CMPU 366 · Natural Language Processing

# Attention and Transformers

3 November 2025

# BuzzFeed

QUIZ

Community · Posted on May 14, 2019

✈ **Subscribe** to Quizzes Newsletter ∨

# Can You Get Mr. Darcy From "Pride And Prejudice" To Propose To You?

If you can manage to avoid Mr. Collins, of course.

by **rsk17**
Community Contributor

✦ 820 points

Approved and edited by
**BuzzFeed Community Team**

💬 View All 14 Comments

# What would you do if a wealthy man moved

Breadth

Individual work

Learn and practice foundational
NLP concepts

Rule-based systems (regular expressions)

Language models

Logistic regression

Word embeddings

Deep learning

Breadth

Individual work

Learn and practice foundational NLP concepts

- Rule-based systems (regular expressions)
- Language models
- Logistic regression
- Word embeddings
- Deep learning

Depth

Group work

Deep and sustained creative problem-solving to build complex projects around one topic

# Special topics

This is a chance to explore a topic we didn't cover yet. It's a good idea for this to be on the general area you'd like your final project to be in – though it doesn't need to be!

Sign up for a topic/day before the start of next class.

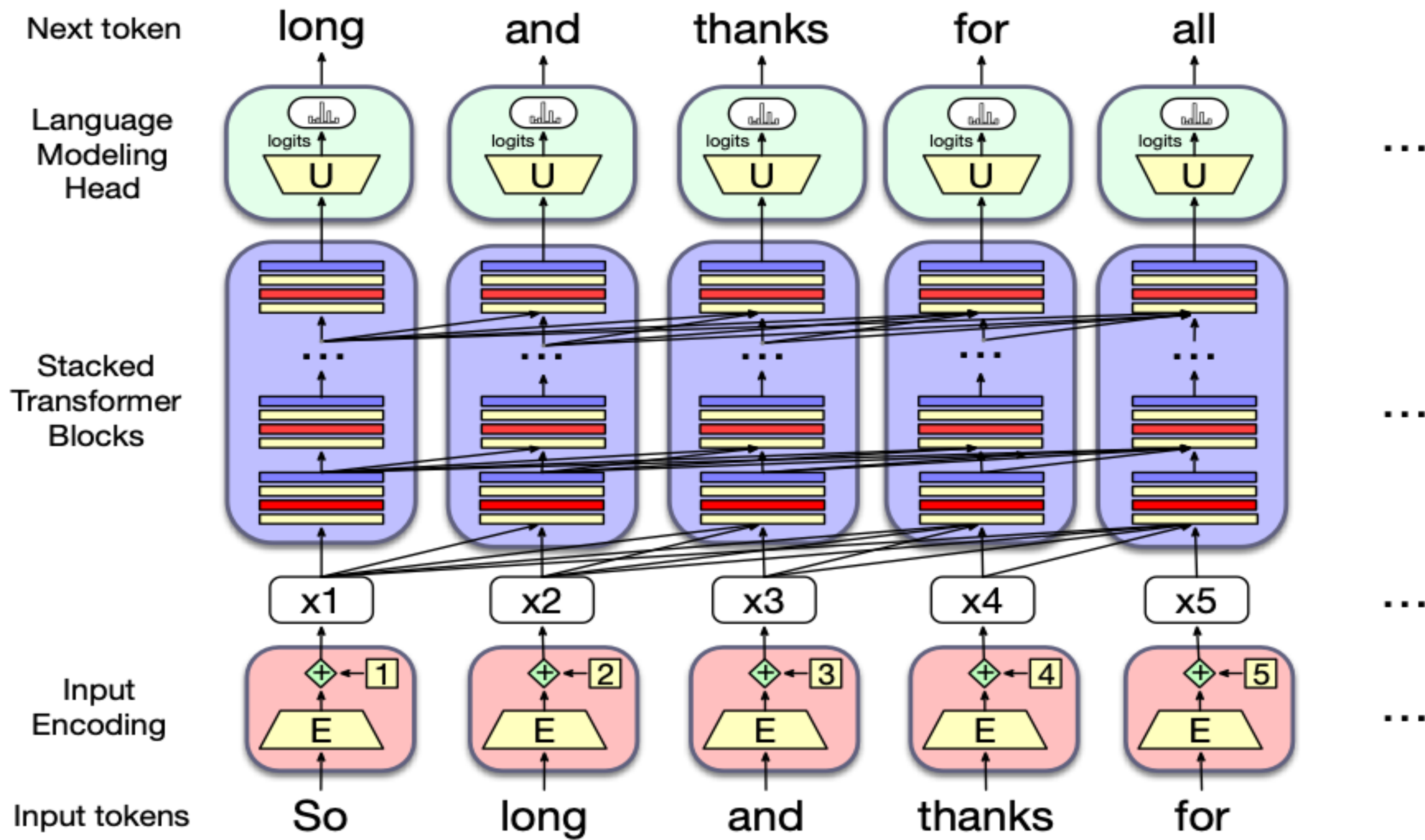# Projects

(Short) proposals will be due next Tuesday.

We'll use class time next Wednesday to start working on projects.

LLMs are built out of transformers.

A *transformer* is a specific kind of network architecture – like a fancier feedforward network, but based on *attention*.

# Attention

Let's start by considering the embeddings for an individual word from a particular layer.

The problem with embeddings like Word2vec is they're *static*; the embedding for a word doesn't reflect how its meaning changes in *context*.

*The chicken didn't cross the road because it was too tired.*

What does *it* mean in a static embedding?

*The chicken didn't cross the road because it was too tired.*

A better representation of the meaning of a word would be different in different contexts!

This is called a *contextual embedding*.

*The chicken didn't cross the road because it was too tired.*

*The chicken didn't cross the road because it was too wide.*

*The chicken didn't cross the road because it …*

A transformer uses *attention* to build a contextual representation of a word by selectively integrating information from all the neighboring words.

columns corresponding to input tokens



Layer k+1

The chicken didn't cross the road because **it** was too tired

**self-attention distribution**

Layer k

The chicken didn't cross the road because it was too tired

Outputs

$\mathbf{a}_1$

Self-attention layer

attention

Inputs

$\mathbf{x}_1$

Outputs

$\mathbf{a}_1$

Self-attention layer

attention

Inputs

$\mathbf{x}_1$    $\mathbf{x}_2$

Outputs

$\mathbf{a}_1$

Self-attention layer

attention        attention

Inputs

$\mathbf{x}_1$        $\mathbf{x}_2$

Outputs

$\mathbf{a}_1$   $\mathbf{a}_2$

Self-attention layer

attention   attention

Inputs

$\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{x}_3$

Outputs

Self-attention layer

Inputs

$\mathbf{a}_1$ $\mathbf{a}_2$

attention attention attention

$\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_3$

Outputs

Self-attention layer

$\mathbf{a}_1$ $\mathbf{a}_2$ $\mathbf{a}_3$

attention attention attention

Inputs

$\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_3$

Outputs

$\mathbf{a}_1$  $\mathbf{a}_2$  $\mathbf{a}_3$

Self-attention layer
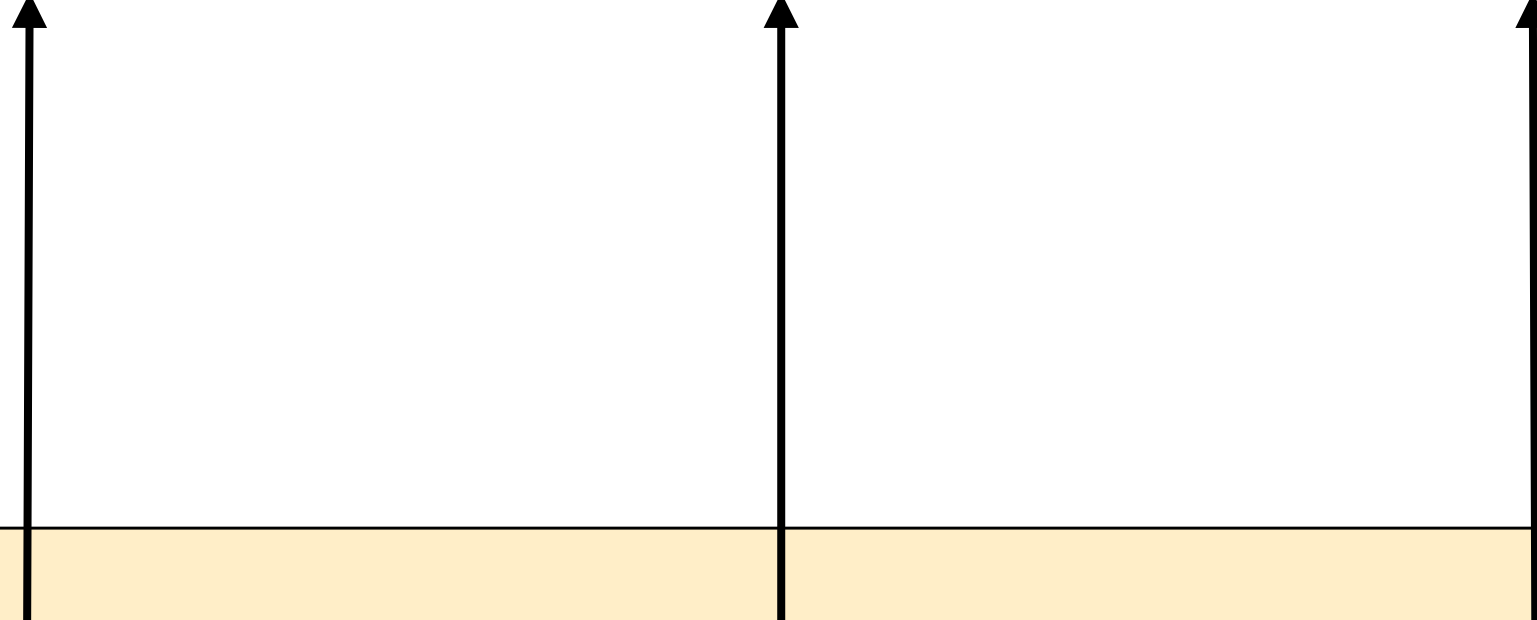
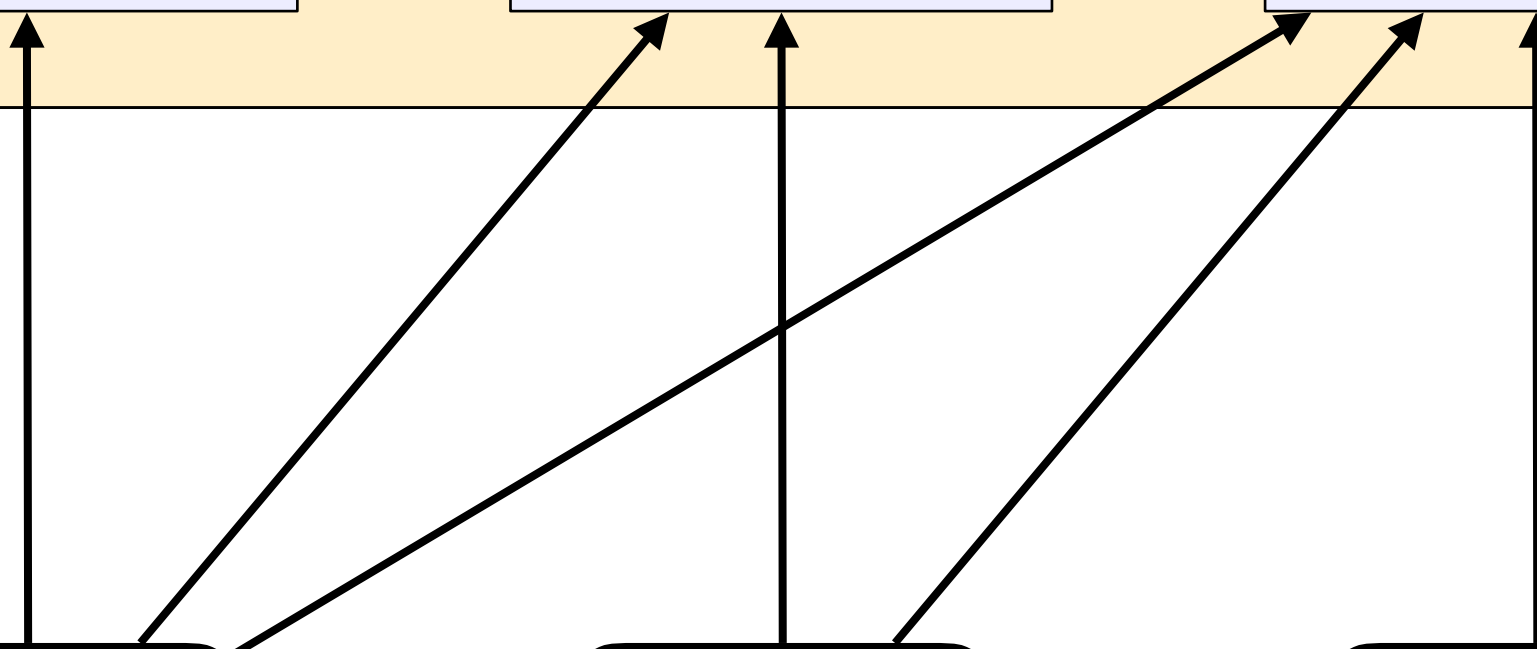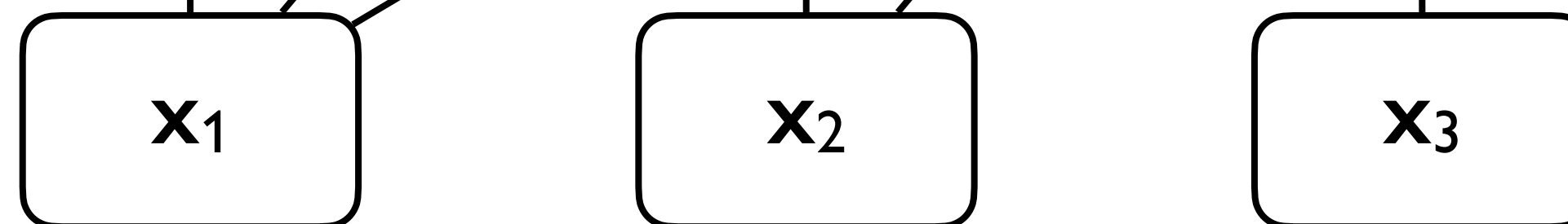attention  attention  attention

