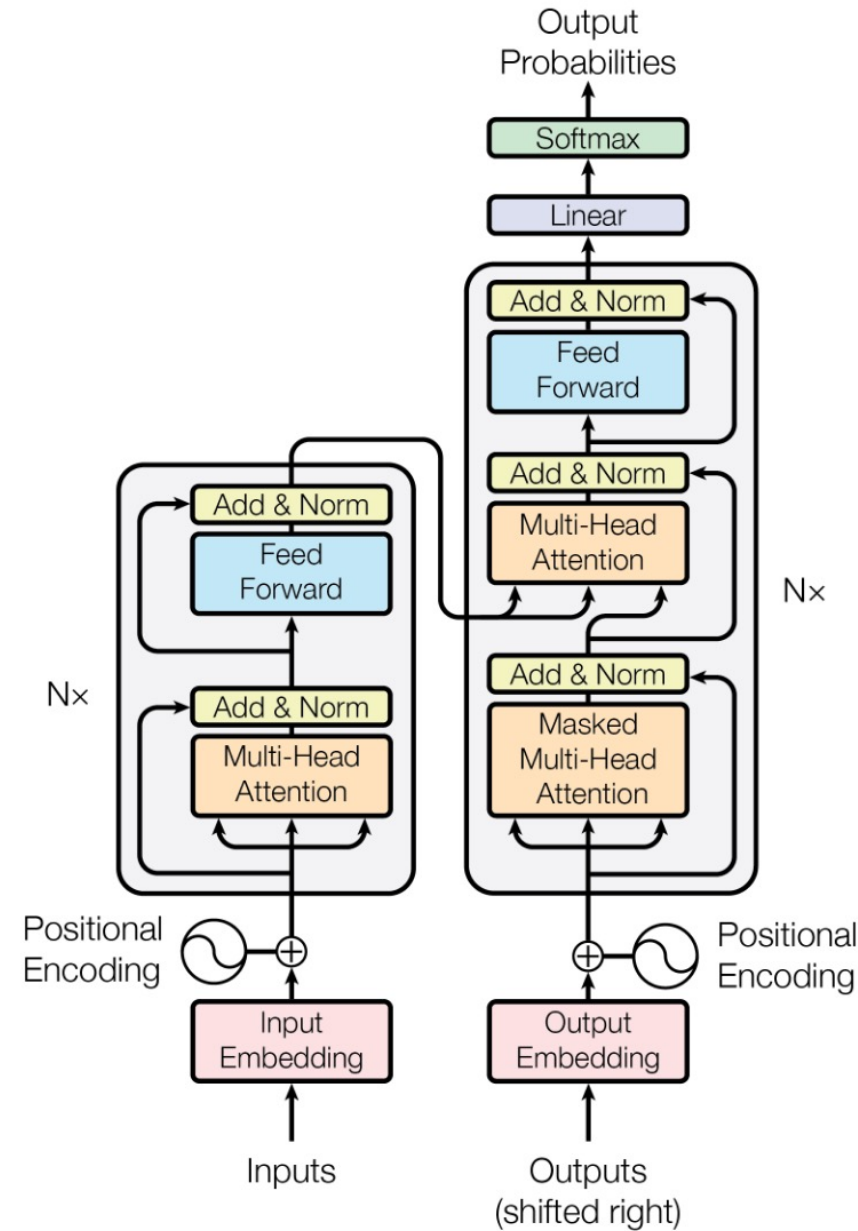# IT5005 Artificial Intelligence

Sirigina Rajendra Prasad
AY2025/2026: Semester 1

## Transformers Contd.

# Output Layer

- **Inference**:
  - Auto-regressive and sequential over time

- At each step:
  - Input to the decoder is the tokens generated so far
    - **Example**: at time $t + 1$, input to decoder is $\widehat{\boldsymbol{y}}_{1:t}$
  - Decoder cross-attends to all encoder outputs
    - No masking at cross-attention layer
  - $Q, K$ and $V$ shapes grow with time
    - "*KV caching*" to save computation time
    - Query only for the newly added position

# KV caching during inference

# Inference: at $t = 1$

$Q \in R^{1\times 2}$

| $x_1$ |
| $x_2$ |
| $x_3$ |

$W_i^Q \in R^{2\times 2}$

$[w_1^q \ w_2^q]$

$QW_i^Q = Q \in R^{1\times 2}$

| $[x_1 w_1^q \quad x_1 w_2^q] = q_1$ |
| $[x_2 w_1^q \quad x_2 w_2^q] = q_2$ |
| $[x_3 w_1^q \quad x_3 w_2^q] = q_3$ |

$X \in R^{1\times 2}$

| $x_1$ |
| $x_2$ |
| $x_3$ |

| Transformatoren |
| sind |
| großartig |

Raw Data      SE+PE

$K \in R^{1\times 2}$

| $x_1$ |
| $x_2$ |
| $x_3$ |

$W_i^K \in R^{2\times 2}$

$[w_1^k \ w_2^k]$

$KW_i^K = K \in R^{1\times 2}$

| $[x_1 w_1^k \quad x_1 w_2^k] = k_1$ |
| $[x_2 w_1^k \quad x_2 w_2^k] = k_2$ |
| $[x_3 w_1^k \quad x_3 w_2^k] = k_3$ |

