# Fantastic Language Models and How to Build Them

**Guest Lecture — CS 224U: Natural Language Understanding**
*Stanford || Zoom || Folks 2x-ing the Recording*
April 12, 2023

Stanford|NLP

iliad
intelligent and interactive autonomous systems

**Siddharth Karamcheti**

$$i\hbar\frac{dc_1}{dt} = E_0 C_1 - A C_2$$

$$i\hbar\frac{dc_2}{dt} = E_0 C_2 - A C_1$$

$-A = H_{12} = H_{21}$

$\langle \cdots$

$(\cdots) = (E-A)(c_1 + c_2)$

$\sim e^{-\frac{i}{\hbar}(E-A)\,t}$

$E-A \qquad C_1 = C_2$

$E+A \qquad C_2 = -C_1$

$|\psi\rangle = \sum_i |i\rangle\langle i|\psi\rangle$

$|\varphi\rangle$ at $t_1$

delay until $t_2$

$\langle \chi | U(t_2,t_1)|\varphi\rangle = \sum_i \langle \cdots$

$|\psi(t)\rangle = \sum_i |i\rangle C_i(t)$

# On the Importance of "Building"

**Today** — a *practical* take on large-scale language models (LLMs).

Whirlwind tour of the full pipeline:

- **Model Architecture** — Evolution of the Transformer

- **Training at Scale** — From 124M to 1T+ Parameters

- **Efficient Finetuning & Inference** — Tips & Tricks

**Punchline**: From "folk knowledge" —> insight / intuition / (re-)discovery!

*Please ask lots of questions! Why is this information useful to <YOU>?*

# Part I: Evolution of the Transformer

"Experiment is the mother of knowledge."
— Madeline L'Engle, *A Wrinkle in Time*

# Recipe for a Good™ Language Model

**Massive amounts of cheap, easy to acquire data...**

X

**... a simple, high-throughput way to *consume* it!**

Natural to scale with data.
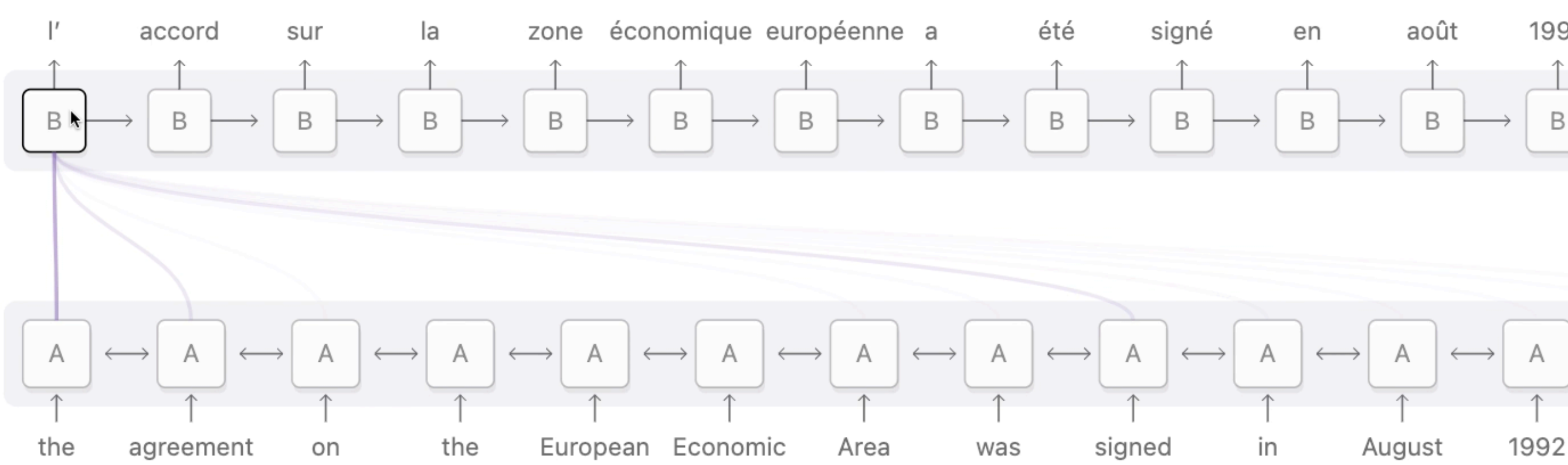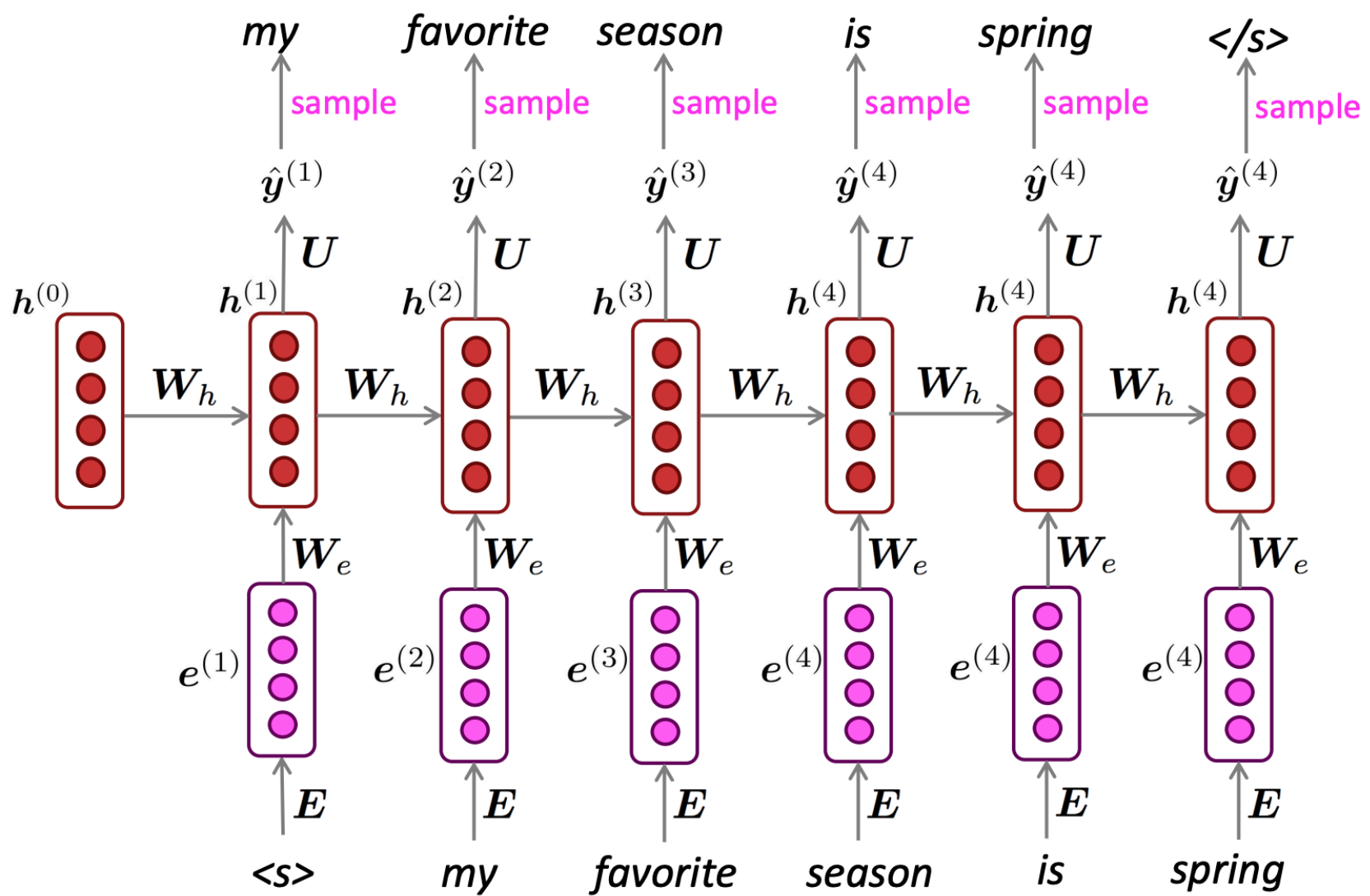
Composable and "general".

Fast & parallelizable training.

High hardware utilization.

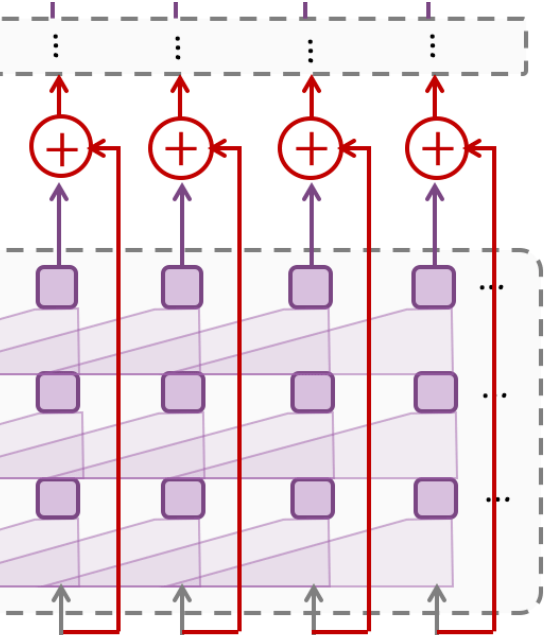Minimal "assumptions" on

relationships between data?

