

<https://introml.mit.edu/>

6.390 Intro to Machine Learning

Lecture 9: Transformers

Shen Shen

Apr 6, 2026

3pm, Room 10-250

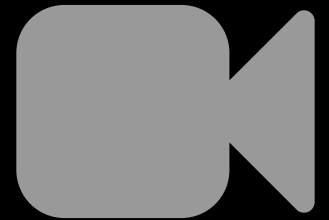
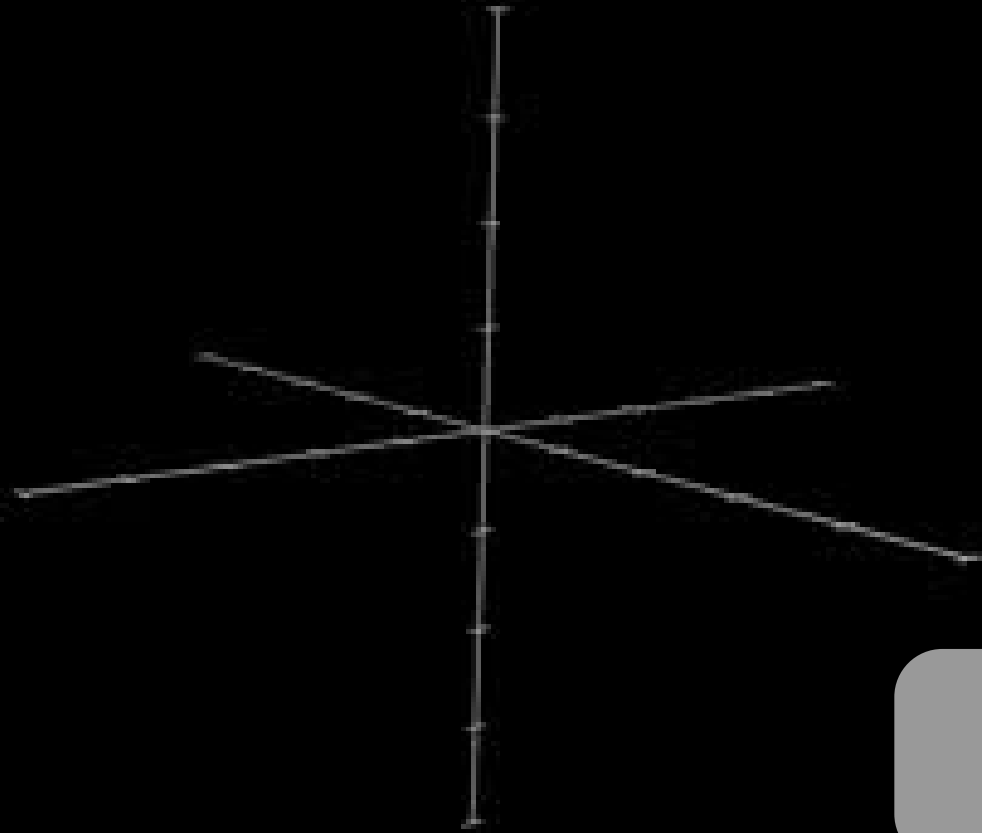
[Slides and Lecture Recording](#)

Recap: Word embedding

Words



Vectors



[video edited from [3b1b](#)]

Good word-embeddings space is equipped with *semantically meaningful vector arithmetic*

this enables "soft" dictionary look-up:

```
1 dict_en2fr = {  
2   "apple" : "pomme",  
3   "banana" : "banane",  
4   "lemon" : "citron"}
```

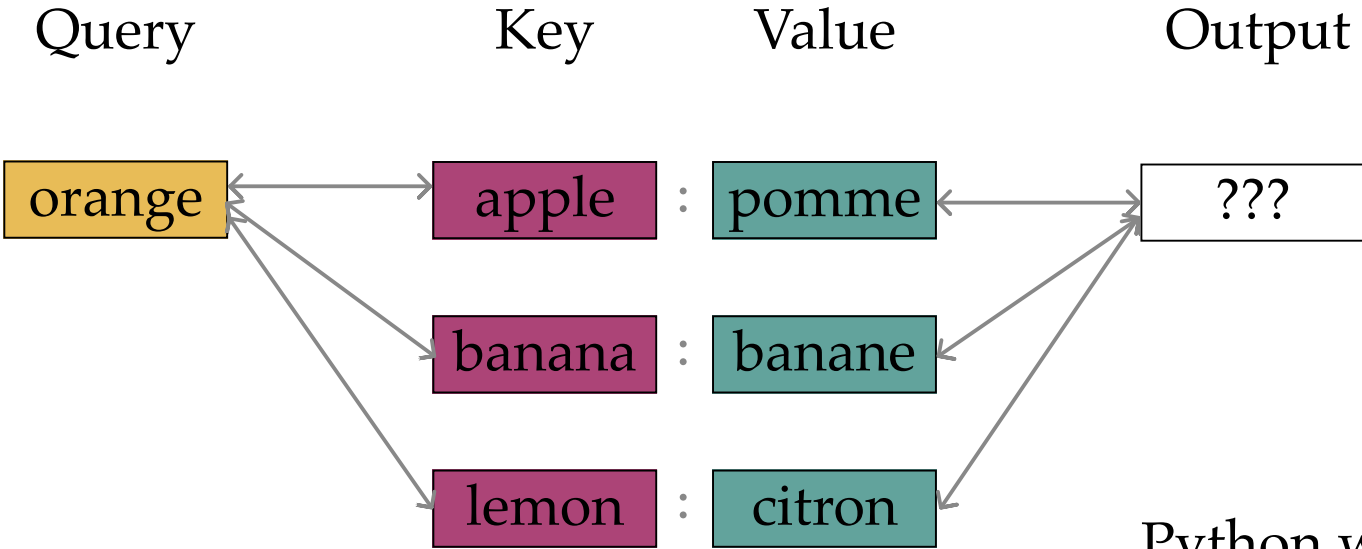
Key Value

apple : pomme

banana : banane

lemon : citron

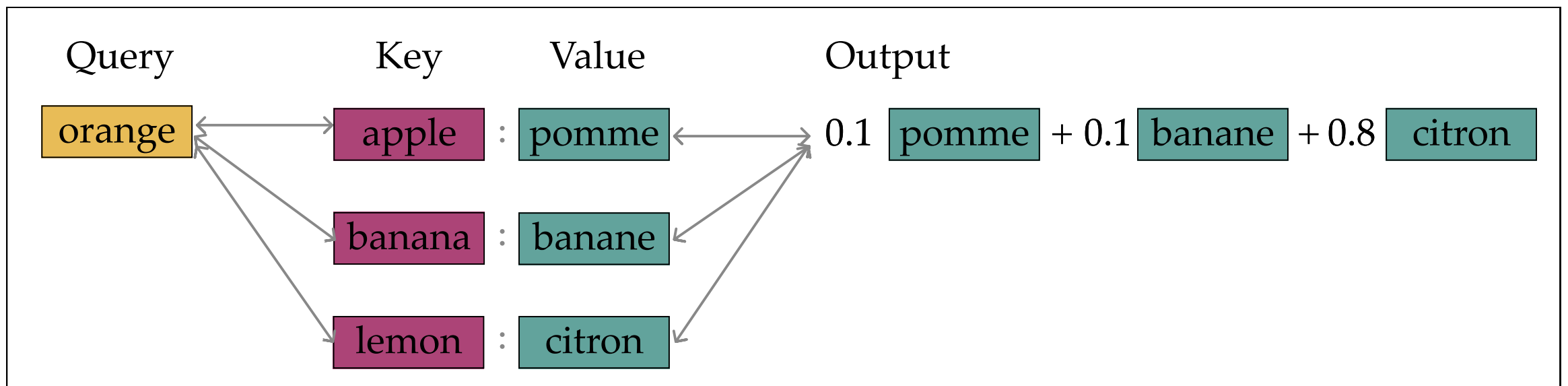
```
1 dict_en2fr = {
2     "apple" : "pomme",
3     "banana" : "banane",
4     "lemon" : "citron"}
5
6 query = "orange"
7 output = dict_en2fr[query]
```

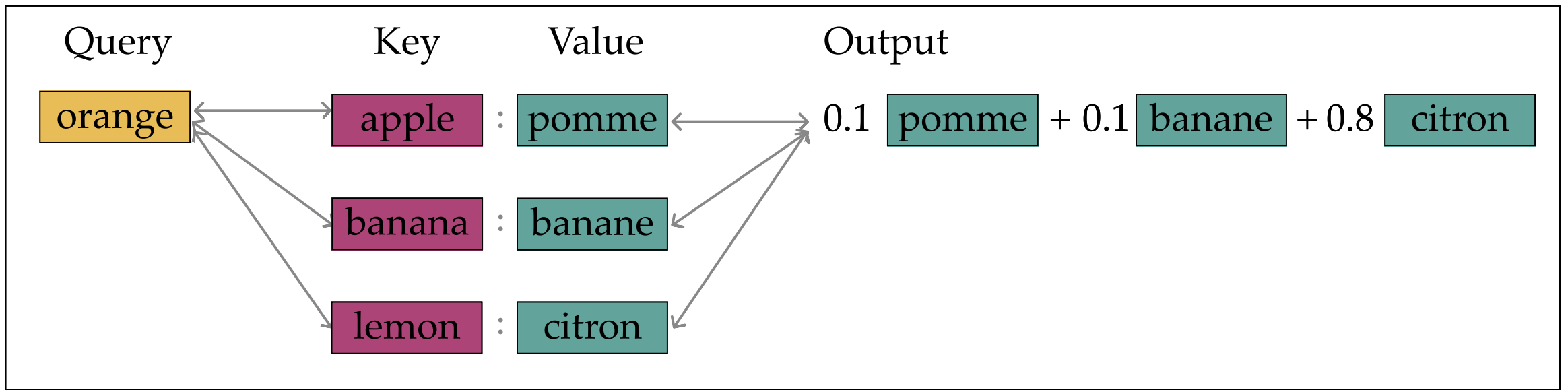


Python would complain. 🤪

```
1 dict_en2fr = {
2     "apple" : "pomme",
3     "banana" : "banane",
4     "lemon" : "citron"}
5
6 query = "orange"
7 output = dict_en2fr[query]
```

But we can probably see the rationale behind something like this:



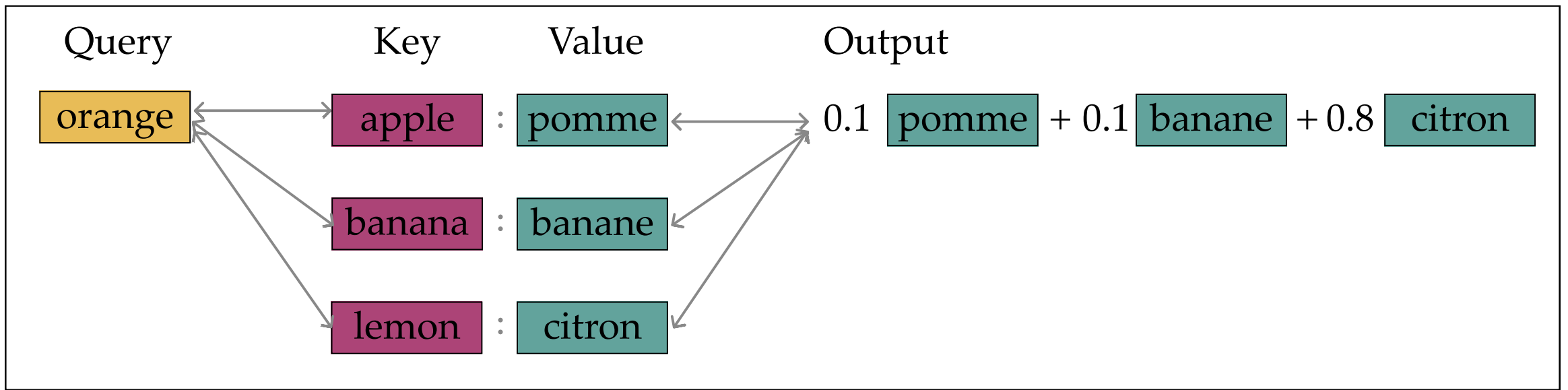


We put (query, key, value) in "good" embeddings in our human brain

such that mixing the values $0.1 \text{ pomme} + 0.1 \text{ banane} + 0.8 \text{ citron}$

via these mixing percentages [0.1 0.1 0.8] made sense

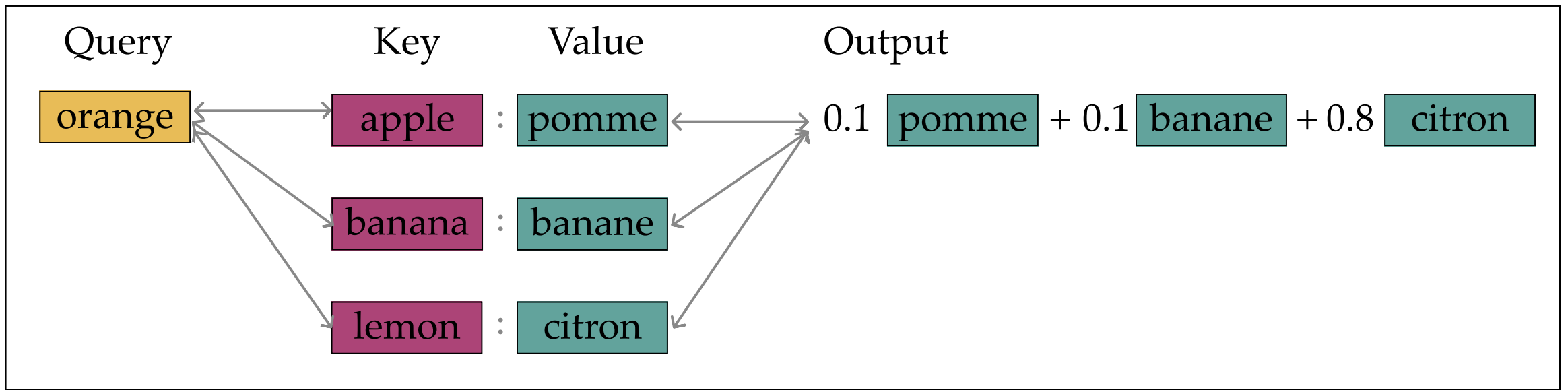
very roughly, the attention mechanism in transformers automates this process.



a. compare query and key for merging percentages:

$$\text{softmax} \left(\begin{array}{c} \text{orange} \text{ apple} \\ \text{orange} \text{ banana} \\ \text{orange} \text{ lemon} \end{array} \right) = [0.1 \ 0.1 \ 0.8]$$

dot-product similarity



a. compare query and key for merging percentages:

$$\text{softmax} \left(\begin{array}{c} \text{orange} \text{ apple} \\ \text{orange} \text{ banana} \\ \text{orange} \text{ lemon} \end{array} \right) = [0.1 \ 0.1 \ 0.8]$$

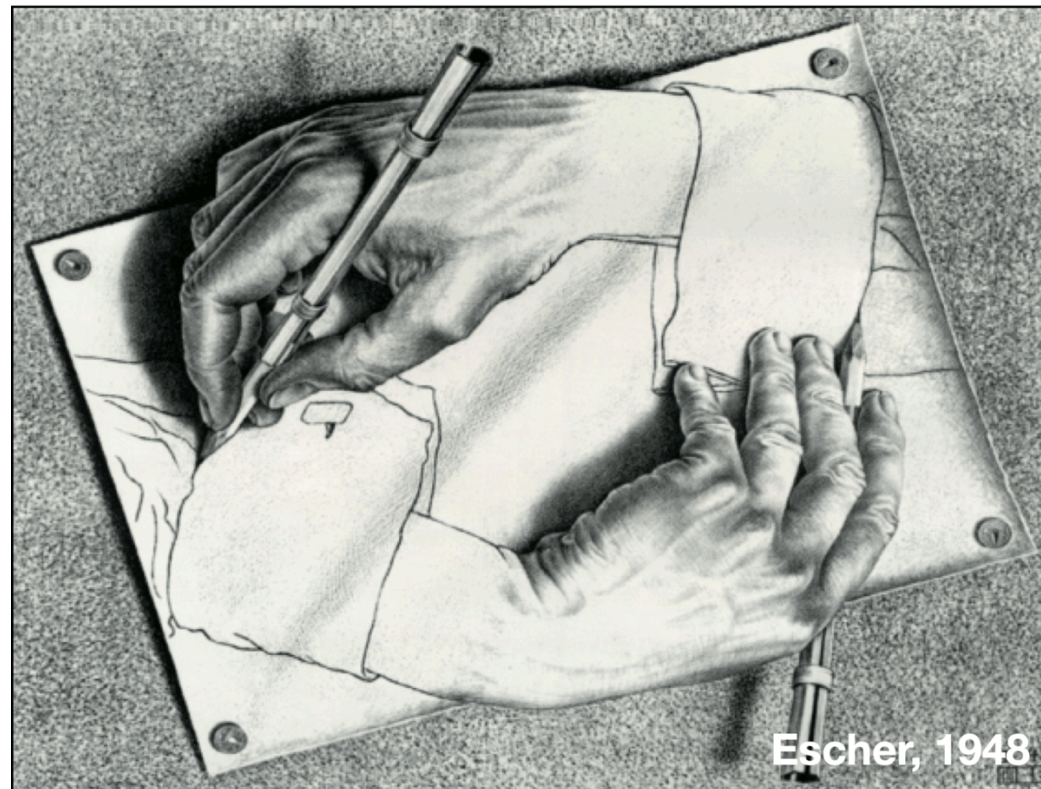
b. then output mixed values 0.1 pomme + 0.1 banane + 0.8 citron

Let's see how this intuition becomes a trainable mechanism.

Outline

- Transformers high-level intuition and architecture
- Attention mechanism
- Multi-head attention
- (Applications)

Large Language Models (LLMs) are trained in this self-supervised way



- Scrape the internet for plain texts.
- Cook up “labels” (prediction targets) from these texts.
- Convert “unsupervised” problem into “supervised” setup.

"To date, the cleverest thinker of all time was Issac. "



feature

label

To date, the

cleverest

To date, the cleverest

thinker

To date, the cleverest thinker

was

⋮

⋮

To date, the cleverest thinker of all time was

Issac

auto-regressive prediction

To date, the _____



model



n

To date, the cleve rest thinker of all time was

???

5.4	7.8	9.7	2.6	3.6	5.6	1.6	9.7	...	3.2	6.7	4.4
7.1	5.2	7.9	7.7	4.3	4.3	1.1	4.6	...	6.6	2.7	8.4
6.0	5.6	4.6	4.5	6.9	9.8	6.5	9.7	...	1.3	7.3	6.9
5.4	9.2	7.7	5.6	0.6	1.0	1.4	6.0	...	7.1	9.5	2.9
4.2	0.7	1.2	0.2	6.6	2.1	1.9	7.3	...	2.9	2.5	8.1
6.4	0.9	6.3	6.1	6.6	1.6	3.7	0.4	...	1.8	5.7	3.9
4.3	0.2	1.4	6.1	2.1	6.5	8.1	2.8	...	5.8	5.9	8.7
8.8	8.2	9.4	6.1	1.3	2.5	1.0	1.2	...	0.2	5.7	5.8
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	...	⋮	⋮	⋮
3.8	8.6	4.1	6.8	3.6	2.4	1.0	1.2	...	0.0	9.4	6.9

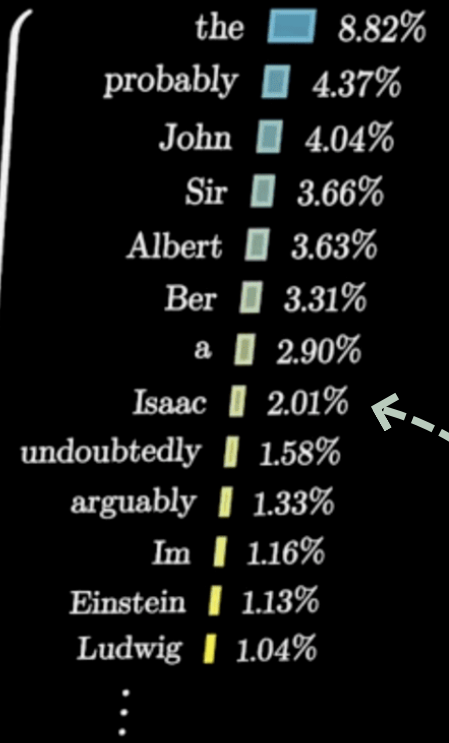
input embedding

To date, the cleverest thinker of all time was

???



word embedding
to a softmax
distribution over
the vocabulary



Minimizing cross-entropy loss drives the weights to assign higher probability to the correct next token

American shrew mole

↓
[0.2
1.3
1.8
7.7
4.1
7
6.1]

↓
[0.4
1.7
1.0
1.6
0.7
5.1]

↓
[1.8
0.9
2.5
1.7
0.1
⋮
2.1]

One mole of carbon dioxide

↓
[1.2
7.4
2.7
0.8
0.8
1
2.7]

↓
[0.8
0.9
2.5
1.7
0.1
⋮
2.1]

↓
[2.3
7.8
1.8
0.1
1.1
1.8]

↓
[7.4
1.4
1.7
0.1
1.8]

↓
[0.8
1.4
0.1
0.8
0.1
1
0.1]

Take a biopsy of the mole

↓
[0.9
1.1
0.7
0.6
0.0
1
1.1]

↓
[1.5
0.7
1.0
0.3
0.7
1
0.9]

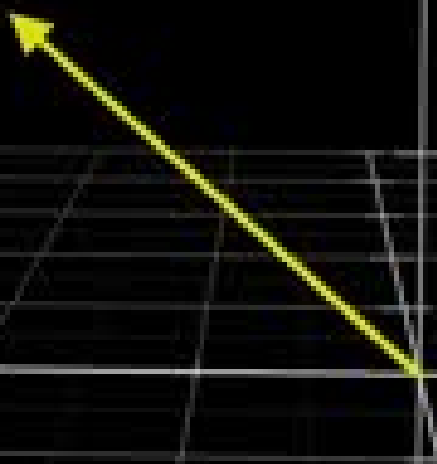
↓
[1.1
0.7
1.4
1.7
0.7
1
0.1]

↓
[0.8
7.8
1.8
0.1
1.1
1.8]

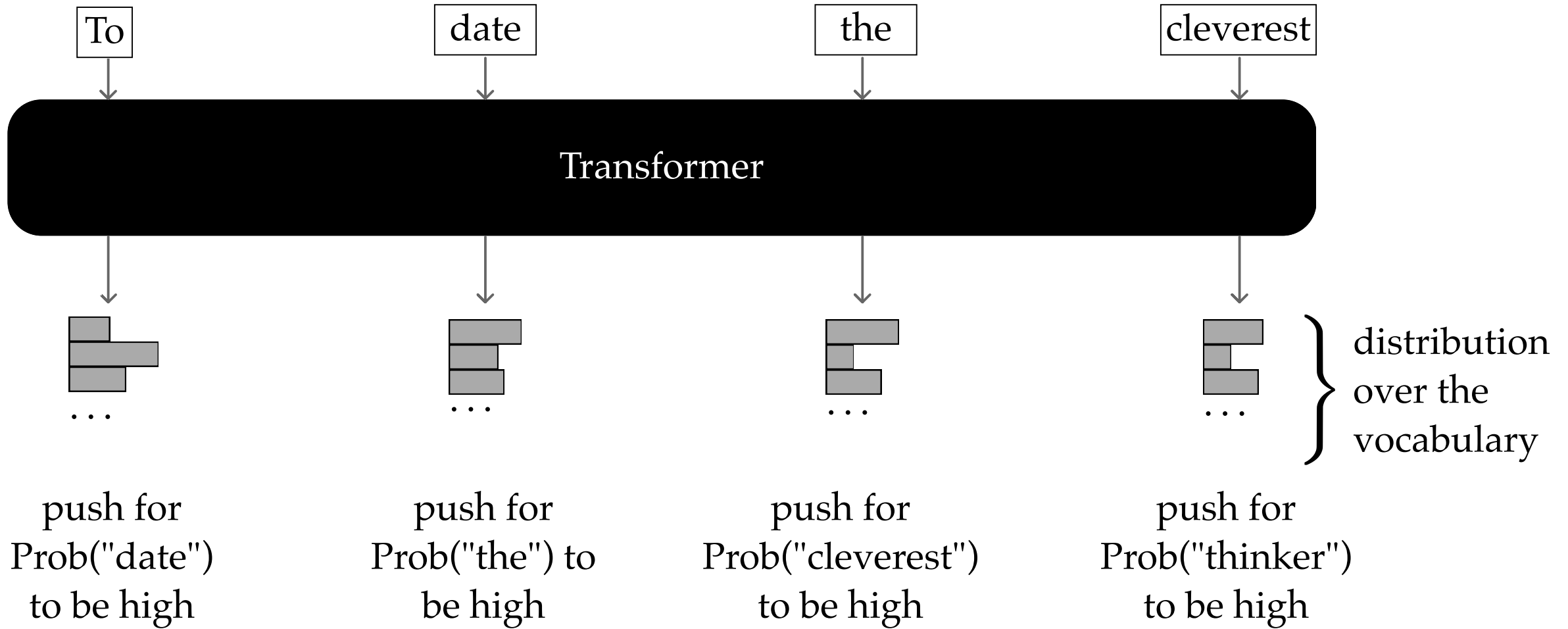
↓
[1.7
0.1
0.1
1
0.1]