Inputs

$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$  $\mathbf{x}_4$

Outputs

Self-attention layer

Inputs

$a_1$  $a_2$  $a_3$  $a_4$

attention  attention  attention  attention

$x_1$  $x_2$  $x_3$  $x_4$

Outputs: $\mathbf{a}_1$, $\mathbf{a}_2$, $\mathbf{a}_3$, $\mathbf{a}_4$

Self-attention layer: attention, attention, attention, attention

Inputs: $\mathbf{x}_1$, $\mathbf{x}_2$, $\mathbf{x}_3$, $\mathbf{x}_4$, $\mathbf{x}_5$

Outputs: $a_1$, $a_2$, $a_3$, $a_4$

Self-attention layer: attention, attention, attention, attention, attention

Inputs: $x_1$, $x_2$, $x_3$, $x_4$, $x_5$
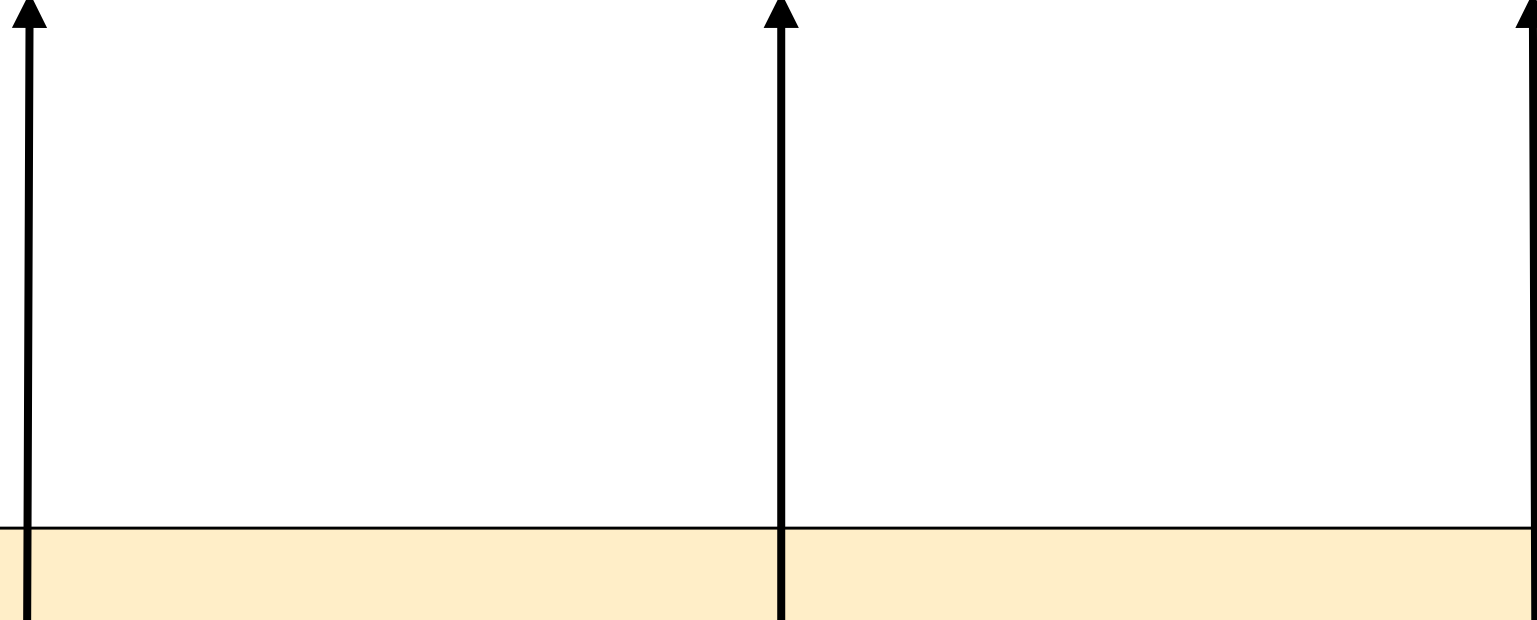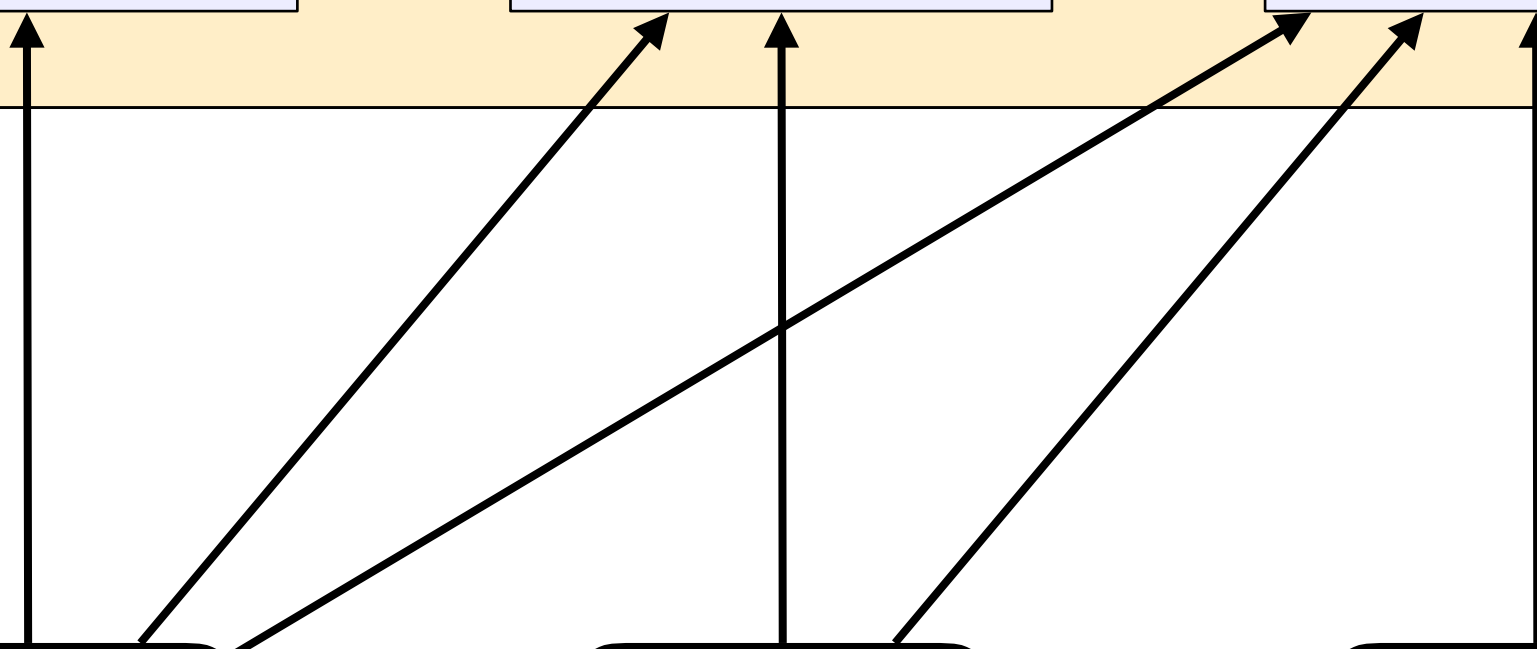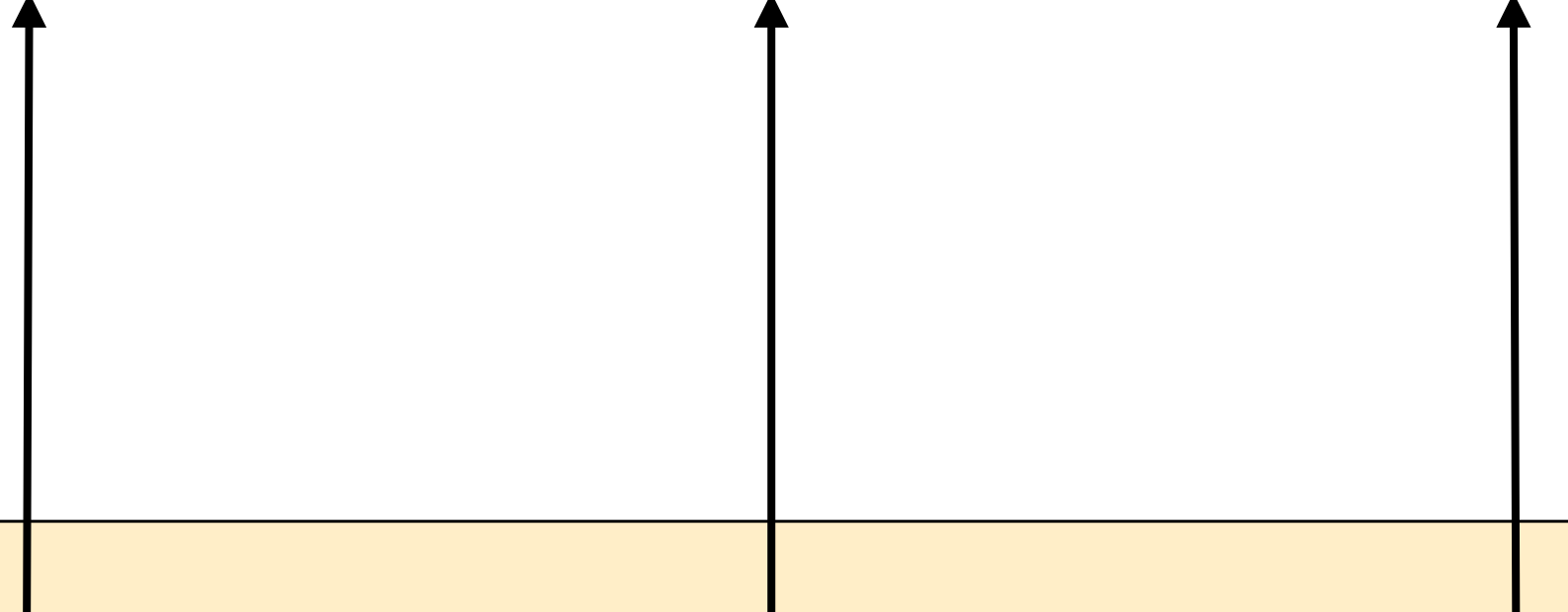
Outputs

Self-attention layer

Inputs

$\mathbf{a}_1$ $\mathbf{a}_2$ $\mathbf{a}_3$ $\mathbf{a}_4$ $\mathbf{a}_5$

attention attention attention attention attention

$\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_3$ $\mathbf{x}_4$ $\mathbf{x}_5$
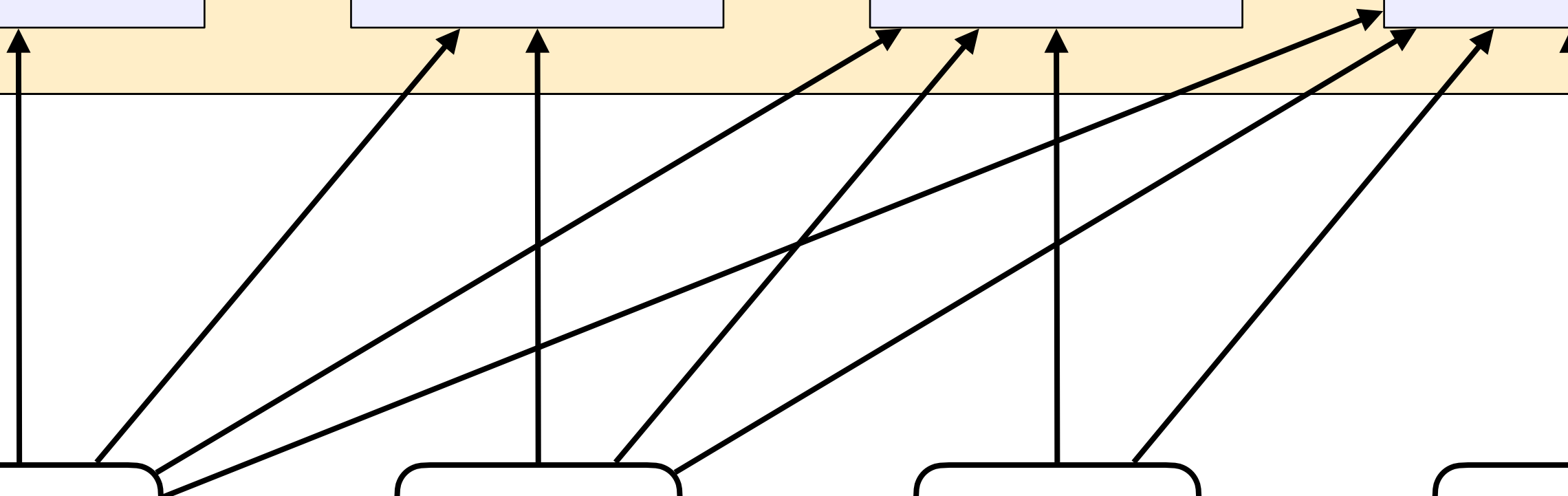
Given a sequence of token embeddings

$$\mathbf{x}_1 \quad \mathbf{x}_2 \quad \mathbf{x}_3 \quad \mathbf{x}_4 \quad \mathbf{x}_5 \quad \mathbf{x}_6 \quad \mathbf{x}_7 \quad \mathbf{x}_i$$

we can produce a sum of the embeddings weighted by their similarity to $\mathbf{x}_i$:

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{x}_j$$

where the weight

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \forall j \leq i,$$

and we can simply define

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j.$$

columns corresponding to input tokens

Layer k+1

The chicken didn't cross the road because **it** was too tired

**self-attention distribution**

Layer k

The chicken didn't cross the road because it was too tired

$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$  $\mathbf{x}_4$  $\mathbf{x}_5$  $\mathbf{x}_6$  $\mathbf{x}_7$  $\mathbf{x}_i$

An actual attention head is slightly more complicated.

Instead of using vectors like $\mathbf{x}_i$ and $\mathbf{x}_4$ directly, we'll represent three separate *roles* each vector $\mathbf{x}_i$ plays:

*Query*: As the *current element* being compared to the preceding inputs

*Key*: As a *preceding input* that is being compared to the current element to determine a similarity

*Value*: a value of a preceding element that gets *weighted and summed*.