**<Story Time>**

# Pre-2017 — Historical Context

**RNNs**



**CNNs**



receptive field

Residual Connection

**RNN Key Ideas:** Long Context, **Attention**

**CNN Key Ideas:**
- Layer: Multiple "Filters" (Views)
- **Scaling Depth** w/ Residuals
- **Parallelizable!**

**< How do I do better? >**

For the RNN top text: my, favorite, season, is, spring, </s>

Reference: "Attention and Augmented Recurrent Neural Networks," Chris Olah and Shan Carter. *Distill, 2016.*
Reference: "Convolutional Neural Networks for Text," Lena Voita. ML for NLP @ YSDA

6

# Formulating the Self-Attention Block



```python
class Attention(nn.Module):
    def __init__(self, embed_dim: int, n_heads: int):
        super().__init__()
        self.n_heads, self.dk = n_heads, (embed_dim // n_heads)
        self.qkv = nn.Linear(embed_dim, 3 * embed_dim)
        self.proj = nn.Linear(embed_dim, embed_dim)

    def forward(self, x: Tensor[bsz, seq, embed_dim]):
        q, k, v = rearrange(
            self.qkv(x),
            "bsz seq (qkv nh dk) -> qkv bsz nh seq dk",
            qkv=3,
            nh=self.n_heads,    # Different "views" (like CNN filters)!
            dk=self.dk,
        ).unbind(0)
```
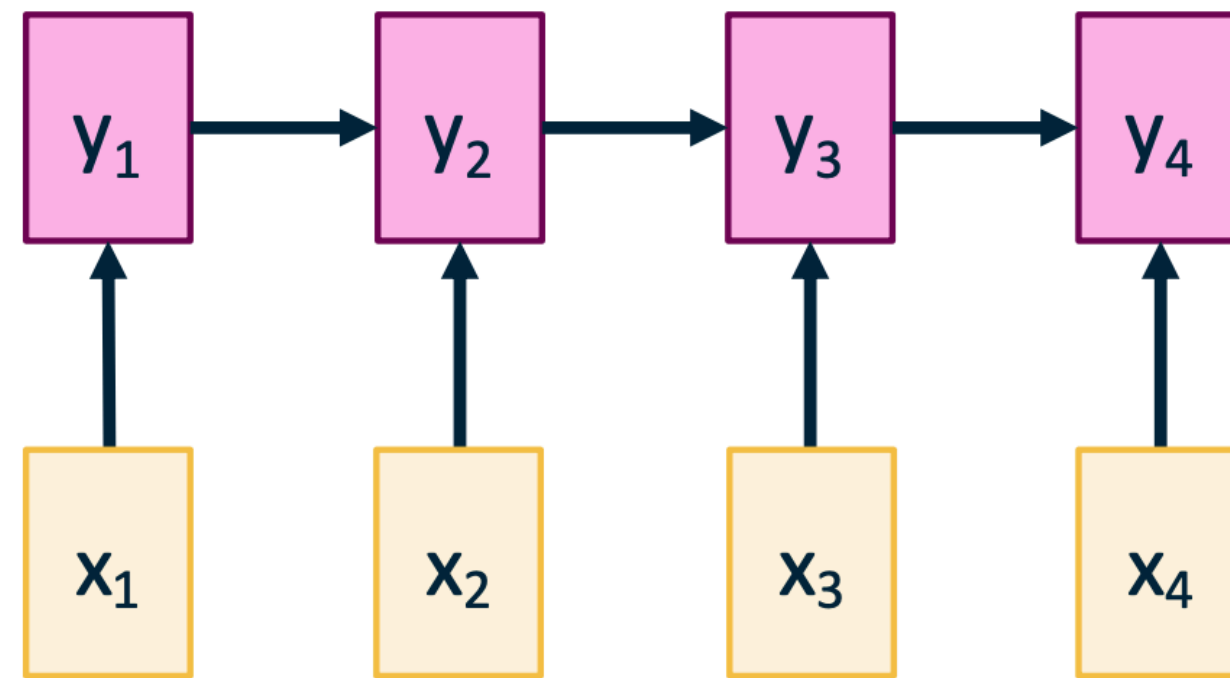
**Self-Attention:** "The" —> **query**, **key,** & **value**

**Multi-Headed:** Different "views" per layer

**< Is this actually better? >**

# Aside — Self-Attention & Parallelization

## Recurrent Neural Network



## 1D Convolution



## Self-Attention
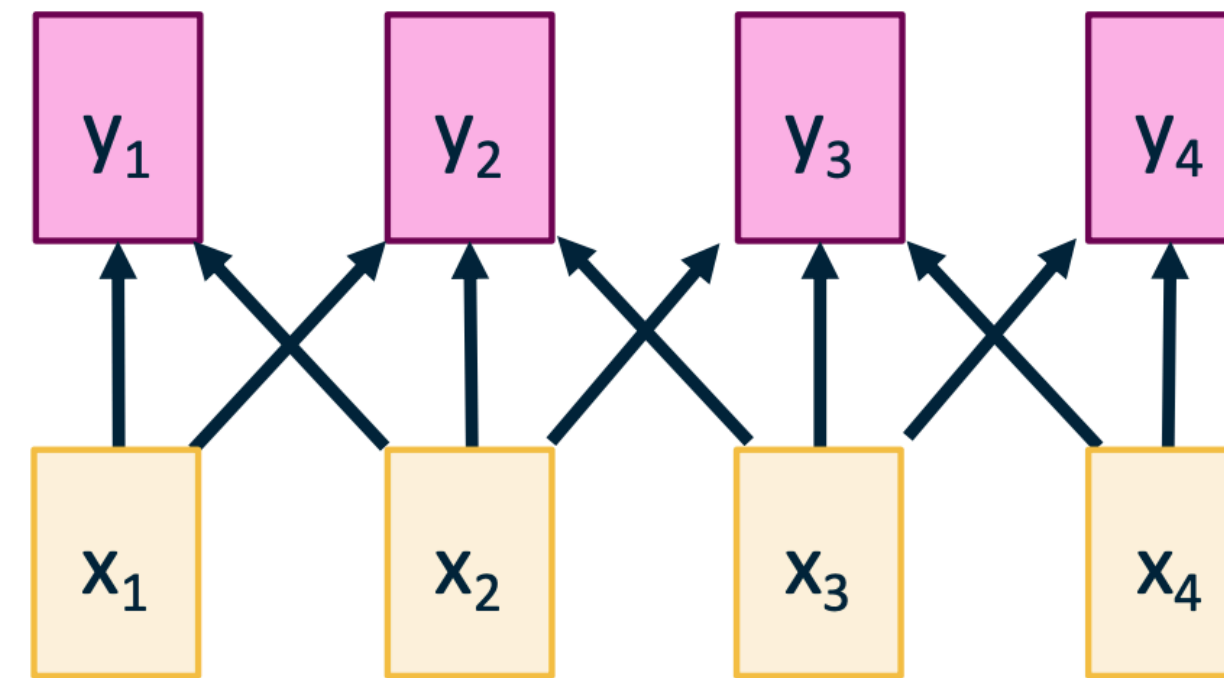


Works on **Ordered Sequences**
(+) Good at long sequences: After one RNN layer, $h_T$ ”sees” the whole sequence
**(-) Not parallelizable: need to compute hidden states sequentially**

Works on **Multidimensional Grids**
**(-) Bad at long sequences: Need to stack many conv layers for outputs to "see" the whole sequence**
(+) Highly parallel: Each output can be computed in parallel

