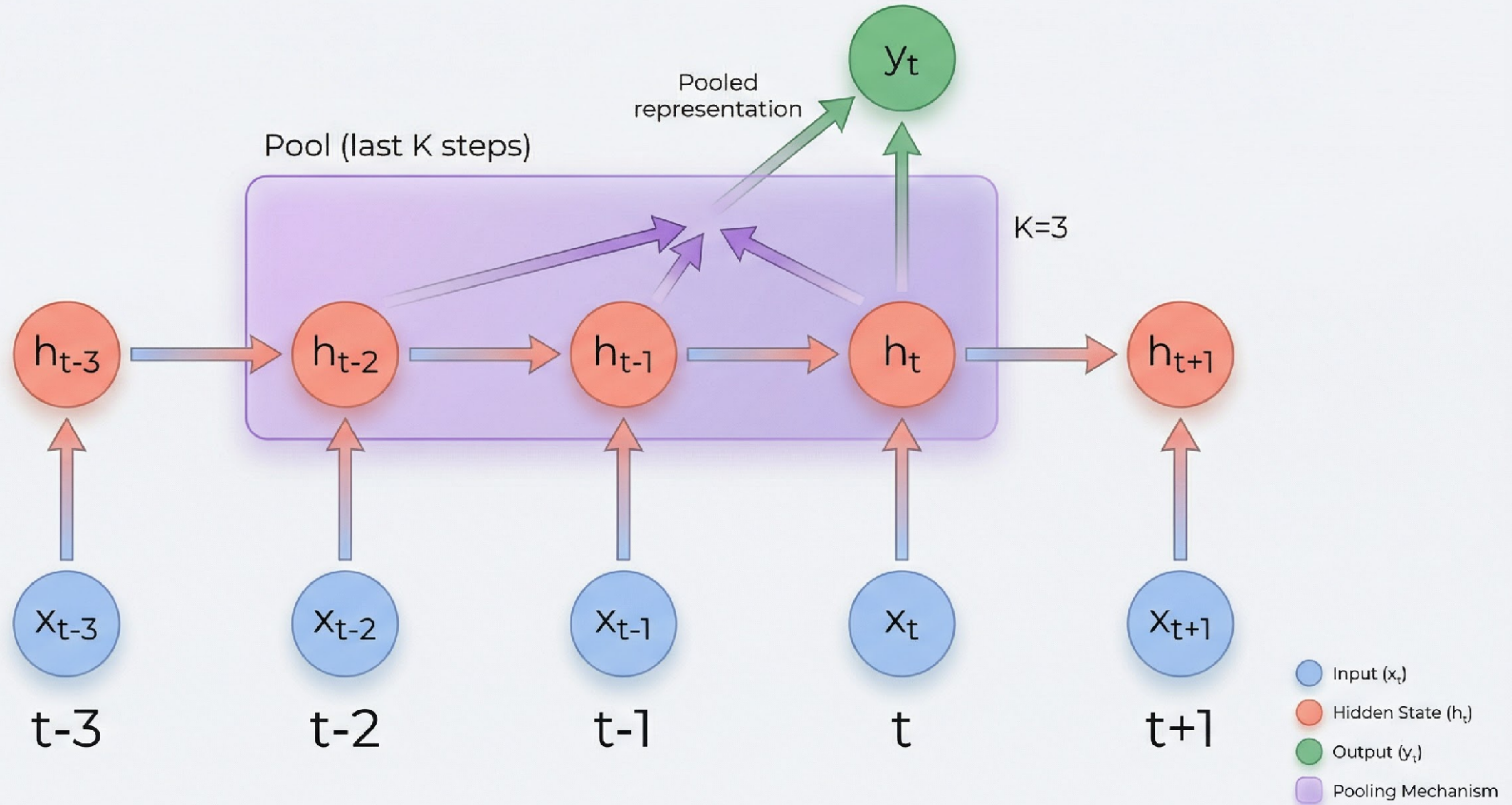# Deep Learning (1470)