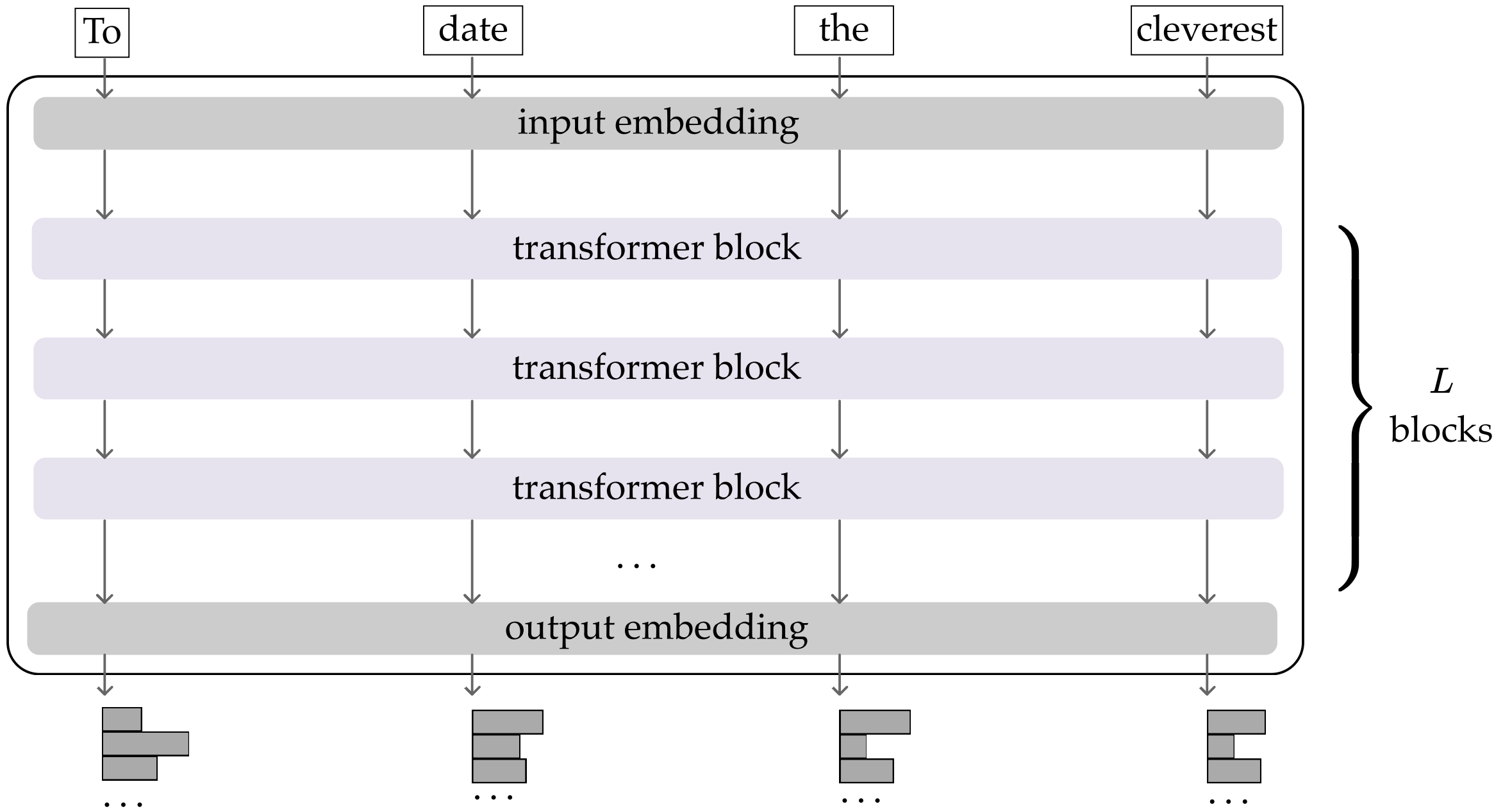
↓
[1.8
0.9
2.5
1.7
0.1
⋮
2.1]

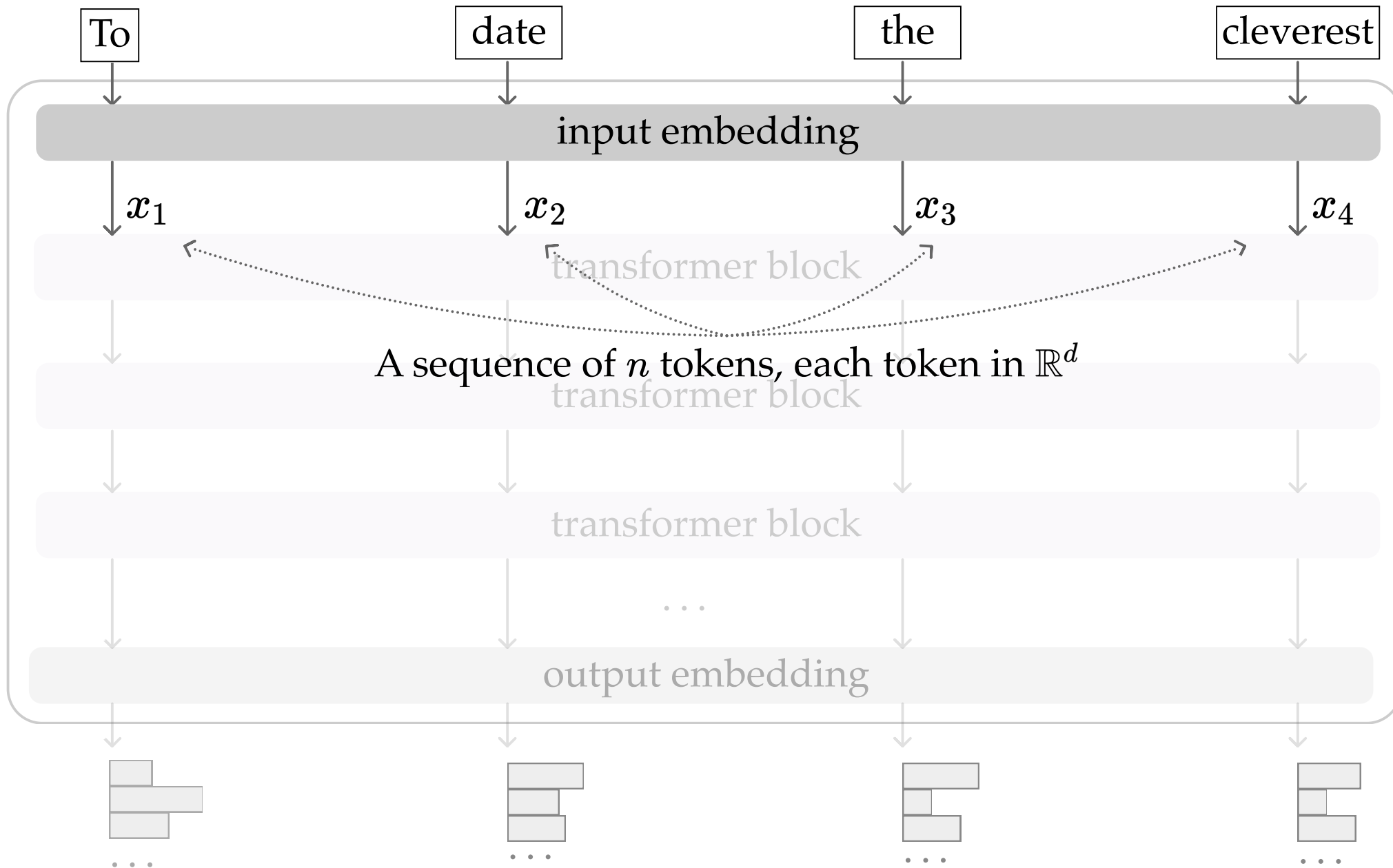
E(mole)

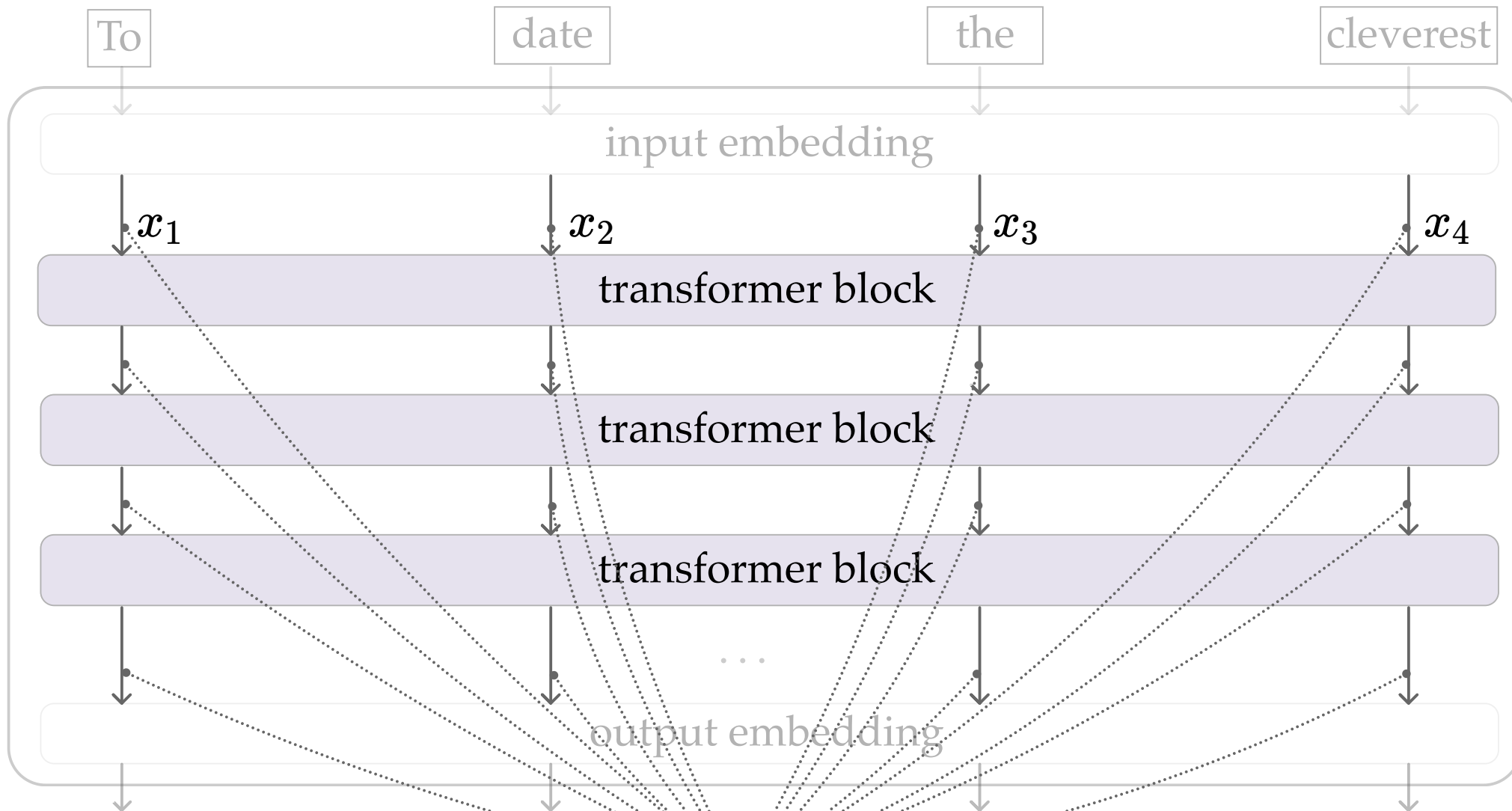


"To date, the cleverest [thinker] of all time was Issac.

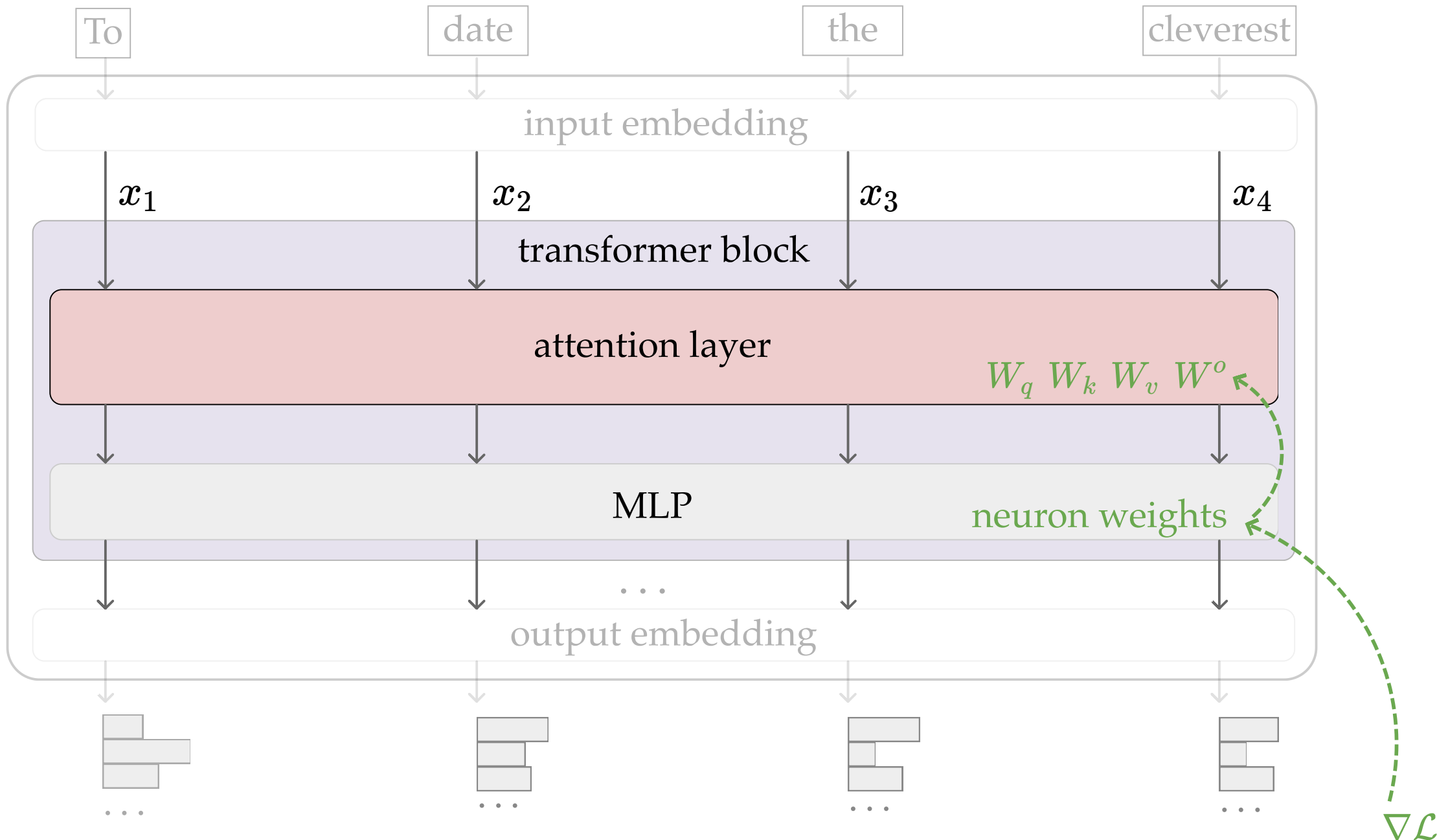


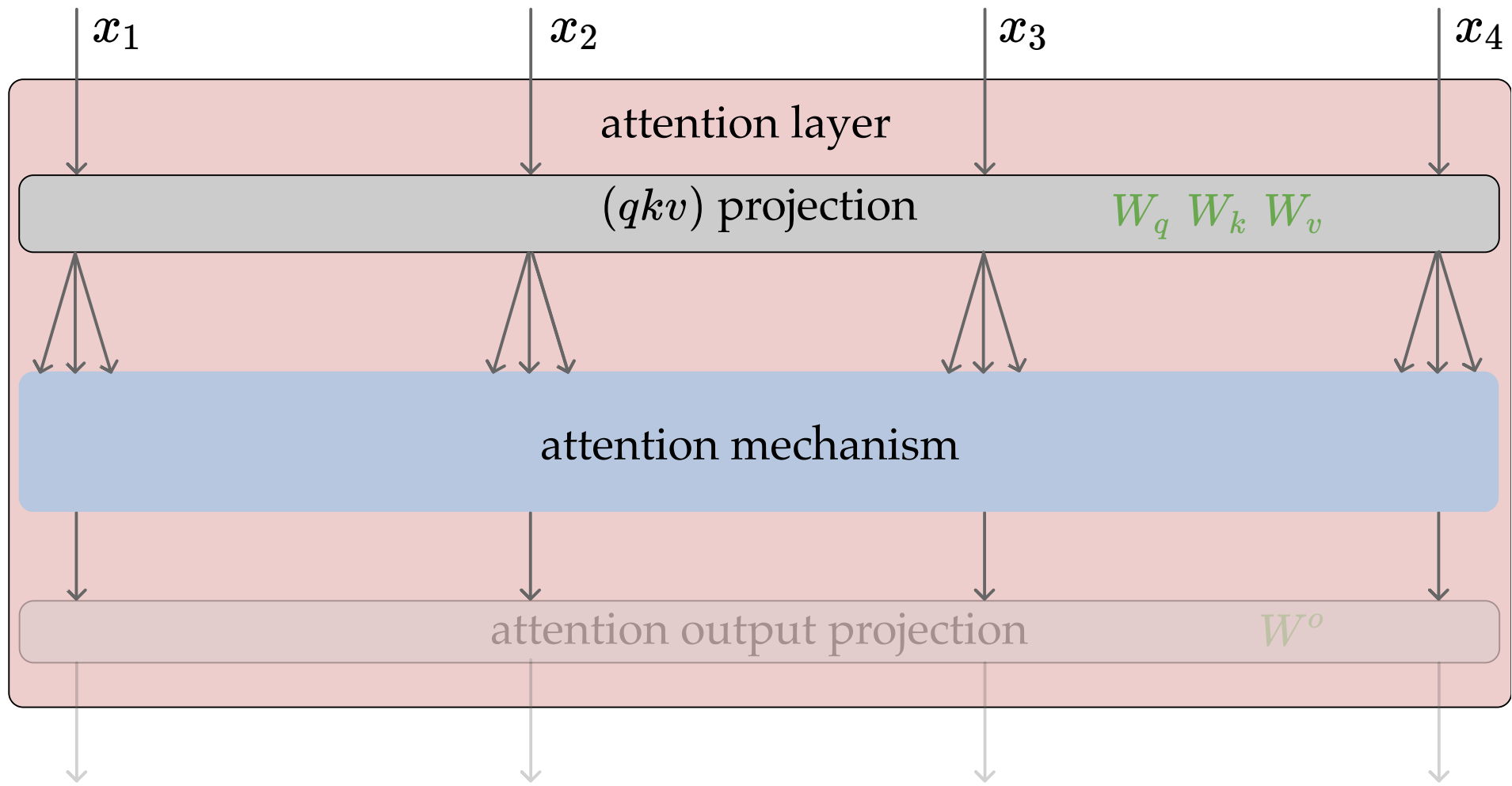






each of the n tokens transformed, block by block
 within a shared d -dimensional word-embedding space.





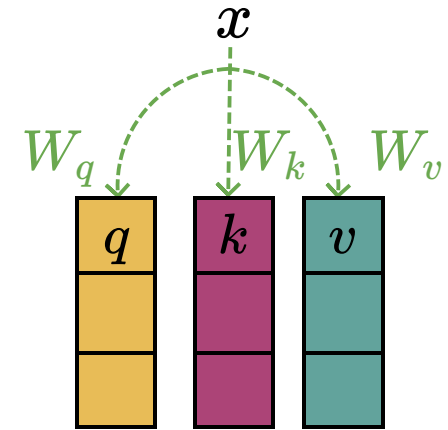
Most important bits in an attention layer:

1. (query, key, value) projection
2. attention mechanism

1. (query, key, value) projection

With *learned* projections, we frame x into:

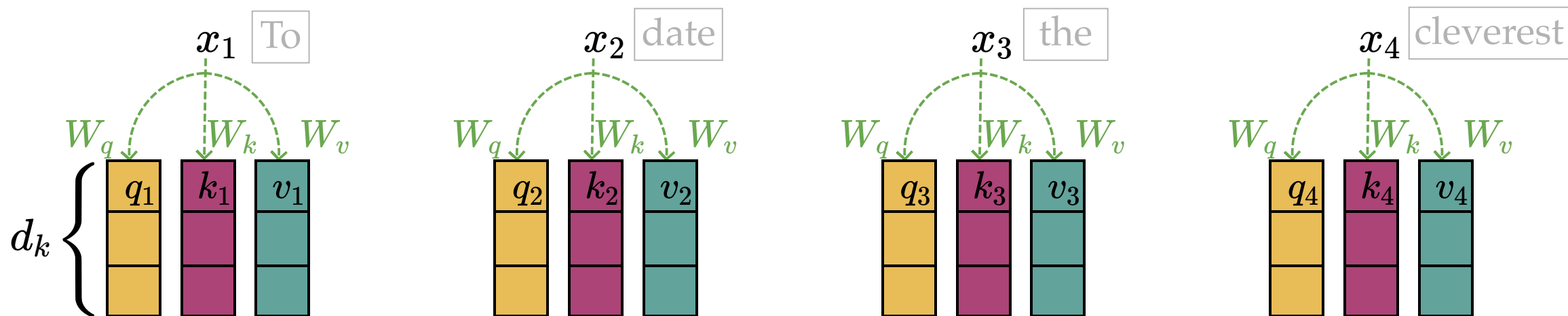
- a query to be the questions
- a key to be compared
- a value to contribute



Why learning these projections:

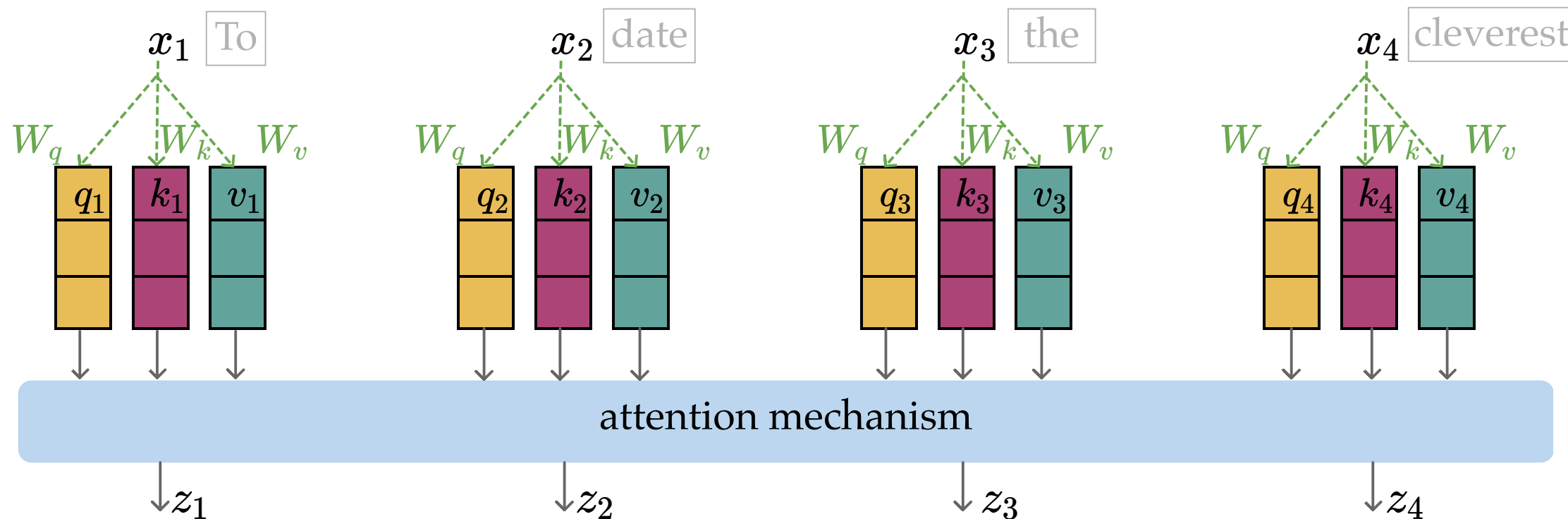
- W_q learns *how* to ask
- W_k learns *how* to listen
- W_v learns *how* to speak

1. (query, key, value) projection



- W_q, W_k, W_v , all in $\mathbb{R}^{d \times d_k}$
- project word embedding x from d -dimensional space to d_k -dimensional (q, k, v) spaces (typically $d_k < d$)
- $q_i = W_q^T x_i, k_i = W_k^T x_i, v_i = W_v^T x_i, \forall i$ — *weight sharing* across positions
- *parallel* and *structurally identical* processing