Works on **Sets of Vectors**
(+) Good at long sequences: after one self-attention layer, each output "sees" all inputs!
(+) Highly parallel: Each output can be computed in parallel
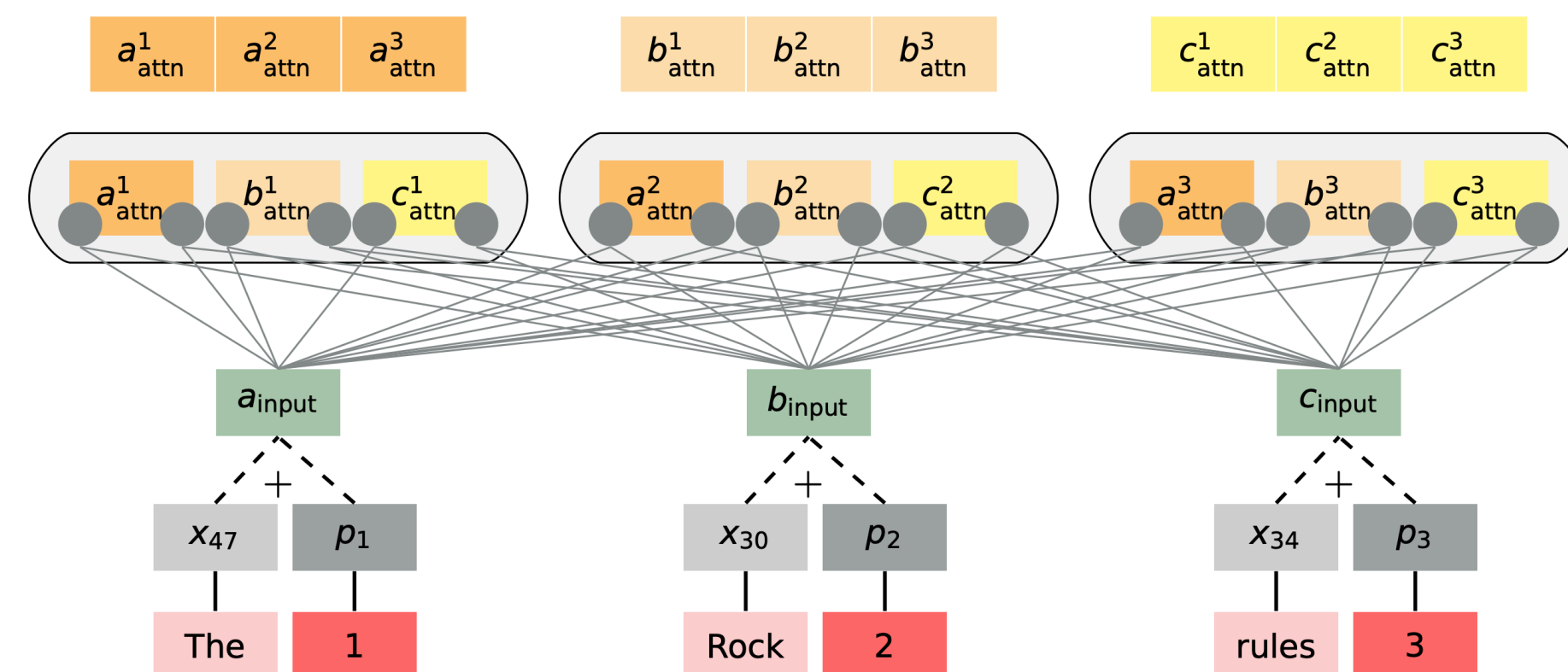**(-) Very memory intensive**

## < Great! But… what am I missing? >

# Formulating the Self-Attention Block

```python
class Attention(nn.Module):
    def __init__(self, embed_dim: int, n_heads: int):
        super().__init__()
        self.n_heads, self.dk = n_heads, (embed_dim // n_heads)
        self.qkv = nn.Linear(embed_dim, 3 * embed_dim)
        self.proj = nn.Linear(embed_dim, embed_dim)

    def forward(self, x: Tensor[bsz, seq, embed_dim]):
        q, k, v = rearrange(
            self.qkv(x),
            "bsz seq (qkv nh dk) -> qkv bsz nh seq dk",
            qkv=3,
            nh=self.n_heads,    # Different "views" (like CNN filters)!
            dk=self.dk,
        ).unbind(0)

        # RNN Attention --> *for each view*
        scores = torch.softmax(
            q @ (k.transpose(-2, -1)),
            dim=-1
        )
        return self.proj(
            rearrange(scores @ v, "b nh seq dk -> b seq (nh dk)")
        )
```



**< Where's my nonlinearity? >**

9

# Expressivity & Nonlinearity

```python
class ExpressiveTransformerBlock(nn.Module):
    def __init__(self, embed_dim: int, n_heads: int, up: int = 4):
        super().__init__()
        self.attn = Attention(embed_dim, n_heads)

        # Project *up* to high-dimension, nonlinear, compress!
        self.mlp = nn.Sequential(
          nn.Linear(embed_dim, up * embed_dim),
          nn.ReLU(),
          nn.Linear(up * embed_dim, embed_dim)
        )

    def forward(self, x: T[bsz, seq, embed_dim]):
        x = x + self.attn(x)
        x = x + self.mlp(x)
        return x
```

**CS 229** —> SVMs & "Implicit Lifting"

Decision surface

**Residual + MLP** —> "Sharpen" + "Forget"

**< New Problem — Activations Blow Up! >**

# Going Deeper —> Activation Instability

```python
class NormalizedTransformerBlock(nn.Module):
    def __init__(self, embed_dim: int, n_heads: int, up: int = 4):
        super().__init__()
        self.attn = Attention(embed_dim, n_heads)
        self.mlp = nn.Sequential(
            nn.Linear(embed_dim, up * embed_dim),
            nn.ReLU(),
            nn.Linear(up * embed_dim, embed_dim)
        )

        # Add Normalization Layers
        self.attn_norm = nn.LayerNorm(embed_dim)
        self.mlp_norm = nn.LayerNorm(embed_dim)

    def forward(self, x: T[bsz, seq, embed_dim]):
        x = self.attn_norm(x + self.attn(x))
        x = self.mlp_norm(x + self.mlp(x))
        return x
```
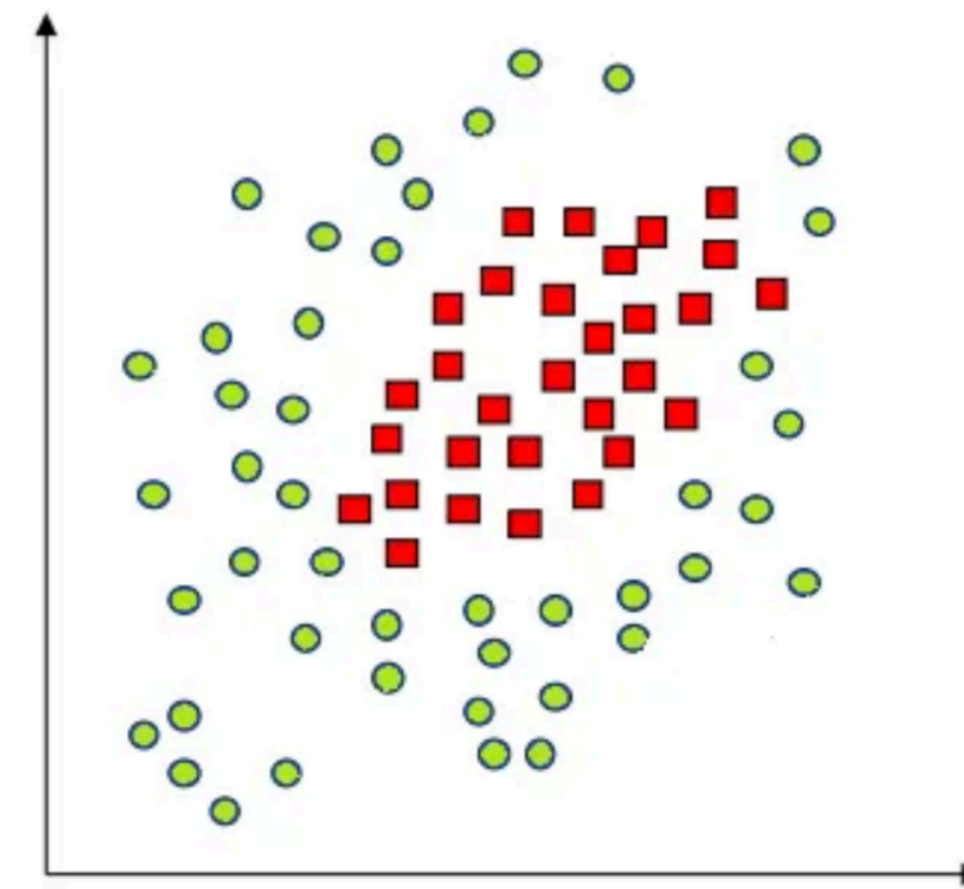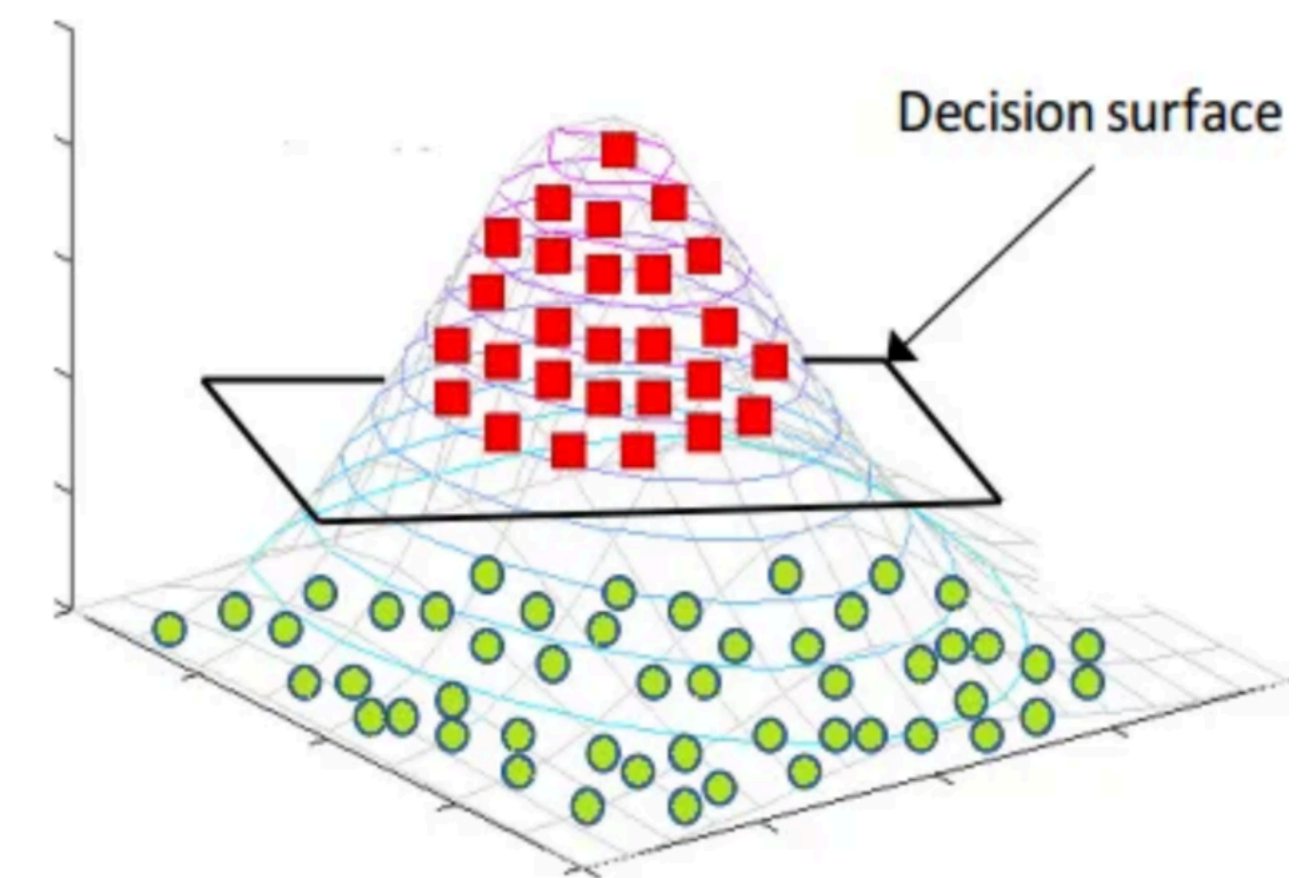


