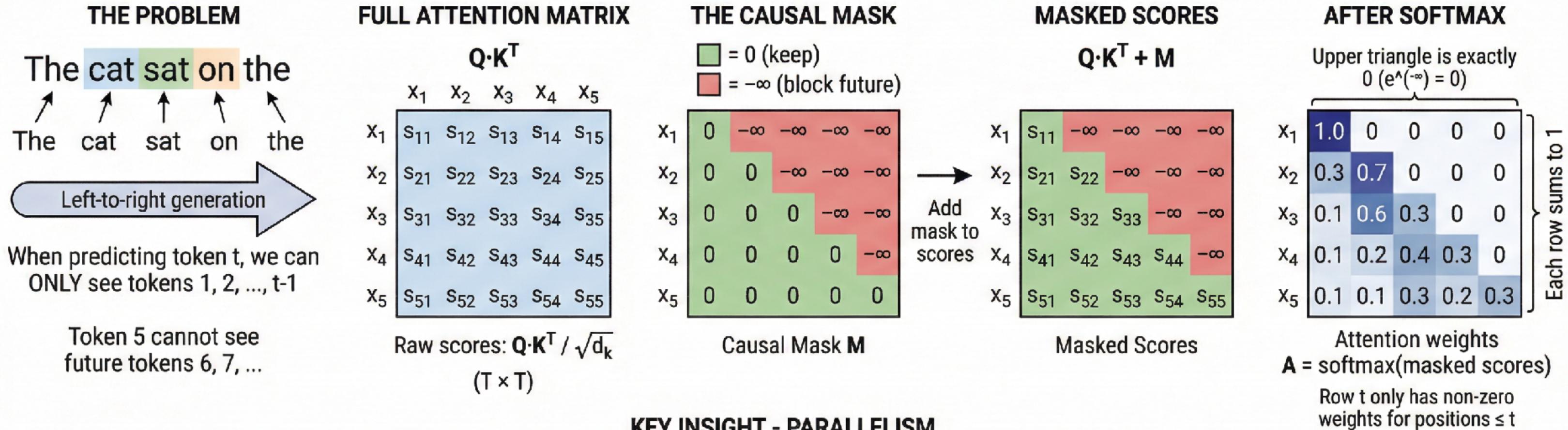# Deep Learning (1470)

## Randall Balestriero

**Class 14: Large Language Models**
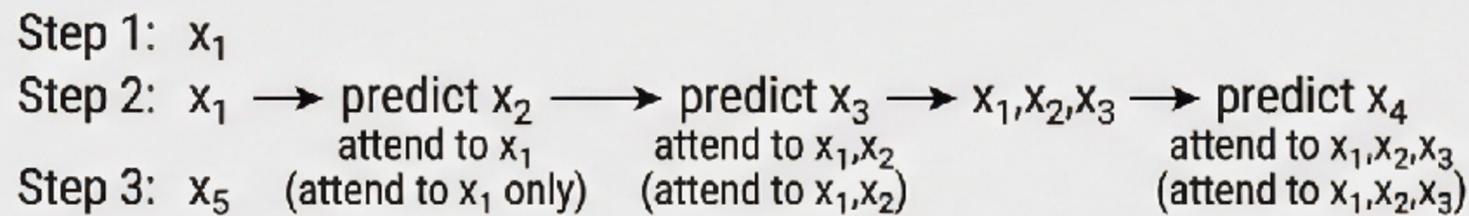
# Recap!

# Causal Self-Attention

## Masked attention for autoregressive models (GPT, LLaMA, etc.)

**THE PROBLEM**

The cat sat on the

The  cat  sat  on  the

Left-to-right generation

When predicting token t, we can ONLY see tokens 1, 2, ..., t-1

Token 5 cannot see future tokens 6, 7, ...

**FULL ATTENTION MATRIX**

$Q \cdot K^T$

|       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| $x_1$ | $s_{11}$ | $s_{12}$ | $s_{13}$ | $s_{14}$ | $s_{15}$ |
| $x_2$ | $s_{21}$ | $s_{22}$ | $s_{23}$ | $s_{24}$ | $s_{25}$ |
| $x_3$ | $s_{31}$ | $s_{32}$ | $s_{33}$ | $s_{34}$ | $s_{35}$ |
| $x_4$ | $s_{41}$ | $s_{42}$ | $s_{43}$ | $s_{44}$ | $s_{45}$ |
| $x_5$ | $s_{51}$ | $s_{52}$ | $s_{53}$ | $s_{54}$ | $s_{55}$ |

Raw scores: $Q \cdot K^T / \sqrt{d_k}$

$(T \times T)$

**THE CAUSAL MASK**

☐ = 0 (keep)
☐ = $-\infty$ (block future)

|       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| $x_1$ | 0 | $-\infty$ | $-\infty$ | $-\infty$ | $-\infty$ |
| $x_2$ | 0 | 0 | $-\infty$ | $-\infty$ | $-\infty$ |
| $x_3$ | 0 | 0 | 0 | $-\infty$ | $-\infty$ |
| $x_4$ | 0 | 0 | 0 | 0 | $-\infty$ |
| $x_5$ | 0 | 0 | 0 | 0 | 0 |