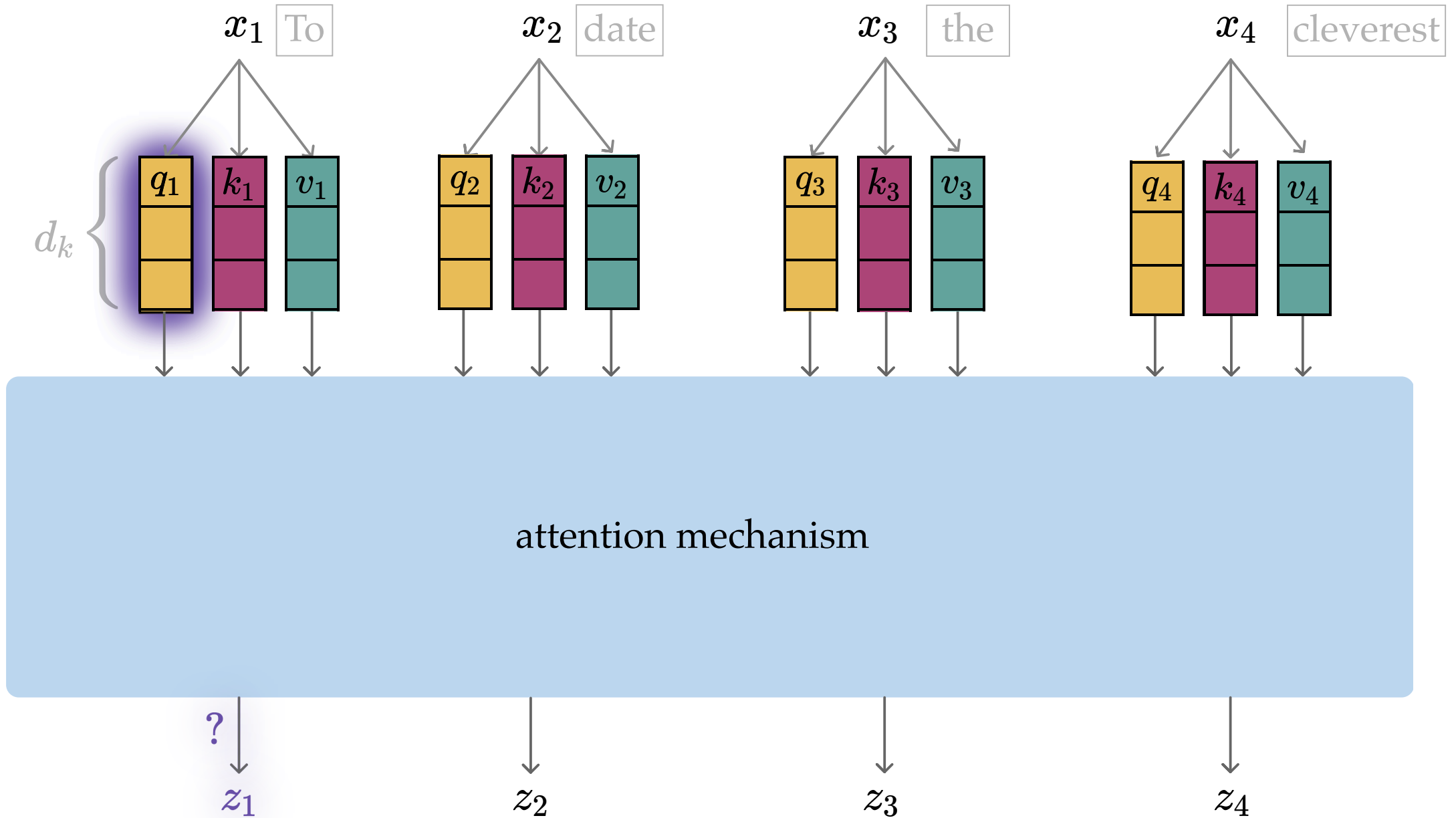
2. Attention mechanism

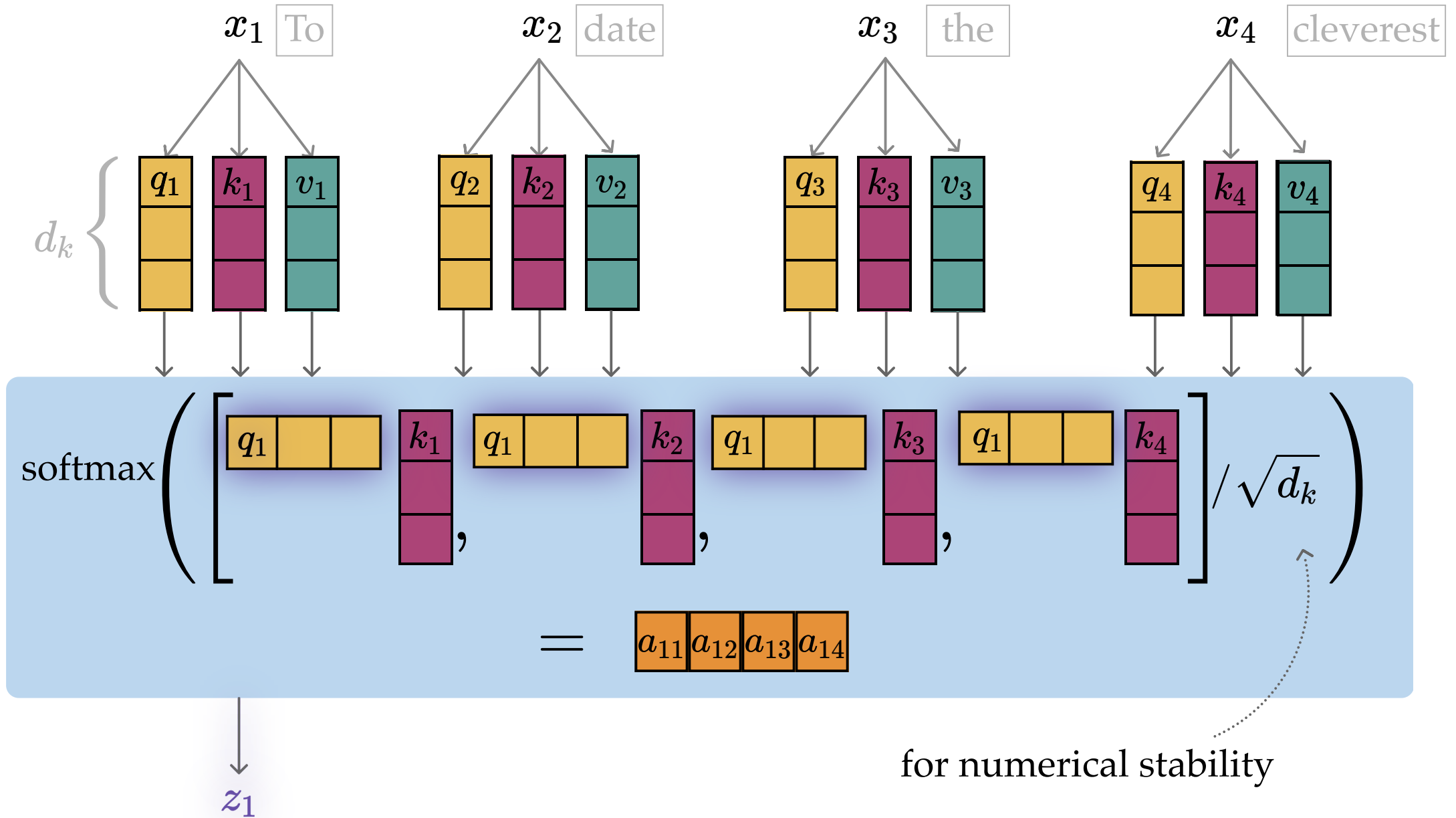


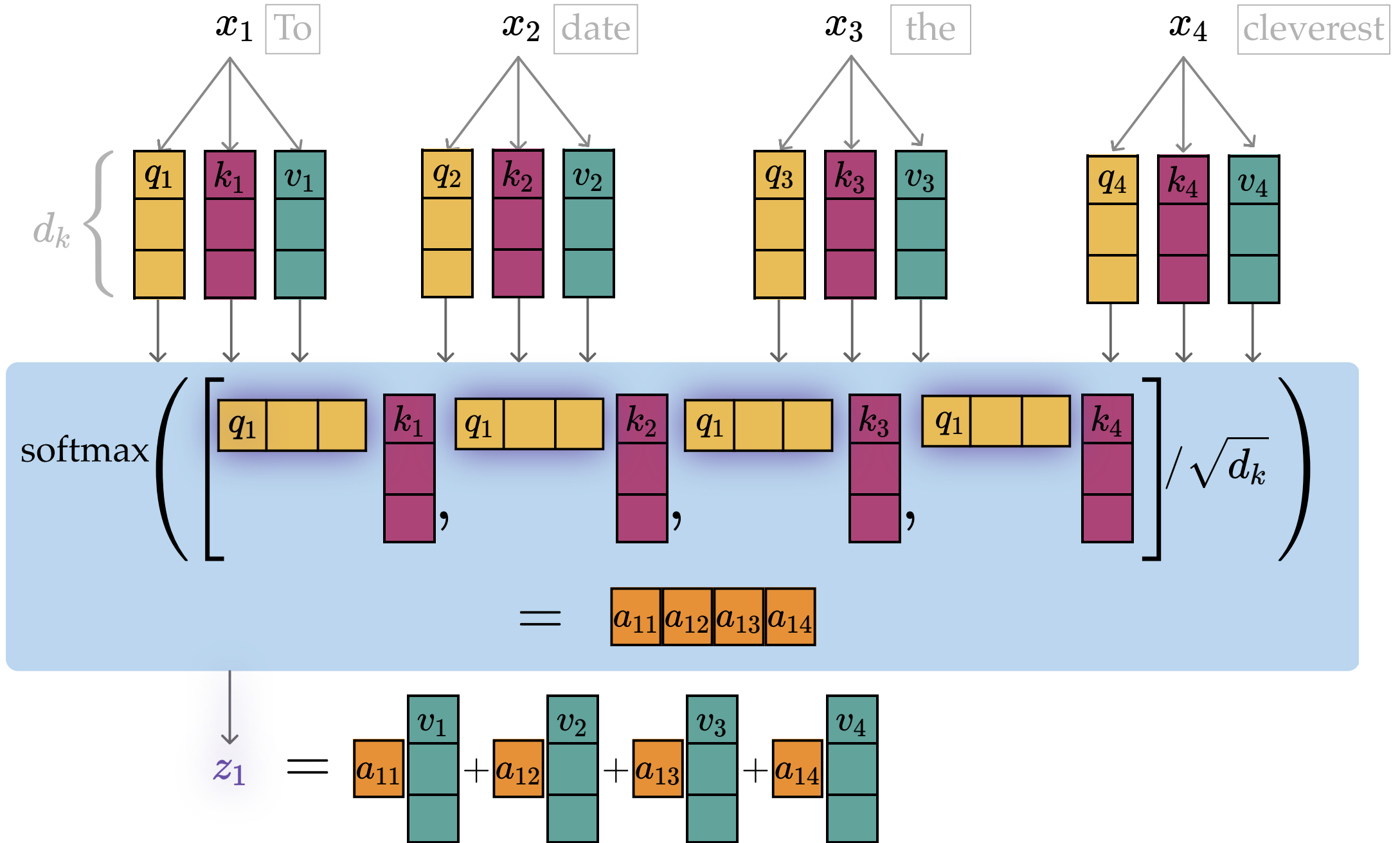
- Attention mechanism turns the projected (q, k, v) into z
- Each z is context-aware: a mixture of everyone's values, weighted by relevance

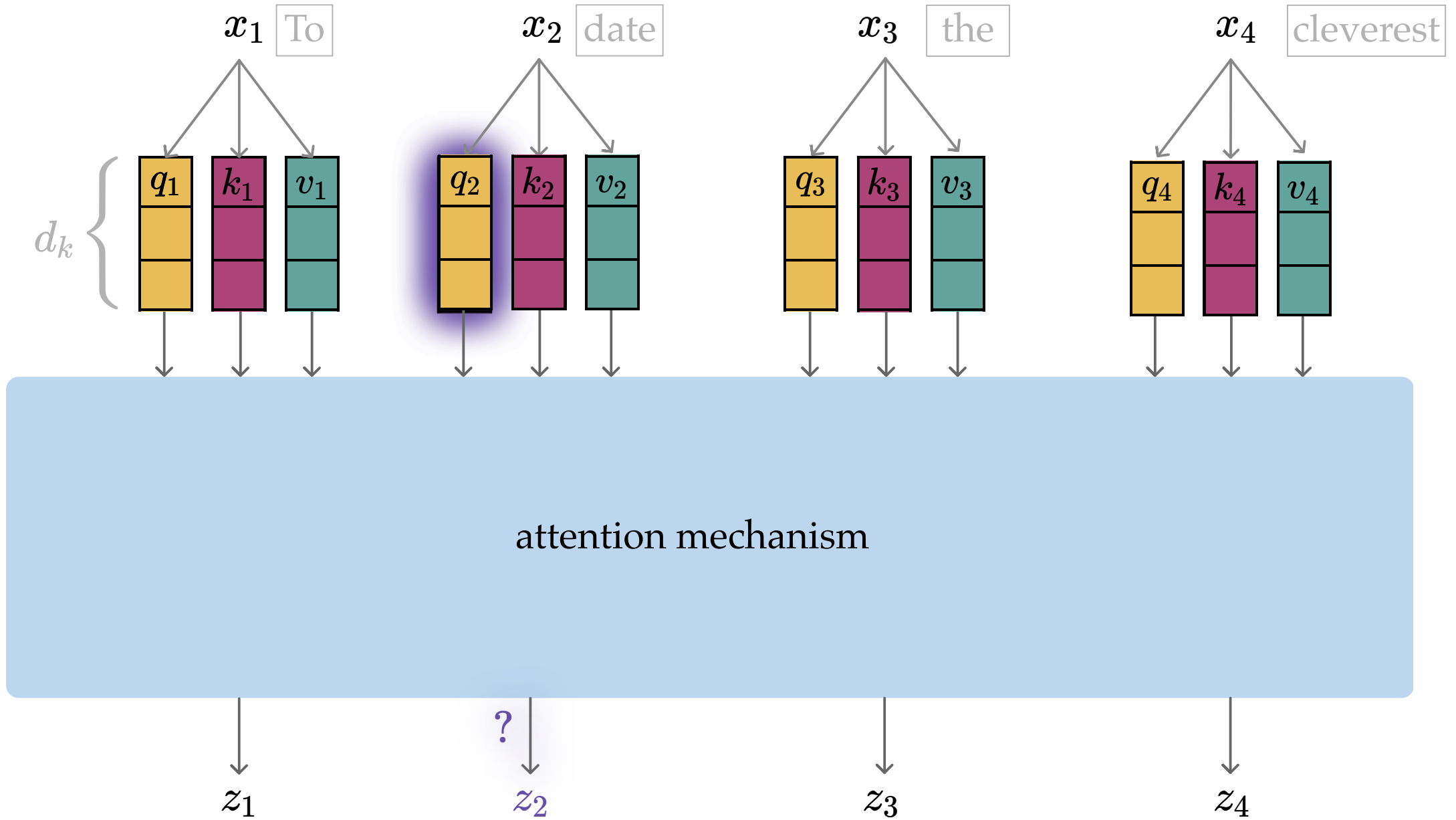
Outline

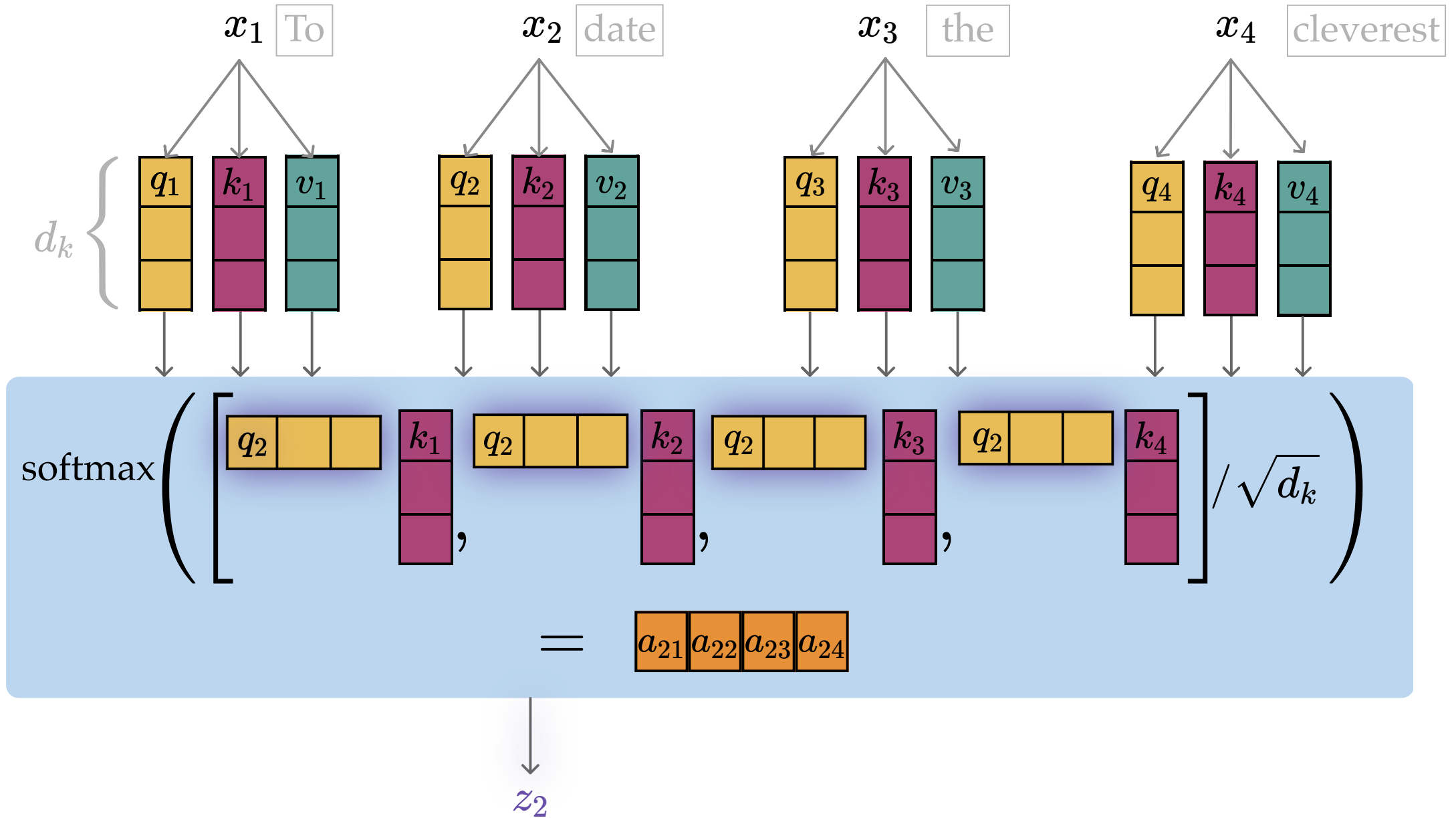
- Transformers high-level intuition and architecture
- Attention mechanism
- Multi-head attention
- (Applications)