**Layer Normalization**

**< And... we're done? >**

# Well, Shucks —> Emergent Optimization Problems

**Typical LR Decay**

[CS 221, CS 229]



Normalized
Transformer

**Transformer Pretraining LR Schedule**

Linear Warmup (5% of Training) then Decay...



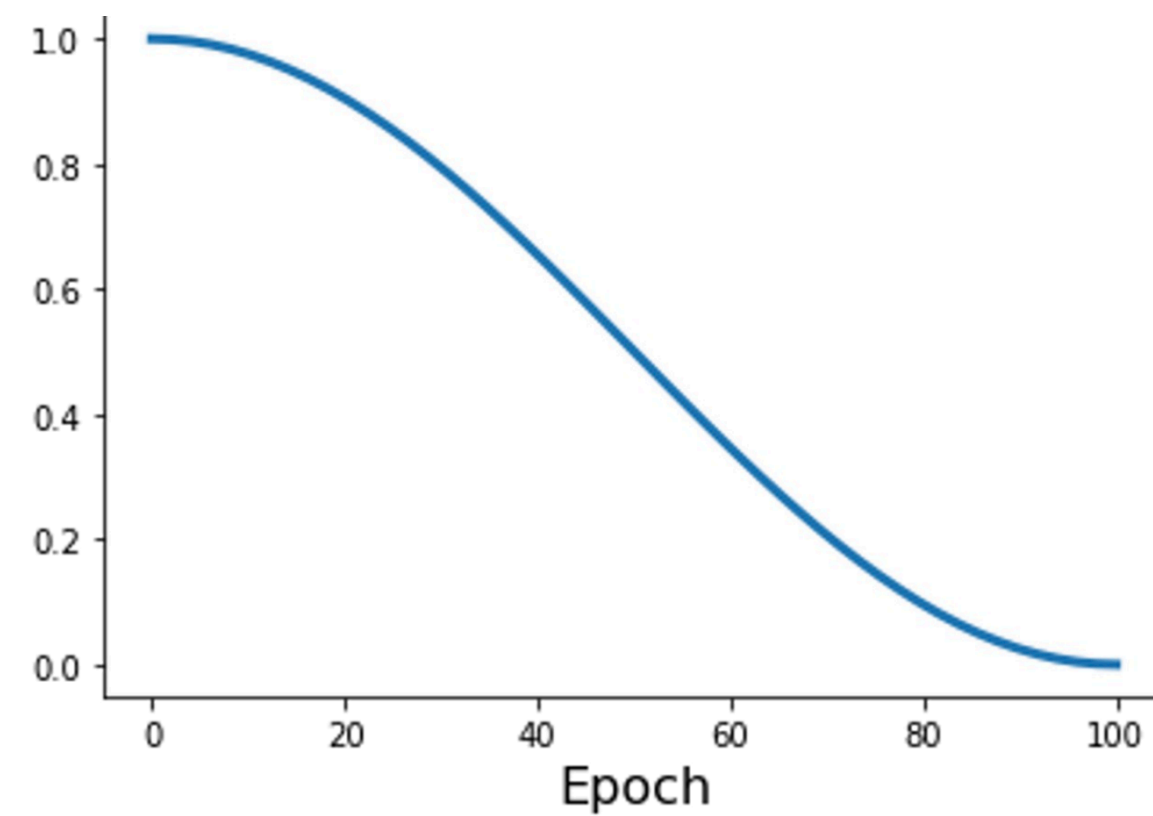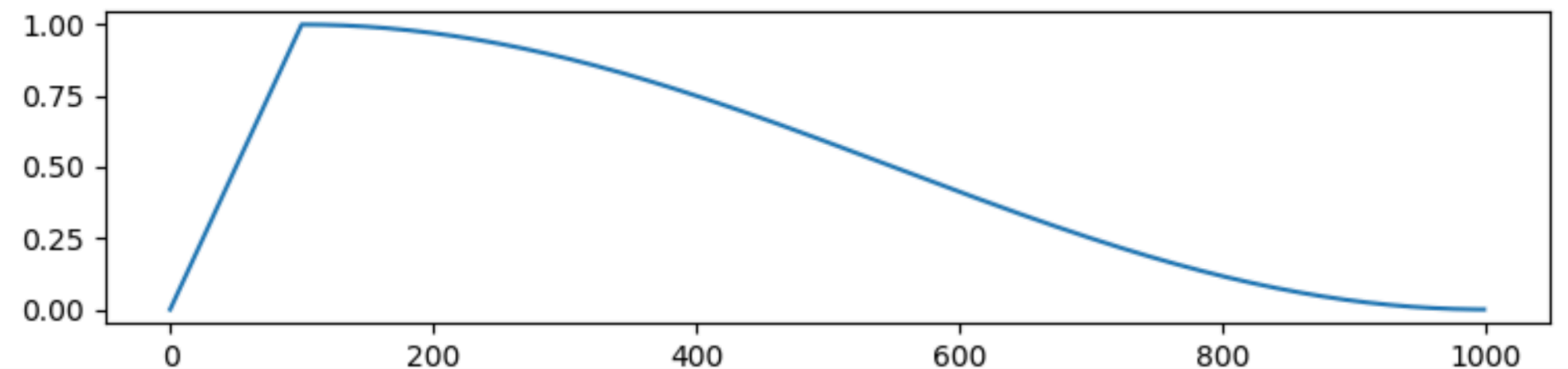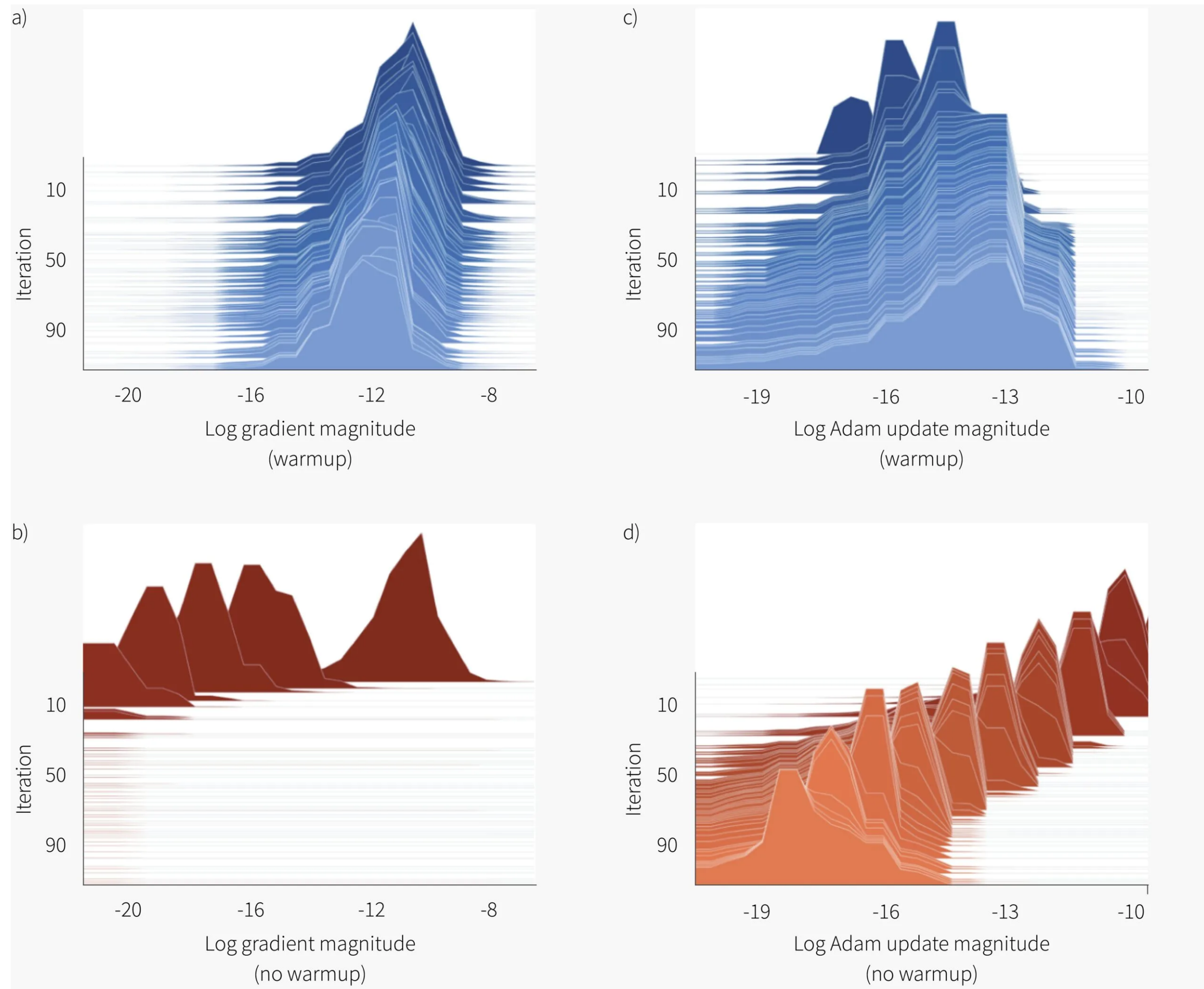**Learning Rate Warmup** —> Breaks conventional machine learning wisdom?