$V \in R^{1\times 2}$

| $x_1$ |
| $x_2$ |
| $x_3$ |

$W_i^V \in R^{2\times 2}$

$[w_1^v \ w_2^v]$

$VW_i^V = V \in R^{1\times 2}$

| $[x_1 w_1^v \quad x_1 w_2^v] = v_1$ |
| $[x_2 w_1^v \quad x_2 w_2^v] = v_2$ |
| $[x_3 w_1^v \quad x_3 w_2^v] = v_3$ |

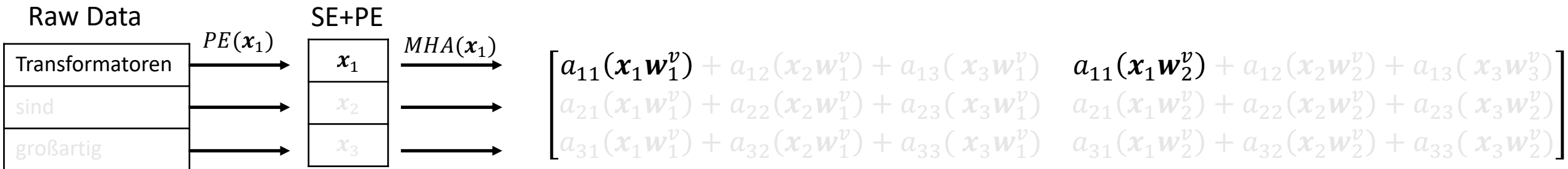$$head_i = softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ M\right)(VW_i^V)$$

|  | Transformers | are | great |
|---|---|---|---|
| Transformers | $a_{11}$ | $a_{12}$ | $a_{13}$ |
| are | $a_{21}$ | $a_{22}$ | $a_{23}$ |
| great | $a_{31}$ | $a_{32}$ | $a_{33}$ |

$$\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} = \begin{bmatrix} \dfrac{\boldsymbol{q_1 k_1^T}}{\sqrt{d_1}} & \dfrac{q_1 k_2^T}{\sqrt{d_1}} & \dfrac{q_1 k_3^T}{\sqrt{d_1}} \\ \dfrac{q_2 k_1^T}{\sqrt{d_2}} & \dfrac{q_2 k_2^T}{\sqrt{d_2}} & \dfrac{q_2 k_3^T}{\sqrt{d_2}} \\ \dfrac{q_3 k_1^T}{\sqrt{d_3}} & \dfrac{q_3 k_2^T}{\sqrt{d_3}} & \dfrac{q_3 k_3^T}{\sqrt{d_3}} \end{bmatrix}$$

$$softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ M\right) = \begin{bmatrix} softmax\left(\begin{bmatrix} \dfrac{\boldsymbol{q_1 k_1^T}}{\sqrt{d_1}} & -\infty & -\infty \end{bmatrix}\right) \\ softmax\left(\begin{bmatrix} \dfrac{\boldsymbol{q_2 k_1^T}}{\sqrt{d_2}} & \dfrac{\boldsymbol{q_2 k_2^T}}{\sqrt{d_2}} & -\infty \end{bmatrix}\right) \\ softmax\left(\begin{bmatrix} \dfrac{\boldsymbol{q_3 k_1^T}}{\sqrt{d_3}} & \dfrac{\boldsymbol{q_3 k_2^T}}{\sqrt{d_3}} & \dfrac{\boldsymbol{q_3 k_3^T}}{\sqrt{d_3}} \end{bmatrix}\right) \end{bmatrix} = \begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

$$head_i = softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ M\right)(VW_i^V)$$

$$softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \circ M\right) = \begin{bmatrix} softmax\left(\left[\frac{q_1k_1^T}{\sqrt{d_1}} \quad -\infty \quad -\infty\right]\right) \\ softmax\left(\left[\frac{q_2k_1^T}{\sqrt{d_2}} \quad \frac{q_2k_2^T}{\sqrt{d_2}} \quad -\infty\right]\right) \\ softmax\left(\left[\frac{q_3k_1^T}{\sqrt{d_3}} \quad \frac{q_3k_2^T}{\sqrt{d_3}} \quad \frac{q_3k_3^T}{\sqrt{d_3}}\right]\right) \end{bmatrix} = \begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

Raw Data     SE+PE

$PE(x_1)$    $MHA(x_1)$

| Transformatoren |
| sind |
| großartig |

| $x_1$ |
| $x_2$ |
| $x_3$ |

$$\begin{bmatrix} a_{11}(x_1w_1^v) + a_{12}(x_2w_1^v) + a_{13}(x_3w_1^v) & a_{11}(x_1w_2^v) + a_{12}(x_2w_2^v) + a_{13}(x_3w_3^v) \\ a_{21}(x_1w_1^v) + a_{22}(x_2w_1^v) + a_{23}(x_3w_1^v) & a_{21}(x_1w_2^v) + a_{22}(x_2w_2^v) + a_{23}(x_3w_2^v) \\ a_{31}(x_1w_1^v) + a_{32}(x_2w_1^v) + a_{33}(x_3w_1^v) & a_{31}(x_1w_2^v) + a_{32}(x_2w_2^v) + a_{33}(x_3w_2^v) \end{bmatrix}$$

Pass this information down the layers and complete the prediction
Use this predicted word as input for the next time instant