**Randall Balestriero**

**Class 13: Transformers**

# Recap!

Recurrent Neural Network with Pooling Architecture

# Transformers anyone?
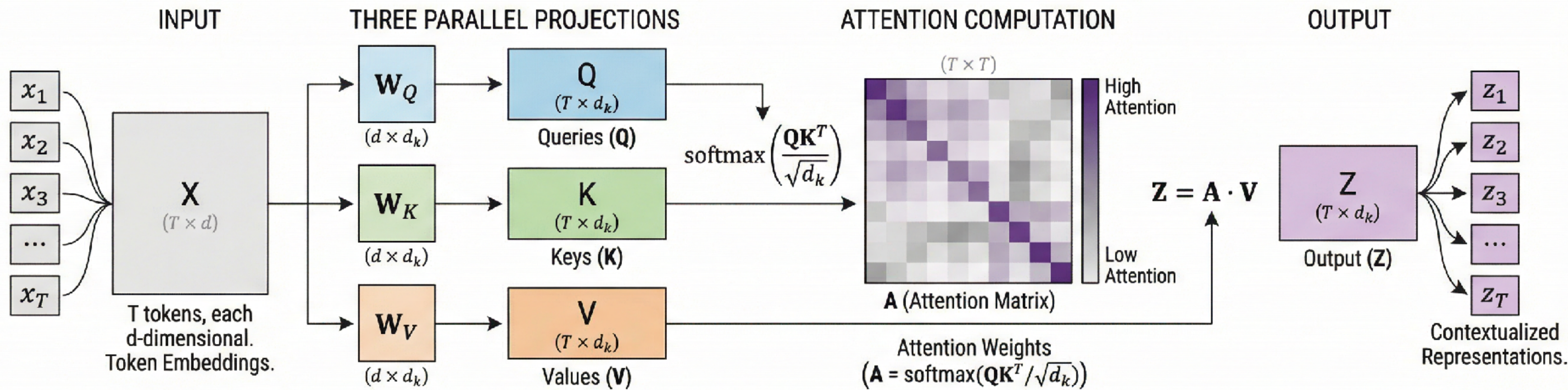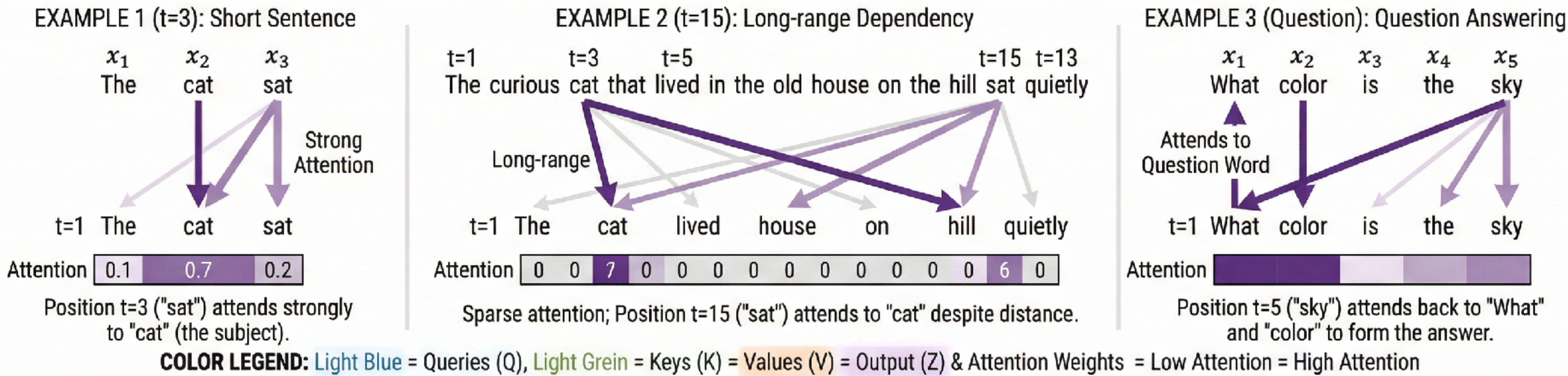
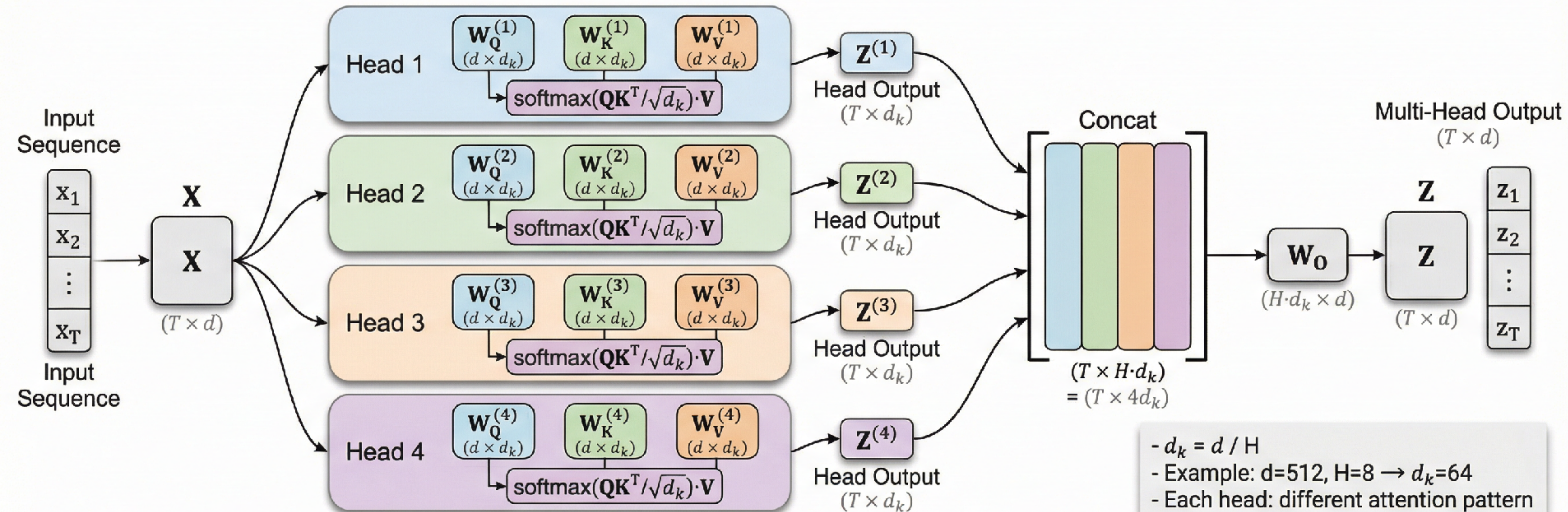# Transformers anyone?

