



Language  
Technologies  
Institute

Carnegie  
Mellon  
University

# Algorithms for NLP

CS 11-711 · Fall 2020

**Lecture 10: Masked language modeling and Transformers**

Emma Strubell

# Announcements

- Share your writeup from P1 (bake-off design and paper summary) on Piazza!
- Friday's recitation will be an overview of P2.

# Previously, in language model pretraining...

## Language models as word representations

- A problem with (most) word2vec-style word embeddings: single representation per word type.
  - the new-look **play** area is due to be completed by early spring 2020
  - gerrymandered congressional districts favor representatives who **play** to the party base
  - the freshman then completed the three-point **play** for a 66-63 lead
- Instead, want **contextualized** word embeddings: token representations that differ depending on the context.
  - But, still want to leverage self-supervised training on massive data.
  - Solution: use hidden representations from a neural language model.

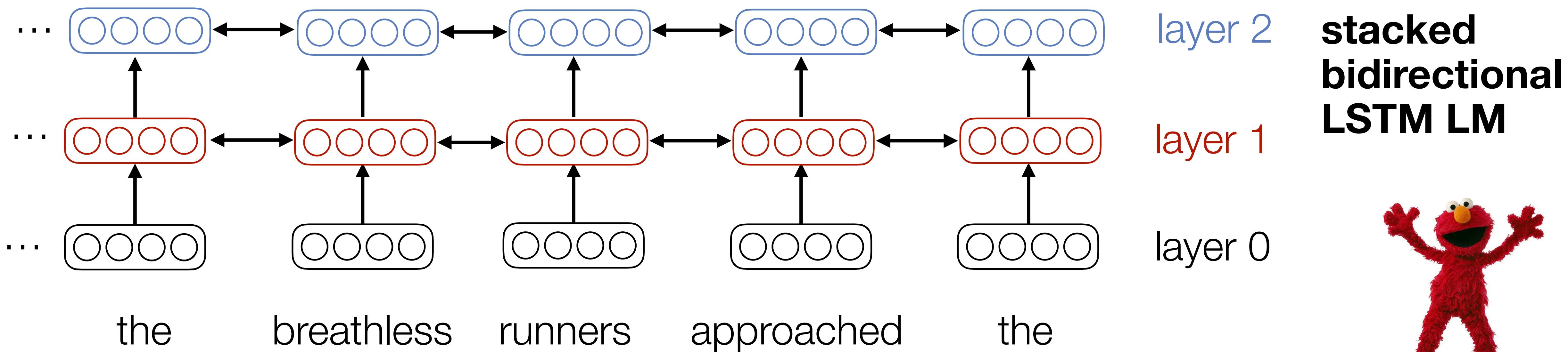
# ELMo: Deep contextualized word embeddings

- Key idea: context-dependent embedding for each word interpolates representations for that word from each layer

$$\text{breathless} = \gamma ( s_1 \cdot \text{blue box} + s_2 \cdot \text{red box} + s_3 \cdot \text{black box} )$$

$$\sum_{j=1}^L s_j = 1$$

- Interpolation weights are task-specific (fine-tuned on supervised data.)



# ELMo: Deep contextualized word embeddings

- Adding ELMo to existing state-of-the-art models provides significant improvement on essentially all NLP tasks.

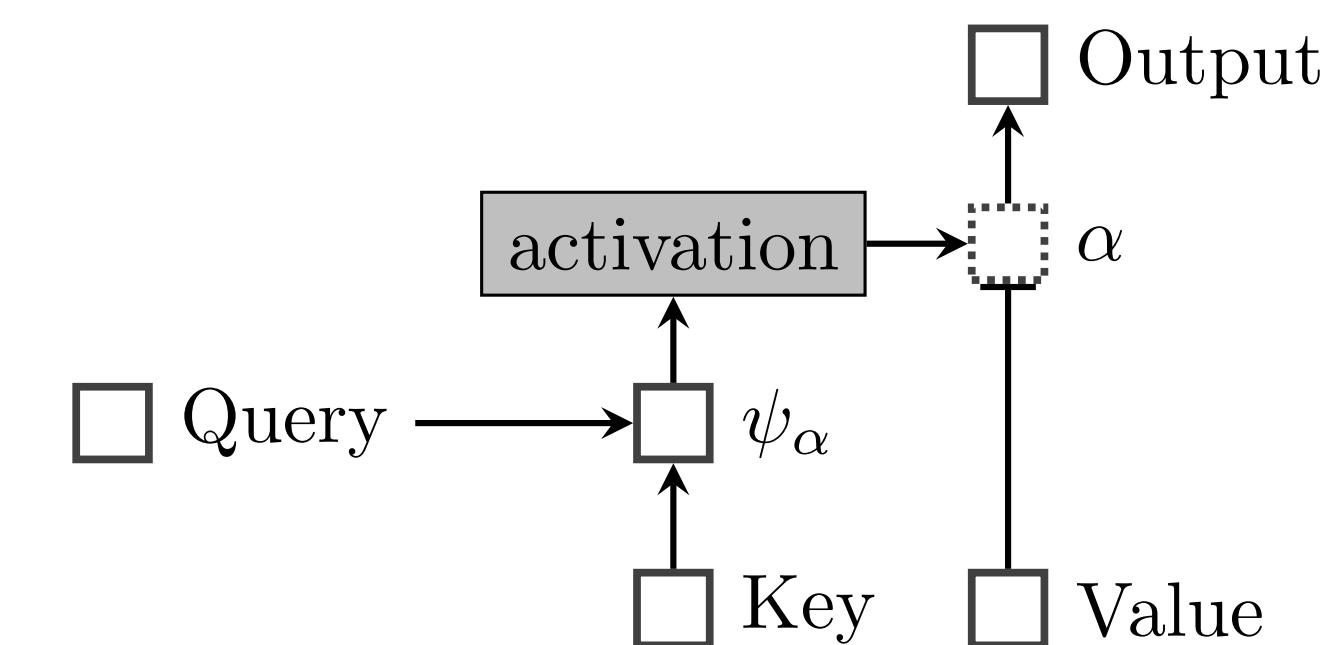
TASK	PREVIOUS SOTA	OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/RELATIVE)
question answering	SQuAD Liu et al. (2017)	84.4	81.1	85.8
natural language inference	SNLI Chen et al. (2017)	88.6	88.0	$88.7 \pm 0.17$
semantic role labeling	SRL He et al. (2017)	81.7	81.4	84.6
coreference	Coref Lee et al. (2017)	67.2	67.2	70.4
named entity recognition	NER Peters et al. (2017)	$91.93 \pm 0.19$	90.15	$92.22 \pm 0.10$
sentiment analysis	SST-5 McCann et al. (2017)	53.7	51.4	$54.7 \pm 0.5$

# Attention in neural networks

# Attention in neural networks

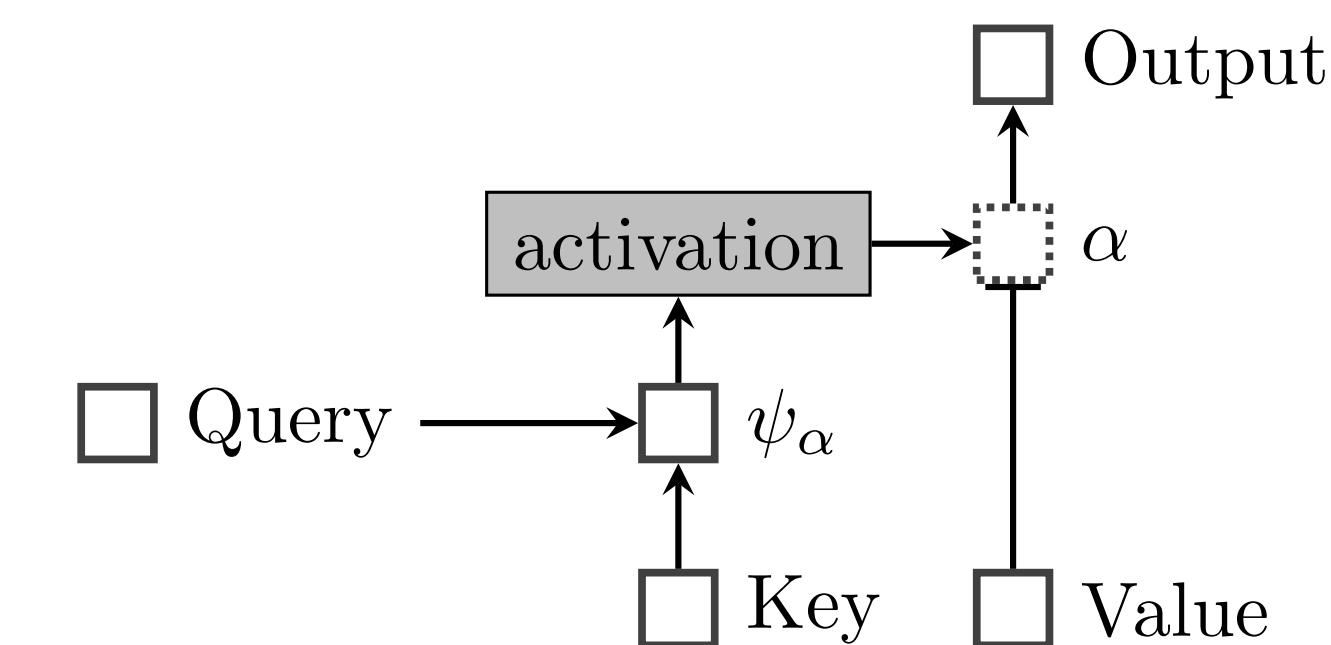
# Attention in neural networks

- **Neural attention mechanism:** a weighted average over (word) representations, where the weights are learned, computed as a function of the input.



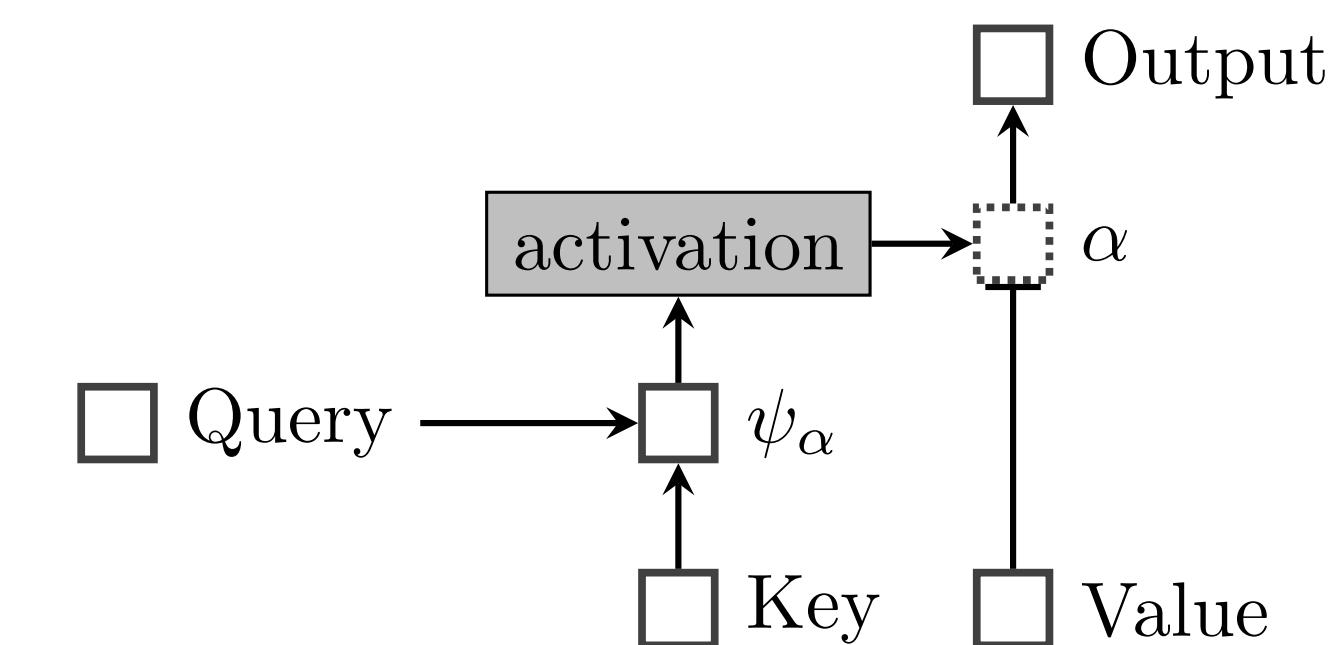
# Attention in neural networks

- **Neural attention mechanism:** a weighted average over (word) representations, where the weights are learned, computed as a function of the input.
- Additive attention:  $\psi$  is a feed-forward layer with concatenated  $[k;q]$  as input.



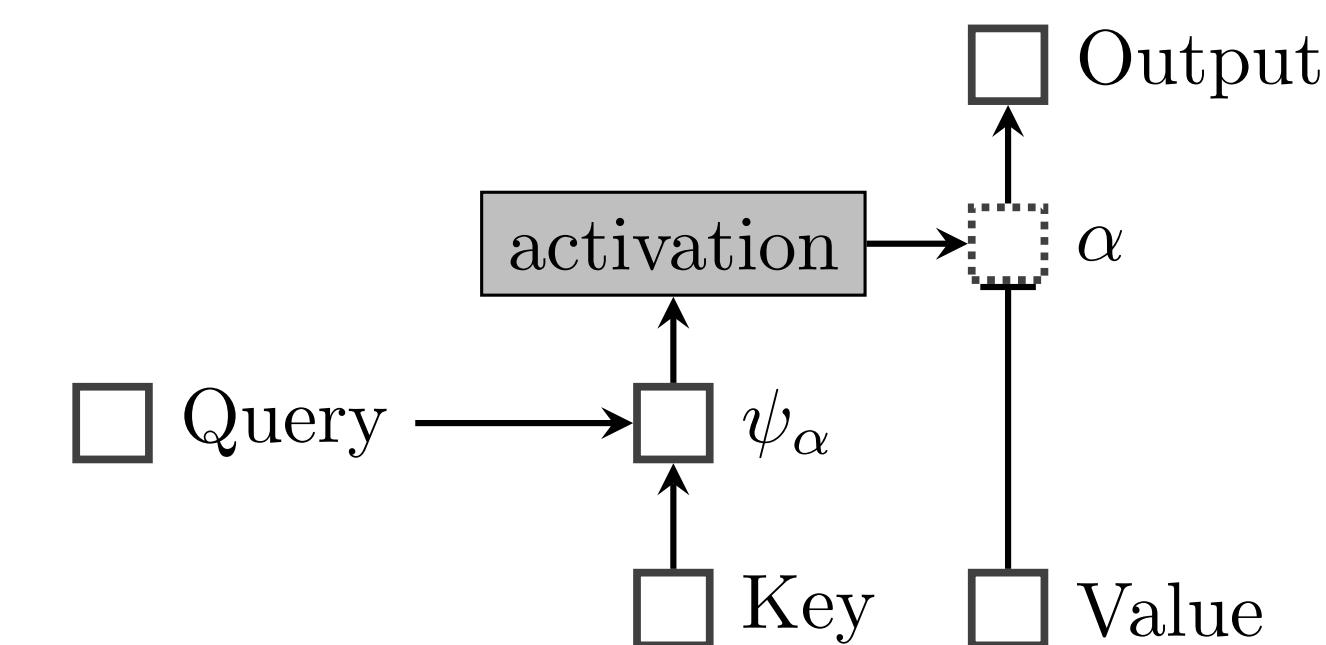
# Attention in neural networks

- **Neural attention mechanism:** a weighted average over (word) representations, where the weights are learned, computed as a function of the input.
  - Additive attention:  $\psi$  is a feed-forward layer with concatenated  $[k;q]$  as input.
  - Multiplicative attention:  $\psi$  is a dot product of linear projections of  $k, q$ .