# Inference: at $t = 2$

$Q \in R^{3 \times 2}$

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

$[w_1^q \ w_2^q]$

| $[x_1 w_1^q \quad x_1 w_2^q] = q_1$ |
|---|
| $[x_2 w_1^q \quad x_2 w_2^q] = q_2$ |
| $[x_3 w_1^q \quad x_3 w_2^q] = q_3$ |

Cached
New Data

$K \in R^{3 \times 2}$

$W_i^K \in R^{2 \times 2}$

$KW_i^K = K \in R^{2 \times 2 \times 1}$

| Transformatoren |
|---|
| sind |
| großartig |

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

$[w_1^k \ w_2^k]$

| $[x_1 w_1^k \quad x_1 w_2^k] = k_1$ |
|---|
| $[x_2 w_1^k \quad x_2 w_2^k] = k_2$ |
| $[x_3 w_1^k \quad x_3 w_2^k] = k_3$ |

Raw Data

SE+PE

$V \in R^{3 \times 2}$

$W_i^V \in R^{2 \times 2}$

$VW_i^V = V \in R^{2 \times 2 \times 1}$

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

$[w_1^v \ w_2^v]$

| $[x_1 w_1^v \quad x_1 w_2^v] = v_1$ |
|---|
| $[x_2 w_1^v \quad x_2 w_2^v] = v_2$ |
| $[x_3 w_1^v \quad x_3 w_2^v] = v_3$ |

Now the decoder sees a sequence of length 2

$$head_i = softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ M\right)(VW_i^V)$$

$$softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \circ M\right)\begin{bmatrix} x_1w_1^v & x_1w_2^v \\ x_2w_1^v & x_2w_2^v \\ x_3w_1^v & x_3w_2^v \end{bmatrix} = \begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}\begin{bmatrix} x_1w_1^v & x_1w_2^v \\ x_2w_1^v & x_2w_2^v \\ x_3w_1^v & x_3w_2^v \end{bmatrix}$$

Cached
New Data

Raw Data

| Transformatoren |
| sind |
| großartig |

$PE(x_{1:2})$

SE+PE

| $x_1$ |
| $x_2$ |
| $x_3$ |

$MHA(x_{1:2})$

$$\begin{bmatrix} a_{11}(x_1w_1^v) + a_{12}(x_2w_1^v) + a_{13}(x_3w_1^v) & a_{11}(x_1w_2^v) + a_{12}(x_2w_2^v) + a_{13}(x_3w_3^v) \\ a_{21}(x_1w_1^v) + a_{22}(x_2w_1^v) + a_{23}(x_3w_1^v) & a_{21}(x_1w_2^v) + a_{22}(x_2w_2^v) + a_{23}(x_3w_2^v) \\ a_{31}(x_1w_1^v) + a_{32}(x_2w_1^v) + a_{33}(x_3w_1^v) & a_{31}(x_1w_2^v) + a_{32}(x_2w_2^v) + a_{33}(x_3w_2^v) \end{bmatrix}$$

Pass this information down the layers and complete the prediction
Use this predicted word as input for the next time instant

# Why KV Caching

$$Q \in R^{3 \times 2}$$

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

$[w_1^q \quad w_2^q]$

| $[x_1 w_1^q \quad x_1 w_2^q] = q_1$ |
|---|
| $[x_2 w_1^q \quad x_2 w_2^q] = q_2$ |
| $[x_3 w_1^q \quad x_3 w_2^q] = q_3$ |

$$X \in R^{3 \times 2} \qquad K \in R^{3 \times 2} \qquad W_i^K \in R^{2 \times 2} \qquad KW_i^K = K \in R^{3 \times 2 \times 1}$$

| Transformatoren |
|---|
| sind |
| großartig |

Raw Data

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

SE+PE

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

$[w_1^k \quad w_2^k]$

| $[x_1 w_1^k \quad x_1 w_2^k] = k_1$ |
|---|
| $[x_2 w_1^k \quad x_2 w_2^k] = k_2$ |
| $[x_3 w_1^k \quad x_3 w_2^k] = k_3$ |

$$V \in R^{3 \times 2} \qquad W_i^V \in R^{2 \times 2} \qquad VW_i^V = V \in R^{3 \times 2 \times 1}$$

| $x_1$ |
|---|
| $x_2$ |
| $x_3$ |

$[w_1^v \quad w_2^v]$

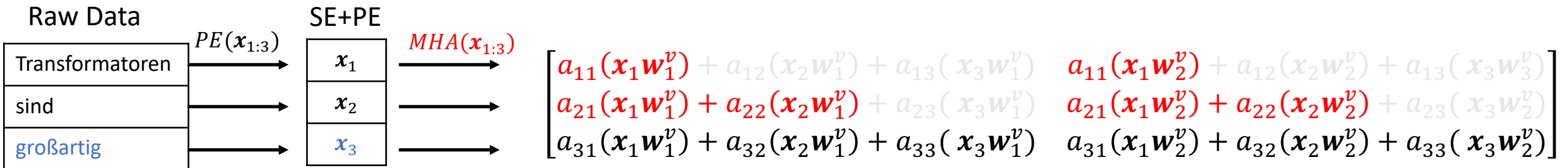| $[x_1 w_1^v \quad x_1 w_2^v] = v_1$ |
|---|
| $[x_2 w_1^v \quad x_2 w_2^v] = v_2$ |
| $[x_3 w_1^v \quad x_3 w_2^v] = v_3$ |