# SINGLE HEAD SELF-ATTENTION MECHANISM



INPUT

$x_1$
$x_2$
$x_3$
$\cdots$
$x_T$

$X$
$(T \times d)$

T tokens, each
d-dimensional.
Token Embeddings.

THREE PARALLEL PROJECTIONS

$W_Q$
$(d \times d_k)$

$W_K$
$(d \times d_k)$

$W_V$
$(d \times d_k)$

$Q$
$(T \times d_k)$
Queries (Q)

$K$
$(T \times d_k)$
Keys (K)

$V$
$(T \times d_k)$
Values (V)

$\text{softmax}\left(\dfrac{QK^T}{\sqrt{d_k}}\right)$

ATTENTION COMPUTATION

$(T \times T)$

High Attention

Low Attention

A (Attention Matrix)

Attention Weights
$(A = \text{softmax}(QK^T/\sqrt{d_k}))$

$Z = A \cdot V$

OUTPUT

$Z$
$(T \times d_k)$
Output (Z)

$z_1$
$z_2$
$z_3$
$\cdots$
$z_T$

Contextualized
Representations.

---

# CONCRETE EXAMPLES OF ATTENTION PATTERNS FOR LANGUAGE MODELING

EXAMPLE 1 (t=3): Short Sentence

$x_1$     $x_2$     $x_3$
The      cat      sat

Strong Attention

t=1  The    cat    sat

Attention | 0.1 | 0.7 | 0.2 |

Position t=3 ("sat") attends strongly
to "cat" (the subject).

EXAMPLE 2 (t=15): Long-range Dependency

t=1        t=3      t=5                                    t=15  t=13
The curious cat that lived in the old house on the hill sat quietly

Long-range

t=1  The    cat    lived   house    on    hill    quietly

Attention | 0 | 0 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 |

Sparse attention; Position t=15 ("sat") attends to "cat" despite distance.

EXAMPLE 3 (Question): Question Answering

$x_1$     $x_2$     $x_3$     $x_4$     $x_5$
What     color     is       the      sky

Attends to
Question Word

t=1  What   color   is    the    sky

Attention

Position t=5 ("sky") attends back to "What"
and "color" to form the answer.

**COLOR LEGEND:** Light Blue = Queries (Q), Light Grein = Keys (K) = Values (V) = Output (Z) & Attention Weights = Low Attention = High Attention

# Multi-Head Self-Attention

**Input Sequence**

$x_1$
$x_2$
$\vdots$
$x_T$

**Input Sequence**

$\mathbf{X}$

$\mathbf{X}$

$(T \times d)$

**Head 1**

$\mathbf{W_Q^{(1)}}$ $(d \times d_k)$   $\mathbf{W_K^{(1)}}$ $(d \times d_k)$   $\mathbf{W_V^{(1)}}$ $(d \times d_k)$

$\text{softmax}(\mathbf{QK^T}/\sqrt{d_k}) \cdot \mathbf{V}$

$\mathbf{Z^{(1)}}$

Head Output $(T \times d_k)$

**Head 2**

$\mathbf{W_Q^{(2)}}$ $(d \times d_k)$   $\mathbf{W_K^{(2)}}$ $(d \times d_k)$   $\mathbf{W_V^{(2)}}$ $(d \times d_k)$

$\text{softmax}(\mathbf{QK^T}/\sqrt{d_k}) \cdot \mathbf{V}$

$\mathbf{Z^{(2)}}$

Head Output $(T \times d_k)$

**Head 3**

$\mathbf{W_Q^{(3)}}$ $(d \times d_k)$   $\mathbf{W_K^{(3)}}$ $(d \times d_k)$   $\mathbf{W_V^{(3)}}$ $(d \times d_k)$

$\text{softmax}(\mathbf{QK^T}/\sqrt{d_k}) \cdot \mathbf{V}$

$\mathbf{Z^{(3)}}$

Head Output $(T \times d_k)$

**Head 4**

$\mathbf{W_Q^{(4)}}$ $(d \times d_k)$   $\mathbf{W_K^{(4)}}$ $(d \times d_k)$   $\mathbf{W_V^{(4)}}$ $(d \times d_k)$