# Attention in neural networks

- **Neural attention mechanism:** a weighted average over (word) representations, where the weights are learned, computed as a function of the input.
  - Additive attention:  $\psi$  is a feed-forward layer with concatenated  $[k;q]$  as input.
  - Multiplicative attention:  $\psi$  is a dot product of linear projections of  $k, q$ .
- First popularized in neural machine translation, where the network would attend over the source sentence to make each prediction in the target sentence.

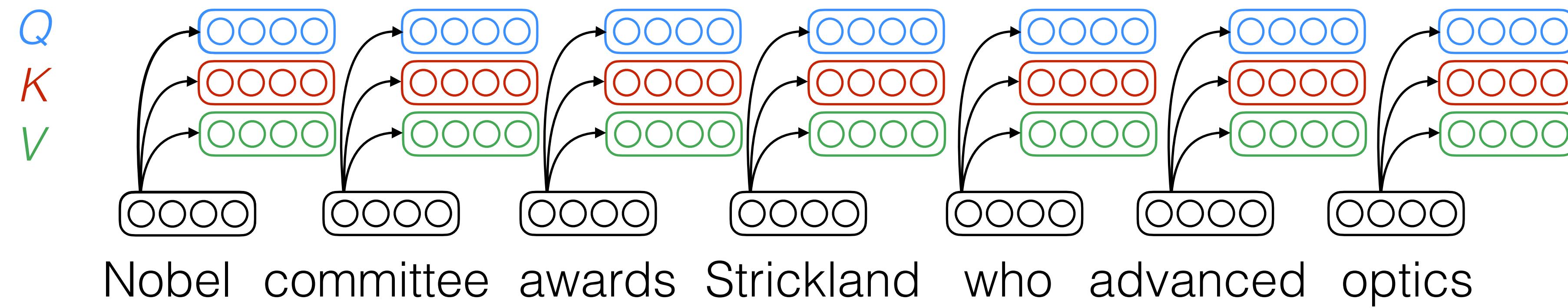


# Transformer self-attention

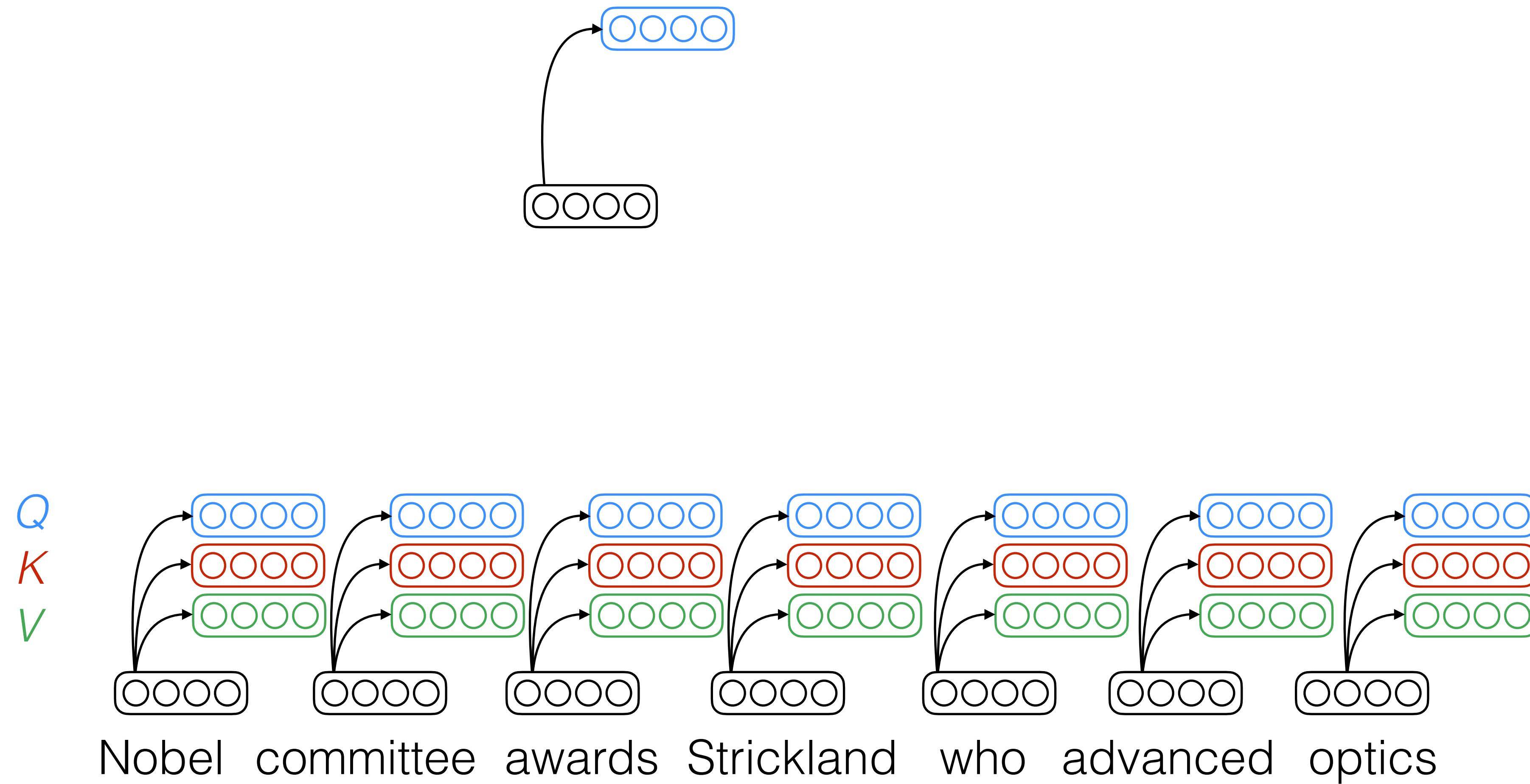
oooo oooo oooo oooo oooo oooo oooo

Nobel committee awards Strickland who advanced optics

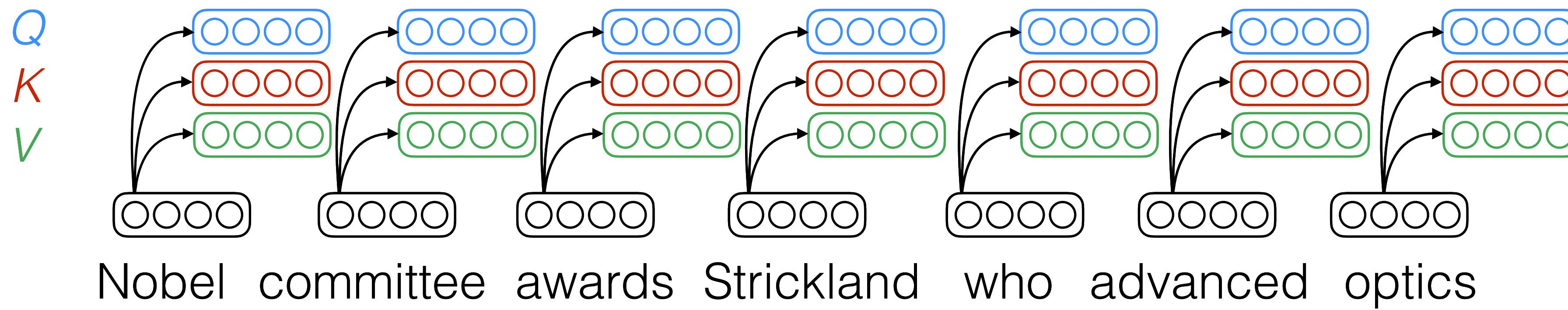
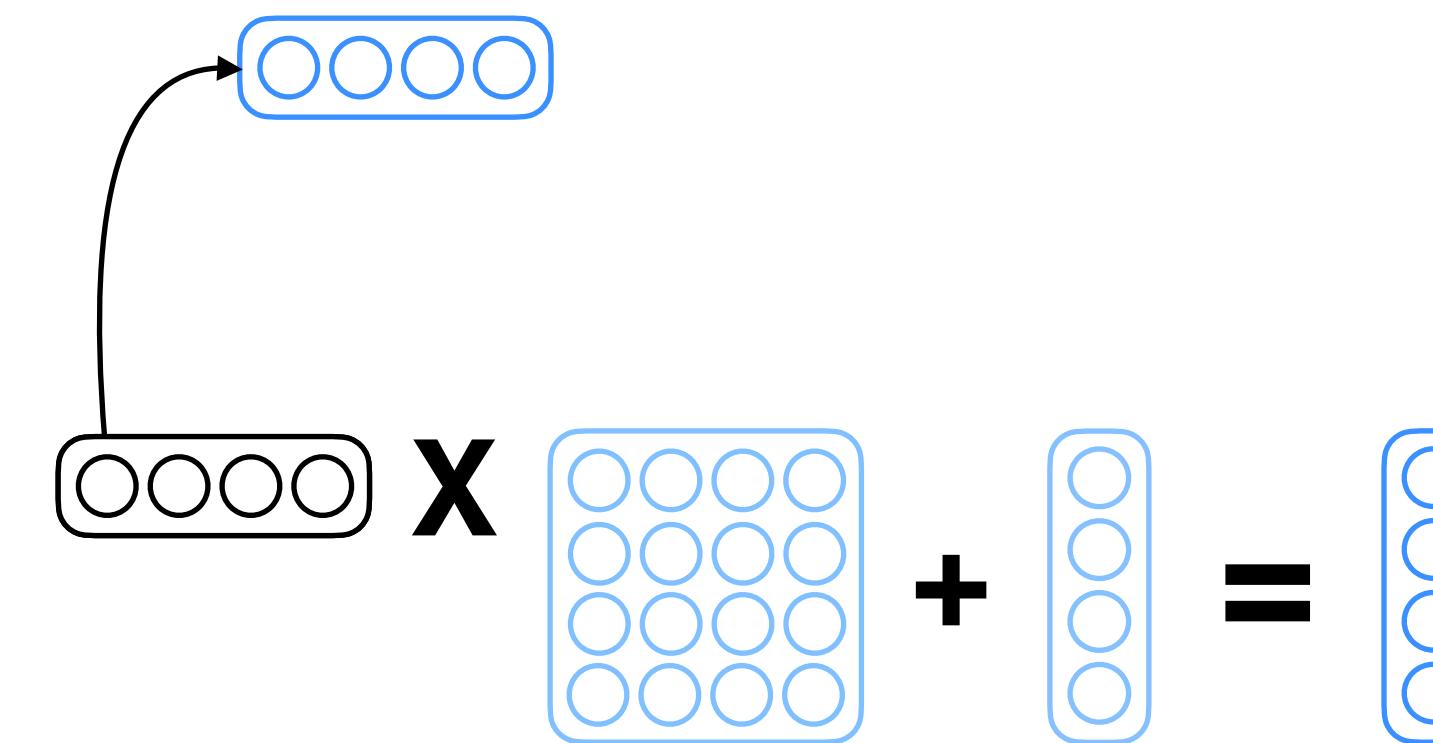
# Transformer self-attention



