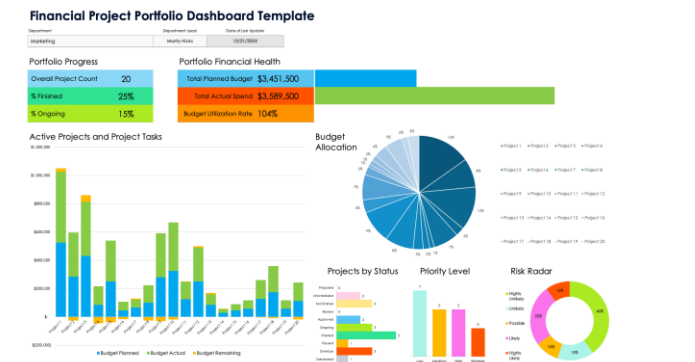
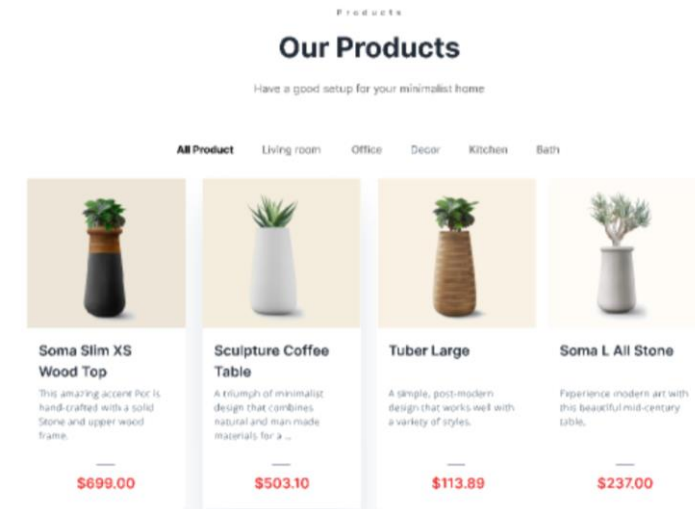
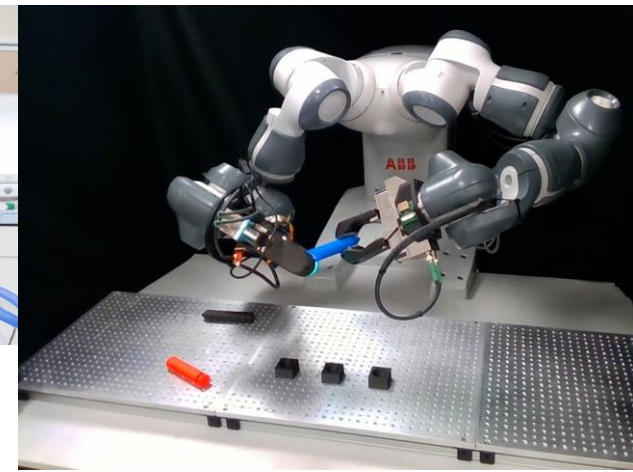


Offline RL

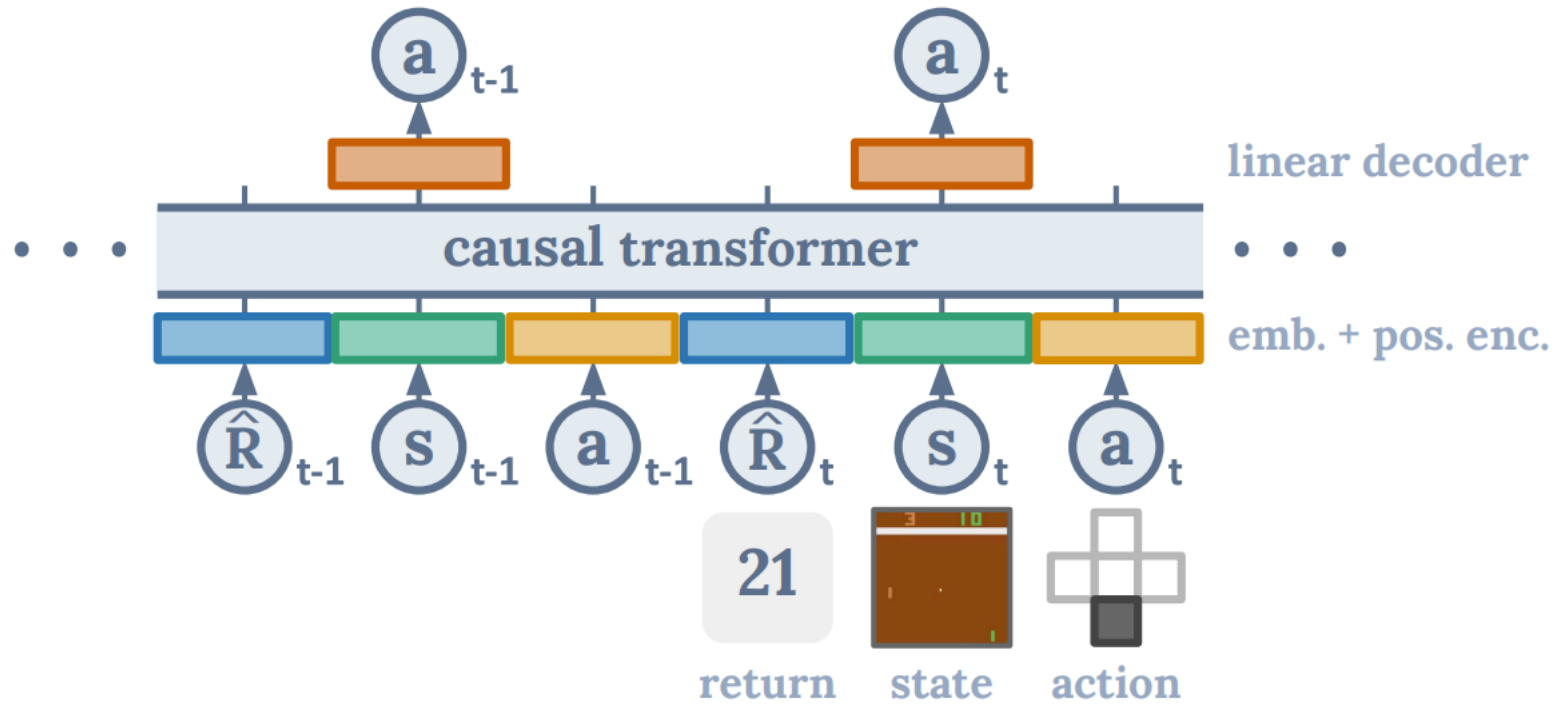
The problem of Offline RL

Applications

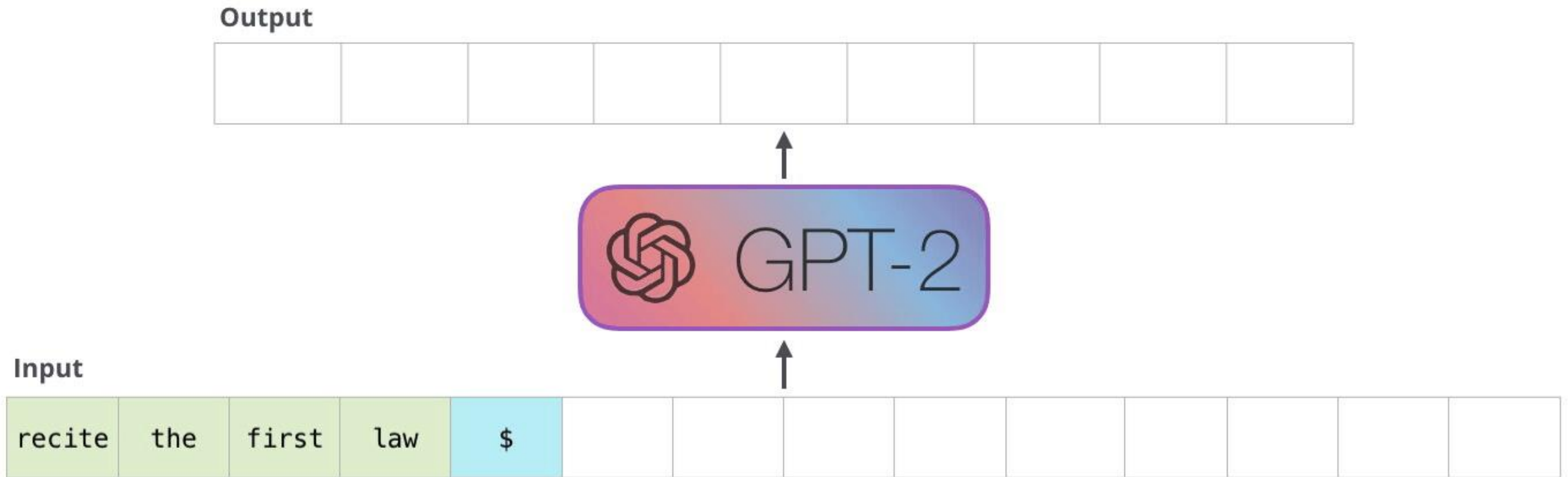
- Healthcare and Personalized Medicine
- Autonomous Driving
- Logistics and warehouse operations
- Robotics (Manipulation & Assistive Robots)
- Recommender Systems
- Smart Grids & Energy Systems
- Education / Intelligent Tutoring
- Finance & Trading



Decision Transformer



GPT-Style Transformers



The challenge of modeling sequences



A sequence may be arbitrarily long. We only have a finite number of parameters.

1. How do we model what comes next using only a finite number of parameters?
2. How do we build “expressive enough” models without cutting off the history at an arbitrary point (i.e. a Markov assumption)?
3. Can we build models that can encode items in a sequence in parallel?

The challenge of modeling sequences

A sequence may be arbitrarily long. We only have a finite number of parameters.

1. How do we model what comes next using only a finite number of parameters?

✓ Markov models ✓ Recurrent Neural Networks ✓ Transformers

2. How do we build “expressive enough” models without cutting off the history at an arbitrary point (i.e. a Markov assumption)?

X Markov models ✓ Recurrent Neural Networks ✓ Transformers

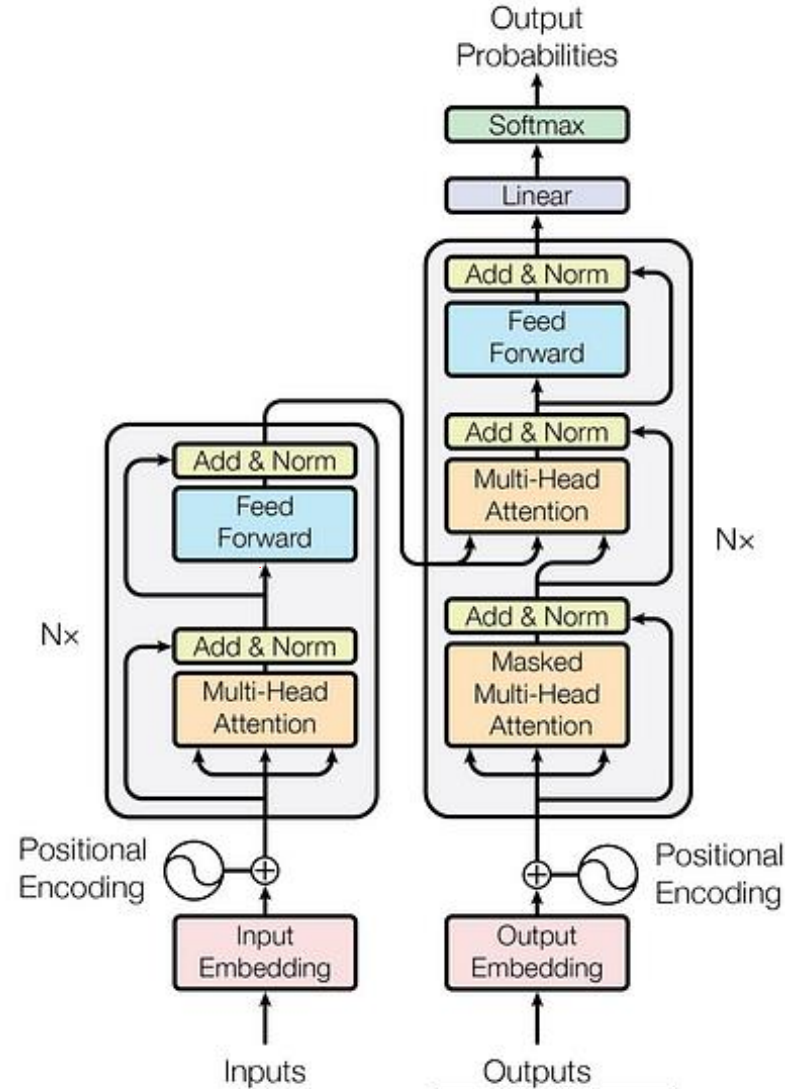
3. Can we build models that can encode items in a sequence in parallel?

X Markov models X Recurrent Neural Networks ✓ Transformers

Transformer architecture: Vaswani et al 2017

BERT

Encoder



GPT

Decoder

Output (labels, a sequence of words)



BERT, RoBERTa, GPT-2, GPT-3, ...



The fat cat sat on the mat

Output (labels, a sequence of words)



BERT, RoBERTa, GPT-2, GPT-3, ...



The fat cat sat on the mat

Output (labels, a sequence of words)



BERT, RoBERTa, GPT-2, GPT-3, ...

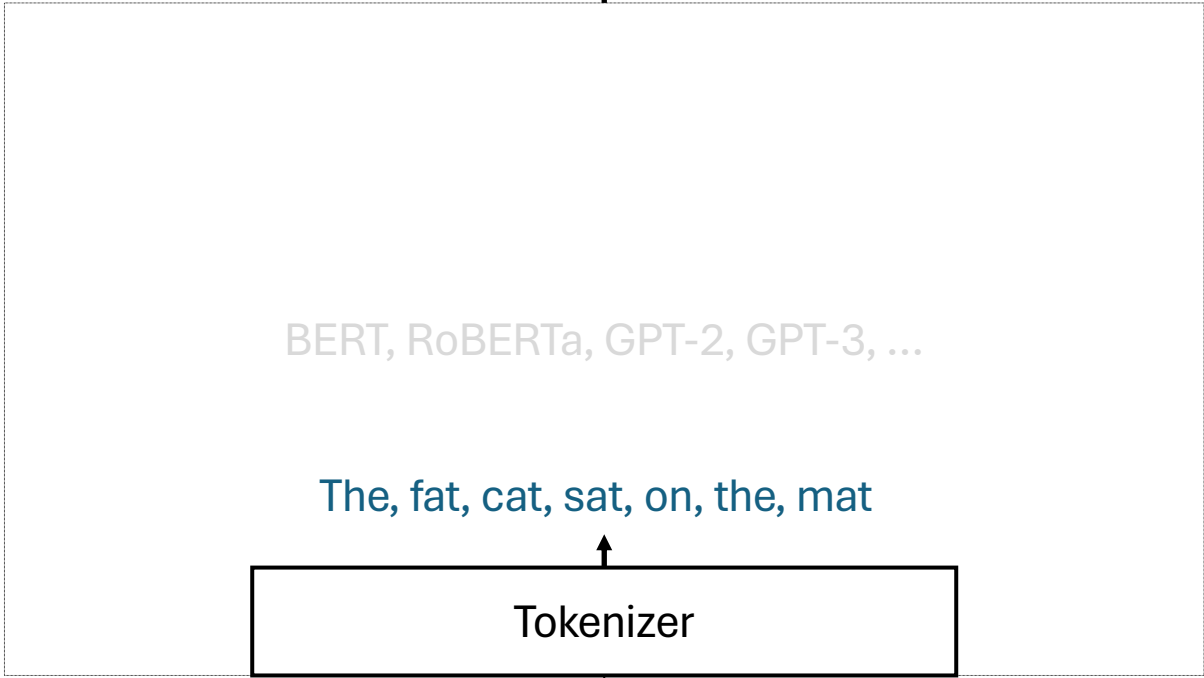
The, fat, cat, sat, on, the, mat



Tokenizer



The fat cat sat on the mat



Output (labels, a sequence of words)



BERT, RoBERTa, GPT-2, GPT-3, ...

The, fat, cat, sat, on, the, mat

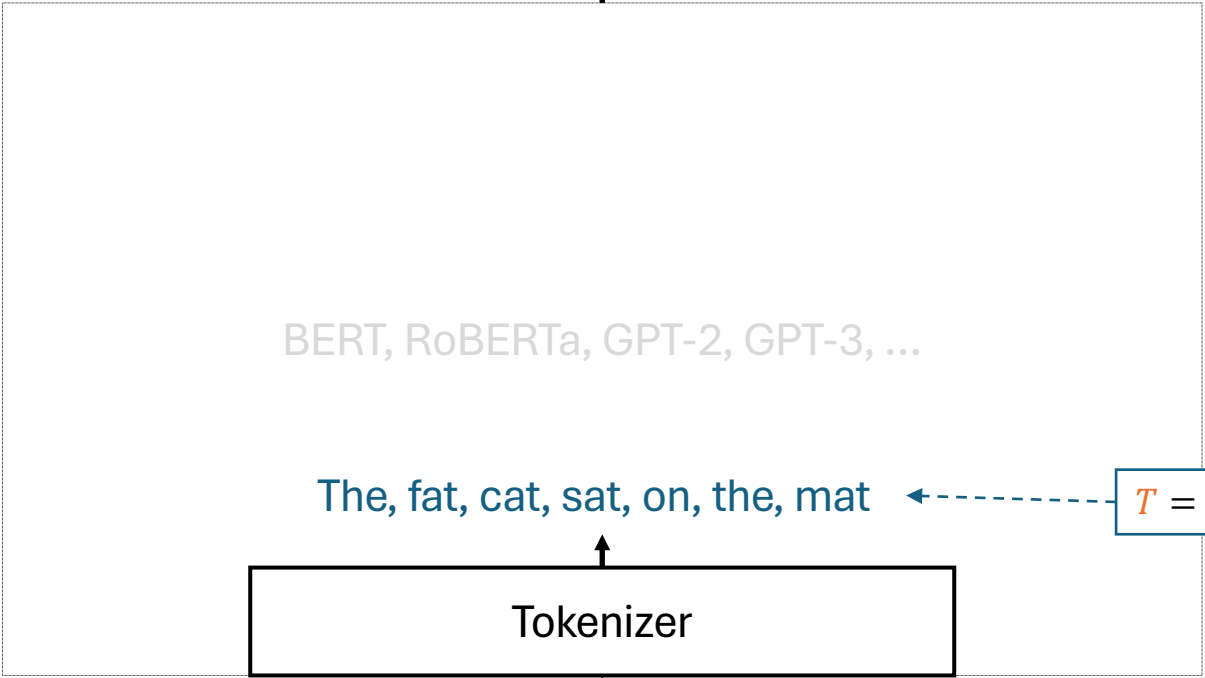


$T = 7$ tokens in this sequence.