$\text{softmax}(\mathbf{QK^T}/\sqrt{d_k}) \cdot \mathbf{V}$

$\mathbf{Z^{(4)}}$

Head Output $(T \times d_k)$

**Concat**

$(T \times H \cdot d_k)$
$= (T \times 4d_k)$

$\mathbf{W_O}$

$(H \cdot d_k \times d)$

**Multi-Head Output** $(T \times d)$

$\mathbf{Z}$

$\mathbf{Z}$

$(T \times d)$

$z_1$
$z_2$
$\vdots$
$z_T$

- $d_k = d / H$
- Example: d=512, H=8 → $d_k$=64
- Each head: different attention pattern

$$\text{MultiHead}(\mathbf{X}) = \text{Concat}(\text{head}_1, \text{head}_2, \ldots, \text{head}_H) \cdot \mathbf{W_O} \quad \text{where head}_h = \text{Attention}(\mathbf{XW_Q^{(h)}}, \mathbf{XW_K^{(h)}}, \mathbf{XW_V^{(h)}})$$

**Why Multiple Heads?**

Head 1: Syntactic

$T \times T$ / $T \times T$

Head 2: Semantic

$T \times T$ / $T \times T$

Head 3: Local

$T \times T$ / $T \times T$

Head 4: Long-range

$T \times T$ / $T \times T$

Light Blue: Query (Q)
Light Green: Key (K)
Light Orange: Value (V)
Light Purple: Attention
Gray Text: Shapes

# Break: Remember batch-norm?

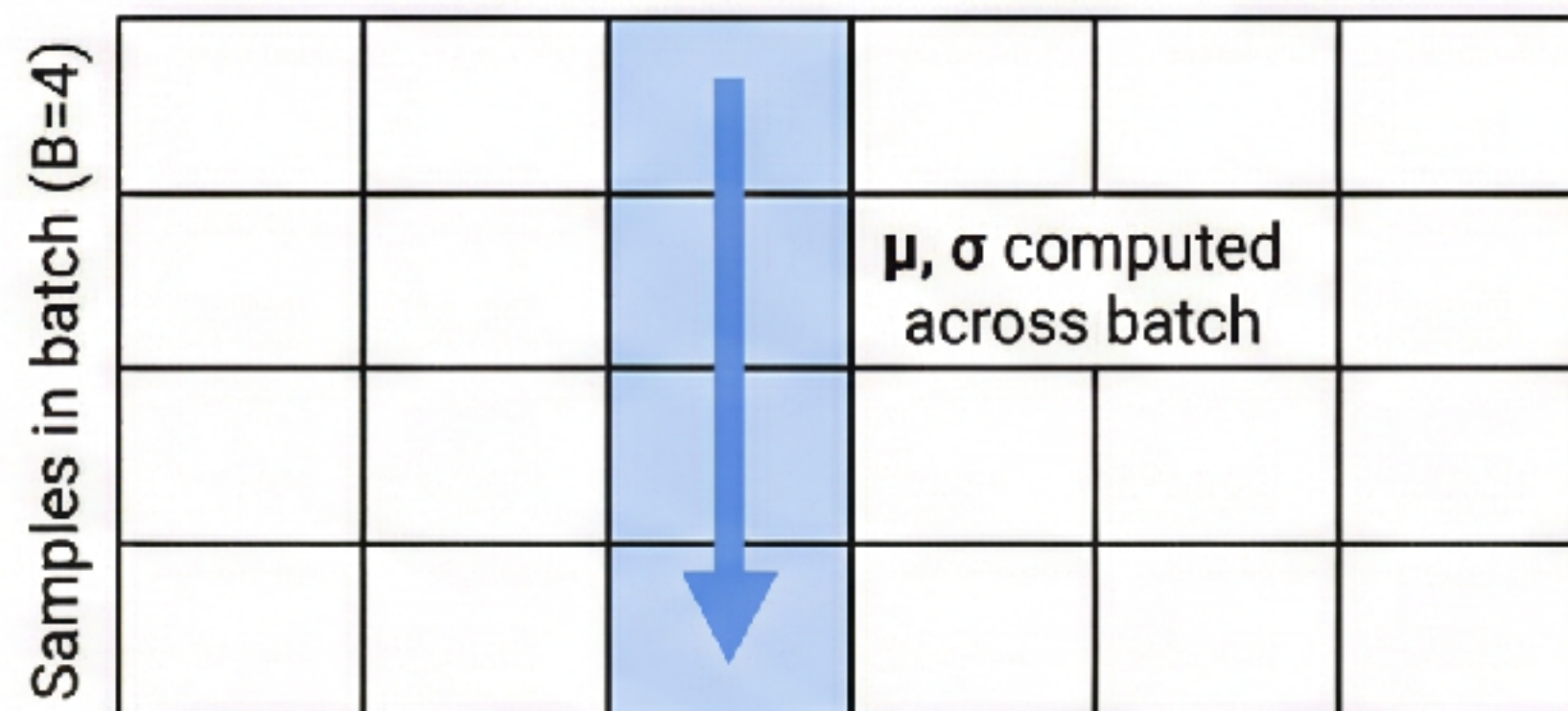# Batch Normalization vs Layer Normalization