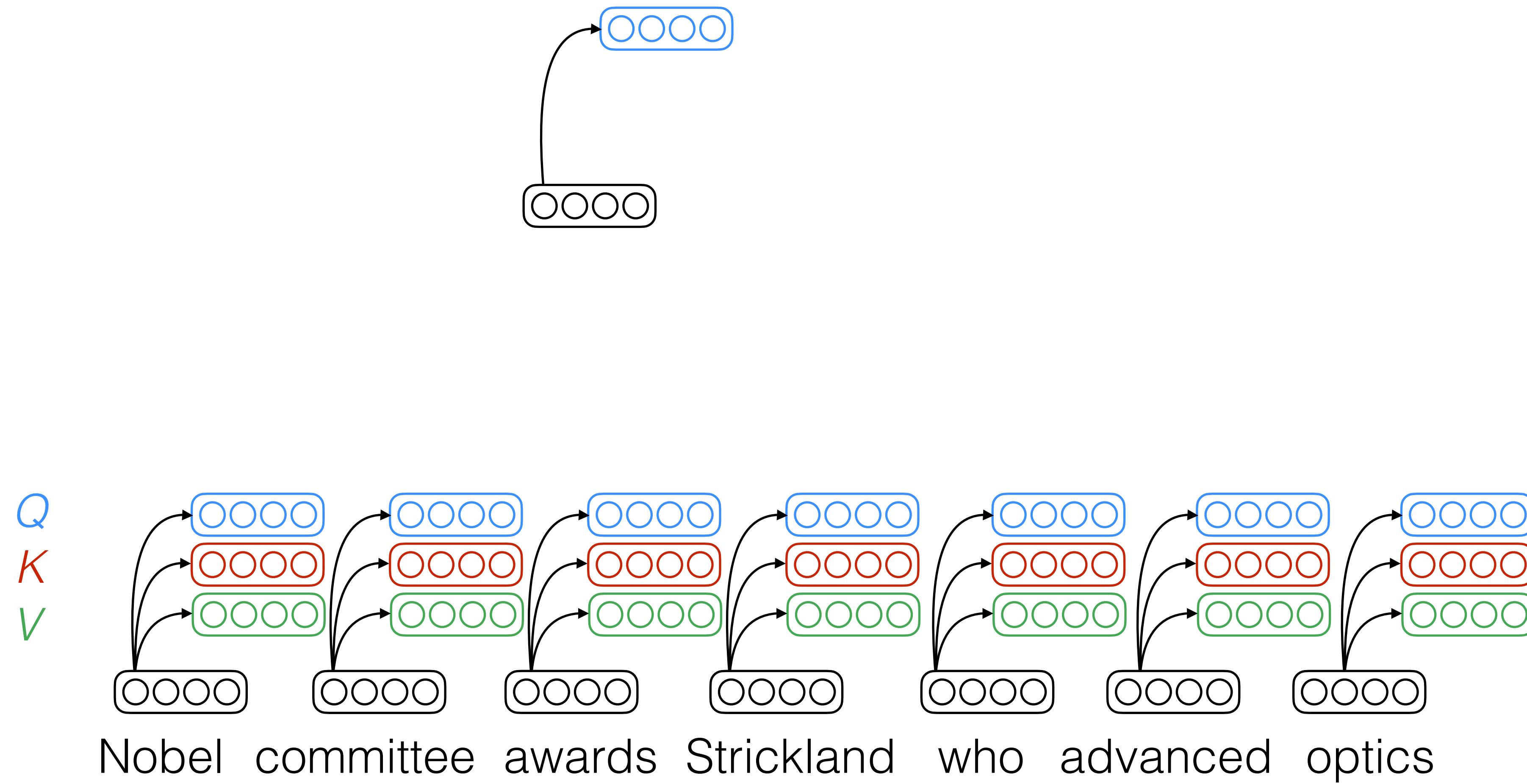
# Transformer self-attention



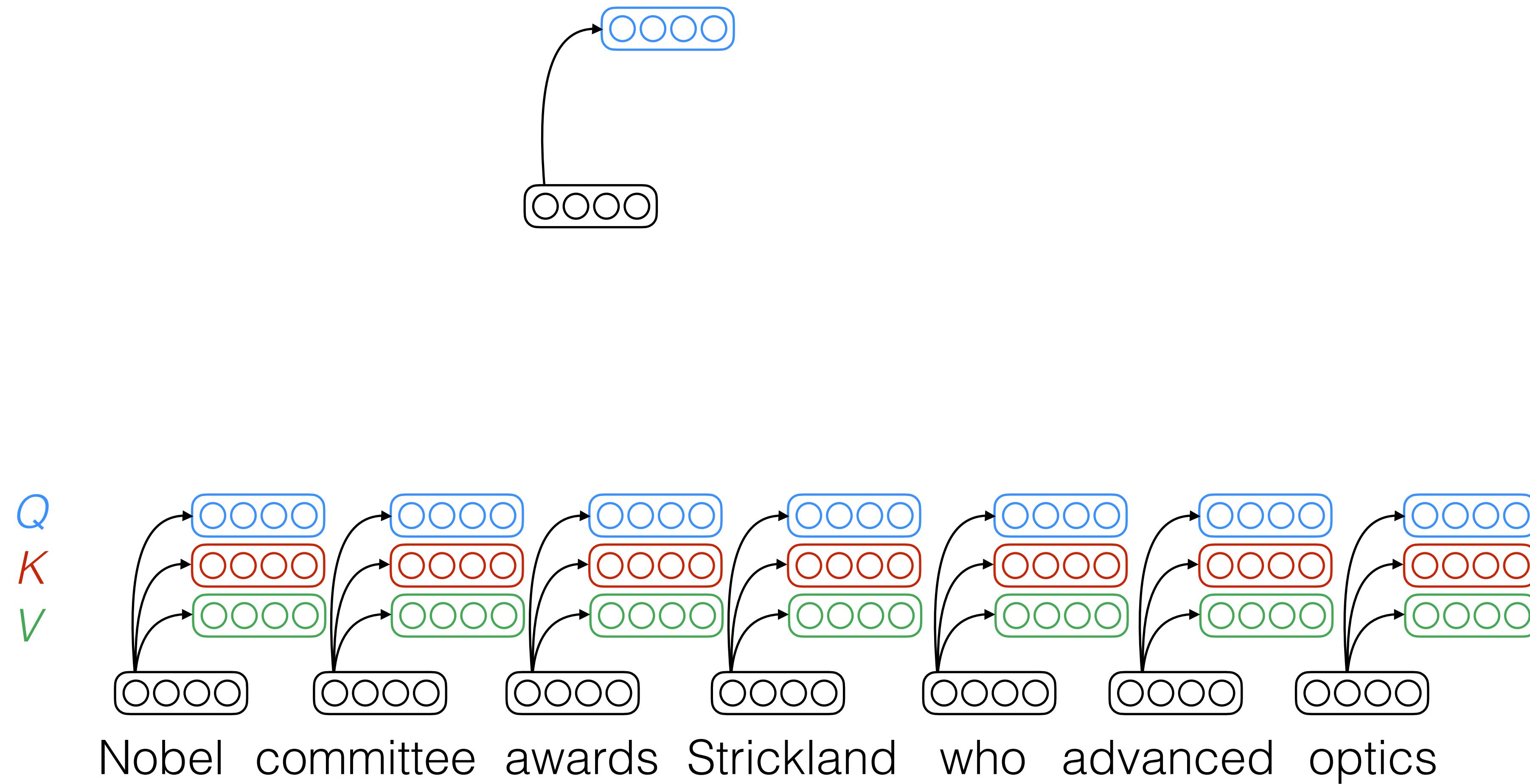
# Transformer self-attention



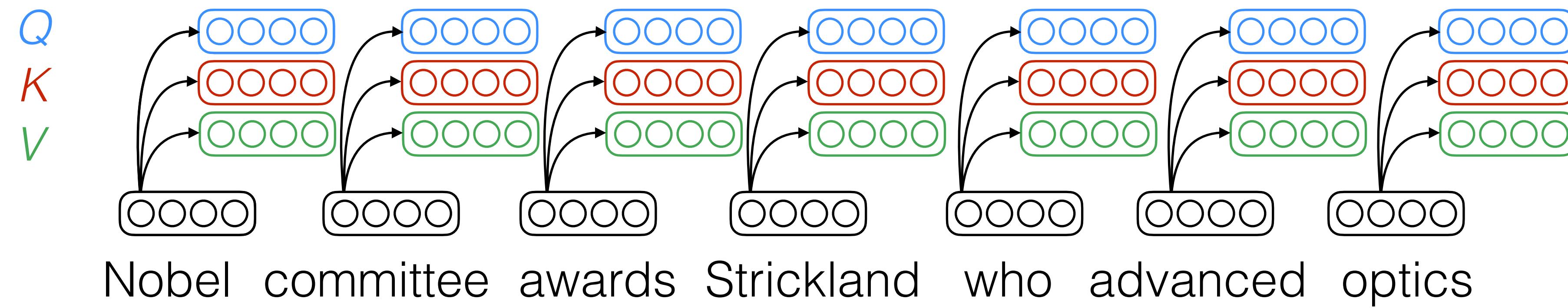
# Transformer self-attention



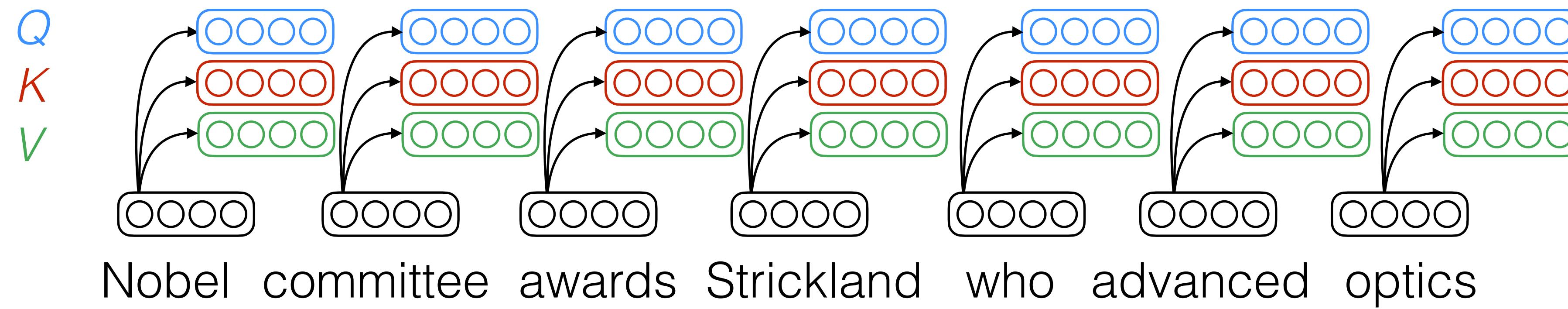
# Transformer self-attention



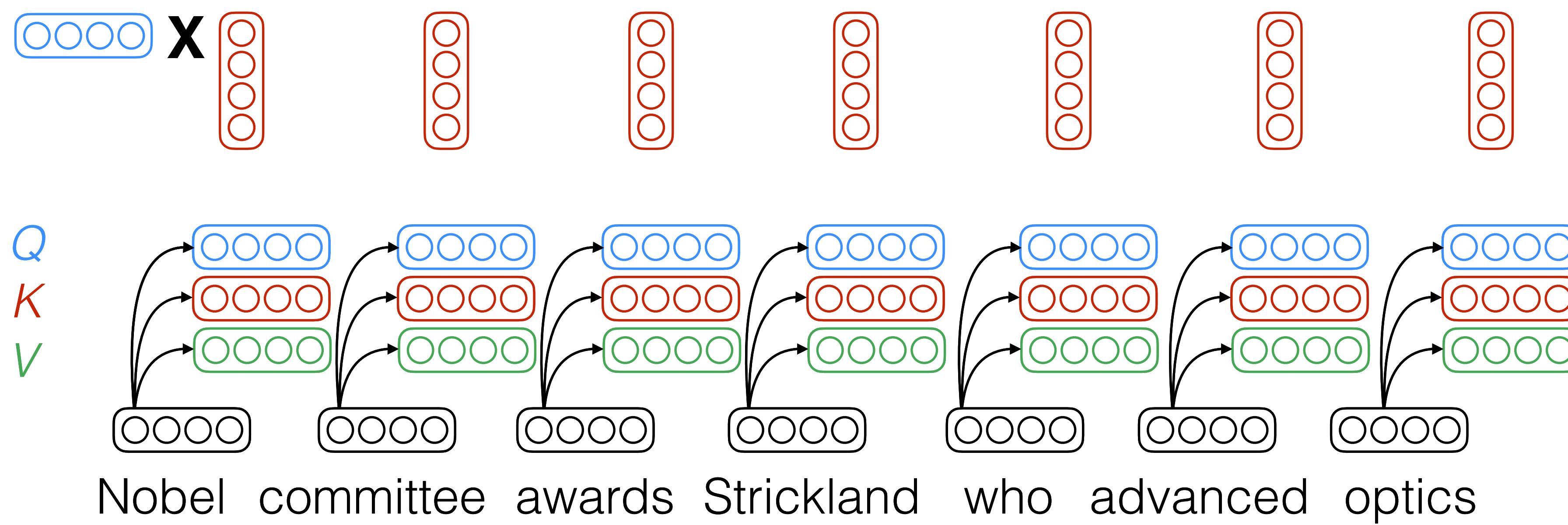
# Transformer self-attention



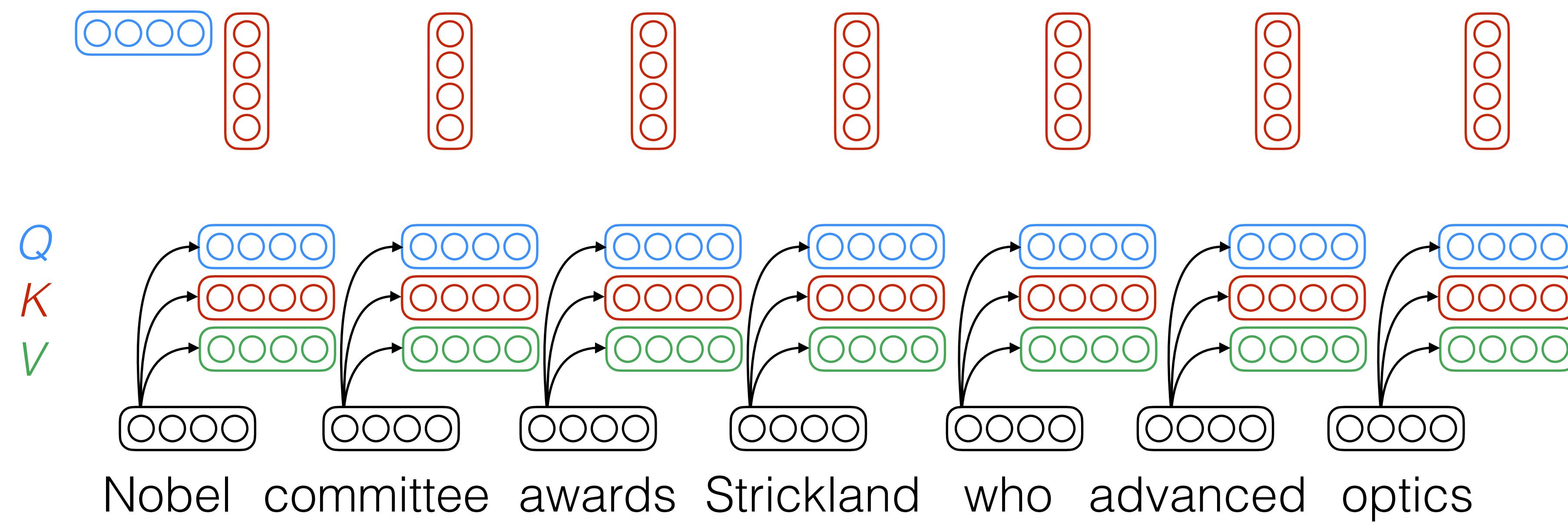
# Transformer self-attention



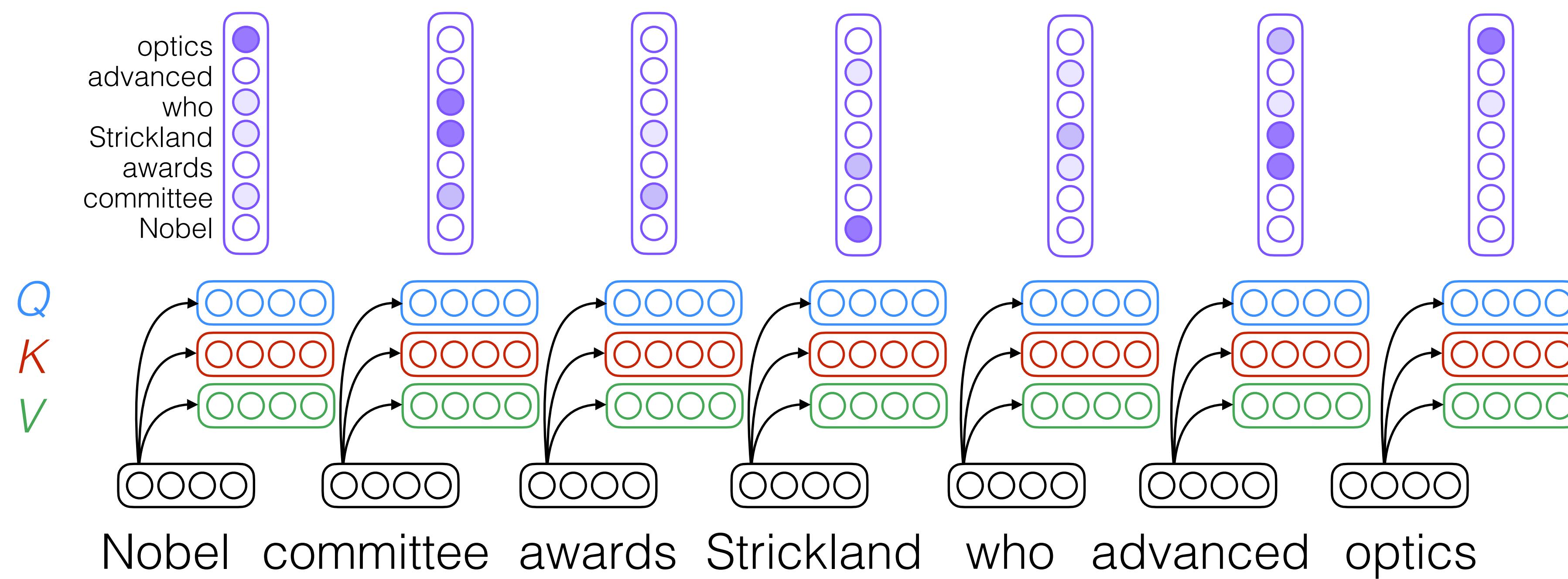
# Transformer self-attention



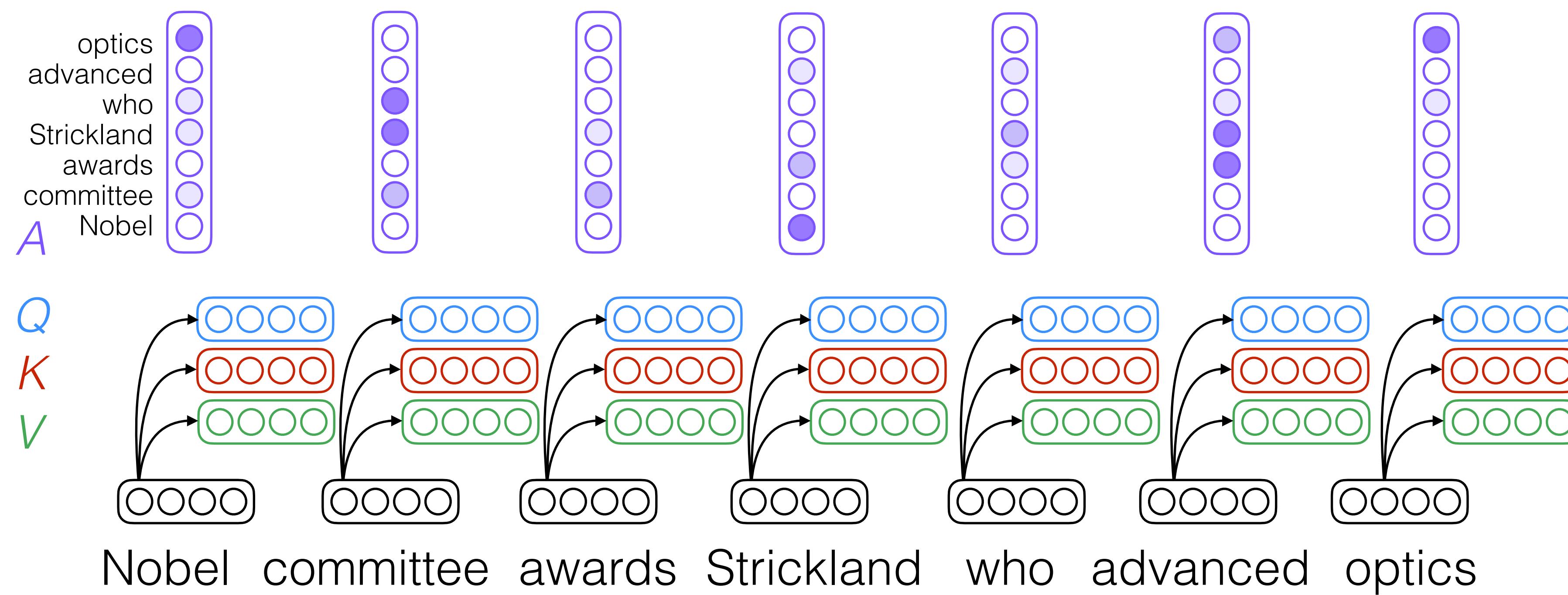
# Transformer self-attention



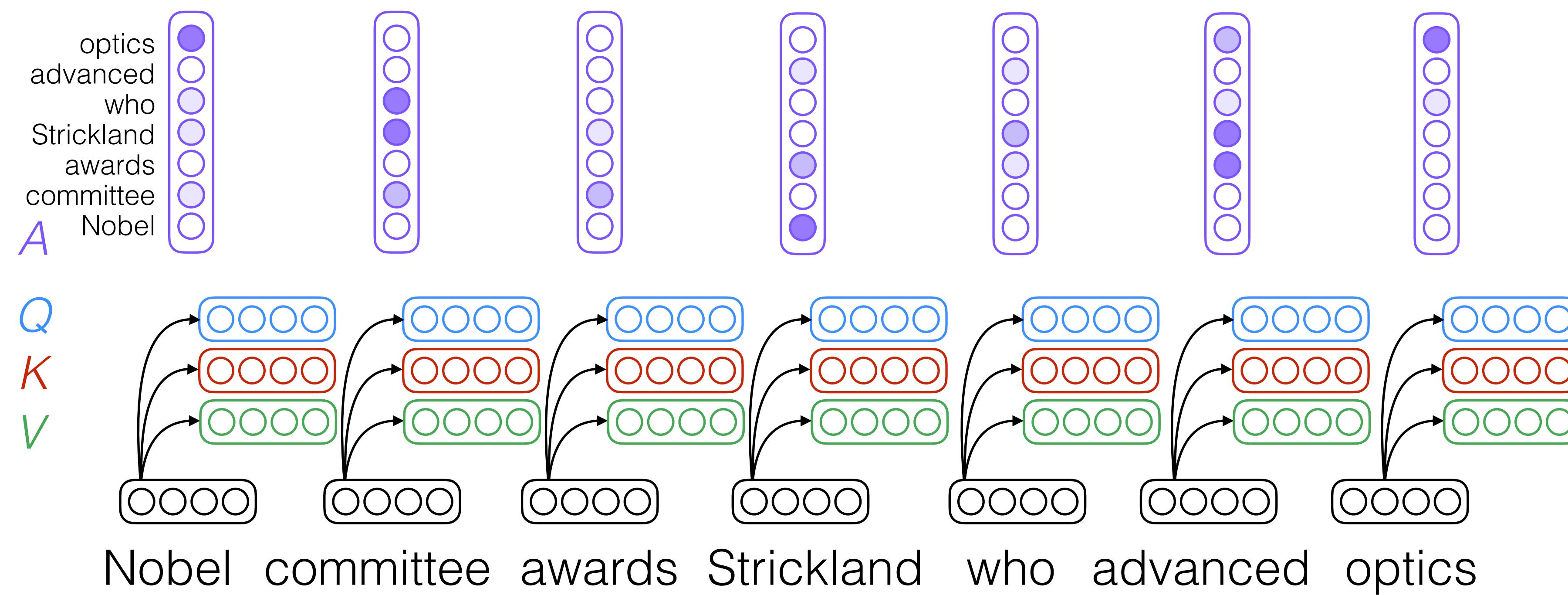
# Transformer self-attention



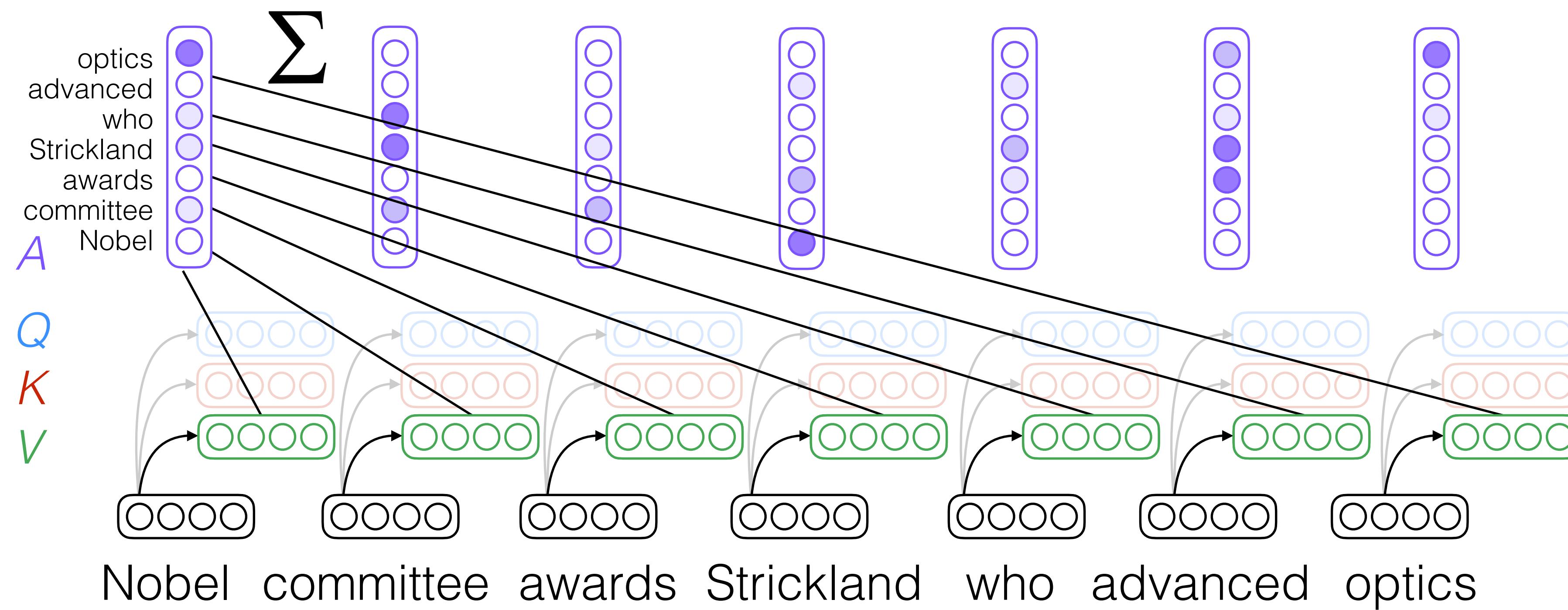