Layer k+1

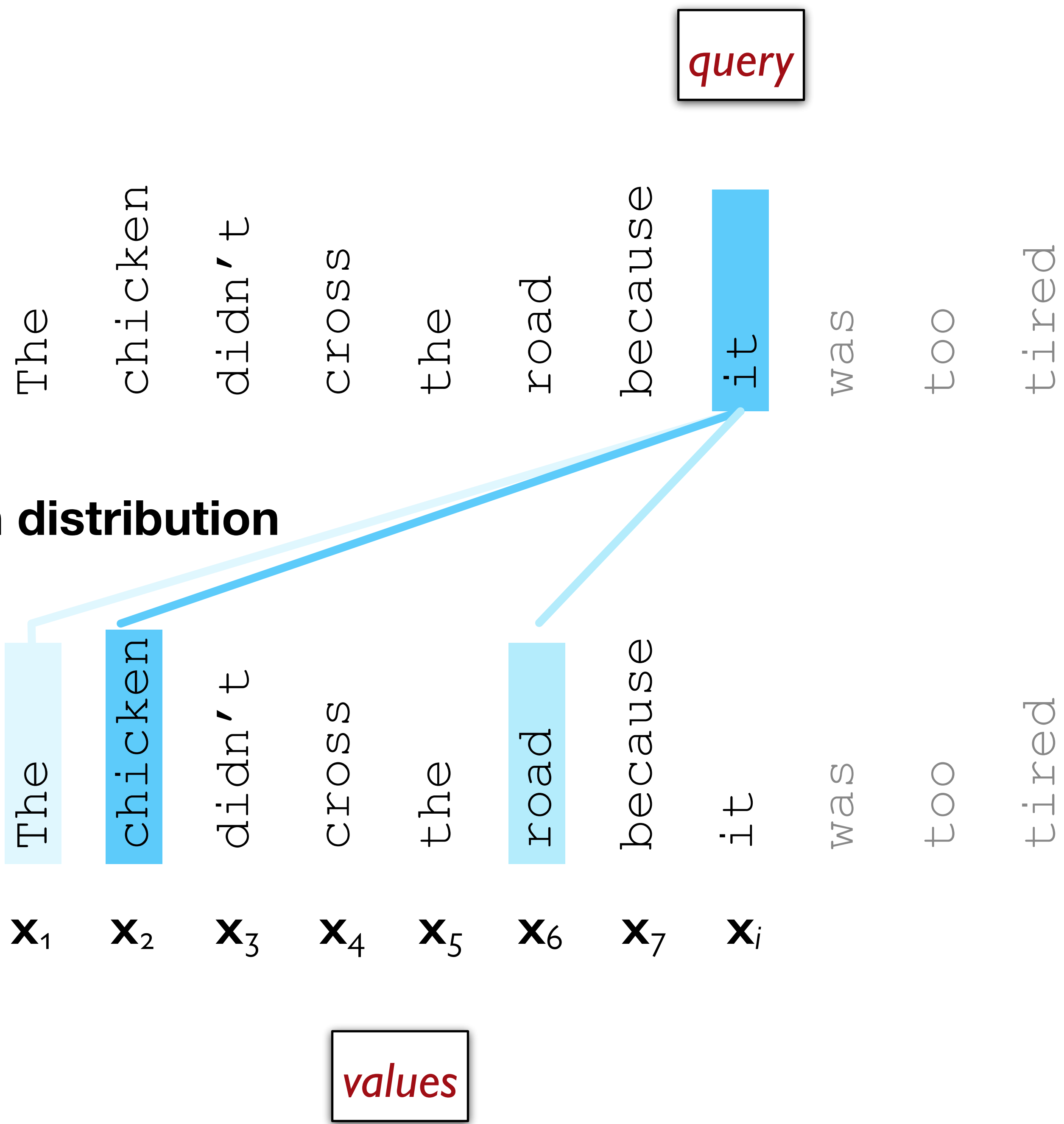The chicken didn't cross the road because **it** was too tired
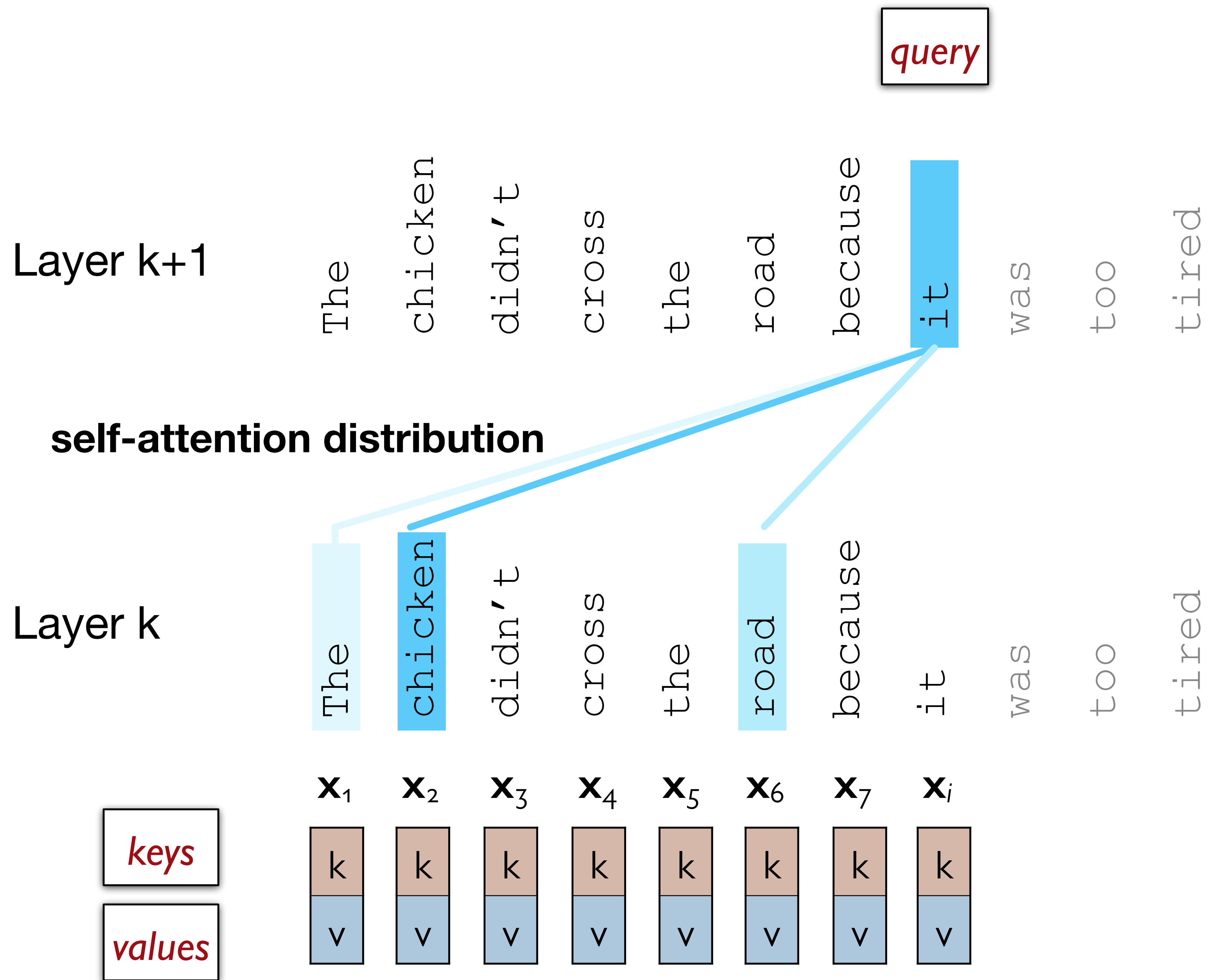
*query*

**self-attention distribution**

Layer k

The chicken didn't cross the road because it was too tired

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$  $x_6$  $x_7$  $x_i$

*values*

query

Layer k+1

The chicken didn't cross the road because it was too tired

**self-attention distribution**

Layer k

The chicken didn't cross the road because it was too tired

$\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_3$ $\mathbf{x}_4$ $\mathbf{x}_5$ $\mathbf{x}_6$ $\mathbf{x}_7$ $\mathbf{x}_i$

keys

k k k k k k k k

values

v v v v v v v v

We'll use weight matrices $\mathbf{W^Q}$, $\mathbf{W^K}$, and $\mathbf{W^V}$ to project each vector $\mathbf{x}_i$ into a representation of its role as a

query:    $\mathbf{q}_i = \mathbf{x}_i\mathbf{W^Q}$

key:      $\mathbf{k}_i = \mathbf{x}_i\mathbf{W^K}$

value:    $\mathbf{v}_i = \mathbf{x}_i\mathbf{W^V}$

To compute the similarity of the current element $\mathbf{x}_i$ with some prior element $\mathbf{x}_j$, we'll use the dot product between the projections $\mathbf{q}_i$ and $\mathbf{k}_j$.

And instead of summing up $\mathbf{x}_j$, we'll sum up the projection $\mathbf{v}_j$.

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^{\mathbf{Q}}$$

$$\mathbf{k}_j = \mathbf{x}_j \mathbf{W}^{\mathbf{K}}$$

$$\mathbf{v}_j = \mathbf{x}_j \mathbf{W}^{\mathbf{V}}$$

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j \qquad\qquad \text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{q}_i \cdot \mathbf{k}_j$$

$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \ \forall j \leq i \qquad\qquad \alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \ \forall j \leq i$$

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{x}_j \qquad\qquad\qquad \mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W^Q}$$

$$\mathbf{k}_j = \mathbf{x}_j \mathbf{W^K}$$

$$\mathbf{v}_j = \mathbf{v}_j \mathbf{W^V}$$

*$d_k$ is the dimensionality of the query and key vectors*

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}} \longleftarrow \text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{q}_i \cdot \mathbf{k}_j$$

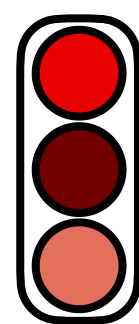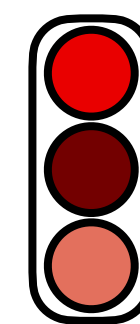$$\alpha_{ij} = \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \ \ \forall j \leq i$$

*This is a practical change to avoid numerical issues during training.*

$$\mathbf{a}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{x}_1$

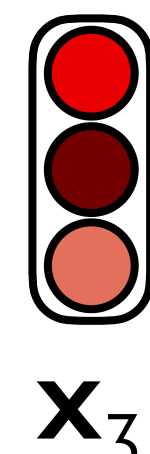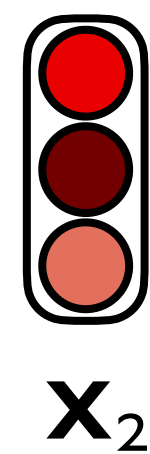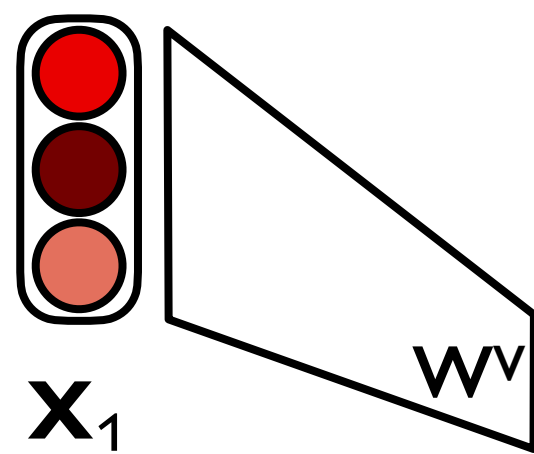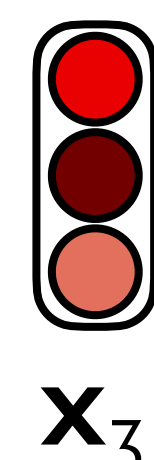$\mathbf{x}_2$

$\mathbf{x}_3$

Project each $\mathbf{x}_i$ into key, query, and value vectors

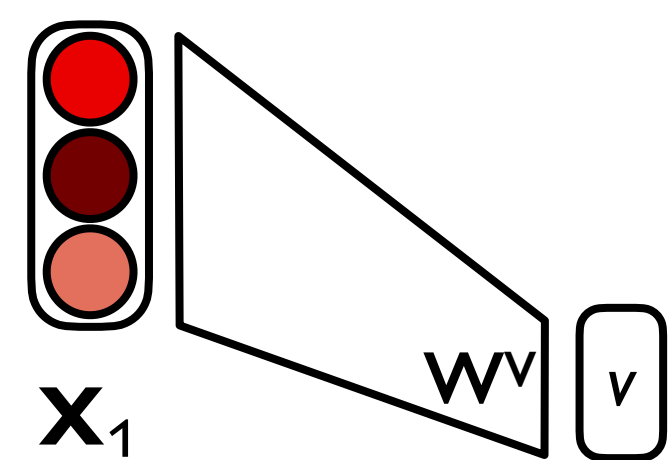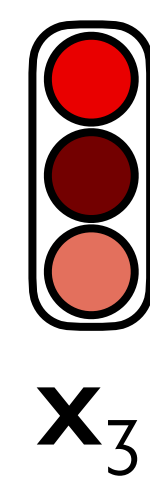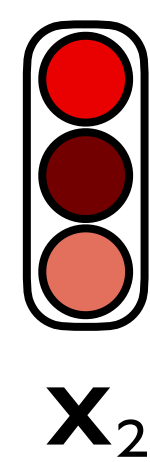$\mathbf{x}_1$     $\mathbf{x}_2$     $\mathbf{x}_3$

*Project each $\mathbf{x}_i$ into key, query, and value vectors*
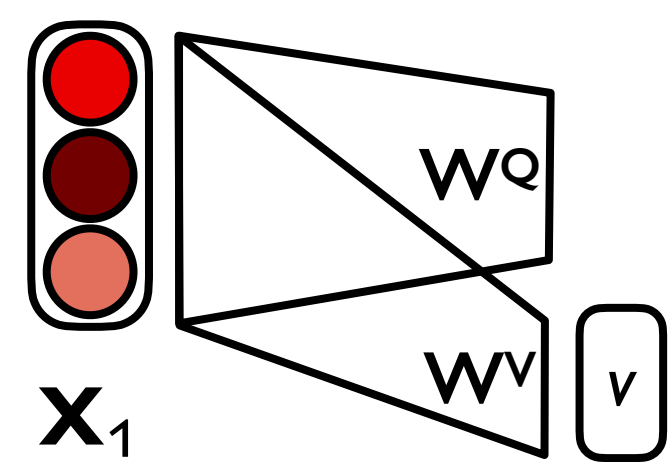
$\mathbf{x}_1$ $W^v$ $v$

$\mathbf{x}_2$

$\mathbf{x}_3$

Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{x}_1$  W$^Q$  W$^V$  v

$\mathbf{x}_2$

$\mathbf{x}_3$

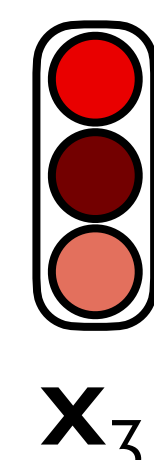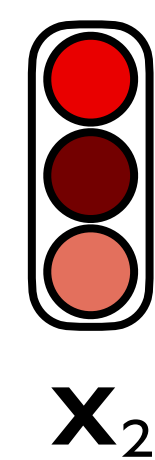Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{x}_1$ $W^Q$ $q$ $W^V$ $v$

