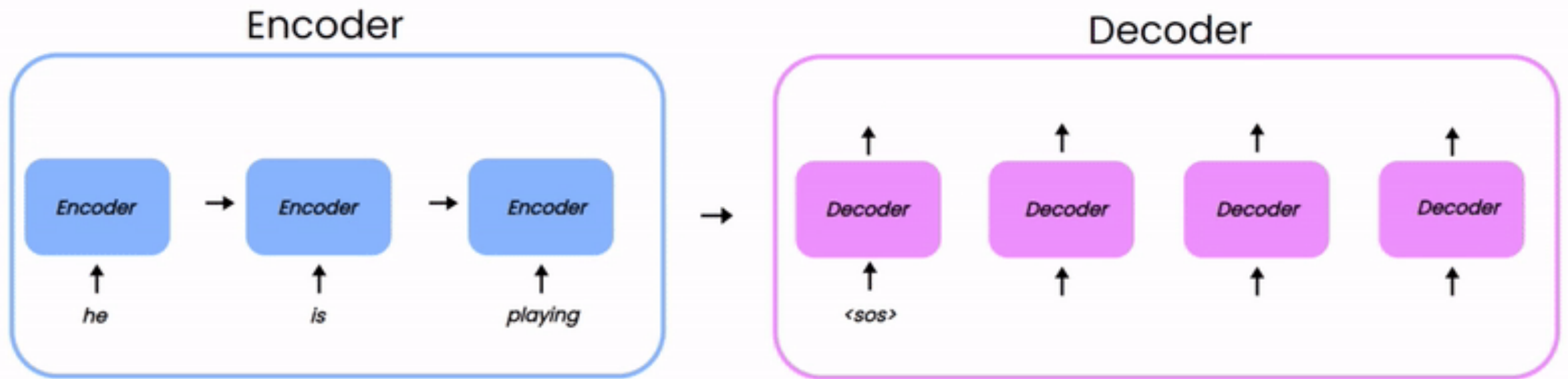# ENM 5310: Data-driven Modeling and Probabilistic Scientific Computing
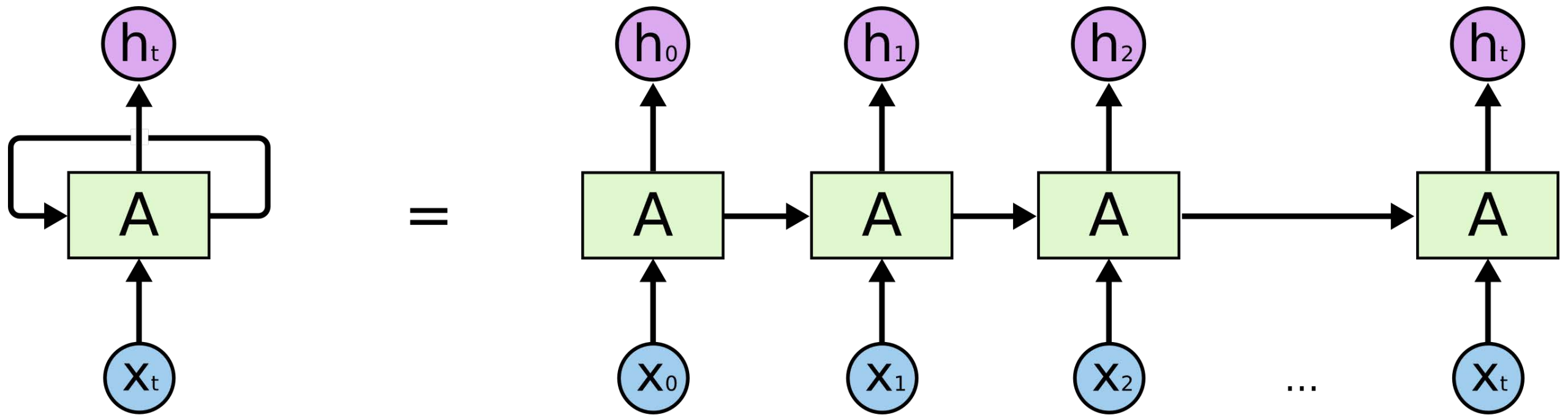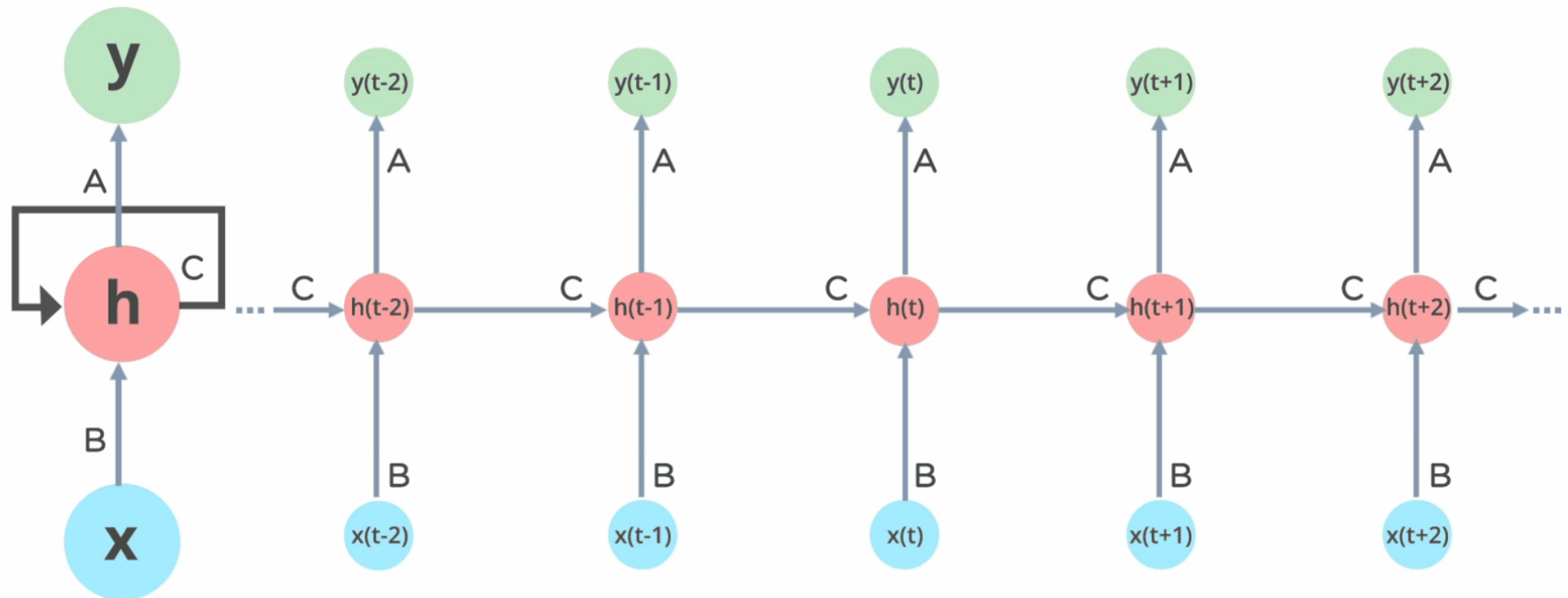
## Lecture #15: Transformers

# Modeling of sequence data

# RNNs



An unrolled recurrent neural network.