Tokenizer



The fat cat sat on the mat



Output (labels, a sequence of words)



BERT, RoBERTa, GPT-2, GPT-3, ...

The, fat, cat, sat, on, the, mat

$T = 7$ tokens in this sequence.

Tokenizer

- More than whitespace splitting to handle unknown/very long words
- Common approach: **Byte-pair encoding**

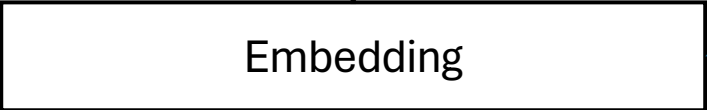
The fat cat sat on the mat



Output (labels, a sequence of words)



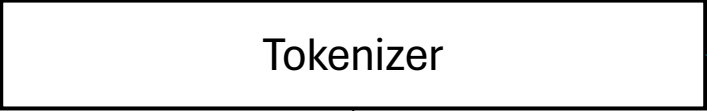
$A T \times d$ matrix



Embedding



The, fat, cat, sat, on, the, mat



Tokenizer

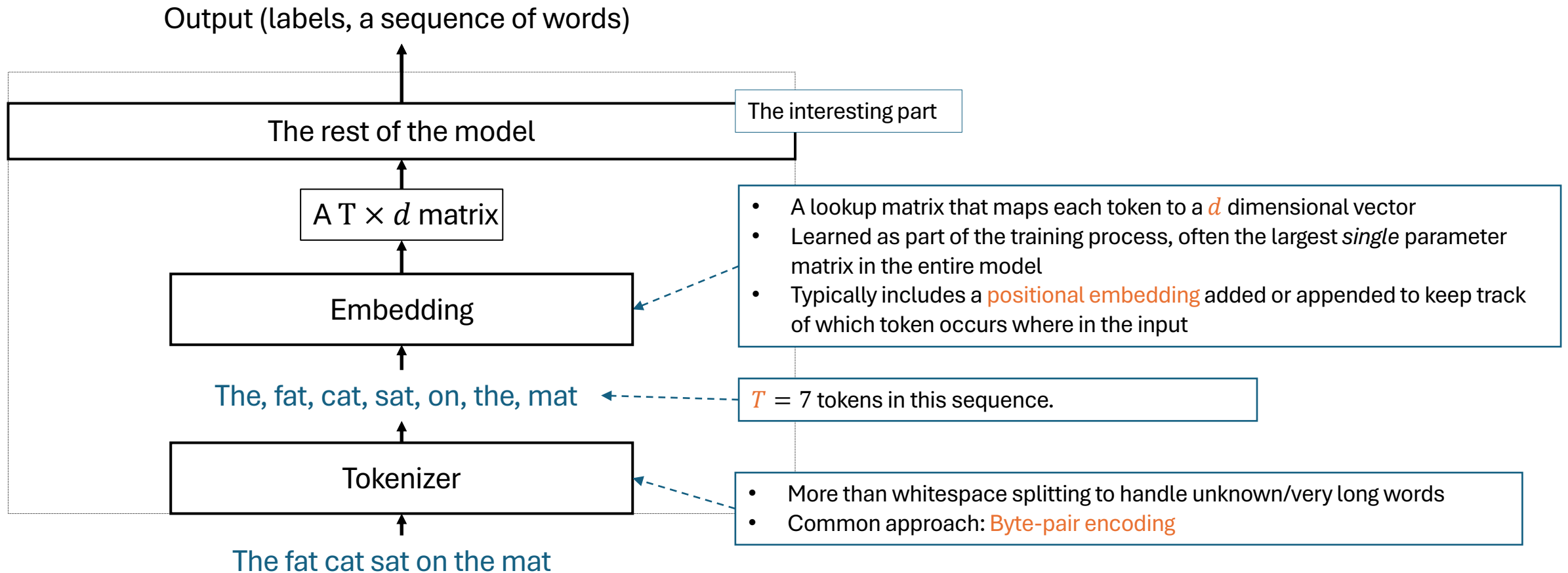


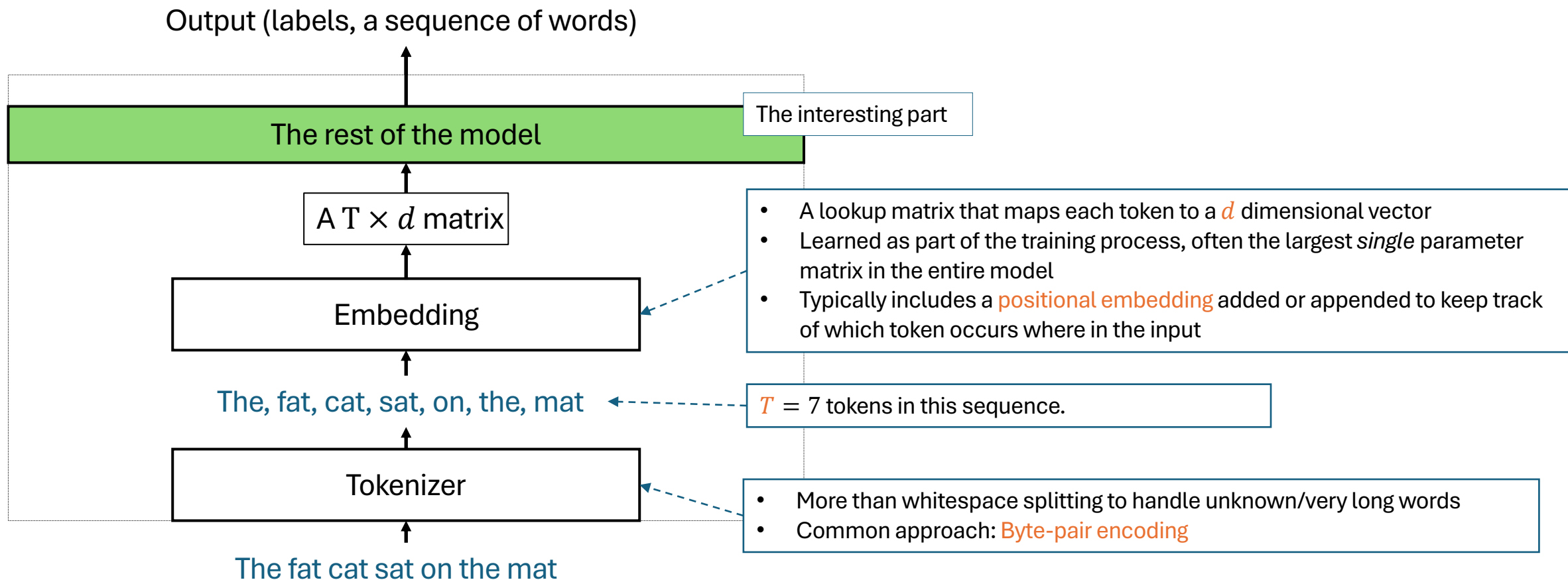
The fat cat sat on the mat

- A lookup matrix that maps each token to a d dimensional vector
- Learned as part of the training process, often the largest *single* parameter matrix in the entire model
- Typically includes a **positional embedding** added or appended to keep track of which token occurs where in the input

$T = 7$ tokens in this sequence.

- More than whitespace splitting to handle unknown/very long words
- Common approach: **Byte-pair encoding**





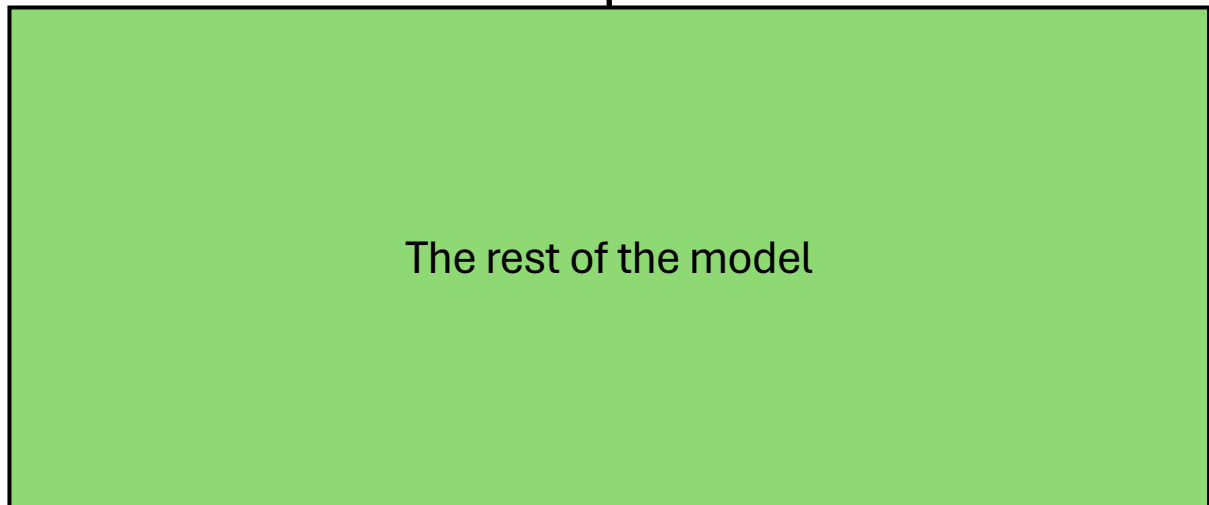
Output (labels, a sequence of words)



The rest of the model



$T \times d$ matrix after embedding tokens



Output (labels, a sequence of words)



The rest of the model



$T \times d$ matrix after embedding tokens

T: Sequence length
 d : Embedding size

Output (labels, a sequence of words)



The rest of the model

Transformer Layer



$T \times d$ matrix after embedding tokens

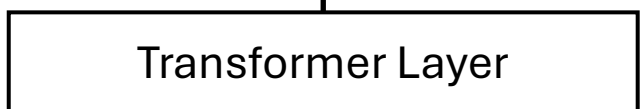
T: Sequence length

d : Embedding size

Output (labels, a sequence of words)



⋮



$T \times d$ matrix after embedding tokens

T: Sequence length

d : Embedding size

Output (labels, a sequence of words)



Transformer Layer

⋮

Transformer Layer



Transformer Layer



$T \times d$ matrix after embedding tokens

T: Sequence length

d : Embedding size

Output (labels, a sequence of words)

A small model (typically linear + softmax) that produces the desired probabilities

Transformer Layer

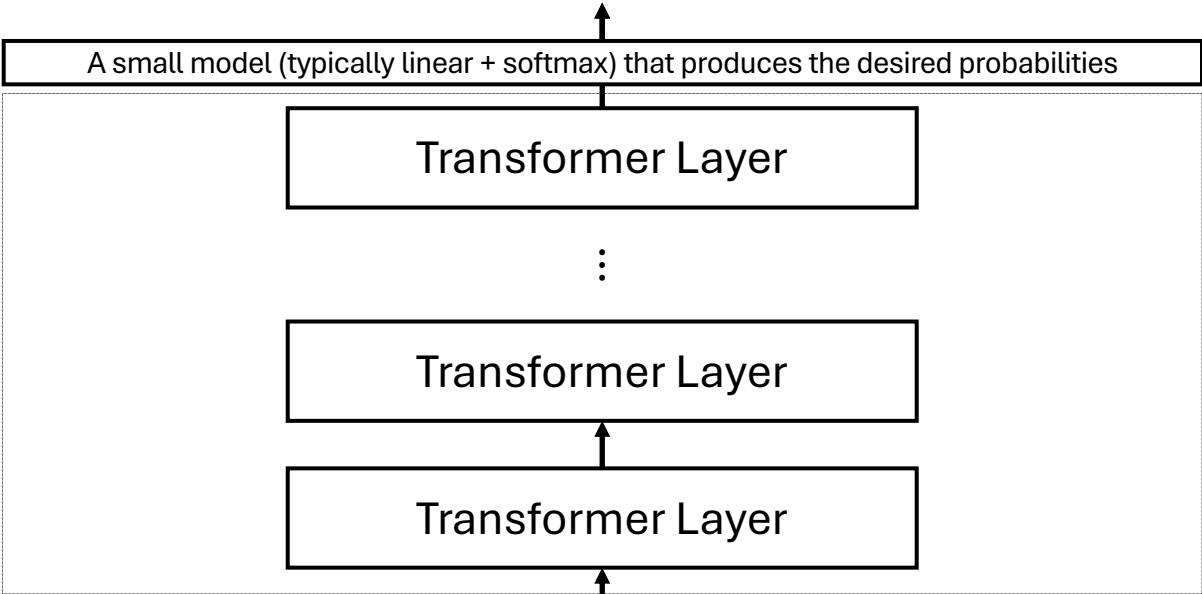
⋮

Transformer Layer

Transformer Layer

$T \times d$ matrix after embedding tokens

T: Sequence length
 d : Embedding size



Output (labels, a sequence of words)

A small model (typically linear + softmax) that produces the desired probabilities

Transformer Layer

⋮

Transformer Layer

Transformer Layer

Each transformer layer converts a $T \times d$ matrix into a "transformed" $T \times d$ matrix

Transformer layers are structurally identical, but have their own parameters

Two different types of transformers in the original paper:

- Encoder: for BERT, etc whose goal is to embed text
- Decoder: for GPT etc whose goal is to generate text
- Minor differences between them

$T \times d$ matrix after embedding tokens

T: Sequence length
 d : Embedding size

Output (labels, a sequence of words)

A small model (typically linear + softmax) that produces the desired probabilities

Transformer Layer

⋮

Transformer Layer

Transformer Layer

$T \times d$ matrix after embedding tokens

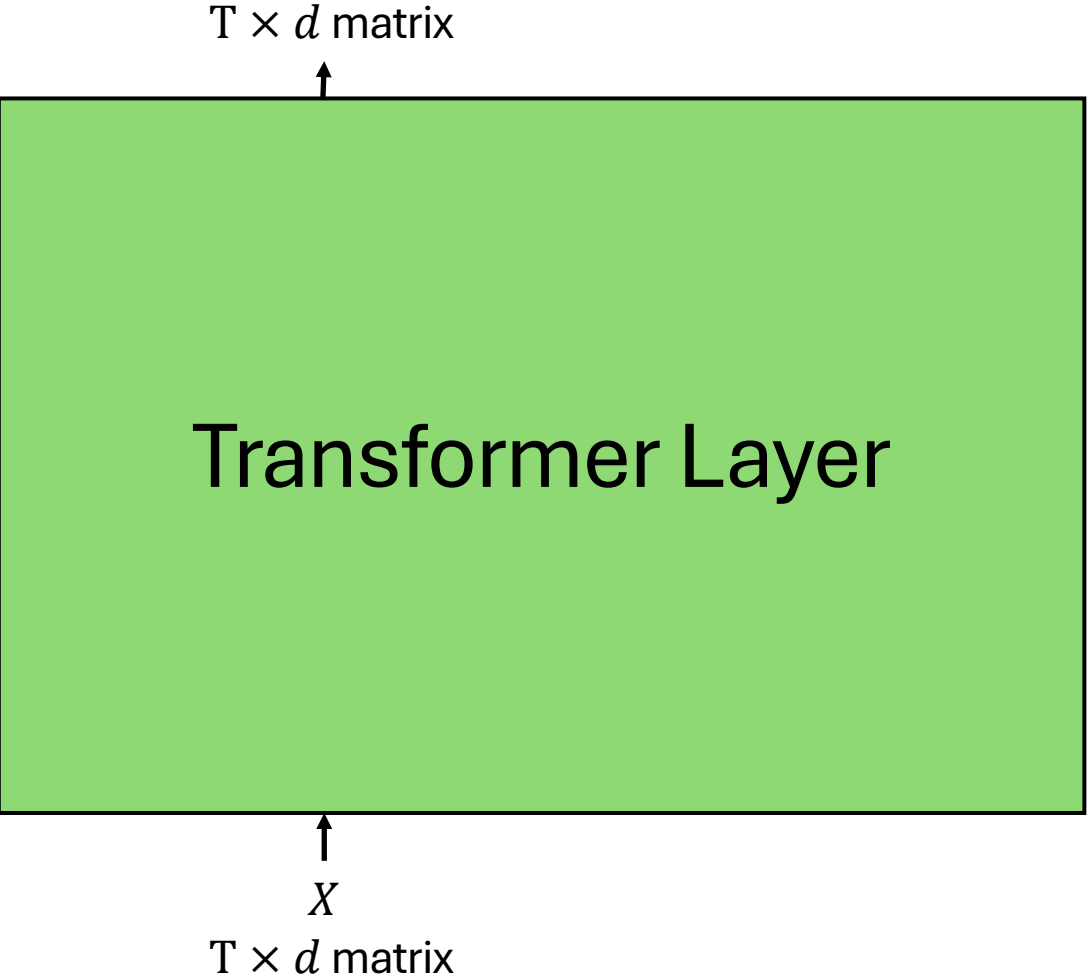
T: Sequence length
 d : Embedding size

Each transformer layer converts a $T \times d$ matrix into a "transformed" $T \times d$ matrix