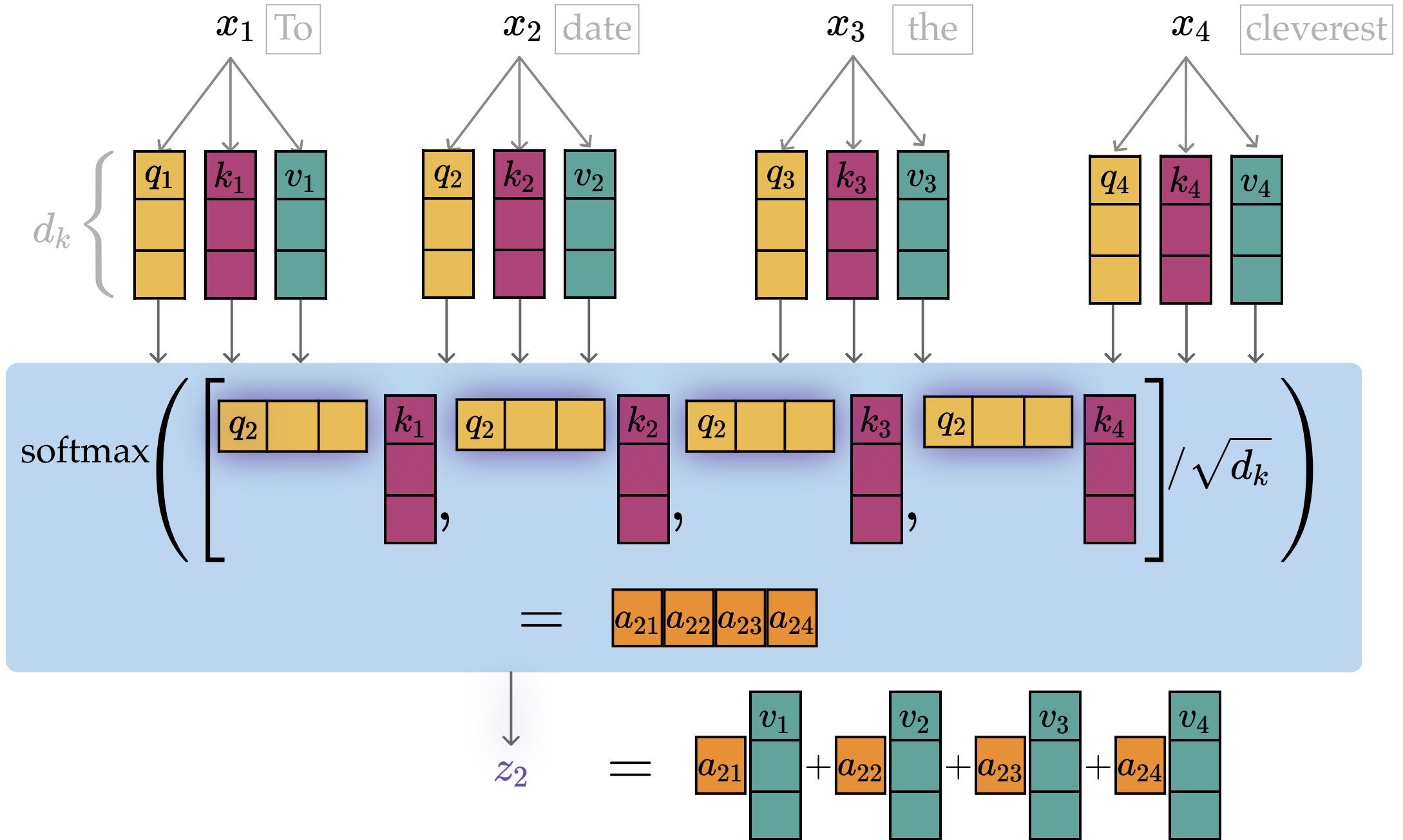


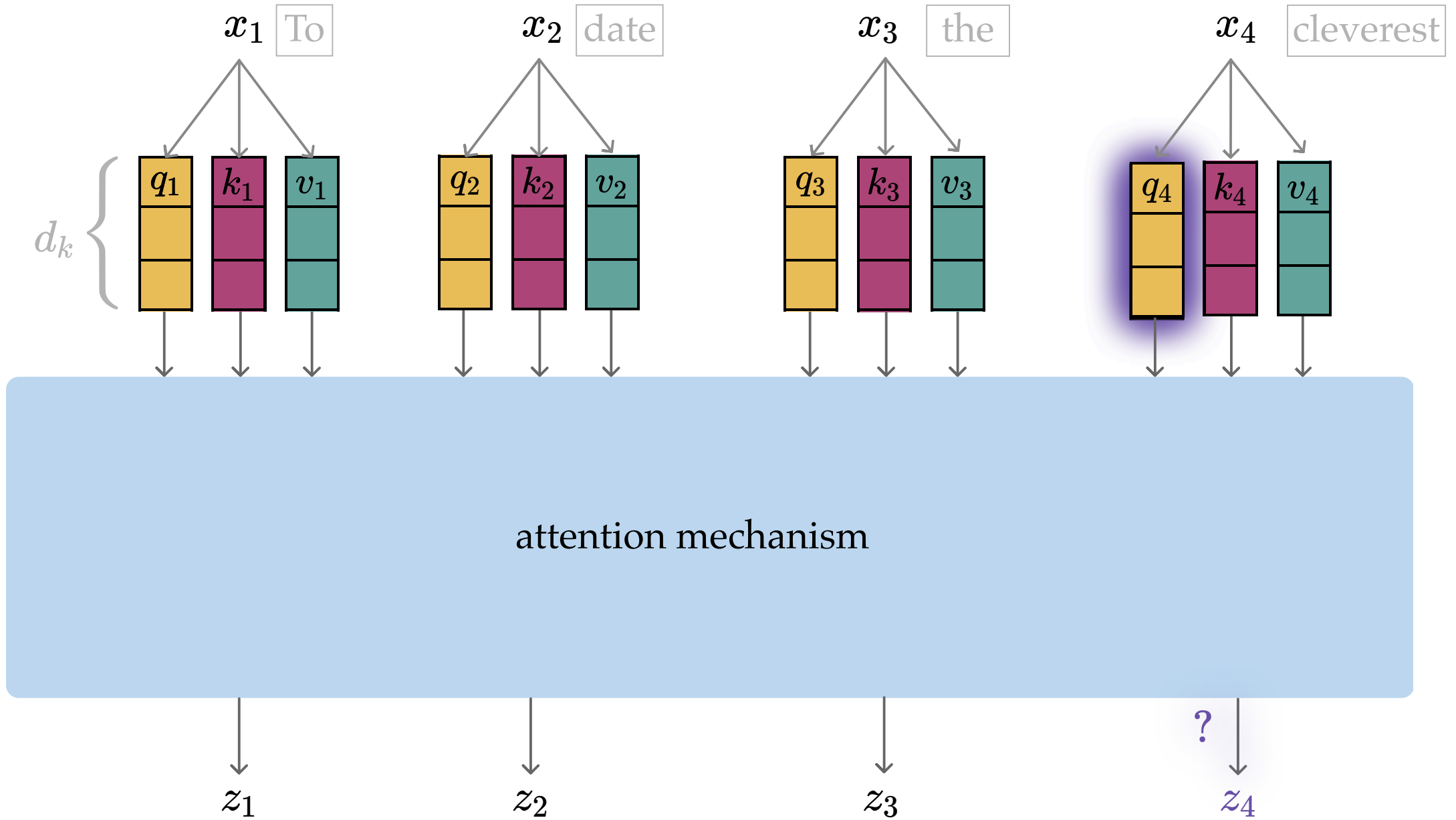


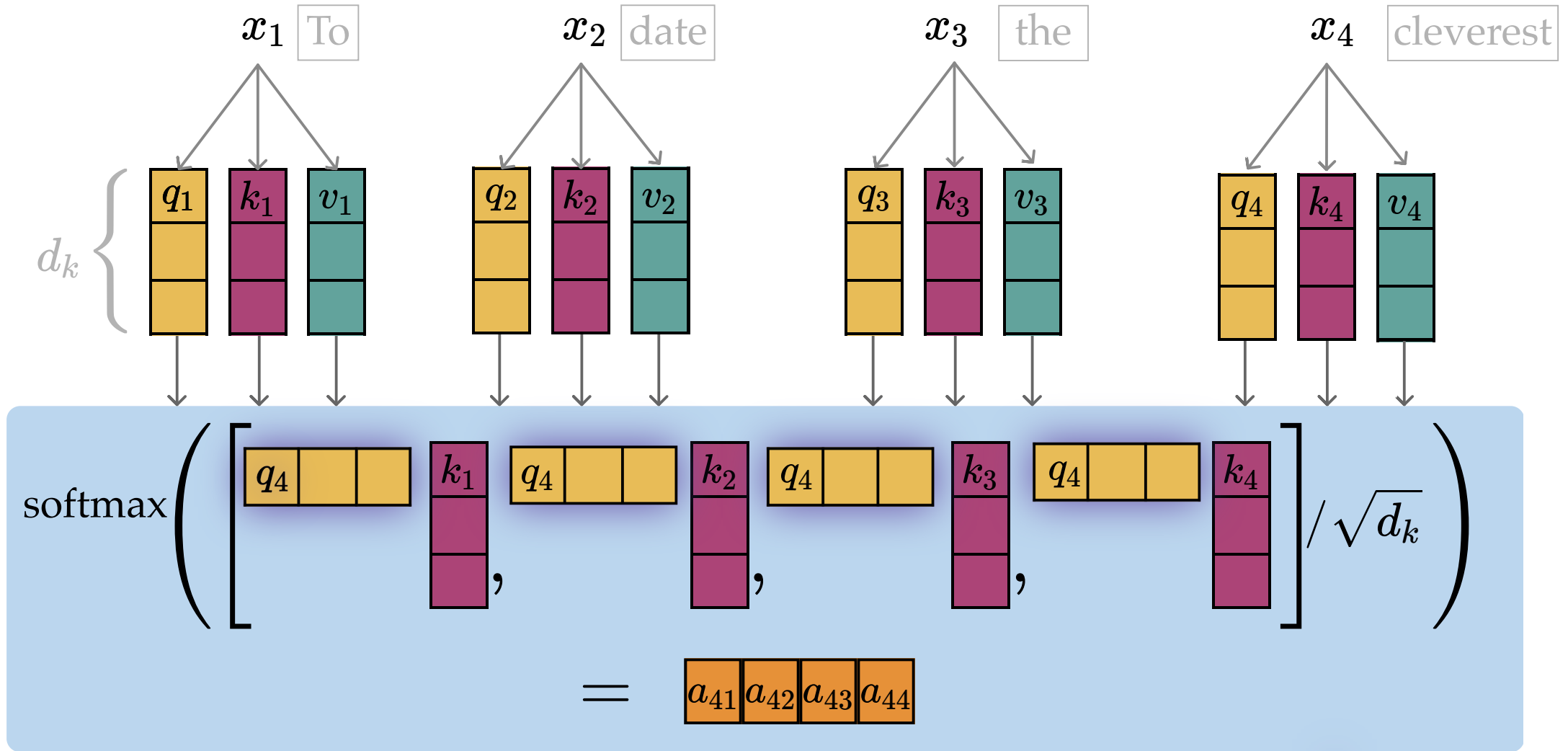




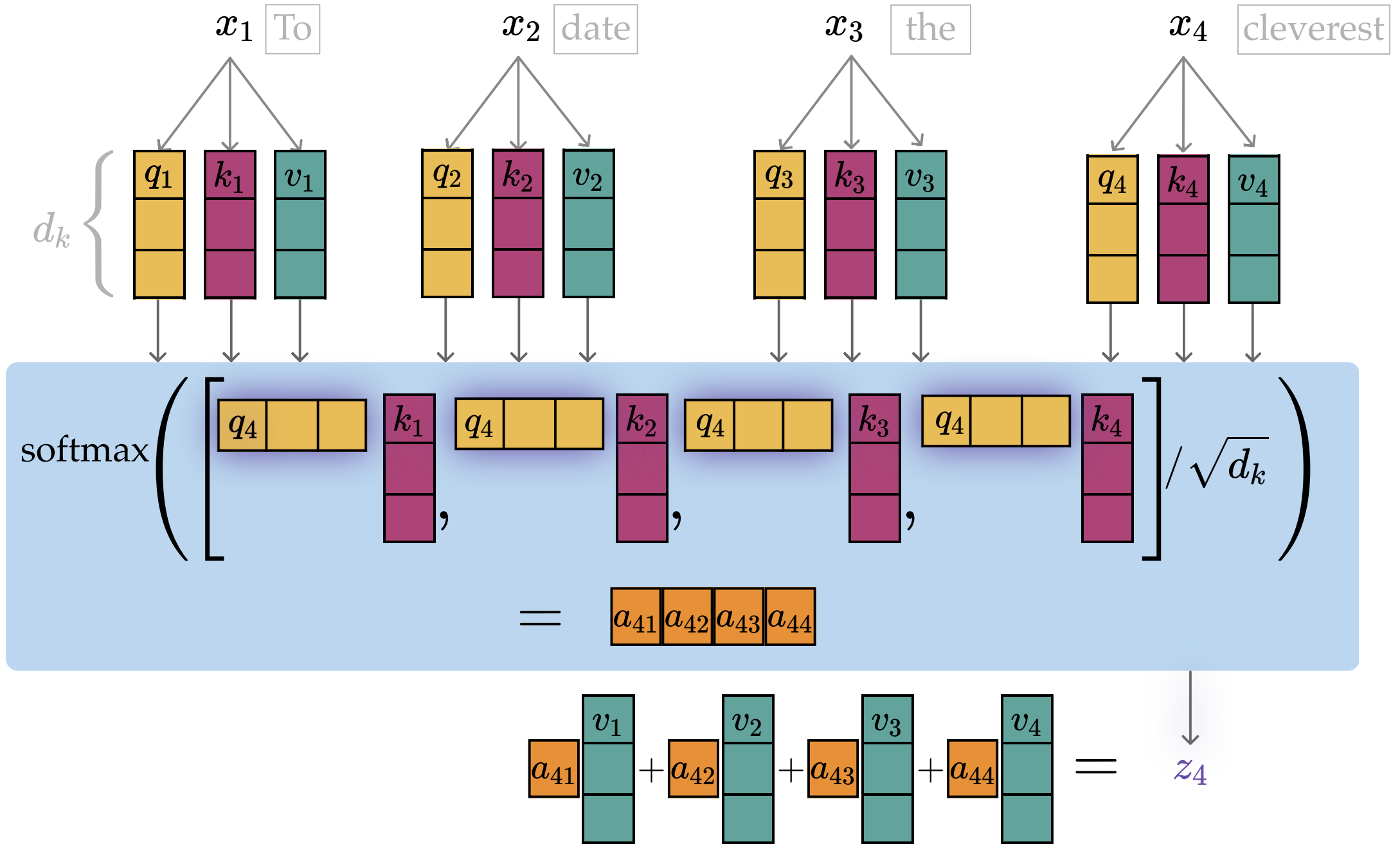


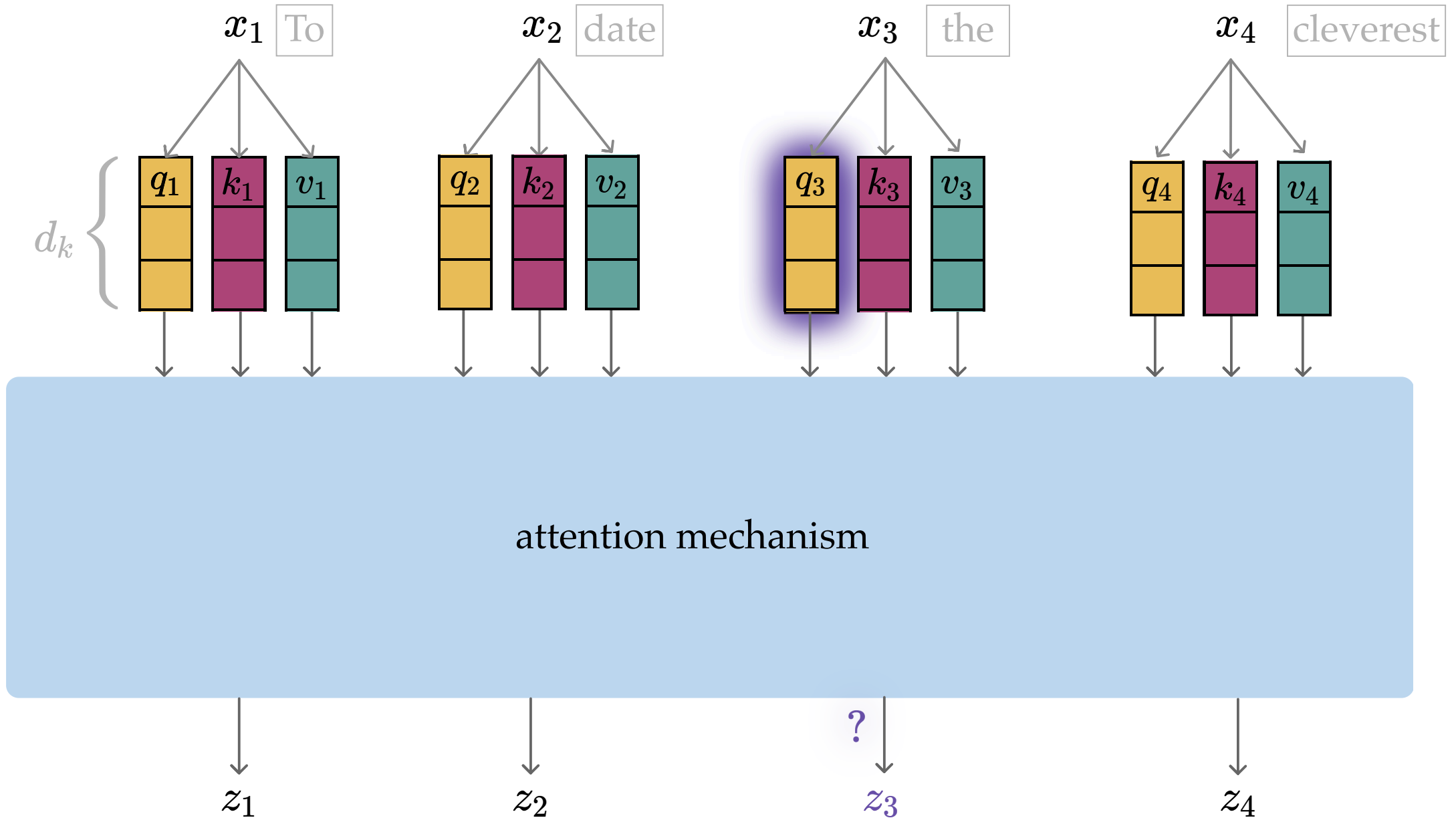


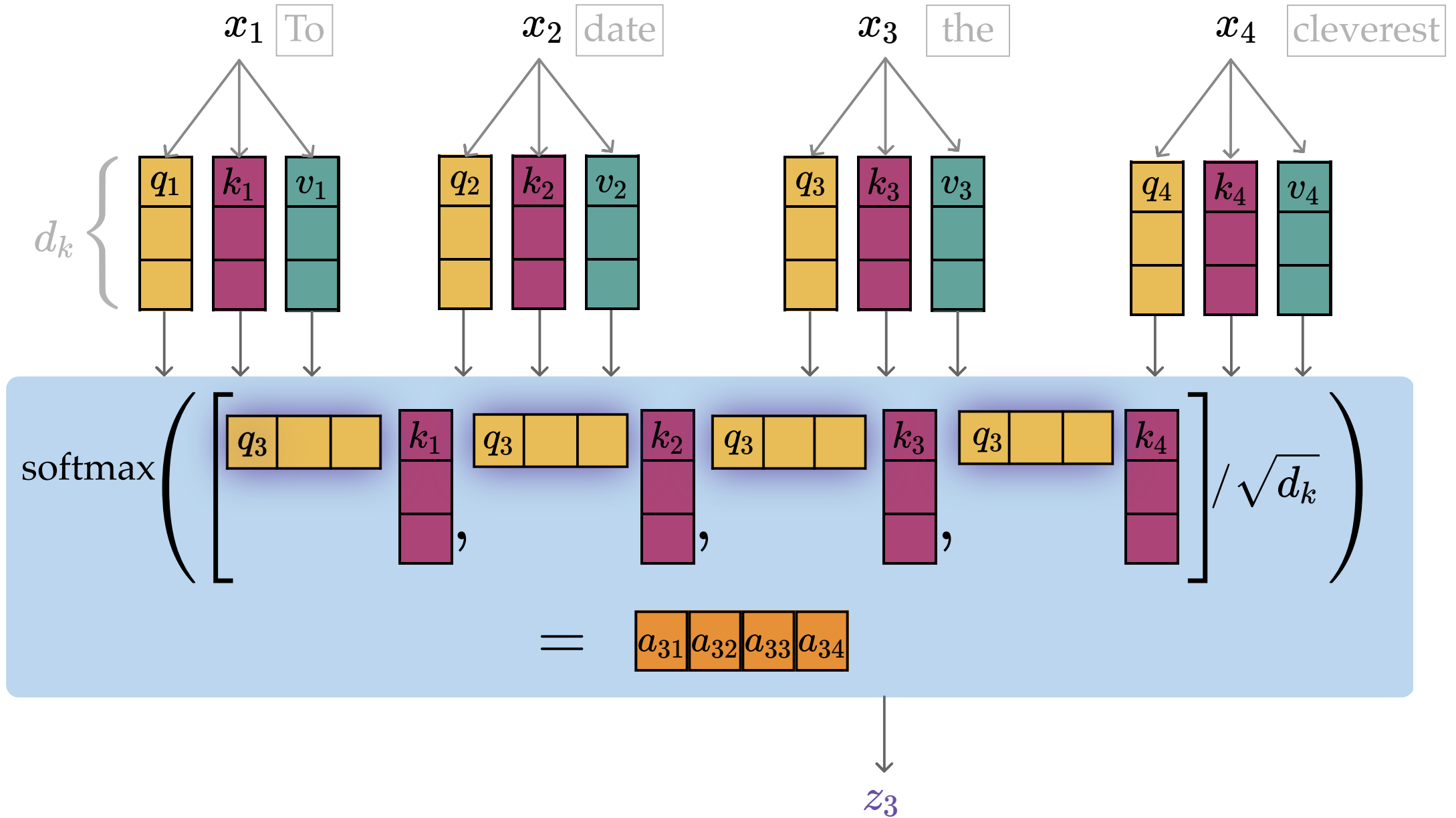


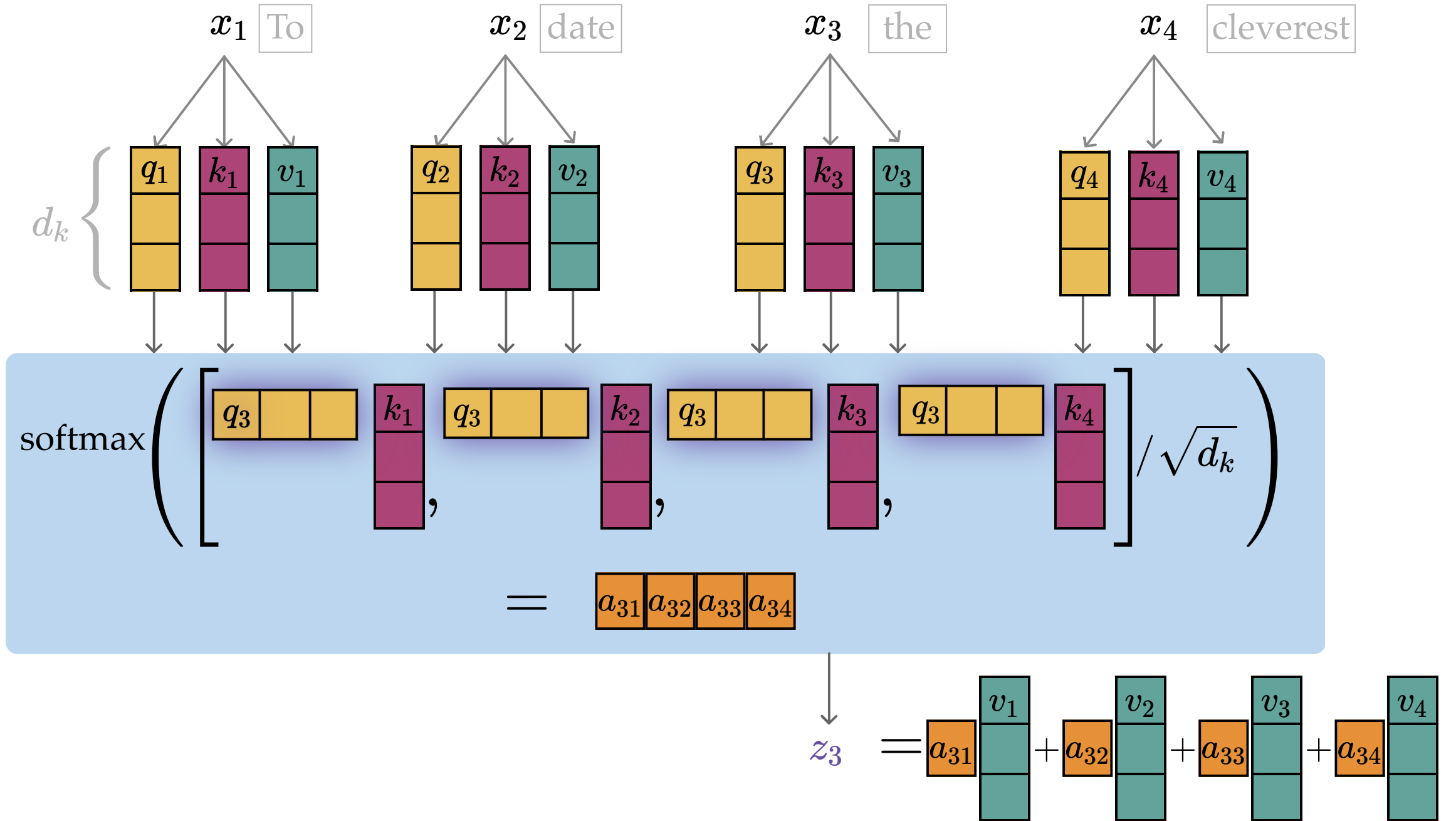


parallel and structurally identical processing
 can calculate z_4 without z_3

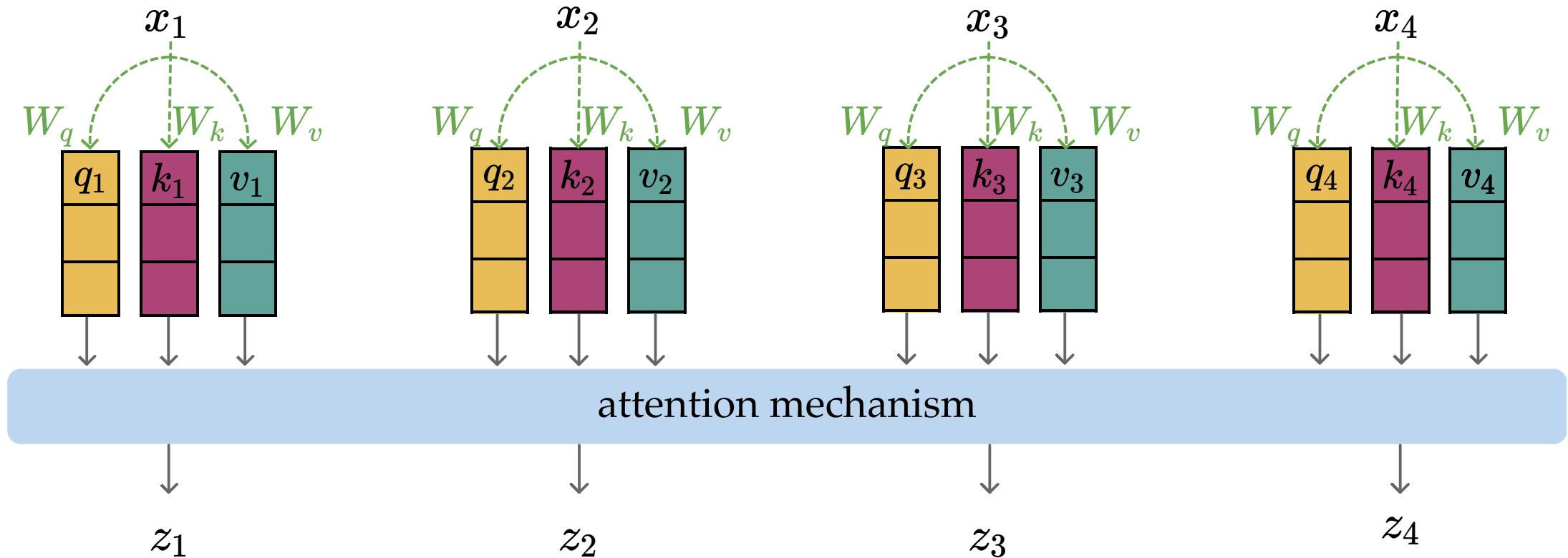








Attention head

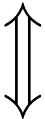
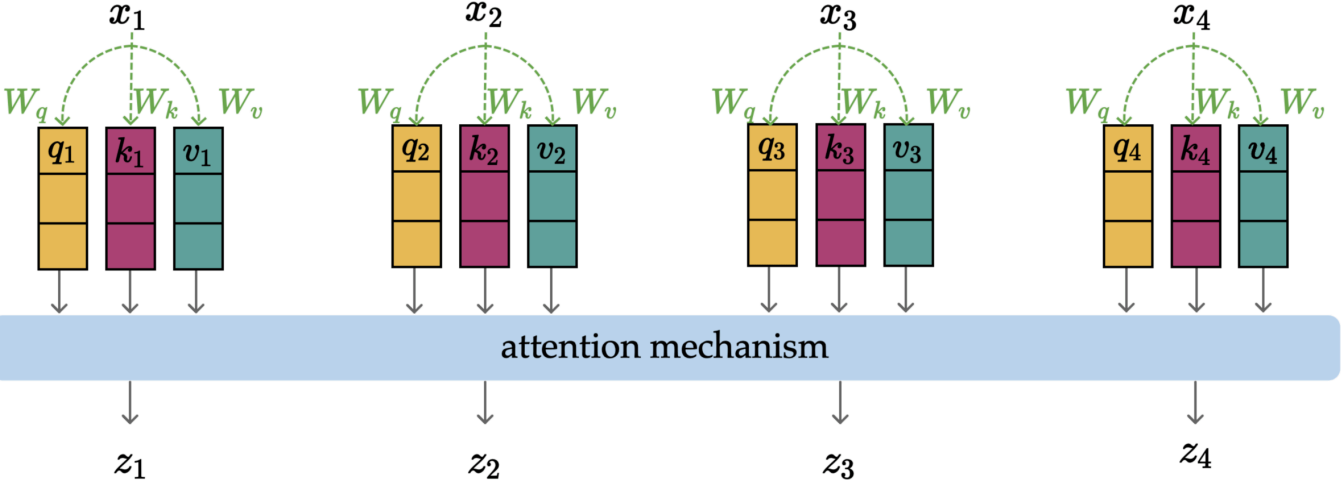


maps sequence of x to sequence of z :

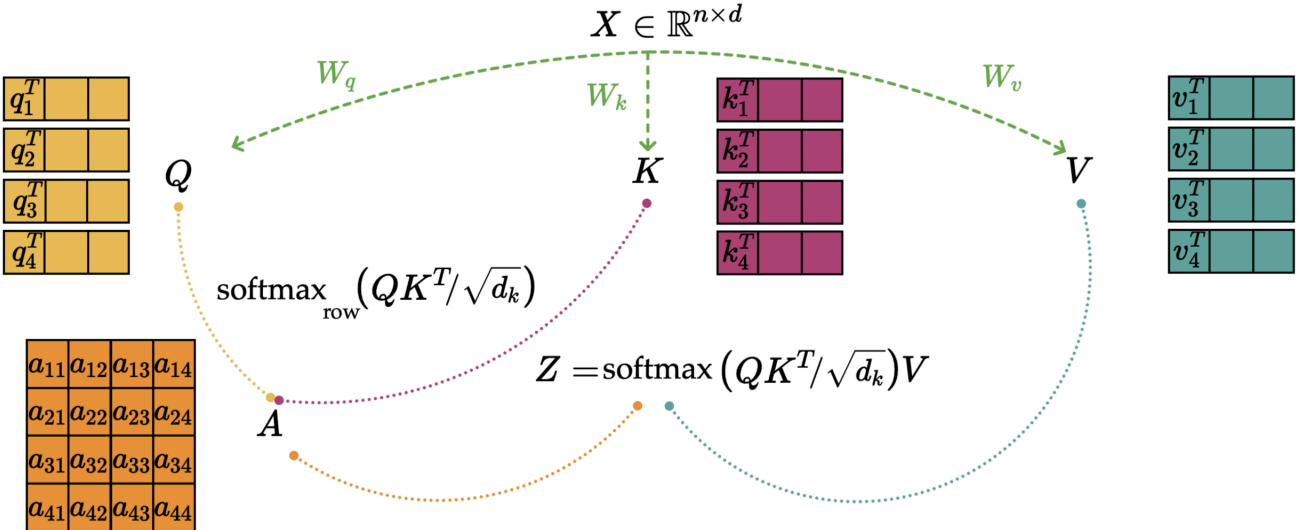
1. (query, key, value) projection
2. attention mechanism