# RNN limitations

- Sequential processing is slow (can't parallelize)
- Long-range dependency issues (vanishing gradients)
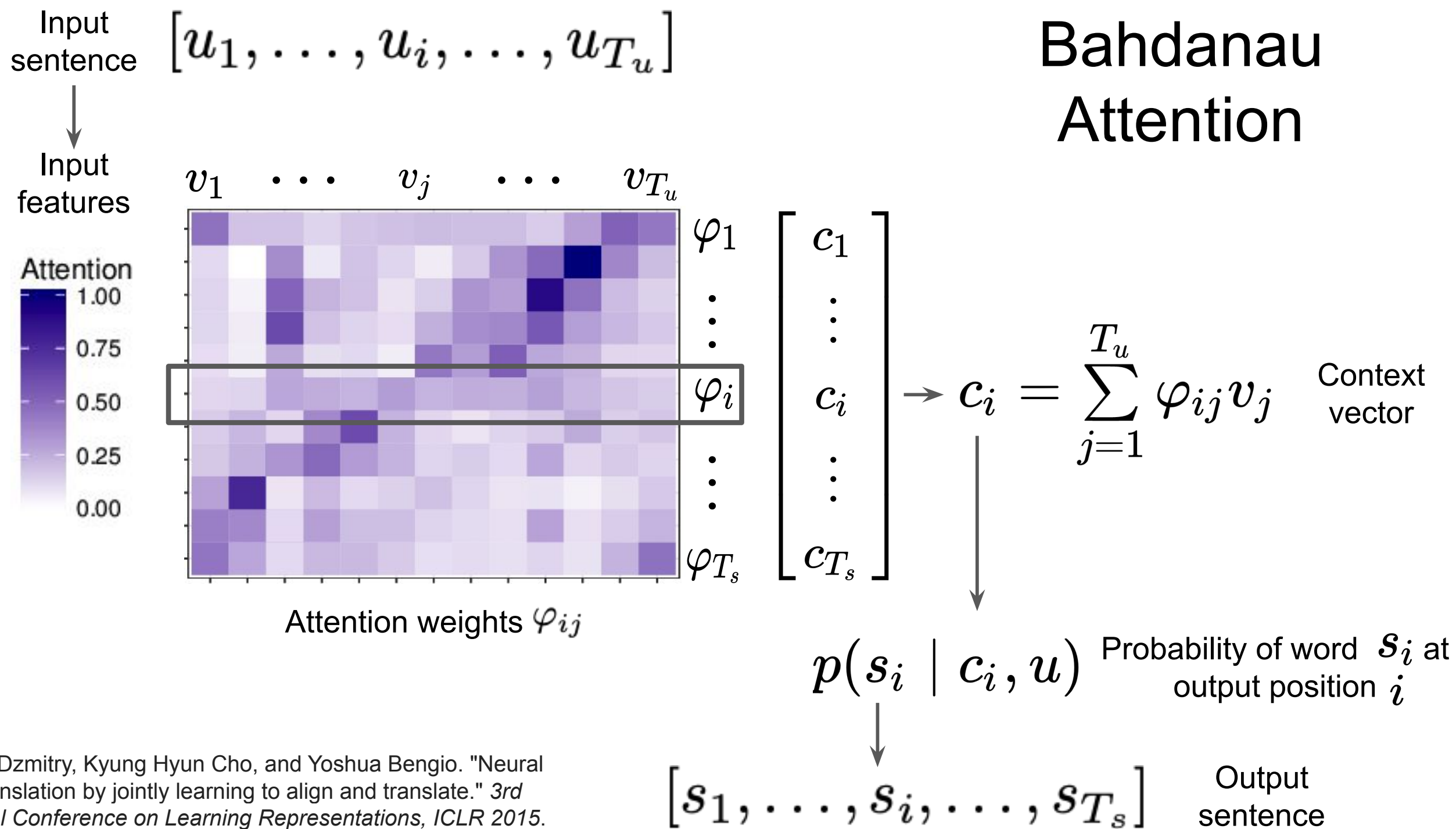- Limited context window in practice

   *"The cat is on the table"*

   *"The cat, who belongs to my mother, is on the table".*

  E.g.: when translating we need to remember "cat" is the subject even after processing many words.

- A need for better models:
  - Language translation requiring full sentence context
  - Document summarization
  - Question answering

  - Protein Structure Prediction (long amino acid sequences)
  - Molecular Property Prediction (multi-atom interactions)
  - Climate Science (long timescales)
  - Astronomical Data Analysis (vast spatial and time scales)
  -

# Key Innovation: Attention

Input sentence $[u_1, \ldots, u_i, \ldots, u_{T_u}]$

Bahdanau Attention

Input features

$v_1 \quad \cdots \quad v_j \quad \cdots \quad v_{T_u}$



Attention
1.00
0.75
0.50
0.25
0.00

$\varphi_1$
$\vdots$
$\varphi_i$
$\vdots$
$\varphi_{T_s}$

$\begin{bmatrix} c_1 \\ \vdots \\ c_i \\ \vdots \\ c_{T_s} \end{bmatrix}$

$$c_i = \sum_{j=1}^{T_u} \varphi_{ij} v_j$$

Context vector

Attention weights $\varphi_{ij}$

$p(s_i \mid c_i, u)$

Probability of word $s_i$ at output position $i$

$[s_1, \ldots, s_i, \ldots, s_{T_s}]$

Output sentence

Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *3rd International Conference on Learning Representations, ICLR 2015.* 2015.