**< Ok but... why? >**

**Reference:** "The Annotated Transformer," Sasha Rush. *Harvard NLP (2018)*

12

a)

Iteration

Log gradient magnitude
(warmup)

c)

Iteration

Log Adam update magnitude
(warmup)

b)

Iteration

Log gradient magnitude
(no warmup)

d)

Iteration

Log Adam update magnitude
(no warmup)

## 3.1. Problem in Transformer Optimization

In this section we demonstrate that the requirement for warmup comes from a combined effect of high variance in the Adam optimizer and backpropagation through layer normalization. Liu et al. (2020) showed that at the begin-

of the input. Specifically, the gradient has the following property:

$$\left\| \frac{\partial \text{LN}(\boldsymbol{x})}{\partial \boldsymbol{x}} \right\| = O\left( \frac{\sqrt{d}}{\|\boldsymbol{x}\|} \right) \tag{1}$$

where $\boldsymbol{x}$ is the input to layer normalization and $d$ is the embedding dimension. If input norm $\|\boldsymbol{x}\|$ is larger than $\sqrt{d}$ then backpropagation through layer normalization has a down scaling effect that reduces gradient magnitude for lower layers. Compounding across multiple layers this can quickly lead to gradient vanishing.

**< Ok, now we're done...? >**

**Reference:** "Transformers III Training; Tricks for Training Transformers," Borealis AI — 8/6/2021.
**Reference:** "Improving Transformer Optimization through Better Initialization," Huang et. al., *ICML 2020.*

# The Modern Transformer (March 2023)

```python
class ModernTransformerBlock(nn.Module):
    def __init__(self, embed_dim: int, n_heads: int, up: int = 4):
        super().__init__()
        self.attn = Attention(embed_dim, n_heads, qk_bias=False)
        self.mlp = nn.Sequential(
            SwishGLU(embed_dim, up * embed_dim),
            nn.Linear(up * embed_dim, embed_dim)
        )

        # Post-Norm --> *Pre-Norm*
        self.pre_attn_norm = RMSNorm(embed_dim)
        self.pre_mlp_norm = RMSNorm(embed_dim)

    def forward(self, x: T[bsz, seq, embed_dim]):
        x = x + self.attn(self.pre_attn_norm(x))
        x = x + self.mlp(self.pre_mlp_norm(x))
        return x
```

```python
# SwishGLU -- A Gated Linear Unit (GLU) with Swish Activation
class SwishGLU(nn.Module):
    def __init__(self, in_dim: int, out_dim: int):
        super().__init__()
        self.swish = nn.SiLU()
        self.project = nn.Linear(in_dim, 2 * out_dim)

    def forward(self, x: T[bsz, seq, embed_dim]):
        projected, gate = self.project(x).tensor_split(2, dim=-1)
        return projected * self.swish(gate)


# RMSNorm -- Simple Alternative to LayerNorm
class RMSNorm(nn.Module):
    def __init__(self, dim: int, eps: float = 1e-8):
        super().__init__()
        self.scale, self.eps = dim**-0.5, eps
        self.g = nn.Parameter(torch.ones(dim))

    def forward(self, x: T[bsz, seq, embed_dim]):
        norm = torch.norm(x, dim=-1, keepdim=True) * self.scale
        return x / norm.clamp(min=self.eps) * self.g
```
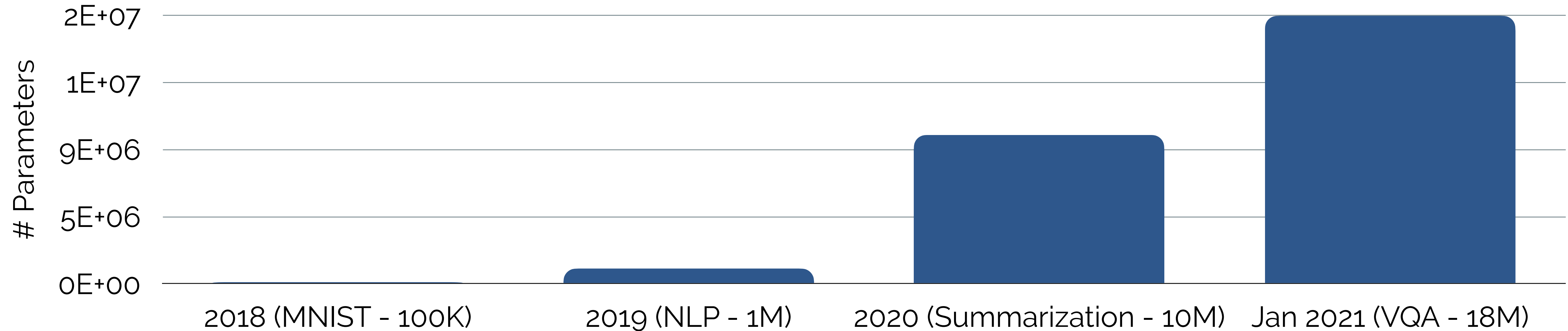
< Fin />

# Part II: Training at Scale

"Nothing in life is to be feared. It is only to be understood."
— Marie Curie

# Short Story — My Deep Learning Trajectory



- "Standard Pipeline": Train on 1 GPU (e.g., on Colab) —> ~**max of a few hours.**

- Let's train a GPT-2 Small (124M)!
  - **Problem:** Batch > 4 goes OOM on a decent GPU = > **12 GB** of GPU RAM
  - Simple Trick —> *Gradient Accumulation!*
    - But... **99.63 Days** to train on Single GPU (400K Steps)

GPT-2 Training Clock

99.63 D

# Shortening the Clock —> The Scaling Toolbox

## GPT-2 Training Clock

99.63 D

**Goal:** 100 Days on 1 GPU —> ~**4 Days on 16 GPUs**

- **Data Parallelism** — Scaling *across* GPUs & Nodes

- **Mixed Precision** — Bits, Bytes, and TensorCores

- **ZeRO Redundancy** — Minimizing Memory Footprint

**Later...** Model Parallelism — Hardware Limitations — Software Optimization

*Even if you're not training big models...* ***understanding breeds innovation!***