## Why transformers use LayerNorm
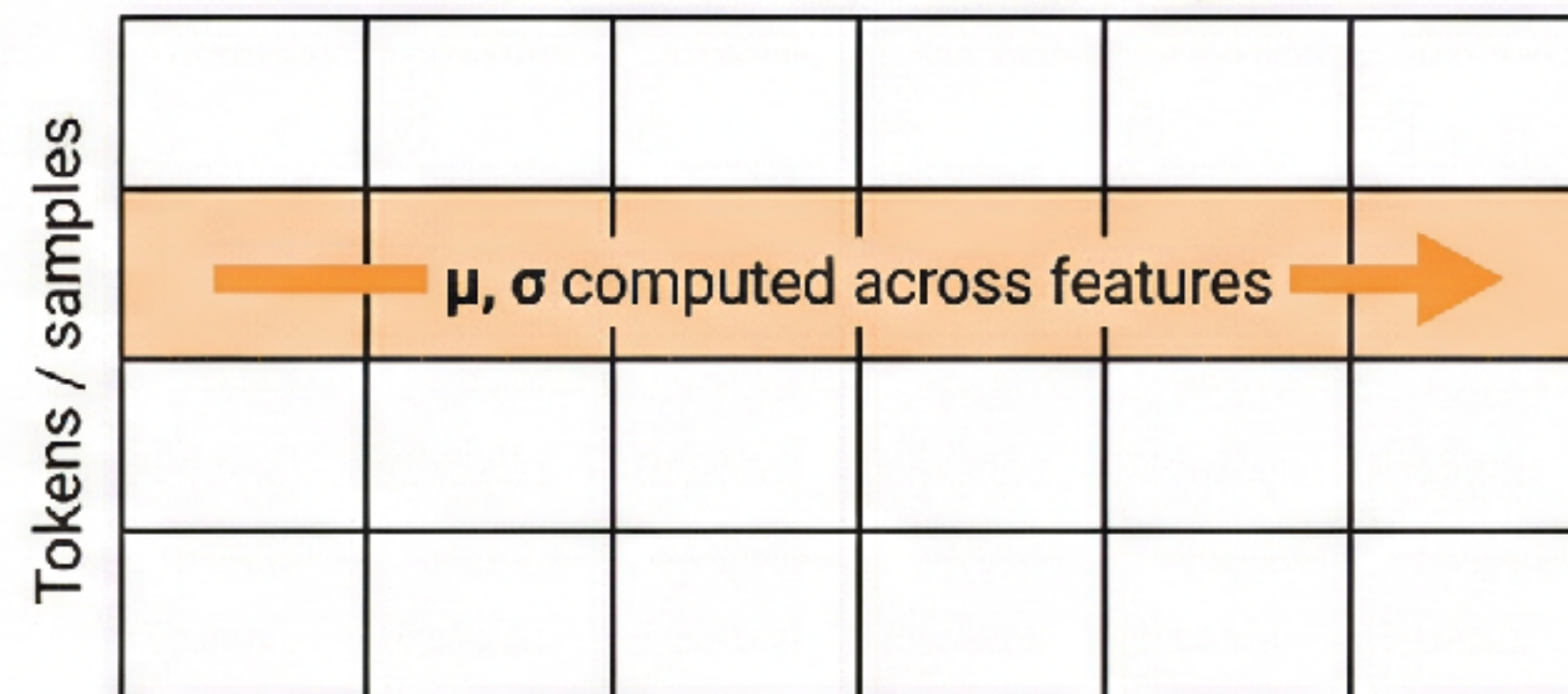
BATCH NORMALIZATION

(Batch/Tokens × Features)

Samples in batch (B=4)

$\mu$, $\sigma$ computed across batch

BatchNorm    LayerNorm

LAYER NORMALIZATION

(Batch/Tokens × Features)

Tokens / samples

$\mu$, $\sigma$ computed across features

**Batch Normalization**

$$\hat{\mathbf{x}} = \frac{\mathbf{x} - \mu_B}{\sigma_B}$$

Per feature, across batch

**Layer Normalization**

$$\hat{\mathbf{x}} = \frac{\mathbf{x} - \mu_d}{\sigma_d}$$

Per token, across features

|  | **Batch Norm** | **Layer Norm** |
|---|---|---|
| Normalizes | Across batch ($\downarrow$) | Across features ($\rightarrow$) |
| Statistics | Per feature | Per token |
| Batch size 1 | ✗ Fails | ✓ Works |
| Train=Inference | ✗ Different | ✓ Same |
| Transformers | ✗ Not used | ✓ Standard |

**LayerNorm**: each token normalized independently → perfect for variable-length sequences & autoregressive generation

$$\text{LayerNorm}(\mathbf{x}) = \gamma \cdot \frac{(\mathbf{x} - \mu)}{\sqrt{\sigma^2 + \epsilon}} + \beta$$