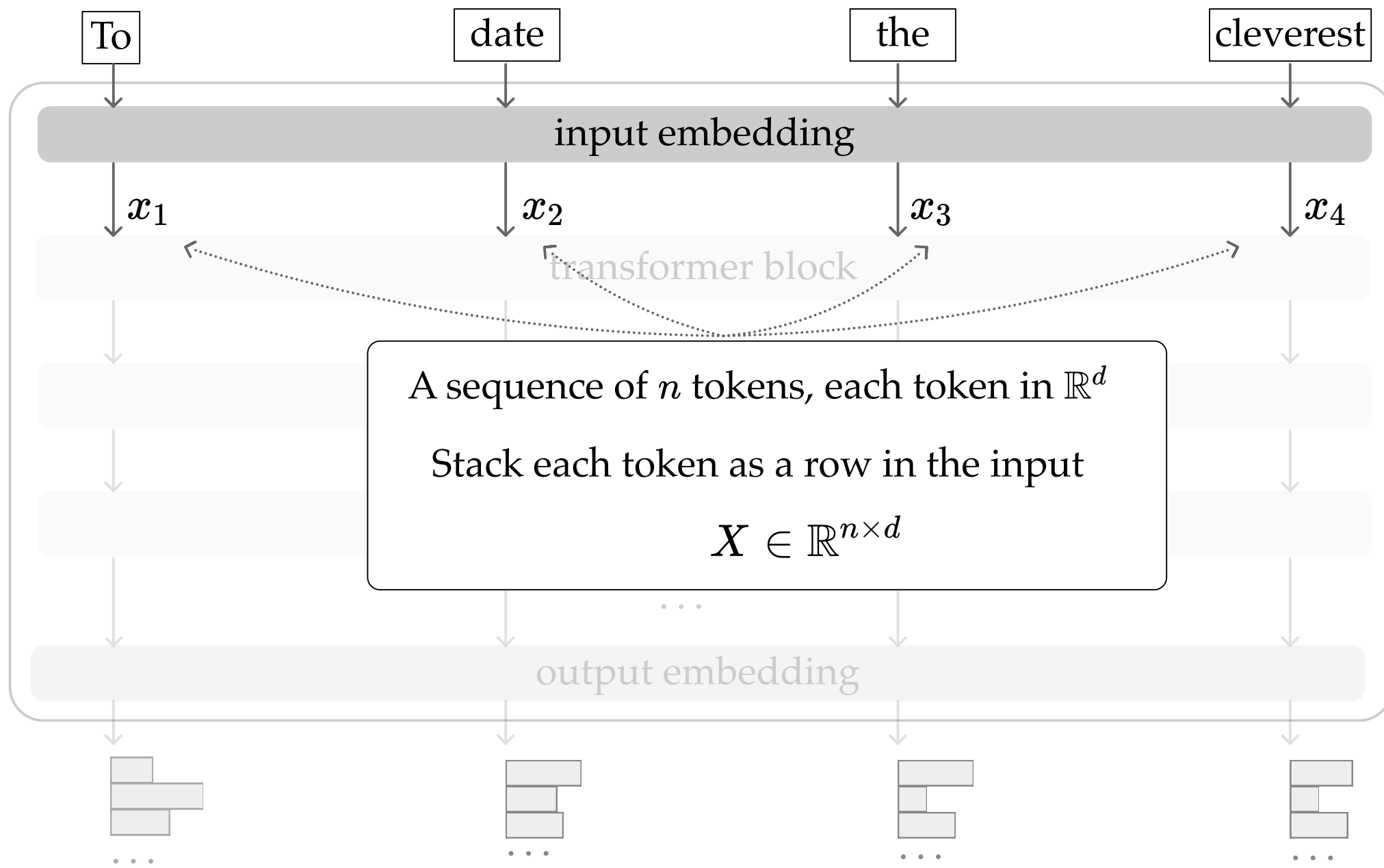
parallel and *structurally identical* processing

Attention head - compact matrix form

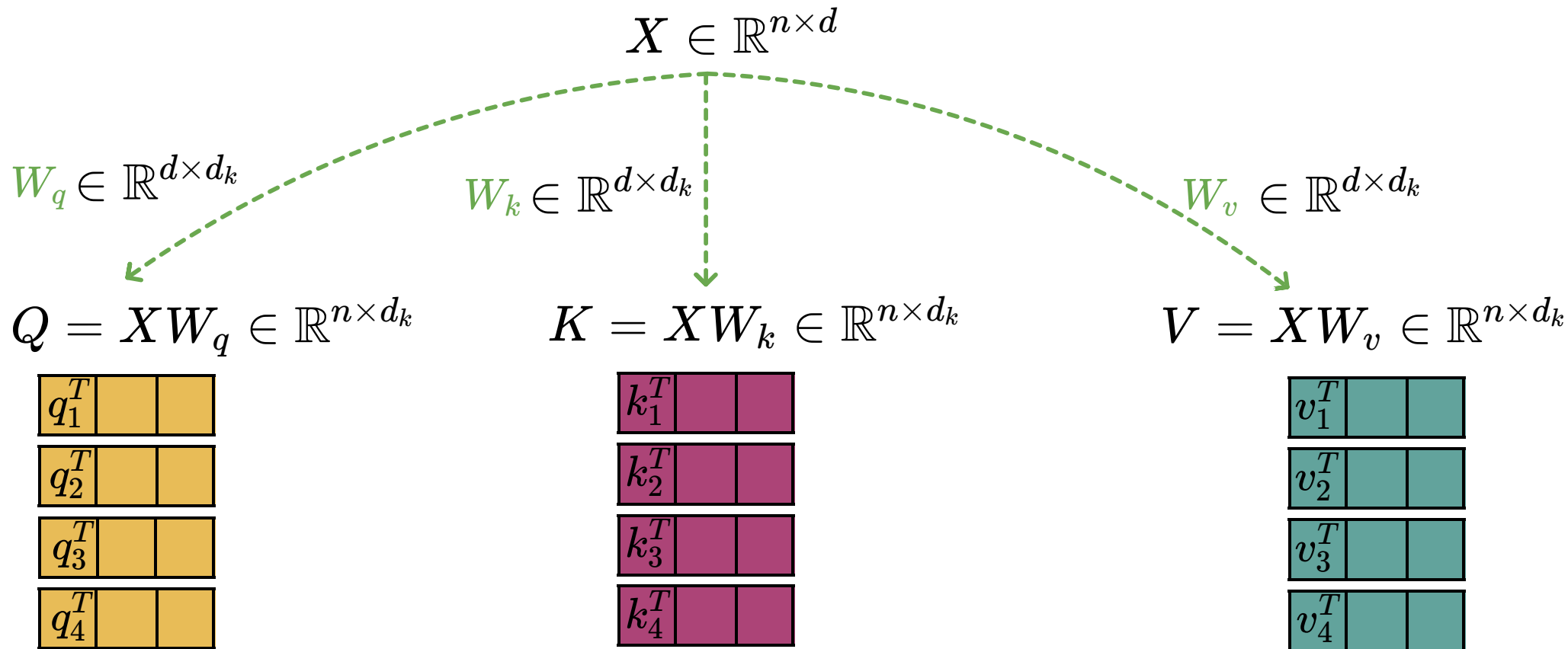


By stacking each individual vector in the sequence as a row





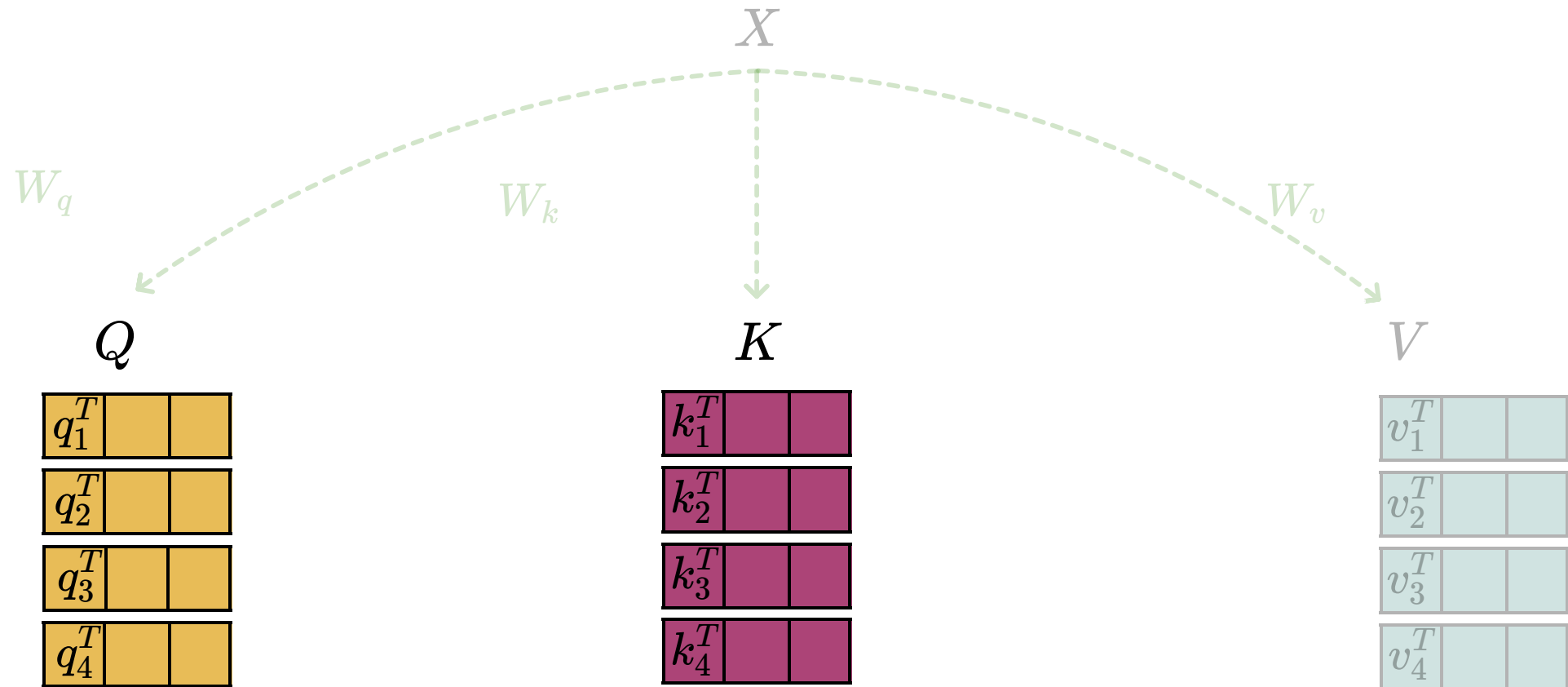
1. (query, key, value) projection



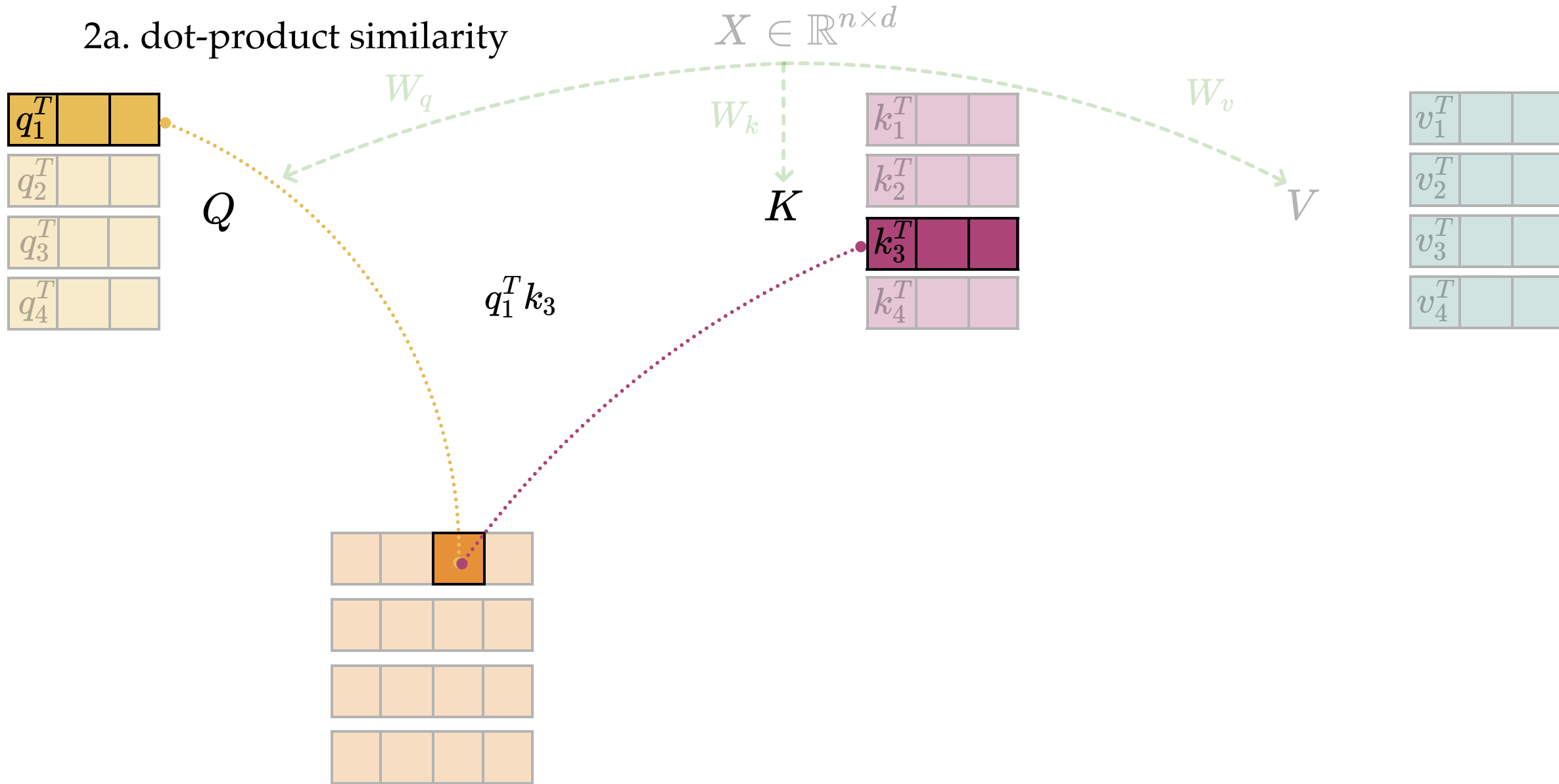
2a. dot-product similarity

compare q_i and k_j

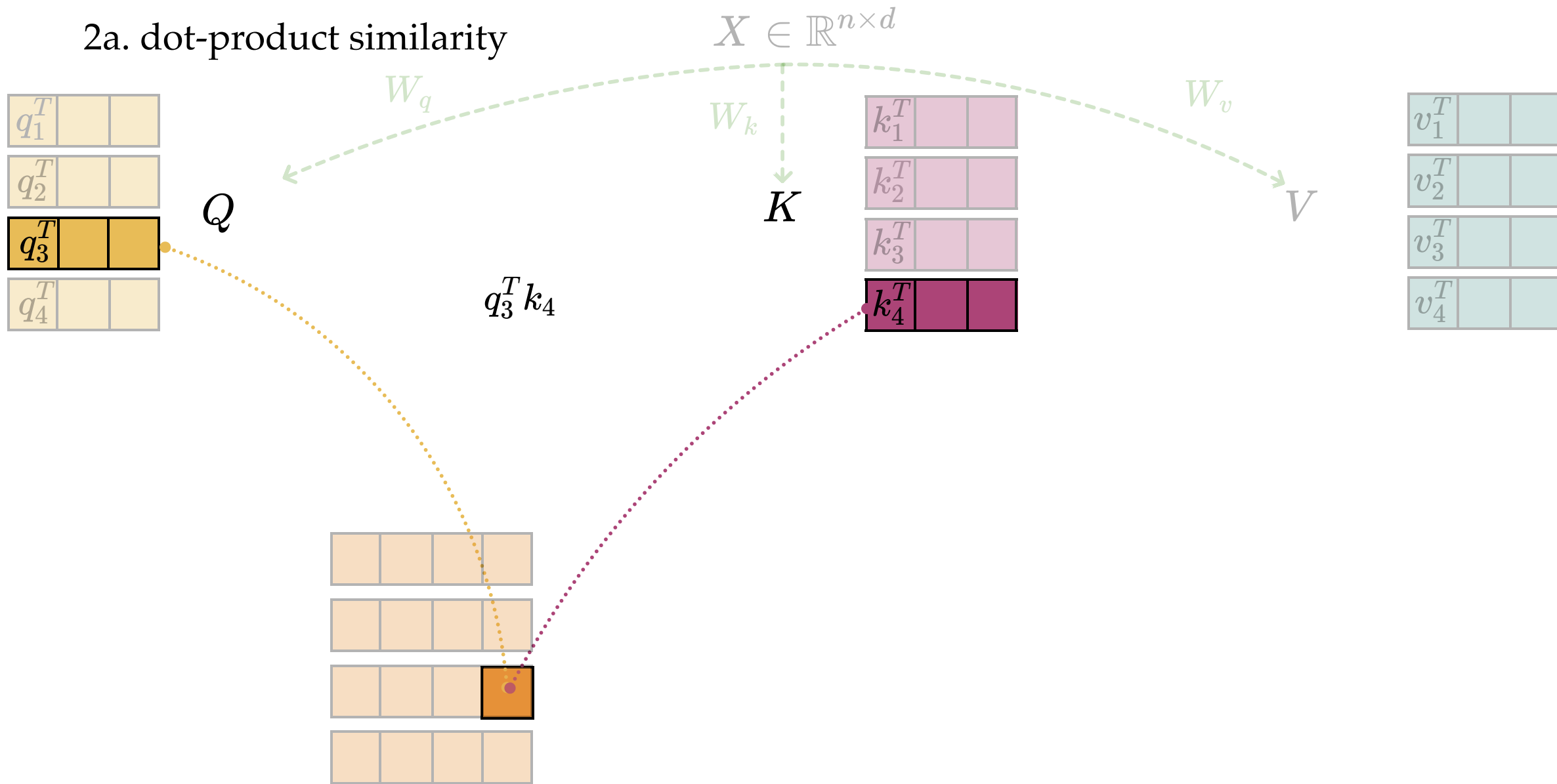
assemble the $n \times n$ similarities so rows correspond to query



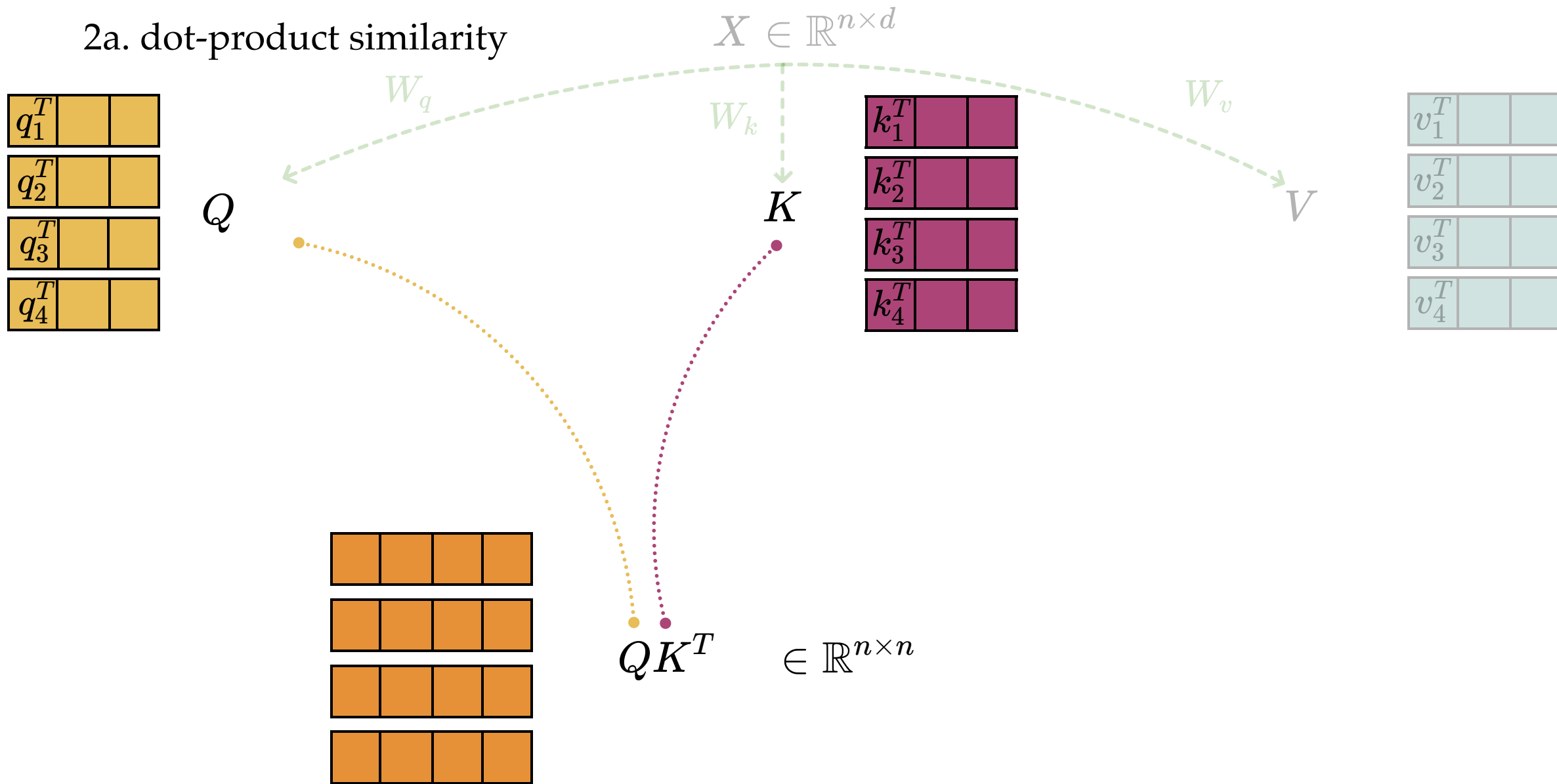
2a. dot-product similarity



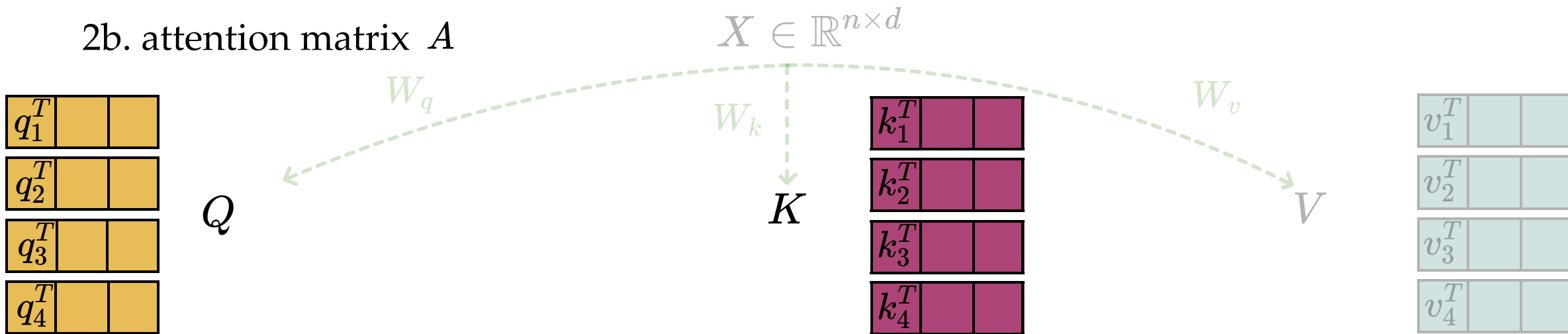
2a. dot-product similarity



2a. dot-product similarity



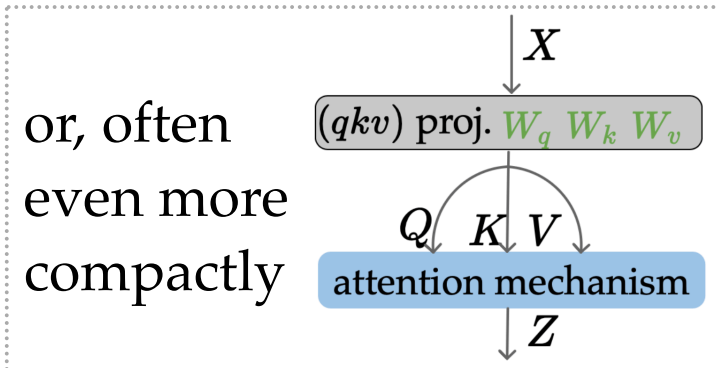
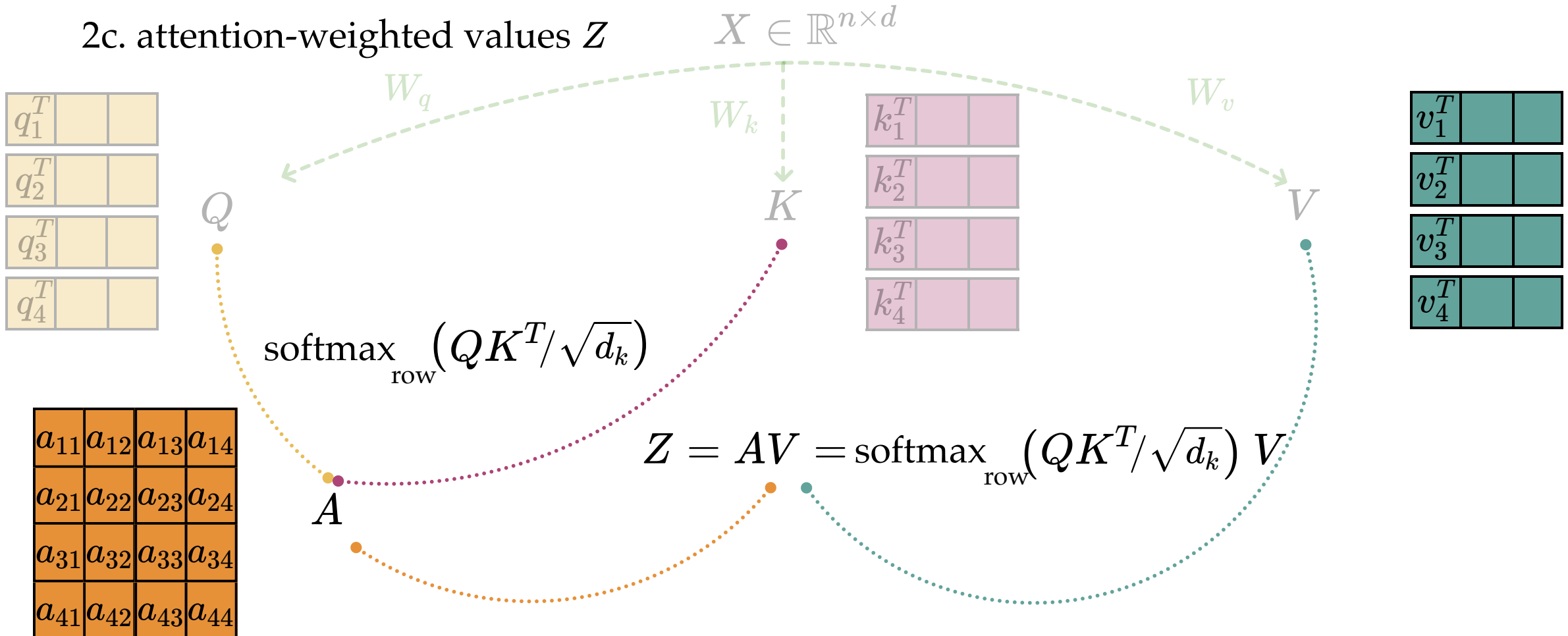
2b. attention matrix A