Causal Mask **M**

Add mask to scores

**MASKED SCORES**

$Q \cdot K^T + M$

|       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| $x_1$ | $s_{11}$ | $-\infty$ | $-\infty$ | $-\infty$ | $-\infty$ |
| $x_2$ | $s_{21}$ | $s_{22}$ | $-\infty$ | $-\infty$ | $-\infty$ |
| $x_3$ | $s_{31}$ | $s_{32}$ | $s_{33}$ | $-\infty$ | $-\infty$ |
| $x_4$ | $s_{41}$ | $s_{42}$ | $s_{43}$ | $s_{44}$ | $-\infty$ |
| $x_5$ | $s_{51}$ | $s_{52}$ | $s_{53}$ | $s_{54}$ | $s_{55}$ |

Masked Scores

**AFTER SOFTMAX**

Upper triangle is exactly 0 ($e^{(-\infty)} = 0$)

|       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|-------|-------|-------|-------|-------|-------|
| $x_1$ | 1.0 | 0 | 0 | 0 | 0 |
| $x_2$ | 0.3 | 0.7 | 0 | 0 | 0 |
| $x_3$ | 0.1 | 0.6 | 0.3 | 0 | 0 |
| $x_4$ | 0.1 | 0.2 | 0.4 | 0.3 | 0 |
| $x_5$ | 0.1 | 0.1 | 0.3 | 0.2 | 0.3 |

Each row sums to 1

Attention weights
**A** = softmax(masked scores)

Row t only has non-zero weights for positions ≤ t

## KEY INSIGHT - PARALLELISM

**Sequential Generation (Inference)**

Step 1:  $x_1$

Step 2:  $x_1$ → predict $x_2$ → predict $x_3$ → $x_1, x_2, x_3$ → predict $x_4$
attend to $x_1$         attend to $x_1, x_2$         attend to $x_1, x_2, x_3$
(attend to $x_1$ only)    (attend to $x_1, x_2$)       (attend to $x_1, x_2, x_3$)

Step 3:  $x_5$

...

T sequential steps

Masking enables parallel training while maintaining causal property

≡

**Parallel Training**

All T positions computed simultaneously
Single matrix multiplication with mask
**Same result as sequential, but parallel!**
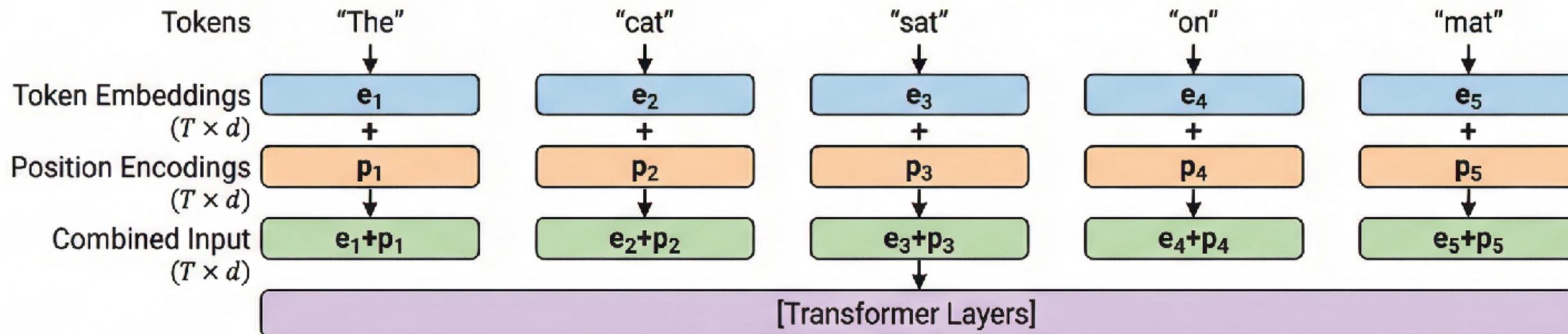
1 parallel step (same result!)

Causal Attention: $\quad \mathbf{A} = \text{softmax}\left(\dfrac{\mathbf{Q} \cdot \mathbf{K}^T + \mathbf{M}}{\sqrt{d_k}}\right) \quad \mathbf{Z} = \mathbf{A} \cdot \mathbf{V} \quad$ where $M_{ij} = 0$ if $j \leq i$, else $-\infty$

# How to encode position?

# Positional Encodings

## Injecting position information into the transformer
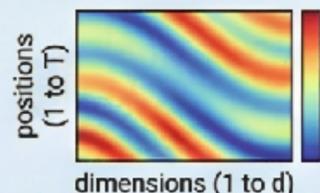
## SECTION 1: WHERE TO ADD

| | "The" | "cat" | "sat" | "on" | "mat" |
|---|---|---|---|---|---|
| Tokens | ↓ | ↓ | ↓ | ↓ | ↓ |
| Token Embeddings $(T \times d)$ | $e_1$ | $e_2$ | $e_3$ | $e_4$ | $e_5$ |
| | + | + | + | + | + |
| Position Encodings $(T \times d)$ | $p_1$ | $p_2$ | $p_3$ | $p_4$ | $p_5$ |
| Combined Input $(T \times d)$ | $e_1+p_1$ | $e_2+p_2$ | $e_3+p_3$ | $e_4+p_4$ | $e_5+p_5$ |

[Transformer Layers]

Position encoding added to token embeddings **BEFORE** transformer layers

## SECTION 2: TYPES OF POSITIONAL ENCODINGS

### Sinusoidal (Original Transformer)

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$



positions (1 to T) / dimensions (1 to d)