$\mathbf{x}_2$

$\mathbf{x}_3$

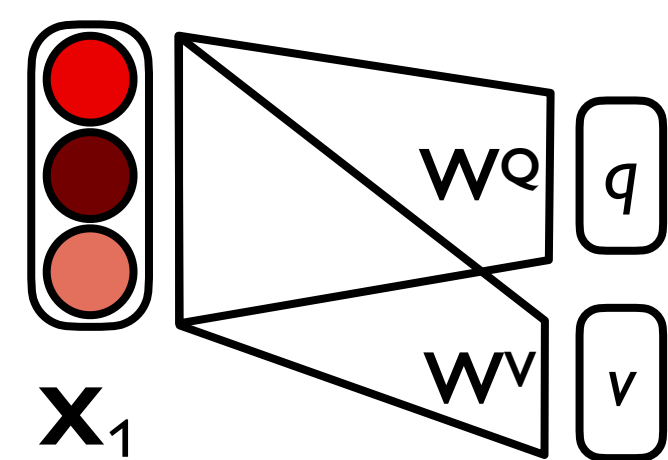Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{x}_1$  $W^K$  $W^Q$  $q$  $W^V$  $v$  $\mathbf{x}_2$  $\mathbf{x}_3$

Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{W}^K$   $k$

$\mathbf{W}^Q$   $q$

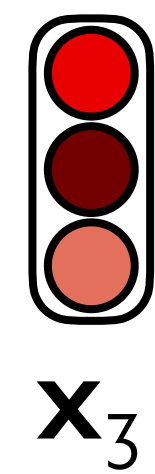$\mathbf{W}^V$   $v$

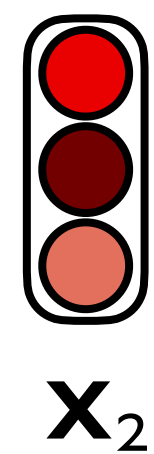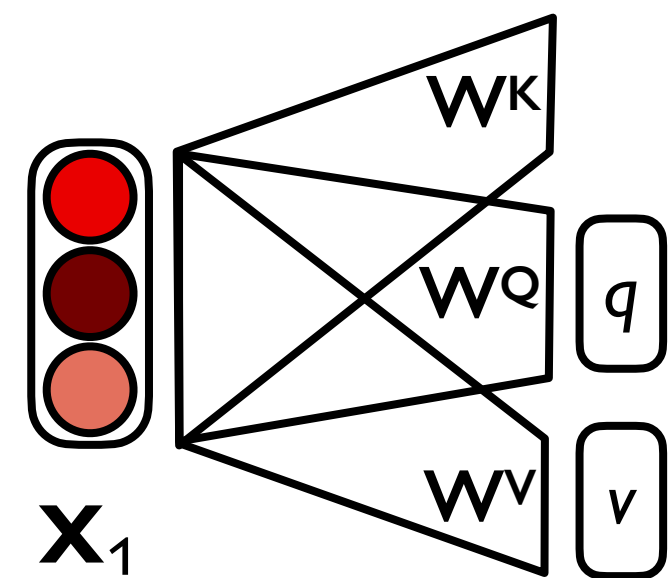$\mathbf{x}_1$         $\mathbf{x}_2$         $\mathbf{x}_3$

*Project each $\mathbf{x}_i$ into key, query, and value vectors*

$\mathbf{x}_1$  $W^K$  $k$  $W^Q$  $q$  $W^V$  $v$

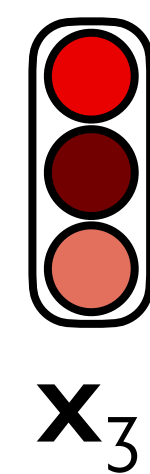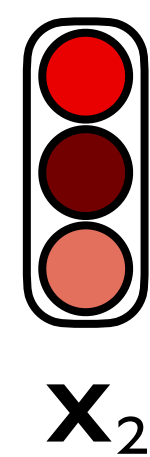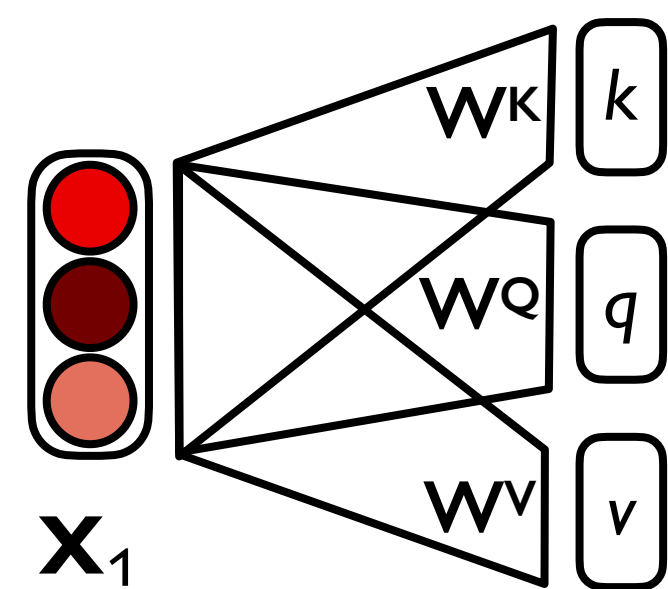$\mathbf{x}_2$  $W^K$  $k$  $W^Q$  $q$  $W^V$  $v$

$\mathbf{x}_3$

*Project each $\mathbf{x}_i$ into key, query, and value vectors*

$\mathbf{x}_1$    $W^K$  $k$    $W^Q$  $q$    $W^V$  $v$

$\mathbf{x}_2$    $W^K$  $k$    $W^Q$  $q$    $W^V$  $v$

$\mathbf{x}_3$    $W^K$  $k$    $W^Q$  $q$    $W^V$  $v$

*Compare **x**$_3$'s query with the keys for **x**$_1$, **x**$_2$, and **x**$_3$.*

*Project each **x**$_i$ into key, query, and value vectors*

**x**$_1$   $W^K$  $k$   $W^Q$  $q$   $W^V$  $v$

**x**$_2$   $W^K$  $k$   $W^Q$  $q$   $W^V$  $v$

**x**$_3$   $W^K$  $k$   $W^Q$  $q$   $W^V$  $v$

Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.

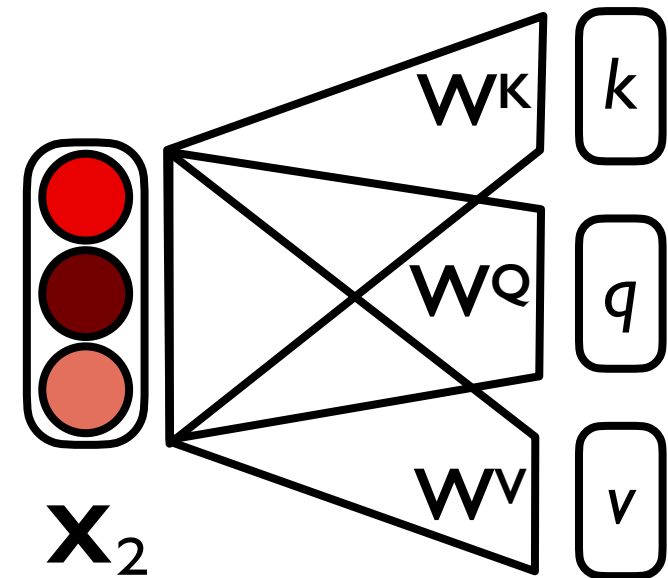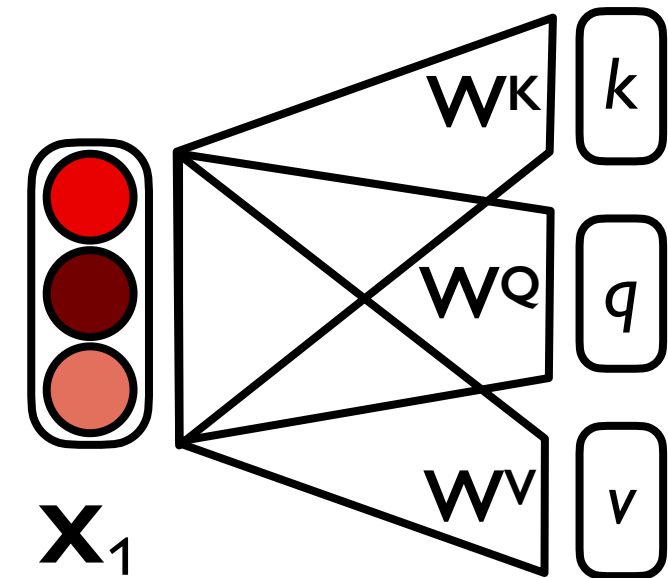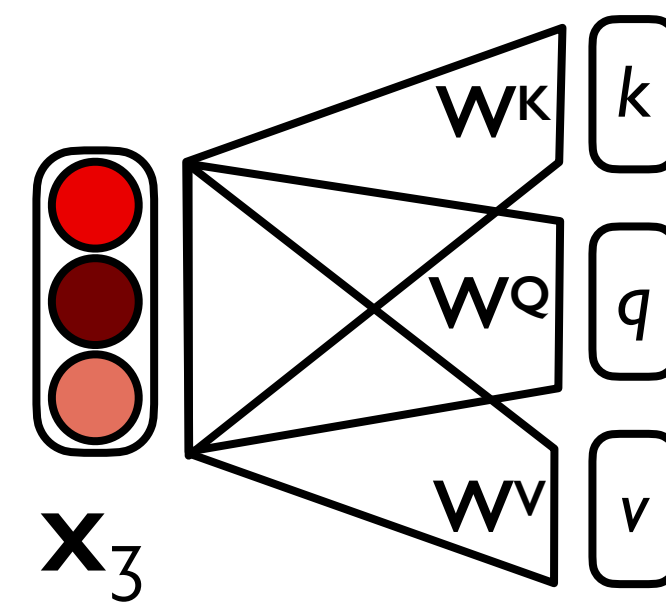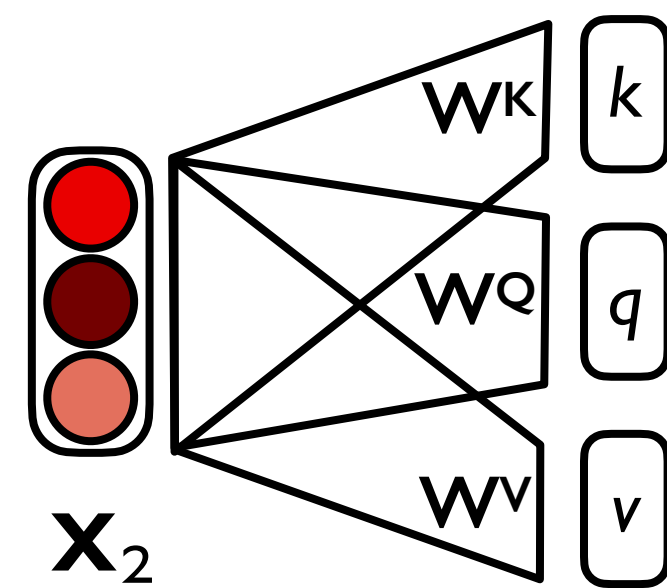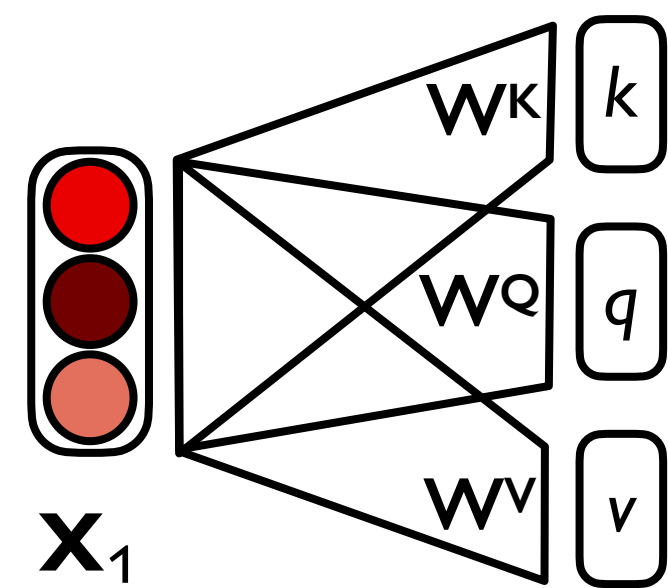Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{x}_1$

$\mathbf{x}_2$

$\mathbf{x}_3$

Compare **x**₃'s query with the keys for **x**₁, **x**₂, and **x**₃.

Project each **x**ᵢ into key, query, and value vectors

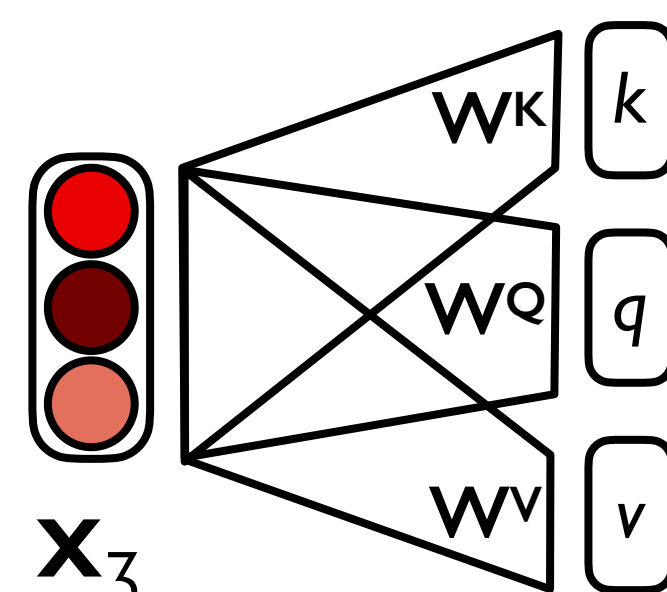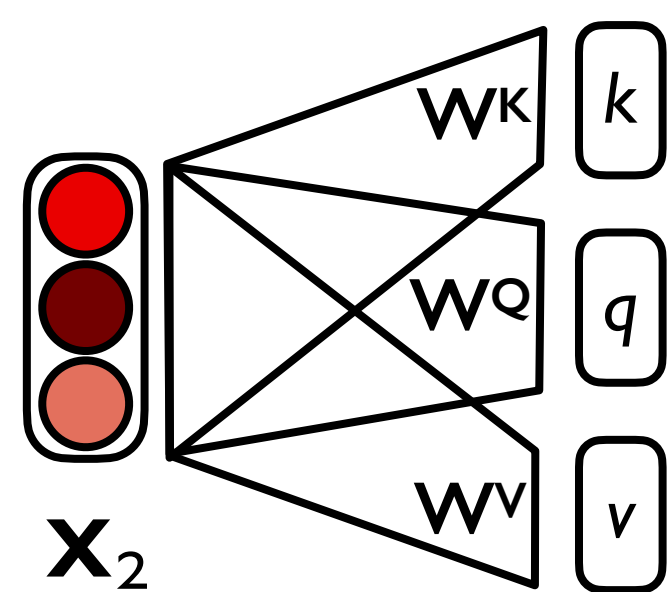*Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.*

*Project each $\mathbf{x}_i$ into key, query, and value vectors*

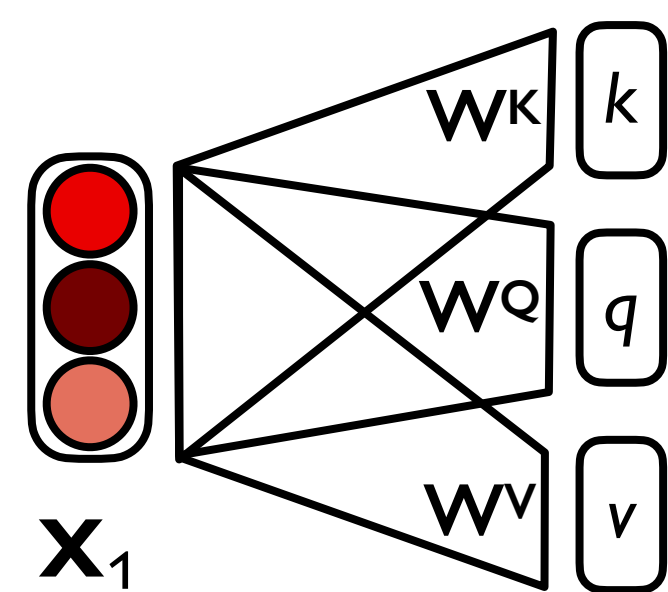Divide score by $\sqrt{d_k}$

Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.