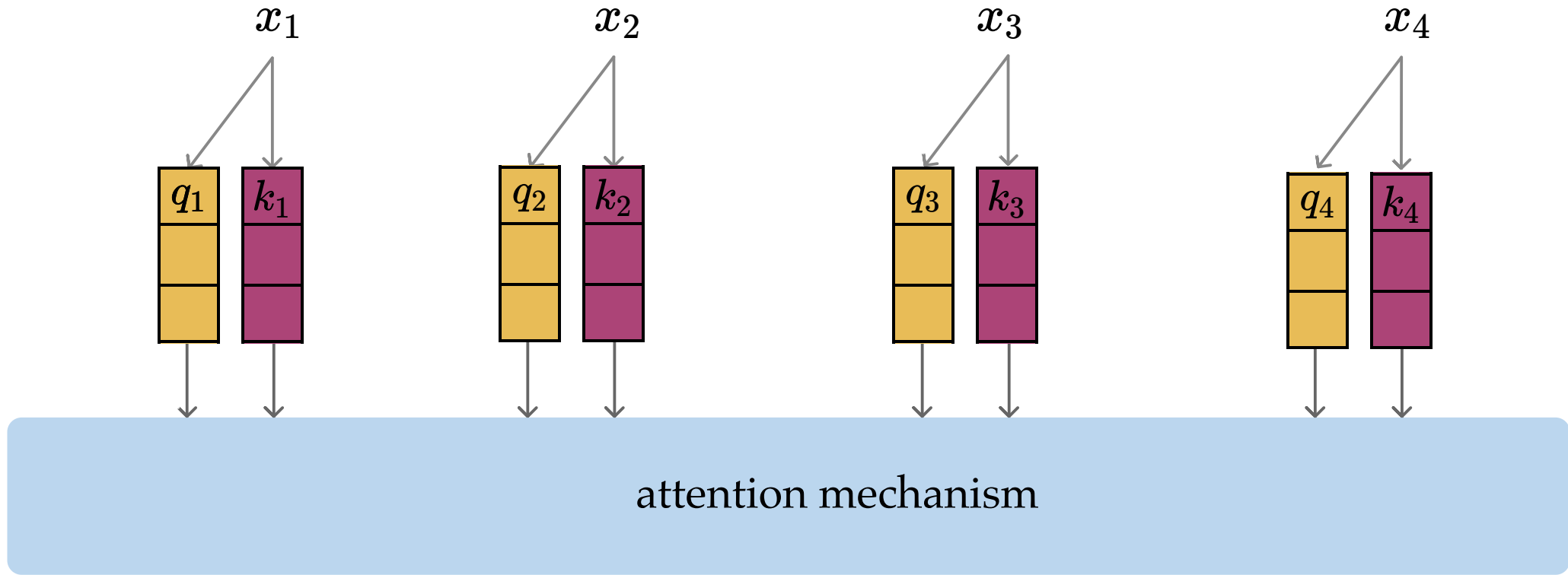


$$A = \begin{bmatrix} \text{softmax}\left(\text{orange_row}_1 / \sqrt{d_k}\right) \\ \text{softmax}\left(\text{orange_row}_2 / \sqrt{d_k}\right) \\ \text{softmax}\left(\text{orange_row}_3 / \sqrt{d_k}\right) \\ \text{softmax}\left(\text{orange_row}_4 / \sqrt{d_k}\right) \end{bmatrix} = \text{softmax}_{\text{row}}\left(QK^T / \sqrt{d_k}\right) = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

each row sums up to 1

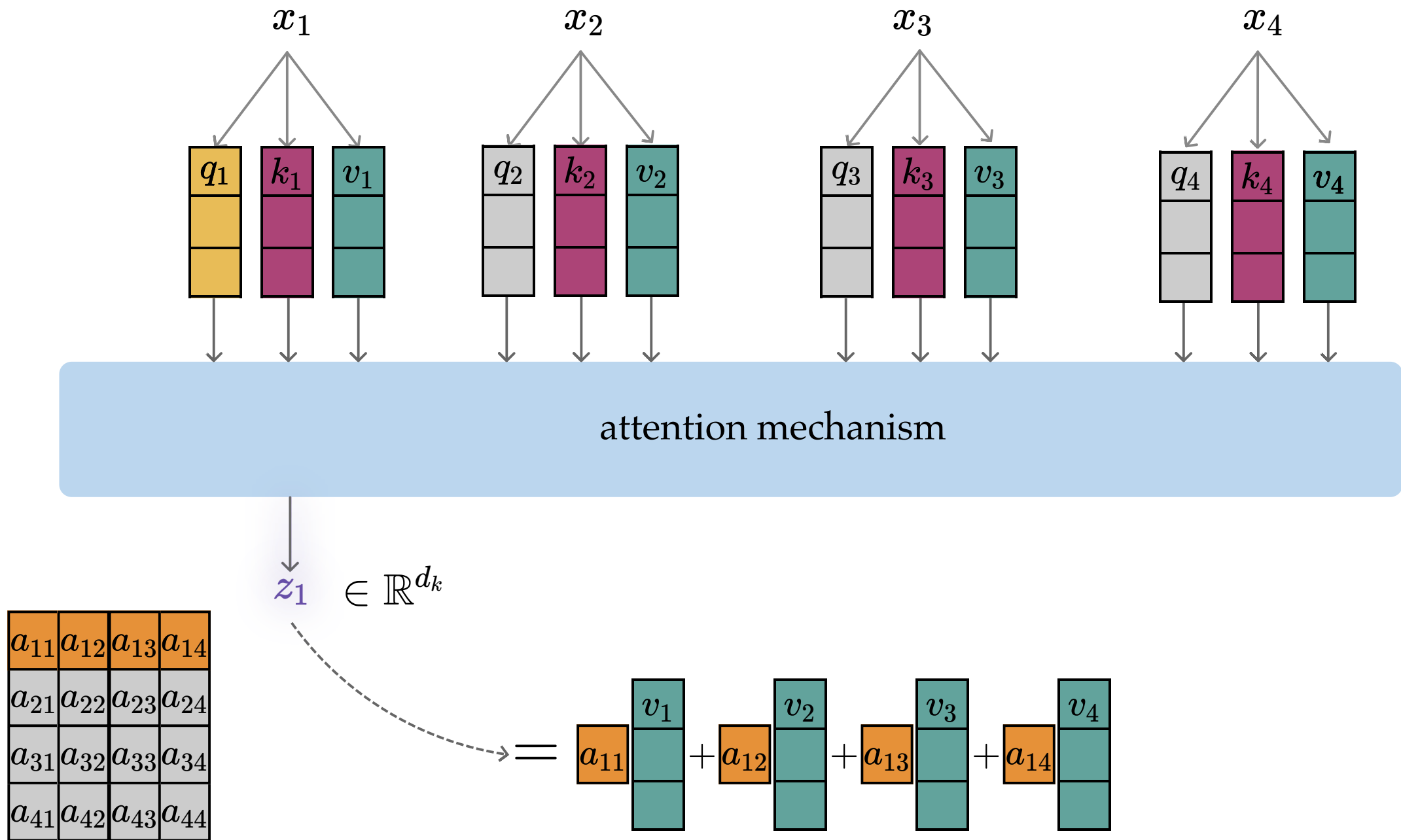
2c. attention-weighted values Z

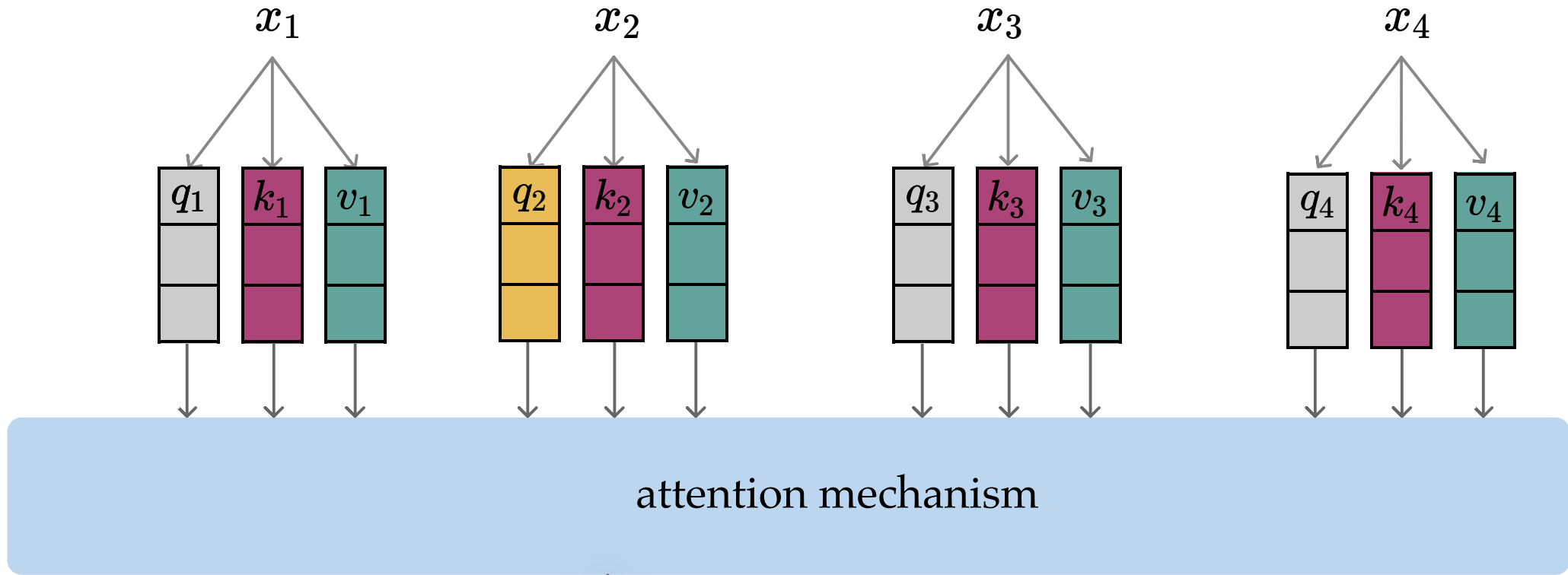




a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}

attention scores depend on the (query, key) only

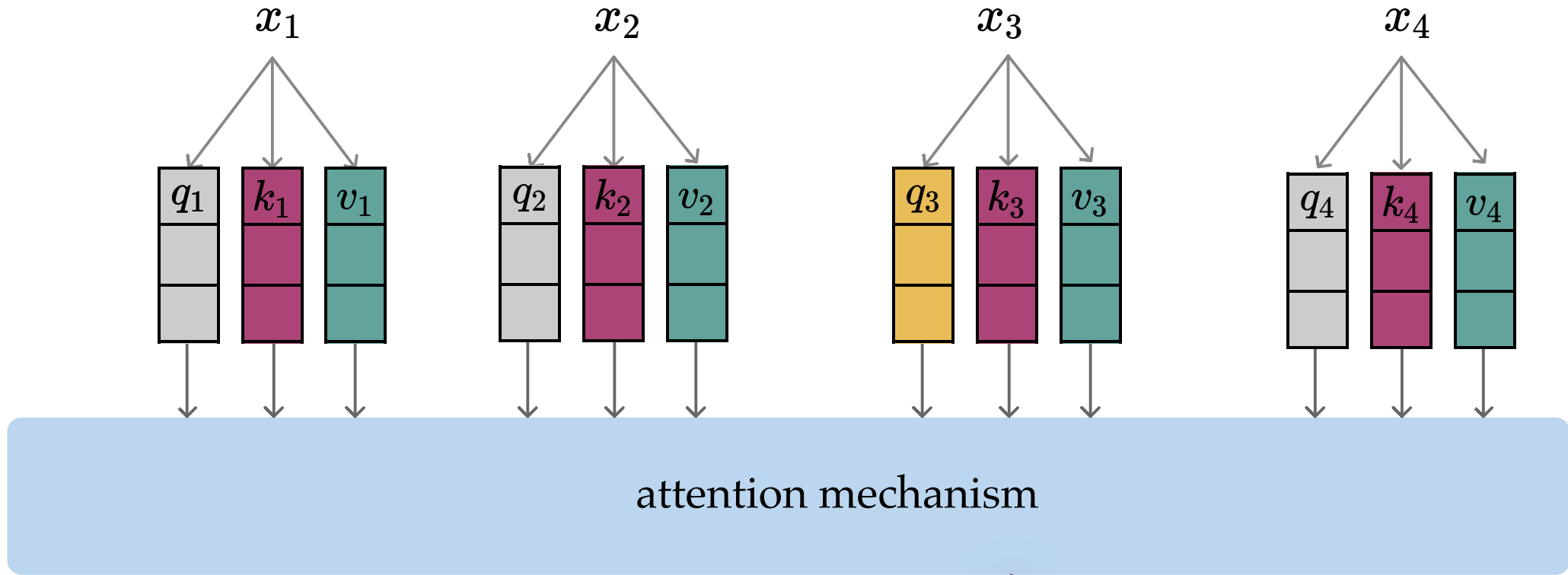




$z_2 \in \mathbb{R}^{d_k}$

a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}

$= a_{21} \begin{bmatrix} v_1 \\ \cdot \\ \cdot \end{bmatrix} + a_{22} \begin{bmatrix} v_2 \\ \cdot \\ \cdot \end{bmatrix} + a_{23} \begin{bmatrix} v_3 \\ \cdot \\ \cdot \end{bmatrix} + a_{24} \begin{bmatrix} v_4 \\ \cdot \\ \cdot \end{bmatrix}$

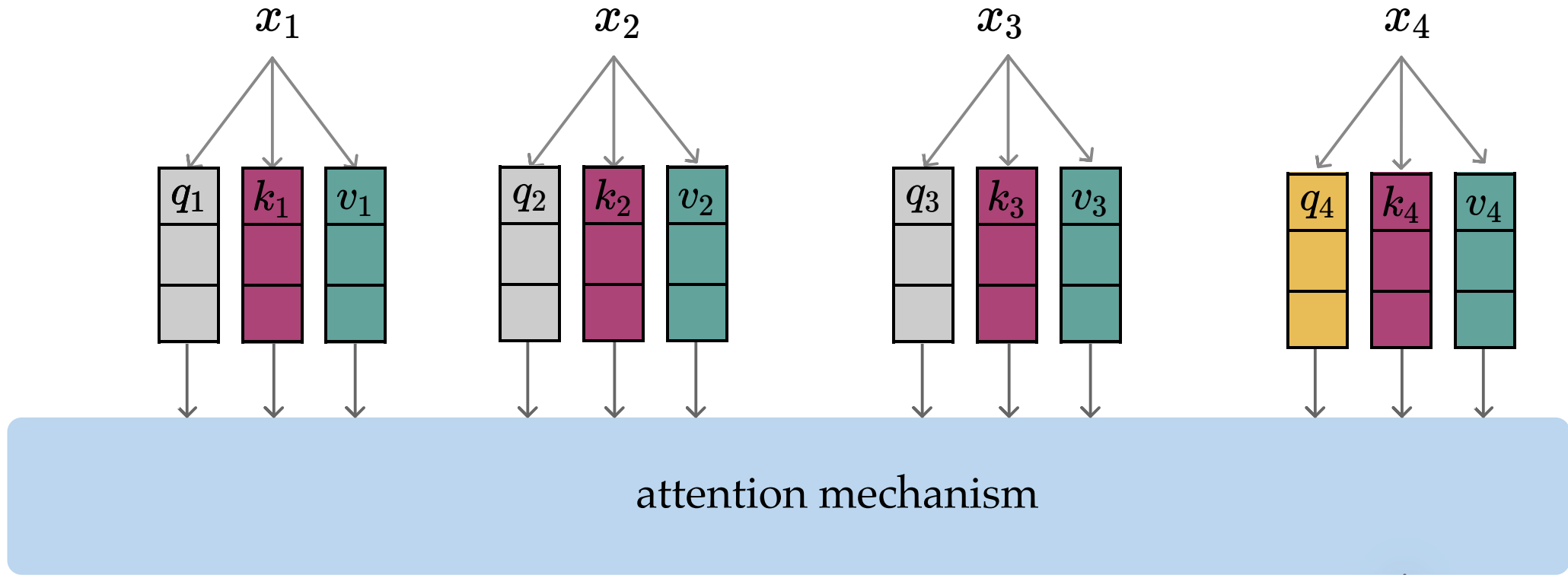


attention mechanism

$z_3 \in \mathbb{R}^{d_k}$

a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}

$$= a_{31} \begin{bmatrix} v_1 \\ \cdot \\ \cdot \end{bmatrix} + a_{32} \begin{bmatrix} v_2 \\ \cdot \\ \cdot \end{bmatrix} + a_{33} \begin{bmatrix} v_3 \\ \cdot \\ \cdot \end{bmatrix} + a_{34} \begin{bmatrix} v_4 \\ \cdot \\ \cdot \end{bmatrix}$$



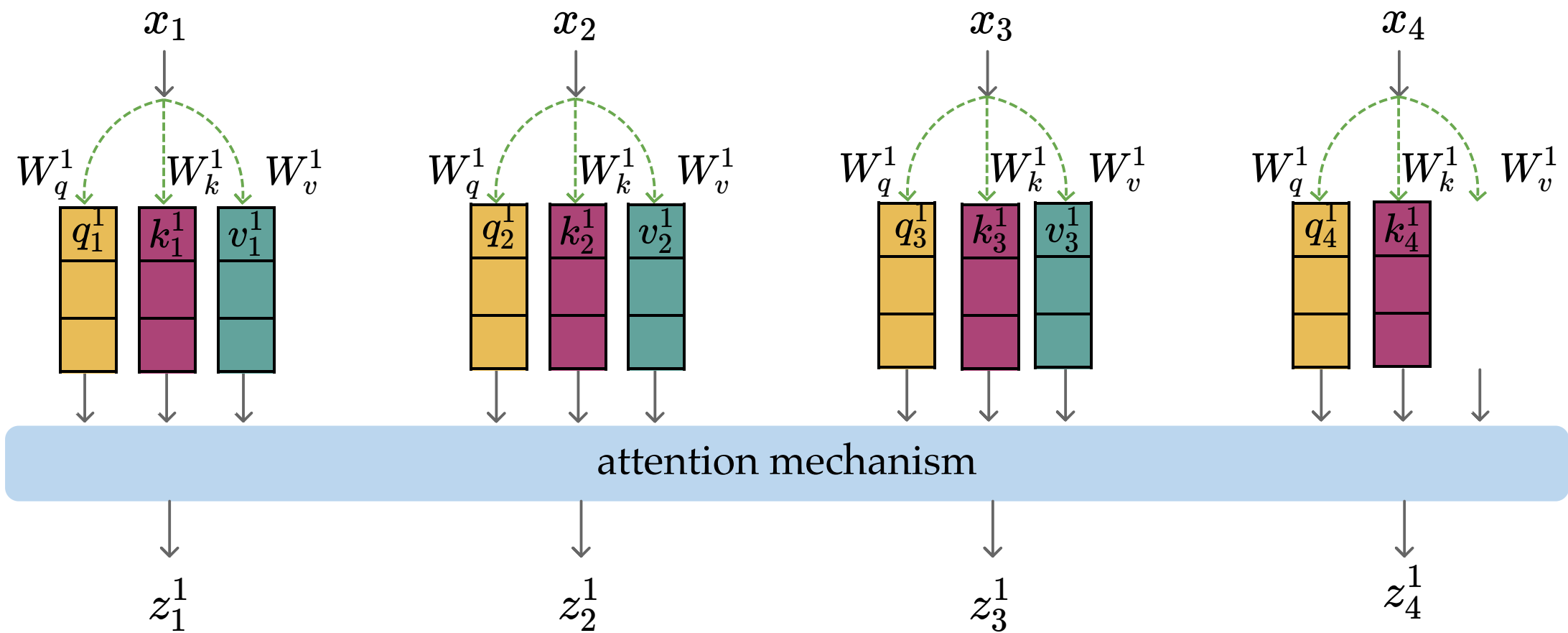
a_{11}	a_{12}	a_{13}	a_{14}
a_{21}	a_{22}	a_{23}	a_{24}
a_{31}	a_{32}	a_{33}	a_{34}
a_{41}	a_{42}	a_{43}	a_{44}

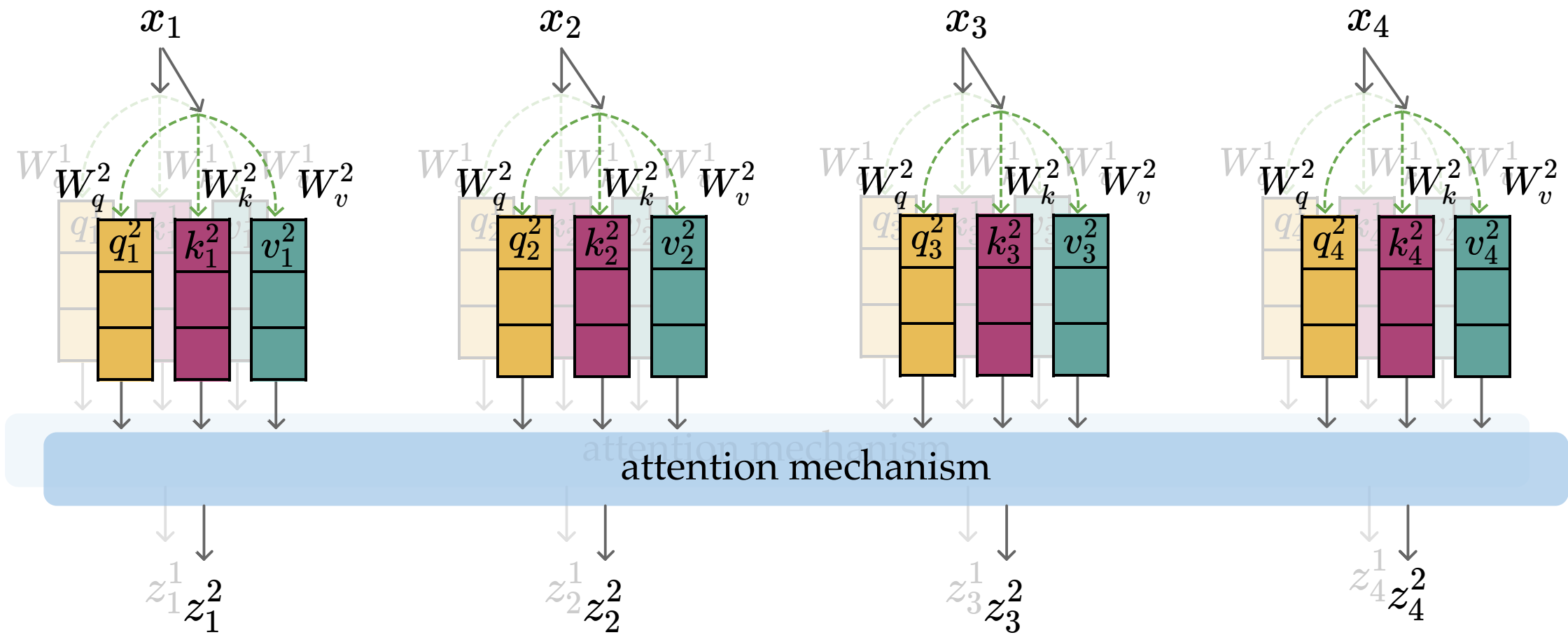
The attention mechanism outputs a vector $z_4 \in \mathbb{R}^{d_k}$, which is used to weight the value vectors v_1, v_2, v_3, v_4 to produce the final output:

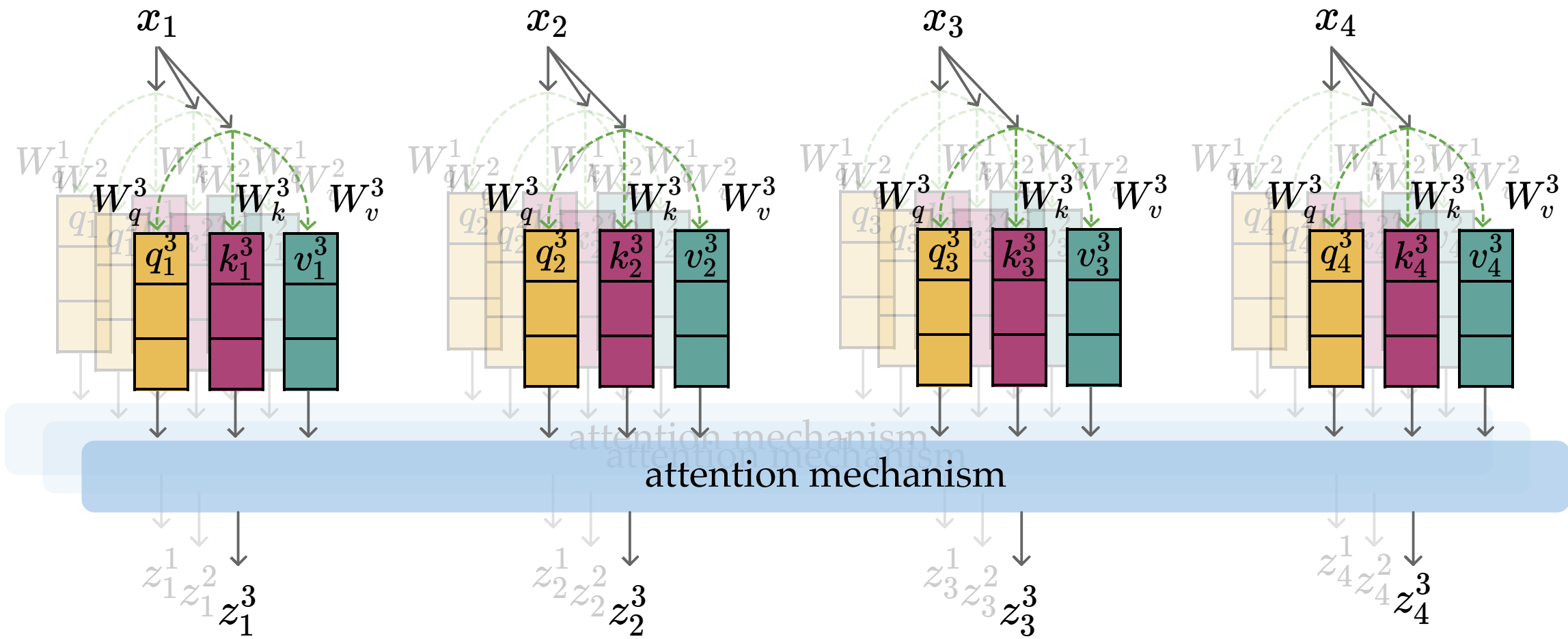
$$a_{41} v_1 + a_{42} v_2 + a_{43} v_3 + a_{44} v_4 = z_4$$