# From Sequential to Parallel Processing

- Compare how humans read vs how RNNs process text
  - Humans: Quick glances at relevant parts
  - RNNs: Must process word by word
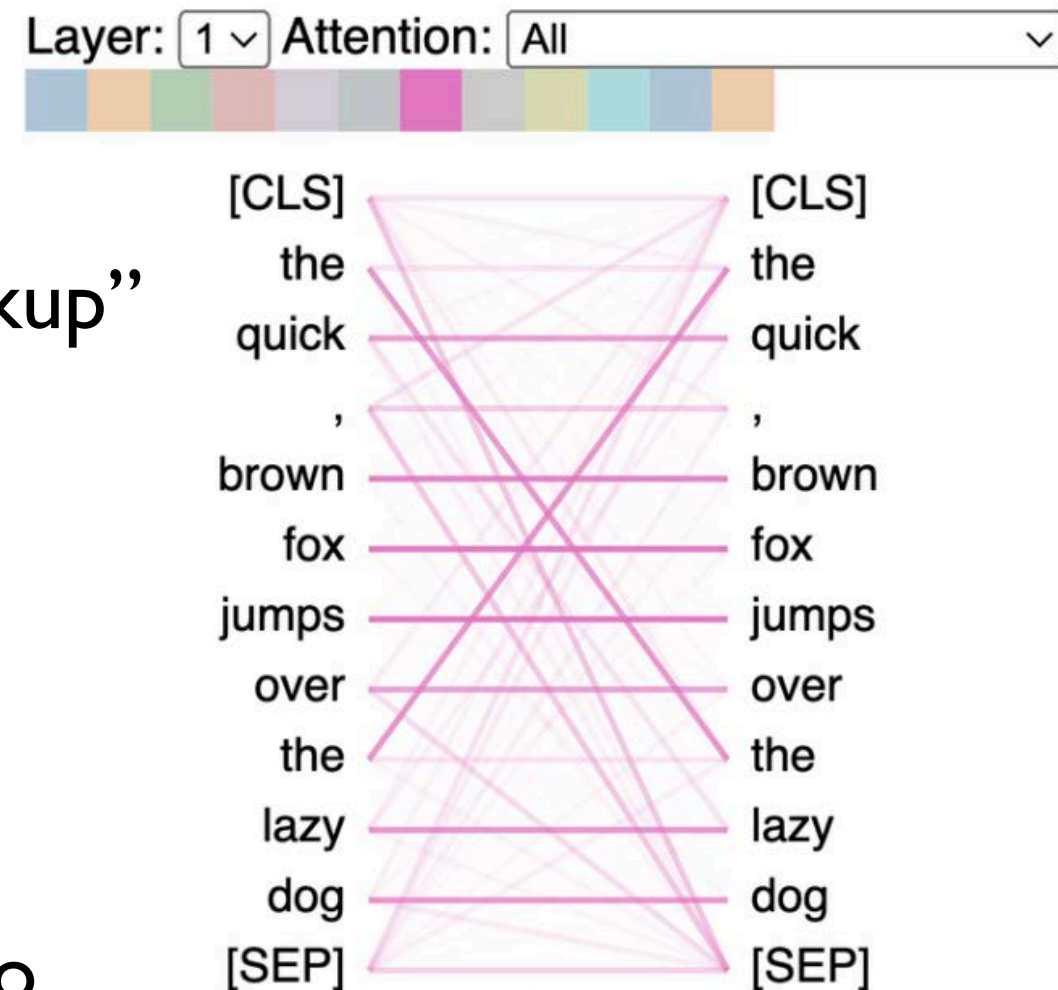- Attention is "looking" at all words simultaneously

Key innovation: Attention
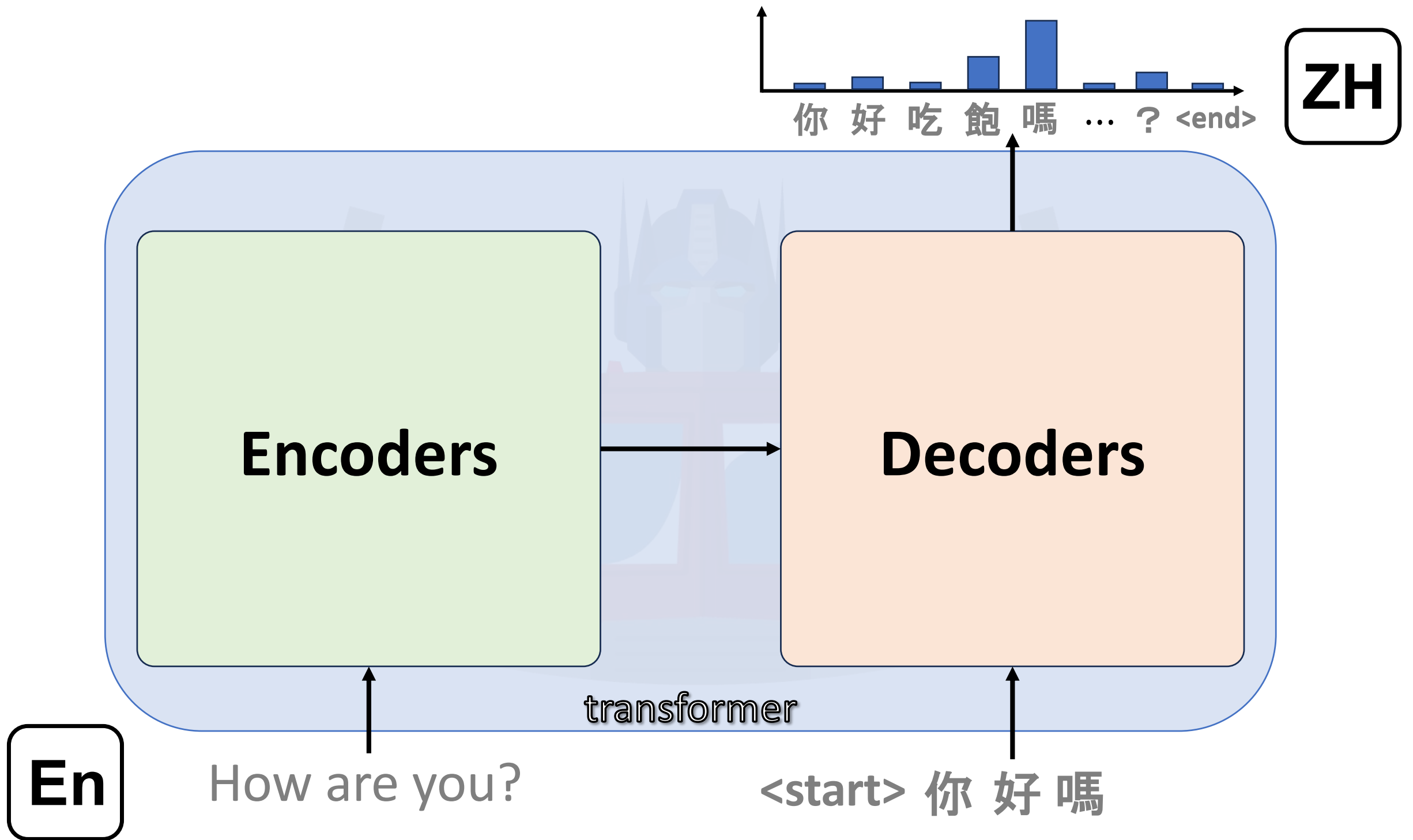
Think of attention as a "smart dictionary lookup"

- Keys: What you're looking for
- Values: What you get back
- Queries: What you're asking for

Simple example: Translation attention
- Each word in target language might need to "look at" multiple source words