$$head_i = softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_K}} \circ M\right)(VW_i^V)$$

$$softmax\left(\frac{(QW_i^Q)(KW_i^K)^T}{\sqrt{d_k}} \circ M\right)\begin{bmatrix} \boldsymbol{x}_1\boldsymbol{w}_1^v & \boldsymbol{x}_1\boldsymbol{w}_2^v \\ \boldsymbol{x}_2\boldsymbol{w}_1^v & \boldsymbol{x}_2\boldsymbol{w}_2^v \\ \boldsymbol{x}_3\boldsymbol{w}_1^v & \boldsymbol{x}_3\boldsymbol{w}_2^v \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}\begin{bmatrix} \boldsymbol{x}_1\boldsymbol{w}_1^v & \boldsymbol{x}_1\boldsymbol{w}_2^v \\ \boldsymbol{x}_2\boldsymbol{w}_1^v & \boldsymbol{x}_2\boldsymbol{w}_2^v \\ \boldsymbol{x}_3\boldsymbol{w}_1^v & \boldsymbol{x}_3\boldsymbol{w}_2^v \end{bmatrix}$$

Cached

New Data

Raw Data    SE+PE

| Transformatoren |
| sind |
| großartig |

$PE(\boldsymbol{x}_{1:3})$

| $\boldsymbol{x}_1$ |
| $\boldsymbol{x}_2$ |
| $\boldsymbol{x}_3$ |

$MHA(\boldsymbol{x}_{1:3})$

$$\begin{bmatrix} a_{11}(\boldsymbol{x}_1\boldsymbol{w}_1^v) + a_{12}(\boldsymbol{x}_2\boldsymbol{w}_1^v) + a_{13}(\boldsymbol{x}_3\boldsymbol{w}_1^v) & a_{11}(\boldsymbol{x}_1\boldsymbol{w}_2^v) + a_{12}(\boldsymbol{x}_2\boldsymbol{w}_2^v) + a_{13}(\boldsymbol{x}_3\boldsymbol{w}_3^v) \\ a_{21}(\boldsymbol{x}_1\boldsymbol{w}_1^v) + a_{22}(\boldsymbol{x}_2\boldsymbol{w}_1^v) + a_{23}(\boldsymbol{x}_3\boldsymbol{w}_1^v) & a_{21}(\boldsymbol{x}_1\boldsymbol{w}_2^v) + a_{22}(\boldsymbol{x}_2\boldsymbol{w}_2^v) + a_{23}(\boldsymbol{x}_3\boldsymbol{w}_2^v) \\ a_{31}(\boldsymbol{x}_1\boldsymbol{w}_1^v) + a_{32}(\boldsymbol{x}_2\boldsymbol{w}_1^v) + a_{33}(\boldsymbol{x}_3\boldsymbol{w}_1^v) & a_{31}(\boldsymbol{x}_1\boldsymbol{w}_2^v) + a_{32}(\boldsymbol{x}_2\boldsymbol{w}_2^v) + a_{33}(\boldsymbol{x}_3\boldsymbol{w}_2^v) \end{bmatrix}$$
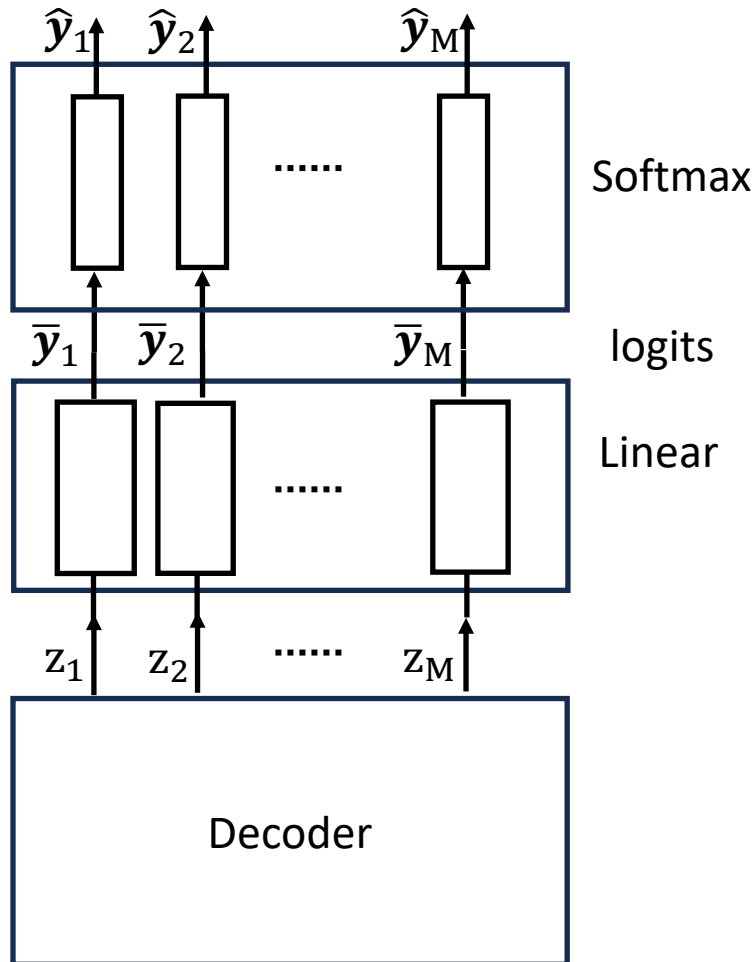
More details here:

