

Transformers (II)

CS 685, Fall 2025

Advanced Natural Language Processing

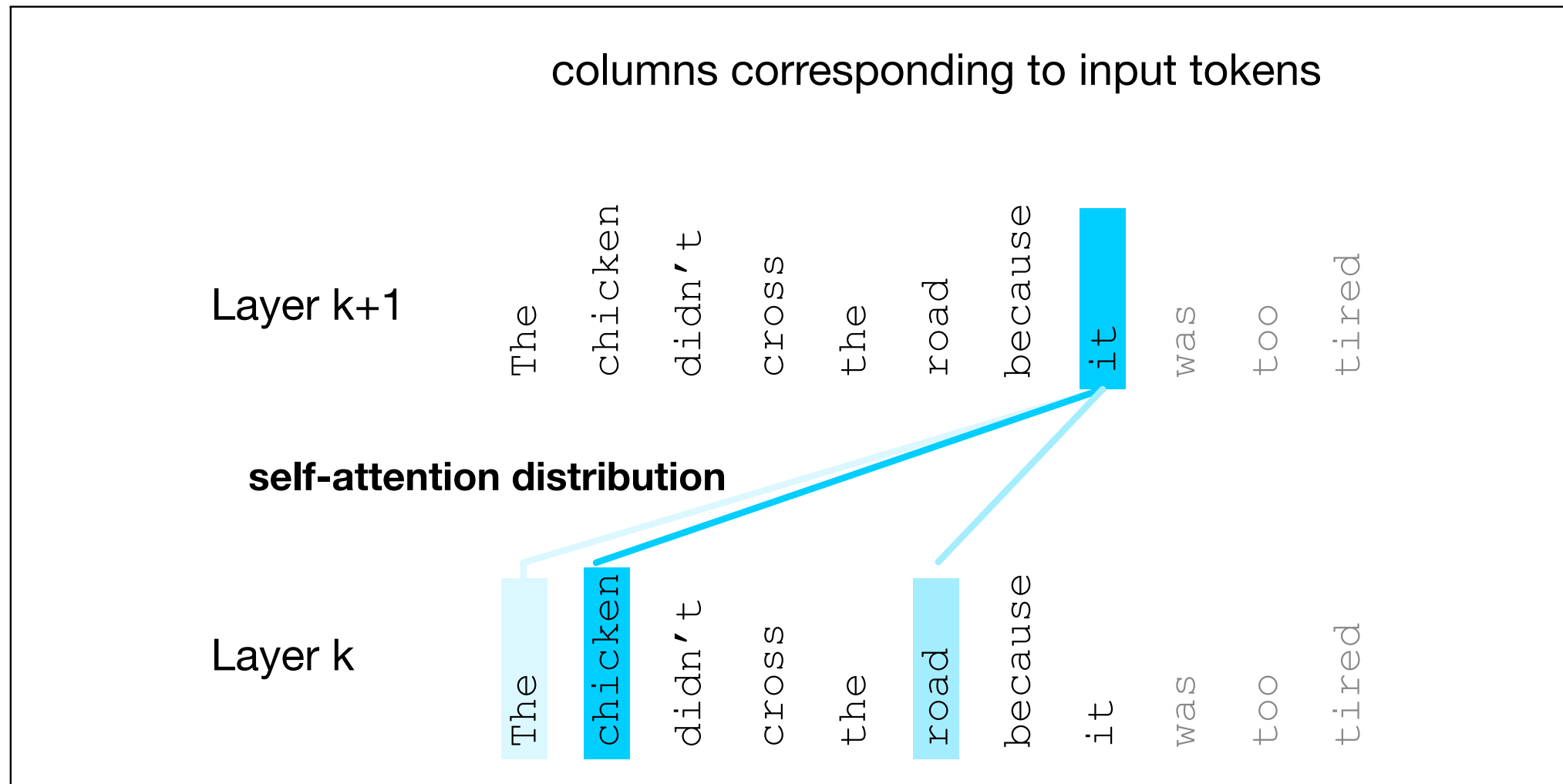
https://people.cs.umass.edu/~brenocon/cs685_f25/

Brendan O'Connor

College of Information and Computer Sciences

University of Massachusetts Amherst

- Practice midterm -- posted to Piazza



$$\mathbf{q}_i = \mathbf{x}_i \mathbf{W}^Q; \quad \mathbf{k}_j = \mathbf{x}_j \mathbf{W}^K; \quad \mathbf{v}_j = \mathbf{x}_j \mathbf{W}^V$$

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

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

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

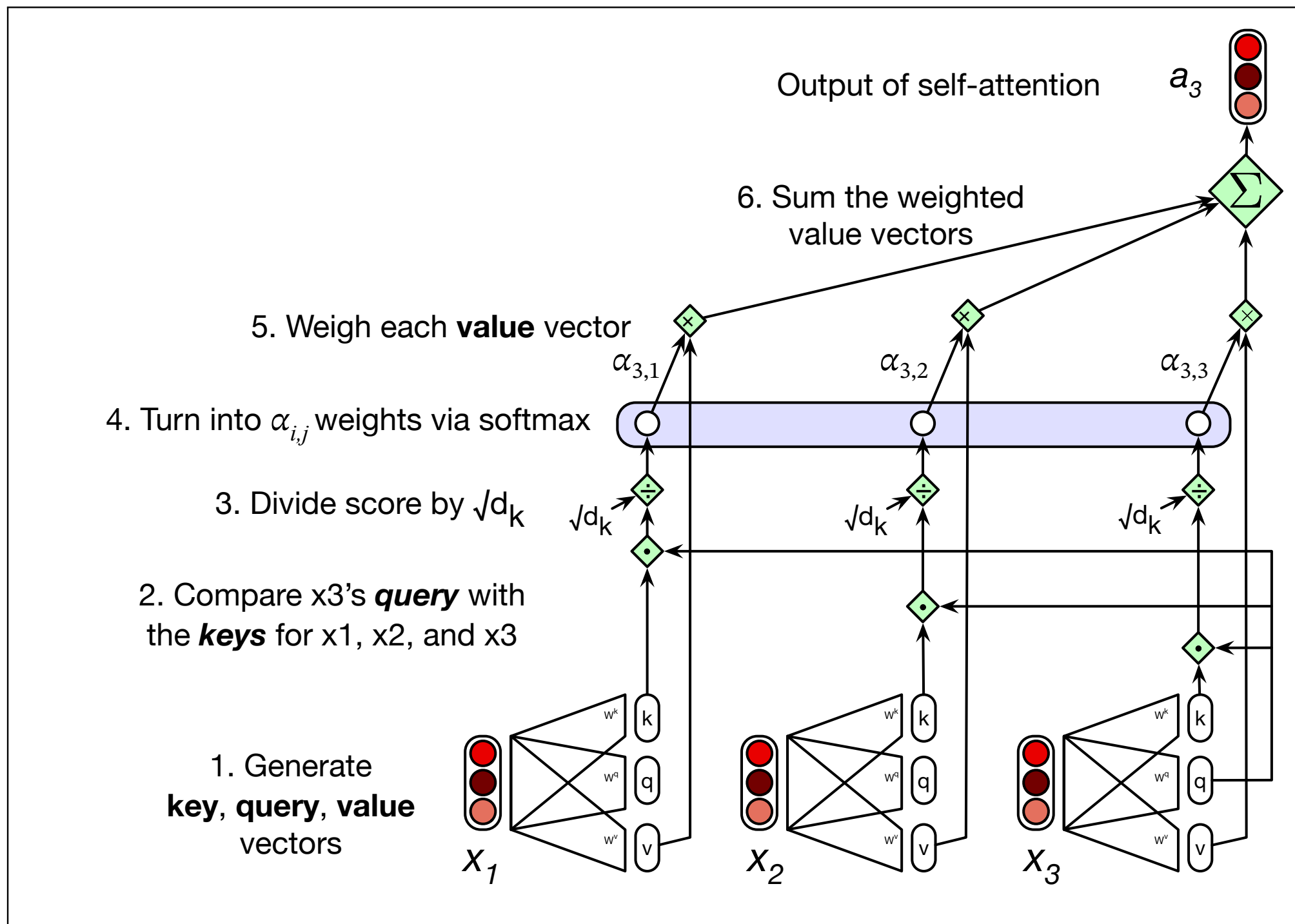


Figure 9.4 Calculating the value of a_3 , the third element of a sequence using causal (left-to-right) self-attention.

- Mathematical comparisons
 - Associate array, but soft
 - Kernel functions, but learned
 - Sequence automaton (choosing, copying, etc.)
- Parallelization of query-key products
 - Comp. advantage vs. RNNs
 - Need to mask out future information