Project each $\mathbf{x}_i$ into key, query, and value vectors

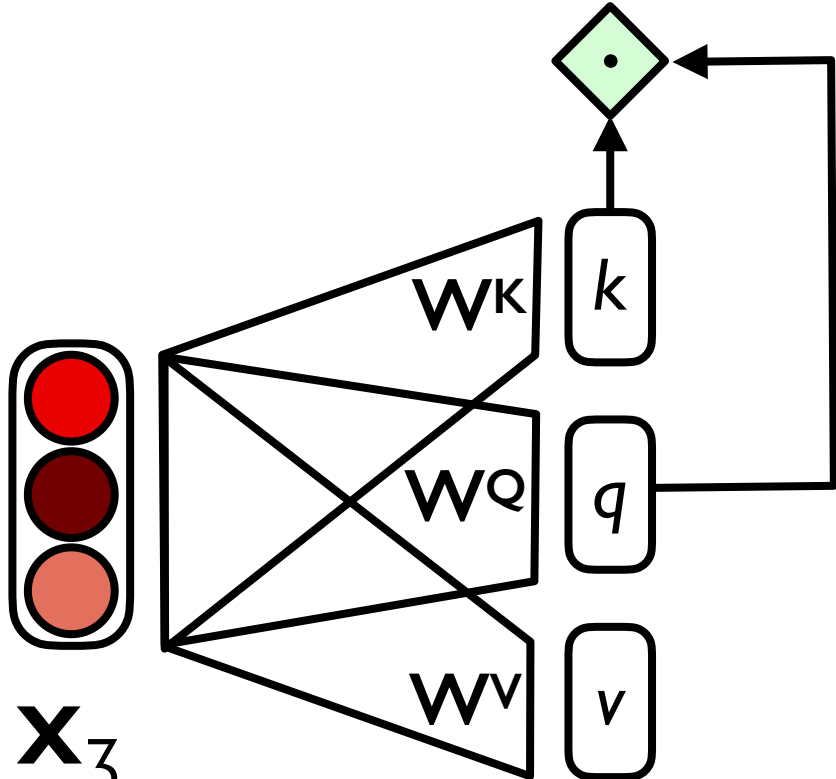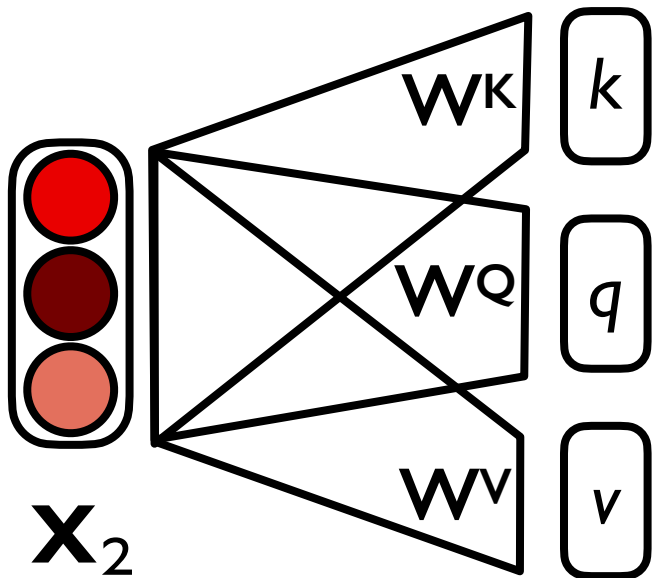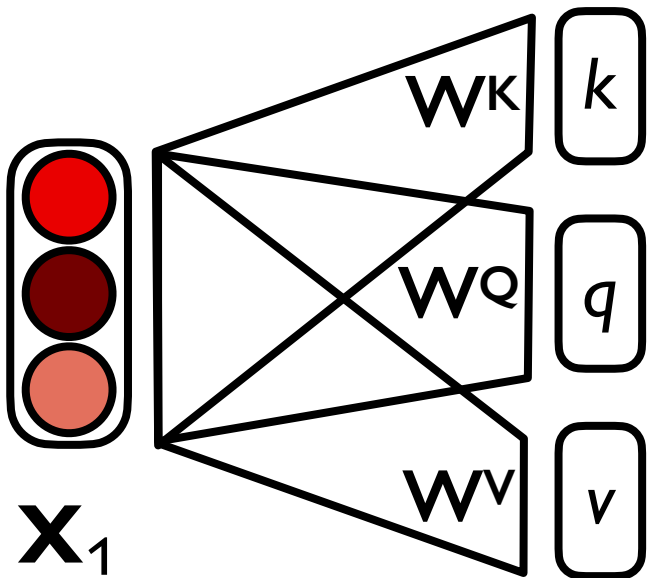Turn into $a_{ij}$ weights via softmax

Divide score by $\sqrt{d_k}$

Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.

Project each $\mathbf{x}_i$ into key, query, and value vectors

$\sqrt{d_k}$

$\sqrt{d_k}$

$\sqrt{d_k}$

$W^K$   $k$

$W^Q$   $q$

$W^V$   $v$

$\mathbf{x}_1$

$\mathbf{x}_2$

$\mathbf{x}_3$

Weigh each value vector

Turn into $a_{ij}$ weights via softmax

Divide score by $\sqrt{d_k}$

Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.

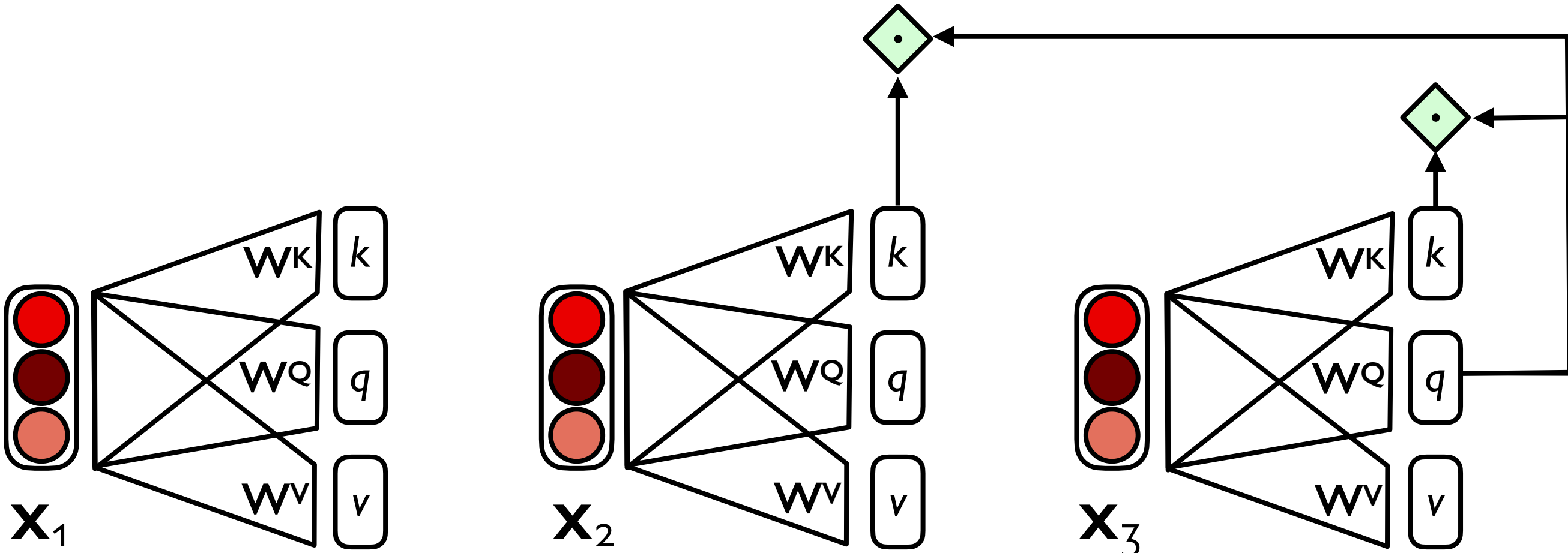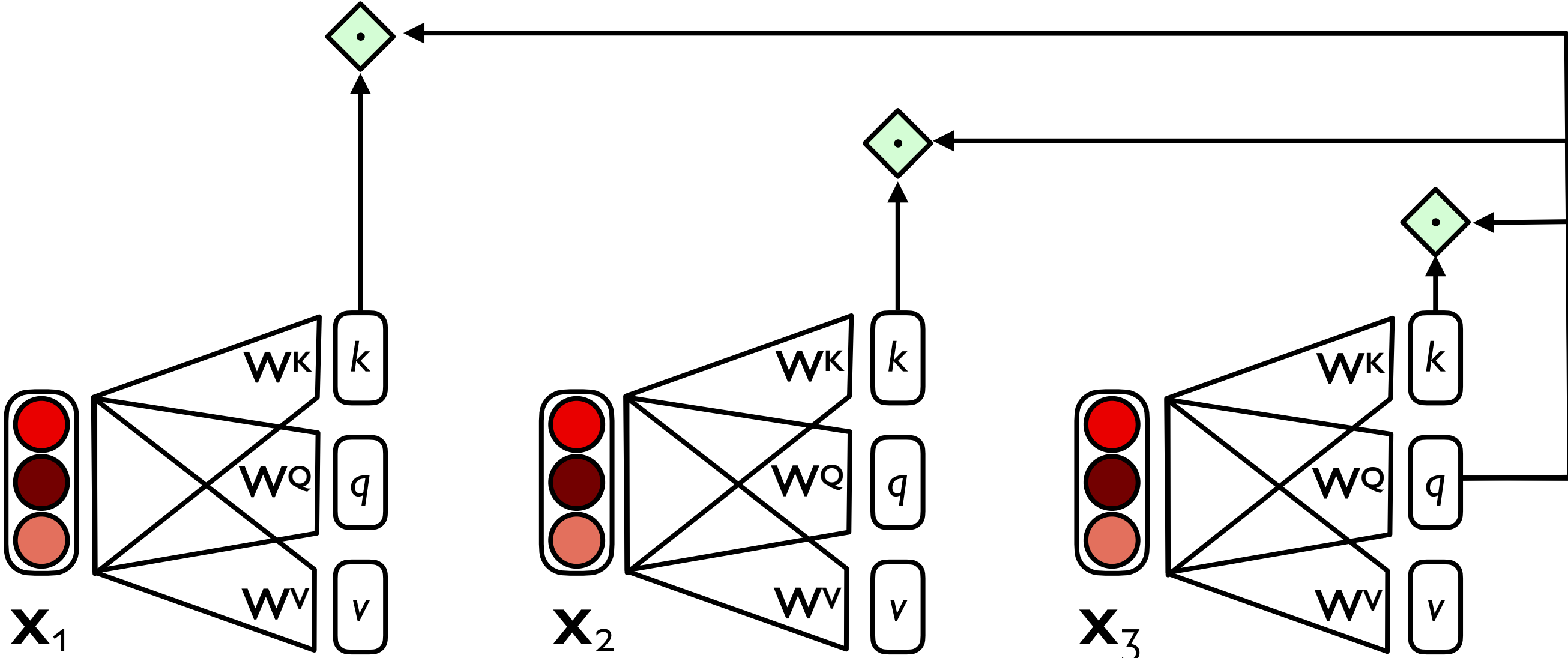Project each $\mathbf{x}_i$ into key, query, and value vectors