# Transformer Layer
## Stacked N× times



**Self-Attention Sublayer:**
$$\mathbf{X'} = \text{LayerNorm}(\mathbf{X} + \text{MultiHeadAttention}(\mathbf{X}))$$

**Feed-Forward Sublayer:**
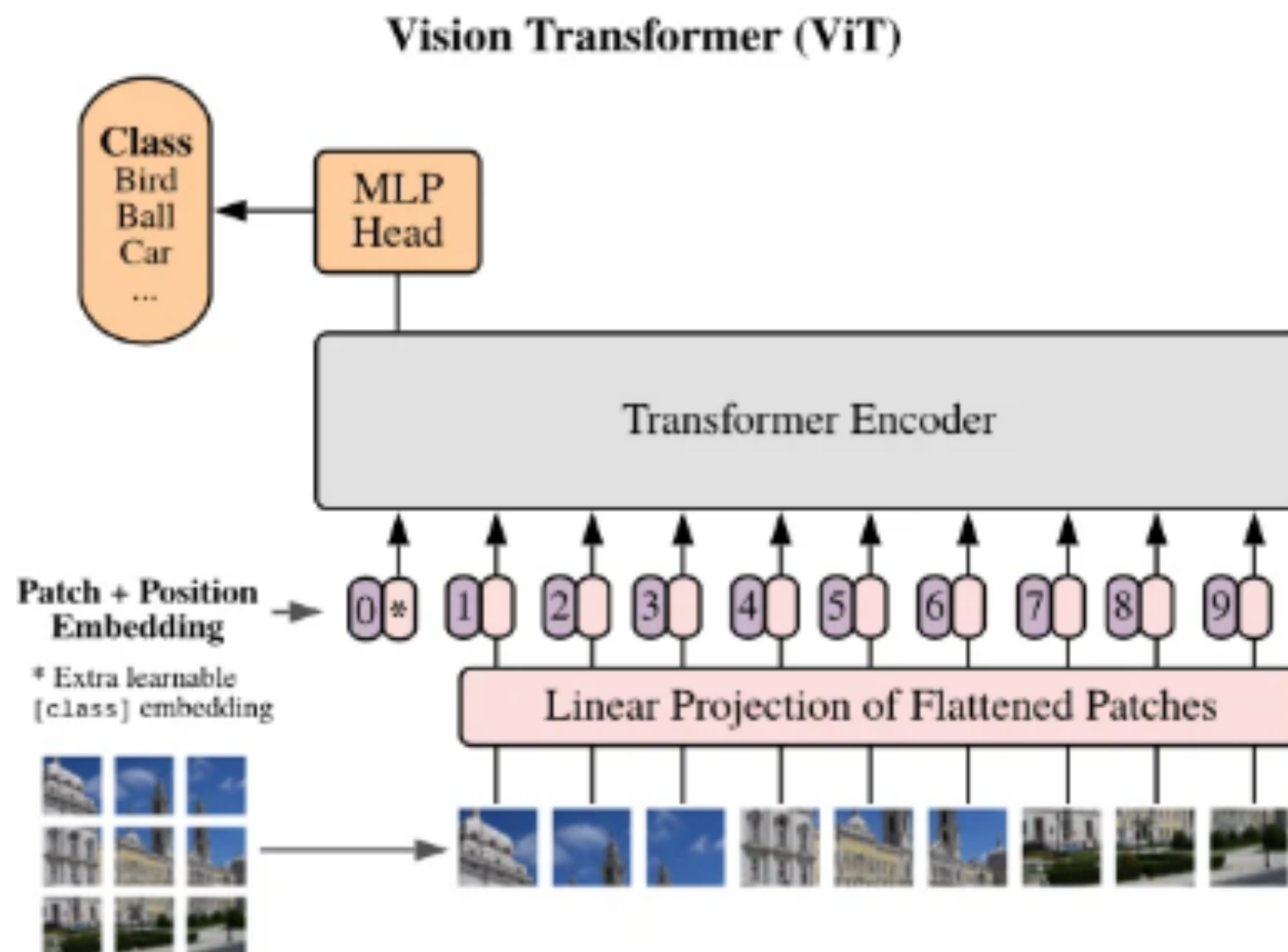$$\mathbf{X''} = \text{LayerNorm}(\mathbf{X'} + \text{FFN}(\mathbf{X'}))$$
$$\text{FFN}(\mathbf{x}) = W_2 \cdot \text{ReLU}(W_1\mathbf{x} + b_1) + b_2$$