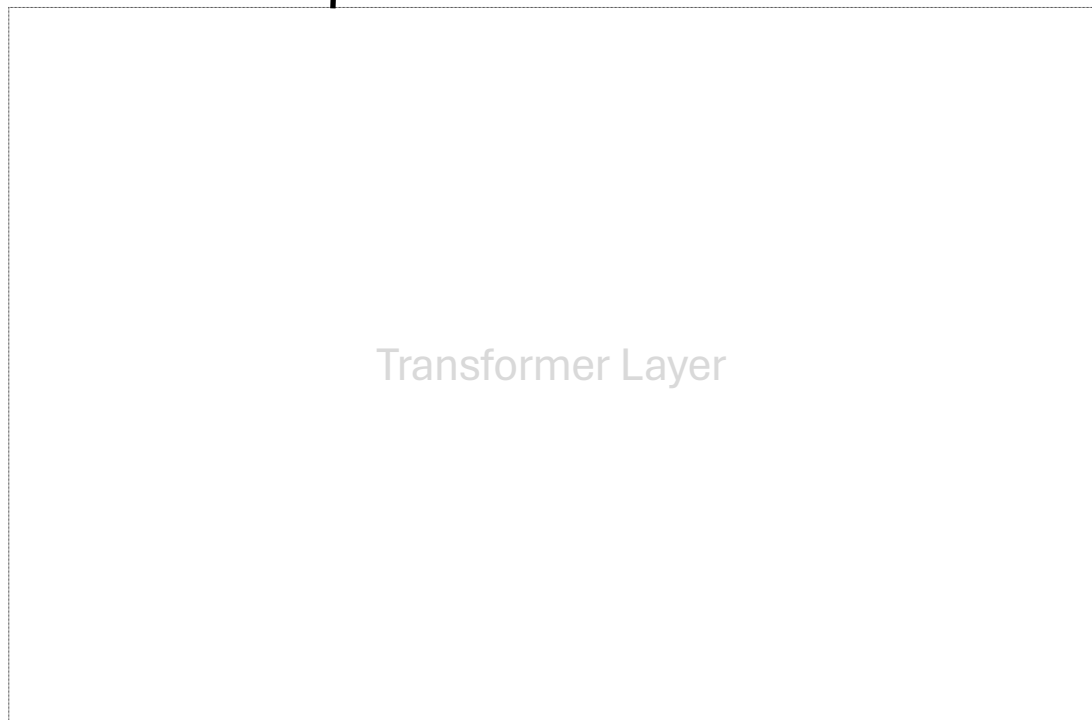
Transformer layers are structurally identical, but have their own parameters

Two different types of transformers in the original paper:

- Encoder: for BERT, etc whose goal is to embed text
- Decoder: for GPT etc whose goal is to generate text
- Minor differences between them