$a_{3,1}$  $a_{3,2}$  $a_{3,3}$

$\sqrt{d_k}$  $\sqrt{d_k}$  $\sqrt{d_k}$

$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$

$W^K$  $k$
$W^Q$  $q$
$W^V$  $v$

Sum the weighted value vectors

Weigh each value vector

Turn into $a_{ij}$ weights via softmax

Divide score by $\sqrt{d_k}$

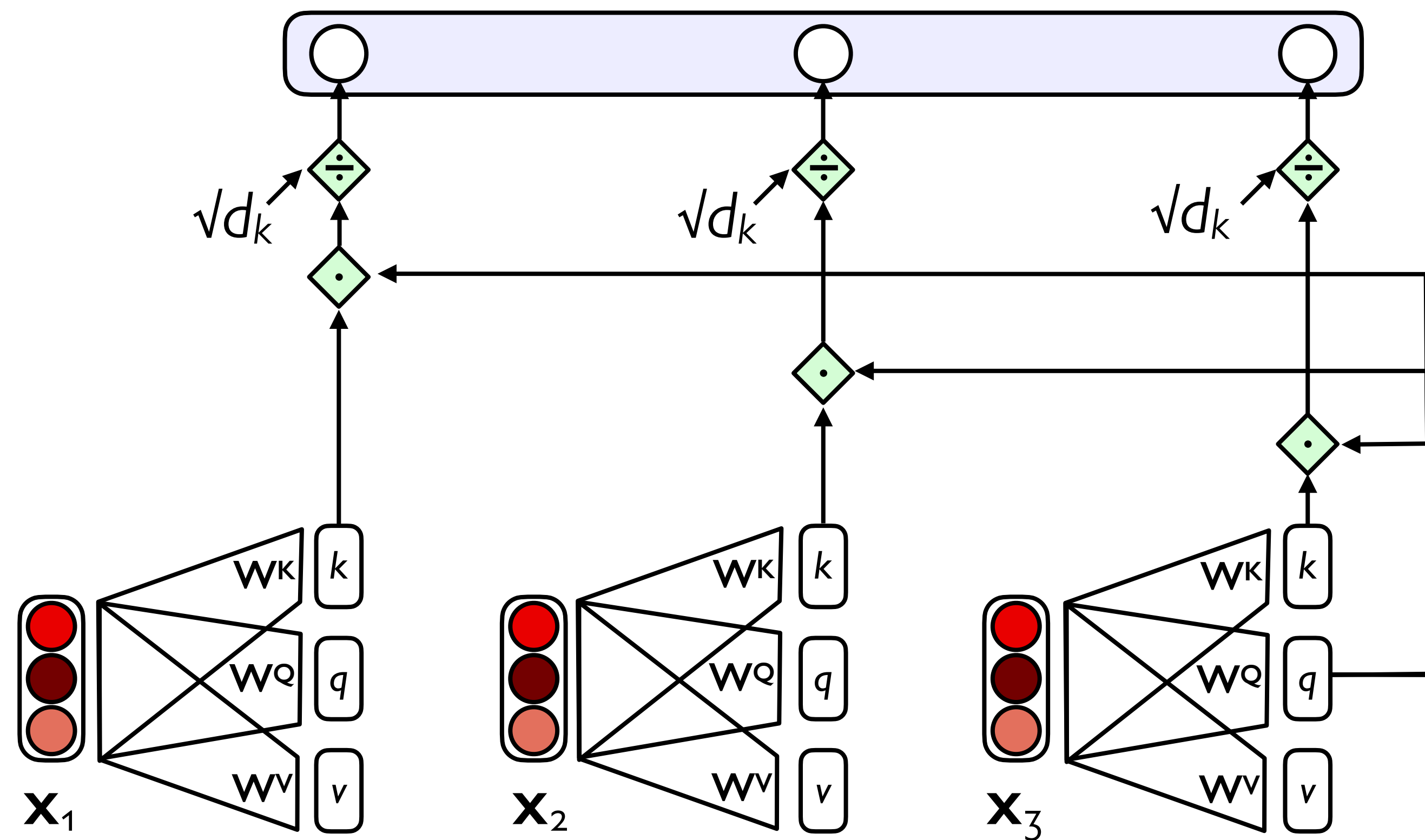Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.

Project each $\mathbf{x}_i$ into key, query, and value vectors

$\Sigma$

$a_{3,1}$   $a_{3,2}$   $a_{3,3}$

$\sqrt{d_k}$   $\sqrt{d_k}$   $\sqrt{d_k}$

$\mathbf{W}^K$   $k$
$\mathbf{W}^Q$   $q$
$\mathbf{W}^V$   $v$

$\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{x}_3$

Output of self-attention

Sum the weighted value vectors

Weigh each value vector

Turn into $a_{ij}$ weights via softmax
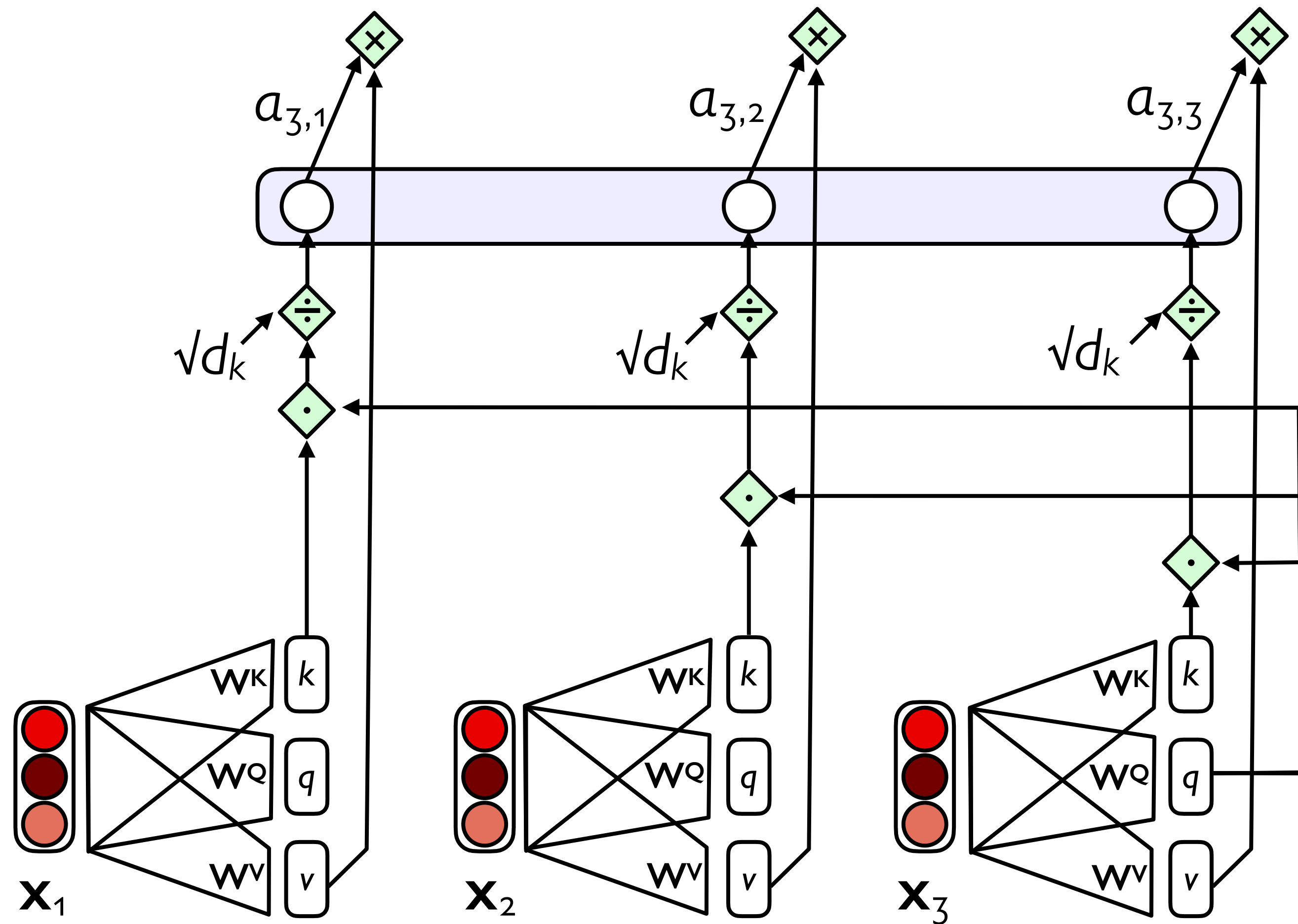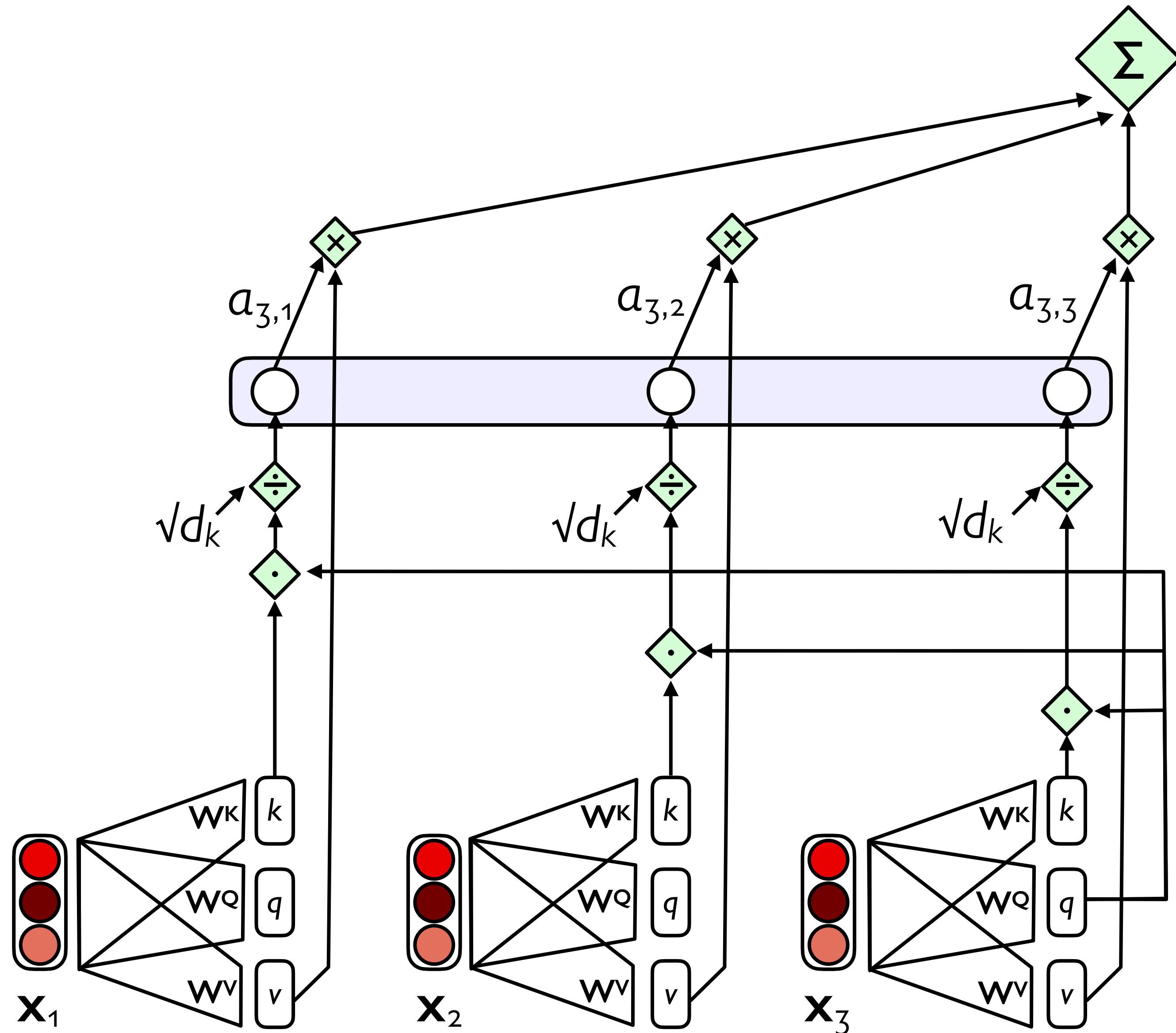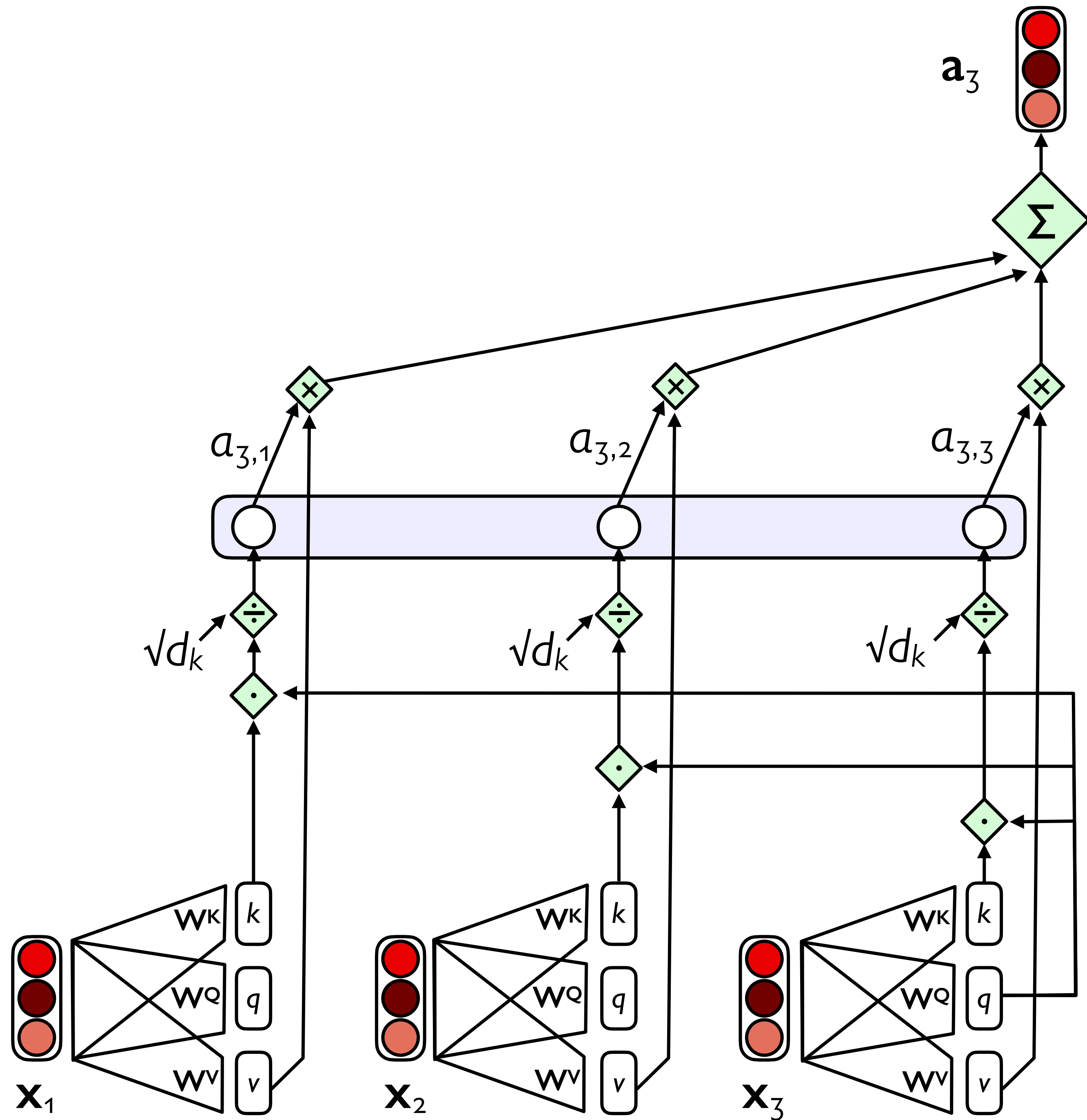
Divide score by $\sqrt{d_k}$

Compare $\mathbf{x}_3$'s query with the keys for $\mathbf{x}_1$, $\mathbf{x}_2$, and $\mathbf{x}_3$.