- Fixed (not learned)
- Deterministic
- Can extrapolate to longer sequences

### Learned (BERT, GPT-2)

$$\mathbf{P} \in \mathbb{R}^{T_{max} \times d}$$

"Lookup table of learnable vectors"

$$\mathbf{P} \begin{array}{c} p_1 \\ p_2 \\ \vdots \\ p_T \end{array}$$

$(T_{max} \times d)$

- Learned during training
- More flexible
- Limited to max sequence length $T_{max}$

### Relative (Transformer-XL, T5)

"Encode relative distance $(i - j)$ not absolute position"

$a_{ij}$ depends on $(i - j)$

|     | -2 | -1 | 0 | +1 | +2 |
|-----|----|----|----|----|----|
| -2  | -2 |    |    |    |    |
| -1  |    | -1 |    |    |    |
| 0   |    |    | 0 |    |    |
| +1  |    |    |    | +1 |    |
| +2  |    |    |    |    | +2 |

- Captures relative distance
- Better for long sequences
- Added in attention computation

### RoPE / Rotary (LLaMA, GPT-NeoX)

"Rotate **Q** and **K** vectors based on position"

$$\mathbf{Q}'_m = \mathbf{R}_m \cdot \mathbf{Q}_m$$

$$\mathbf{K}'_n = \mathbf{R}_n \cdot \mathbf{K}_n$$



- Applied to **Q**, **K** in attention
- Relative position via rotation
- Extrapolates well

## SECTION 3: VISUAL COMPARISON

| Method | Where Added | Learned? | Extrapolation |
|---|---|---|---|
| Sinusoidal | Input | No | ✓ Good |
| Learned | Input | Yes | ✗ Limited |
| Relative | Attention | Yes | ✓ Good |
| RoPE | Q, K | No | ✓ Good |

## EQUATIONS BOX

**Input to transformer:**

$$\mathbf{X} = \text{TokenEmbed(tokens)} + \text{PositionEncode(positions)}$$

**For sinusoidal:**

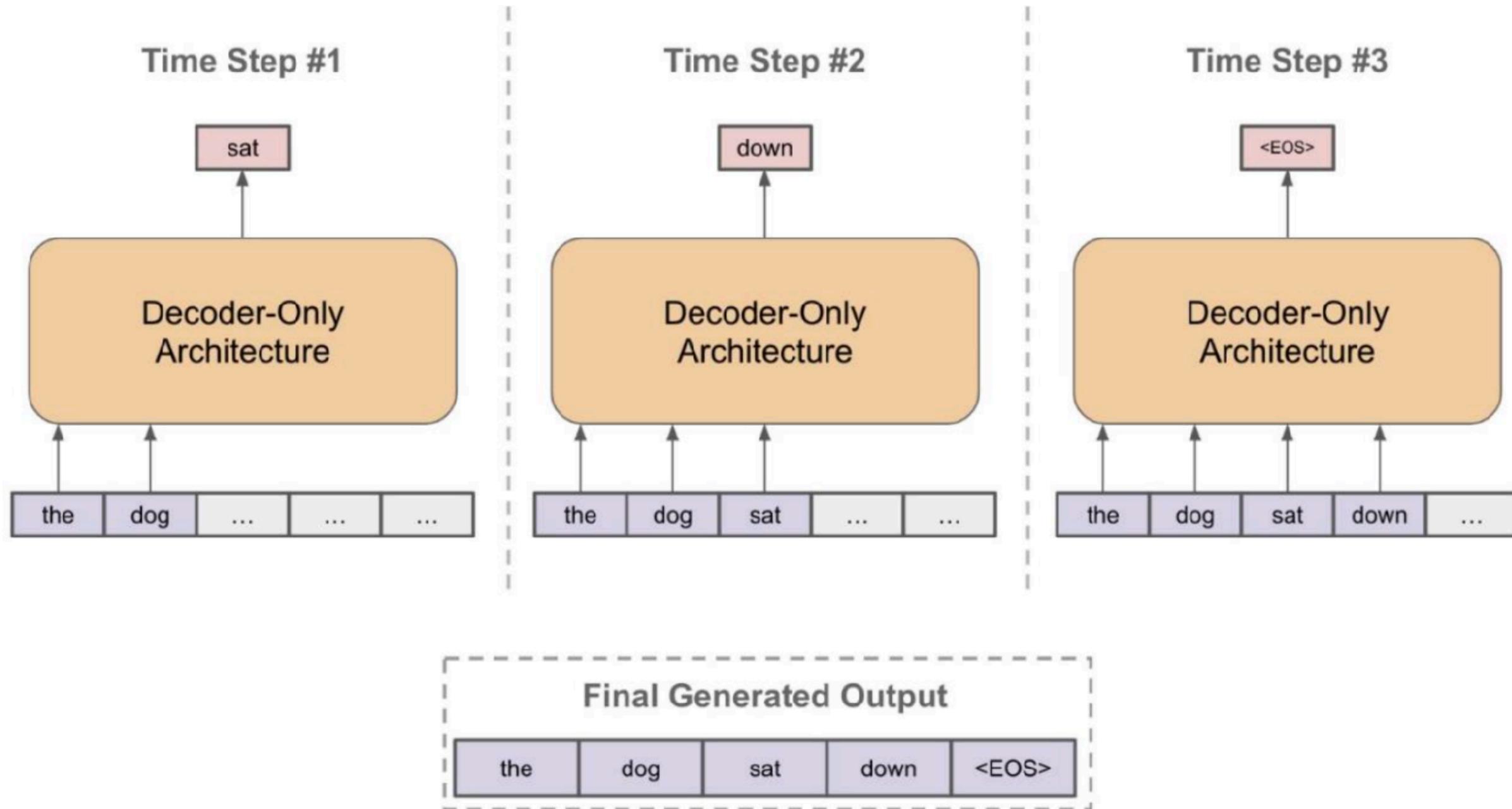$$PE(t, 2i) = \sin\left(\frac{t}{10000^{2i/d}}\right) \qquad PE(t, 2i+1) = \cos\left(\frac{t}{10000^{2i/d}}\right)$$