# Data Parallelism — A Toy Example

## GPT-2 Training Clock

99.63 D

```python
BATCH_SIZE = 128

class MLP(nn.Module):
    def __init__(
      self, n_classes: int = 10, mnist_dim: int = 784, hidden: int = 128
    ):
        super().__init__()
        self.mlp = nn.Sequential(
            nn.Linear(mnist_dim, hidden),
            nn.ReLU(),
            nn.Linear(hidden, hidden),
            nn.ReLU(),
            nn.Linear(hidden, n_classes)
        )

    def forward(self, x: T[bsz, mnist_dim]):
        return self.mlp(x)

# Main Code
dataloader = DataLoader(dataset=torchvision.datasets(...), batch_size=BATCH_SIZE)
model = MLP()

# Train Loop
criterion, opt = nn.CrossEntropyLoss(), optim.AdamW(model.parameters())
for (inputs, labels) in dataloader:
    loss = criterion(model(inputs), labels)
    loss.backward(); opt.step(); opt.zero_grad()
```

**Idea —>** Parallelize?

**SIMD**

Single Instruction, Multiple Data

**SPMD**

Single **Program**, Multiple Data

**< Seems hard? >**

# (Distributed) Data Parallelism — Implementation

## GPT-2 Training Clock

99.63 D

7.2 D — 16 GPUs w/ Data Parallelism (DDP)

```python
from torch.nn.parallel import DistributedDataParallel as DDP
from torch.utils.data.distributed import DistributedSampler

BATCH_SIZE, WORLD_SIZE = 128, 8  # World Size == # of GPUs
class MLP(nn.Module):
    def __init__(
      self, n_classes: int = 10, mnist_dim: int = 784, hidden: int = 128
    ):
        super().__init__()
        self.mlp = nn.Sequential(
            nn.Linear(mnist_dim, hidden),
            nn.ReLU(),
            nn.Linear(hidden, hidden),
            nn.ReLU(),
            nn.Linear(hidden, n_classes)
        )

    def forward(self, x: T[bsz, mnist_dim]):
        return self.mlp(x)

# Main Code
train_set = torchvision.dataset(...)
dist_sampler = DistributedSampler(dataset=train_set)
dataloader = DataLoader(
  train_set, sampler=dist_sampler, batch_size=BATCH_SIZE // WORLD_SIZE
)

model = DDP(
  MLP(),
  device_ids=[os.environ["LOCAL_RANK"]],
  output_device=os.environ["LOCAL_RANK"]
)
```