Transformers

Encoder-Decoder Auto-regressive Models

# Key Building Blocks

- Tokenization

- Token embeddings

- Position embeddings

- Self-attention & Multi-head self-attention
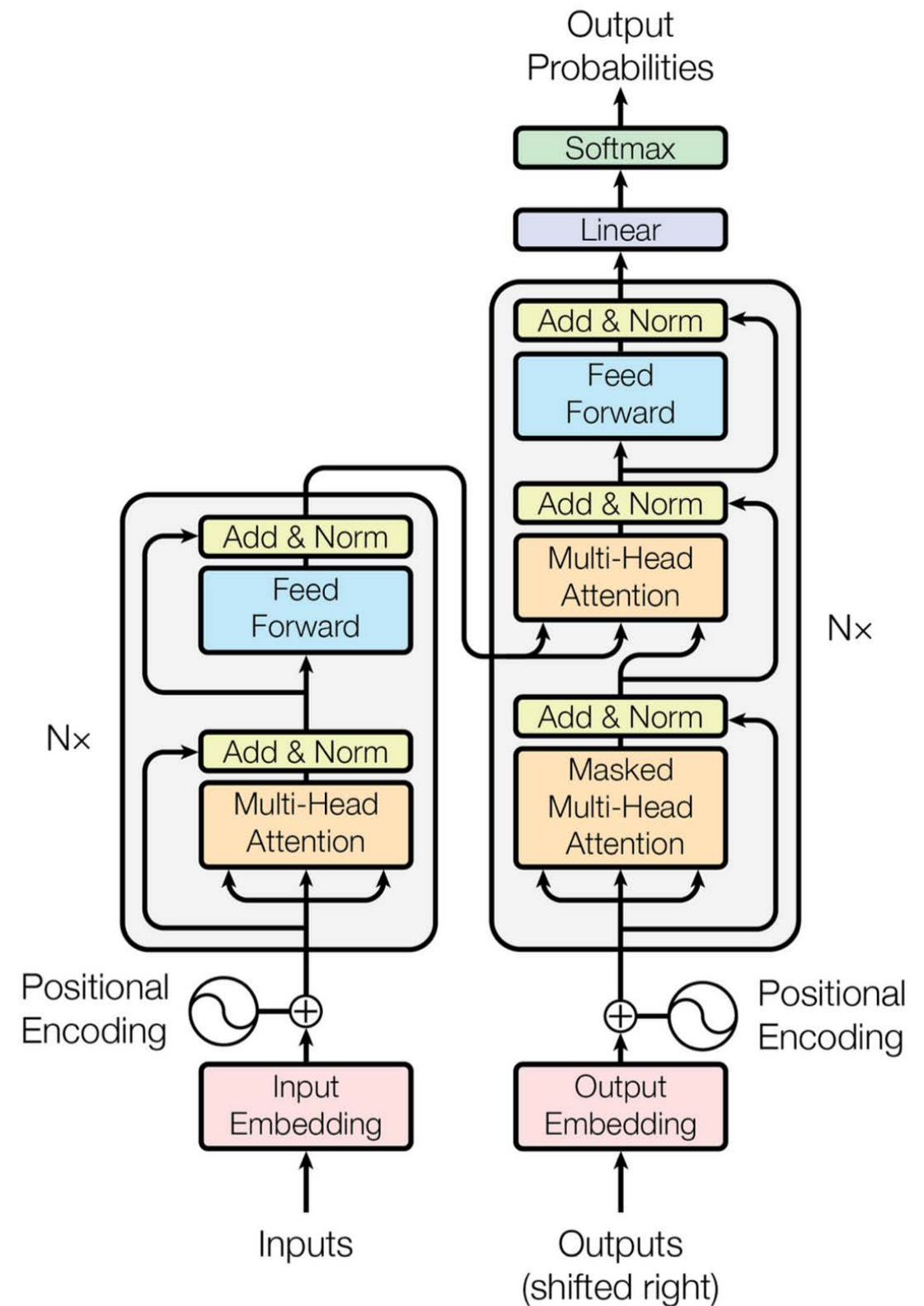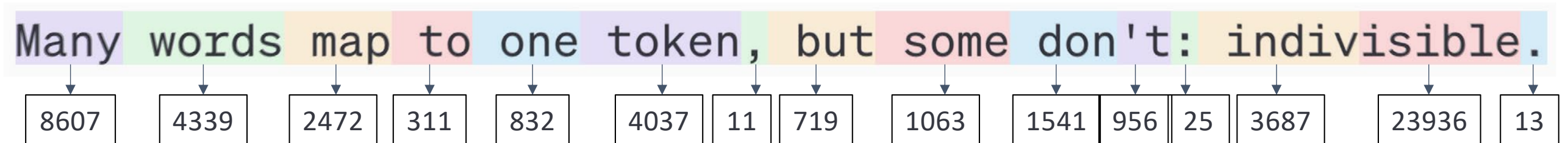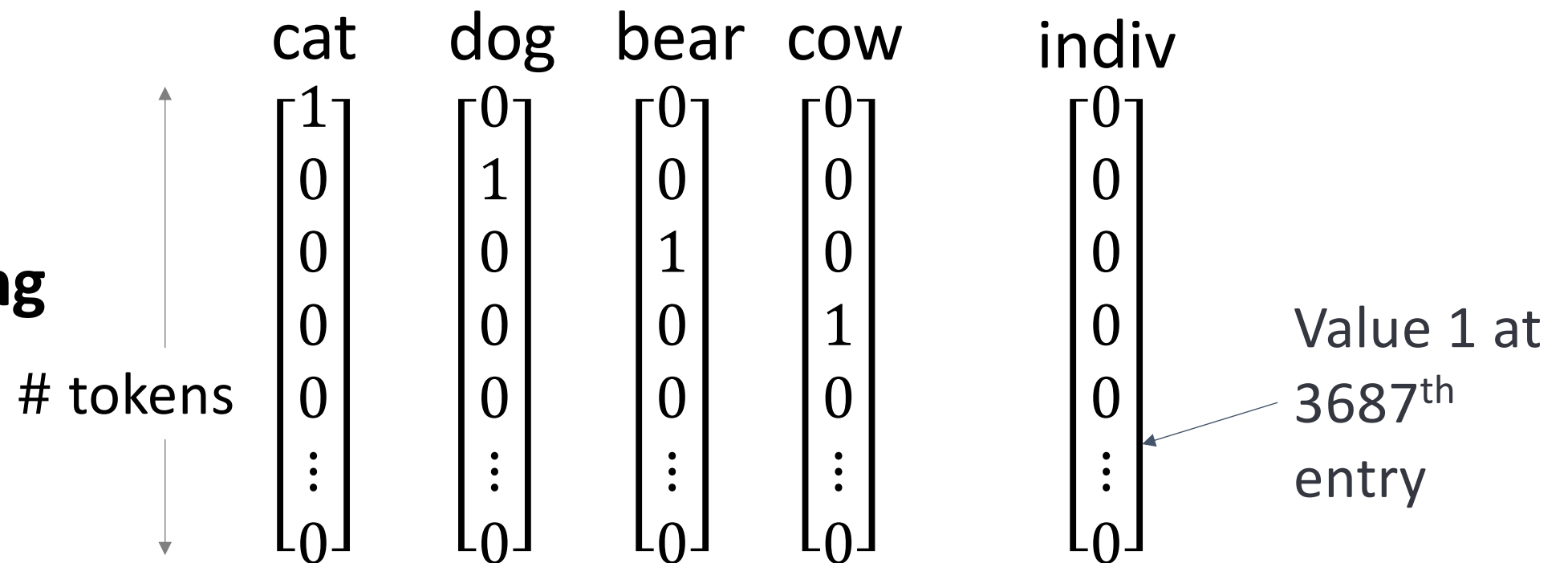
- MLPs