# Transformer self-attention



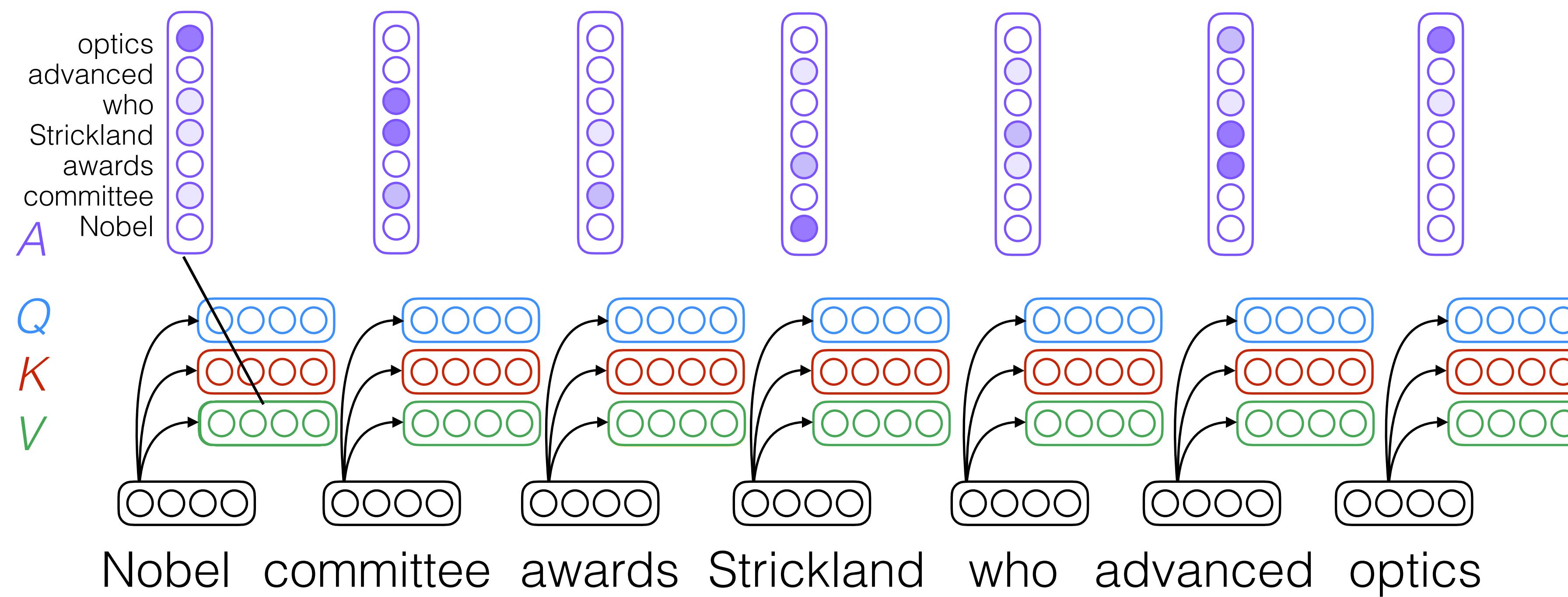
# Transformer self-attention



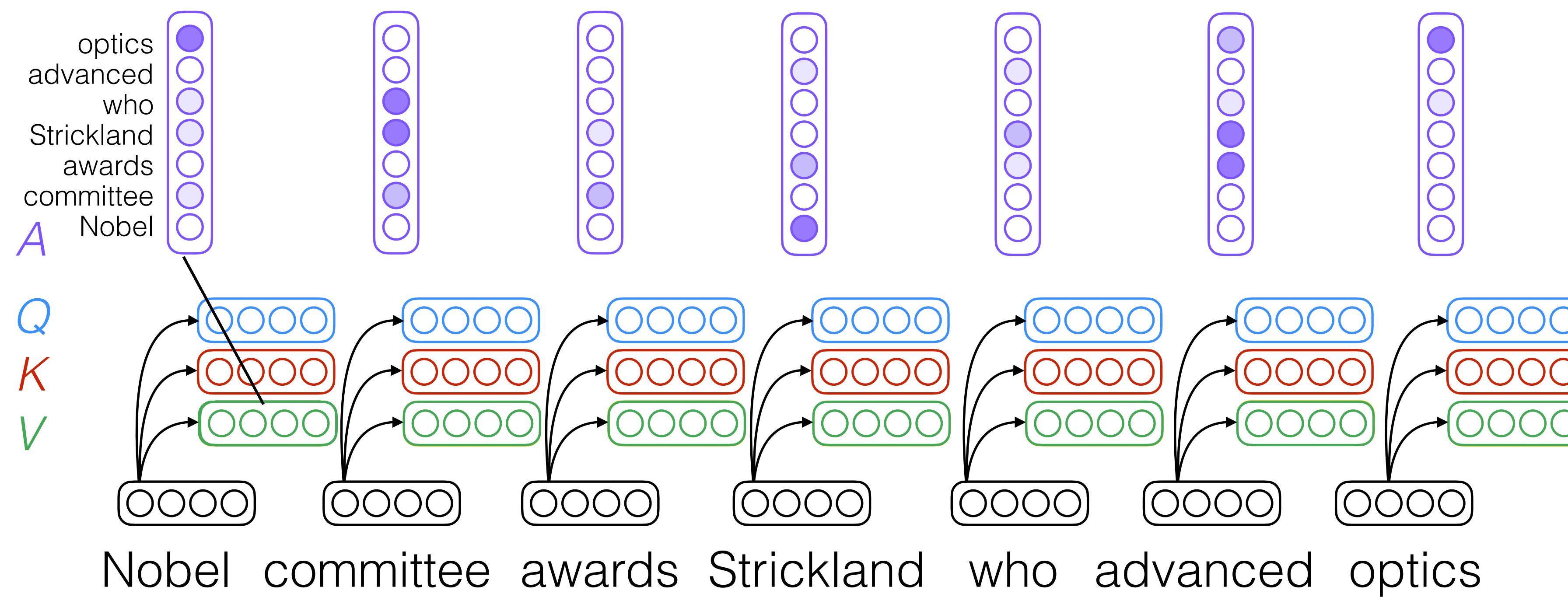
# Transformer self-attention



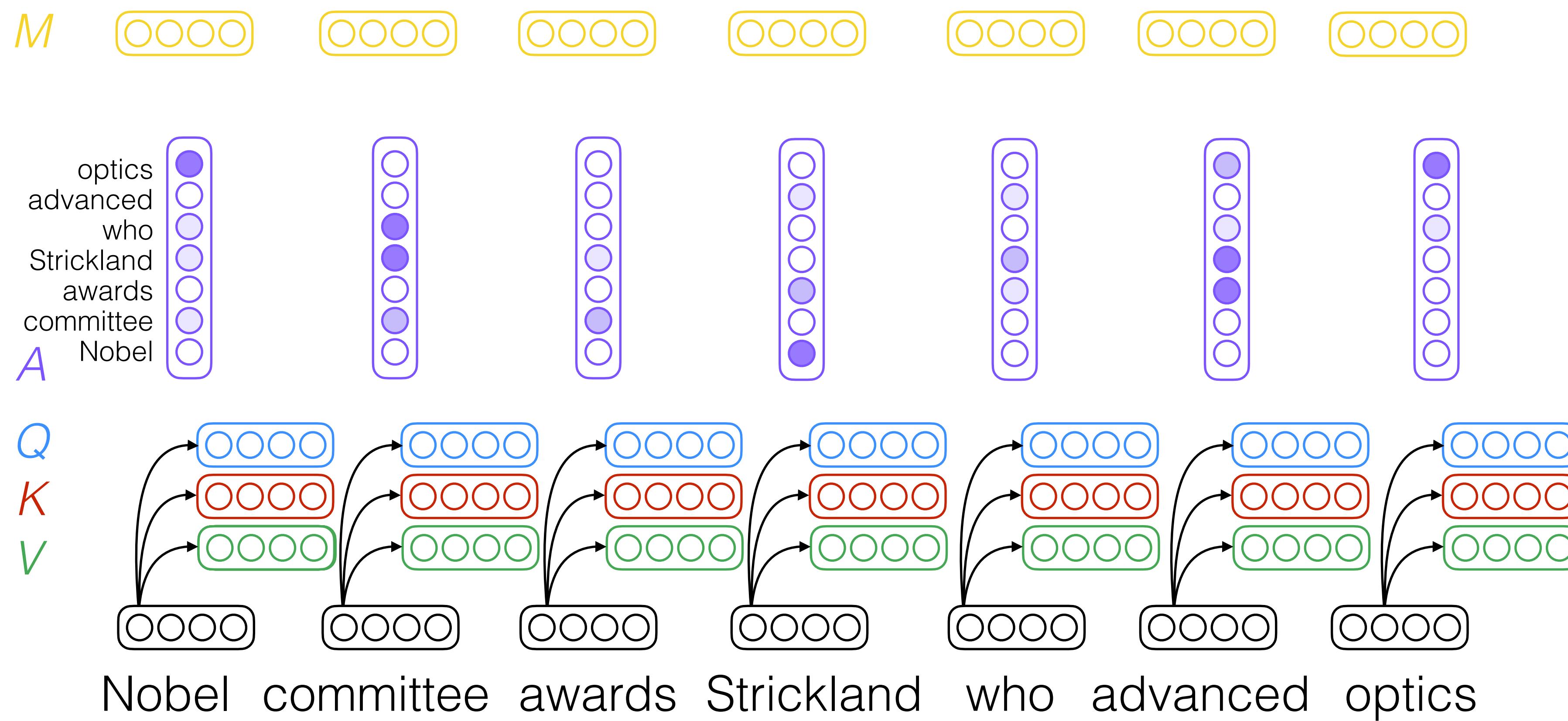
# Transformer self-attention



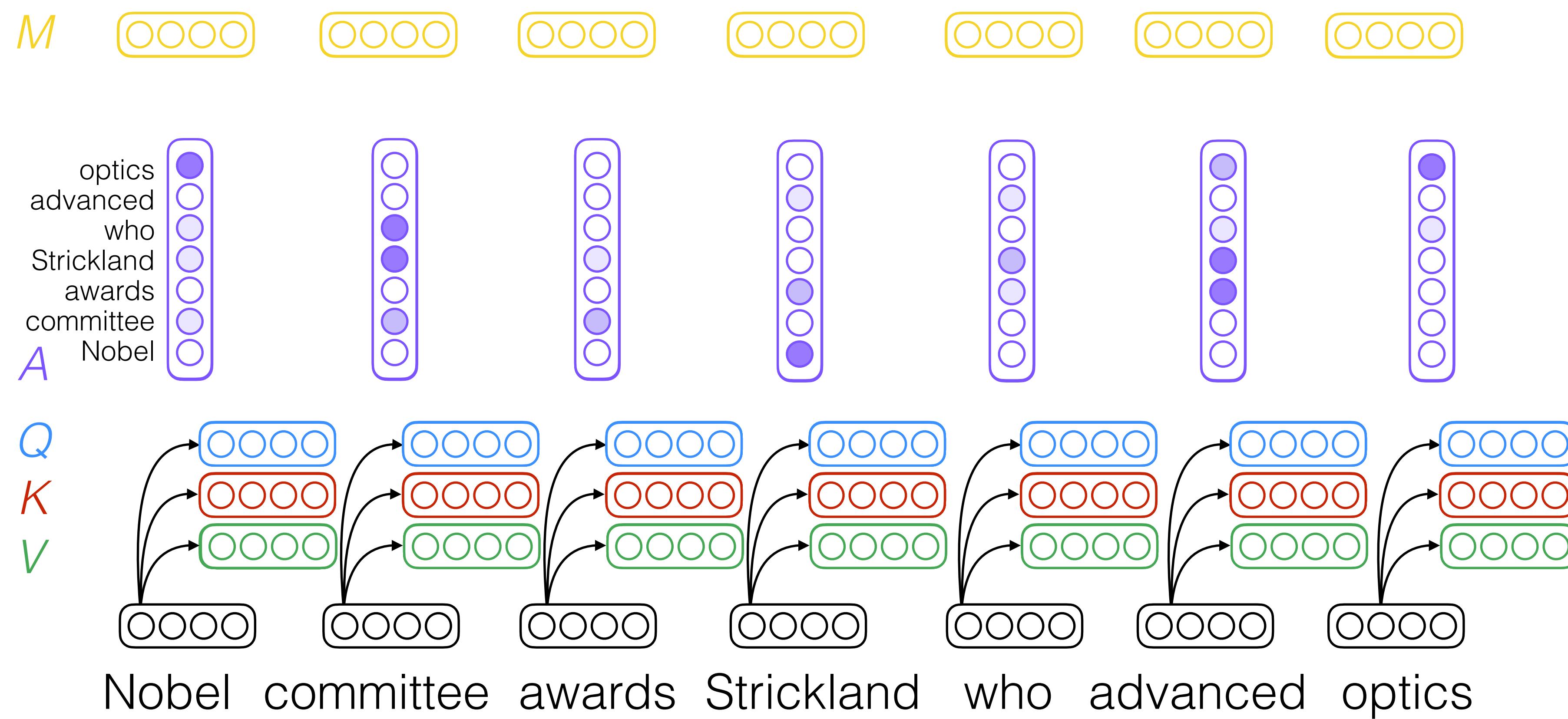
# Transformer self-attention



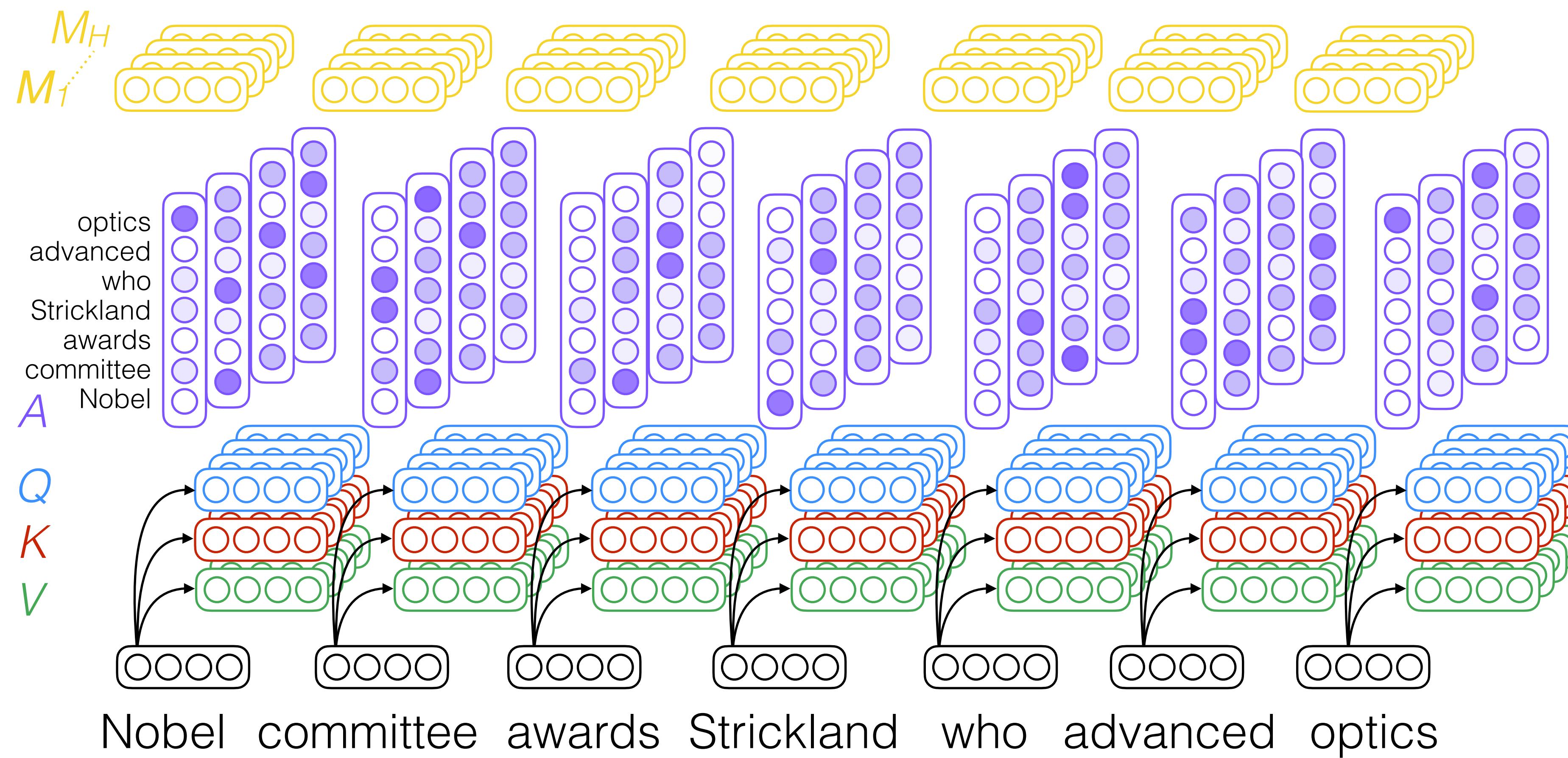
# Transformer self-attention



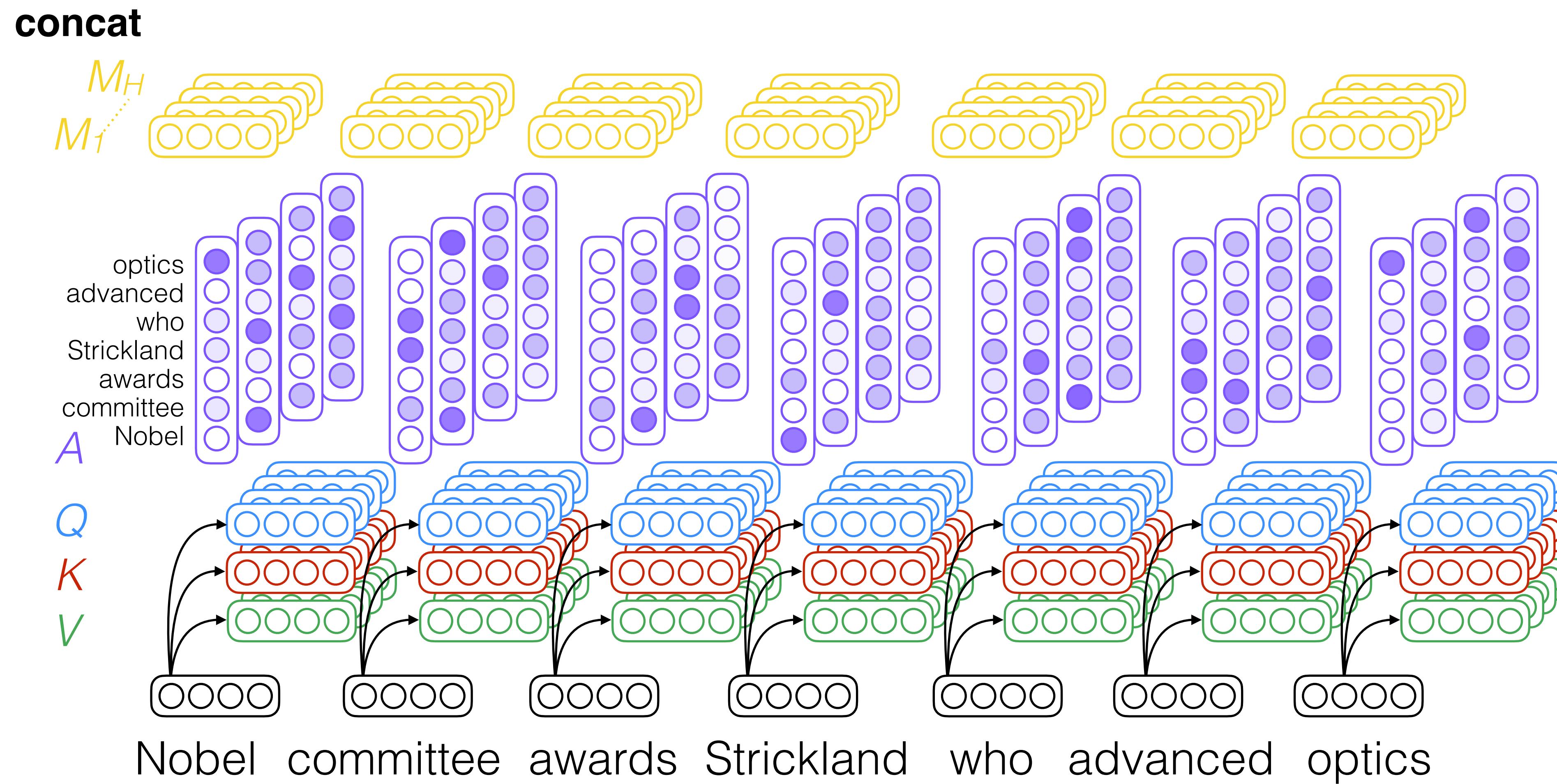
# Transformer self-attention



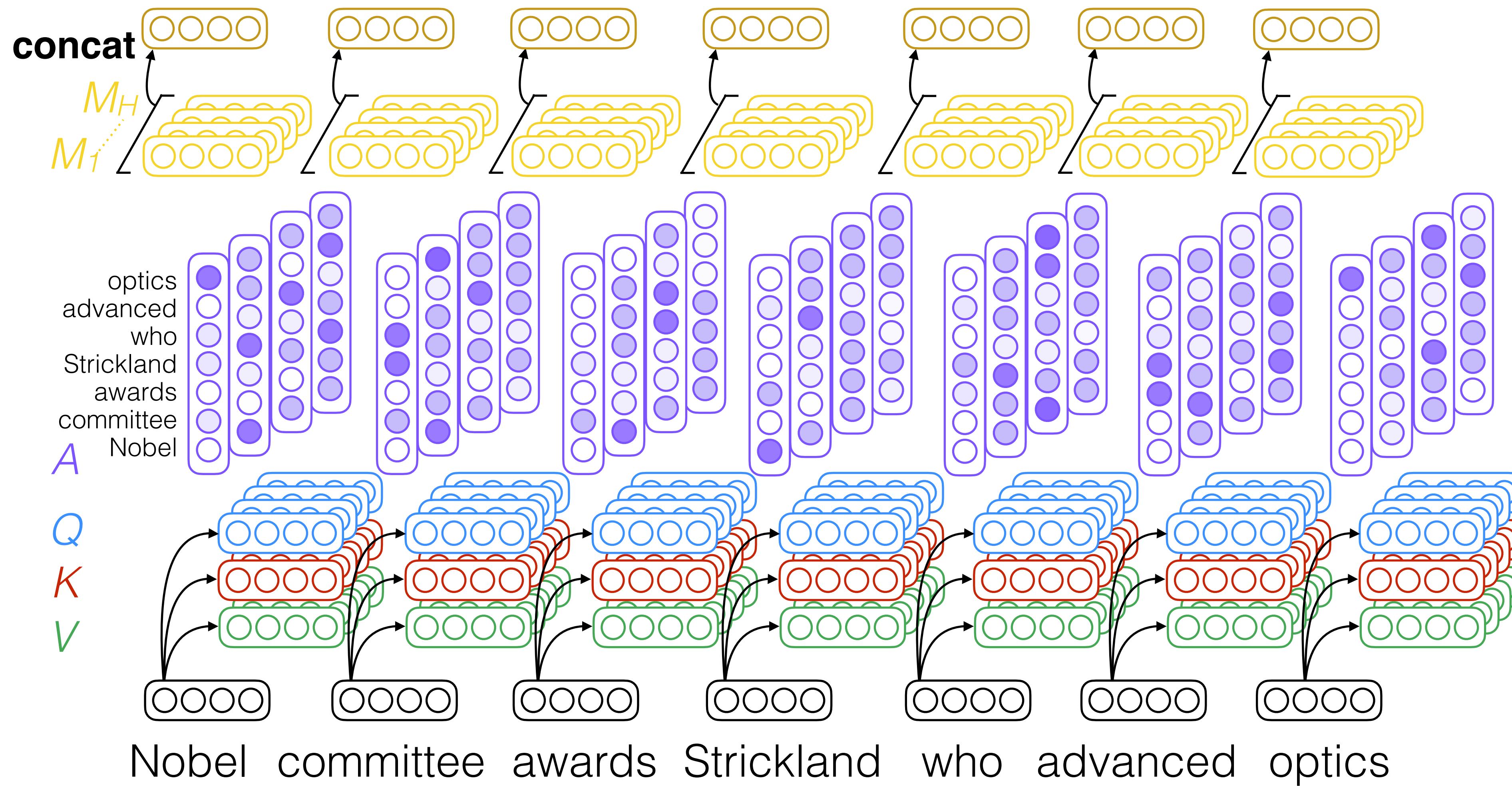
# Multi-head self-attention



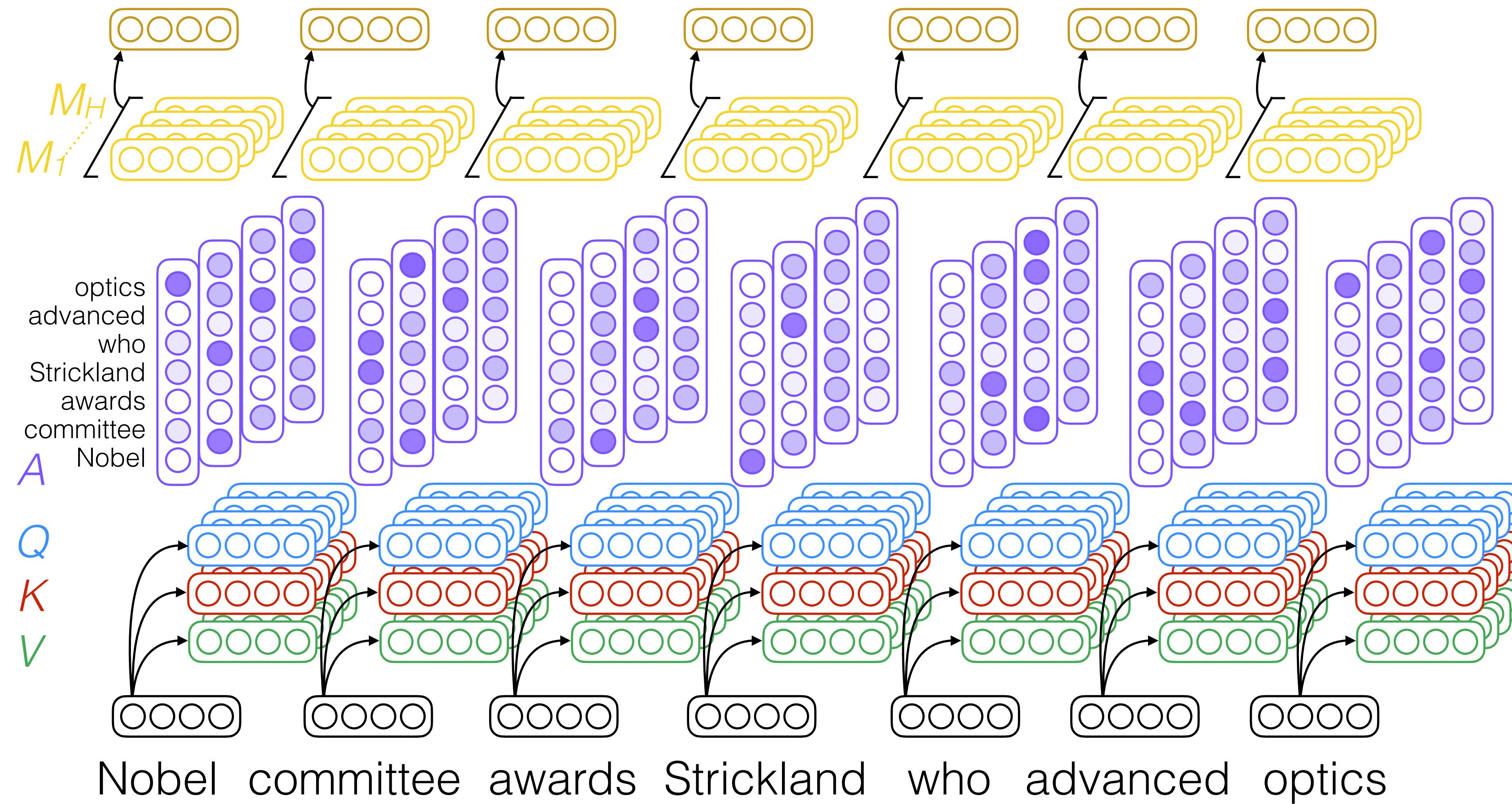
# Multi-head self-attention



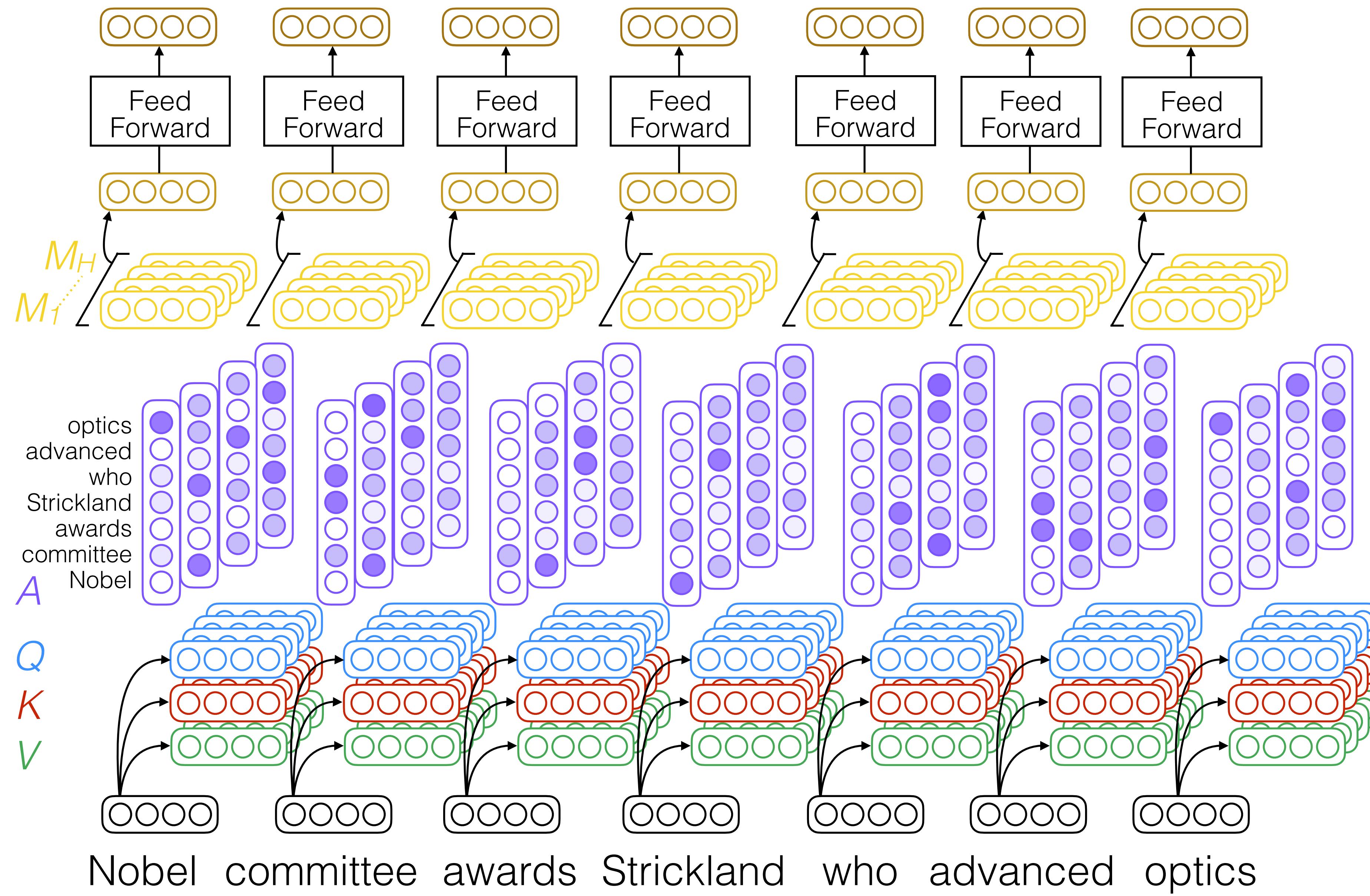
# Multi-head self-attention



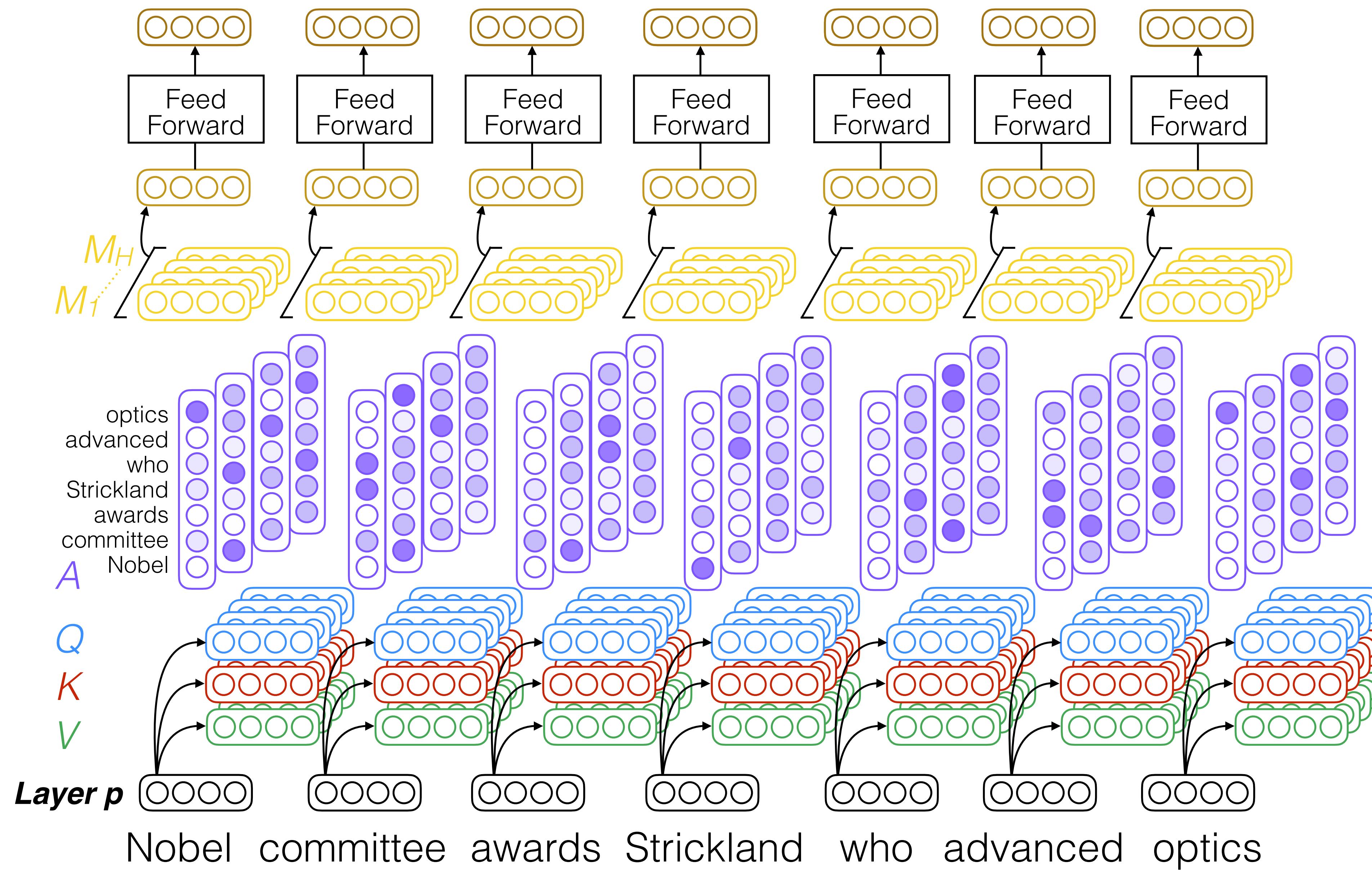