Full Transformer block

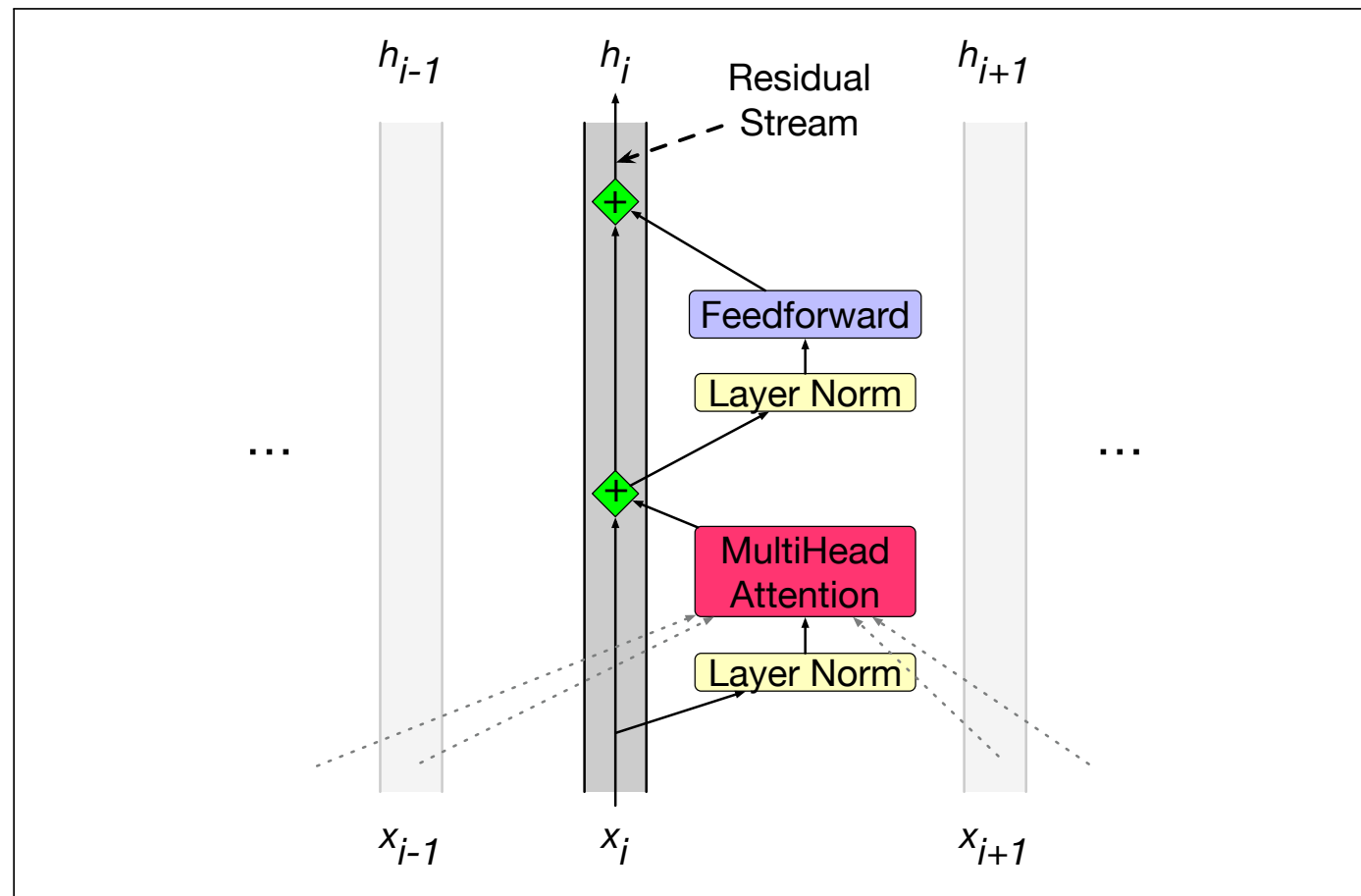


Figure 8.6 The architecture of a transformer block showing the **residual stream**. This figure shows the **prenorm** version of the architecture, in which the layer norms happen before the attention and feedforward layers rather than after.

$$\begin{aligned}
 \mathbf{t}_i^1 &= \text{LayerNorm}(\mathbf{x}_i) \\
 \mathbf{t}_i^2 &= \text{MultiHeadAttention}(\mathbf{t}_i^1, [\mathbf{t}_1^1, \dots, \mathbf{t}_N^1]) \\
 \mathbf{t}_i^3 &= \mathbf{t}_i^2 + \mathbf{x}_i \\
 \mathbf{t}_i^4 &= \text{LayerNorm}(\mathbf{t}_i^3) \\
 \mathbf{t}_i^5 &= \text{FFN}(\mathbf{t}_i^4) \\
 \mathbf{h}_i &= \mathbf{t}_i^5 + \mathbf{t}_i^3
 \end{aligned}$$

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

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

- Multiple layers with MLPs (FF layers)
- Multiple heads

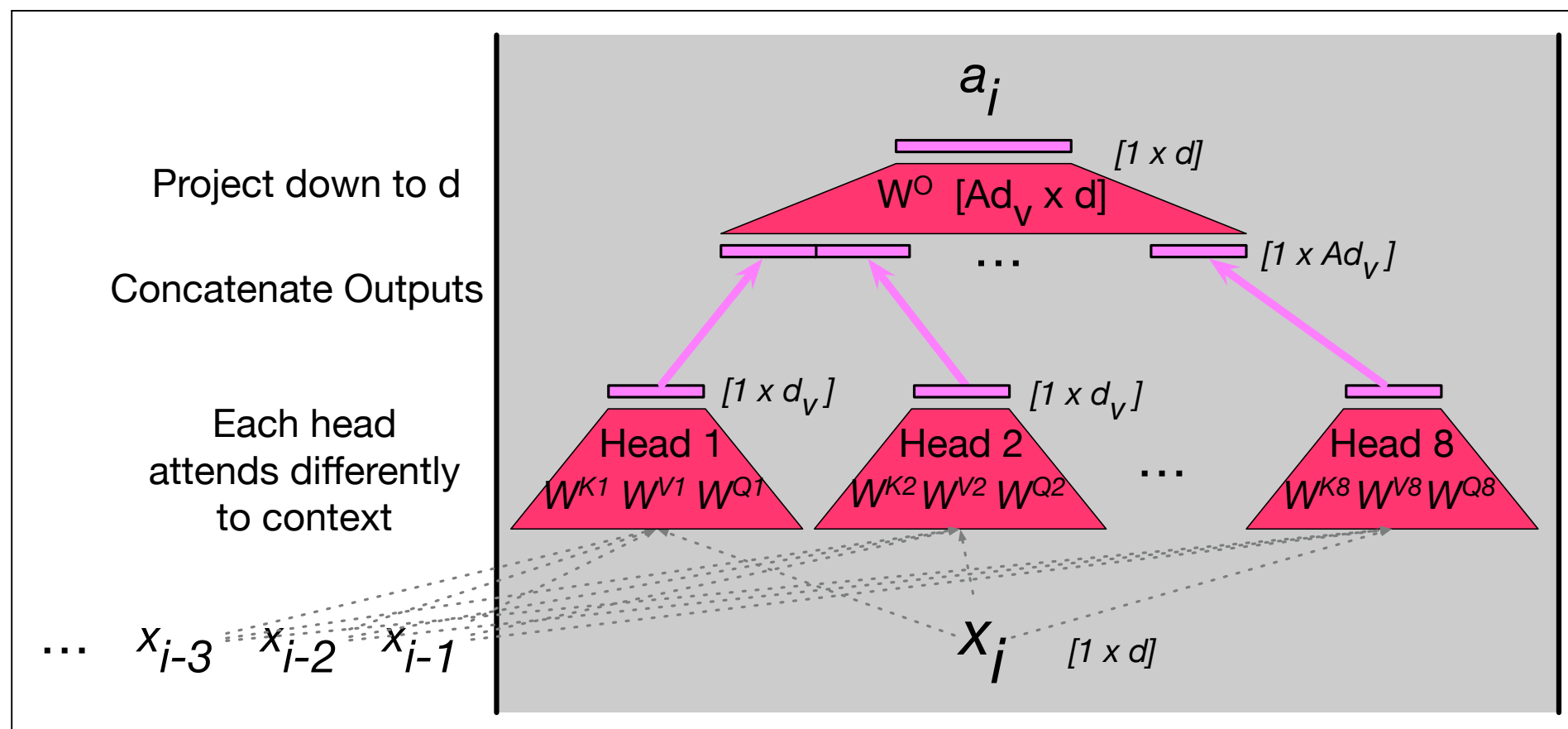


Figure 8.5 The multi-head attention computation for input \mathbf{x}_i , producing output \mathbf{a}_i . A multi-head attention layer has A heads, each with its own query, key, and value weight matrices. The outputs from each of the heads are concatenated and then projected down to d , thus producing an output of the same size as the input.

$$\mathbf{q}_i^c = \mathbf{x}_i \mathbf{W}^{Qc}; \quad \mathbf{k}_j^c = \mathbf{x}_j \mathbf{W}^{Kc}; \quad \mathbf{v}_j^c = \mathbf{x}_j \mathbf{W}^{Vc}; \quad \forall c \quad 1 \leq c \leq A$$

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

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

$$\text{head}_i^c = \sum_{j \leq i} \alpha_{ij}^c \mathbf{v}_j^c$$

$$\mathbf{a}_i = (\text{head}^1 \oplus \text{head}^2 \dots \oplus \text{head}^A) \mathbf{W}^O$$

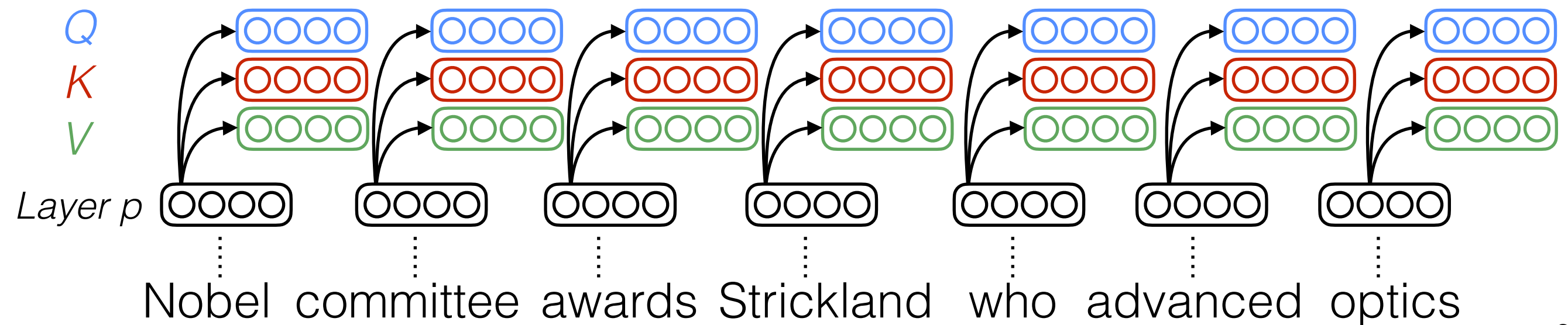
$$\text{MultiHeadAttention}(\mathbf{x}_i, [\mathbf{x}_1, \dots, \mathbf{x}_{i-1}]) = \mathbf{a}_i$$