- Cross-attention



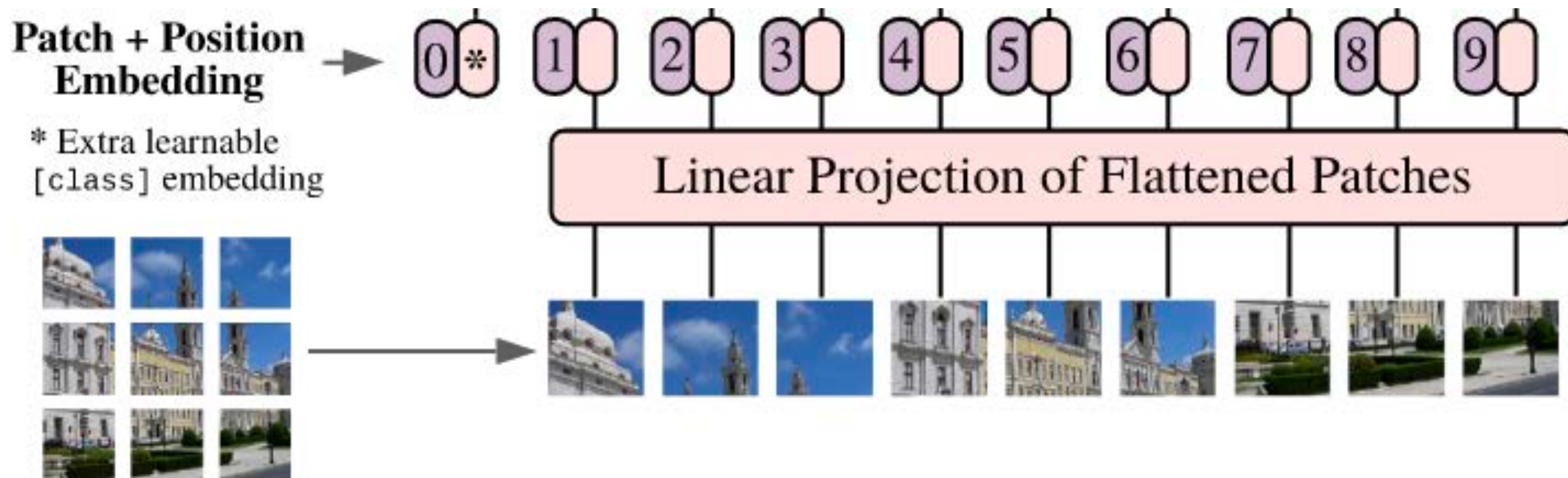Figure 1: The Transformer - model architecture.

Vaswani, A. (2017). Attention is all you need. *Advances in Neural Information Processing Systems.*

# Tokenization in NLP

Many words map to one token, but some don't: indivisible.

| 8607 | 4339 | 2472 | 311 | 832 | 4037 | 11 | 719 | 1063 | 1541 | 956 | 25 | 3687 | 23936 | 13 |

**One-hot encoding**

# tokens

$$
\text{cat} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad
\text{dog} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad
\text{bear} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad
\text{cow} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad
\text{indiv} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
$$

Value 1 at $3687^{th}$ entry

# Tokenization in Vision



Patch + Position Embedding

* Extra learnable [class] embedding

Linear Projection of Flattened Patches

# Tokenization in Point Clouds

# Tokenization in NLP

Many words map to one token, but some don't: indivisible.

8607  4339  2472  311  832  4037  11  719  1063  1541  956  25  3687  23936  13

**One-hot encoding**

# tokens

$$
\text{cat} \quad \text{dog} \quad \text{bear} \quad \text{cow} \quad \quad \text{indiv}
$$

$$
\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
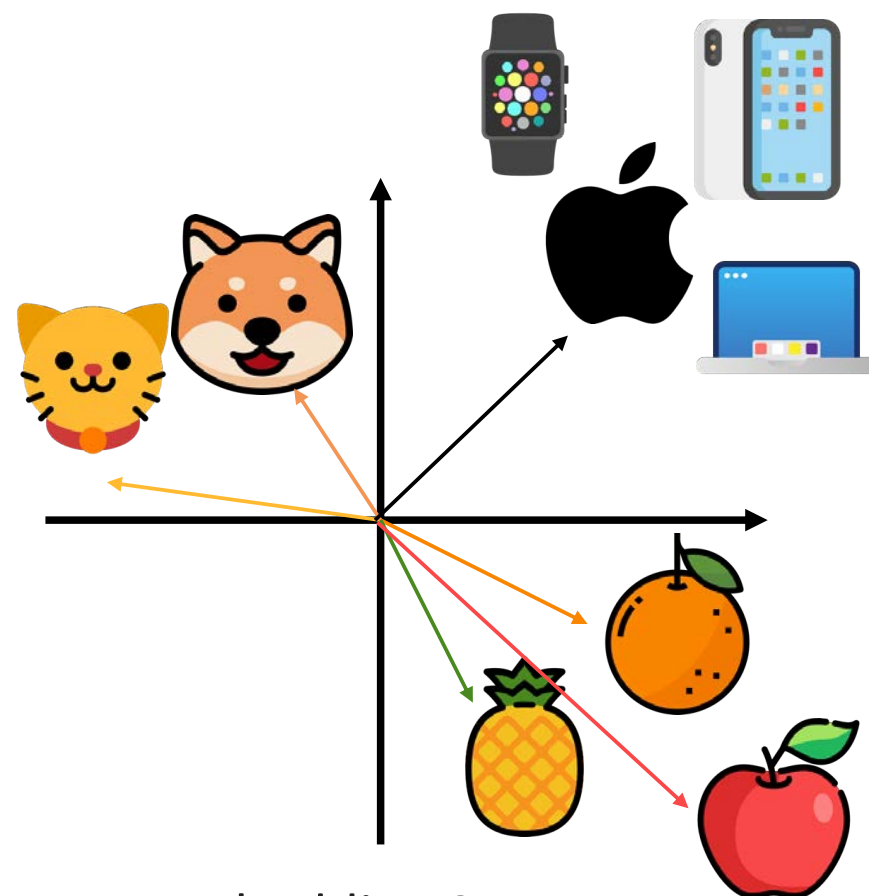$$

Value 1 at 3687[th] entry

# Token Embeddings

cat  dog  bear  cow   indiv

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Value 1 at 3687th entry

I bought an **apple** and an orange.

I bought an **apple** watch.