# Transformer layer



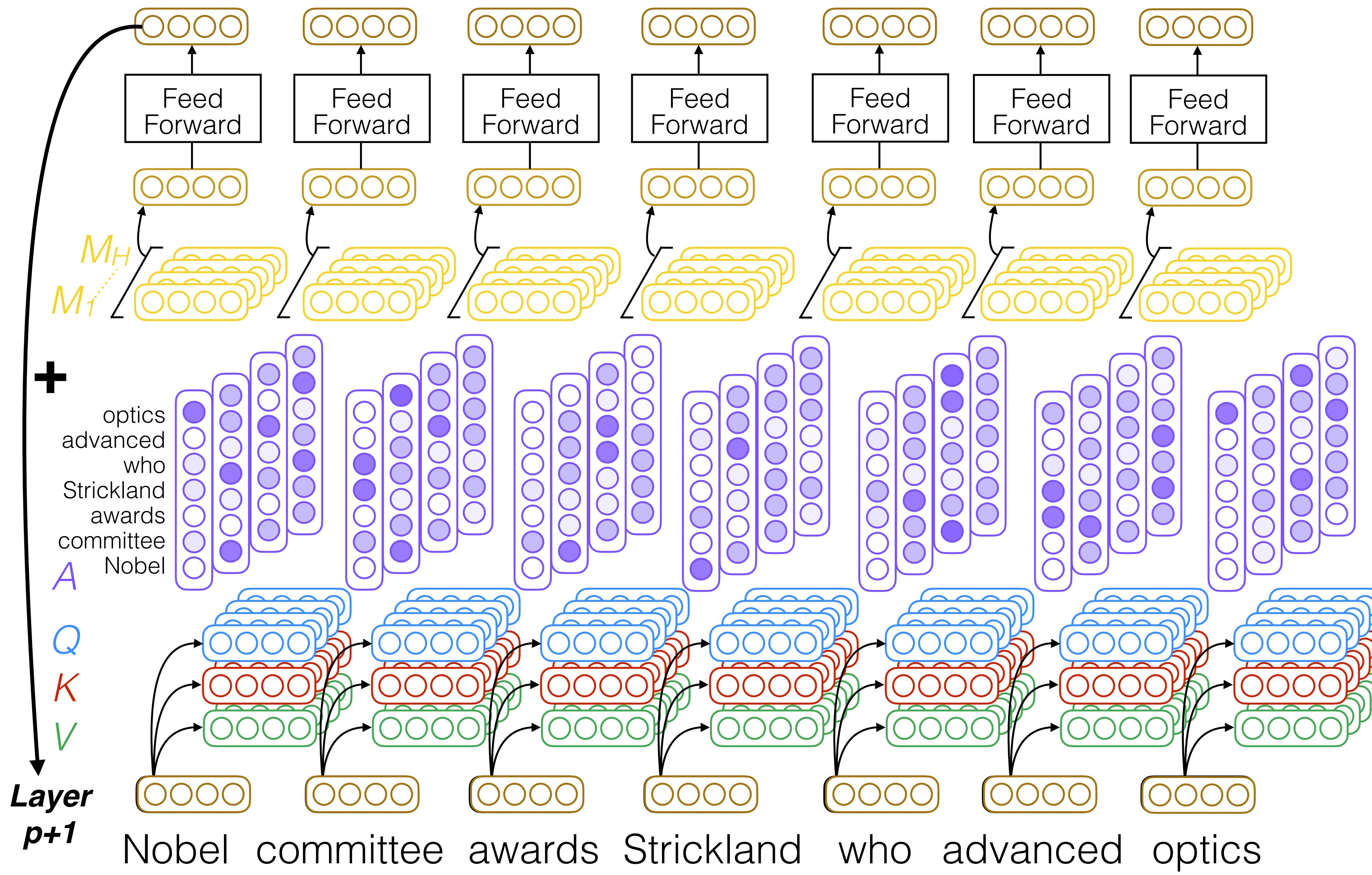
# Transformer layer



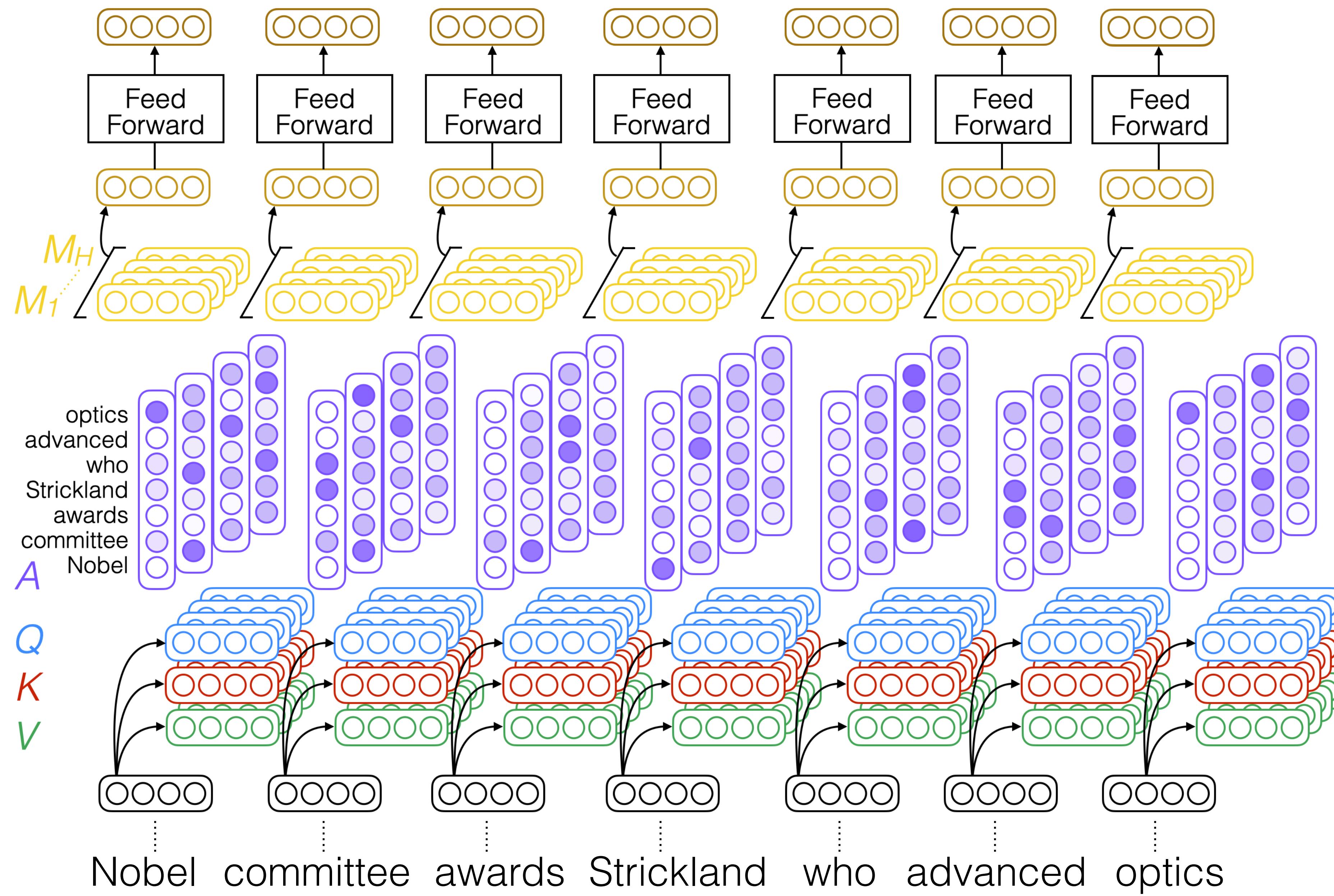
# Transformer layer



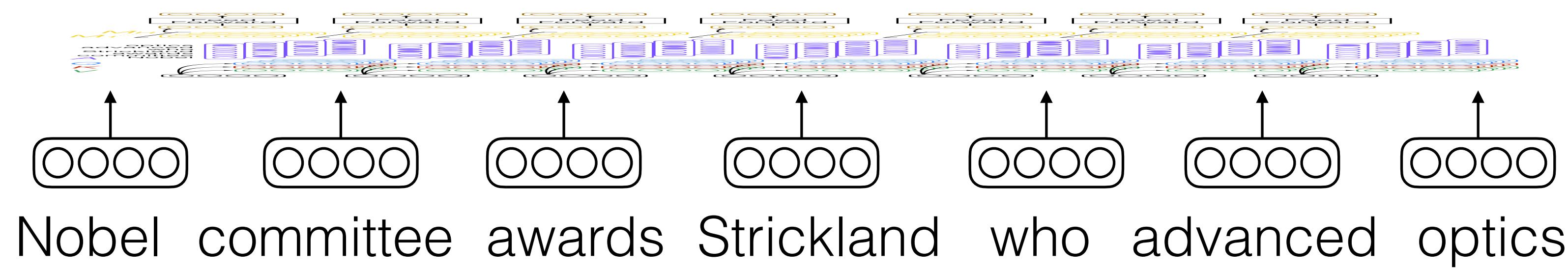
# Transformer layer



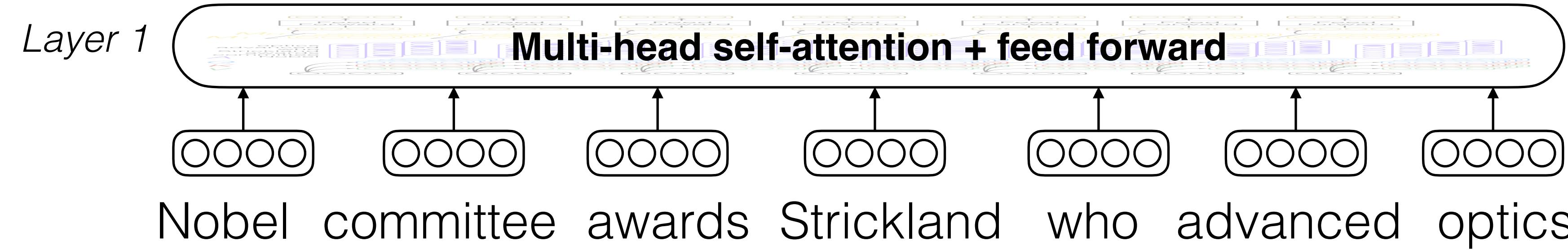
# Transformer layer



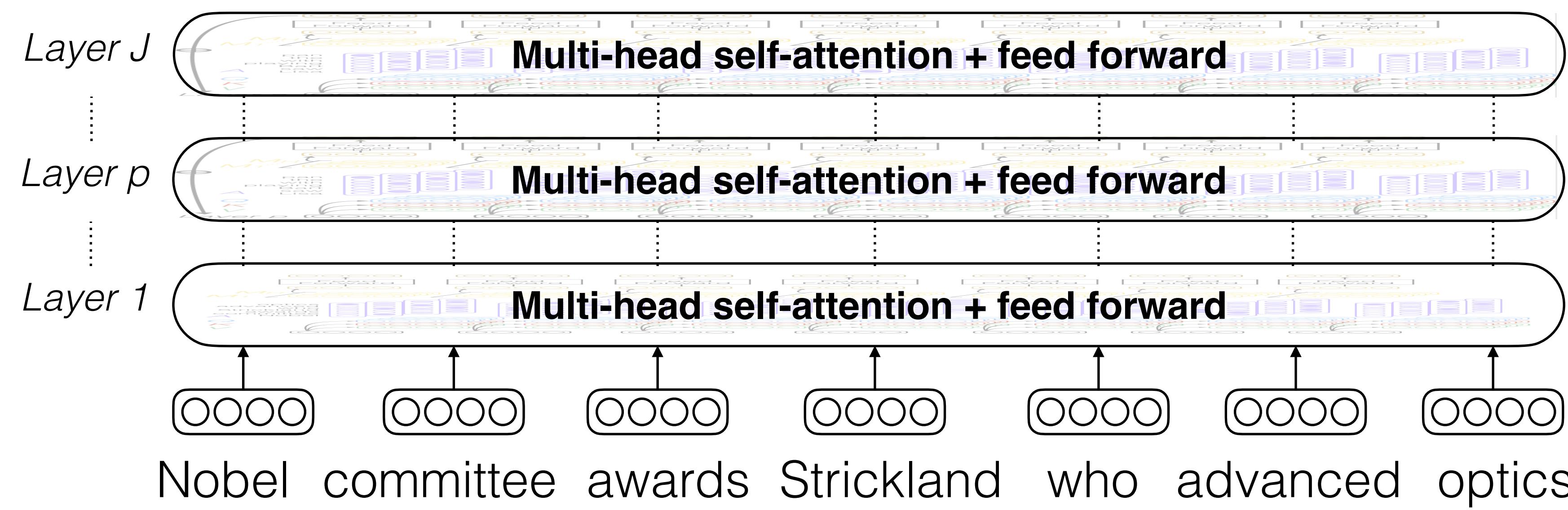
# Transformer architecture



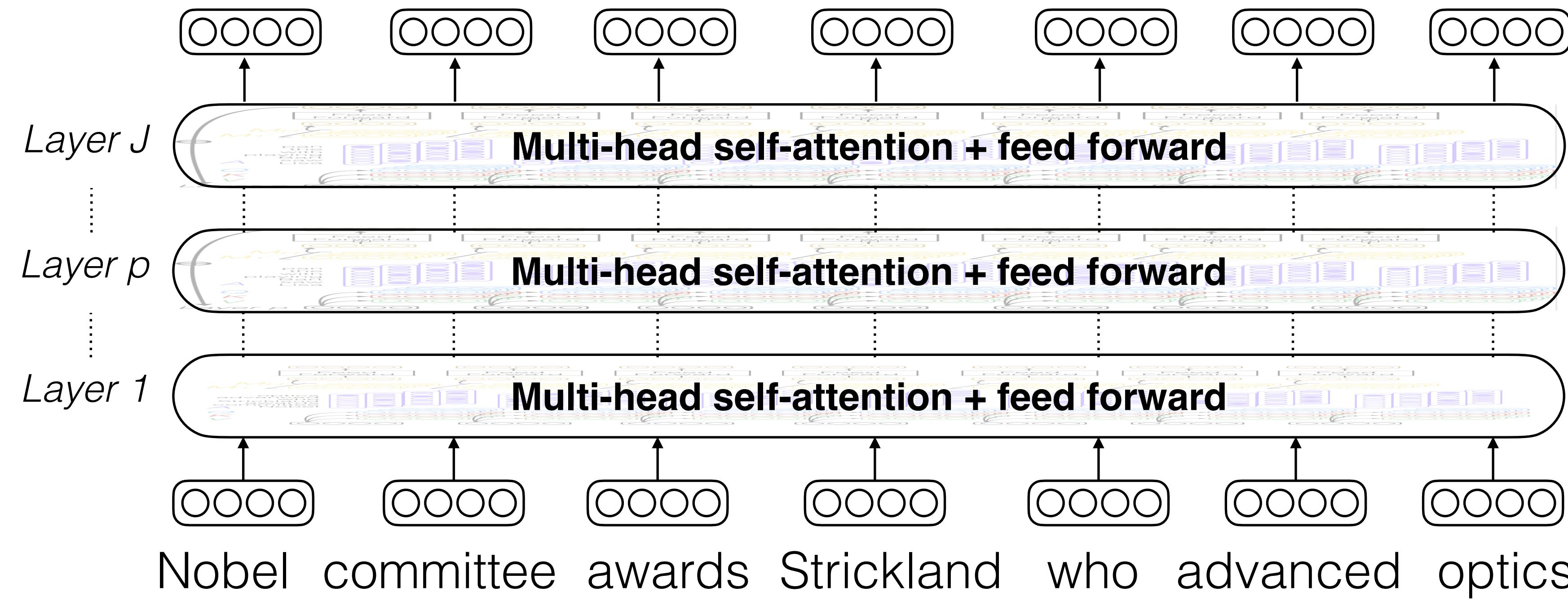
# Transformer architecture



# Transformer architecture



# Transformer architecture



# Modeling context

# Modeling context

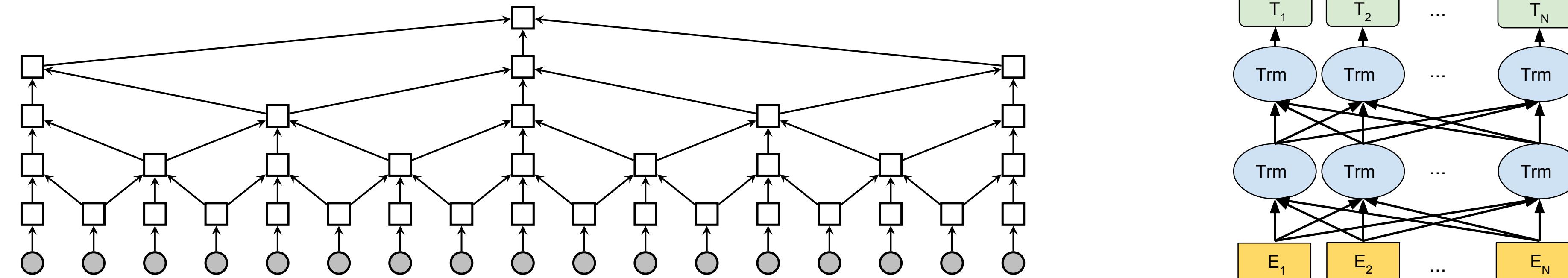
- Like RNNs, each token observes context from the *entire sequence* at each layer.

# Modeling context

- Like RNNs, each token observes context from the *entire* sequence at each layer.