https://huggingface.co/blog/not-lain/kv-caching

# Output Layer

- Let $Z \in R^{M \times d}$ be the output of decoder, where $M$ is the output sequence length and $d$ is the embedding dimension
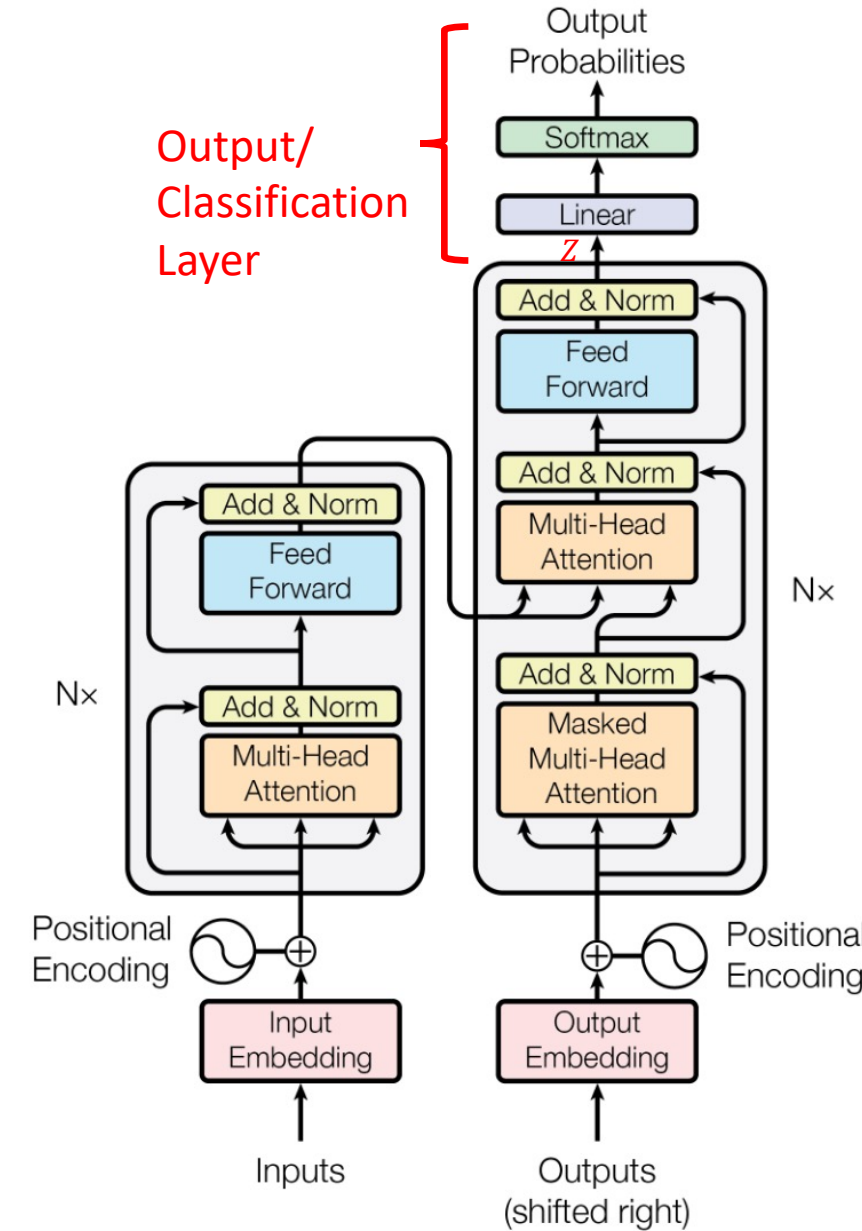  - $Z = [z_1 \quad z_2 \quad \dots \quad z_M]$, where $z \in R^{d \times 1}$



Softmax

logits

Linear

$$\widehat{\boldsymbol{y}}_i = \mathrm{softmax}(\overline{\boldsymbol{y}}_i)$$

$$\widehat{\boldsymbol{y}}_i \in R^{|V| \times 1}$$

$$\overline{\boldsymbol{y}}_i \in R^{|V| \times 1}$$

$$\overline{\boldsymbol{y}}_i = \left(W^O\right)^T z_i$$

$$W^O \in R^{d \times |V|}$$



Output/ Classification Layer

11

# Output Layer

- Output of the $softmax$ is the conditional probability of the next token

- Encoder-Decoder Architecture

$$P(\hat{y}_k | \boldsymbol{x}_{1:N}, \hat{y}_{1:k-1})$$

- Decoder-only Architecture (GPT)

$$P(\hat{y}_k | prompt\_tokens, \hat{y}_{1:k-1})$$

# Sampling from Probability Distributions

- If distribution of data is known, samples (data) can be generated

- Example:

$$\begin{bmatrix} P(X = \text{“Chains”}|Y = \text{“Markov”}) \\ P(X = \text{“Markov”}|Y = \text{“Markov”}) \\ P(X = \text{“Rocks”}|Y = \text{“Markov”}) \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \\ 0.5 \end{bmatrix}$$

$Y \longrightarrow X$

**Most Probable Word**:
$Chains$ and $Rocks$ are equally probable

```
>>> import numpy as np
>>> vocabulary = ["Chains","Markov","Rocks"]
>>> np.random.choice(vocabulary, p=[0.5,0,0.5])
['Chains'
```

| Chains | Rocks |
|:---:|:---:|

0        0.5        1

0.5      0.5

```
>>> import random
>>> random.uniform(0,1)
0.2894913983065621
```

**Generated Sample**: $Chains$

# Greedy Sampling

- Greedy sampling selects the token with the <span style="color:red">highest probability</span> at each step during generation.