$T \times d$ matrix



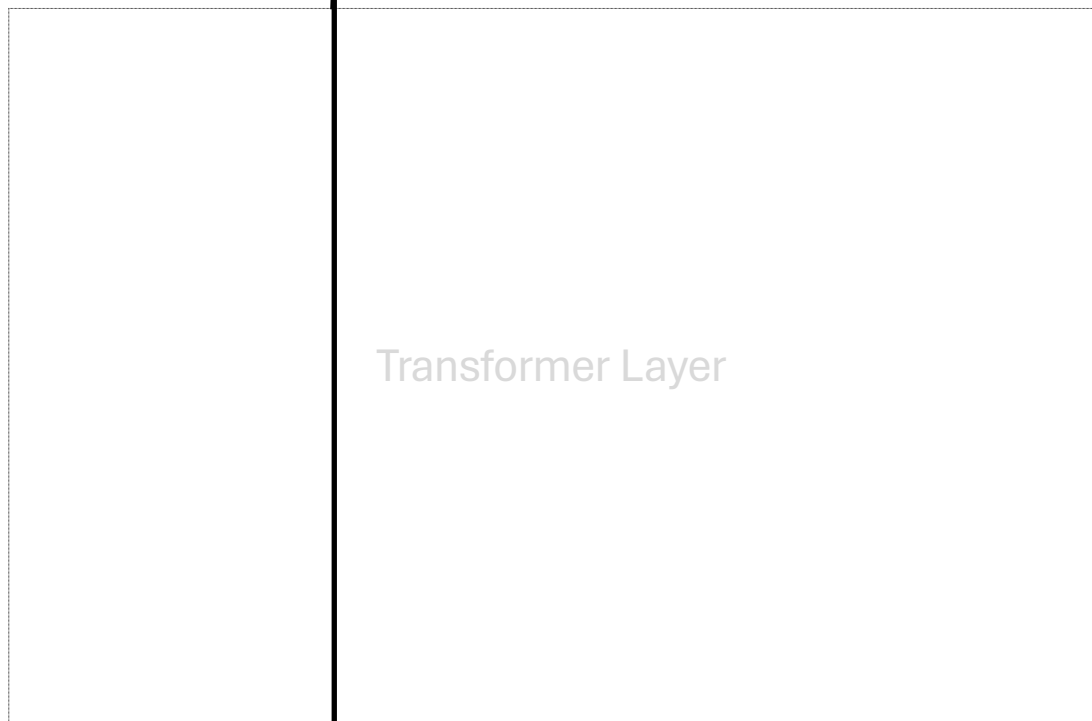
Transformer Layer



X

$T \times d$ matrix

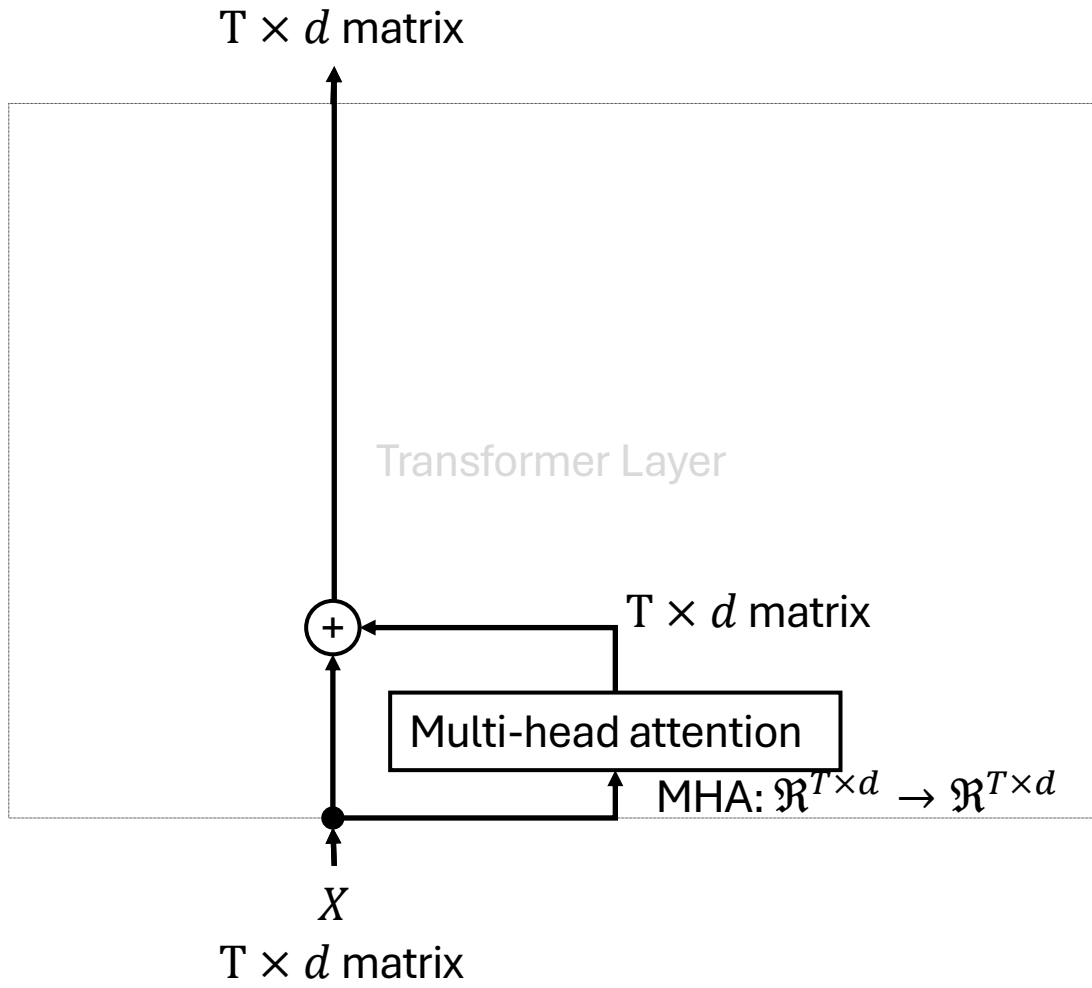
$T \times d$ matrix



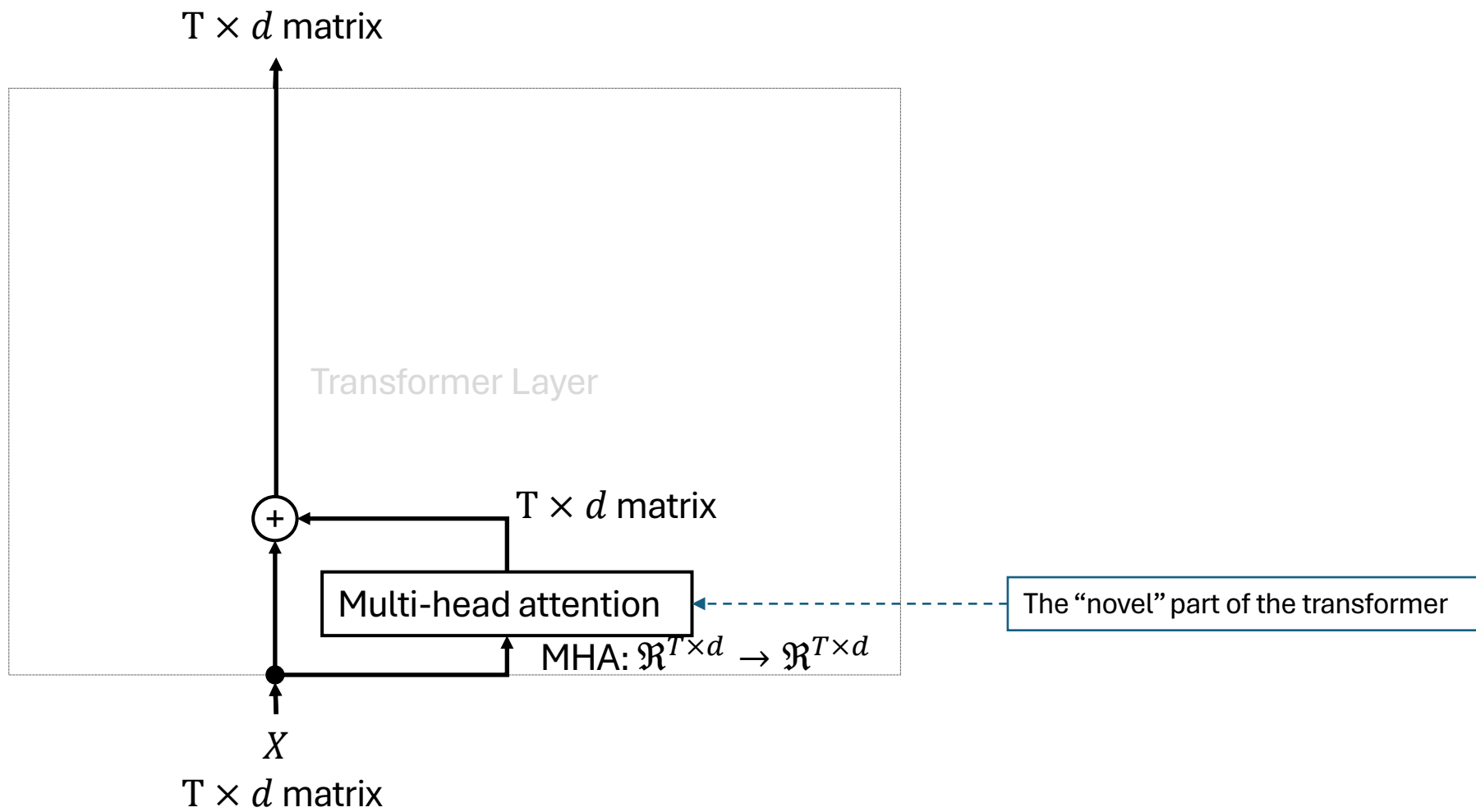
Transformer Layer

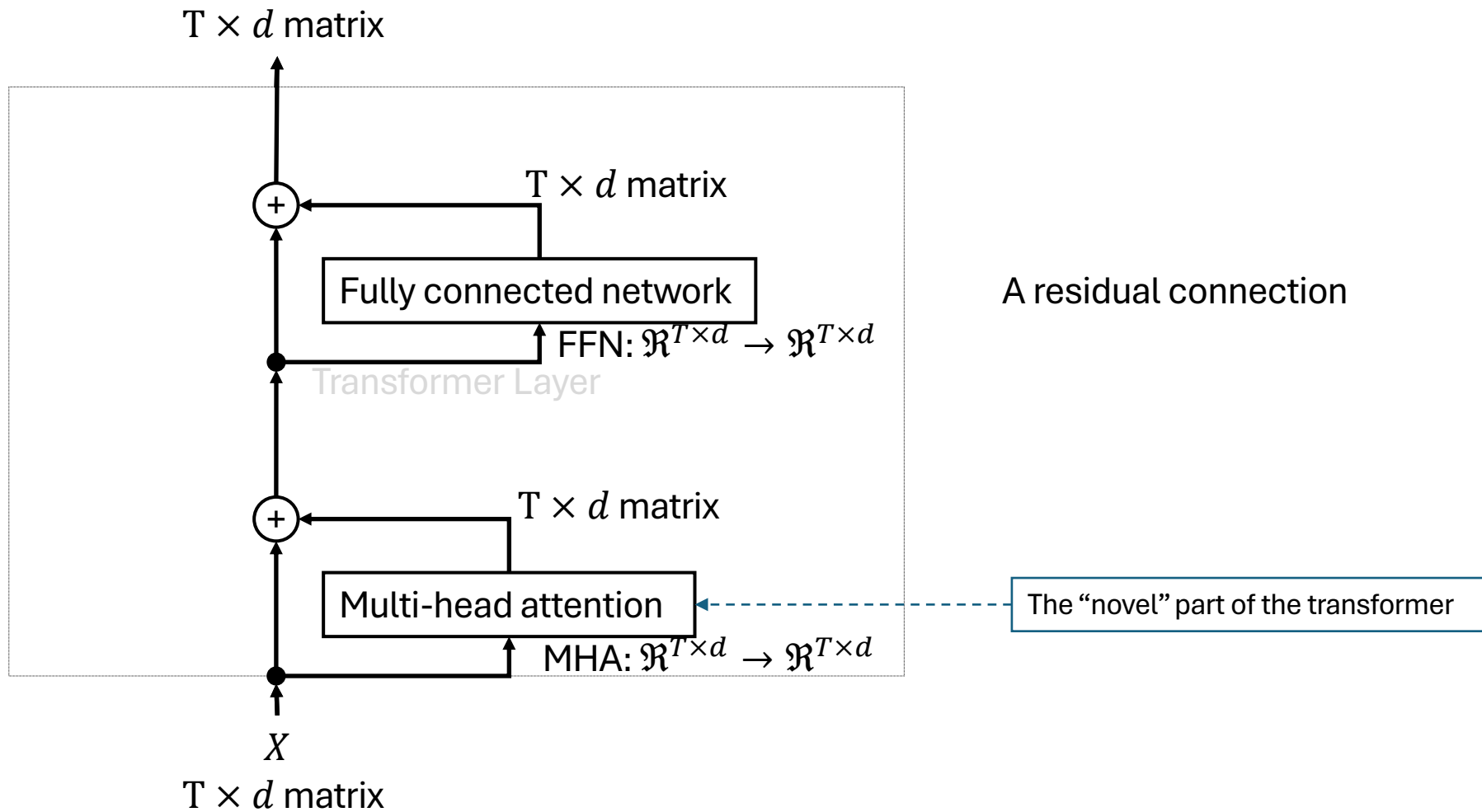
X

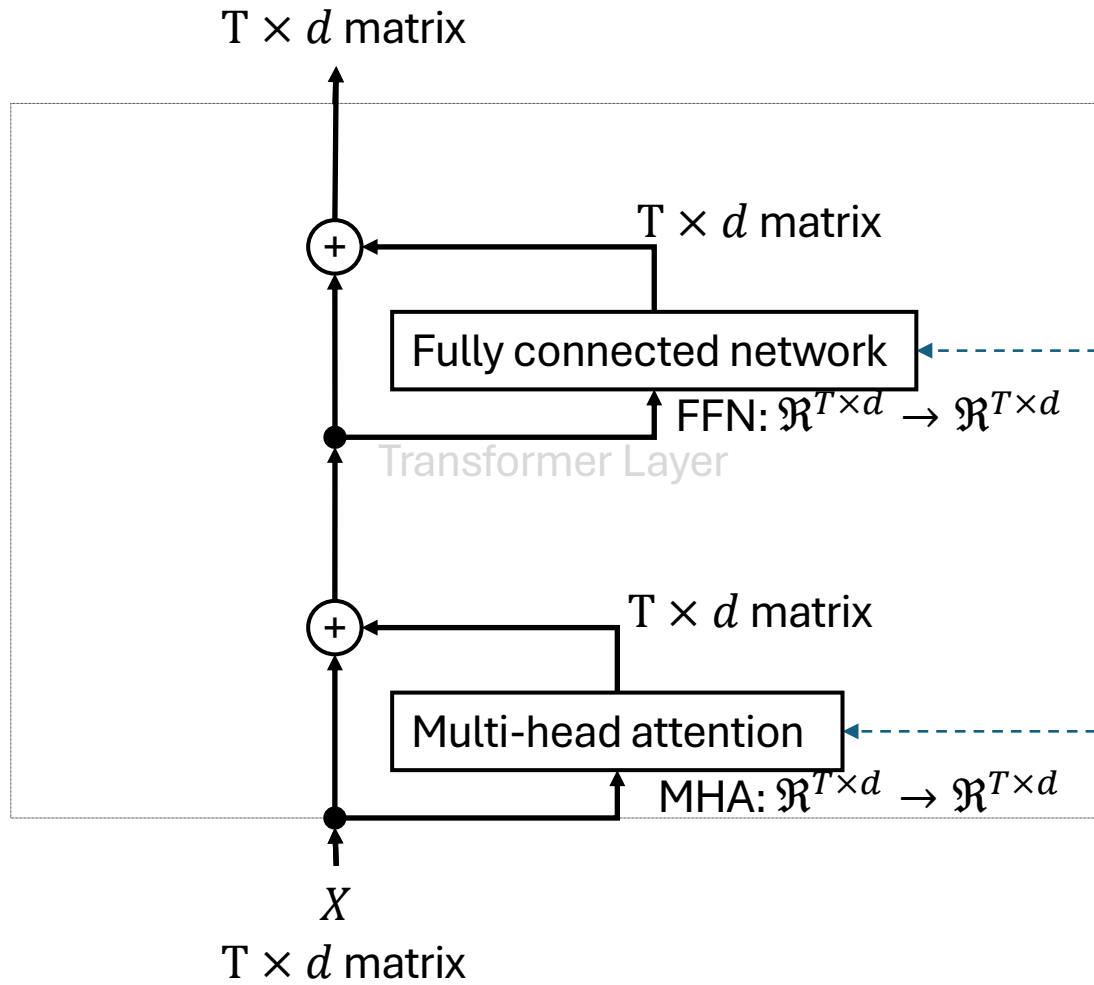
$T \times d$ matrix



A residual connection

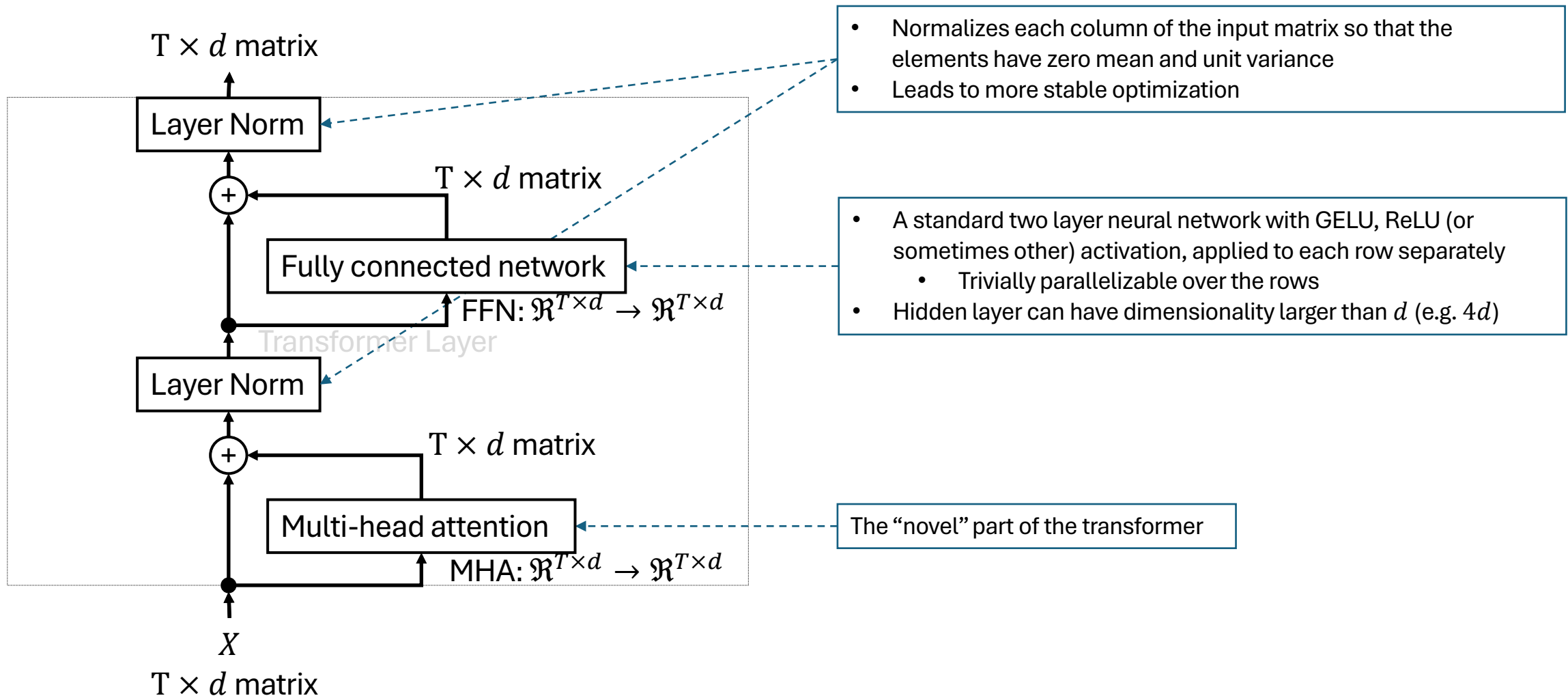






- A standard two layer neural network with GELU, ReLU (or sometimes other) activation, applied to each row separately
 - Trivially parallelizable over the rows
- Hidden layer can have dimensionality larger than d (e.g. $4d$)

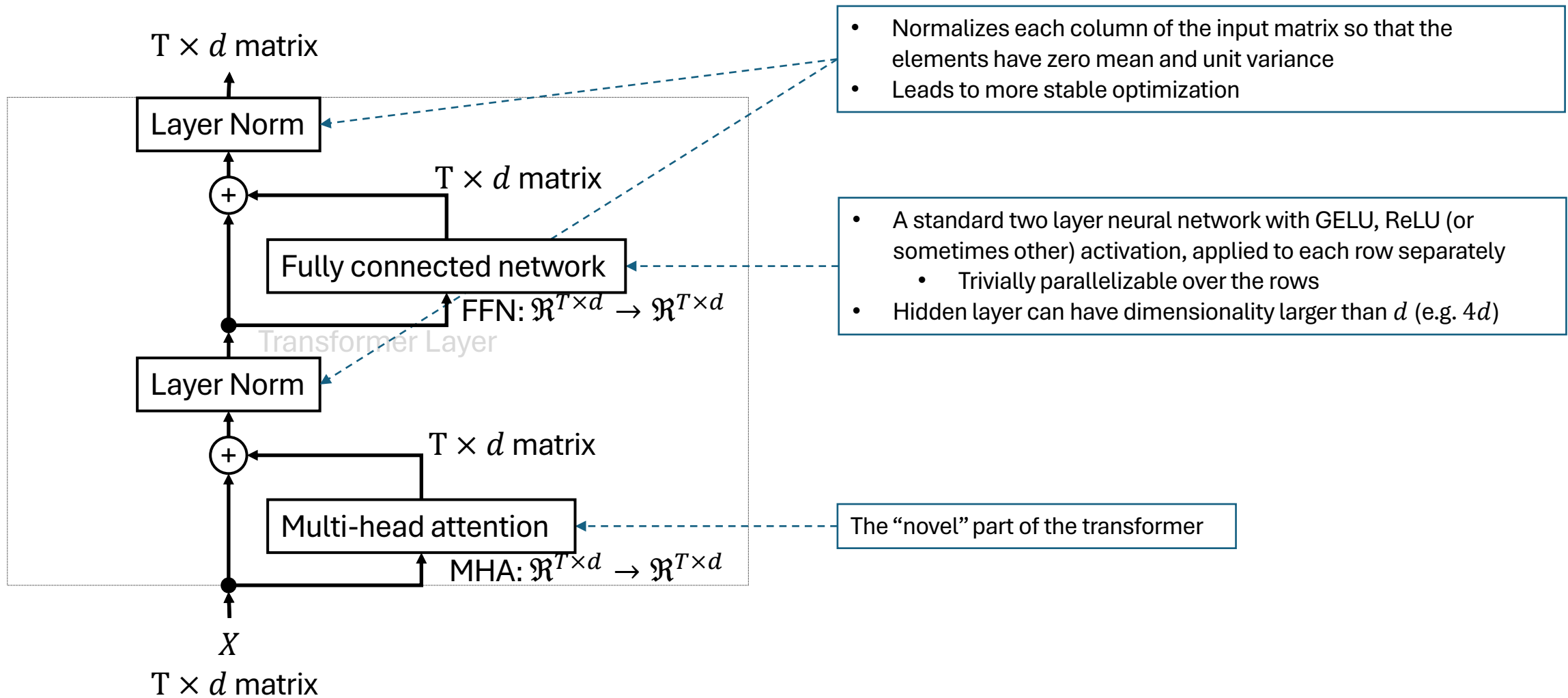
The “novel” part of the transformer



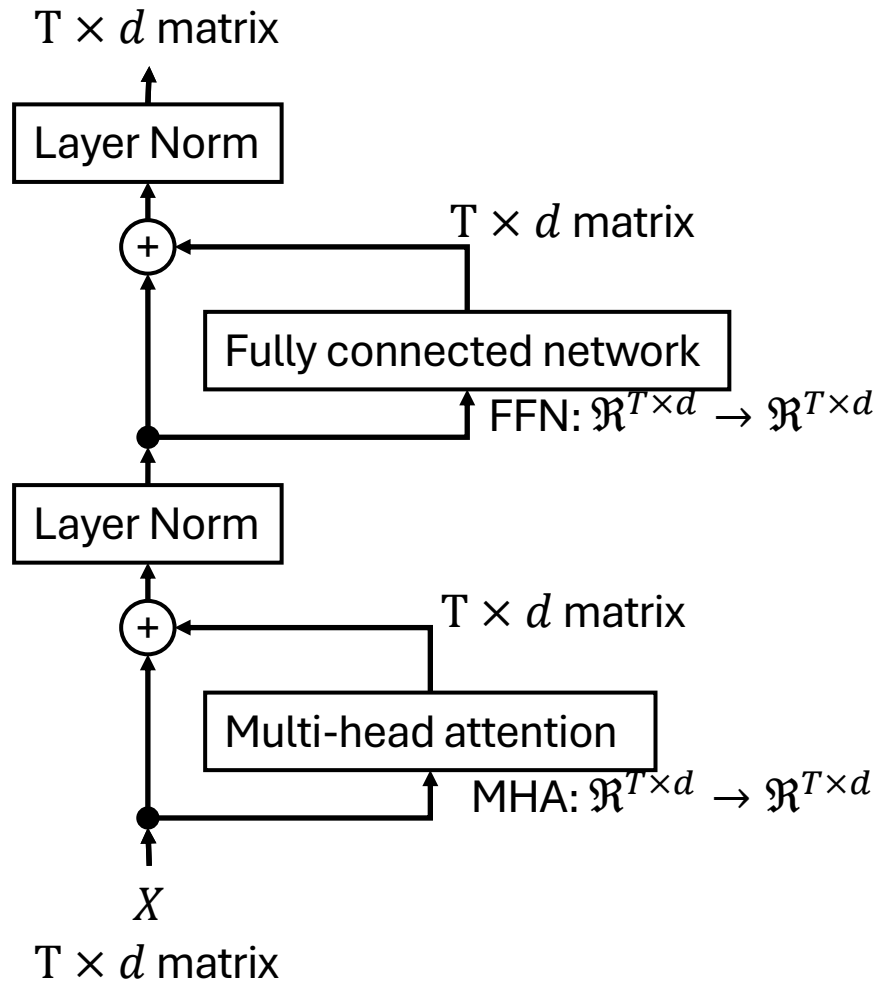
- Normalizes each column of the input matrix so that the elements have zero mean and unit variance
- Leads to more stable optimization

- A standard two layer neural network with GELU, ReLU (or sometimes other) activation, applied to each row separately
 - Trivially parallelizable over the rows
 - Hidden layer can have dimensionality larger than d (e.g. $4d$)

The "novel" part of the transformer



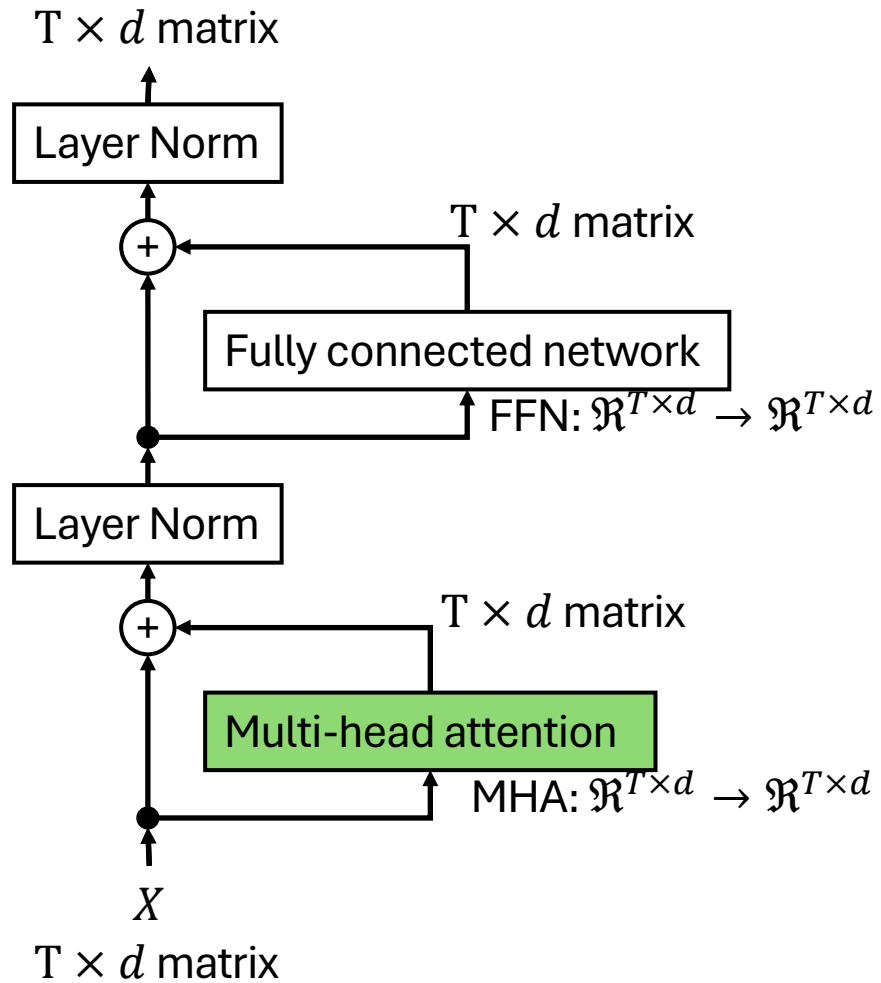
This is the transformer encoder. The decoder transformer has a bit more detail. We will encounter the details later



```
def transformer_layer(X):
    X1 = layer_norm1(X + multi_head_attention(X))
    X2 = layer_norm2(X1 + fully_connected(X1))
    return X2
```

$$X_1 = \text{LayerNorm}(X + \text{MHA}(X))$$

$$\text{Result} = \text{LayerNorm}(X_1 + \text{FFN}(X_1))$$

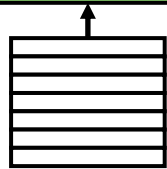
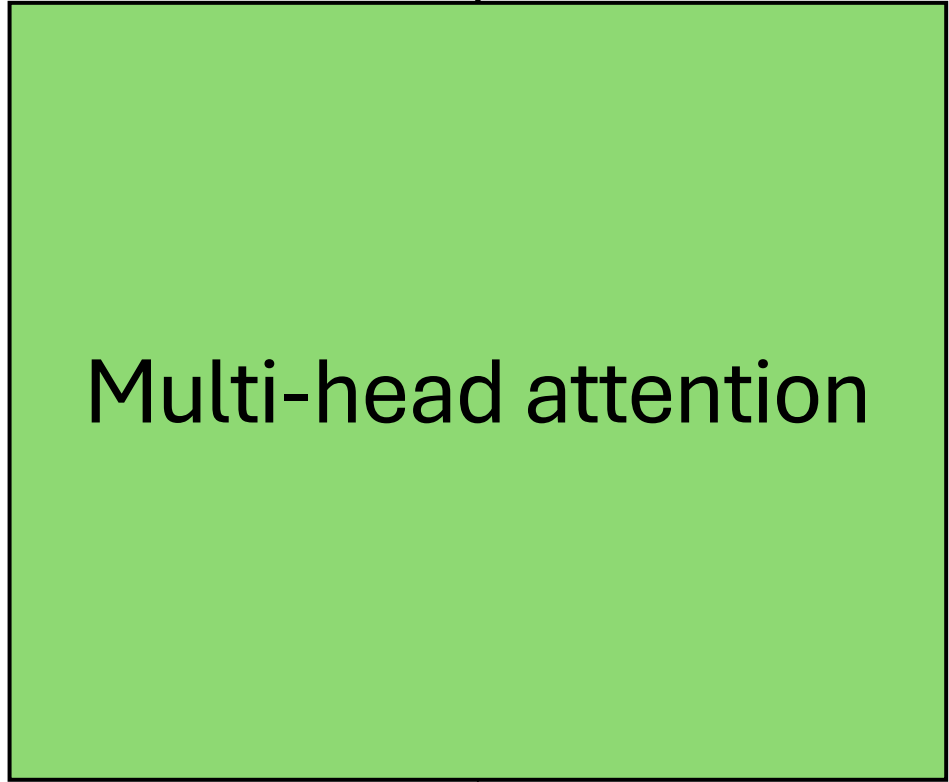


```
def transformer_layer(X):
    X1 = layer_norm1(X + multi_head_attention(X))
    X2 = layer_norm2(X1 + fully_connected(X1))
    return X2
```