Project each $\mathbf{x}_i$ into key, query, and value vectors

$\mathbf{a}_3$

$a_{3,1}$

$a_{3,2}$

$a_{3,3}$

$\sqrt{d_k}$

$\sqrt{d_k}$

$\sqrt{d_k}$

$\mathbf{x}_1$

$\mathbf{x}_2$

$\mathbf{x}_3$

$W^K$ $k$

$W^Q$ $q$

$W^V$ $v$

In practice, instead of one attention head, we'll have lots of them: *multi-head attention*.

Why? Each head might be attending to the context for different purposes – different linguistic relationships or patterns in the context

# Summary

Attention is a method for enriching the representation of a token by incorporating contextual information
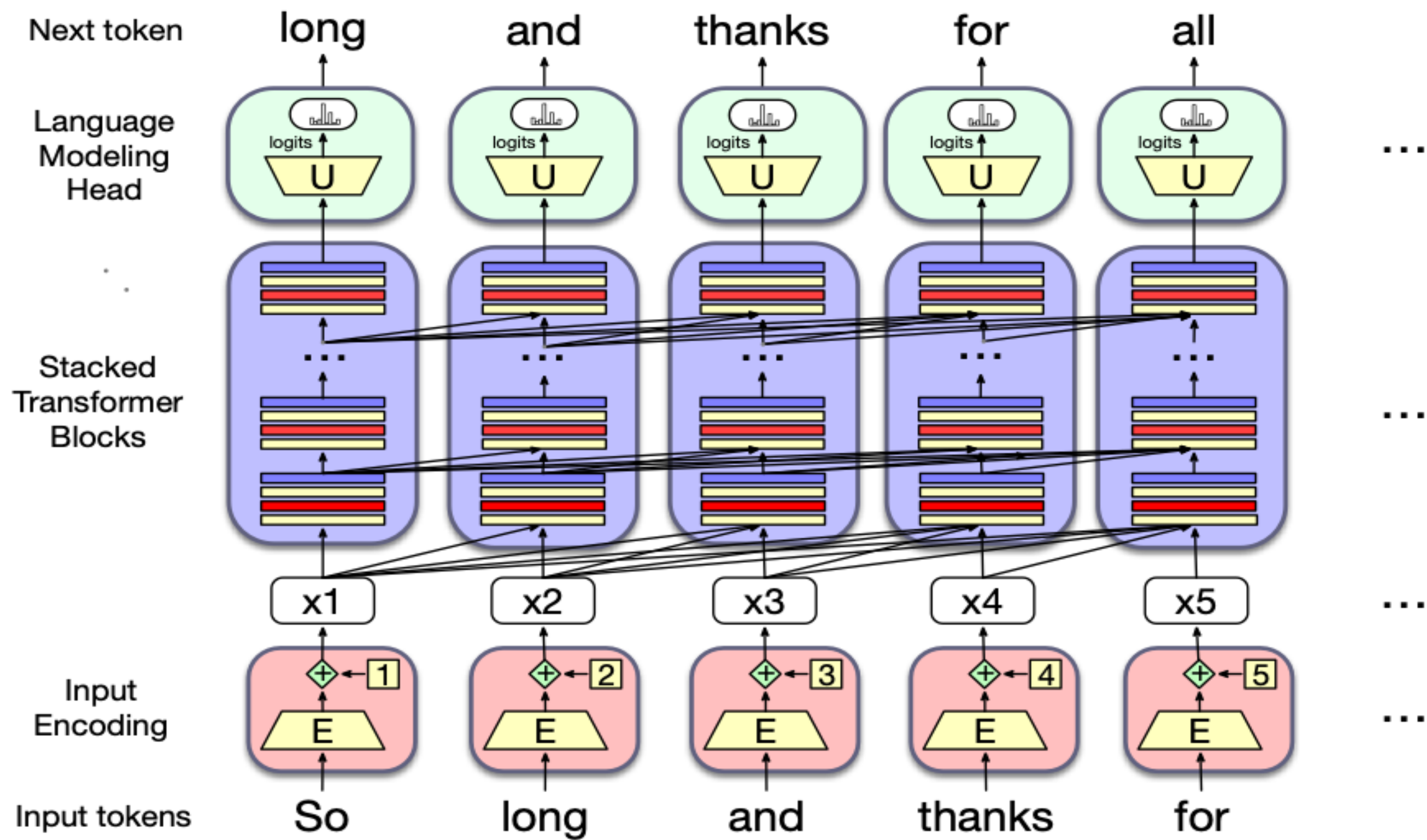
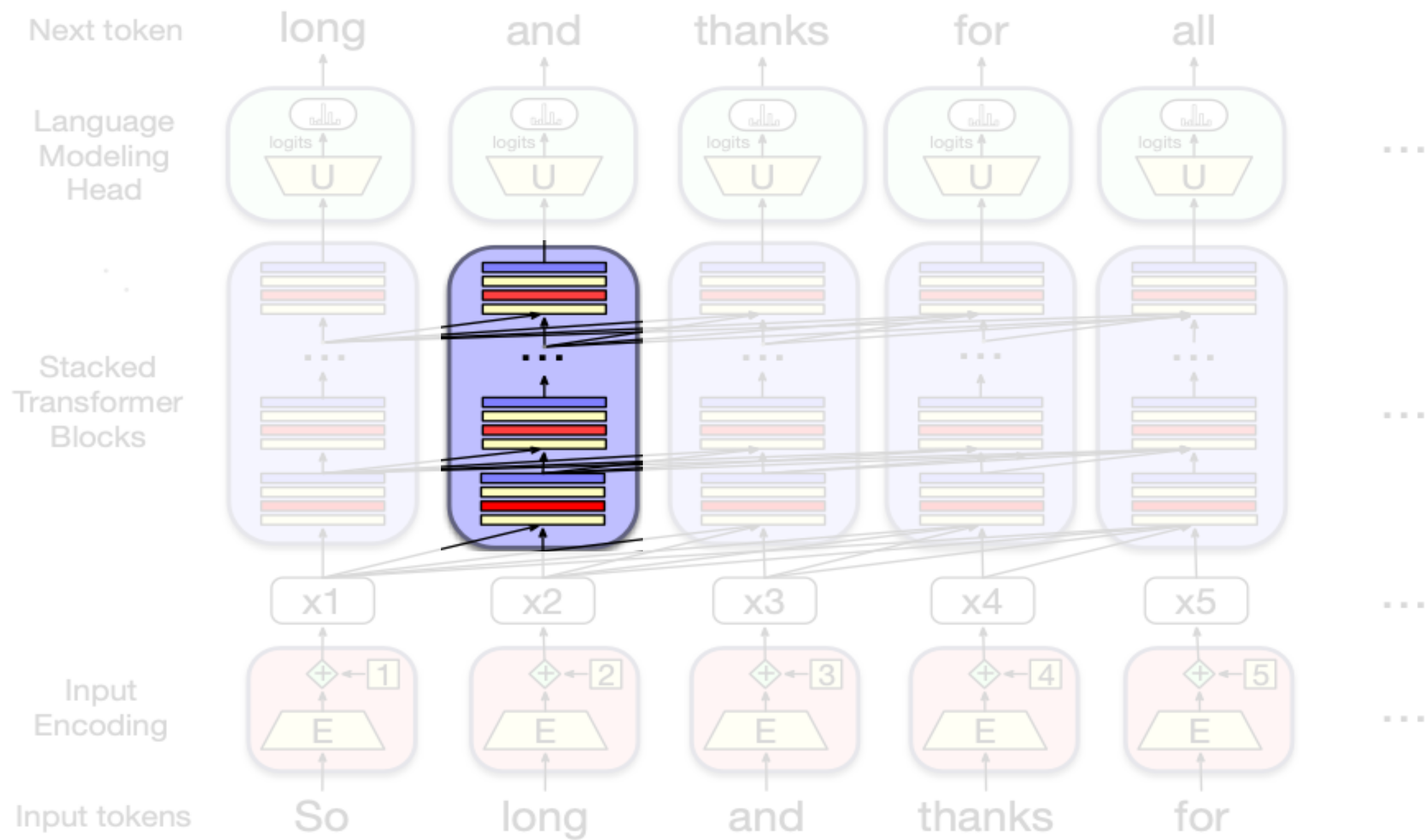The result: the embedding for each word will be different in different contexts!

Contextual embeddings: a representation of word meaning in its context.

We'll see next that attention can also be viewed as a way to move information from one token to another.

# The transformer block

$\mathbf{h}_{i-1}$ $\mathbf{h}_i$ $\mathbf{h}_{i+1}$

Feedforward

Multi-head attention

$\mathbf{x}_{i-1}$ $\mathbf{x}_i$ $\mathbf{x}_{i+1}$

*We'll need nonlinearities, so add a feedforward layer*

$$\mathrm{FFN}(\mathbf{x}_i) = \mathrm{ReLU}(\mathbf{x}_i \mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2$$

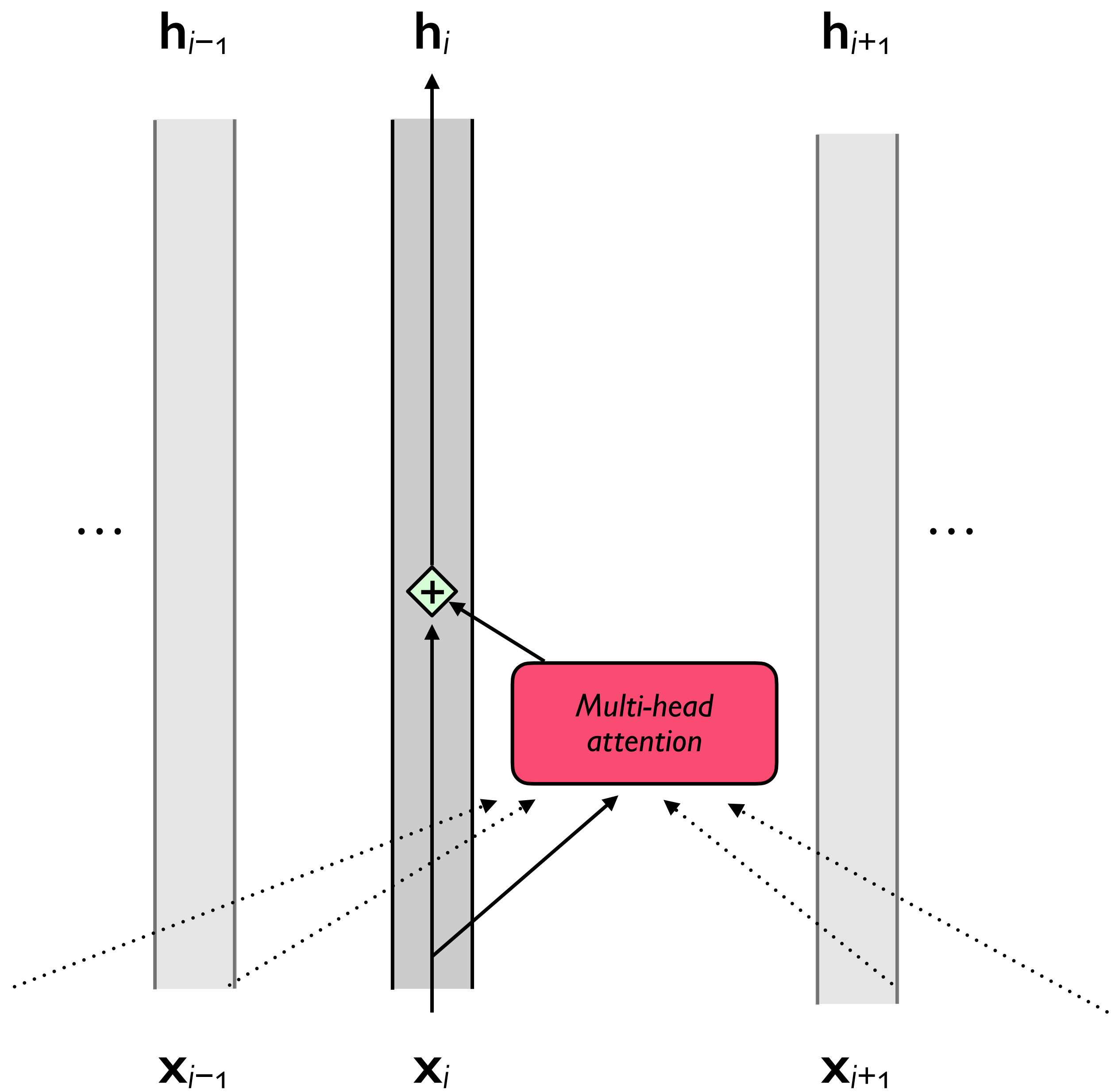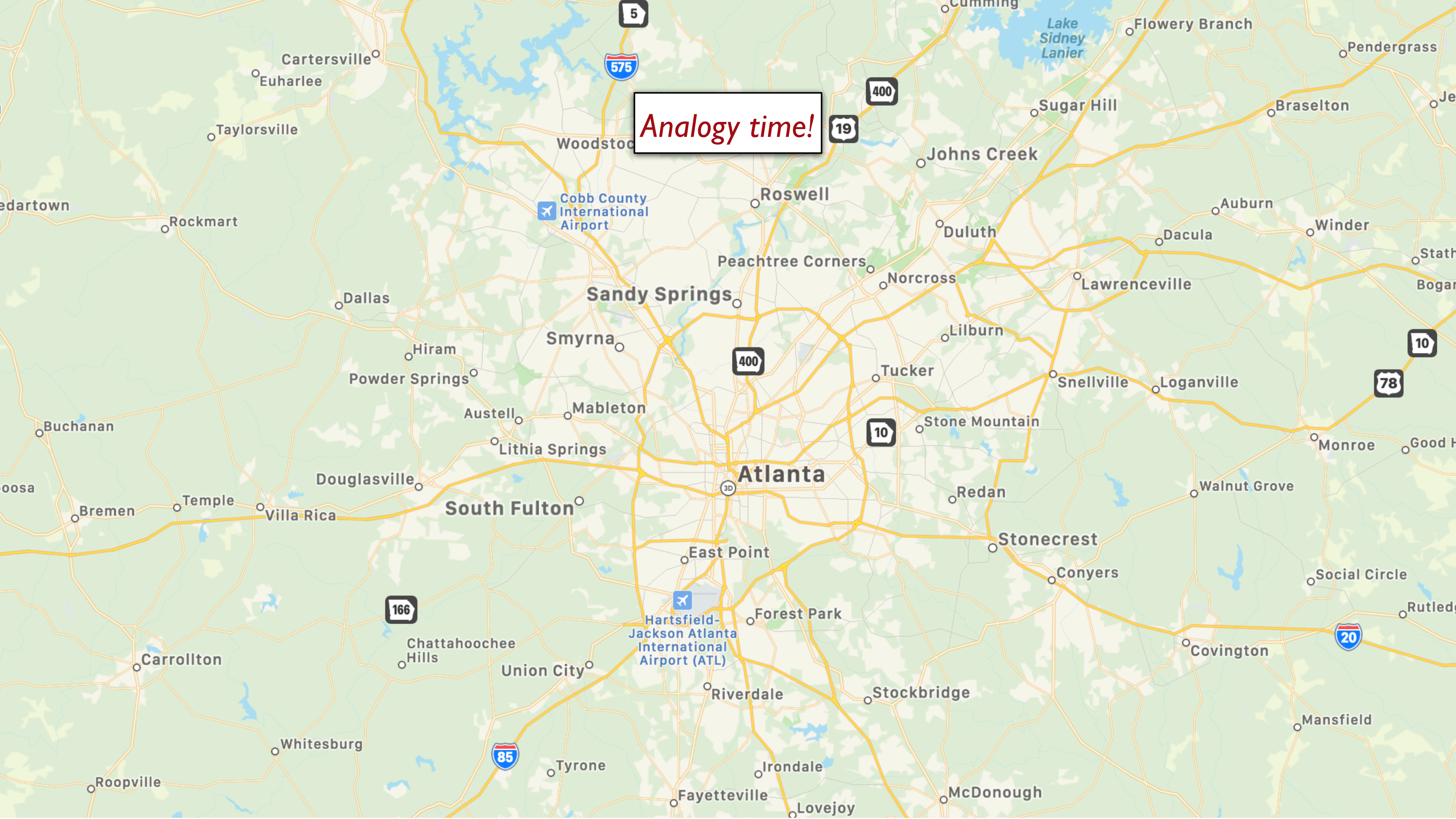$\mathbf{h}_{i-1}$        $\mathbf{h}_i$        $\mathbf{h}_{i+1}$