Apple

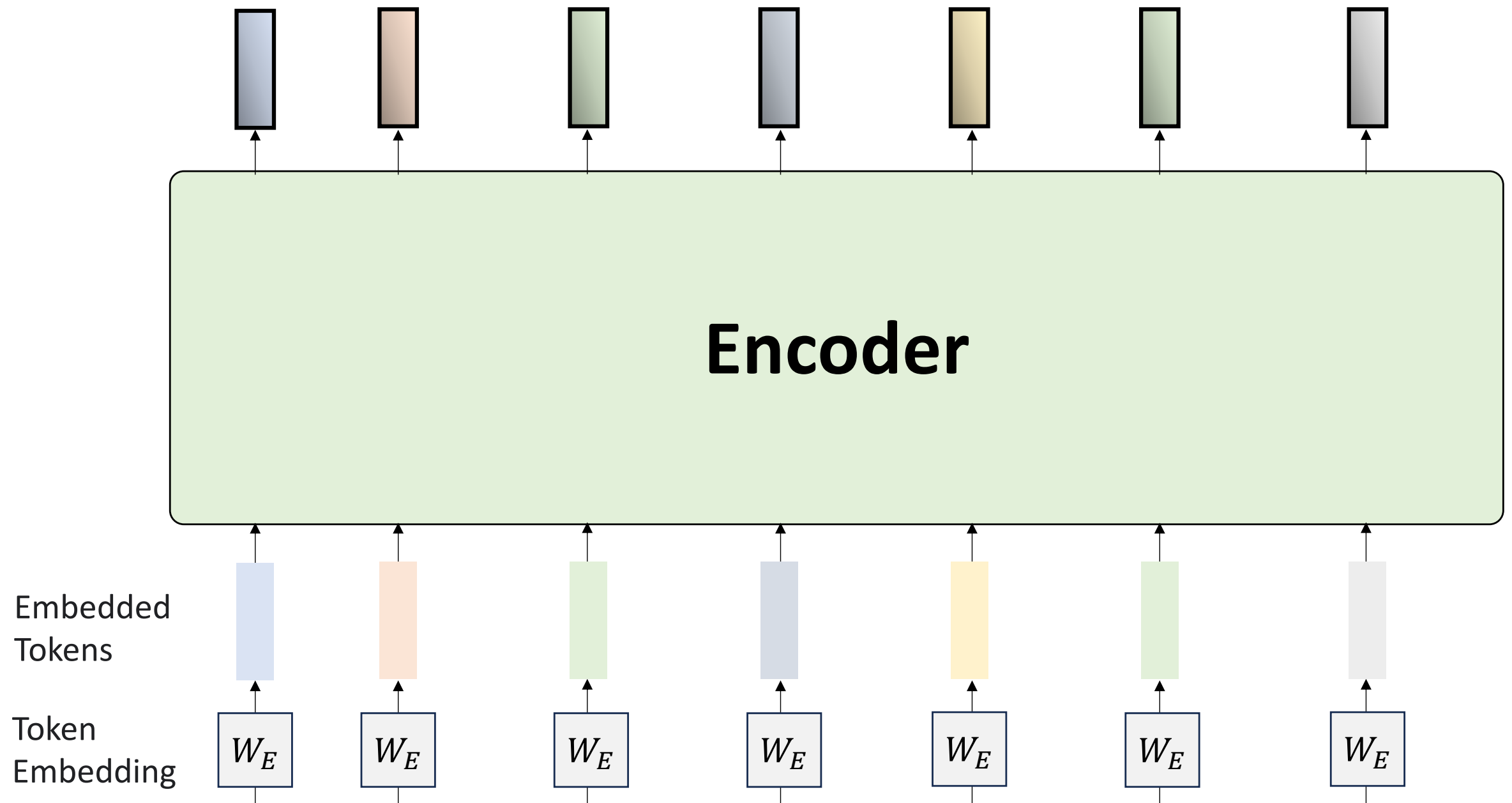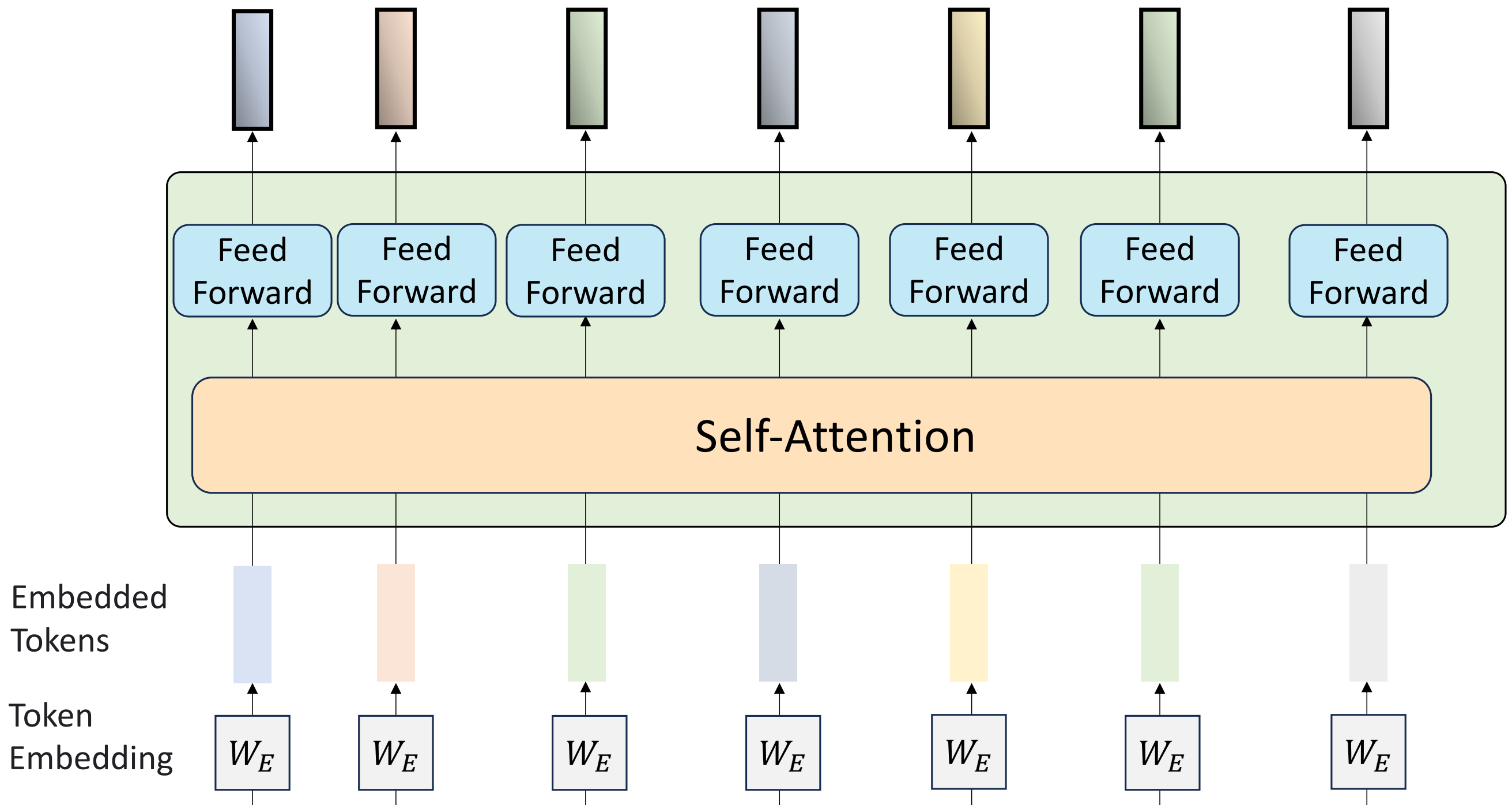$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & & \\ & & & \cdots & & \\ & & & & & \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# tokens

dog

Embedding Space       Embedded token       Embedding Matrix

*Slide by Jia-Bin Huang, University of Maryland College Park*

# Token Embeddings



cat, dog, bear, cow, indiv

$$\text{cat} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{dog} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{bear} \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{cow} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{indiv} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Value 1 at $3687^{th}$ entry

I bought an **apple** and an orange.