- Like CNNs, representations at each layer can be computed *entirely in parallel*.

# Modeling context

- Like RNNs, each token observes context from the *entire* sequence at each layer.
- Like CNNs, representations at each layer can be computed *entirely in parallel*.
  - Dilated CNNs obtain wide (enough) context and are faster than Transformers, but at the cost of expressivity: impose a regular pattern of connectivity between tokens, whereas in Transformer connectivity is learned / each token always has access to all other tokens.



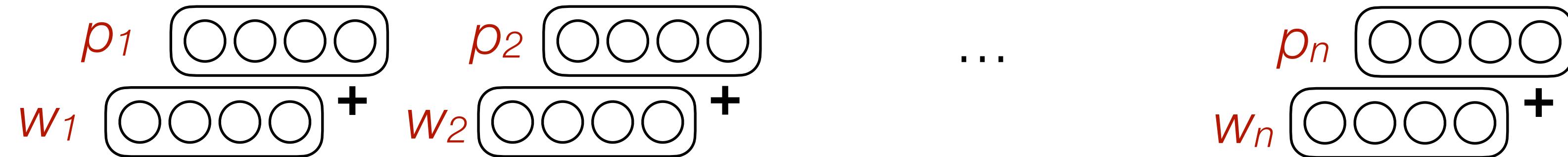
# Positional encoding

# Positional encoding

- **Positional embeddings** encode a token's position in the sentence

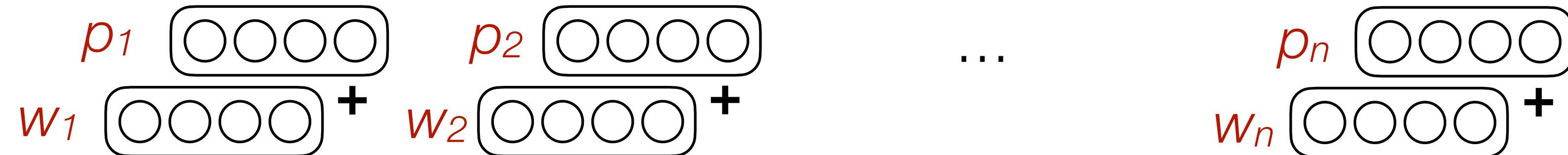
# Positional encoding

- **Positional embeddings** encode a token's position in the sentence



# Positional encoding

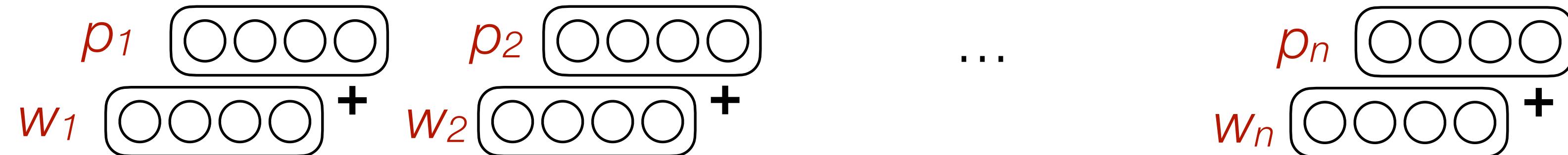
- **Positional embeddings** encode a token's position in the sentence



- **Sinusoidal position encodings** generalize to any sequence length.

# Positional encoding

- **Positional embeddings** encode a token's position in the sentence

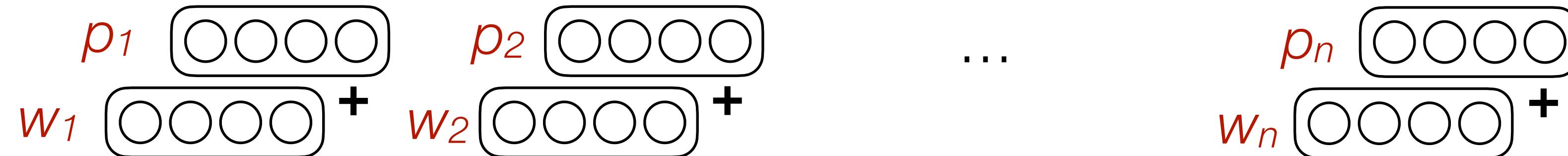


- **Sinusoidal position encodings** generalize to any sequence length.

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \text{with} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

# Positional encoding

- **Positional embeddings** encode a token's position in the sentence