# LLM Hyperparameters

| | OLMo-7B | LLaMA2-7B | OpenLM-7B | Falcon-7B | PaLM-8B |
|---|---|---|---|---|---|
| Dimension | 4096 | 4096 | 4096 | 4544 | 4096 |
| Num heads | 32 | 32 | 32 | 71 | 16 |
| Num layers | 32 | 32 | 32 | 32 | 32 |
| MLP ratio | ~8/3 | ~8/3 | ~8/3 | 4 | 4 |
| Layer norm type | non-parametric | RMSNorm | parametric | parametric | parametric |
| Positional embeddings | RoPE | RoPE | RoPE | RoPE | RoPE |
| Attention variant | full | GQA | full | MQA | MQA |
| Biases | none | none | in LN only | in LN only | none |
| Block type | sequential | sequential | sequential | parallel | parallel |
| Activation | SwiGLU | SwiGLU | SwiGLU | GeLU | SwiGLU |
| Sequence length | 2048 | 4096 | 2048 | 2048 | 2048 |
| Batch size (instances) | 2160 | 1024 | 2048 | 2304 | 512 |
| Batch size (tokens) | ~4M | ~4M | ~4M | ~4M | ~1M |
| Weight tying | no | no | no | no | yes |

Table 2: LM architecture comparison at the 7–8B scale. In the "layer norm type" row, "parametric" and "non-parametric" refer to the usual layer norm implementation with and without adaptive gain and bias, respectively.

# What do we do after training?

# Generating Autoregressive Output

## Time Step #1

sat

Decoder-Only
Architecture

| the | dog | ... | ... | ... |

## Time Step #2

down

Decoder-Only
Architecture

| the | dog | sat | ... | ... |

## Time Step #3

<EOS>

Decoder-Only
Architecture

| the | dog | sat | down | ... |

### Final Generated Output

| the | dog | sat | down | <EOS> |

**Sampling Strategies from Next-Token Distribution**

**Original Distribution**

the 0.35, a 0.25, one 0.18, my 0.12, some 0.07, his 0.03

**Greedy (argmax)**

Always pick max prob!

the 0.35, a 0.25, one 0.18, my 0.12, some 0.07, his 0.03

**Temperature T=0.5 (sharper)**

the 0.52, a 0.26, one 0.14, my 0.06, some 0.02, his 0.00

**Temperature T=2.0 (flatter)**

the 0.26, a 0.22, one 0.18, my 0.15, some 0.12, his 0.08

**Top-k Sampling (k=3)**

the 0.45, a 0.32, one 0.23, my 0.00, some 0.00, his 0.00

**Top-p Nucleus (p=0.8)**

the 0.39, a 0.28, one 0.20, my 0.13, some 0.00, his 0.00

# Can we train ChatGPT now?

# Can we train ChatGPT now?

Yes, but only the poor version

# Language Modeling is not Enough

# Language Modeling is not Enough

- Model is repetitive

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

- Does not feel interactive

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

- Does not feel interactive

- Has unconstrained output (dangerous)

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

- Does not feel interactive

- Has unconstrained output (dangerous)

- We need to "teach it" what to say more precisely!

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

- Does not feel interactive

- Has unconstrained output (dangerous)

- We need to "teach it" what to say more precisely!

  - Supervised finetuning

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

- Does not feel interactive

- Has unconstrained output (dangerous)

- We need to "teach it" what to say more precisely!

  - Supervised finetuning

  - RLHF

# Language Modeling is not Enough

- Model is repetitive

- Does not follow instruction

- Does not feel interactive

- Has unconstrained output (dangerous)

- We need to "teach it" what to say more precisely!

  - Supervised finetuning

  - RLHF

  - ...

# Turning GPT to Chat-GPT

**Step 0: Train GPT**

## Step 1

**Collect demonstration data and train a supervised policy.**

A prompt is sample from our prompt dataset.

Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT

## Step 2

**Collect comparison data and train a reward model.**

A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

A: In reinforcement learning, the agent is...
B: Explain rewards...
C: In machine learning...
D: We give treats and punishments to teach...

A labeler ranks the outputs from best to worst.

D > C > A > B

This data is used to train our reward model.

RM

D > C > A > B

## Step 3

**Optimize a policy against the reward model using the PPO reinforcement learning algorithm.**

A new prompt is sampled from the dataset.

Write a story about otters.

The PPO model is initialized from the supervised policy.

PPO

The policy generates an output.

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.
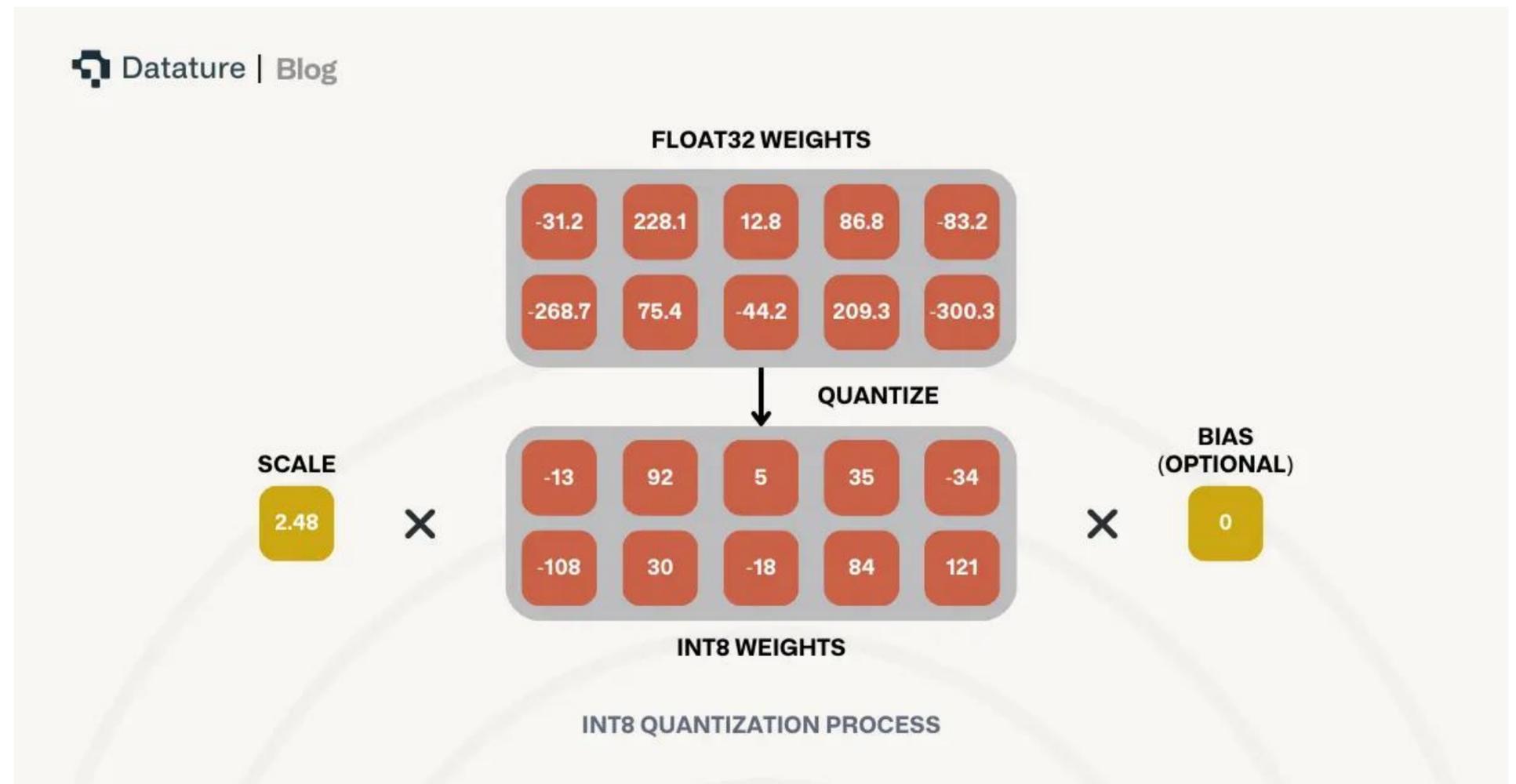
$r_k$

Source: OpenAI

Now that you can generate nice content, how to speed it up?

# Quantization

Can we use smaller representation of parameters?

DeepSeek was able to create distilled and quantized models that only used 4 bits per parameter

https://huggingface.co/neuralmagic/DeepSeek-R1-Distill-Llama-8B-quantized.w4a16
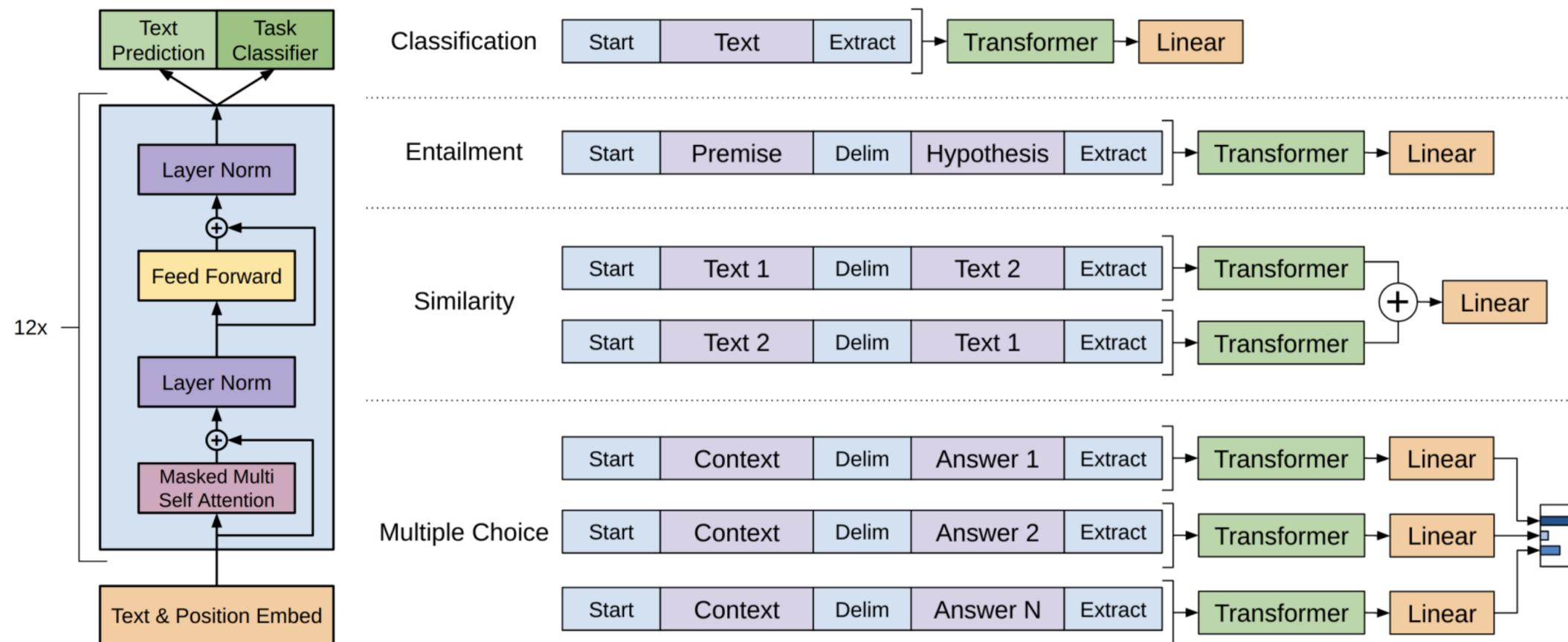
# But you can do much more



Figure 1: **(left)** Transformer architecture and training objectives used in this work. **(right)** Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

Alec Radford et al., Improving Language Understanding by Generative Pre-Training, 2018