- Pros
  - Simple and fast
  - Deterministic, i.e., always produces the same output for a given input

- Cons
  - Can lead to **repetitive** or **generic** outputs
  - May miss better sequences due to lack of exploration
  - Lacks diversity and creativity

- Use Case
  - Best suited for tasks where **precision** and **consistency** are more important than diversity
  - Eg: factual completions or deterministic templates.

# Top-$K$ Sampling

- Restricts the model's output choices to the $K$ most probable tokens at each generation step

- Retains only the top $K$ **highest logit values**

- All other logits are set to $-\infty$ (get **zero probability** after $softmax$)

- Pro:
  - Balances diversity and quality
  - Avoids  low-probability nonsensical tokens
  - Creative than greedy sampling

- Cons:
  - Choice of $K$ is challenging ($K = 1$ is greedy and large $K$ may lead to nonsensical tokens)

# Top-$p$ sampling or Nucleus Sampling

- Dynamically chooses the smallest set of tokens whose cumulative probability exceeds a threshold $p$
  - Sorts token by the predicted probabilities
  - Accumulate probabilities until the sum exceeds $p$
  - Sample the token from this subset
- Pros:
  - Adaptive unlike $top - K$
  - Avoids both repetitive outputs like greedy sampling and nonsensical ones
- Cons:
  - Sensitive to selection of $p$
- Usecase:
  - Chatbots, story generation and creative writing

# Why Temperature in Sampling?



- Logits ($\overline{\boldsymbol{y}}_i$) are scaled by $\frac{1}{T}$

- $\frac{\overline{Y}}{T} = \left[\frac{\overline{\boldsymbol{y}}_1}{T}, \frac{\overline{\boldsymbol{y}}_2}{T}, \frac{\overline{\boldsymbol{y}}_3}{T}, \dots, \frac{\overline{\boldsymbol{y}}_M}{T}\right]$

$$\widehat{\boldsymbol{y}}_i = \text{softmax}\left(\frac{\overline{\boldsymbol{y}}_i}{T}\right)$$

$$\widehat{\boldsymbol{y}}_i = softmax\left(\frac{\overline{\boldsymbol{y}}_i}{T}\right) = \begin{bmatrix} \dfrac{\exp\left(\frac{\overline{\boldsymbol{y}}_{i,1}}{T}\right)}{\sum_i \exp\left(\frac{\overline{\boldsymbol{y}}_i}{T}\right)} \\ \vdots \\ \dfrac{\exp\left(\frac{\overline{\boldsymbol{y}}_{i,|V|}}{T}\right)}{\sum_i \exp\left(\frac{\overline{\boldsymbol{y}}_i}{T}\right)} \end{bmatrix}$$

$$\widehat{\boldsymbol{y}}_i \in R^{|V| \times 1}$$

$\overline{\boldsymbol{y}}_i \in R^{|V| \times 1}$     $\overline{\boldsymbol{y}}_i = \left(W^O\right)^T z_i$

# How temperature impacts sampling?

- Low temperature ($T \rightarrow 0$)
  - Ditribution becomes sharper (?)

- High temperature ($T \rightarrow \infty$)
  - Distribution becomes flatter

$$T = 0.5 \qquad \frac{\bar{\boldsymbol{y}}_i}{0.5} = \begin{bmatrix} 4 \\ 2 \\ 0 \end{bmatrix} \qquad softmax\left(\frac{\bar{\boldsymbol{y}}_i}{0.5}\right) = \begin{bmatrix} 0.868 \\ 0.117 \\ 0.016 \end{bmatrix}$$

Sharper
Deterministic

$$T = 1 \qquad \bar{\boldsymbol{y}}_i = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \qquad softmax(\bar{\boldsymbol{y}}_i) = \begin{bmatrix} 0.665 \\ 0.245 \\ 0.090 \end{bmatrix}$$

$$T = 2 \qquad \frac{\bar{\boldsymbol{y}}_i}{2} = \begin{bmatrix} 4 \\ 2 \\ 0 \end{bmatrix} \qquad softmax\left(\frac{\bar{\boldsymbol{y}}_i}{2}\right) = \begin{bmatrix} 0.506 \\ 0.307 \\ 0.186 \end{bmatrix}$$

$$T = \infty \qquad \frac{\bar{\boldsymbol{y}}_i}{\infty} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \qquad softmax\left(\frac{\bar{\boldsymbol{y}}_i}{2}\right) = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$$