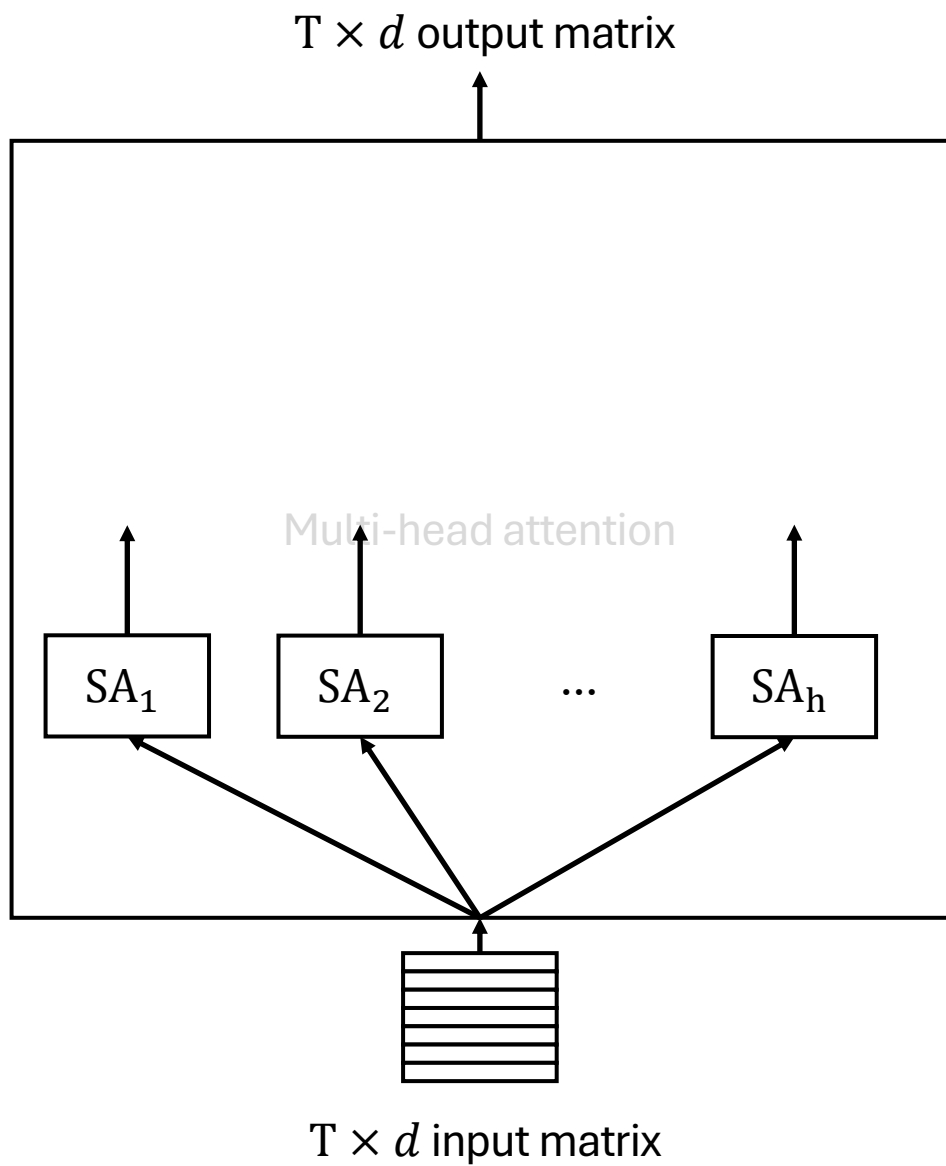
$$X_1 = \text{LayerNorm}(X + \text{MHA}(X))$$

$$\text{Result} = \text{LayerNorm}(X_1 + \text{FFN}(X_1))$$

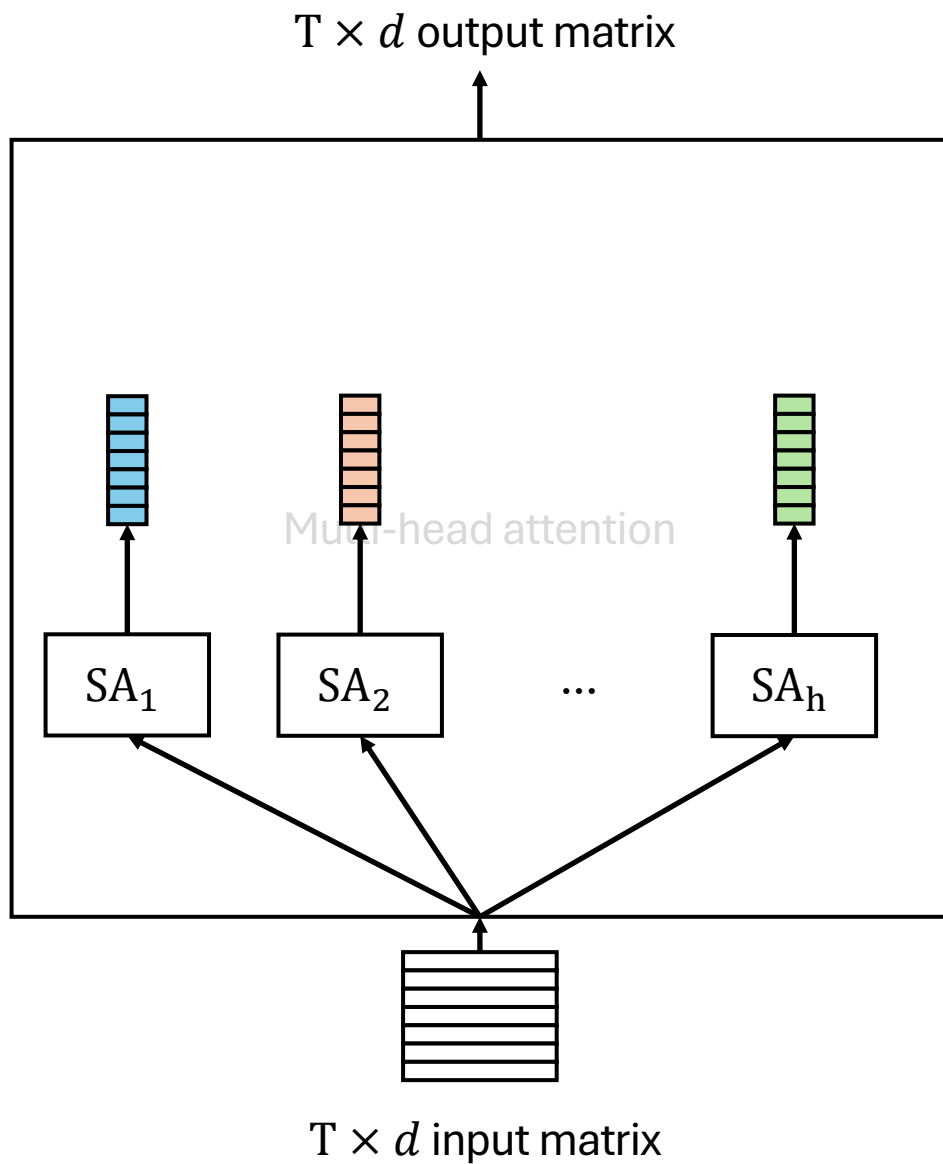
$T \times d$ output matrix



$T \times d$ input matrix

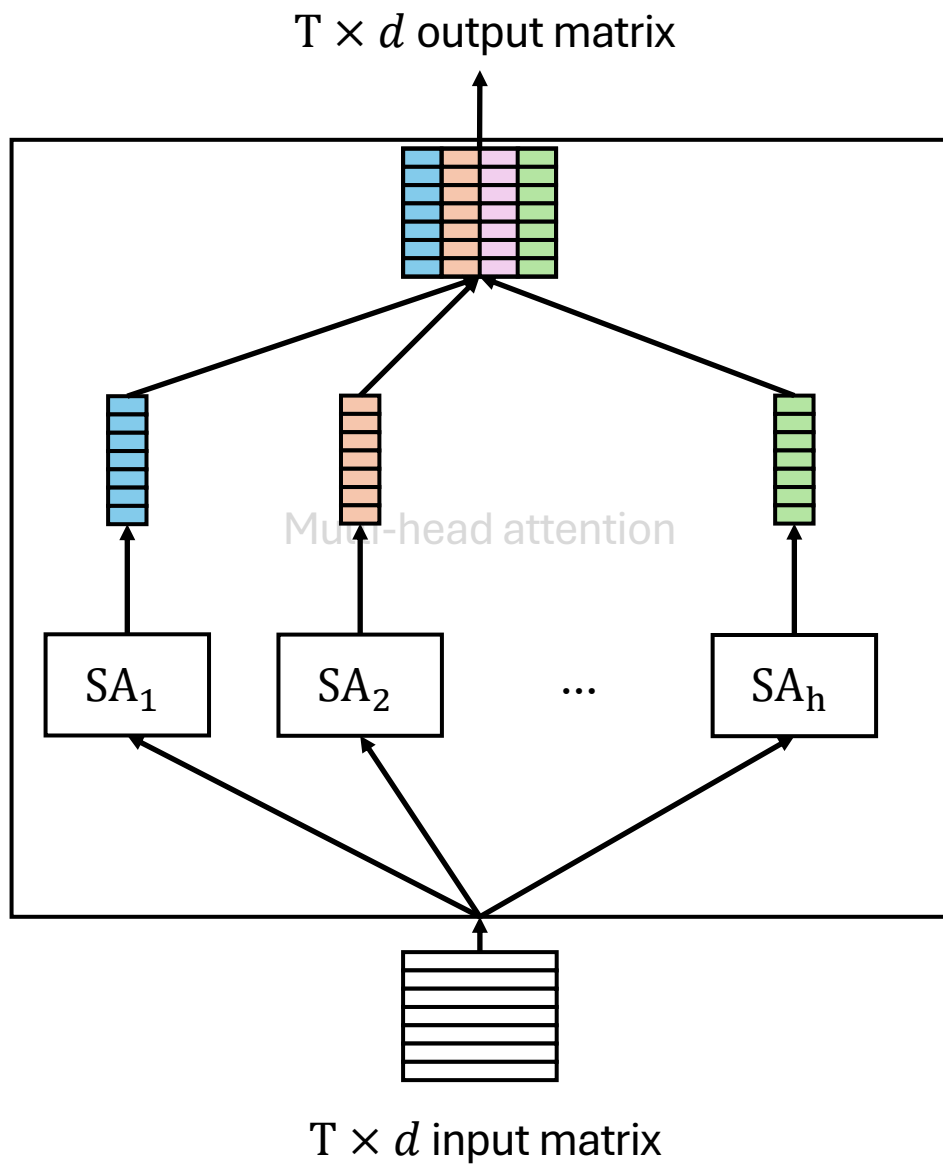


- h self-attention (SA) networks
- Analogous to channels in a CNN



Each produces a matrix of size $T \times \frac{d}{h}$

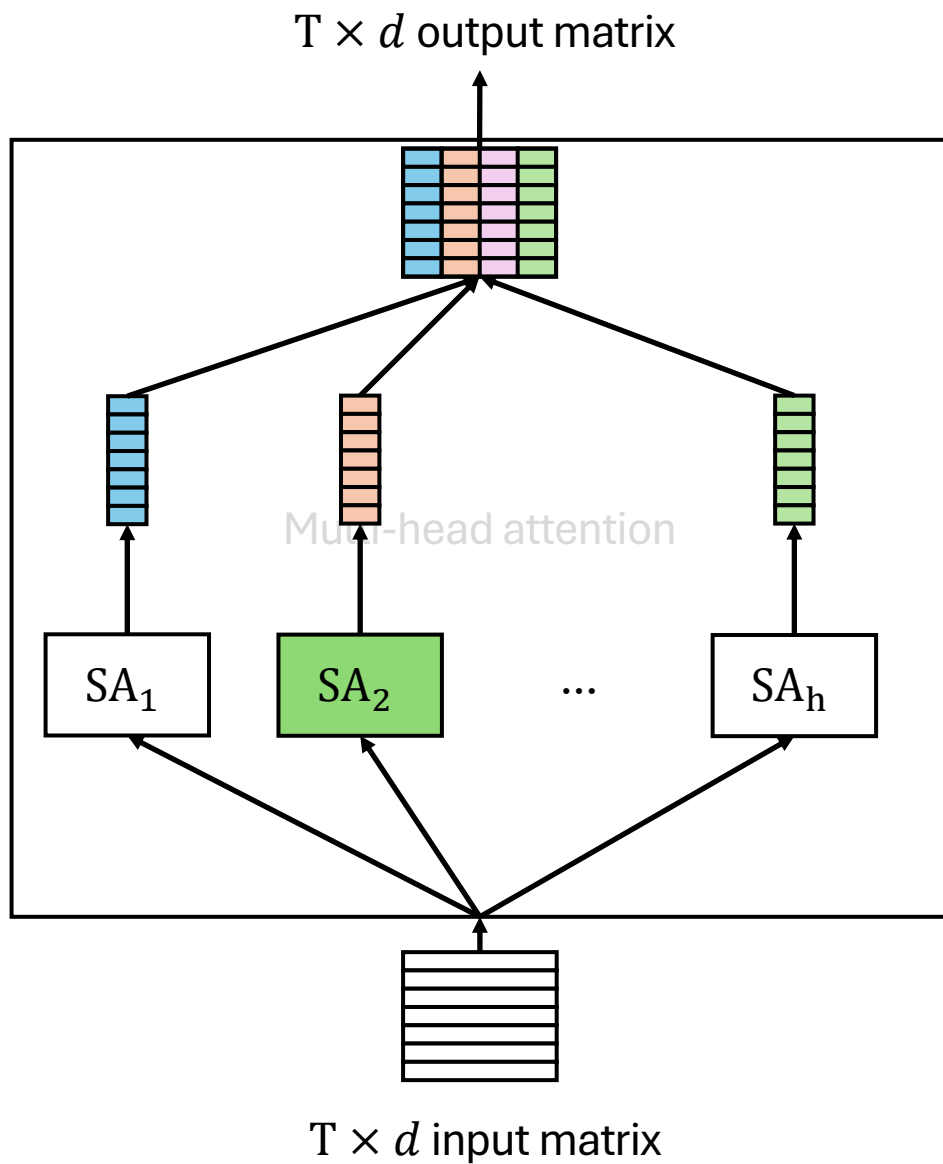
- h self-attention (SA) networks
- Analogous to channels in a CNN