- Multiple layers with MLPs (FF layers)
- Multiple heads

Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

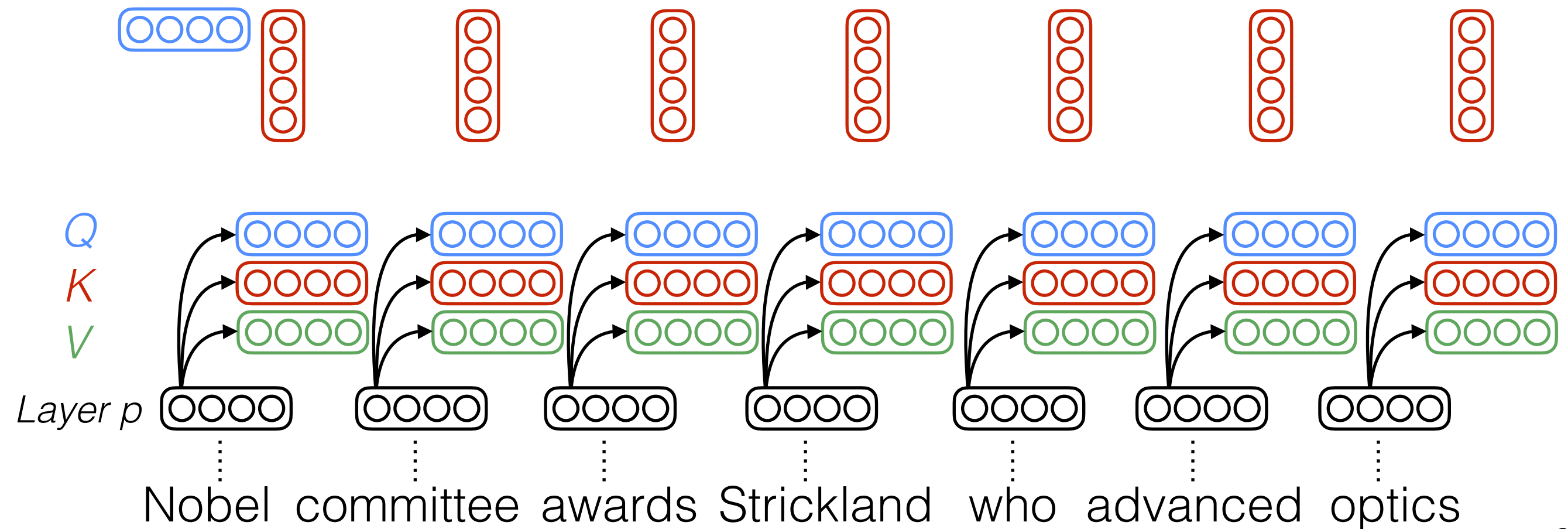
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

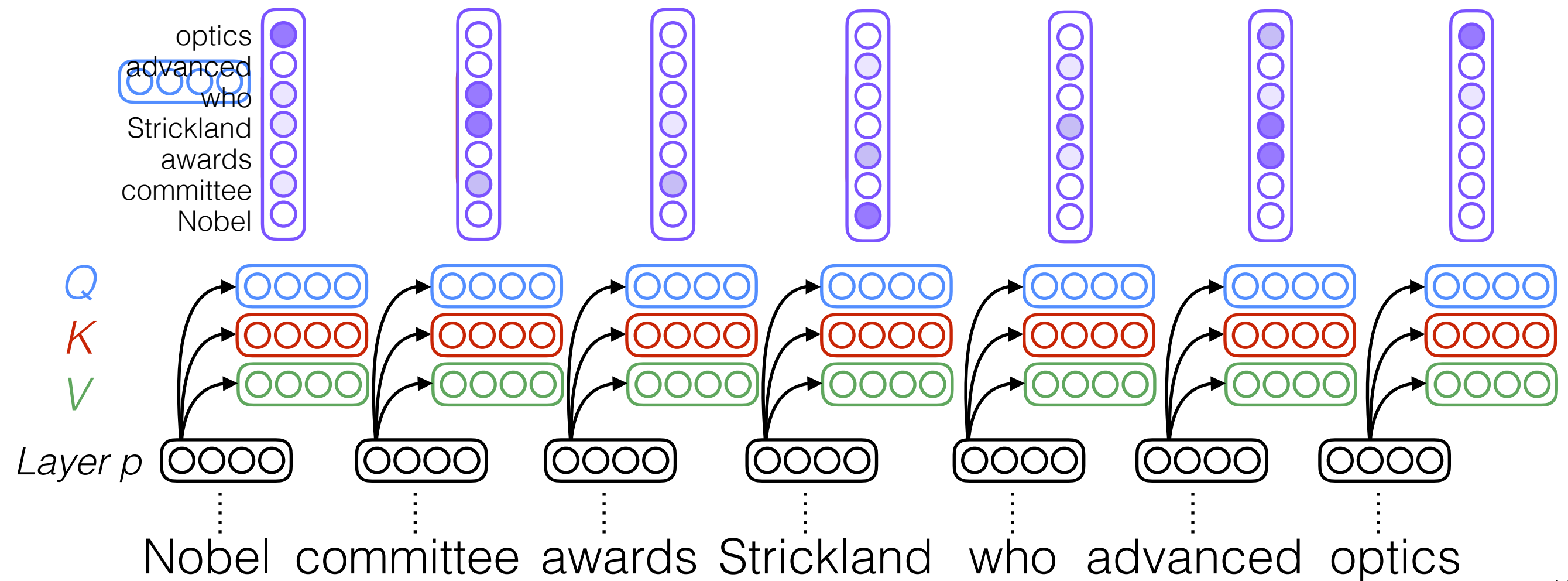
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

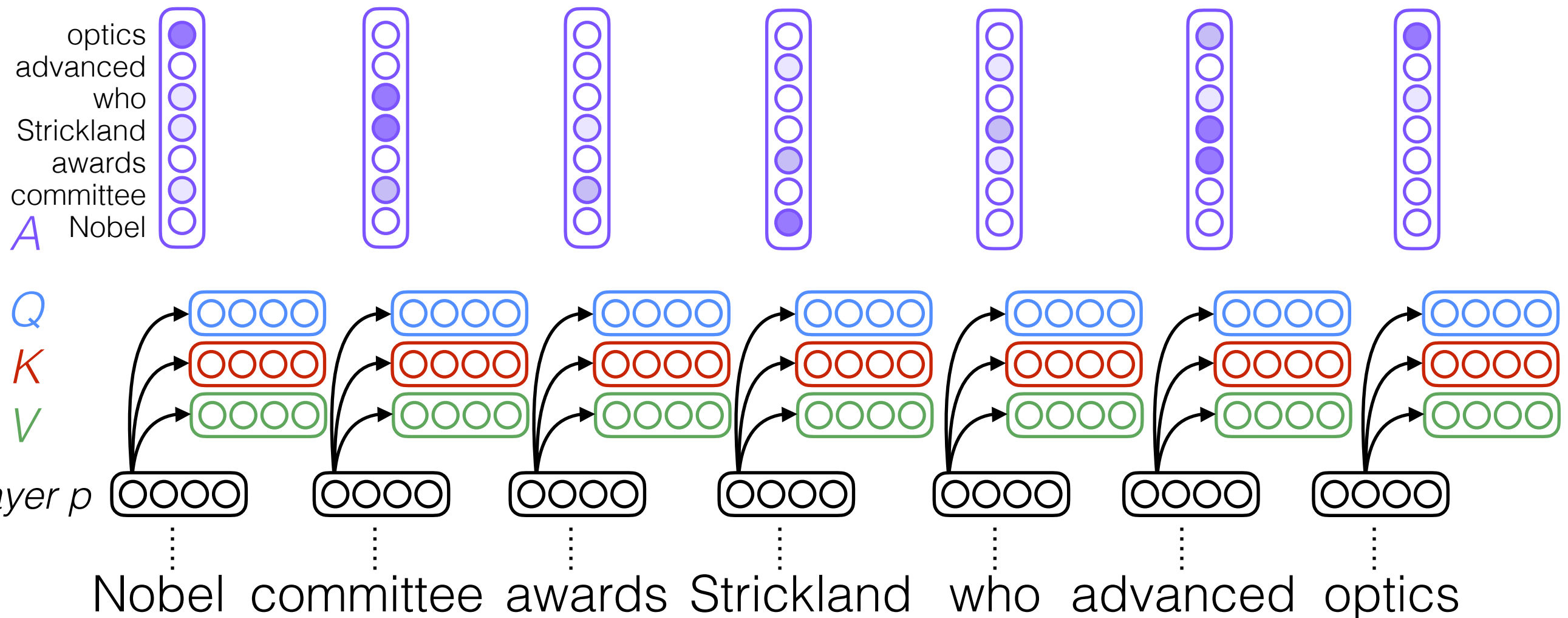
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

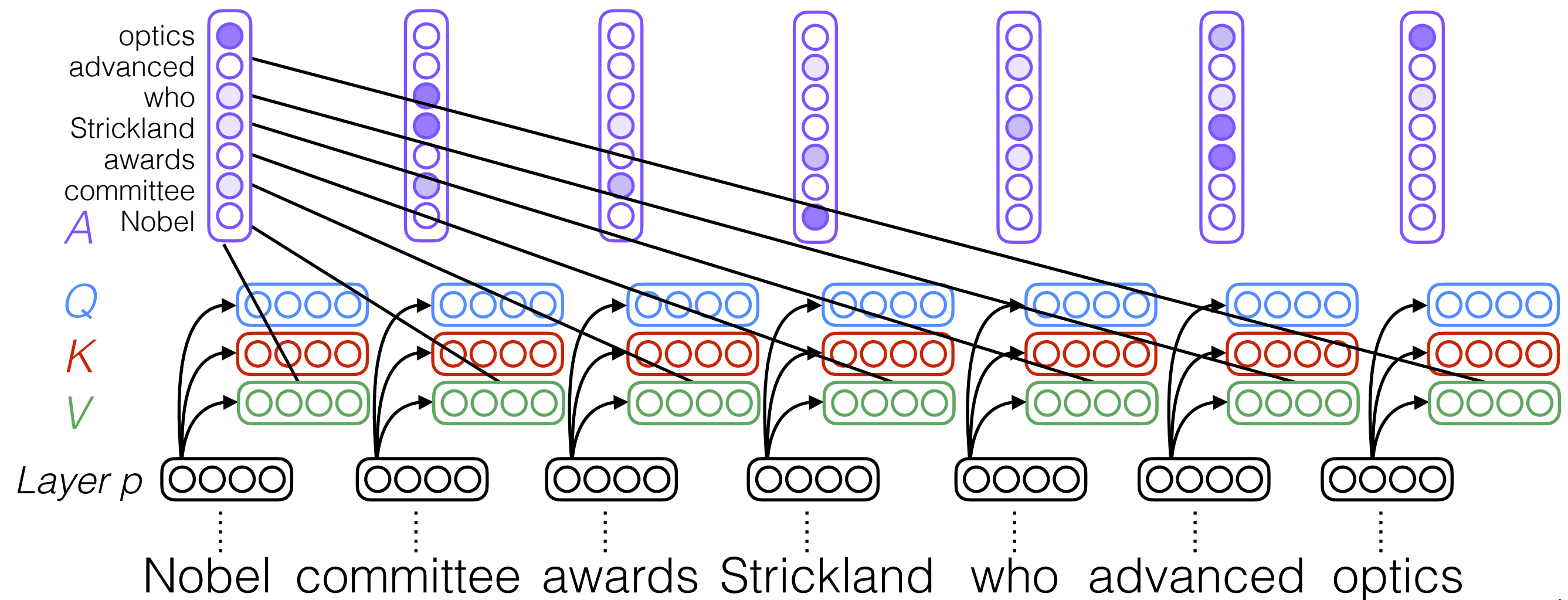
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

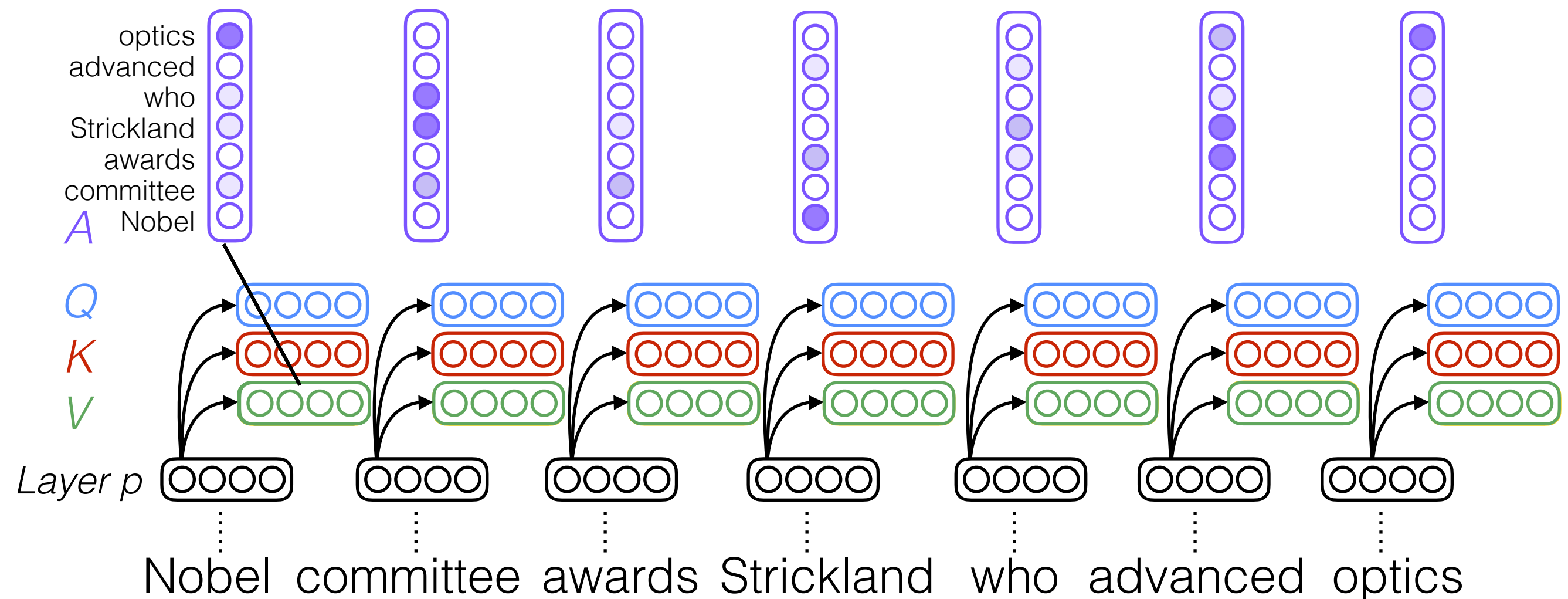
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

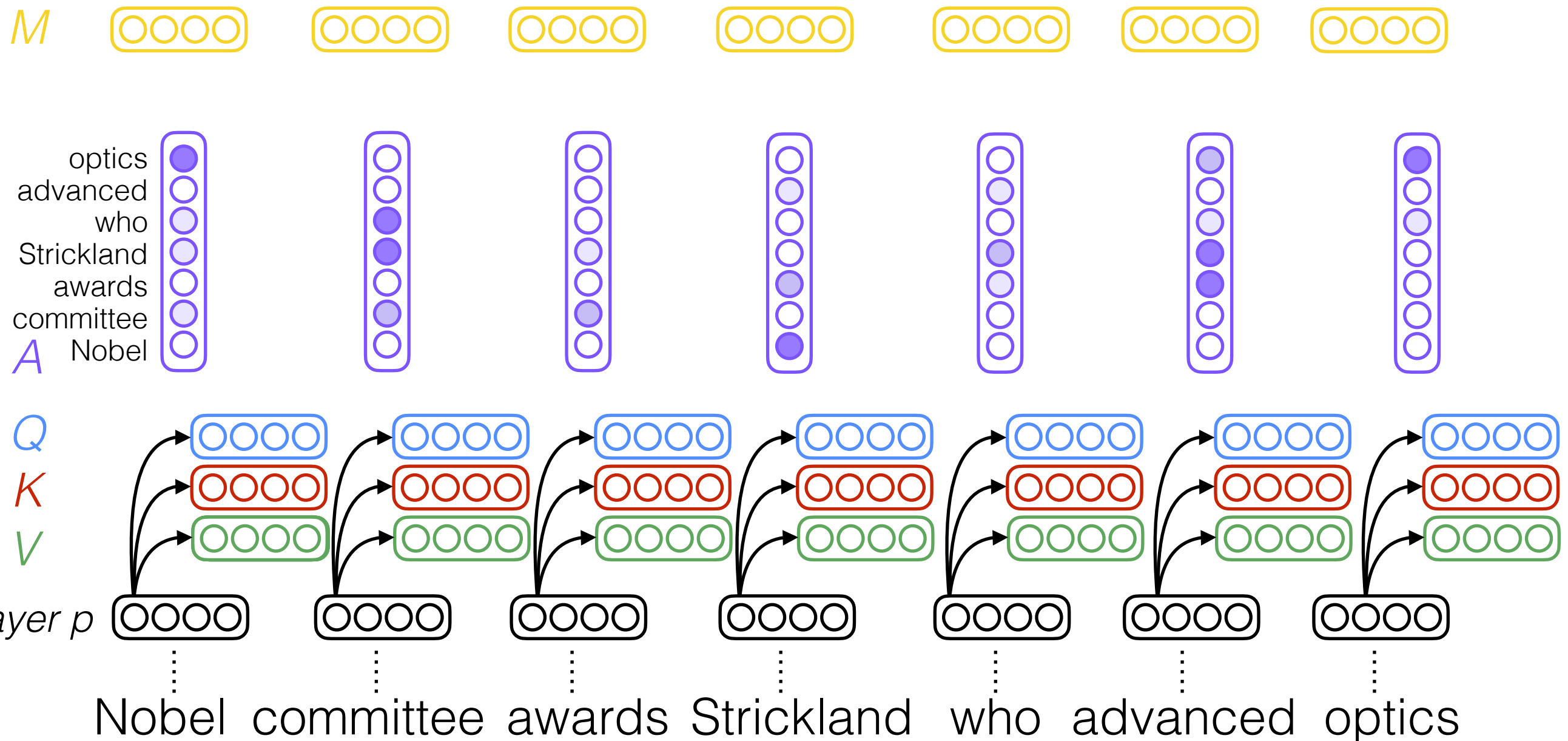
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]

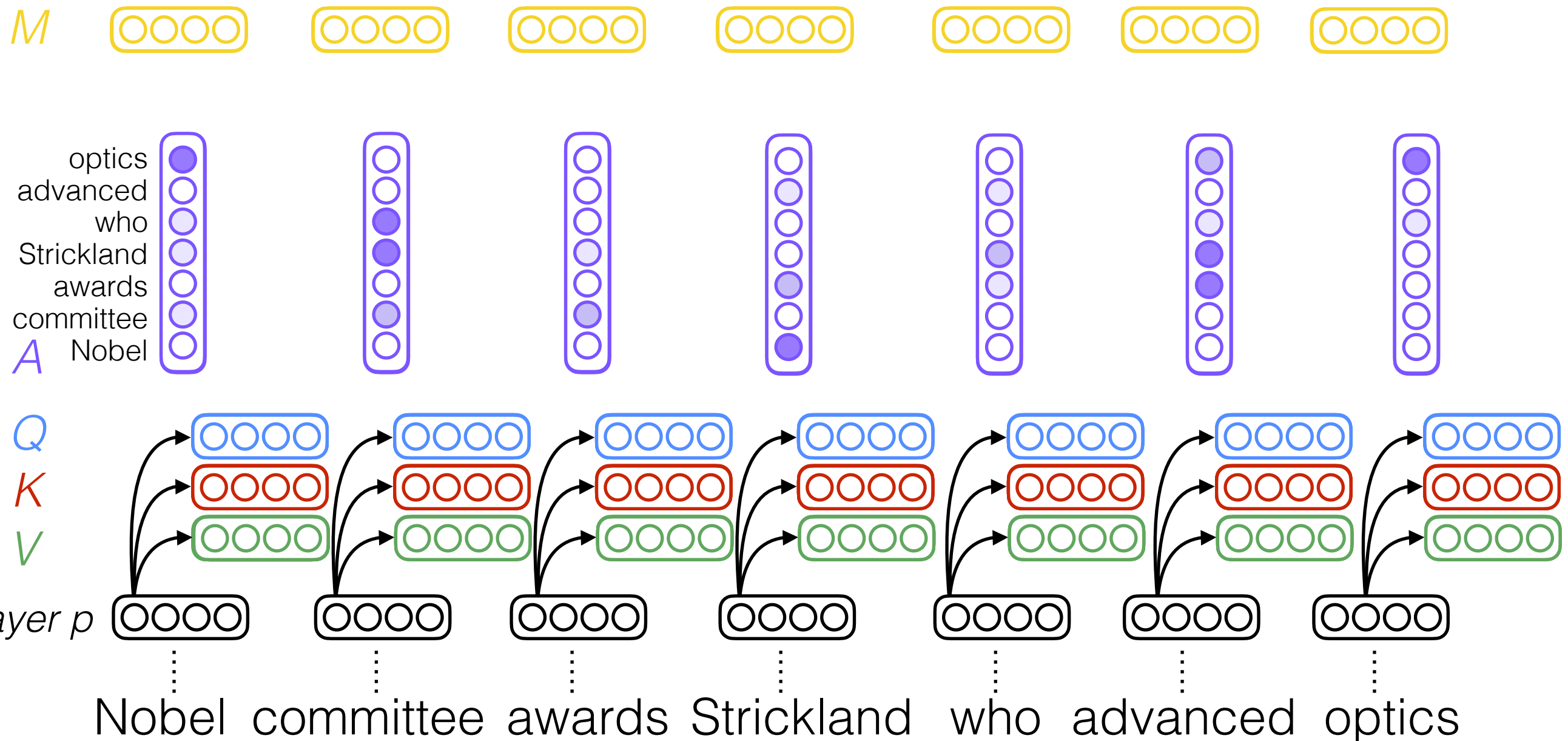
$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Self-attention