Flat
Random generation

# Number of Parameters

- Embedding Layer: $|V| * d_{model}$.  (no bias term in embedding layers)
- Positional Encoding:
  - Fixed PE:  None
  - Learned PE: $L_{seq} * d_{model}$      (no bias terms)
- Per-Encoder:
  - $MHA$: $4 * d_{model} * d_{model} + 4 * d_{model}$. (Q, K, V, W^O) matrix
  - $FFN$: $2 * d_{model} * d_{FF} + d_{FF} + d_{model}$
  - $LN$: $2 * 2 * d_{model}$
- Per-Decoder:
  - $MHA$: $2 * (4 * d_{model} * d_{model} + 4 d_{model})$  (Q, K, V, W^O) matrix
  - $FFN$: $2 * d_{model} * d_{FF} + d_{FF} + d_{model}$
  - $LN$: $3 * 2 * d_{model}$. (no bias terms)
- Output Layer:
  - Projection: $d\_model * |V|$

# FLOPs per Layer

- Attention:
  - Q, K, V Projections
    - $QW^Q, KW^K, VW^V: 3 * 2 * L_{seq} * d_{model}^2$
  - Attention Scores ($A$)
    - $QW^Q(KW^K)^T: 2 * L_{seq}^2 * d_{model}$
    - $softmax: O(L_{seq}^2)$
  - Attention Values
    - $A(VW^V): 2 * L_{seq}^2 * d_{model}$
  - Output Projection
    - $(AVW^V)W_O: 2 * L_{seq} * d_{model}^2$
- FFN
  - $2 * L_{seq} * 2 * d_{model} * d_{FF} \approx 2 * L_{seq} * 2 * d_{model}^2$
- LN:
  - $O(L_{seq} * d_{model})$

Due to multiplications and additions in matrix multiplication

$Q, K, V \in R^{L_{seq} \times d_{model}}$

$W^Q, W^K, W^V \in R^{d_{model} \times d_{model}}$

$QW^Q, KW^K, VW^V \in R^{L_{seq} \times d_{model}}$

$QW^Q(KW^K)^T \in R^{L_{seq} \times L_{seq}}$

$A \in R^{L_{seq} \times L_{seq}}$: Attention Matrix

$W^{[1]} \in R^{d_{model} \times d_{FF}}$

$W^{[2]} \in R^{d_{FF} \times d_{model}}$

$d_{FF} = 4d_{model}$

# FLOPS per Layer

- Quadratic in both sequence length and embedding dimension

- For smaller sequences
  - FFN is the bottleneck

- For larger sequences
  - Attention layer is the bottleneck

- KV Caching makes attention cost linear in $L_{seq}$

# Vision Transformers

## An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale

Alexey Dosovitskiy[*,†], Lucas Beyer[*], Alexander Kolesnikov[*], Dirk Weissenborn[*],
Xiaohua Zhai[*], Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby[*,†]
[*]equal technical contribution, [†]equal advising
Google Research, Brain Team
{adosovitskiy, neilhoulsby}@google.com
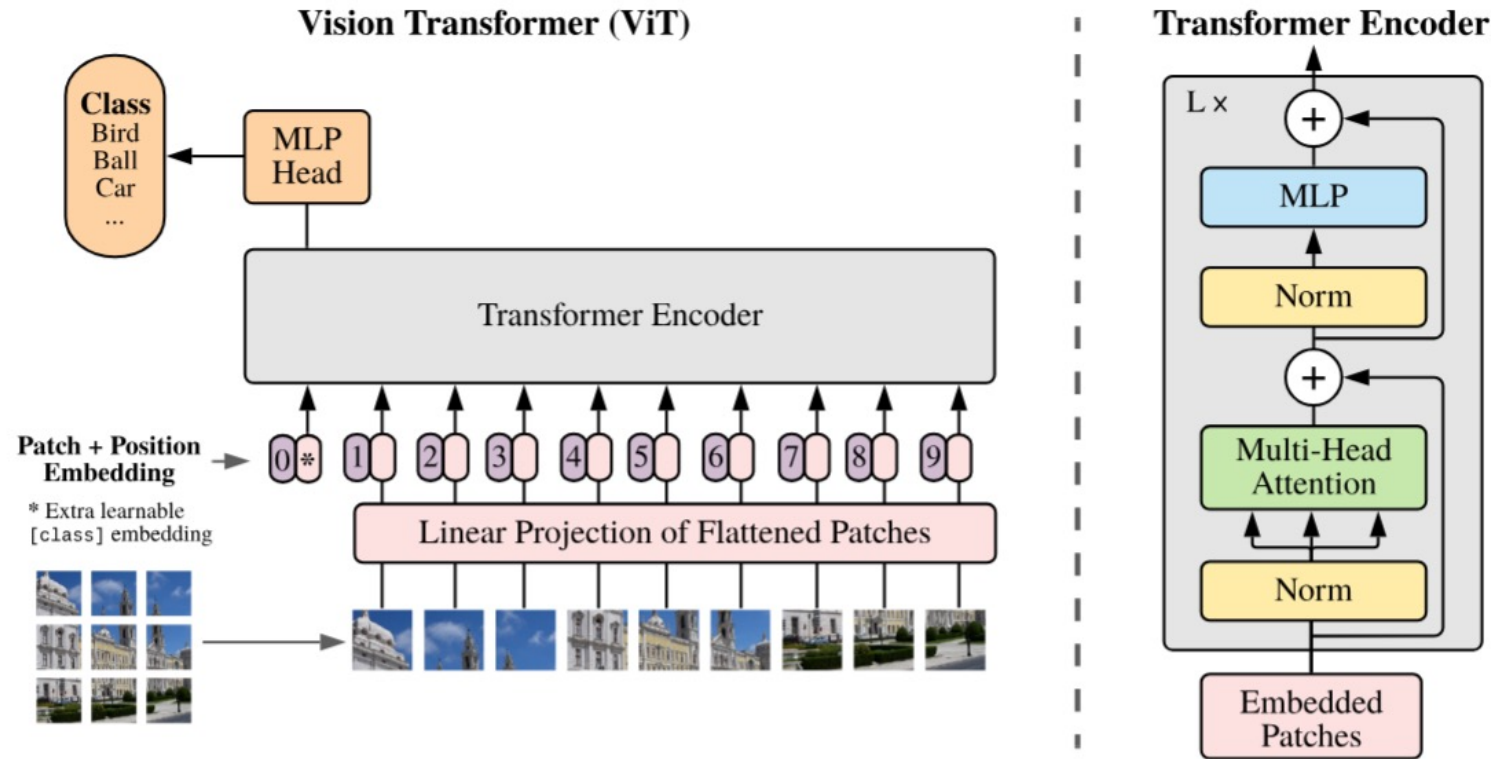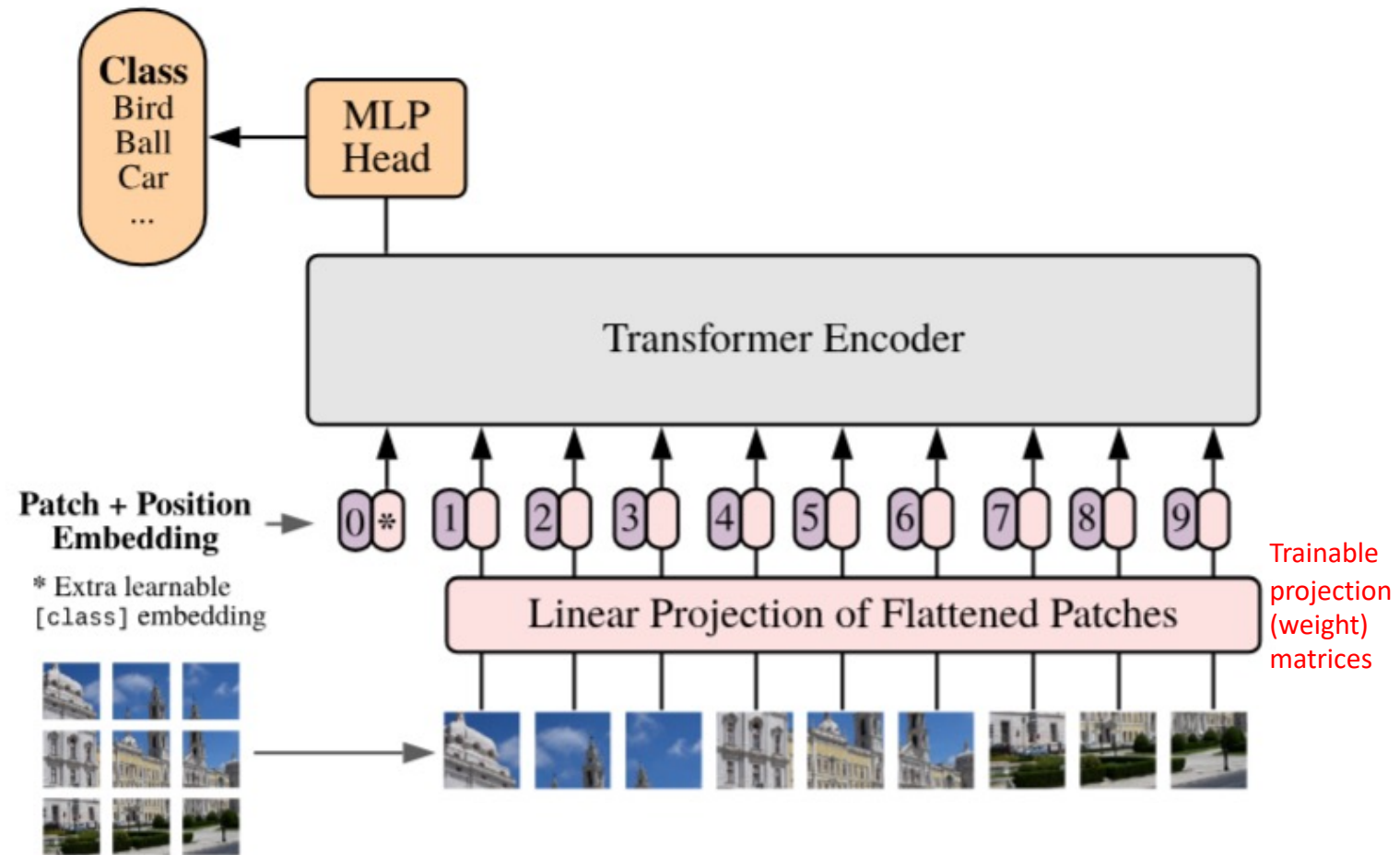
# Vision Transformers



Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable "classification token" to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

# Vision Transformer (ViT)



$x \in R^{H \times W \times 3}$

$x \in R^{N \times (3P^2)}$

$N = \frac{HW}{P^2}$ : Number of patches

$P$: Patch height (width)

Output embedding vector corresponding to * token represents the whole image (similar to feature vector)