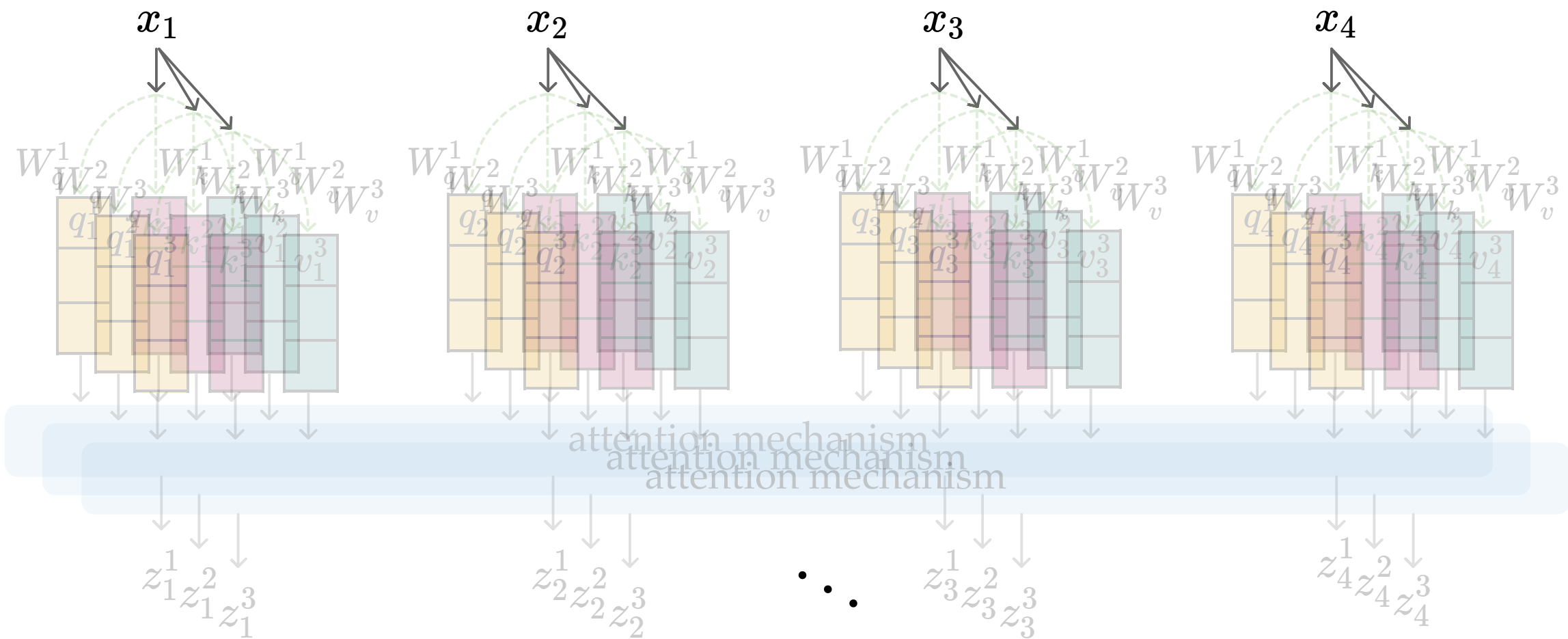
Outline

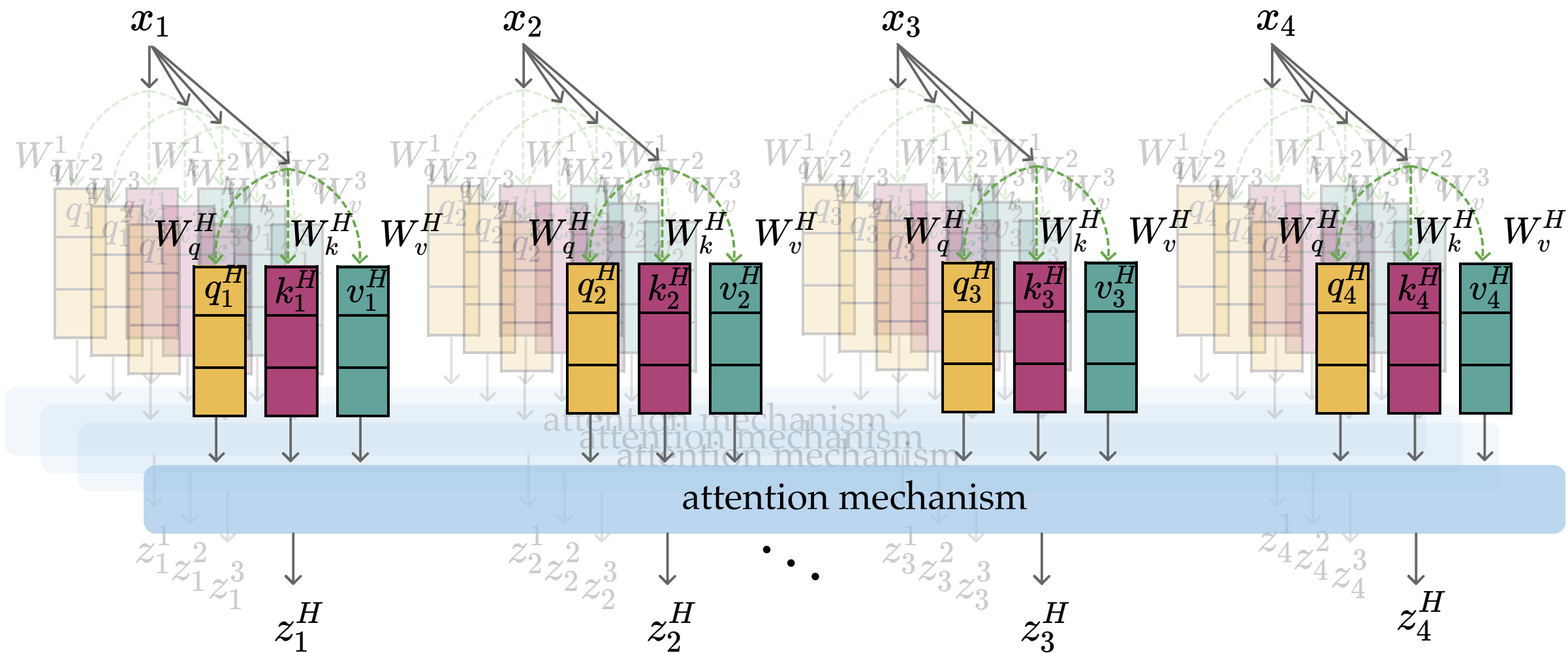
- Transformers high-level intuition and architecture
- Attention mechanism
- Multi-head attention
- (Applications)









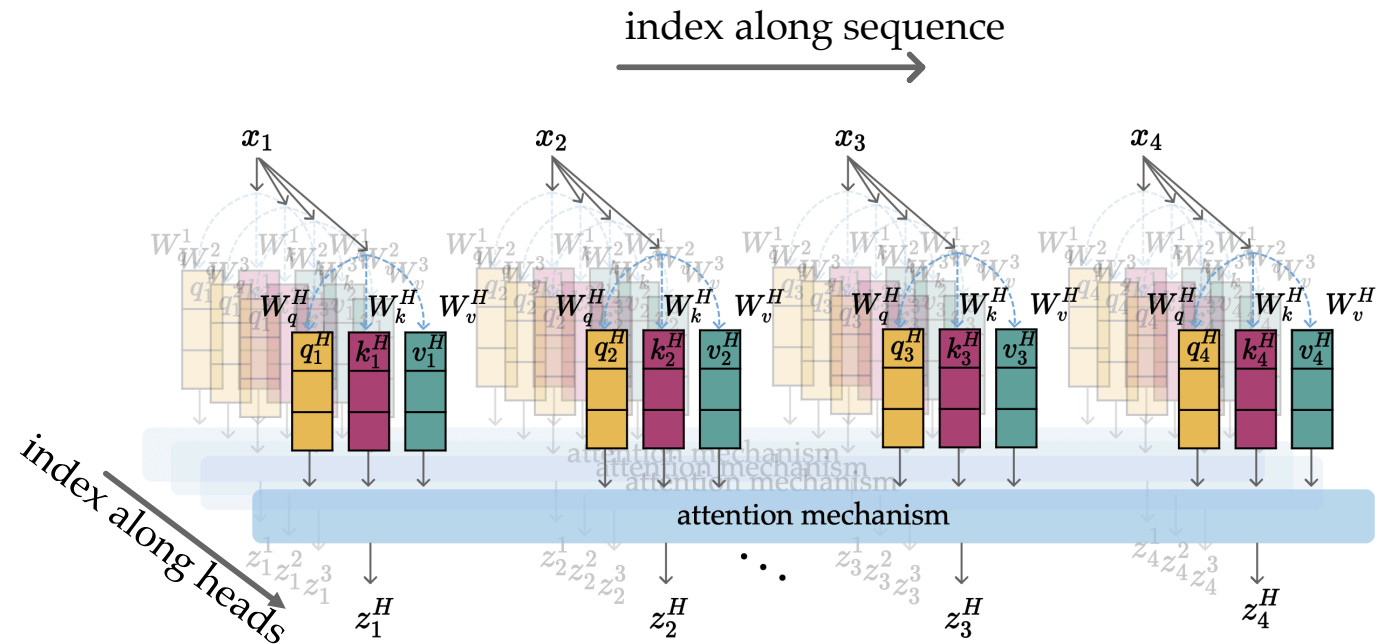


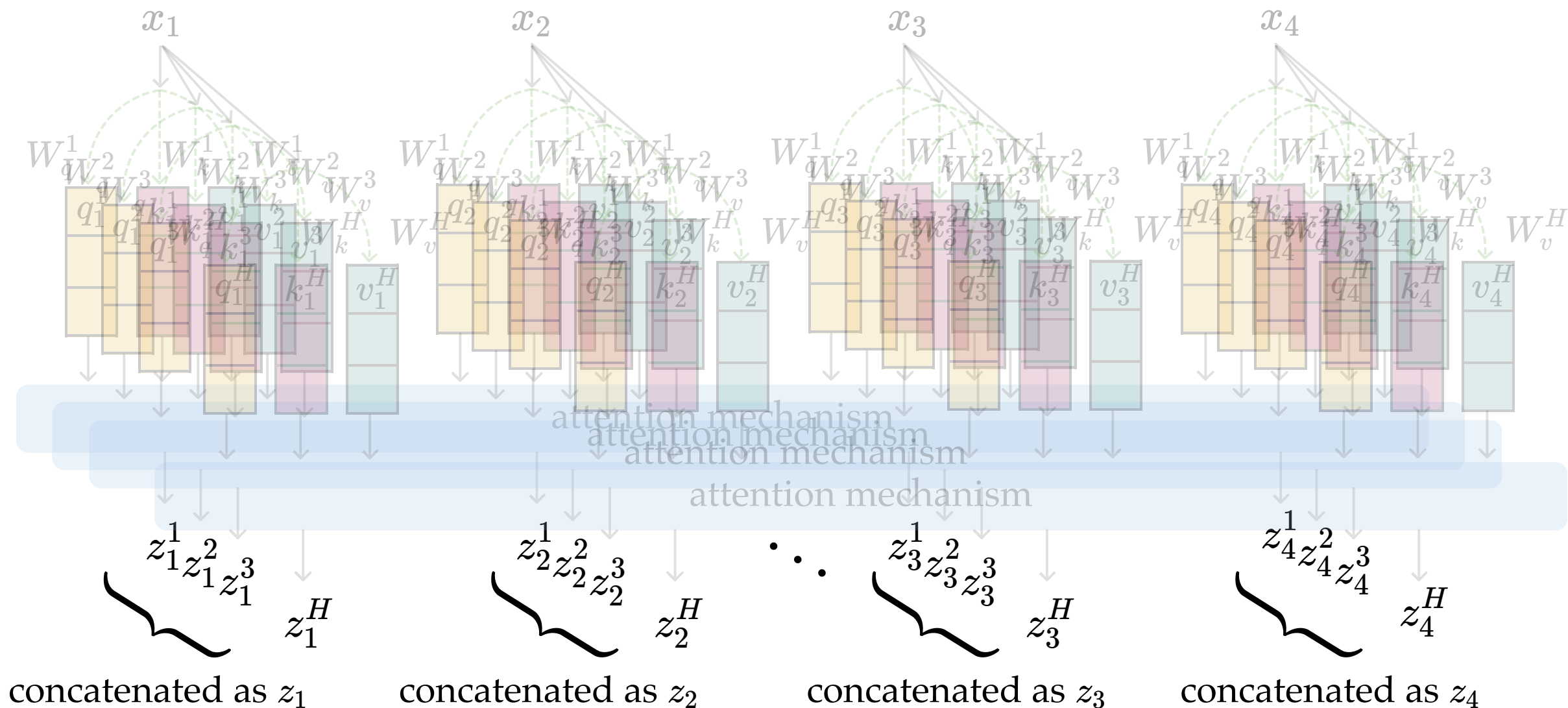
Multi-head Attention

Parallel, and structurally identical processing across all heads and tokens.

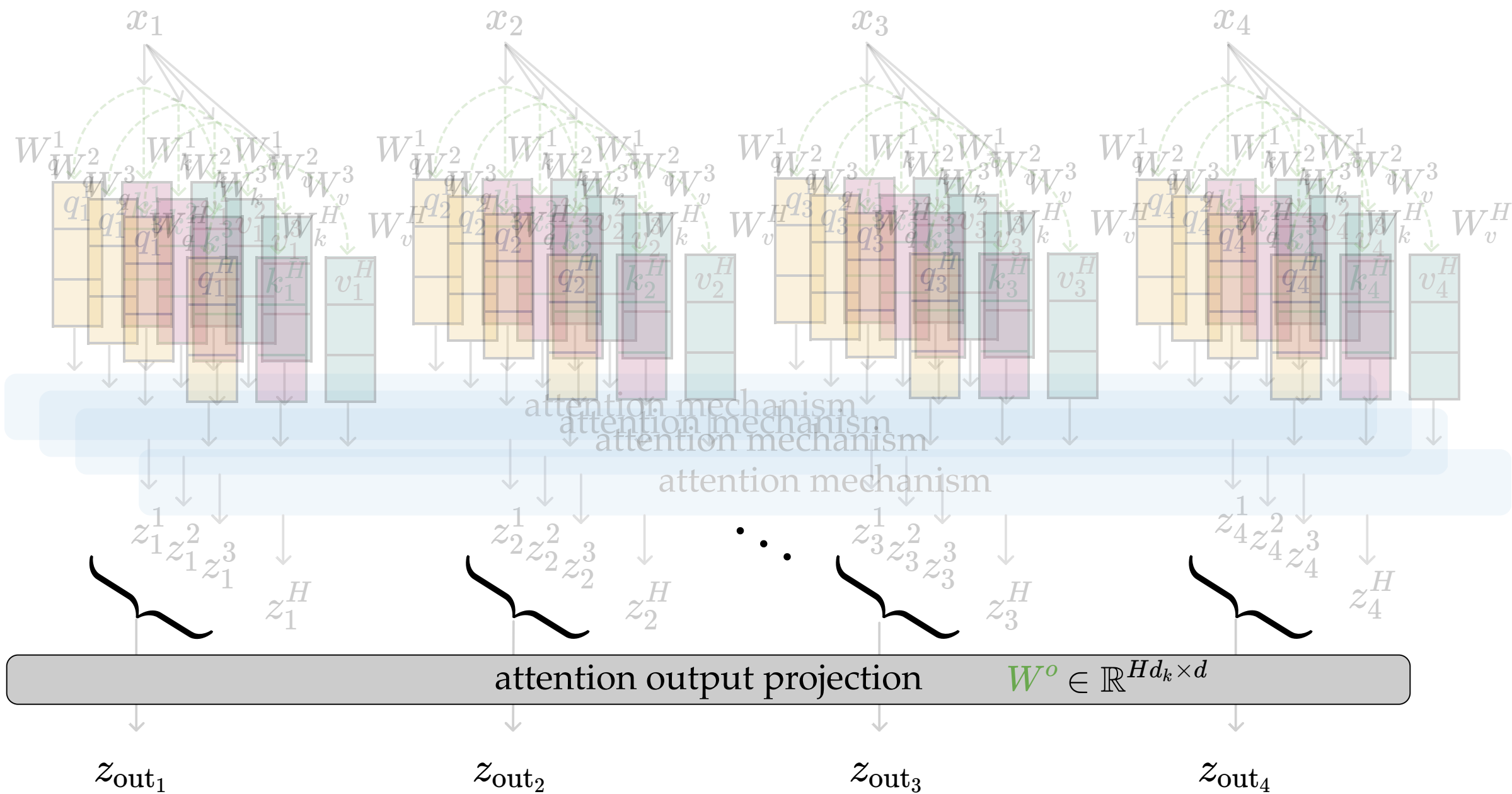
In particular, each head:

- learns its own set of W_q, W_k, W_v
- creates its own projected sequence of (q, k, v)
- computes its own sequence of z
- structurally identical processing
- for each token in the sequence:
 - structurally identical processing

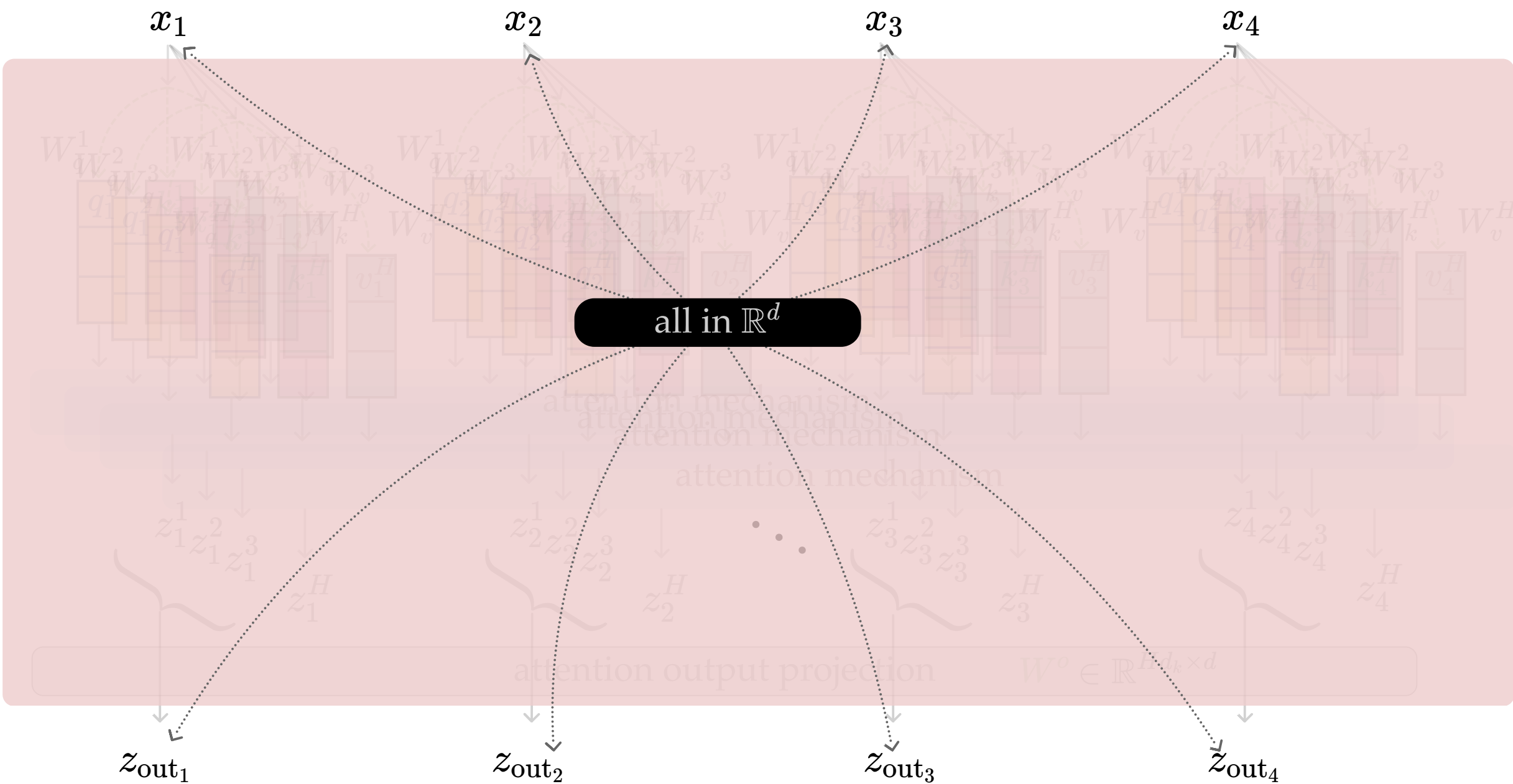




each concatenated $z_i \in \mathbb{R}^{Hd_k}$



multi-head attention

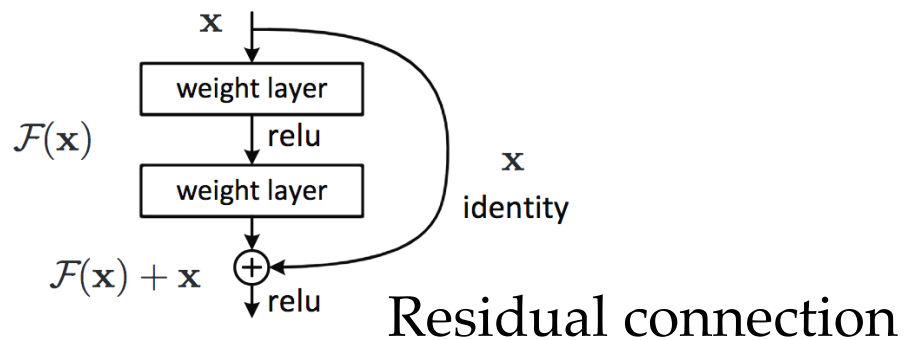


Shape Example:

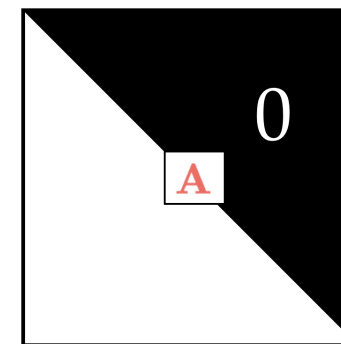
				X	input	$n \times d$	5×6	
			for a single attention head	W_q^h	query proj	$d \times d_k$	6×3	
				W_k^h	key proj	$d \times d_k$	6×3	
				W_v^h	value proj	$d \times d_k$	6×3	
n	num tokens	5			Q^h	query	$n \times d_k$	5×3
d	word-embedding dim	6			K^h	key	$n \times d_k$	5×3
d_k	(qkv) embedding dim	3			V^h	value	$n \times d_k$	5×3
					A^h	attn matrix	$n \times n$	5×5
				Z^h	attn head out	$n \times d_k$	5×3	
H	num heads	2		$\text{concat}(Z^1 \dots Z^H)$		multi-head out	$n \times Hd_k$	5×6
				W^o	output proj	$Hd_k \times d$	6×6	
				Z_{out}	attn layer out	$n \times d$	5×6	

W s are the learned weights

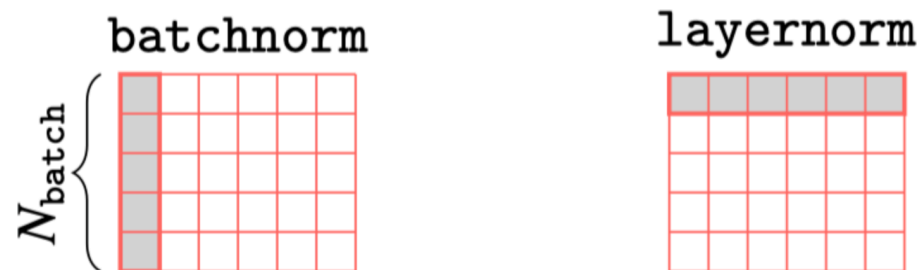
Some practical techniques commonly needed when training auto-regressive transformers:



masking



Positional encoding



Layer normalization



<https://poloclub.github.io/transformer-explainer/>

Outline

- Transformers high-level intuition and architecture
- Attention mechanism
- Multi-head attention
- (Applications)

Generative Boba



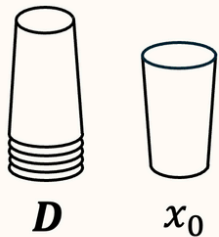
By Boyuan Chen
boyuanc@mit.edu

Intro:

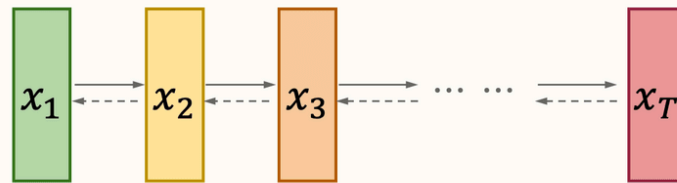
Boba tea is an essential component of Asian students' well-being and productivity. Prior studies [1] have suggested the lack of accessible boba shop hurts MIT's reputation as a top CS school. To address this problem, we open a novel boba shop at MIT - Generative Boba. By providing PhD students with free boba every afternoon, Generative Boba boosts the research productivity of the floor by TBD%. Generative Boba is also looking for Chief Boba Scientists to scale up! Contact me to be a co-author.