Feedforward

Layer norm

Multi-head attention

Layer norm

$\mathbf{x}_{i-1}$        $\mathbf{x}_i$        $\mathbf{x}_{i+1}$

Layer norm

The vector $\mathbf{x}_i$ is normalized twice

*Layer norm* is a variation of the z-score from statistics, applied to a single vector in a hidden layer:

mean

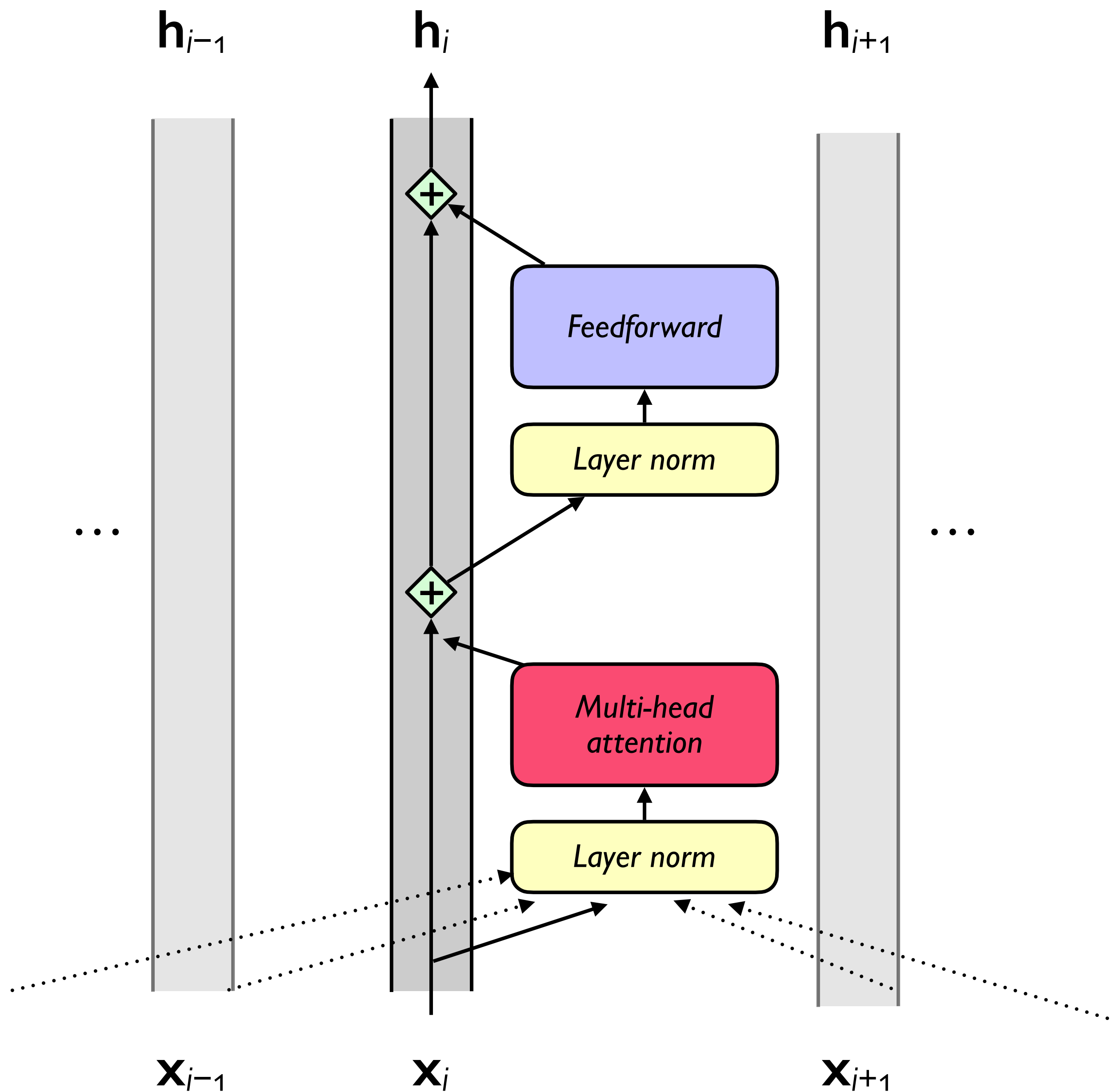$$\mu = \frac{1}{d}\sum_{i=1}^{d} x_i$$

standard deviation

$$\sigma = \sqrt{\frac{1}{d}\sum_{i=1}^{d}(x_i - \mu)^2}$$

normalized

$$\hat{\mathbf{x}} = \frac{(\mathbf{x} - \mu)}{\sigma}$$

normalized in a tunable way

$$\text{LayerNorm}(\mathbf{x}) = \gamma\frac{(\mathbf{x} - \mu)}{\sigma} + \beta$$

$$\mathbf{h}_i = \mathbf{t}_i^5 + \mathbf{t}_i^3$$
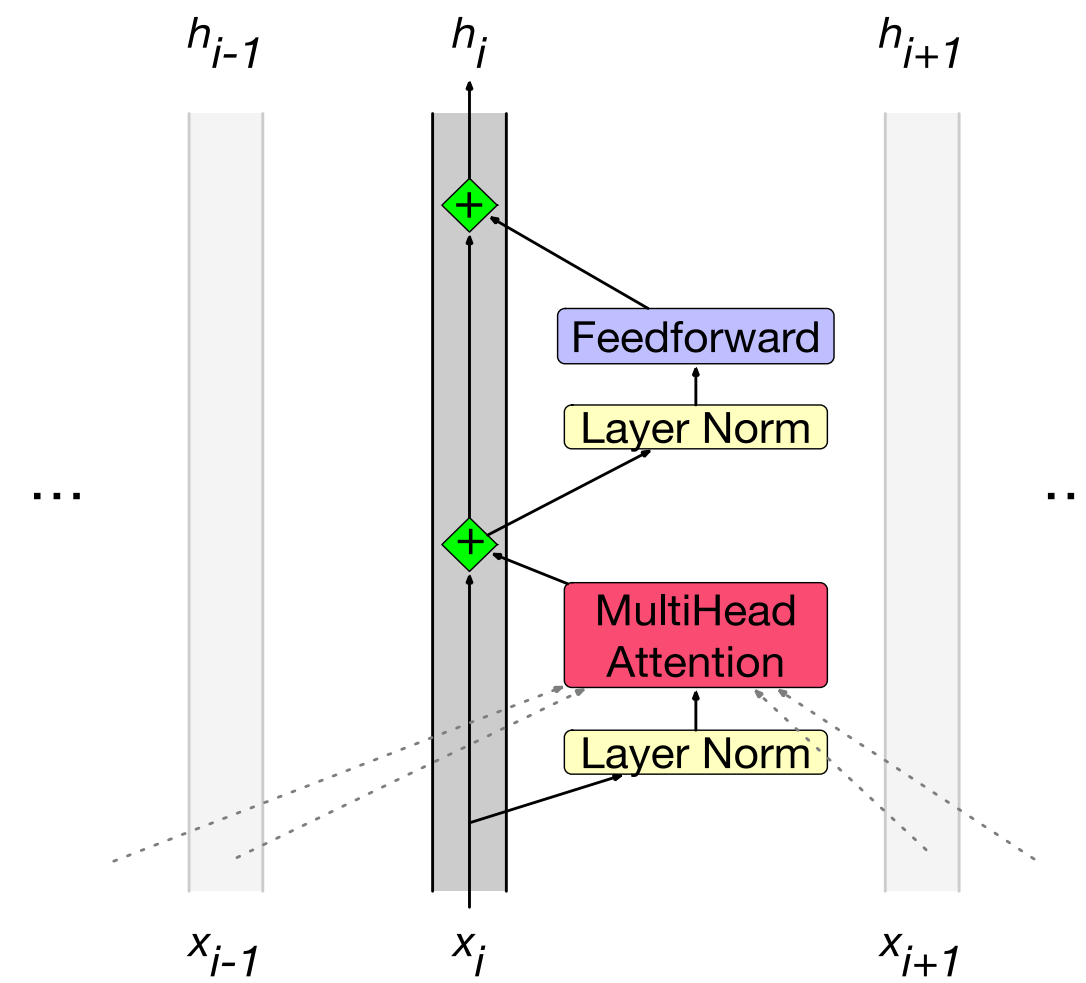
$$\mathbf{t}_i^5 = \text{FFN}(\mathbf{t}_i^4)$$

$$\mathbf{t}_i^4 = \text{LayerNorm}(\mathbf{t}_i^3)$$

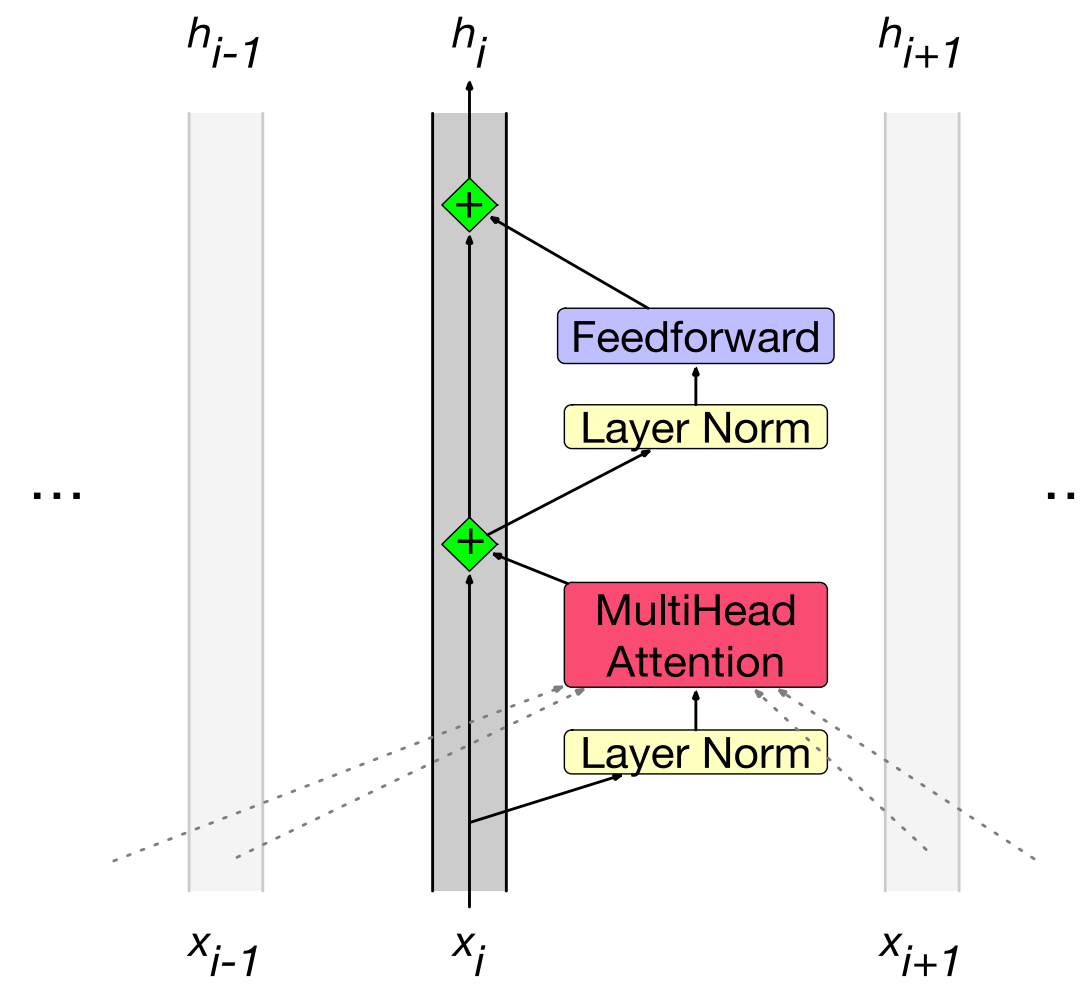$$\mathbf{t}_i^3 = \mathbf{t}_i^2 + \mathbf{x}_i$$

$$\mathbf{t}_i^2 = \text{MultiHeadAttention}(\mathbf{t}_i^1, [\mathbf{t}_1^1, \ldots, \mathbf{t}_N^1])$$

$$\mathbf{t}_i^1 = \text{LayerNorm}(\mathbf{x}_i)$$

Block 2

Block 1

A transformer is a stack of these blocks – so all the vectors are of the same dimensionality d

Notice that all parts of the transformer block apply to 1 residual stream (1 token) – except attention, which takes information from other tokens.

Elhage et al. (2021) show that we can view attention heads as literally moving information from the residual stream of a neighboring token into the current stream:



Token A residual stream

Token B residual stream

# Acknowledgments

The lecture incorporates material from: