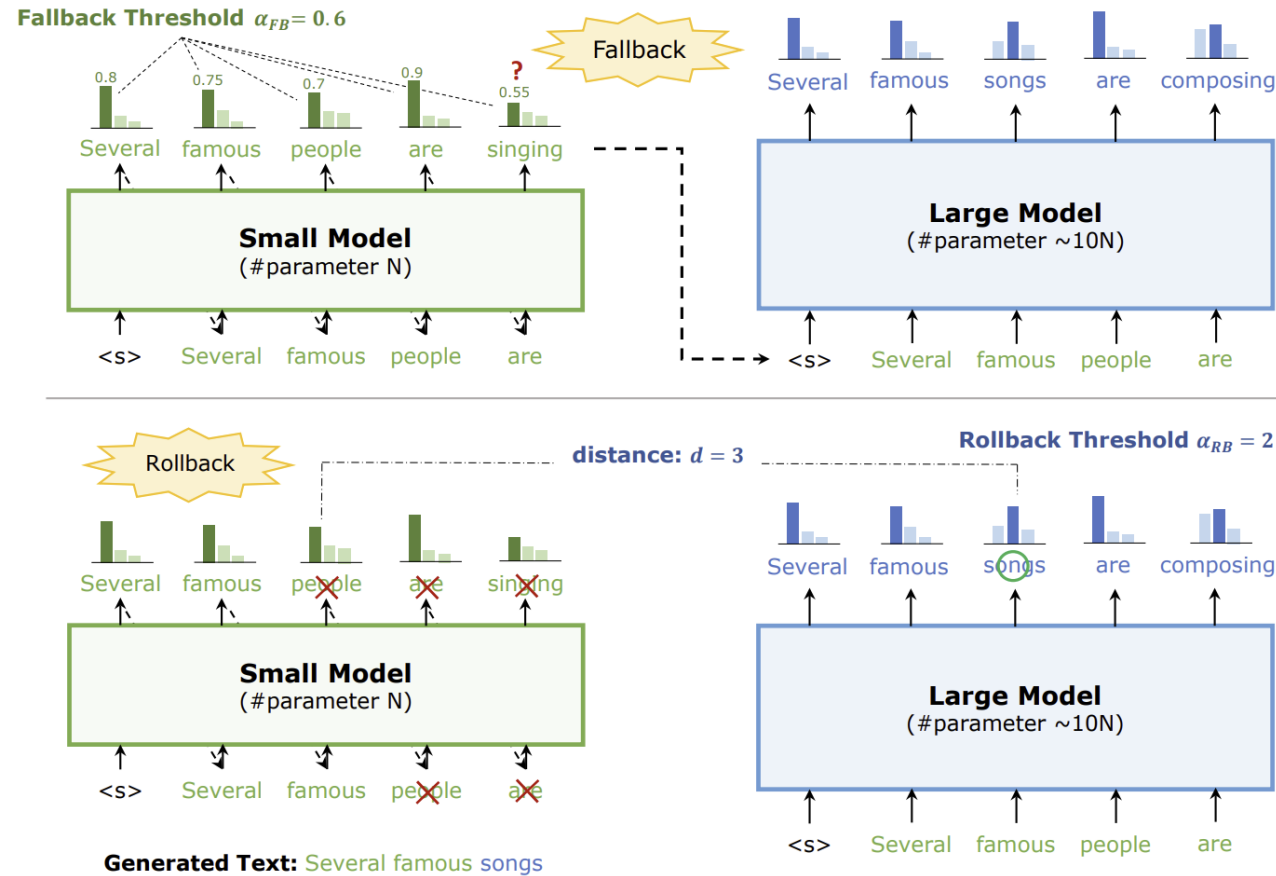
Pruning



Key links:

1: [Pruning Tutorial Torch](#)

Speculative Decoding



Key links:

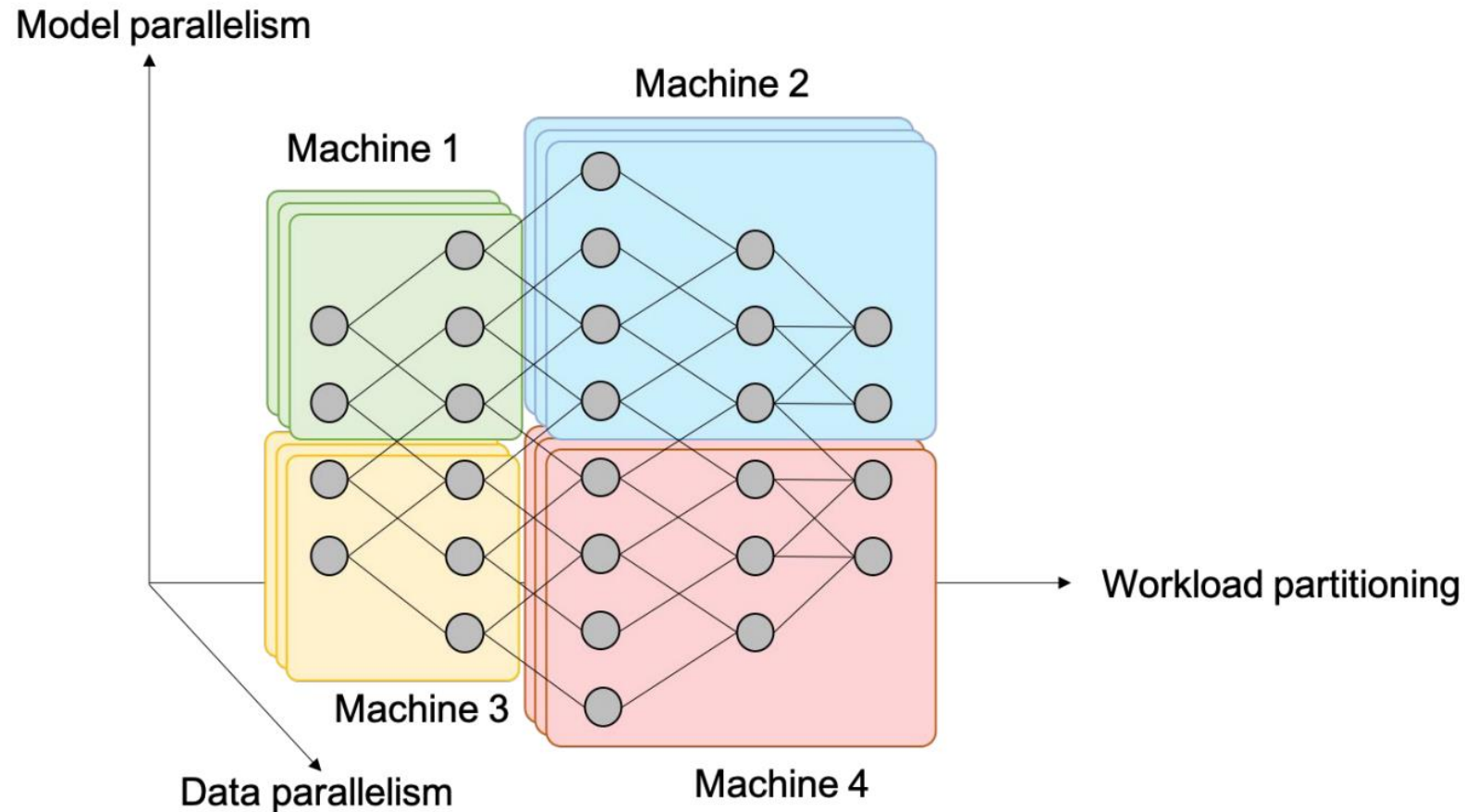
- 1: [Speculative Decoding paper](#)
- 2: [A Hitchhiker's Guide to Speculative Decoding](#)



Table of Content

- Just finished training, now what?
- Engineering tricks
- Architectural tricks
- **Going distributed**

Distributed is painful on 3 main dimensions



Key links:

- 1: [Torch distributed tutorial](#)
- 2: [Torch Distributed Data Parallel](#)
- 3: [Torch on Model Parallelism](#)
- 4: [On Parameter Servers](#)
- 5: [Accelerate by HF](#)
- 6: [Horovod](#)
- 7: [NCCL by Nvidia](#)
- 8: [DeepSpeed by Microsoft](#)
- 9: [Ray](#)