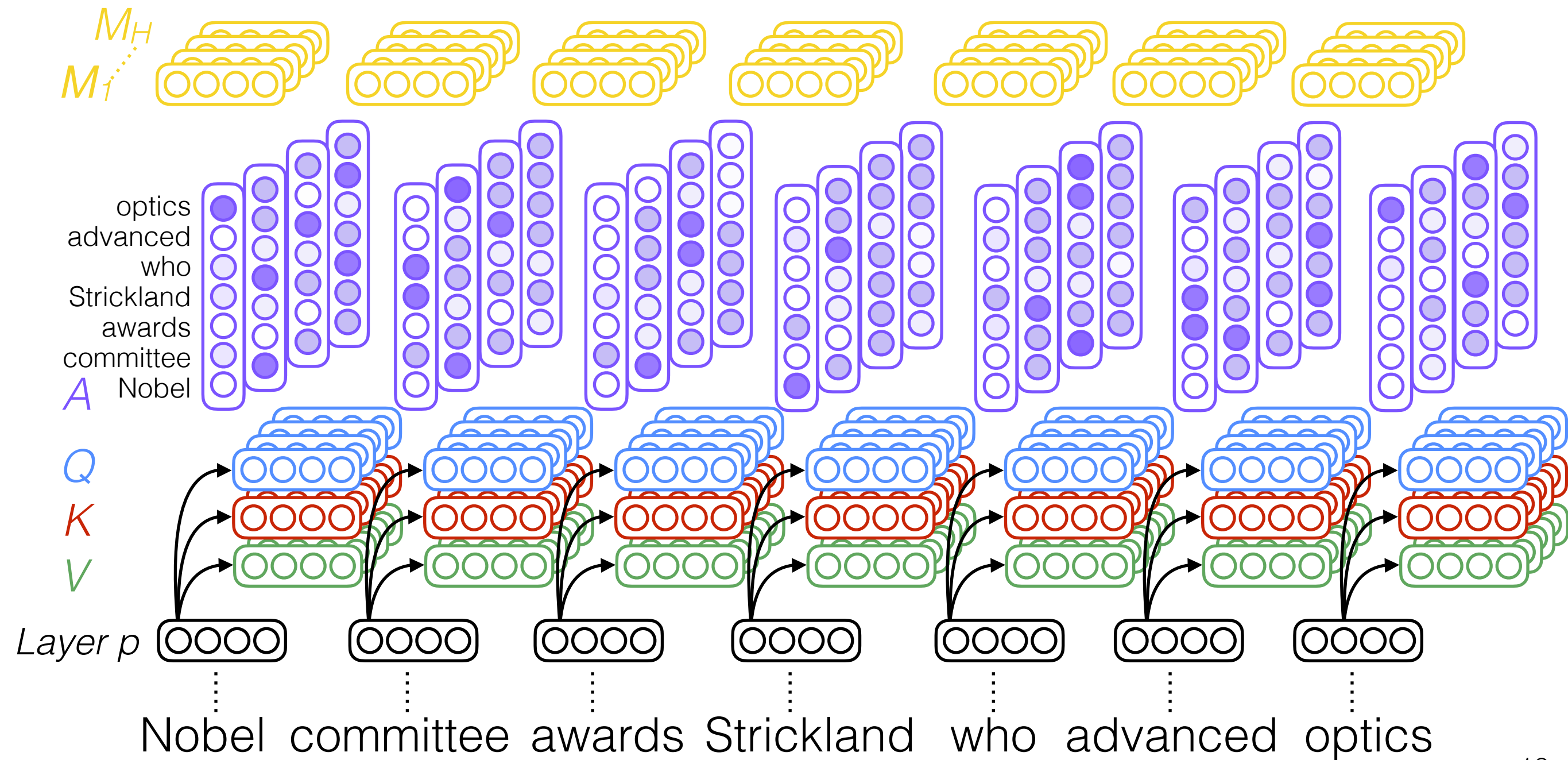
[Vaswani et al. 2017.
slide: Emma Strubell]

$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



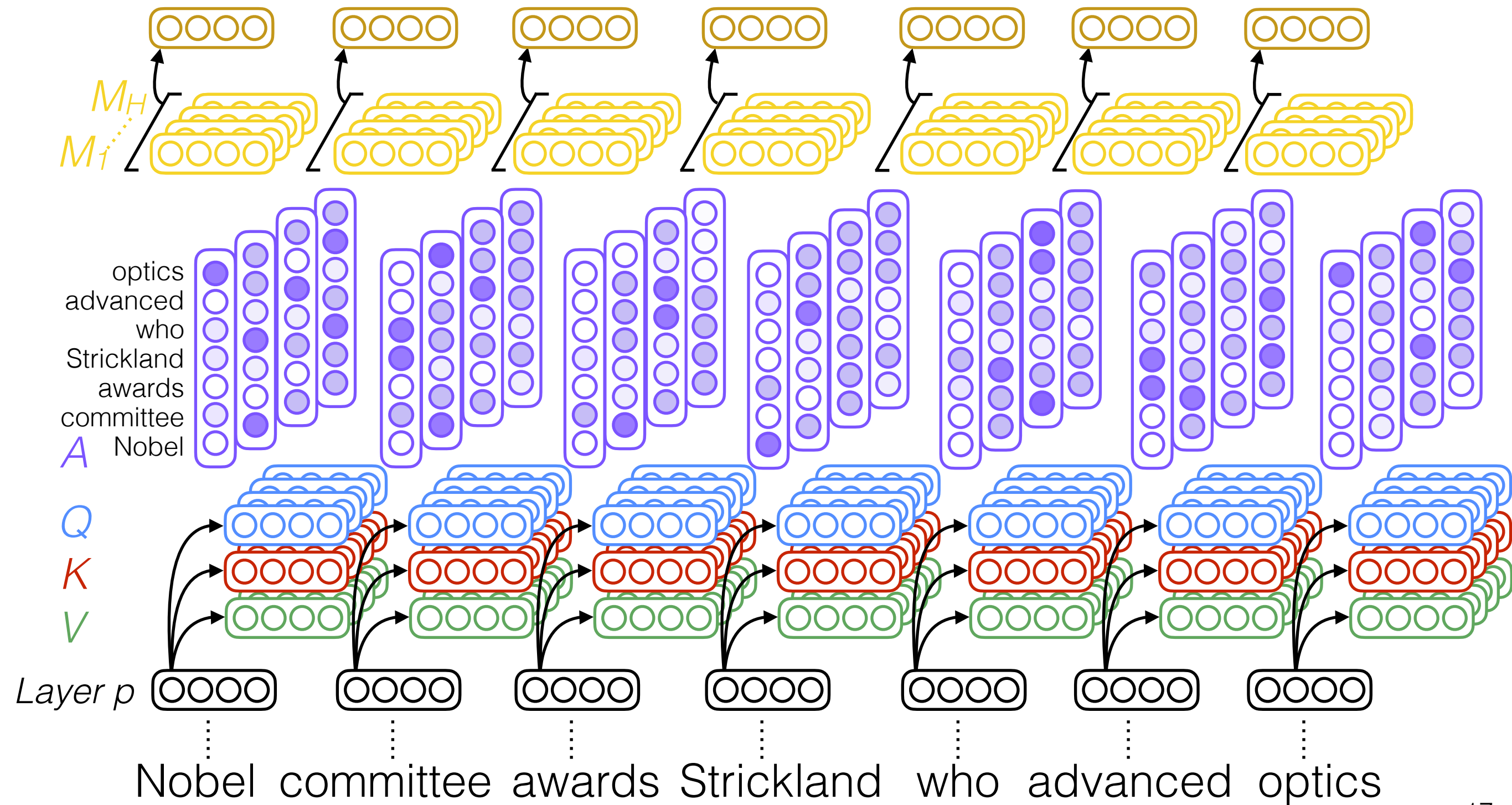
Multi-head self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]



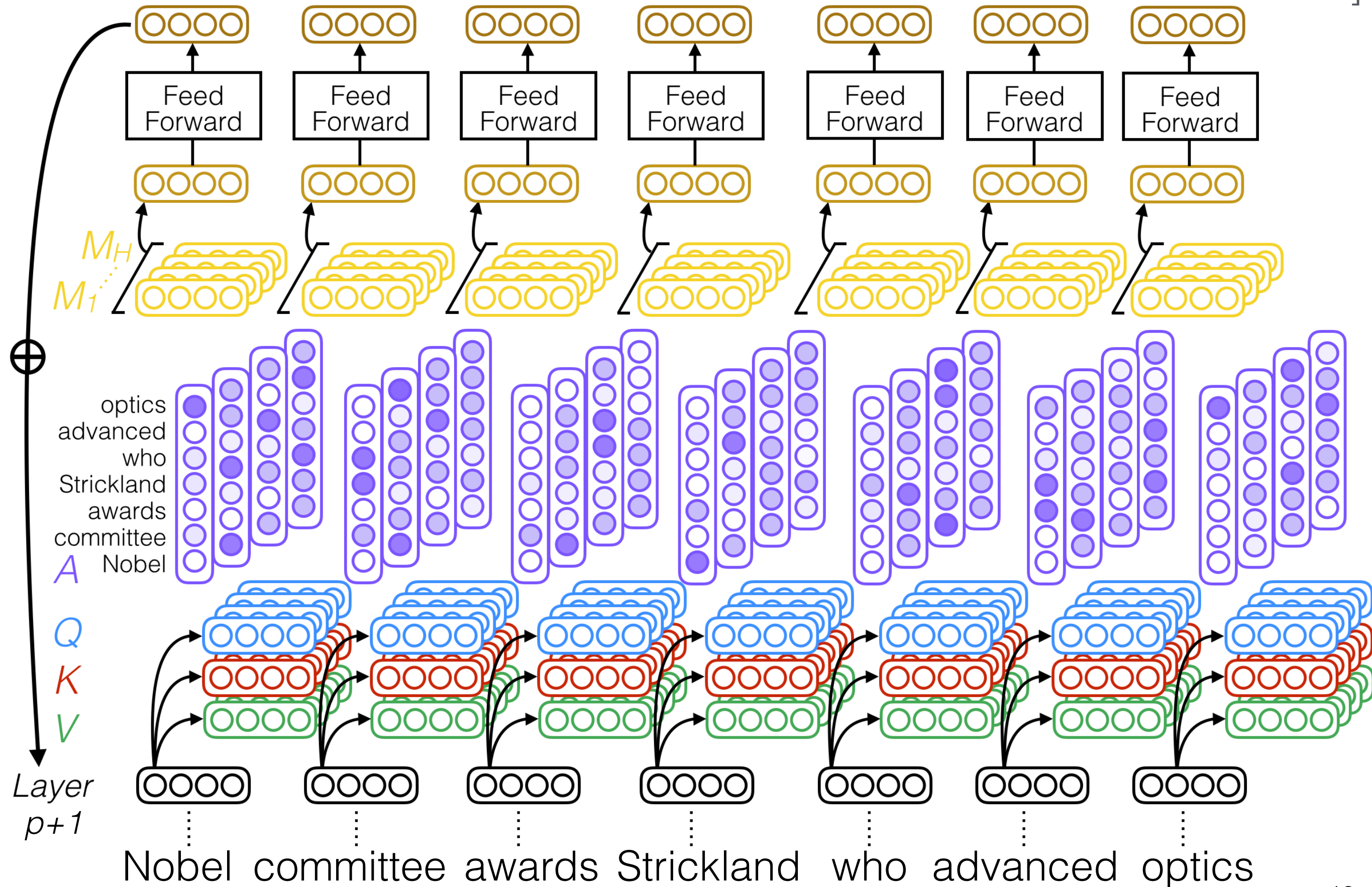
Multi-head self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]



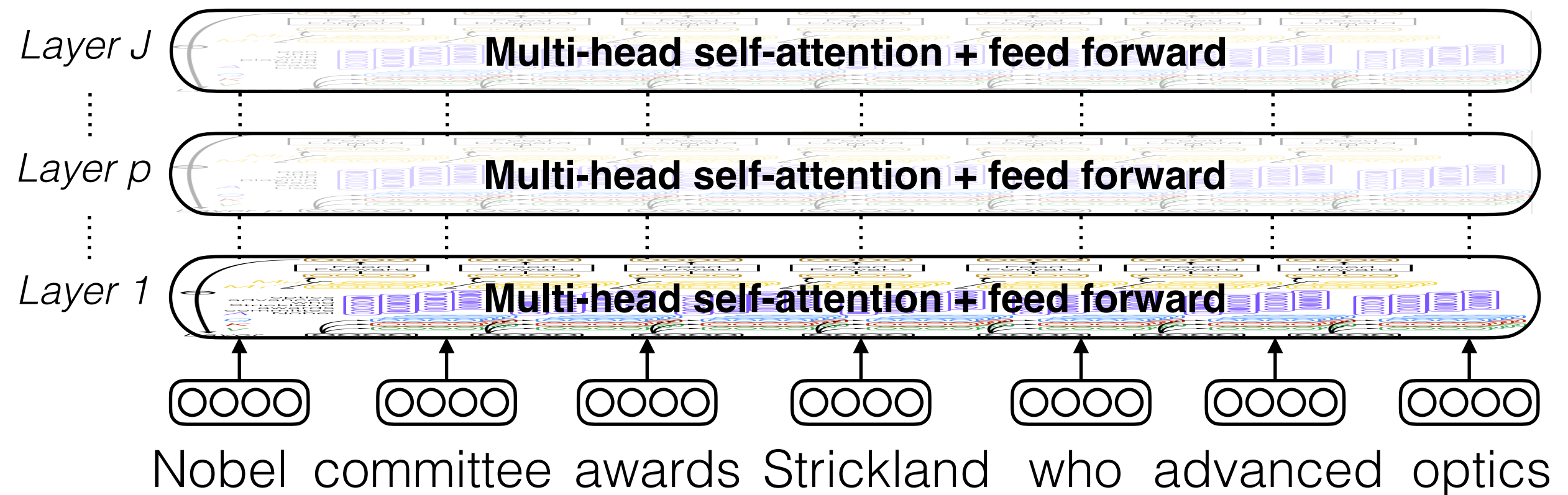
Multi-head self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]



Multi-head self-attention

[Vaswani et al. 2017.
slide: Emma Strubell]



- Parallelized training (switch slides)

Transformer
s

Parallelizing Attention Computation

Parallelizing computation using X

For attention/transformer block we've been computing a **single** output at a **single** time step i in a **single** residual stream.

But we can pack the N tokens of the input sequence into a single matrix X of size $[N \times d]$.

Each row of X is the embedding of one token of the input.

X can have 1K-32K rows, each of the dimensionality of the embedding d (the **model dimension**)

$$Q = XW^Q; \quad K = XW^K; \quad V = XW^V$$

[slide: SLP3]

QK^T

Now can do a single matrix multiply to combine Q and K^T

N	$q1 \cdot k1$	$q1 \cdot k2$	$q1 \cdot k3$	$q1 \cdot k4$
	$q2 \cdot k1$	$q2 \cdot k2$	$q2 \cdot k3$	$q2 \cdot k4$
	$q3 \cdot k1$	$q3 \cdot k2$	$q3 \cdot k3$	$q3 \cdot k4$
	$q4 \cdot k1$	$q4 \cdot k2$	$q4 \cdot k3$	$q4 \cdot k4$
	N			

Parallelizing attention

- Scale the scores, take the softmax, and then multiply the result by V resulting in a matrix of shape $N \times d$
- An attention vector for each input token

$$\mathbf{A} = \text{softmax} \left(\text{mask} \left(\frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right) \right) \mathbf{V}$$

Masking out the future

$$\mathbf{A} = \text{softmax} \left(\text{mask} \left(\frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right) \right) \mathbf{V}$$

- What is this mask function?
 \mathbf{QK}^\top has a score for each query dot every key, *including those that follow the query.*
- Guessing the next word is pretty simple if you already know it!

Masking out the future

$$\mathbf{A} = \text{softmax} \left(\text{mask} \left(\frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right) \right) \mathbf{V}$$

Add $-\infty$ to cells in upper triangle

The softmax will turn it to 0

N

q1•k1	$-\infty$	$-\infty$	$-\infty$
q2•k1	q2•k2	$-\infty$	$-\infty$
q3•k1	q3•k2	q3•k3	$-\infty$
q4•k1	q4•k2	q4•k3	q4•k4

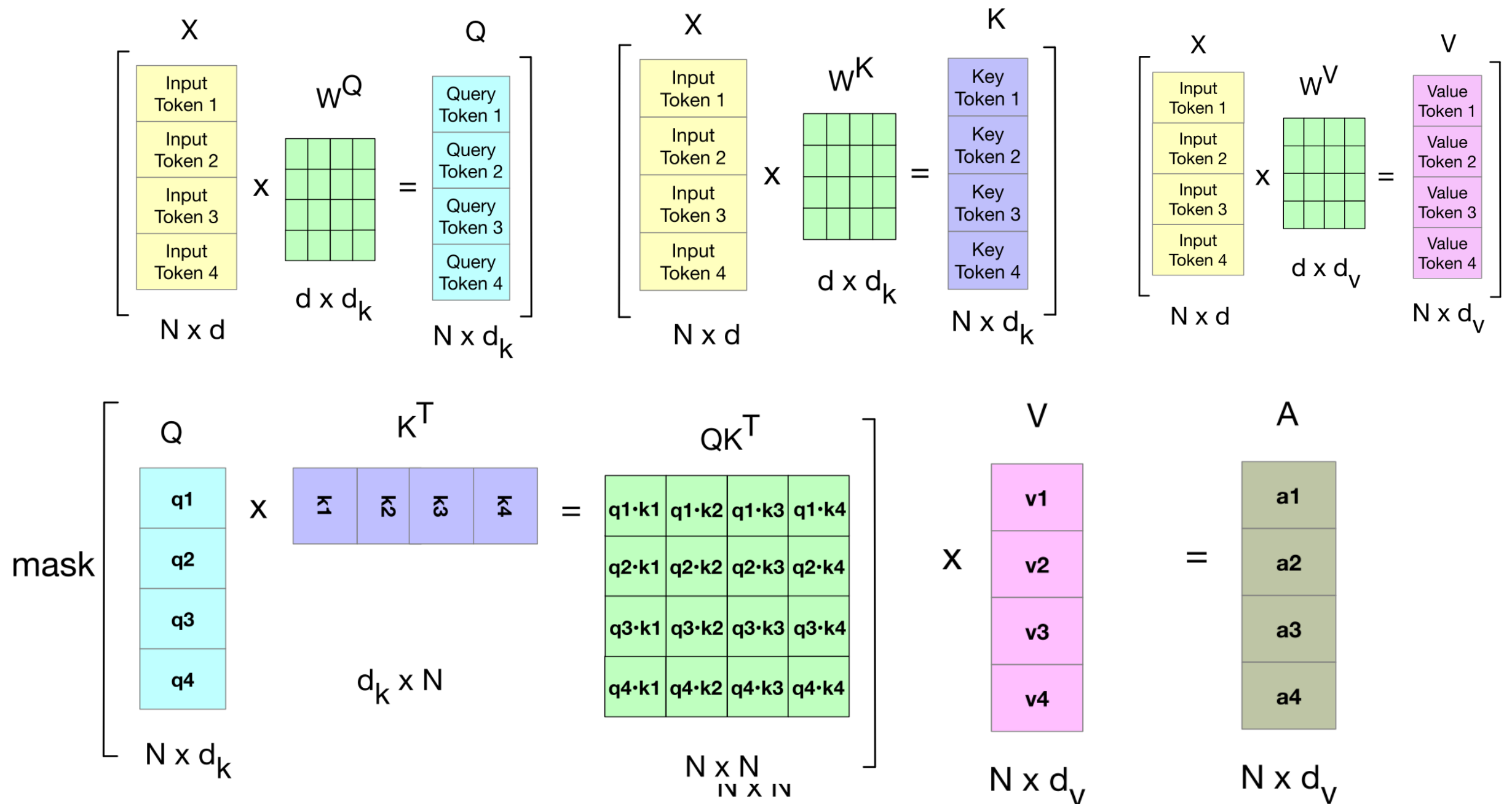
N

Another point: Attention is quadratic in length

$$\mathbf{A} = \text{softmax} \left(\text{mask} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \right) \mathbf{V}$$

N	q1•k1	−∞	−∞	−∞
	q2•k1	q2•k2	−∞	−∞
	q3•k1	q3•k2	q3•k3	−∞
	q4•k1	q4•k2	q4•k3	q4•k4
N				

Attention again



Parallelizing Multi-head Attention

$$\mathbf{Q}^i = \mathbf{XW}^{\mathbf{Q}^i}; \quad \mathbf{K}^i = \mathbf{XW}^{\mathbf{K}^i}; \quad \mathbf{V}^i = \mathbf{XW}^{\mathbf{V}^i}$$

$$\mathbf{head}_i = \text{SelfAttention}(\mathbf{Q}^i, \mathbf{K}^i, \mathbf{V}^i) = \text{softmax} \left(\frac{\mathbf{Q}^i \mathbf{K}^{i\top}}{\sqrt{d_k}} \right) \mathbf{V}^i$$

$$\text{MultiHeadAttention}(\mathbf{X}) = (\mathbf{head}_1 \oplus \mathbf{head}_2 \dots \oplus \mathbf{head}_h) \mathbf{W}^{\mathbf{O}}$$

Parallelizing Multi-head Attention

$$\mathbf{O} = \text{LayerNorm}(\mathbf{X} + \text{MultiHeadAttention}(\mathbf{X}))$$

$$\mathbf{H} = \text{LayerNorm}(\mathbf{O} + \text{FFN}(\mathbf{O}))$$

or

$$\mathbf{T}^1 = \text{MultiHeadAttention}(\mathbf{X})$$

$$\mathbf{T}^2 = \mathbf{X} + \mathbf{T}^1$$

$$\mathbf{T}^3 = \text{LayerNorm}(\mathbf{T}^2)$$

$$\mathbf{T}^4 = \text{FFN}(\mathbf{T}^3)$$

$$\mathbf{T}^5 = \mathbf{T}^4 + \mathbf{T}^3$$

$$\mathbf{H} = \text{LayerNorm}(\mathbf{T}^5)$$

Transformer
s

Parallelizing Attention Computation

Transformer s

Input and output:
Position embeddings
and the Language
Model Head

Token and Position Embeddings

The matrix X (of shape $[N \times d]$) has an embedding for each word in the context.

This embedding is created by adding two distinct embeddings for each input

- token embedding
- positional embedding

Token Embeddings

Embedding matrix E has shape $[|V| \times d]$.

- One row for each of the $|V|$ tokens in the vocabulary.
- Each word is a row vector of d dimensions

Given: string *"Thanks for all the"*

1. Tokenize with BPE and convert into vocab indices

$w = [5, 4000, 10532, 2224]$

2. Select the corresponding rows from E , each row an embedding
- (row 5, row 4000, row 10532, row 2224).

Position Embeddings

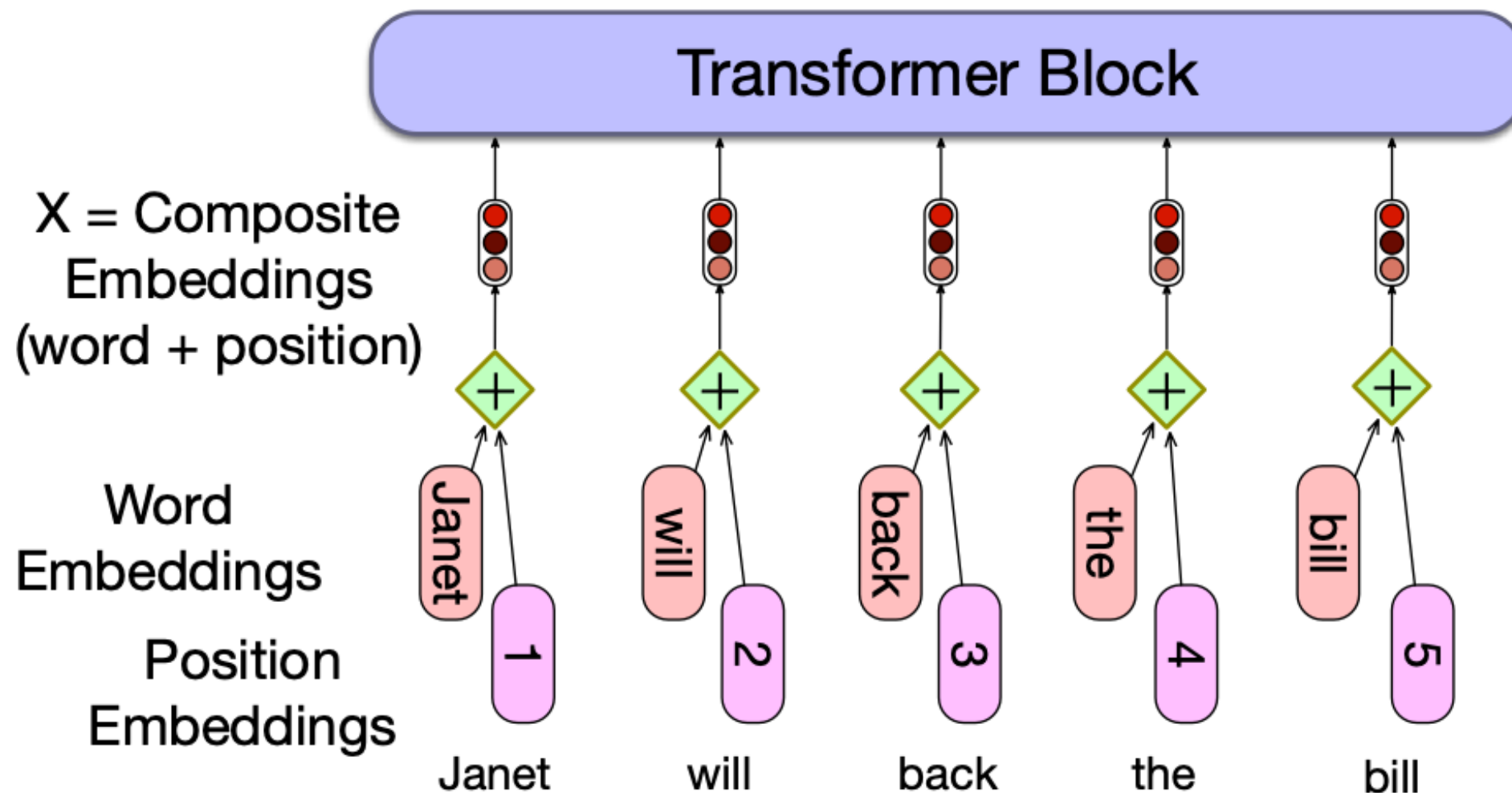
There are many methods, but we'll just describe the simplest: absolute position.

Goal: learn a position embedding matrix E_{pos} of shape $[1 \times N]$.

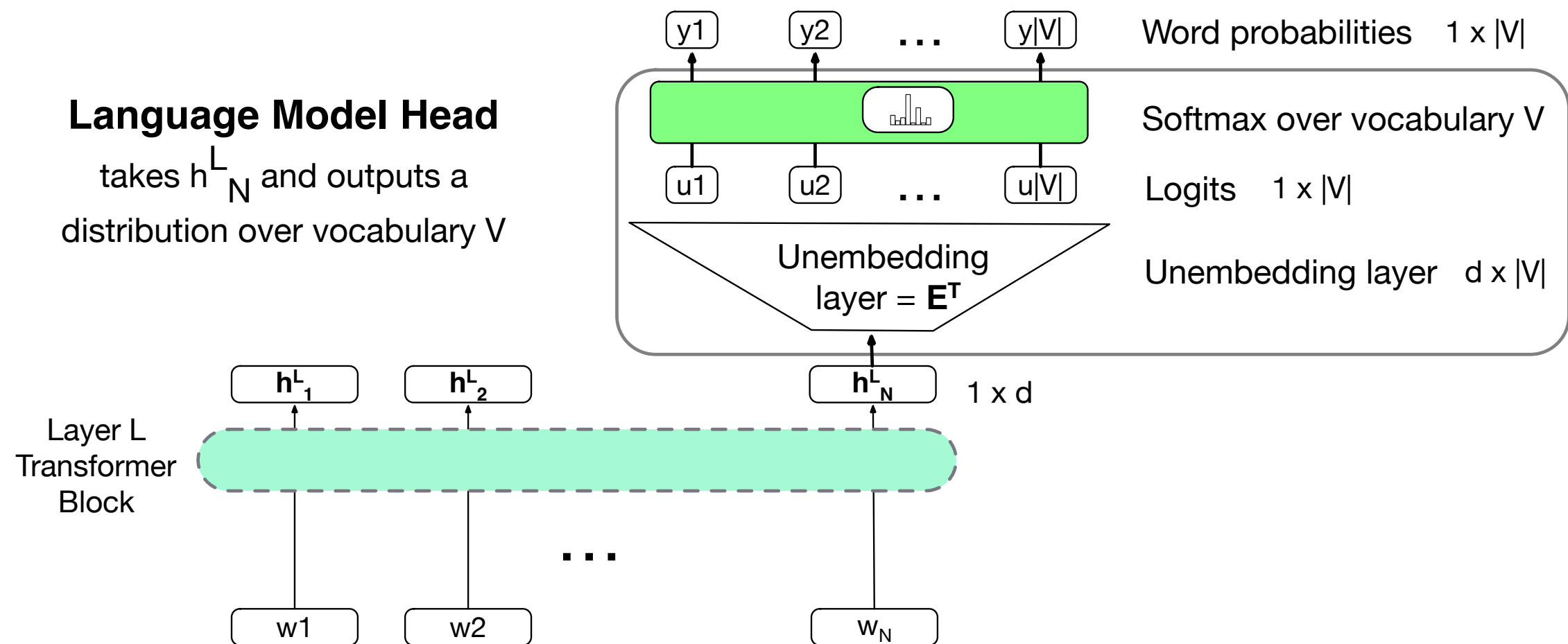
Start with randomly initialized embeddings

- one for each integer up to some maximum length.
- i.e., just as we have an embedding for token *fish*, we'll have an embedding for position 3 and position 17.
- As with word embeddings, these position embeddings are learned along with other parameters during training.

Each x is just the sum of word and position embeddings

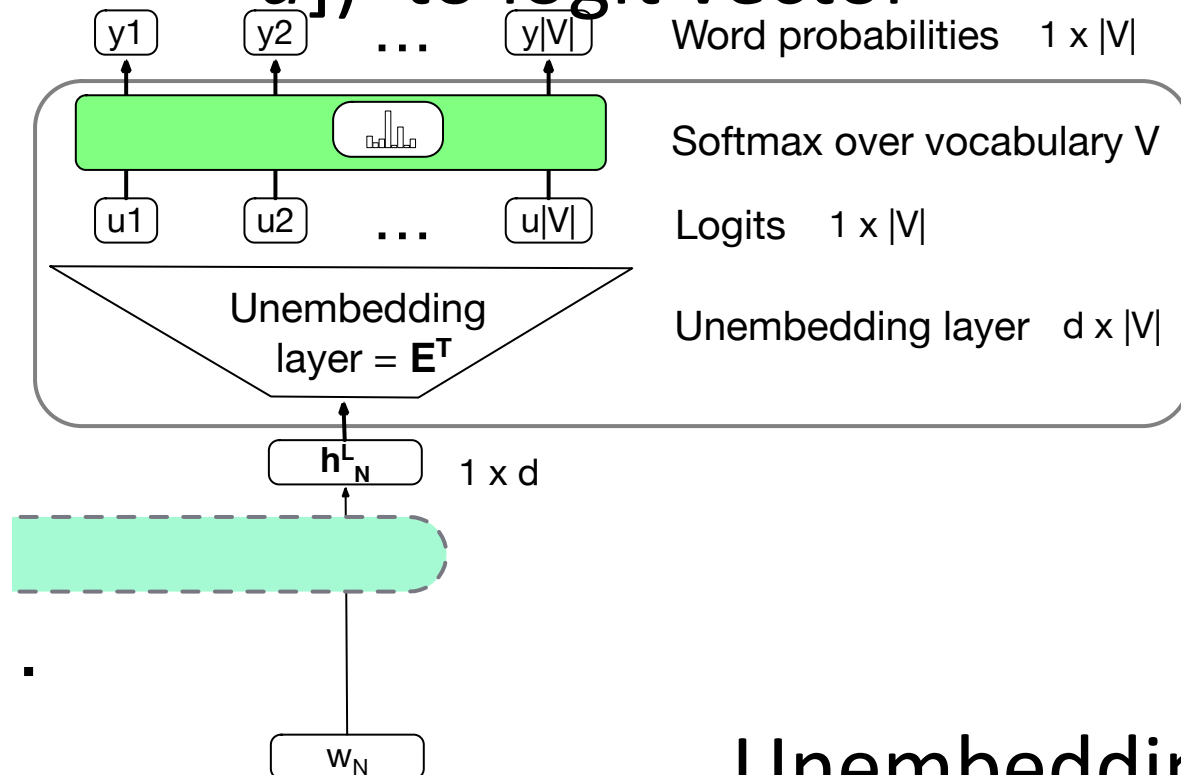


Language modeling head



Language modeling head

Unembedding layer: linear layer projects from h_N^L (shape $[1 \times d]$) to logit vector



Why "unembedding"?

Tied to E^T

Weight tying, we use the same weights for two different matrices

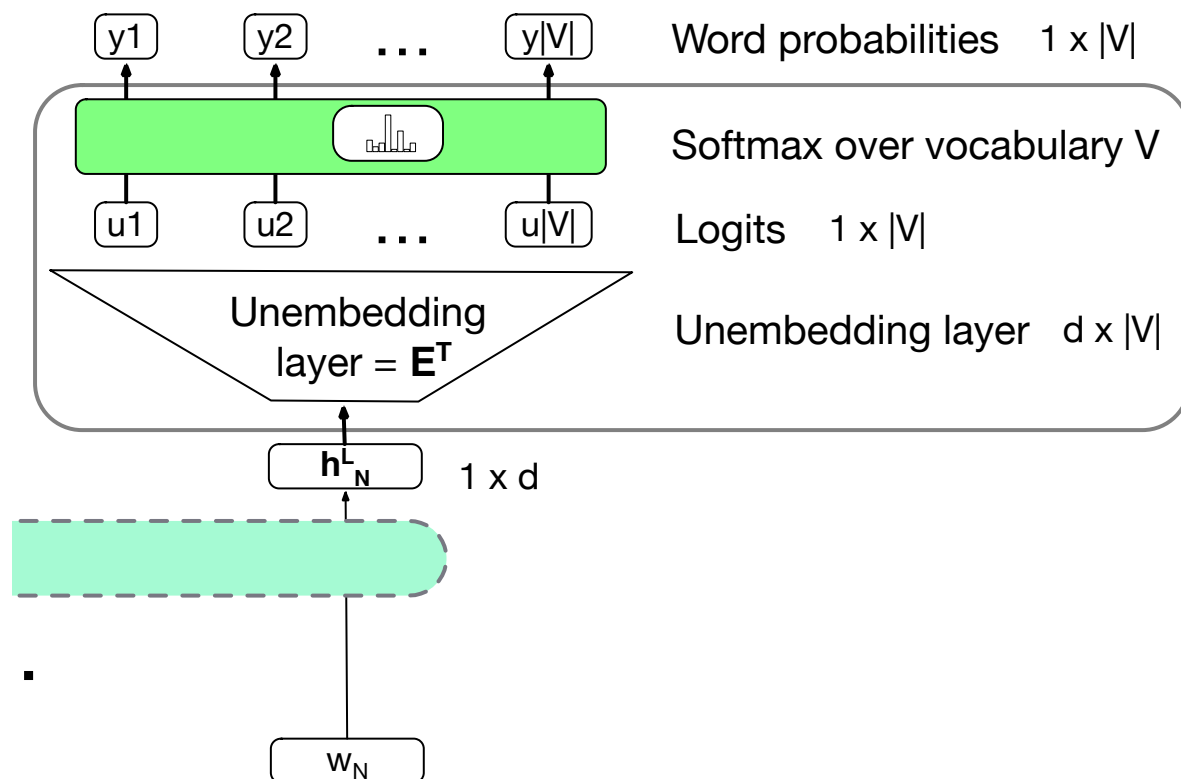
Unembedding layer maps from an embedding to a $1 \times |V|$ vector of logits

Language modeling head

Logits, the score vector u

One score for each of the $|V|$ possible words in the vocabulary V . Shape $1 \times |V|$.

Softmax turns the logits into probabilities over vocabulary. Shape $1 \times |V|$.

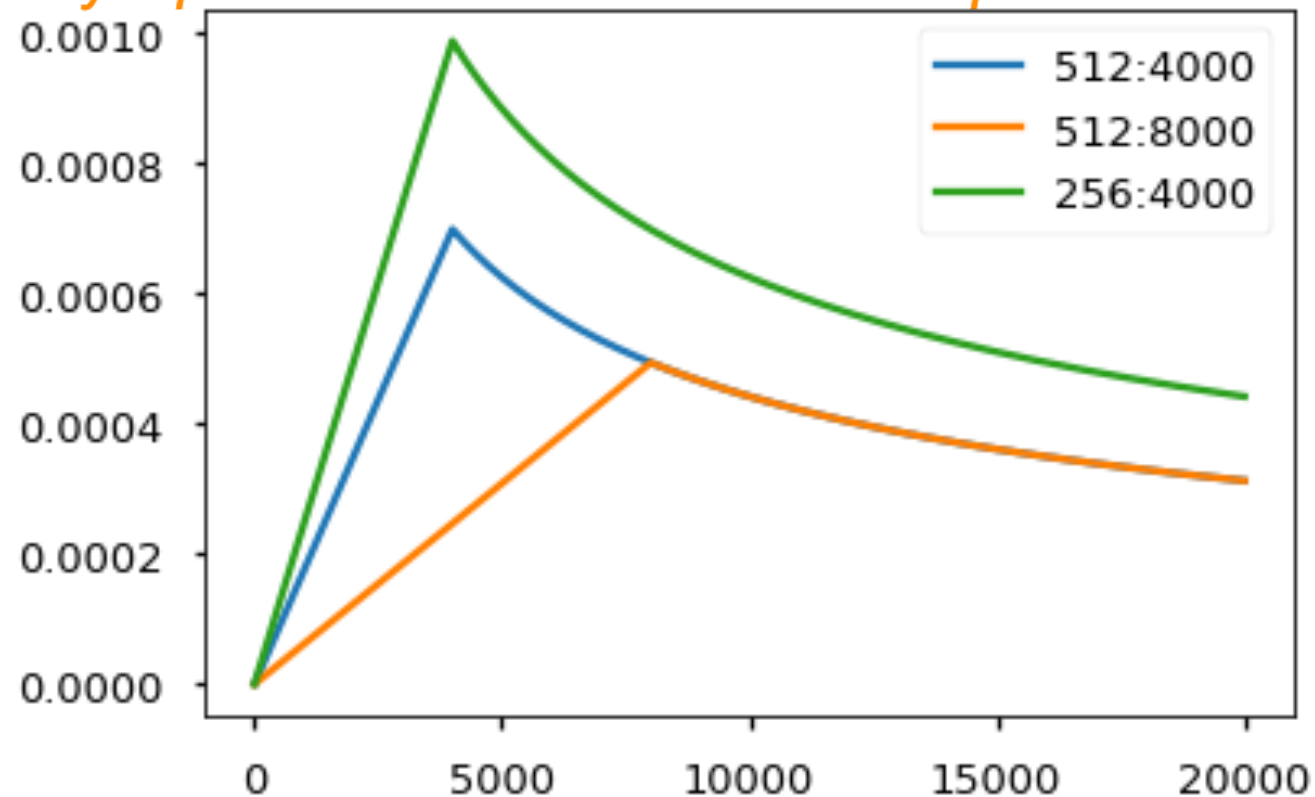


$$u = h_N^L E^T$$
$$y = \text{softmax}(u)$$

Optimizer

We used the Adam optimizer (cite) with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$. We varied the learning rate over the course of training, according to the formula: $lrate = d_{\text{model}}^{-0.5} \cdot \min(step_num^{-0.5}, step_num \cdot warmup_steps^{-1.5})$ This corresponds to increasing the learning rate linearly for the first $warmup_steps$ training steps, and decreasing it thereafter proportionally to the inverse square root of the step number. We used $warmup_steps = 4000$.

Note: This part is very important. Need to train with this setup of the model.



- Training instability is a notorious issue
 - Esp. with many layers, >10 or >20
- Yet something is going right. Not clear why!

Byte pair encoding (BPE)

- Deal with rare words / large vocabulary by instead using **subword** tokenization

system	sentence
source	health research institutes
reference	Gesundheitsforschungsinstitute
WDict	Forschungsinstitute
C2-50k	Fo rs ch un gs in st it ut io ne n
BPE-60k	Gesundheits forsch ungsinstitu ten
BPE-J90k	Gesundheits forsch ungsin stitute
source	asinine situation
reference	dumme Situation
WDict	asinine situation → UNK → asinine
C2-50k	as in in e situation → As in en si tu at io n
BPE-60k	as in ine situation → A in line- Situation
BPE-J90K	as in ine situation → As in in- Situation