I bought an **apple** watch.

Apple

Embedding Space

$$d \begin{bmatrix} 0.5 \\ 2.7 \\ 1.2 \\ \vdots \\ 0.2 \end{bmatrix} = d \begin{bmatrix} & & & & & \\ & & & \cdots & & \\ & & & & & \end{bmatrix} \quad \text{dog} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

# tokens

Embedded token     Embedding Matrix

*Slide by Jia-Bin Huang, University of Maryland College Park*

# Encoder



Slide by Jia-Bin Huang, University of Maryland College Park

# Self-Attention

Updated feature  $\boldsymbol{x}'_4 \quad = \quad \alpha'_{4,1}\,\boldsymbol{x}_1 \quad + \quad \alpha'_{4,2}\,\boldsymbol{x}_2 \quad + \quad \alpha'_{4,3}\,\boldsymbol{x}_3 \quad + \quad \alpha'_{4,4}\,\boldsymbol{x}_4 \quad + \quad \alpha'_{4,5}\,\boldsymbol{x}_5$

Attention Scores  $\alpha'_{4,1}$  0.082  $\alpha'_{4,2}$  0.0495  $\alpha'_{4,3}$  0.0199  $\alpha'_{4,4}$  0.6034  $\alpha'_{4,5}$  0.2452

$$\text{Softmax}$$

$$\alpha'_{4,i} = \frac{\exp(\alpha_{4,i})}{\sum_j \exp(\alpha_{4,j})}$$

Token similarity  $\alpha_{4,1} = \boldsymbol{x}_4^\top \boldsymbol{x}_1 \quad \alpha_{4,2} = \boldsymbol{x}_4^\top \boldsymbol{x}_2 \quad \alpha_{4,3} = \boldsymbol{x}_4^\top \boldsymbol{x}_3 \quad \alpha_{4,4} = \boldsymbol{x}_4^\top \boldsymbol{x}_4 \quad \alpha_{4,5} = \boldsymbol{x}_4^\top \boldsymbol{x}_5$

Embedded Tokens  $\boldsymbol{x}_1 \in R^d \quad \boldsymbol{x}_2 \in R^d \quad \boldsymbol{x}_3 \in R^d \quad \boldsymbol{x}_4 \in R^d \quad \boldsymbol{x}_5 \in R^d$

# Self-Attention

Updated feature $\boldsymbol{x}'_4 = \boxed{W^O} ( \alpha'_{4,1}\, v_1 + \alpha'_{4,2}\, v_2 + \alpha'_{4,3}\, v_3 + \alpha'_{4,4}\, v_4 + \alpha'_{4,5}\, v_5 )$

$$= \sum_i \alpha'_{4,1} \left( \boxed{W^O} \boxed{W^V} \right) \boldsymbol{x}_i$$

$W^O \in R^{d \times d_v}$

$W^Q \in R^{d_k \times d}$

$W^K \in R^{d_k \times d}$

$W^V \in R^{d_v \times d}$



$\alpha'_{4,1} \quad \alpha'_{4,2} \quad \alpha'_{4,3} \quad \alpha'_{4,4} \quad \alpha'_{4,5}$

Softmax

$\alpha_{4,1} = k_1^\top q_4 \quad \alpha_{4,2} = k_2^\top q_4 \quad \alpha_{4,3} = k_3^\top q_4 \quad \alpha_{4,4} = k_4^\top q_4 \quad \alpha_{4,5} = k_5^\top q_4$

$k_1 \quad v_1 \qquad k_2 \quad v_2 \qquad k_3 \quad v_3 \quad q_4 \quad k_4 \quad v_4 \qquad k_5 \quad v_5$

$W^K \quad W^V \qquad W^K \quad W^V \qquad W^K \quad W^V \quad W^Q \quad W^K \quad W^V \qquad W^K \quad W^V$

Embedded Tokens

$\boldsymbol{x}_1 \qquad \boldsymbol{x}_2 \qquad \boldsymbol{x}_3 \qquad \boldsymbol{x}_4 \qquad \boldsymbol{x}_5$

# Self-Attention

**Single-head** attention

$$\text{Attention}(Q, K, V) = V \, \text{softmax}\left(\frac{K^\top Q}{\sqrt{d_k}}\right)$$

$$Q = \boxed{W^Q} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$$K = \boxed{W^K} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

$$V = \boxed{W^V} \; \boxed{x_1} \; \boxed{x_2} \; \boxed{x_3} \; \boxed{x_4} \; \boxed{x_5}$$

# MLPs

# Positional Encoding

| Position $k$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $2^3$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | ← Slow oscillating
| $2^2$ | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| $2^1$ | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| $2^0$ | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | ← Fast oscillating

Dimension

Positional embedding  $d$

Embedded Tokens  $d$   $\boldsymbol{x}_1$   $\boldsymbol{x}_2$   $\boldsymbol{x}_3$   $\boldsymbol{x}_4$   $\boldsymbol{x}_5$

# Positional Encoding

$$\mathbf{S}_i = \left[ \sin(\omega_1 i), \cos(\omega_1 i), \sin(\omega_2 i), \cos(\omega_2 i), \ldots, \sin(\omega_{e/2} i), \cos(\omega_{e/2} i) \right]$$

Position $k$

Angular frequency

$w_i = N^{-2i/d}$

$N = 100{,}000$

$$d \begin{bmatrix} \sin(w_0 k) \\ \cos(w_0 k) \\ \sin(w_1 k) \\ \cos(w_1 k) \\ \vdots \\ \vdots \\ \sin\left(w_{\frac{d}{2}-1} k\right) \\ \cos\left(w_{\frac{d}{2}-1} k\right) \end{bmatrix}$$

← Fast oscillating

← Slow oscillating

👍 **Normalized Range**

👍 **Unique identifier, unlimited length**

👍 **Relative positions as linear transform**

# Positional Encoding

# Transformer Encoder



Encoder #2

Encoder #1

Embedded Tokens

$d$

$x_1$ $P_1$  $x_2$ $P_2$  $x_3$ $P_3$  $x_4$ $P_4$  $x_5$ $P_5$

# Transformer Encoder



💡 **Residual connection**

💡 **Layer normalization**

$$\mathrm{LayerNorm}(\boldsymbol{x}) =$$

$$\gamma \left( \frac{\boldsymbol{x} - \mathrm{mean}(\boldsymbol{x})}{\sqrt{\mathrm{Variance}(\mathbf{x}) + \epsilon}} \right) + \beta$$