Auto-Partitions Data across Processes

Simple Wrapper around nn.Module()

```python
# Train Loop
criterion, opt = nn.CrossEntropyLoss(), optim.AdamW(model.parameters())
for (inputs, labels) in dataloader:
    loss = criterion(model(inputs), labels)
    loss.backward(); opt.step(); opt.zero_grad()

# Run: `torchrun --nnodes 1 --nproc_per_node=8 main.py`
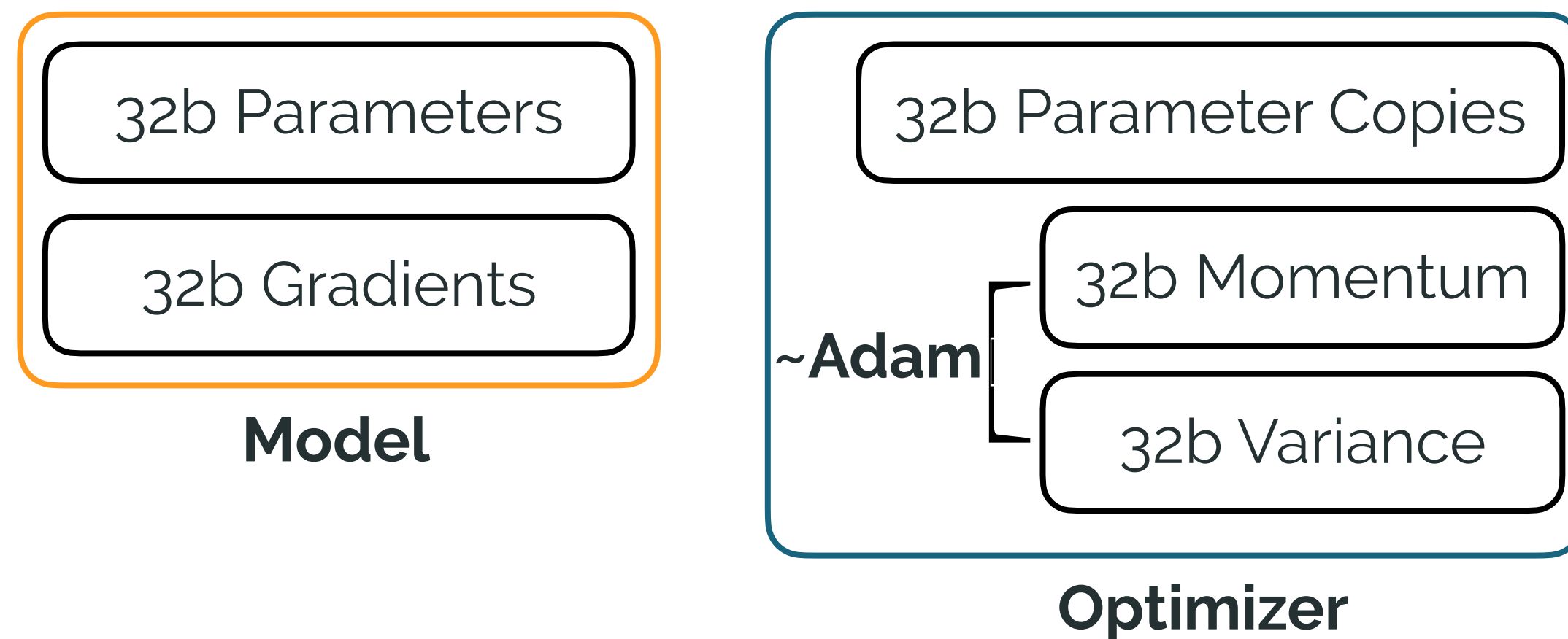```

Nifty Utility —> Spawns Processes

# **Important** — Memory Footprint of Training?

99.63 D

7.2 D — 16 GPUs w/ Data Parallelism (DDP)

## **Standard (Float 32) Memory Footprint**

[Excludes Activations + Temporary Buffers]

| 32b Parameters |
|---|
| 32b Gradients |

**Model**

| 32b Parameter Copies |
|---|

~Adam
| 32b Momentum |
|---|
| 32b Variance |

**Optimizer**

Lower Bound on "Static" Memory (w/ Adam):

= **# Parameters * 20 Bytes**

Activation Memory >> Static Memory

## **Training Implications**

- 1B Parameters —> 18 GB (~**31 GB w/ BSZ = 1**)

- 175B Parameters —> *3 TB (w/o activations!)*

## **Facts about Floating Points**

- Float32 — Standard defined in IEEE-754
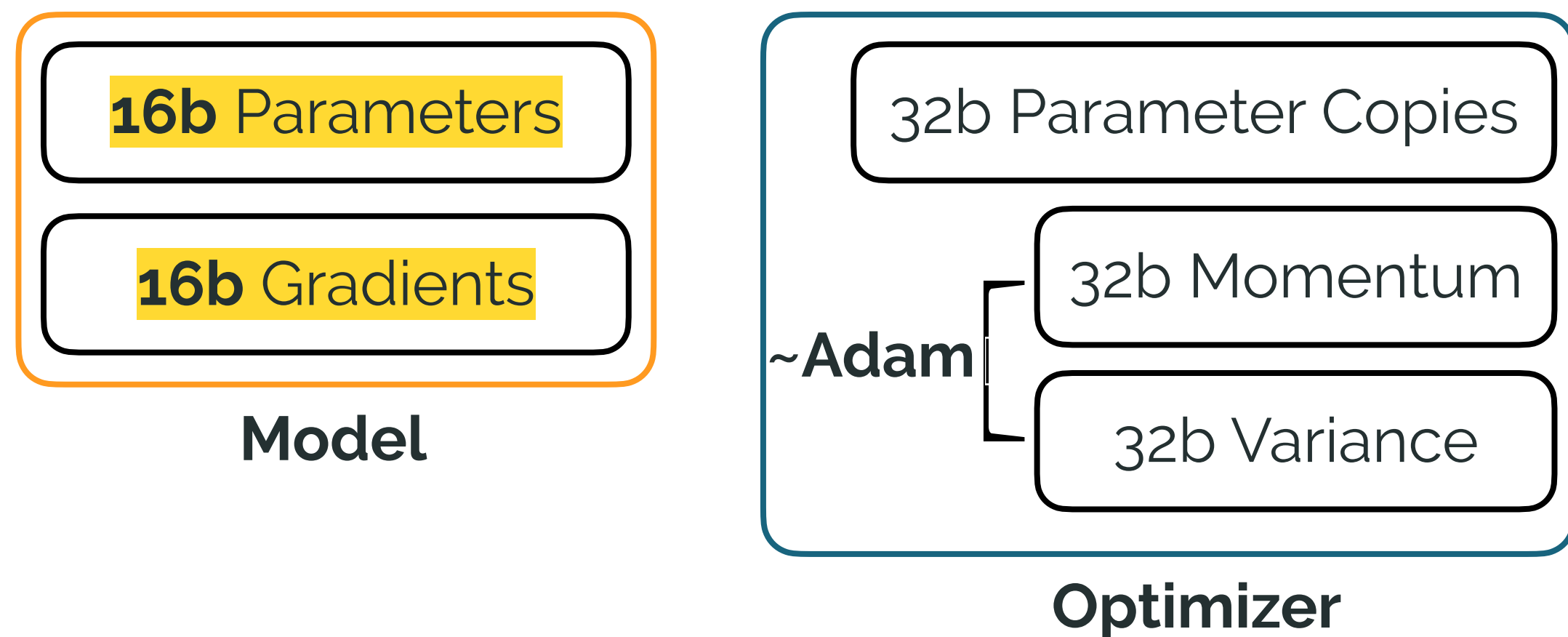  - Sign (1) — Exponent (8) — Significand (23)
  - Wide Range —> up to 1e38

## **< Do we need \*all\* 32 bits? >**

**Reference**: "ZeRO: Memory Optimizations Toward Training Trillion Parameter Models," Rajbhandari, Rasley, Ruwase, and He. *SC 2020.*

# Mixed Precision Training

## GPT-2 Training Clock

99.63 D

7.2 D — 16 GPUs w/ Data Parallelism (DDP)

6.01 D — 16 GPUs w/ DDP, FP16

## Mixed Precision (FP16) Memory Footprint
[Excludes Activations + Temporary Buffers]

**16b** Parameters

**16b** Gradients

**Model**

32b Parameter Copies

~**Adam**
32b Momentum

32b Variance

**Optimizer**

Lower Bound on "Static" Memory (w/ Adam):
= **# Parameters** * **16 Bytes**

Activation Memory —> *halved!*

## Hmm... Optimizer Memory?

FP16 **does not mean *everything*** is FP16.

**Real Gain:** NVIDIA Tensor Core Speedup!



TENSOR CORES

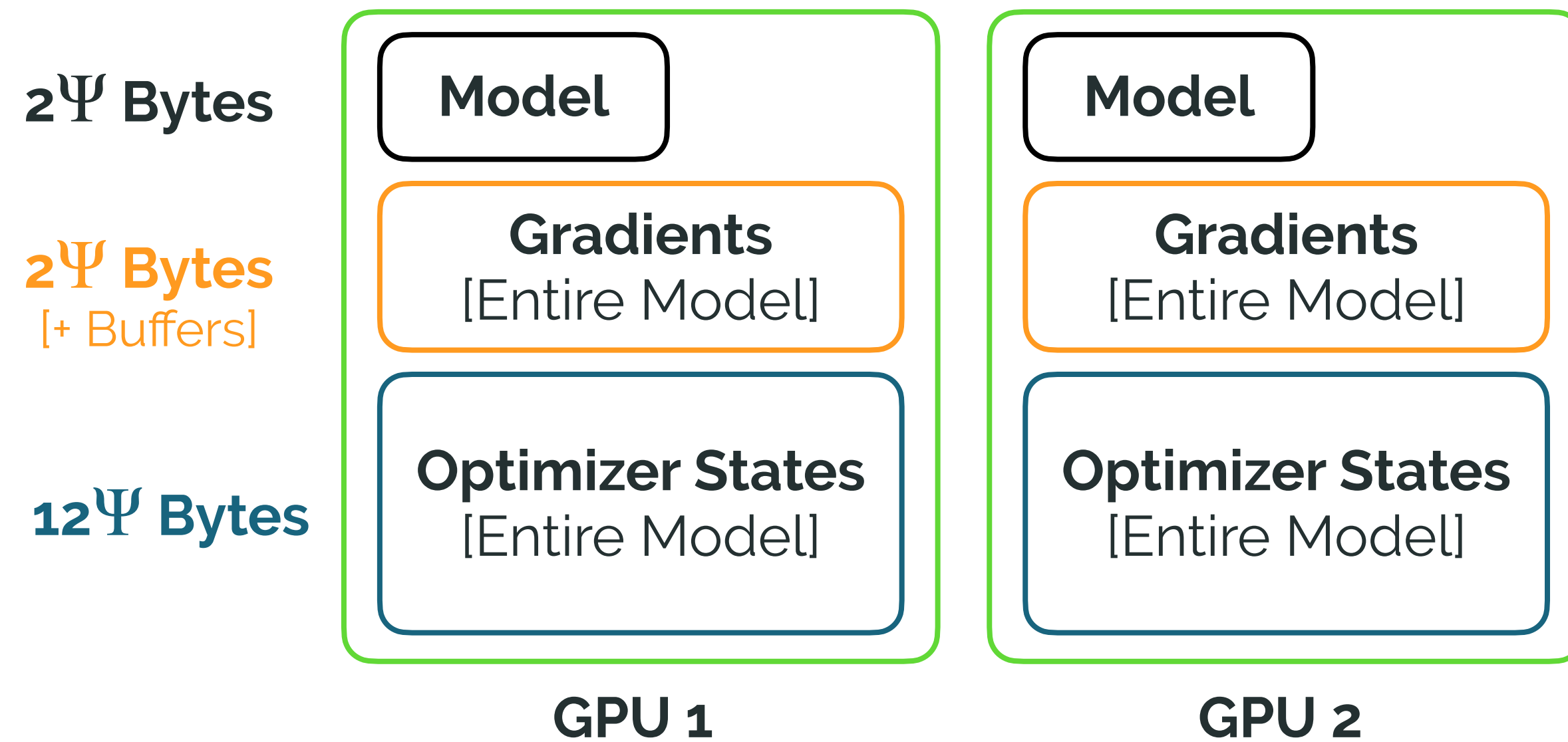**Reference**: "ZeRO: Memory Optimizations Toward Training Trillion Parameter Models," Rajbhandari, Rasley, Ruwase, and He. *SC 2020.*

# Eliminate Redundancies —> ZeRO

## GPT-2 Training Clock

99.63 D