# The Era of Foundation Models

# Foundation Models



https://dataforest.ai/blog/ai-foundation-models-for-big-business-innovation

# Large Language Model Scaling "Laws"

The bigger the better



Larger models require **fewer samples** to reach the same performance
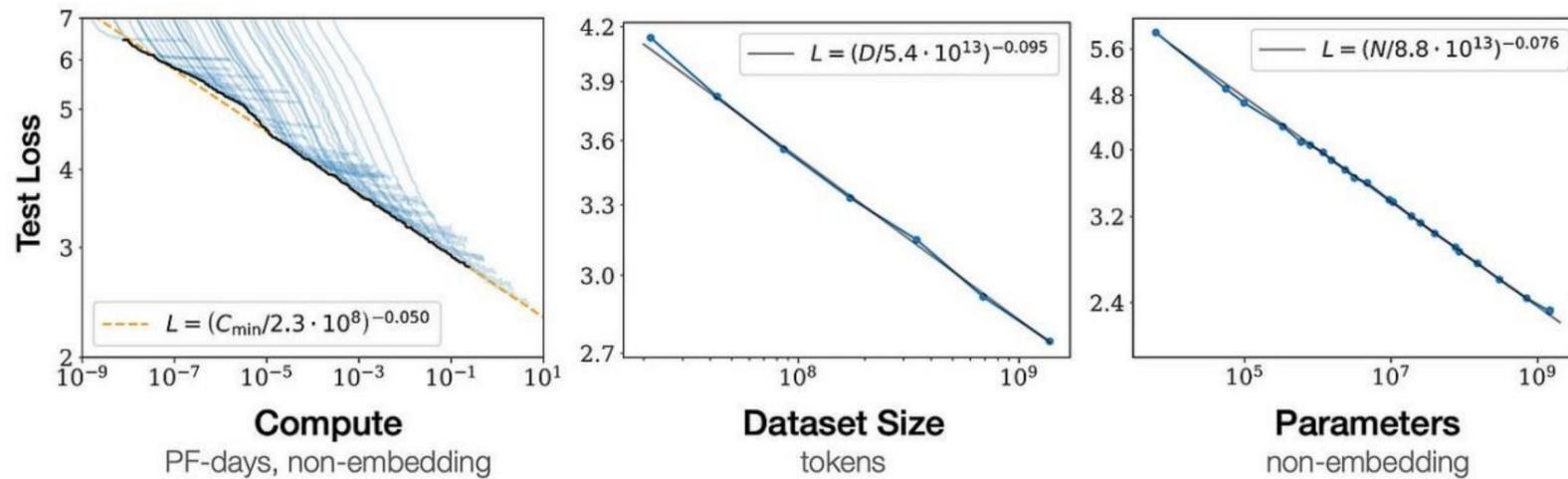
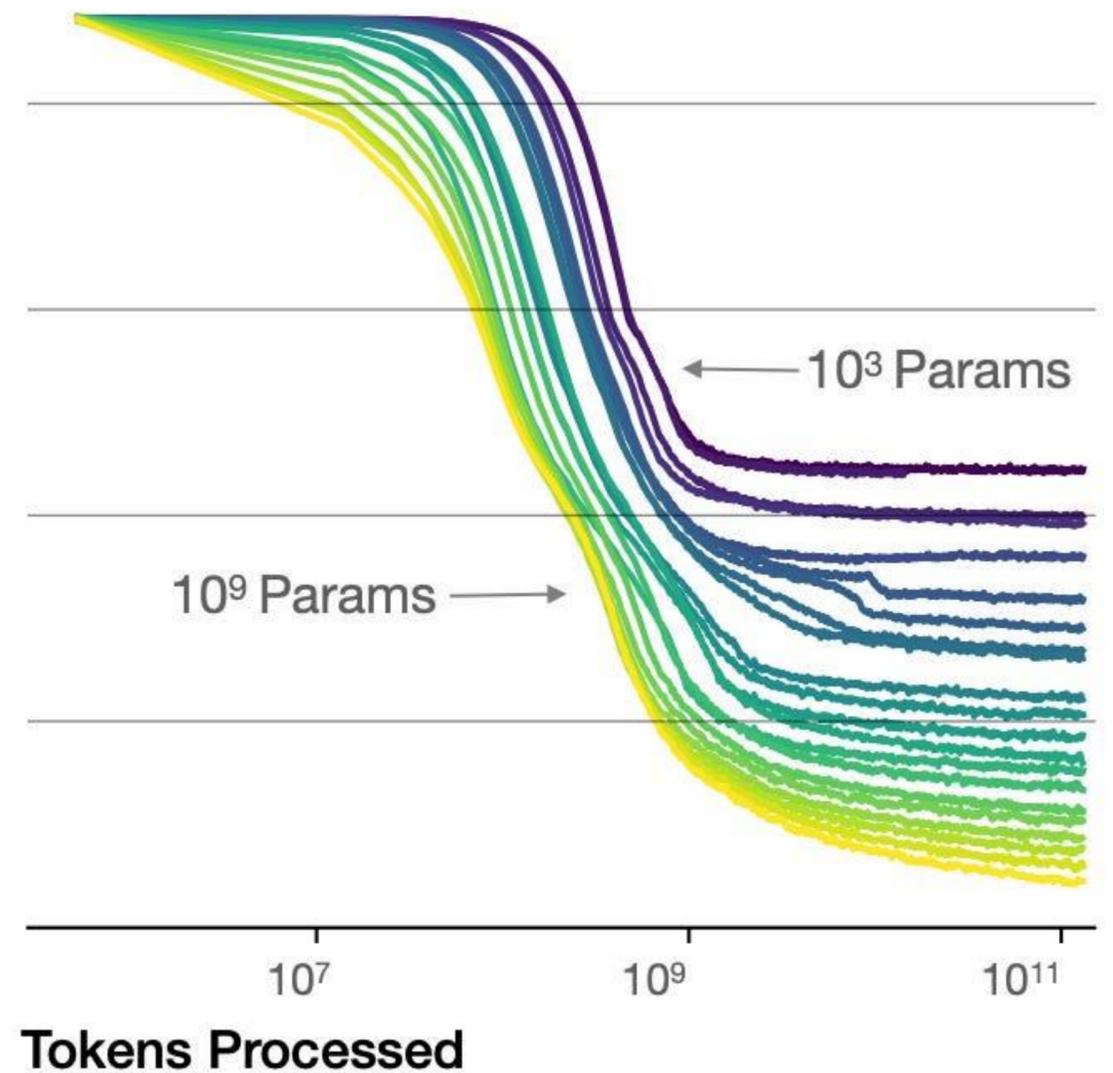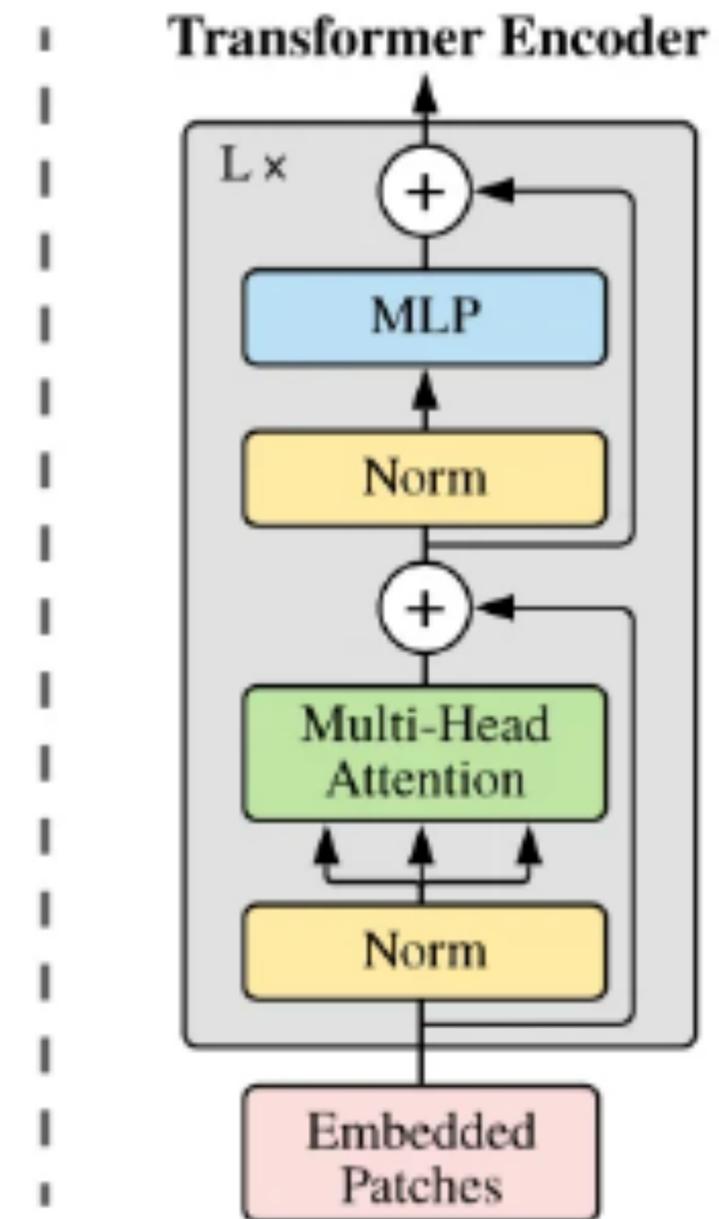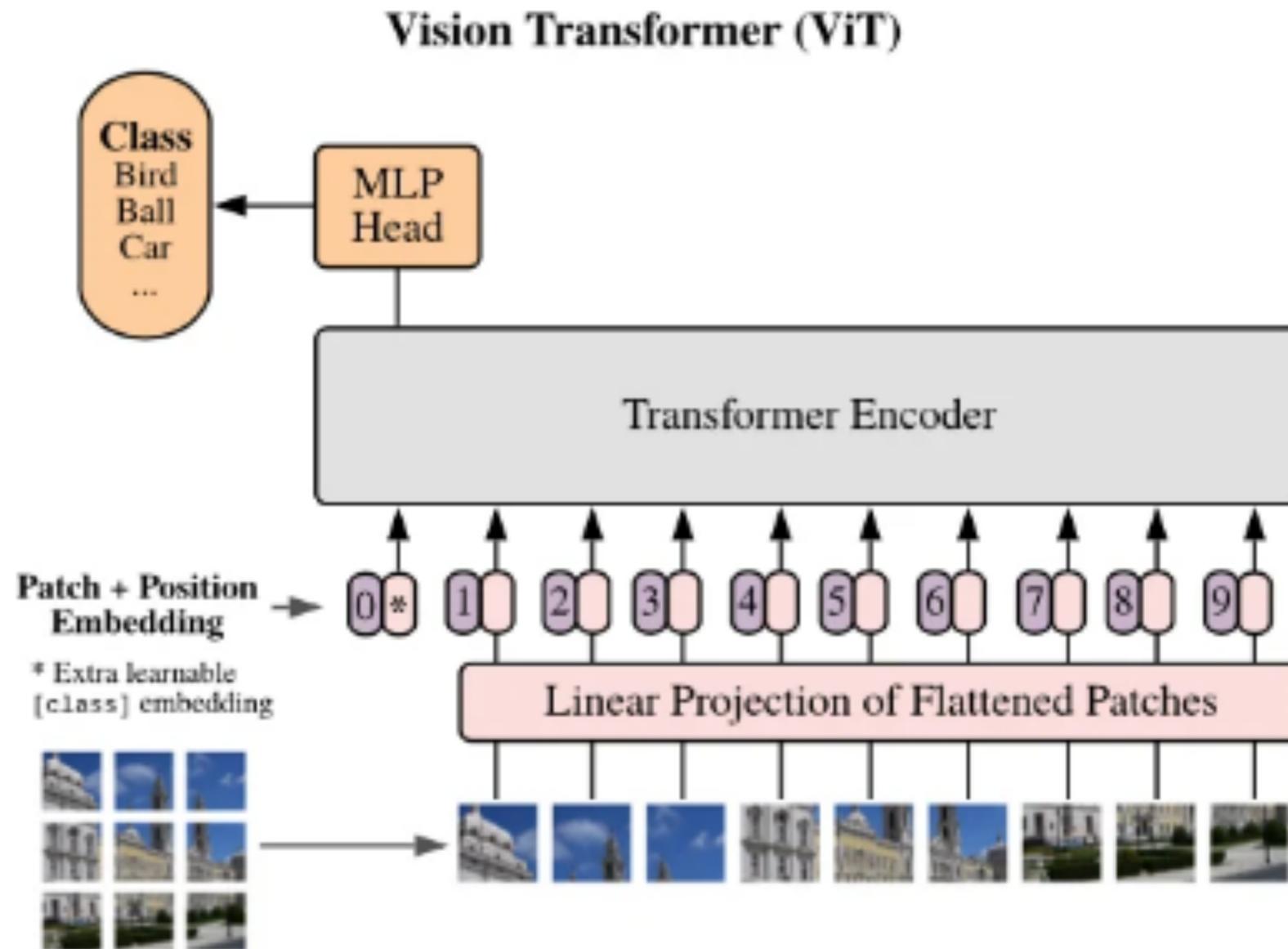Test Loss

$10^3$ Params

$10^9$ Params

Tokens Processed

Figure 1 Language modeling performance improves smoothly as we increase the model size, datasetset size, and amount of compute[2] used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

Kaplan et al. "Scaling Laws for Neural Language Models"

# How to move Beyond Language?

# Vision Transformers!

See you on Wednesday!