$\gamma, \beta \in R$
Learnable parameters

*Slide by Jia-Bin Huang, University of Maryland College Park*

# Decoders

# Decoders

# Vision Transformers



*Dosovitskiy, A. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.*

# Advantages & Limitations

- Parallel processing
- Better at capturing long-range dependencies
- No vanishing gradient problem
- Scales well with data and compute
- Modality-agnostic



- Self-attention cost and memory scales quadratically with seq. length
- Needs a lot of data
- Expensive and difficult to train

# Generative Pretrained Transformers
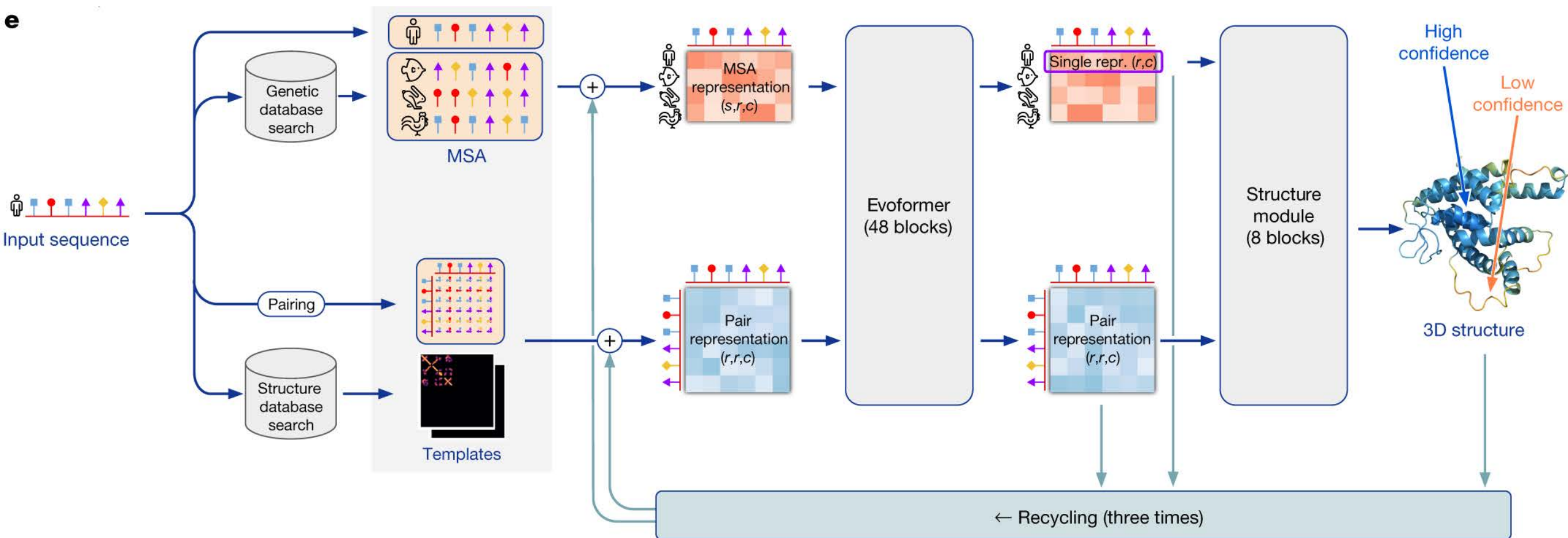
# Multi-modal Transformers
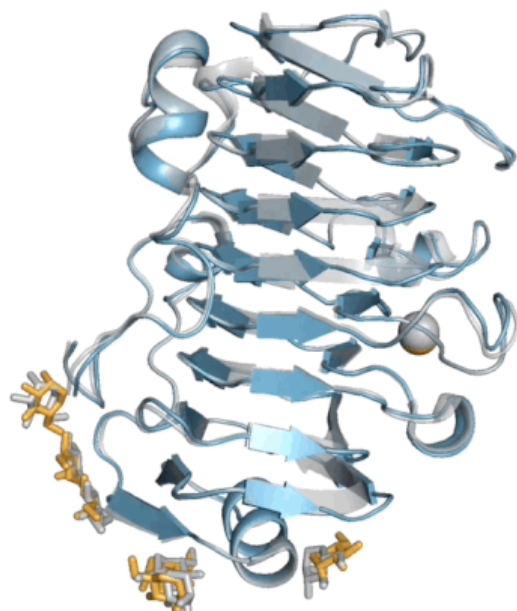
# Multi-modal Transformers



Driving scenes generated by Wayve's GAIA-1, a new generative AI model that creates realistic driving videos and offers fine-grained control over ego-vehicle behaviour and scene features.
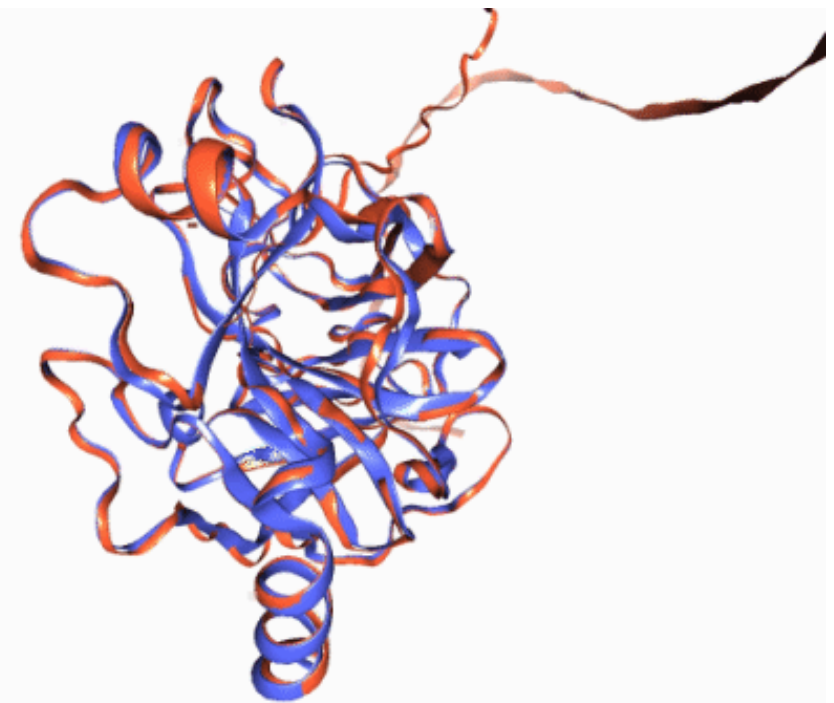
# AlphaFold 3



7BBV

Ground truth shown in grey

Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., ... & Jumper, J. M. (2024). Accurate structure prediction of biomolecular interactions with AlphaFold 3. Nature, 1-3.

# UvA Tutorial Notebooks

# Build your own GPT from scratch