$d$ = model dim (512, 768)
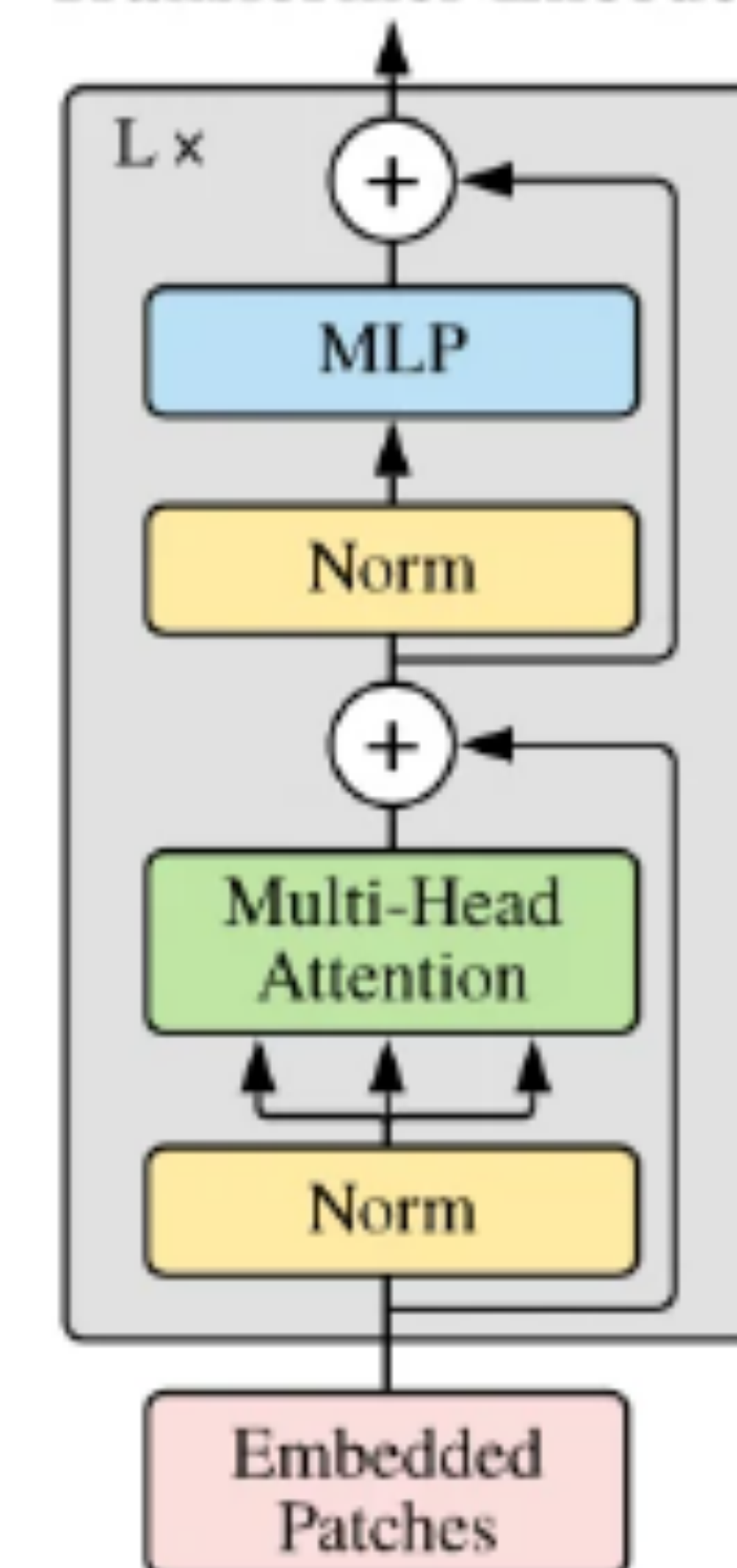$d_{ff}$ = 4d (2048, 3072)
$H$ = heads (8, 12)

Light Blue: Attention Components
Light Green: LayerNorm
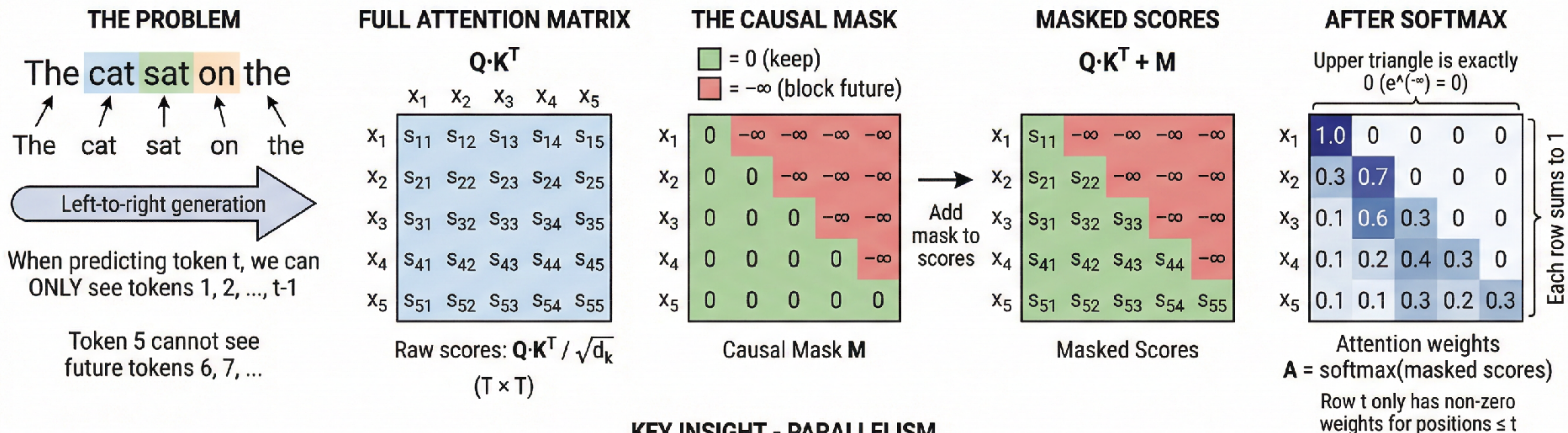Light Orange: Feed-Forward Network