These matrices are simply stacked to produce the output

Each produces a matrix of size $T \times \frac{d}{h}$

- h self-attention (SA) networks
- Analogous to channels in a CNN



These matrices are simply stacked to produce the output

Each produces a matrix of size $T \times \frac{d}{h}$

- h self-attention (SA) networks
- Analogous to channels in a CNN

Self attention: An example



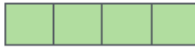
Self attention: An example

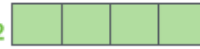
Input

Thinking

Machines

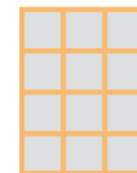
Embedding

x_1 

x_2 



W^Q



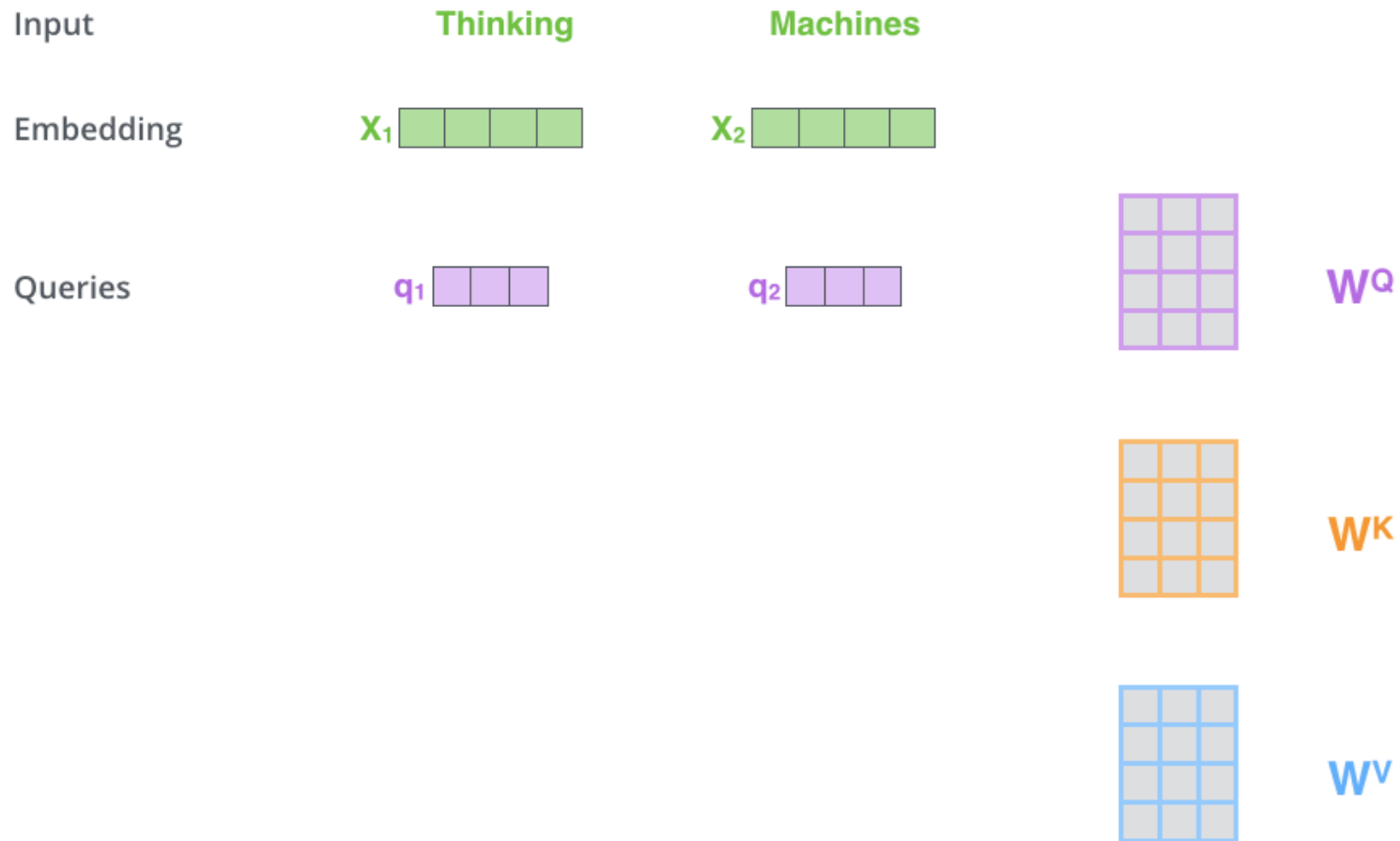
W^K



W^V

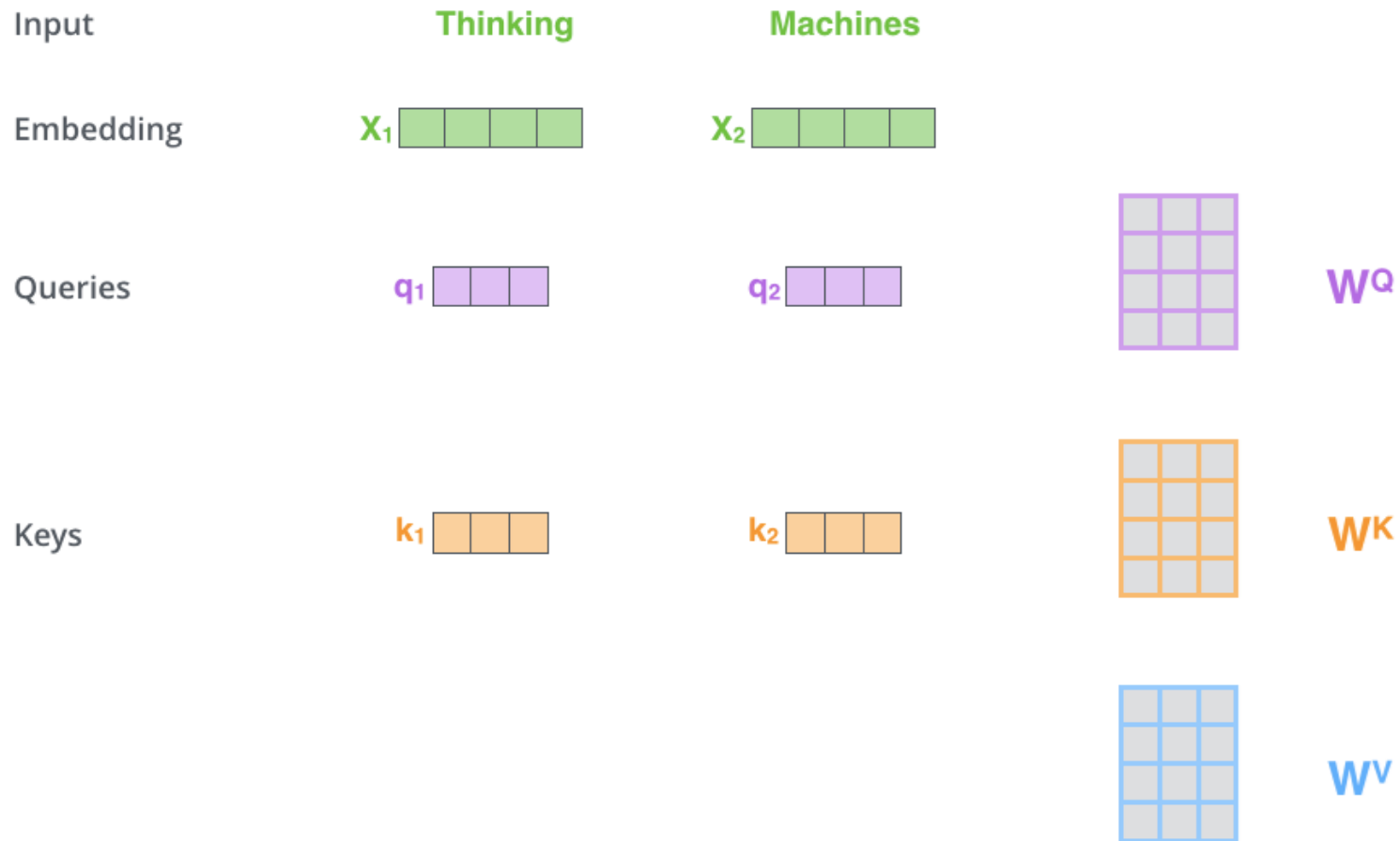
Three parameter matrices associated with this self attention block

Self attention: An example



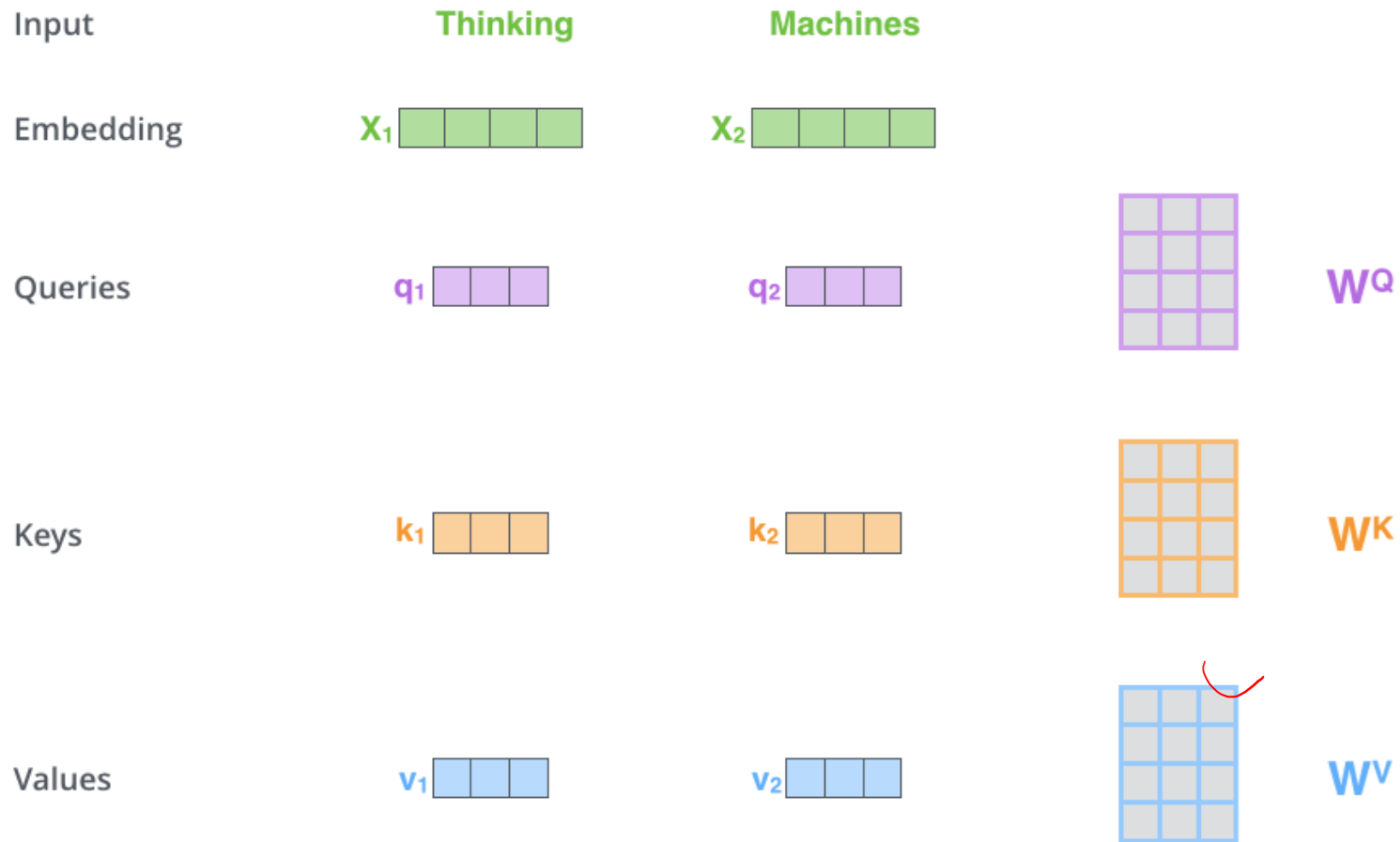
Three parameter matrices associated with this self attention block

Self attention: An example



Three parameter matrices associated with this self attention block

Self attention: An example



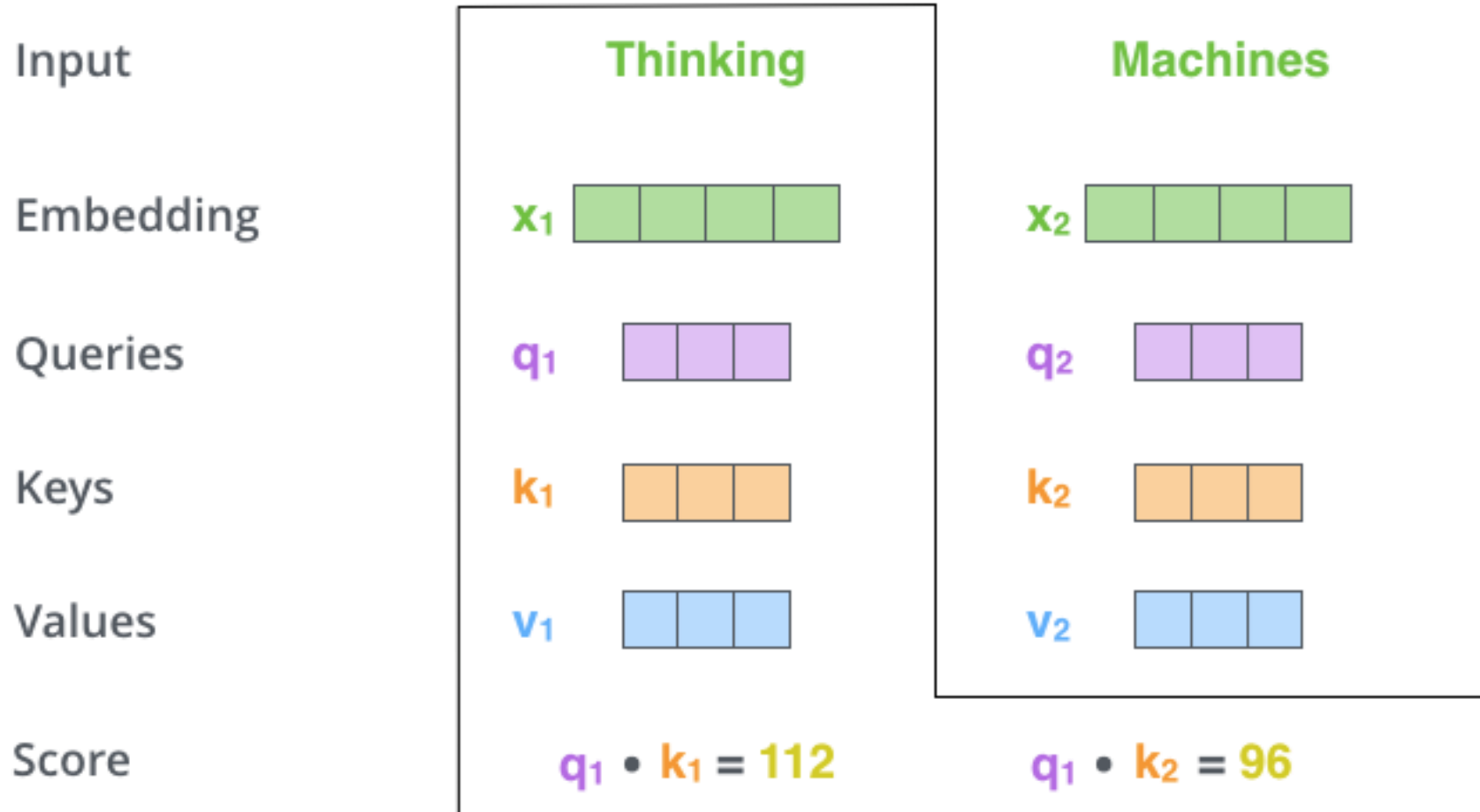
Three parameter matrices associated with this self attention block

For every element of the sequence, create three vectors that are called its *query*, *key* and *value* vectors.

Self attention: An example

For each word, compute the self attention.

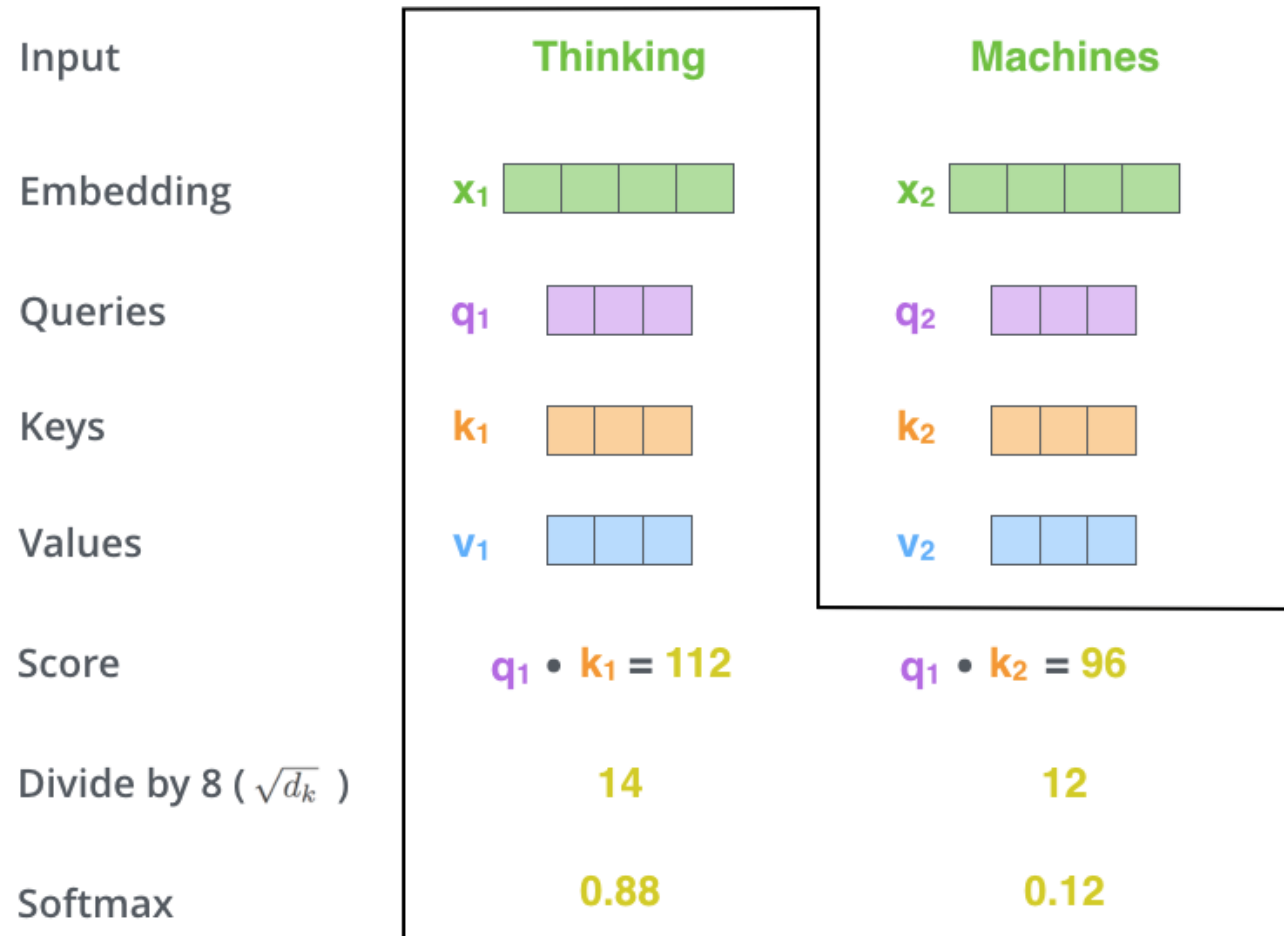
First compute the dot product of its query vector with the key vector of all words in the sentence



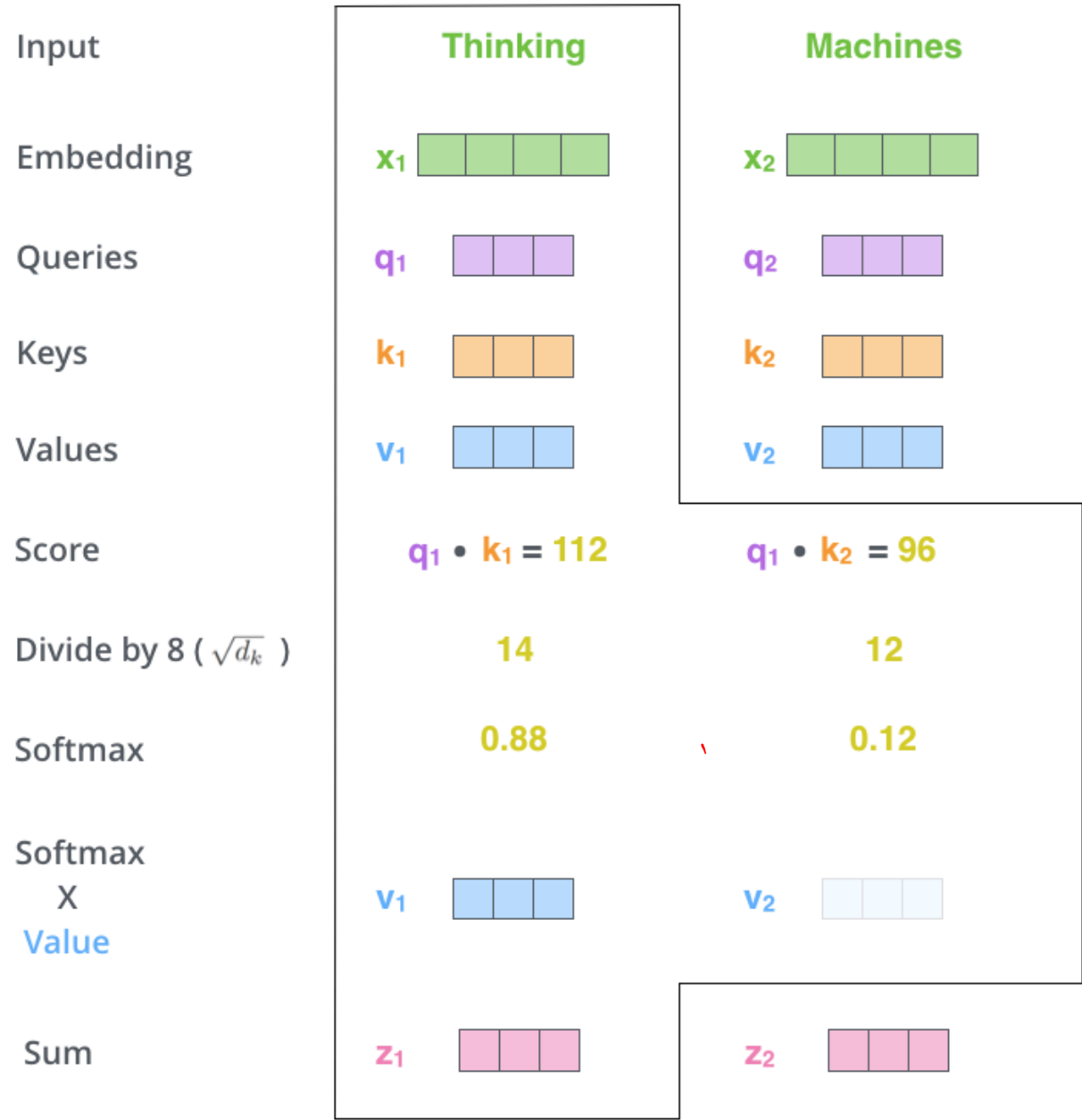
Self attention: An example

For each word, compute the self attention.