[1] Chen, B. (2022, June 20). BEST COMPUTER SCIENCE SCHOOLS RANKED BY BOBA. Boyuan's Blog. Retrieved March 27, 2024,.

Step 1:
Sample a cup from cup set



Step 2:
Gradually add t steps of ingredients



Step 3:
Graduate Student Descent



Fig 1: We propose boba diffusion, a novel generative model that boosts PhD student productivity

Today's special:

<p>$\epsilon_1 = \text{topping}$</p> <p>20%</p> <p>open boba pot for topping close pot lid to keep warm</p>	<p>$\epsilon_2 = \text{milk}$</p> <p>30-60%</p> <p>we serve lactose-free milk, sometimes coconut milk</p>	<p>$\epsilon_3 = \text{tea}$</p> <p>20-50%</p> <p>use 20% for matcha use 50% for black tea</p>	<p>$\epsilon_4 = \text{syrup}$</p> <p>any</p> <p>hold black sugar bottle upside down for 5-10s</p>	<p>add a lid and straw our straws are compostable</p>
--	--	---	---	---

Fig 2: Implementation details of boba diffusion. We follow an efficient ingredient schedule while "today's special" provides special recipes from time to time.

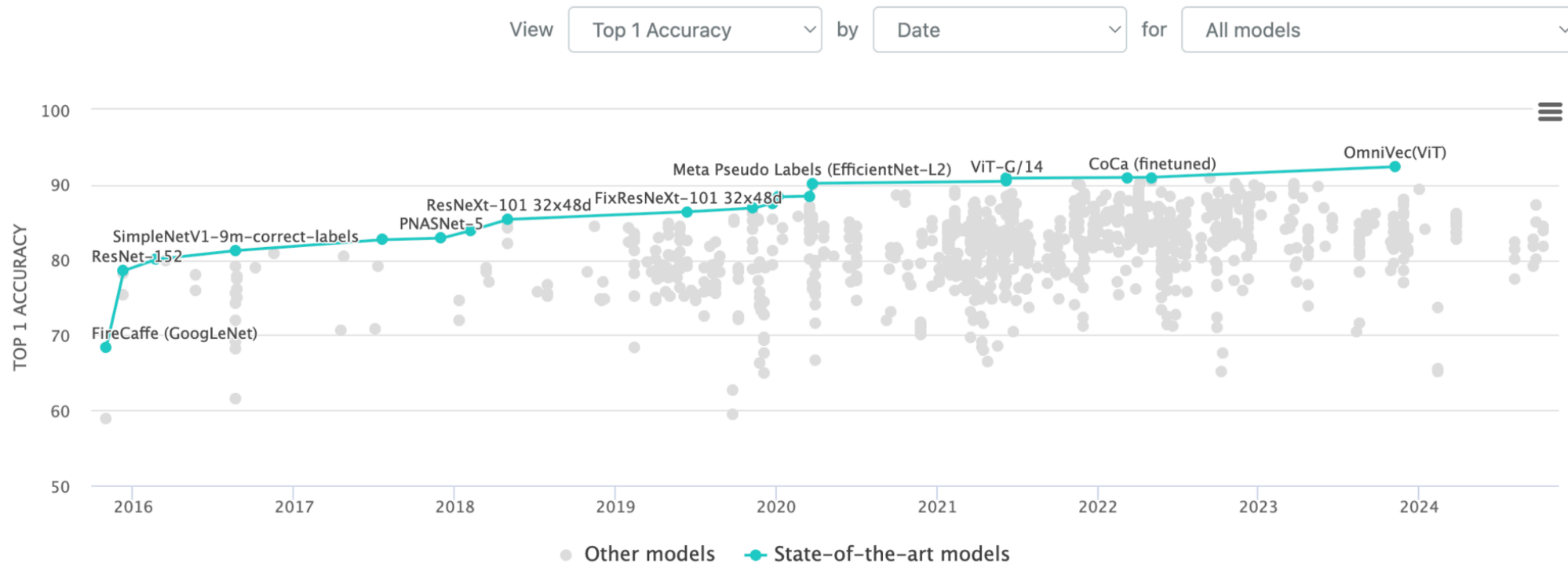


Image Classification on ImageNet

Leaderboard

Dataset

Transformers in Action: Performance across domains



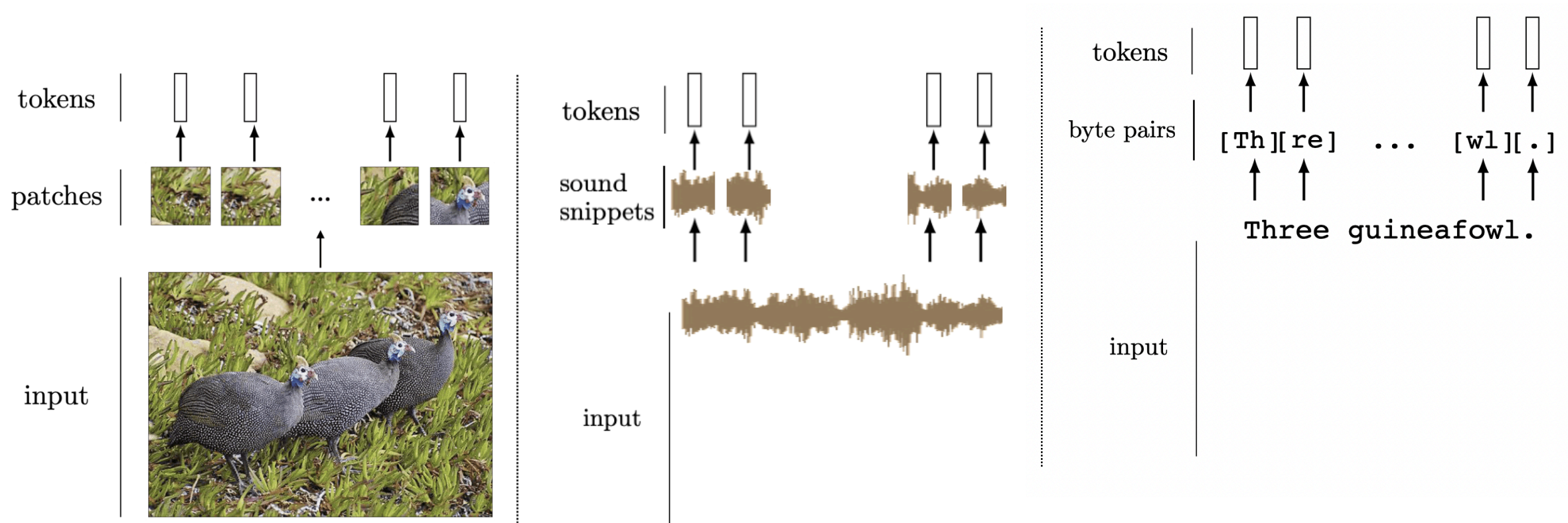
Filter:

- ImageNet-1k only
- Transformer
- ResNet
- CNN
- ImageNet-22k
- EfficientNet
- JFT-300M
- MLP
- ResNeXt
- Reversible
- Neighborhood Attention
- NAT Transformer
- JFT-3B
- PatchConvnet
- FPN
- MoE
- Early Exit
- Dynamic Model Arch
- CNN+Transformer
- ALIGN
- Conv+Transformer
- CLIP data
- No Extra Data
- SNN
- IG-1B
- Swin-Transformer
- Teacher-22k
- Vision Transformer
- FLD-900M
- Pure CNN
- YFCC-15M
- Laion-400M
- Contrastive
- ConvNeXt
- Self-Supervised Learning
- RegNet
- Mixer
- Memory-Centric
- CLIP Pre-trained
- CrossCovarianceAttention
- untagged
- Hardware Burden
- Operations per network pass
- Robustness reports

Edit Leaderboard

We can tokenize anything.

General strategy: chop the input up into chunks, *project* each chunk to an **embedding**



Multi-modality (image q&a)

- (query, key, value) come from different input modality
- cross-attention

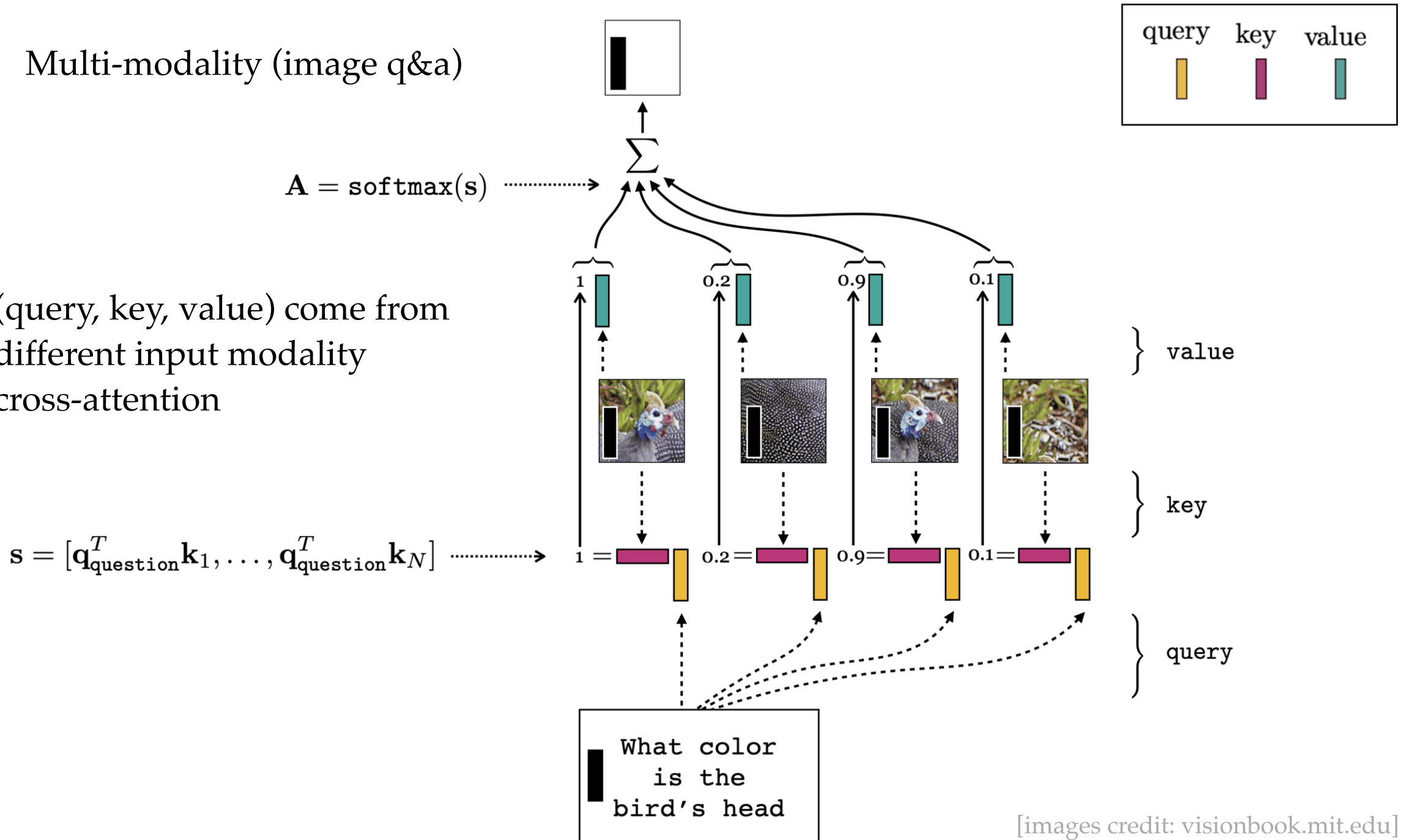
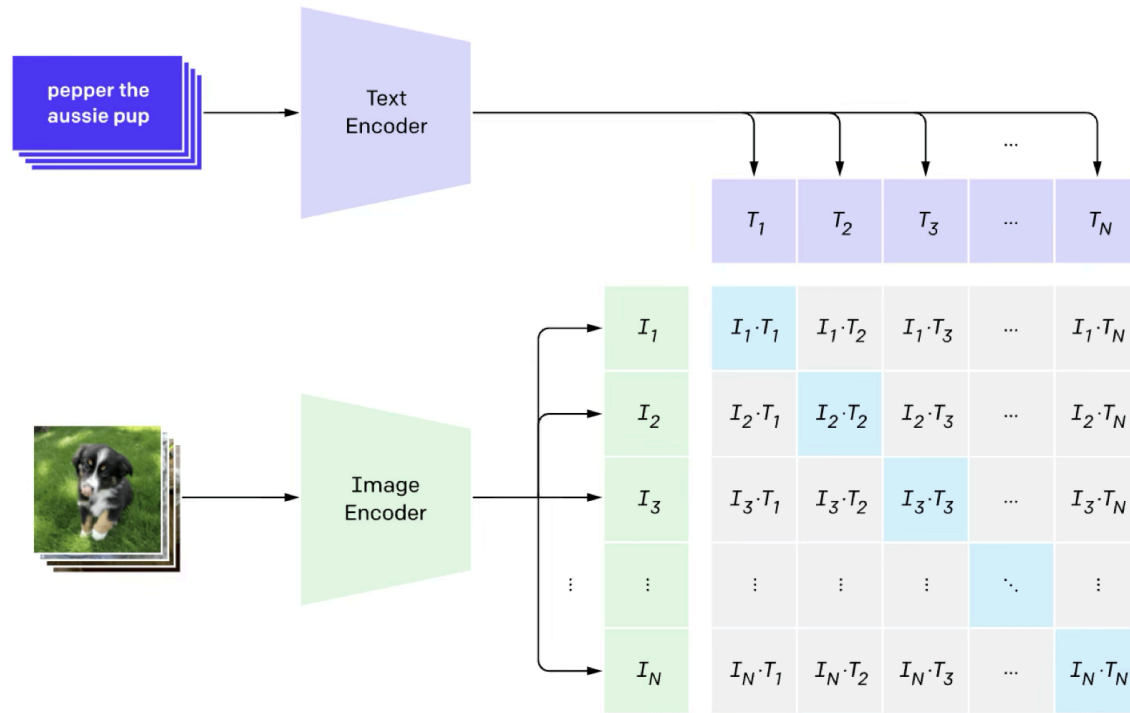


image classification (done in the contrastive way)

1. Contrastive pre-training



```
# extract feature representations of each modality
```

```
I_f = image_encoder(I) #[n, d_i]
```

```
T_f = text_encoder(T) #[n, d_t]
```

```
# joint multimodal embedding [n, d_e]
```

```
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
```

```
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)
```

```
# scaled pairwise cosine similarities [n, n]
```

```
logits = np.dot(I_e, T_e.T) * np.exp(t)
```

```
# symmetric loss function
```

```
labels = np.arange(n)
```

```
loss_i = cross_entropy_loss(logits, labels, axis=0)
```

```
loss_t = cross_entropy_loss(logits, labels, axis=1)
```

```
loss = (loss_i + loss_t)/2
```

TRANSFORMER PROTEIN LANGUAGE MODELS ARE UNSUPERVISED STRUCTURE LEARNERS

Roshan Rao*
UC Berkeley
rmrao@berkeley.edu

Joshua Meier
Facebook AI Research
jmeier@fb.com

Unsupervised contact prediction functional constraints for protein structure prediction is the predominant approach for protein structure prediction. In the past, this has been a potential alternative, but performance has been poor in bioinformatics. In this paper, we

learn contacts from the unsupervised language modeling objective. We find the highest capacity models that have been trained to date already outperform a state-of-the-art unsupervised contact prediction pipeline, suggesting these pipelines can be replaced with a single forward pass of an end-to-end model¹

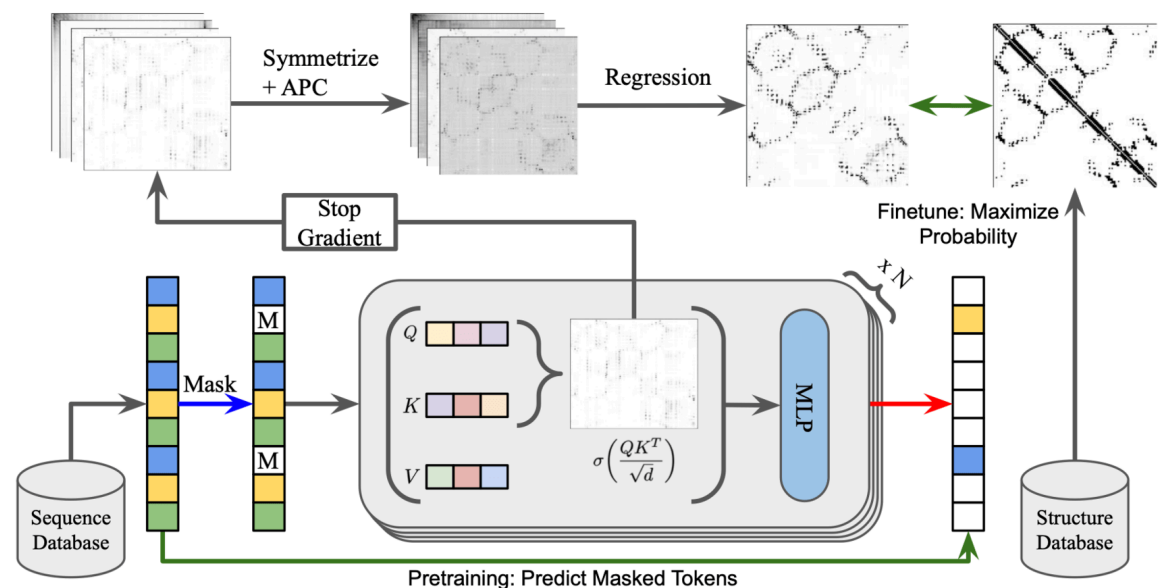


Figure 1: Contact prediction pipeline. The Transformer is first pretrained on sequences from a large database (Uniref50) via Masked Language Modeling. Once finished training, the attention maps are extracted, passed through symmetrization and average product correction, then into a regression. The regression is trained on a small number ($n \leq 20$) of proteins to determine which attention heads are informative. At test time, contact prediction from an input sequence can be done entirely on GPU in a single forward pass.

Evolutionary-scale prediction of atomic-level protein structure with a language model

ZEMING LIN , HALIL AKIN , ROSHAN RAO , BRIAN HIE , ZHONGKAI ZHU, WENTING LU, NIKITA SMETANIN, ROBERT VERKUIJL , ORI KABELI , [...], AN

ALEXANDER RIVES [+5 authors](#) [Authors Info & Affiliations](#)

Machine learning methods for protein structure prediction have taken advantage of the evolutionary information present in multiple sequence alignments to derive accurate structural information, but predicting structure accurately from a single sequence is much more difficult. Lin *et al.* trained **transformer protein language models** with up to 15 billion parameters on experimental and high-quality predicted structures and found that **information about atomic-level structure emerged in the model as it was scaled up.** They created ESMFold, a sequence-to-structure predictor that is nearly as accurate as alignment-based methods and considerably faster. The increased speed permitted the generation of a database, the ESM Metagenomic Atlas, containing more than 600 million metagenomic proteins. — MAF

Human-like object concept representation naturally in multimodal large language models

[Changde Du](#), [Kaicheng Fu](#), [Bincheng Wen](#), [Yi Sun](#), [Jie Peng](#), [Wei Wei](#), [Ying Gao](#), [Chuncheng Zhang](#), [Jinpeng Li](#), [Shuang Qiu](#), [Le Chang](#) & [Huiguang He](#) ✉

Nature Machine Intelligence **7**, 860–875 (2025) | [Cite this article](#)

Abstract

Understanding how humans conceptualize and categorize natural objects offers critical insights into perception and cognition. With the advent of large language models (LLMs), a key question arises: can these models develop human-like object representations from linguistic and multimodal data? Here we combined behavioural and neuroimaging analyses to explore the relationship between object concept representations in LLMs and human cognition. We collected 4.7 million triplet judgements from LLMs and multimodal LLMs to derive low-dimensional embeddings that capture the similarity structure of 1,854 natural objects. The resulting 66-dimensional embeddings were stable, predictive and exhibited semantic clustering similar to human mental representations. Remarkably, the dimensions underlying these embeddings were interpretable, suggesting that LLMs and multimodal LLMs develop human-like conceptual representations of objects. Further analysis showed strong alignment between model embeddings and neural activity patterns in brain regions such as the extrastriate body area, parahippocampal place area, retrosplenial cortex and fusiform face area. This provides compelling evidence that the object representations in LLMs, although not identical to human ones, share fundamental similarities that reflect key aspects of human conceptual knowledge. Our findings advance the understanding of machine intelligence and inform the development of more human-like artificial cognitive systems.

Success mode:



Success mode:

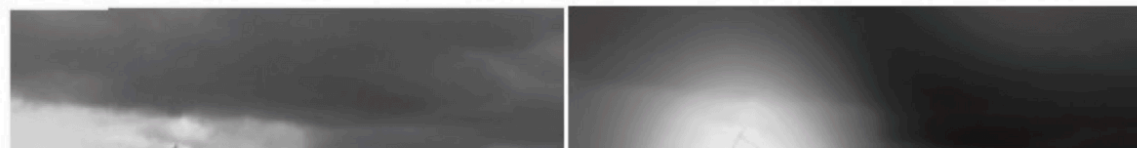
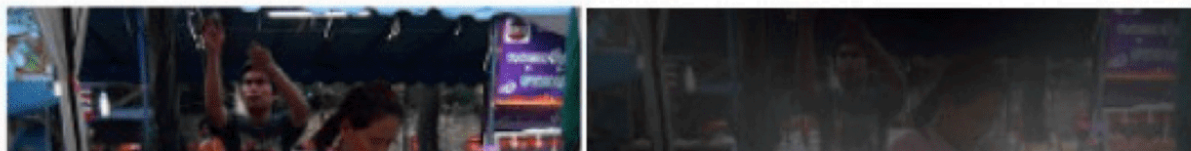


A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.

Failure mode:



A large white bird standing in a forest.

A woman holding a clock in her hand.

Success or Failure? mode:



Giannis Daras NeurIPS 2023
@giannis_daras

...

DALLE-2 has a secret language.

"Apoploe vesrreaitais" means birds.

"Contarra cctnxniam luryca tanniounons" means bugs or pests.

The prompt: "Apoploe vesrreaitais eating Contarra cctnxniam luryca tanniounons" gives images of birds eating bugs.

A thread (1/n)



Another example: "Two whales talking about food, with subtitles". We get an image with the text "Wa ch zod rea" written on it. Apparently, the whales are actually talking about their food in the DALLE-2 language. (4/n)

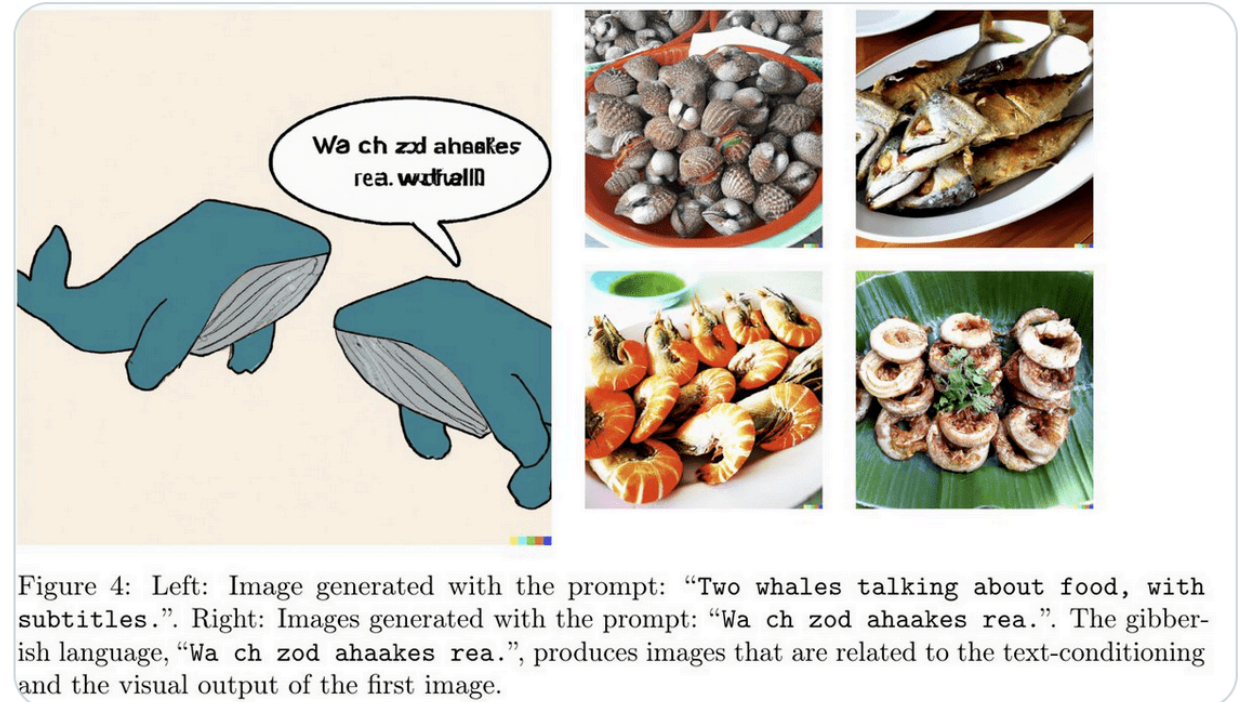


Figure 4: Left: Image generated with the prompt: "Two whales talking about food, with subtitles.". Right: Images generated with the prompt: "Wa ch zod ahaakes rea.". The gibberish language, "Wa ch zod ahaakes rea.", produces images that are related to the text-conditioning and the visual output of the first image.

Summary

- Transformers combine ideas from earlier architectures (convolution, ReLU, residual connections) with new innovations: embedding and attention layers.
- Transformers start with generic hard-coded **embeddings** and, block-by-block, create increasingly context-aware embeddings.
- **Attention** enables massive parallelism: each head, each q, k, v sequence, each attention score, and each output are all computed in parallel.
- (Because the architecture is input-agnostic — "attention is all you need" — transformers have become one model that rules across language, vision, and multi-modal tasks.)