# Vision Transformers!

# Break: efficient implementation for sequence modeling?

# 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)
■ = −∞ (block future)

| | | | | | |
|---|---|---|---|---|---|
| $x_1$ | 0 | −∞ | −∞ | −∞ | −∞ |
| $x_2$ | 0 | 0 | −∞ | −∞ | −∞ |
| $x_3$ | 0 | 0 | 0 | −∞ | −∞ |
| $x_4$ | 0 | 0 | 0 | 0 | −∞ |
| $x_5$ | 0 | 0 | 0 | 0 | 0 |

Causal Mask **M**

Add mask to scores

**MASKED SCORES**

$Q \cdot K^T + M$

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

Masked Scores

**AFTER SOFTMAX**

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

| | | | | | |
|---|---|---|---|---|---|
| $x_1$ | 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: $\mathbf{A} = \text{softmax}\left(\dfrac{\mathbf{Q} \cdot \mathbf{K}^T + \mathbf{M}}{\sqrt{d_k}}\right)$    $\mathbf{Z} = \mathbf{A} \cdot \mathbf{V}$    where $M_{ij} = 0$ if $j \leq i$, else $-\infty$
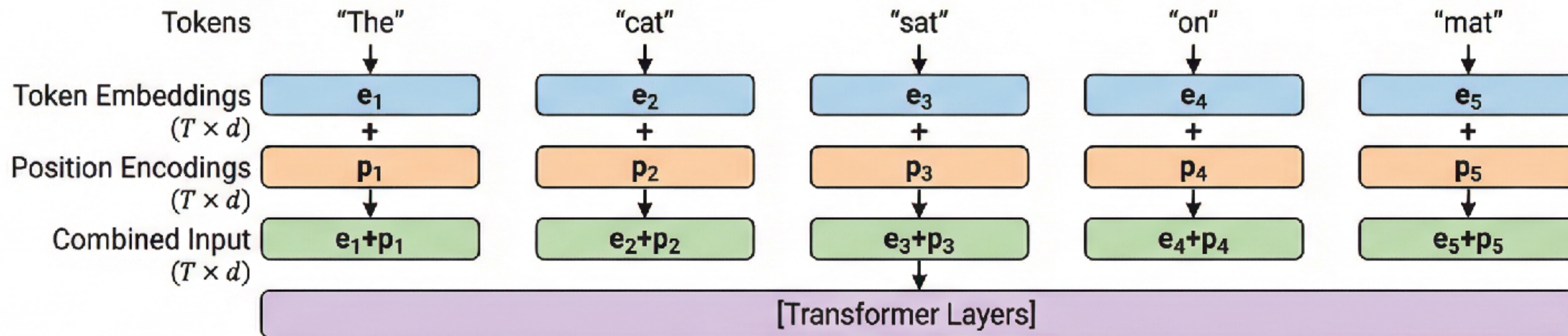
# Break: How to encode position?

# Positional Encodings

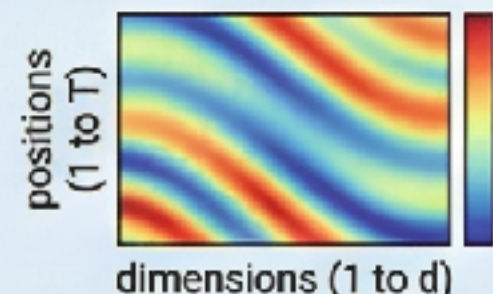Injecting position information into the transformer

## SECTION 1: WHERE TO ADD

| Tokens | "The" | "cat" | "sat" | "on" | "mat" |
|---|---|---|---|---|---|
| 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"

$\mathbf{P}$
$p_1$
$p_2$
$\vdots$
$p_T$

$(T_{max} \times d)$

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

### Relative (Transformer-XL, T5)

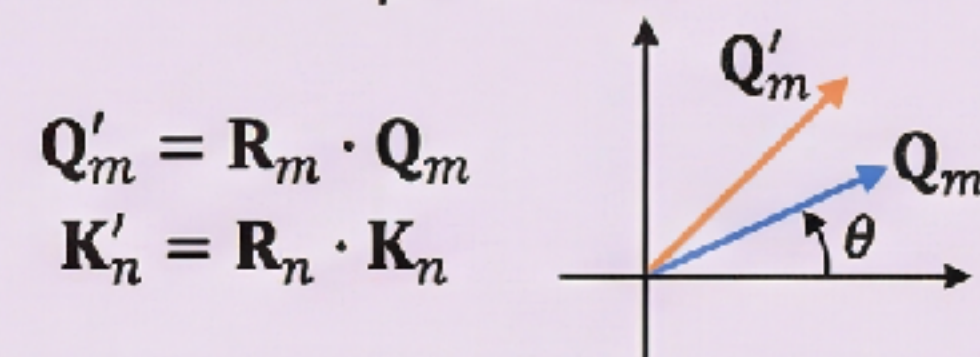"Encode relative distance $(i - j)$ not absolute position"

$a_{ij}$ 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"

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

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

$\mathbf{Q}'_m$, $\mathbf{Q}_m$, $\theta$

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

# See you Monday!