Then, normalize



Use the attention probabilities to weigh the value vectors to produce the output vector for that word



Self attention: Illustrated

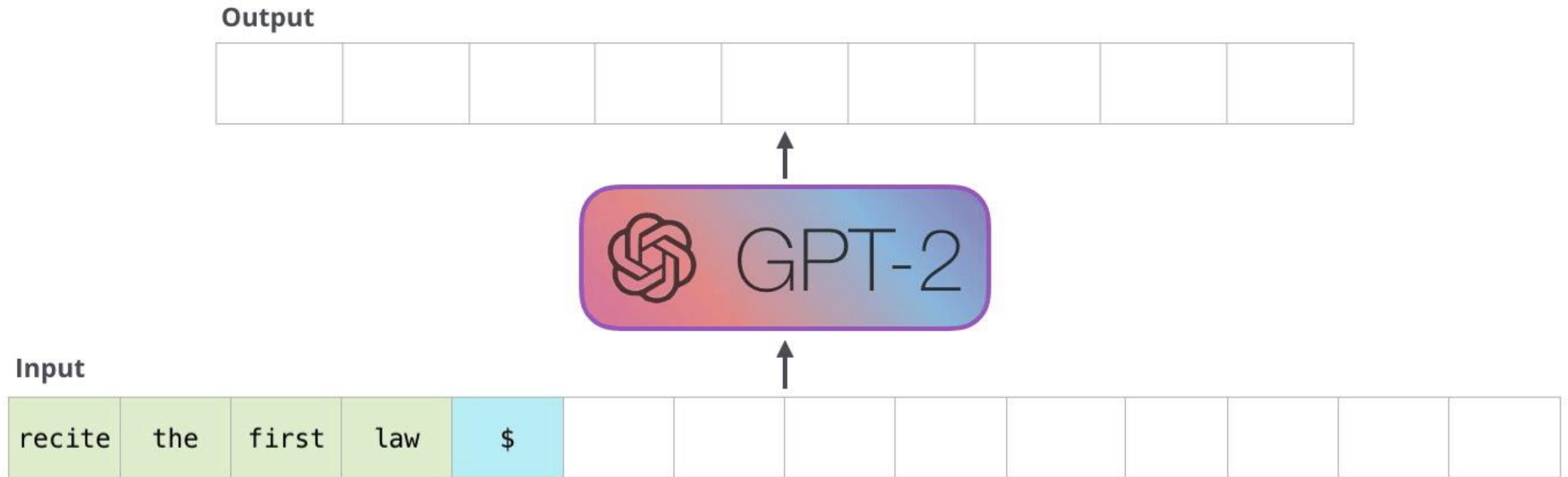
$$X \times W^Q = Q$$

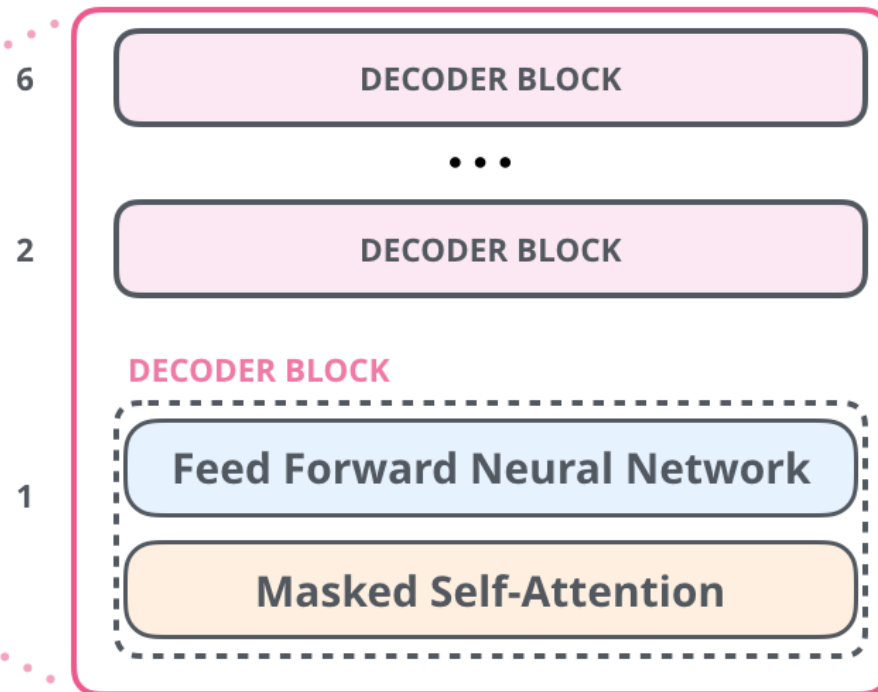
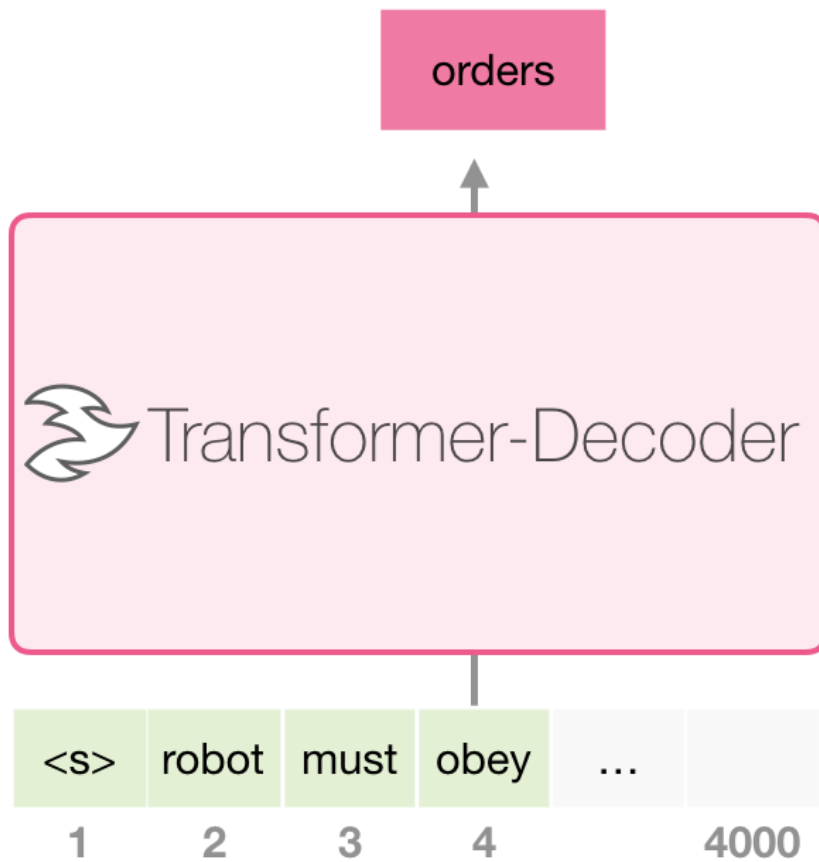
$$X \times W^K = K$$

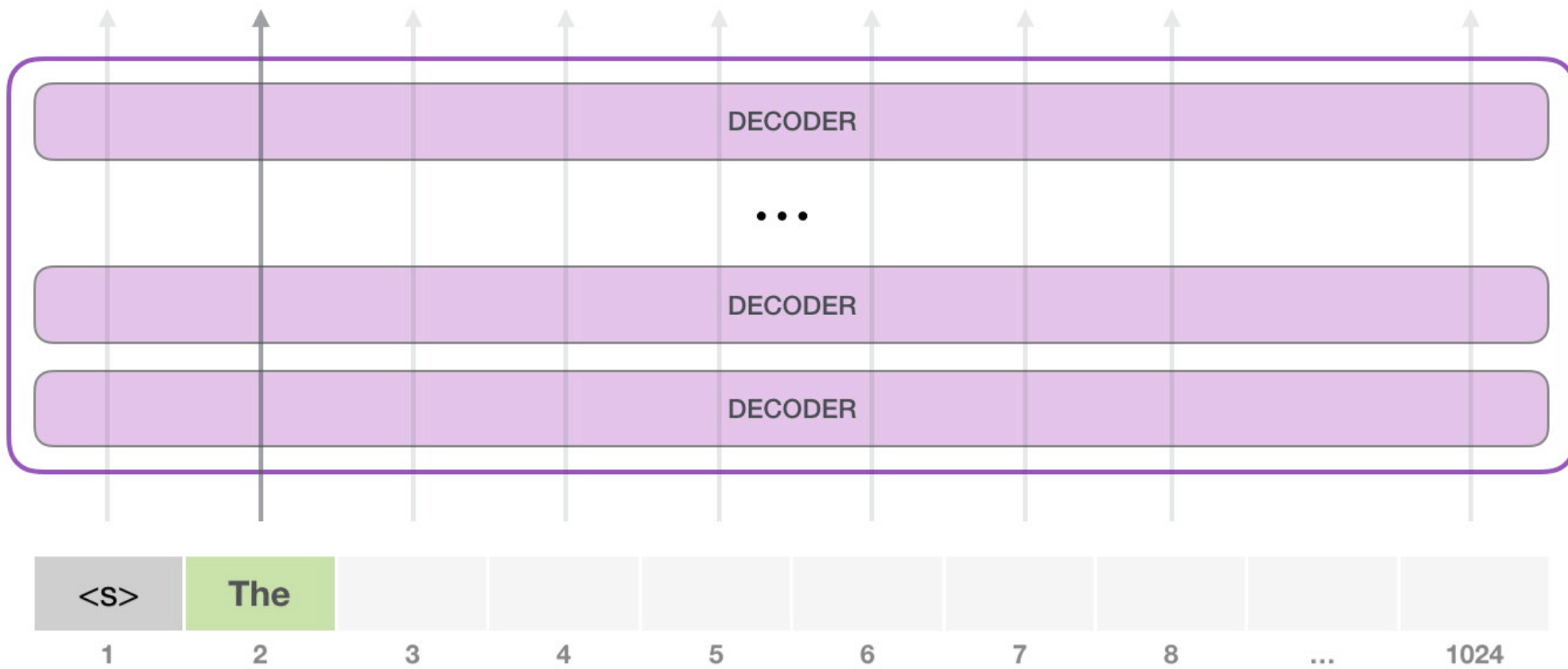
$$X \times W^V = V$$

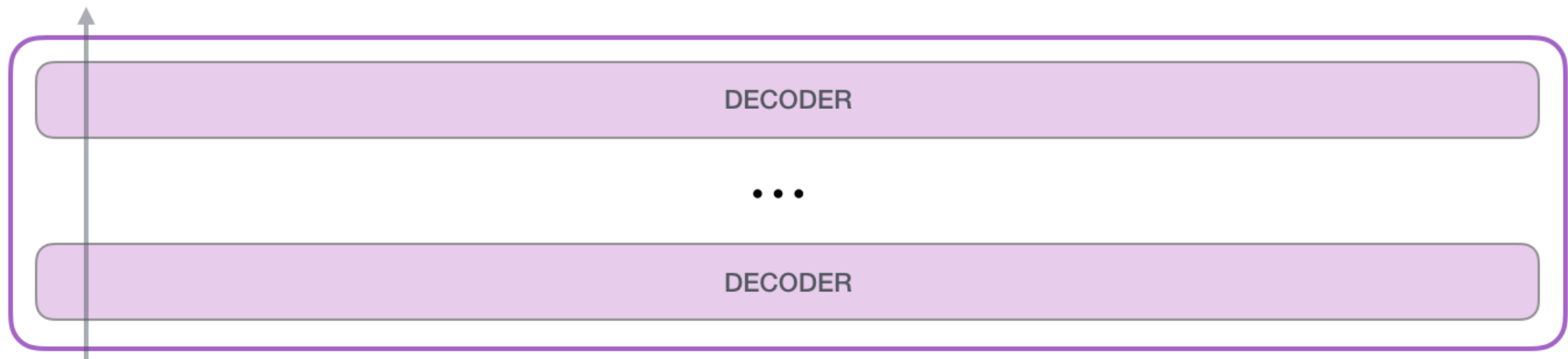
$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) \times V = Z$$

Large (Transformer-Based) Language Models

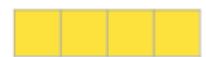






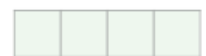


=



Positional encoding for token #1

+



Token embedding of <s>



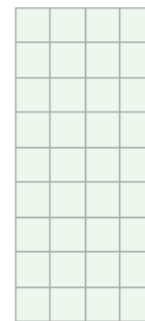
1

2

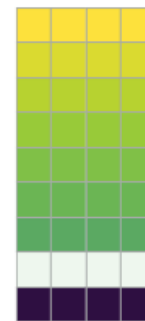
...

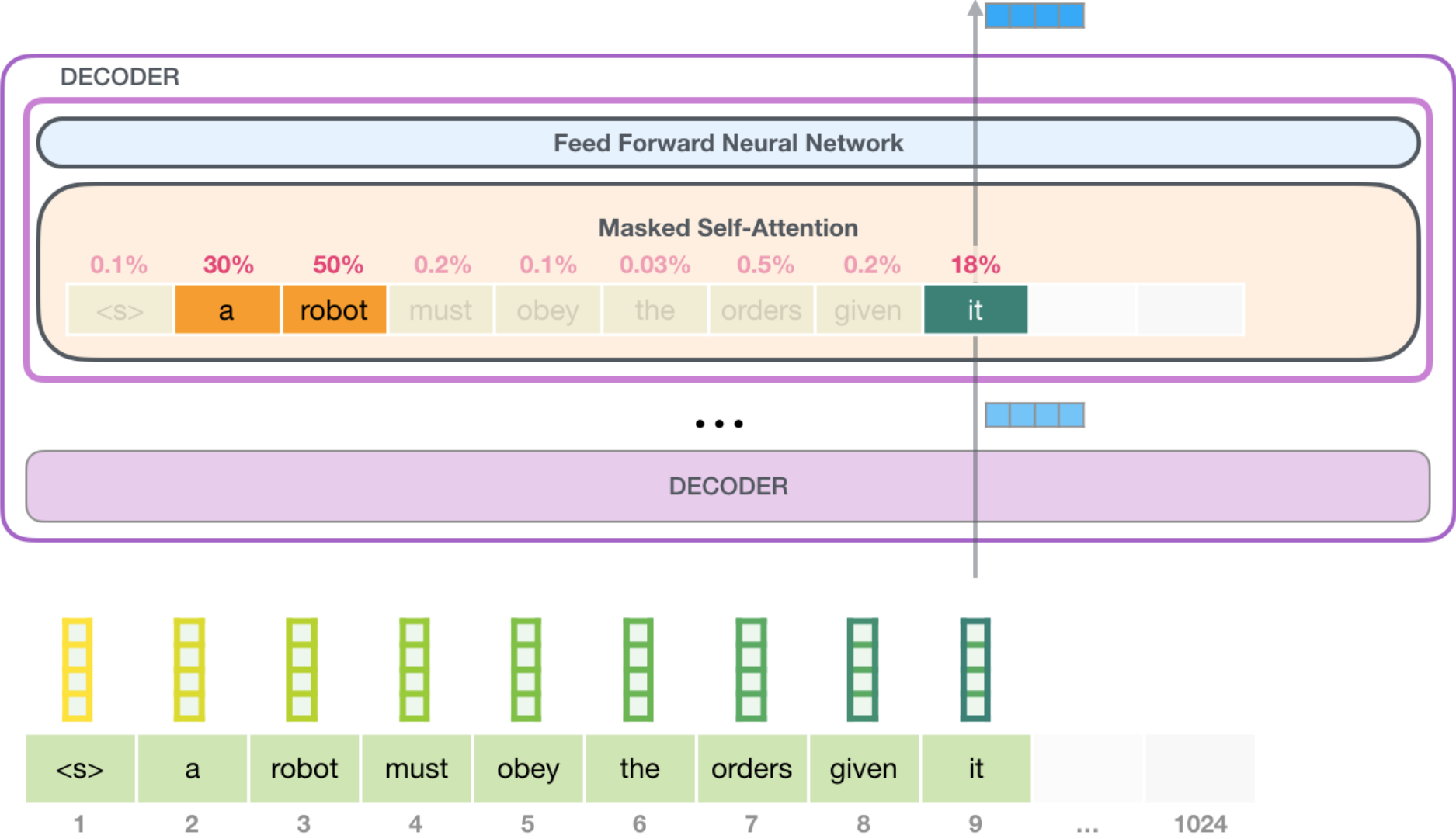
1024

Token
Embeddings



Positional
Encodings



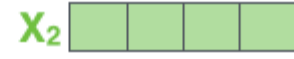
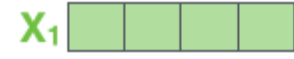


Input

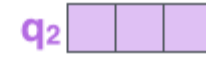
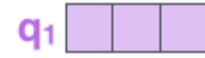
Thinking

Machines

Embedding

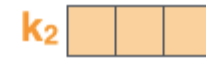
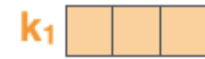


Queries



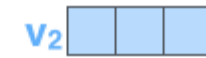
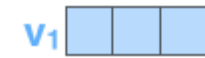
W^Q