7.2 D — 16 GPUs w/ Data Parallelism (DDP)

6.01 D — 16 GPUs w/ DDP, FP16

3.37 D — 16 GPUs w/ DDP, FP16, ZeRO

**Punchline:** "Shards" Memory by # of GPUs!

### Standard Data Parallelism
"Replicate everything but the data!"

### ZeRO Data Parallelism
"Replicate only what you need"

| | Standard | | ZeRO | | |
|---|---|---|---|---|---|
| $2\Psi$ Bytes | **Model** | **Model** | **Model** | **Model** | $2\Psi$ Bytes |
| $2\Psi$ Bytes [+ Buffers] | **Gradients** [Entire Model] | **Gradients** [Entire Model] | **Gradients** *[Layers 1-6]* | **Gradients** *[Layers 7-12]* | $(2\Psi / W)$ Bytes [+ Buffers] |
| $12\Psi$ Bytes | **Optimizer States** [Entire Model] | **Optimizer States** [Entire Model] | **Opt. States** *[Layers 1-6]* | **Opt. States** *[Layers 7-12]* | $(12\Psi / W)$ Bytes |
| | **GPU 1** | **GPU 2** | **GPU 1** | **GPU 2** | **W = # of GPUs** |

$\Psi$ = # of Parameters

**Reference**: "ZeRO: Memory Optimizations Toward Training Trillion Parameter Models," Rajbhandari, Rasley, Ruwase, and He. *SC 2020.*
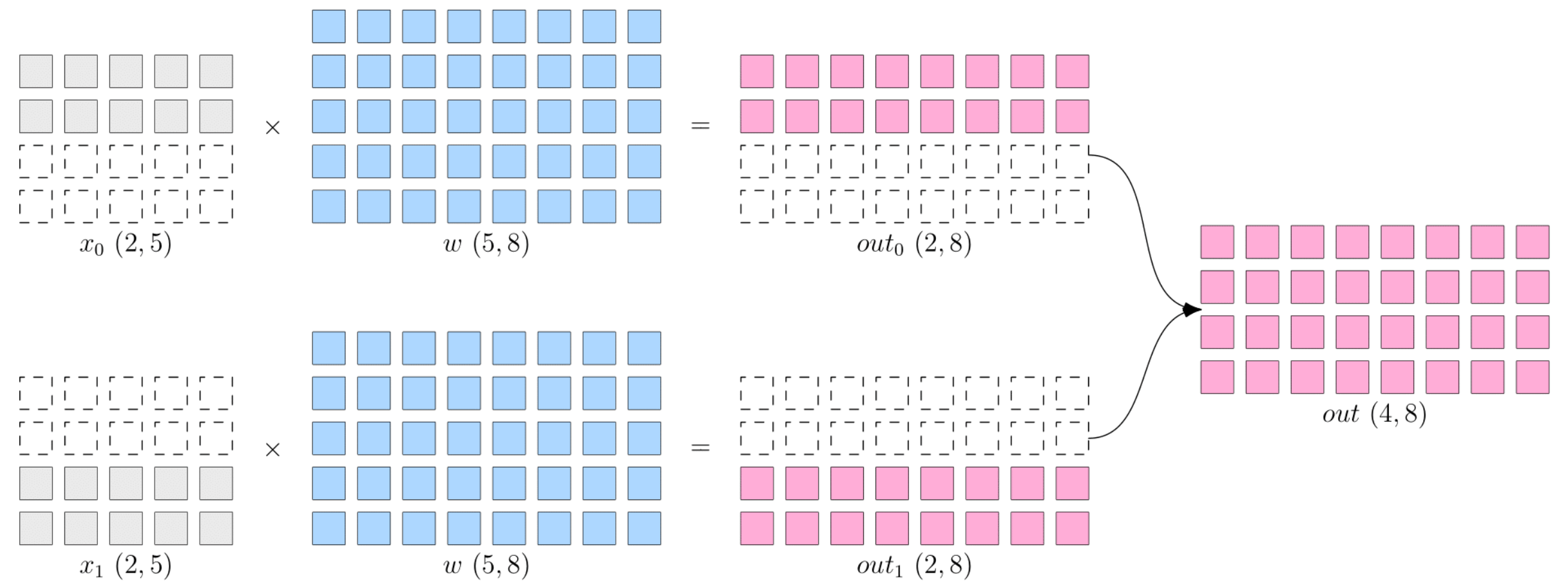
# Alas — Hitting a (Communication) Wall

**Problem** — At some point, communication cost between nodes is too much!

**Answers:**

Exploit Matrix Multiplication...



Schedule Backwards Pass Wisely...



**< Harder to implement, model-specific... still miles to go! >**

# Part III: Fine-Tuning and Inference

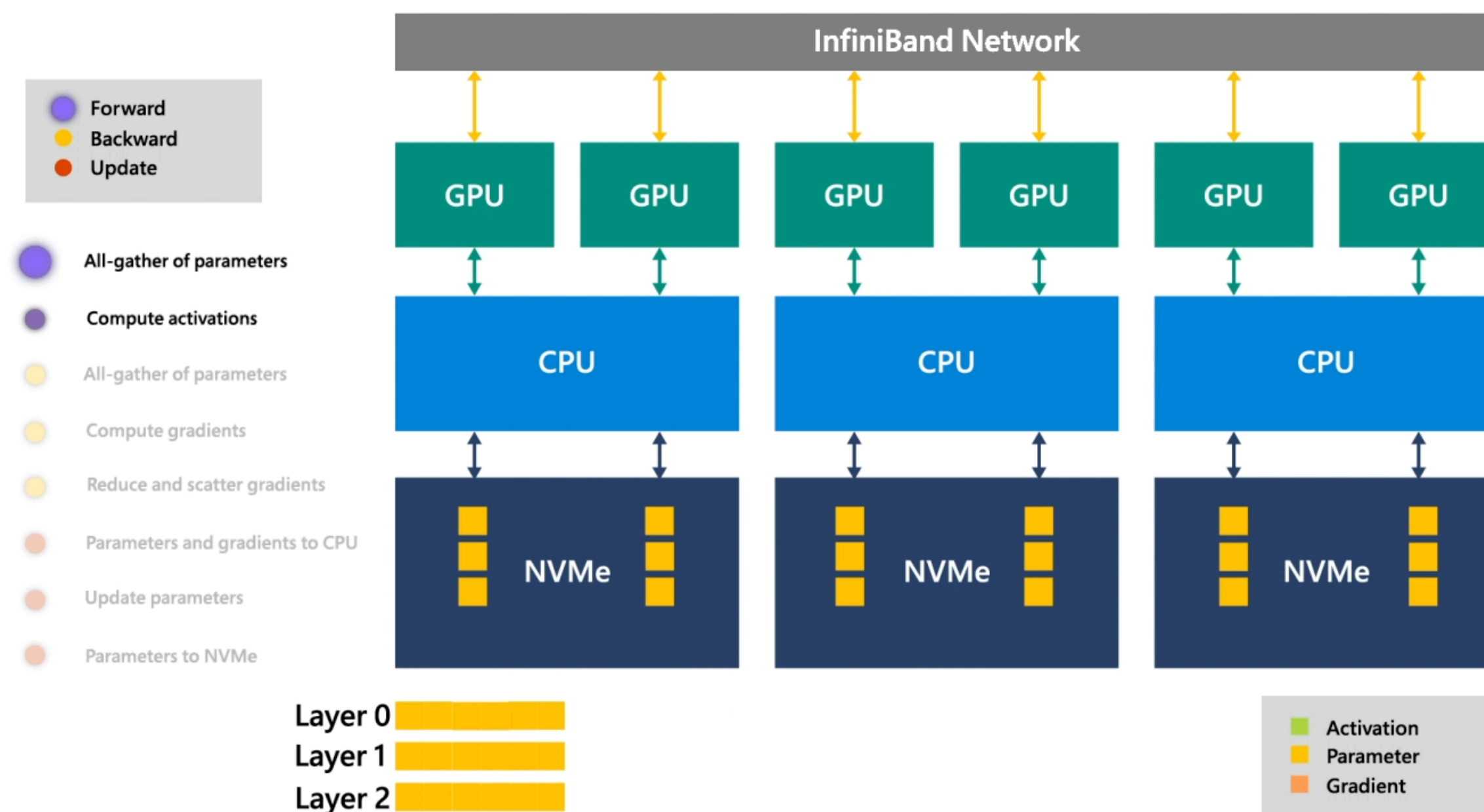"It's such a happiness, when good people get together."
— Jane Austen, *Emma*

# Tools for Training —> Tools for Fine-Tuning

**Silver Lining** — Learning to scale training —> informs *fine-tuning & inference!*
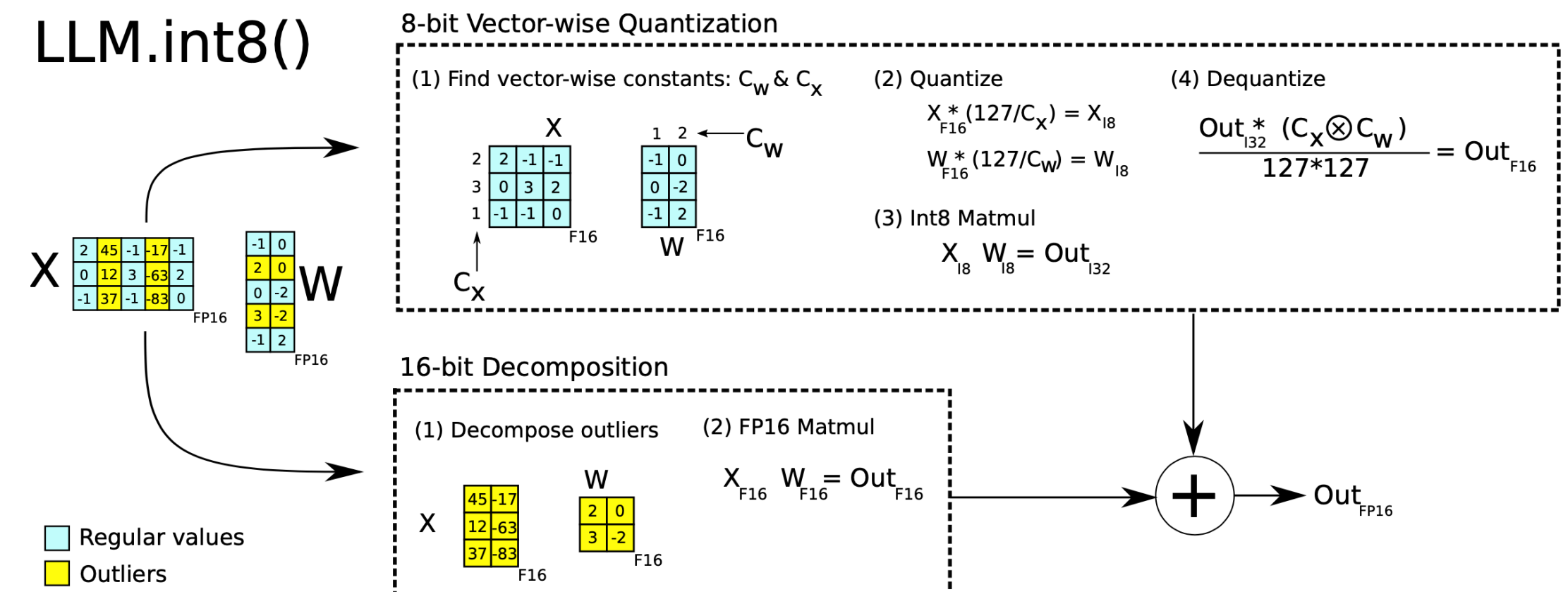
### ZeRO Data Parallelism

### Mixed Precision (FP16)

### ZeRO Infinity —> CPU/NVMe Offloading

### 8-Bit Quantization



*Powers `llama.cpp` and more!*

**Reference**: "LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale," Dettmers, Lewis, Belkada, and Zettlemoyer. *NeurIPS 2022.*

# Teaser for Later —> Parameter-Efficient Fine-Tuning



**LoRA (Low-Rank Adaptation)**

**adaLN (Adapted LayerNorm)**

**...and more!**

# That's all Folks!

"This wind, it is not an ending…"
— Robert Jordan, *A Memory of Light*