Keys



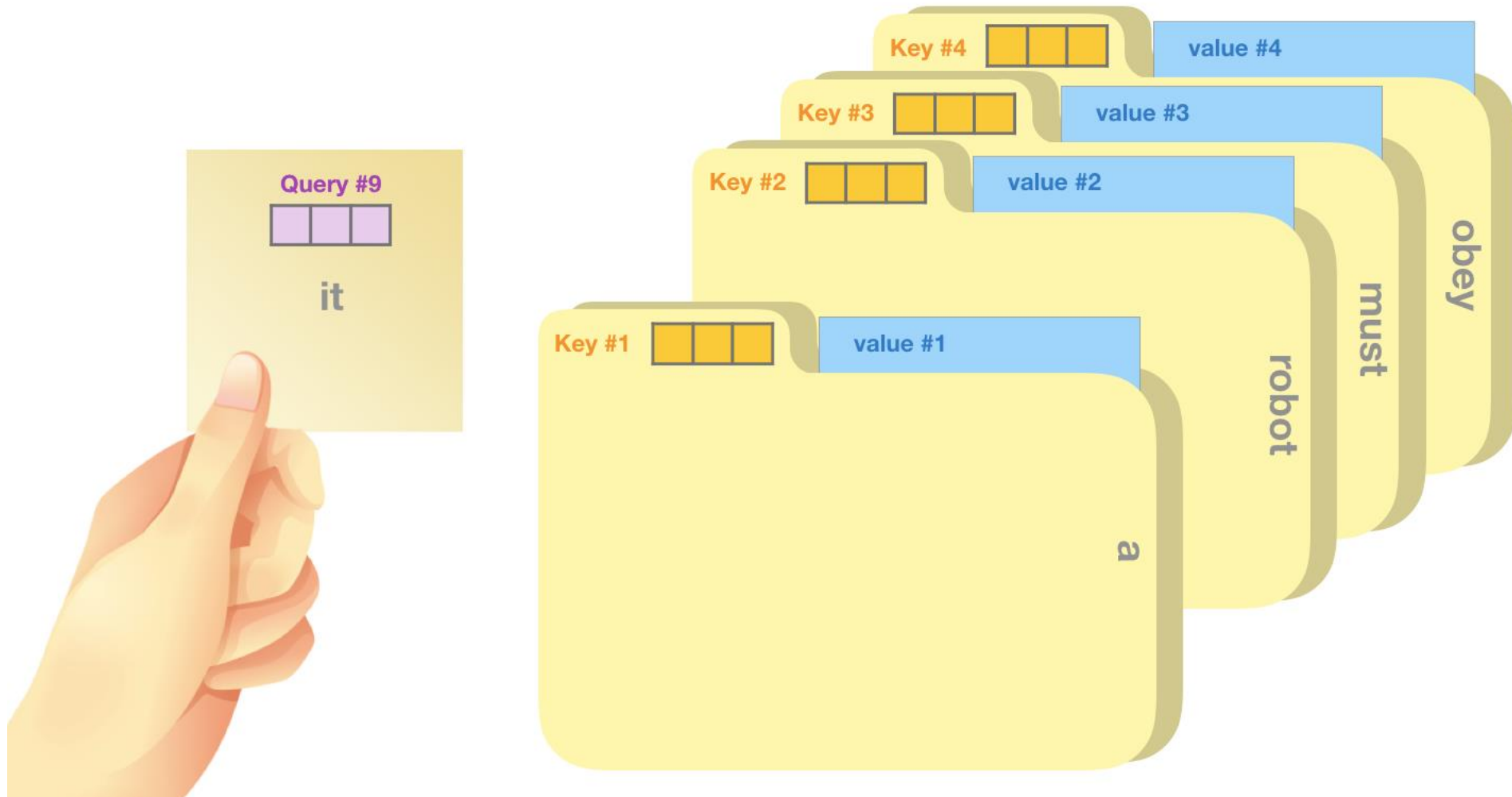
W^K

Values

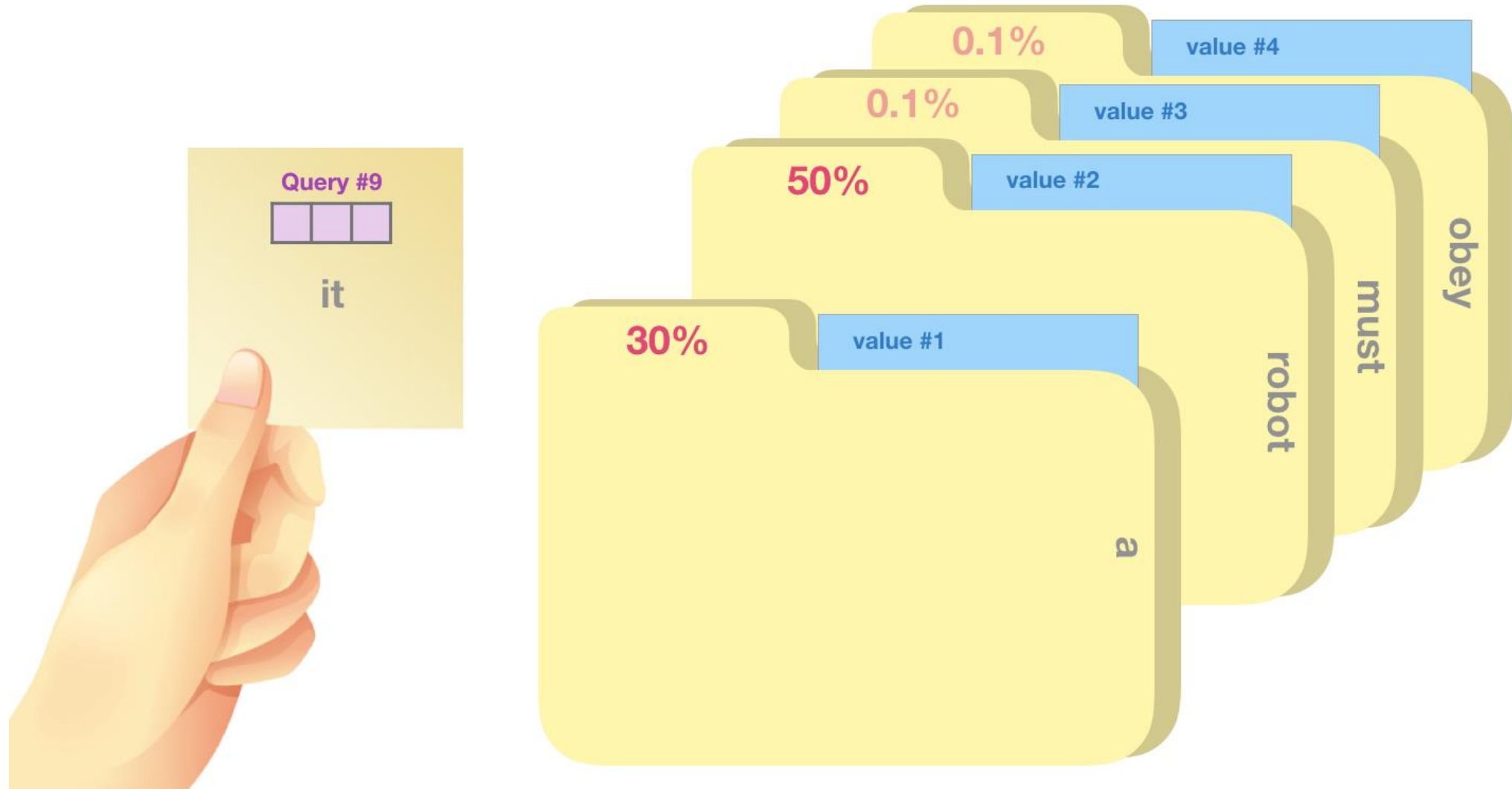



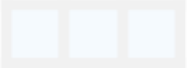





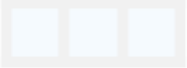



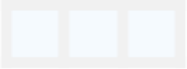







W^V

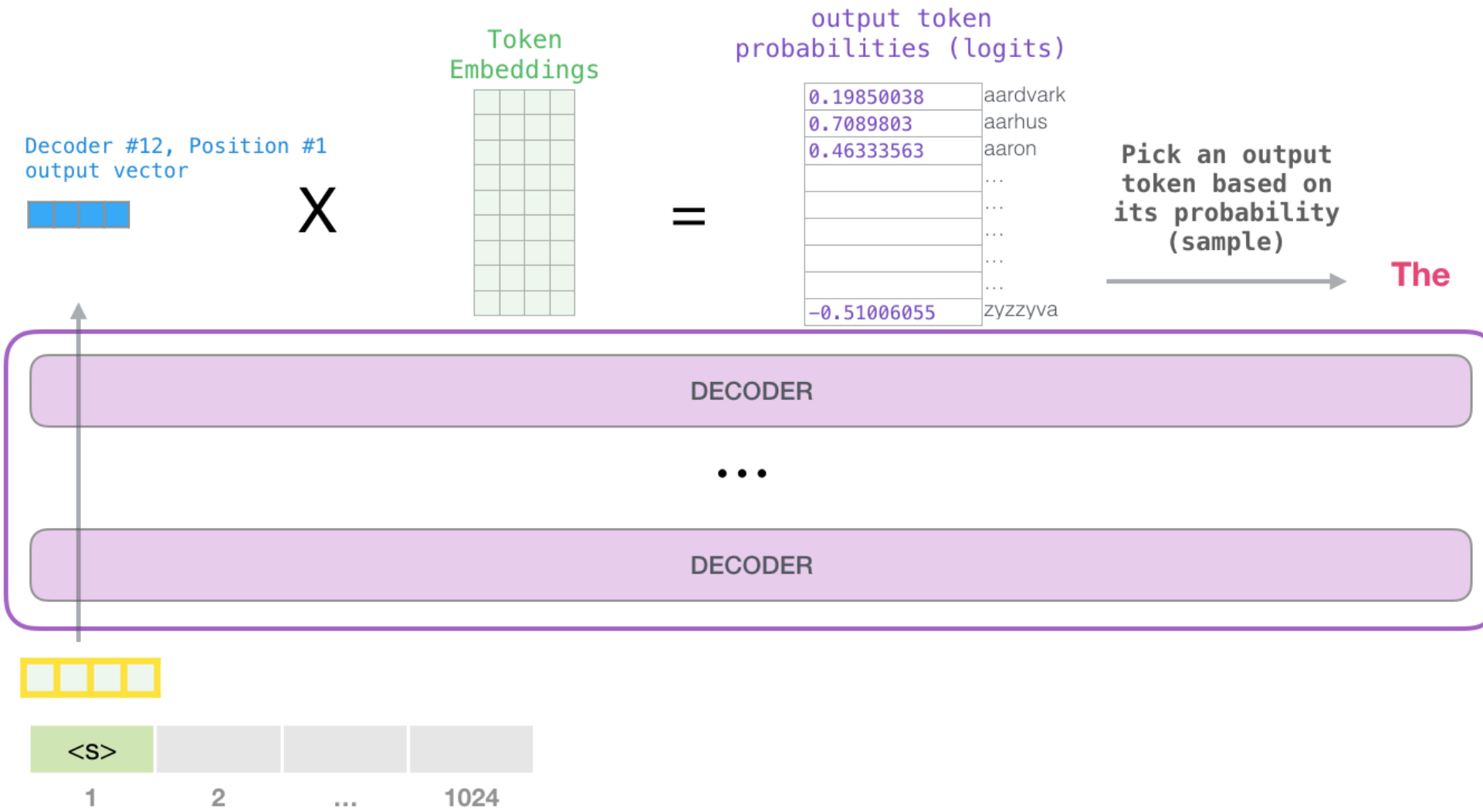
Perform dot product between query and all keys to get a raw score for each previous word (including current word).



Normalize these scores via a softmax to get a probability distribution. Then return a weighted sum of the values.



Word	Value vector	Score	Value X Score
<S>		0.001	
a		0.3	
robot		0.5	
must		0.002	
obey		0.001	
the		0.0003	
orders		0.005	
given		0.002	
it		0.19	
		Sum:	



Math from the paper

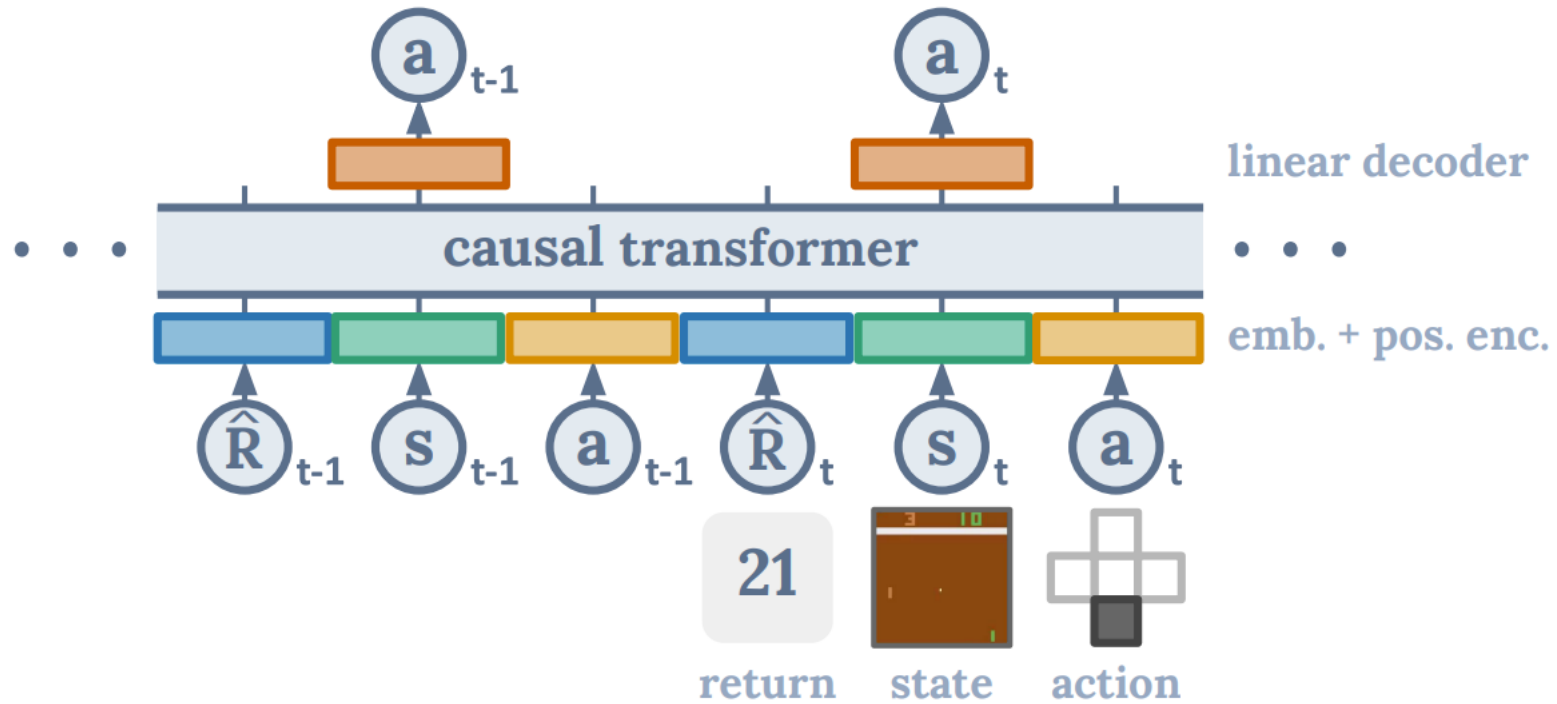
- Given n inputs we create n embeddings x_i
- Map these to queries, keys, and values
- The i th output is then given by

$$z_i = \sum_{j=1}^n \text{softmax}(\{\langle q_i, k_{j'} \rangle\}_{j'=1}^n)_j \cdot v_j$$

We hope this can learn “credit assignment” via attention

Decision Transformer uses a GPT architecture (each token only pays attention to previous tokens)

Decision Transformer



Interesting design decisions

- We want the model to output actions based on desired future returns
- Rather than (s, a, r, s') sequences

$$\tau = \left(\hat{R}_1, s_1, a_1, \hat{R}_2, s_2, a_2, \dots, \hat{R}_T, s_T, a_T \right)$$

- During training we predict actions to generate training loss
- At test time we input desired future return and starting state and then generate the action

Reinforcement Learning Upside Down: Don't Predict Rewards - Just Map Them to Actions

NNAISENSE/IDSIA Technical Report

Jürgen Schmidhuber
The Swiss AI Lab, IDSIA, USI & SUPSI
NNAISENSE, Lugano, Switzerland

23 June 2020 (based on version v1 of 5 Dec 2019)

Earlier drafts: 21 Dec, 31 Dec 2017, 20 Jan, 4 Feb, 9 Mar, 20 Apr, 16 Jul 2018

Questions

- Does the Decision Transformer make any Markov assumptions?
- How is it similar/dissimilar to behavior cloning?
- When might we expect it to perform well/poorly?

DT is closer to:

- behavior cloning + conditioning
- sequence modeling
- generative modeling

DT is *not*:

- value-based RL
- policy improvement via Bellman updates
- exploration-driven learning

Conservative Q-Learning (CQL)

- Core problem with Offline RL

$$\max_a Q(s, a)$$

- Offline data is limited so we often overestimate OOD actions and policy exploits errors -> catastrophic failure at test time
- Core idea: pessimism in the face of uncertainty

$$\mathcal{L}_{\text{CQL}} = \underbrace{\mathcal{L}_{\text{TD}}}_{\text{Bellman error}} + \alpha \left(\underbrace{\log \sum_a \exp(Q(s, a))}_{\text{push all actions down}} - \underbrace{\mathbb{E}_{a \sim \mathcal{D}}[Q(s, a)]}_{\text{pull data actions up}} \right)$$

Implicit Q-Learning (IQL)

- Offline RL dilemma: You want to **improve the policy**, but you **can't trust values of unseen actions**
- Possible solutions
 - Regularize/penalize values of unseen actions (e.g., CQL)
 - Avoid values entirely (e.g., Decision Transformer)
- IQL
 - Do dynamic programming *without ever evaluating unseen actions*

IQL Key Ideas

- Only learn values using dataset actions
- Approximate best action via expectile regression

Mechanism:

1. Learn value function via **upper expectile**:

$$V(s) \approx \text{“high-value actions in dataset”}$$

2. Do TD updates:

$$Q(s, a) \leftarrow r + \gamma V(s')$$

3. Extract policy via **advantage-weighted BC**

Expectile Regression

- Normal regression minimizes $\mathbb{E}[(y - m)^2]$

- Expectile regression minimizes an asymmetric squared loss

$$L_{\tau}(u) = |\tau - \mathbf{1}(u < 0)| u^2$$

where:

- $u = y - m$
- $\tau \in (0, 1)$ controls asymmetry
- If $\tau = 0.5$:
 - symmetric \rightarrow mean
- If $\tau > 0.5$:
 - penalize underestimation more
 - \rightarrow pushes estimate upward
- If $\tau < 0.5$:
 - pushes estimate downward