- **Sinusoidal position encodings** generalize to any sequence length.

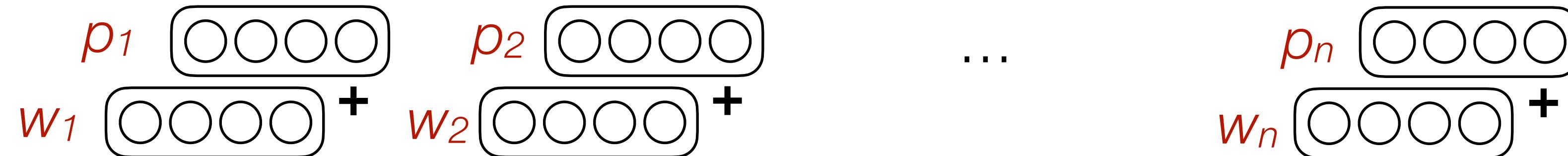
$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$

with  $\omega_k = \frac{1}{10000^{2k/d}}$

$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

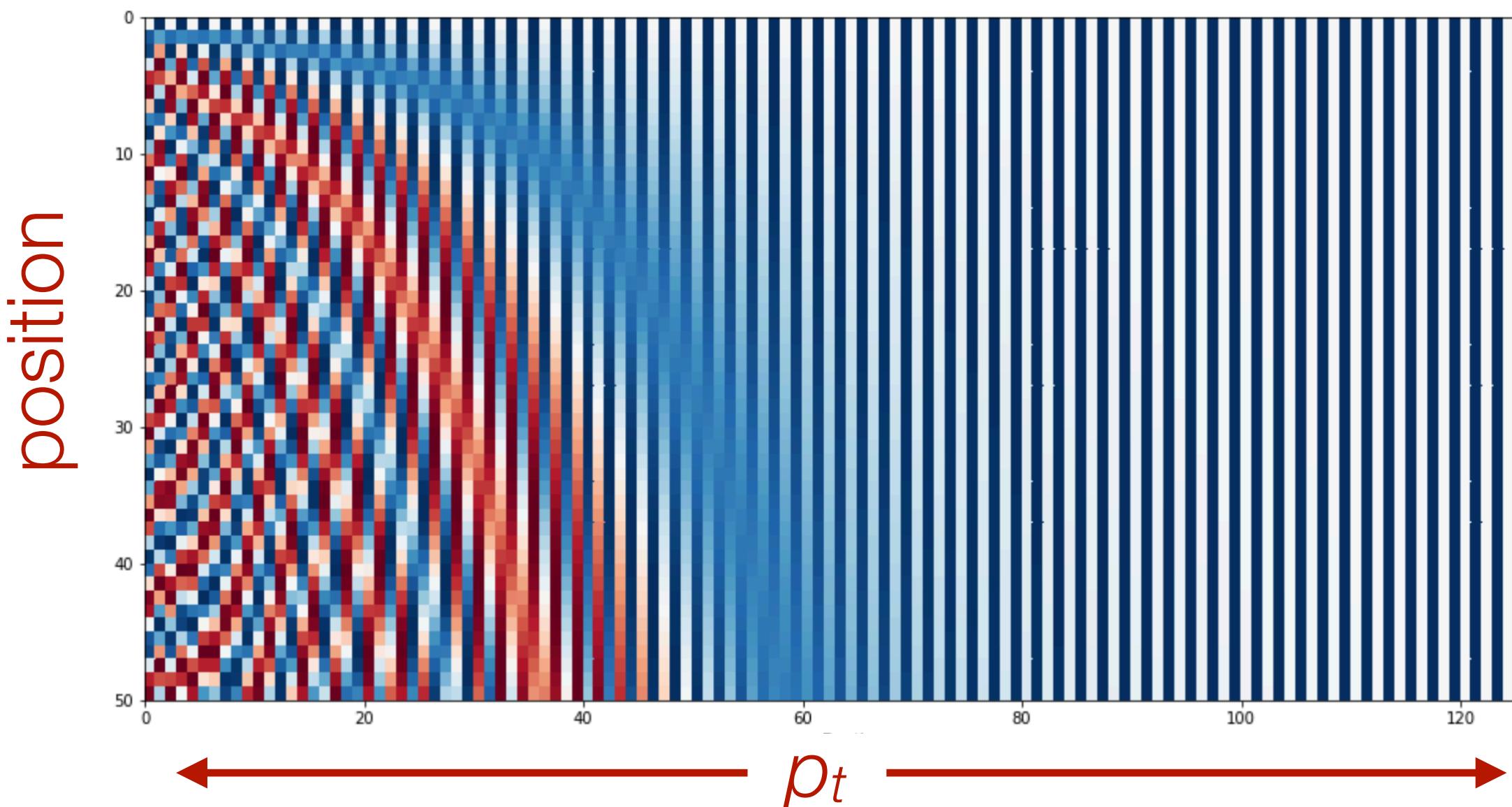
# Positional encoding

- **Positional embeddings** encode a token's position in the sentence

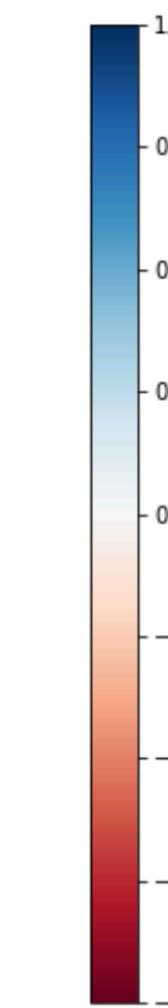


- **Sinusoidal position encodings** generalize to any sequence length.

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$



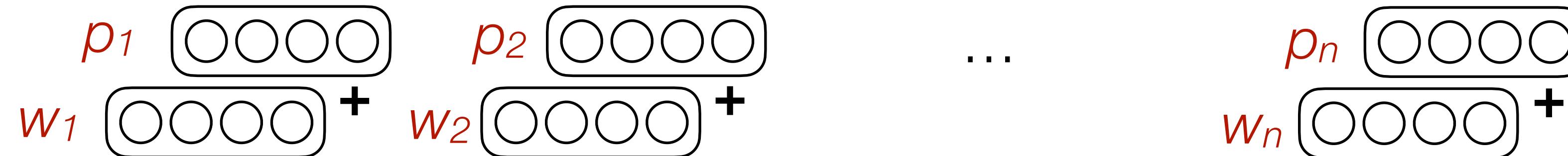
with  $\omega_k = \frac{1}{10000^{2k/d}}$



$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

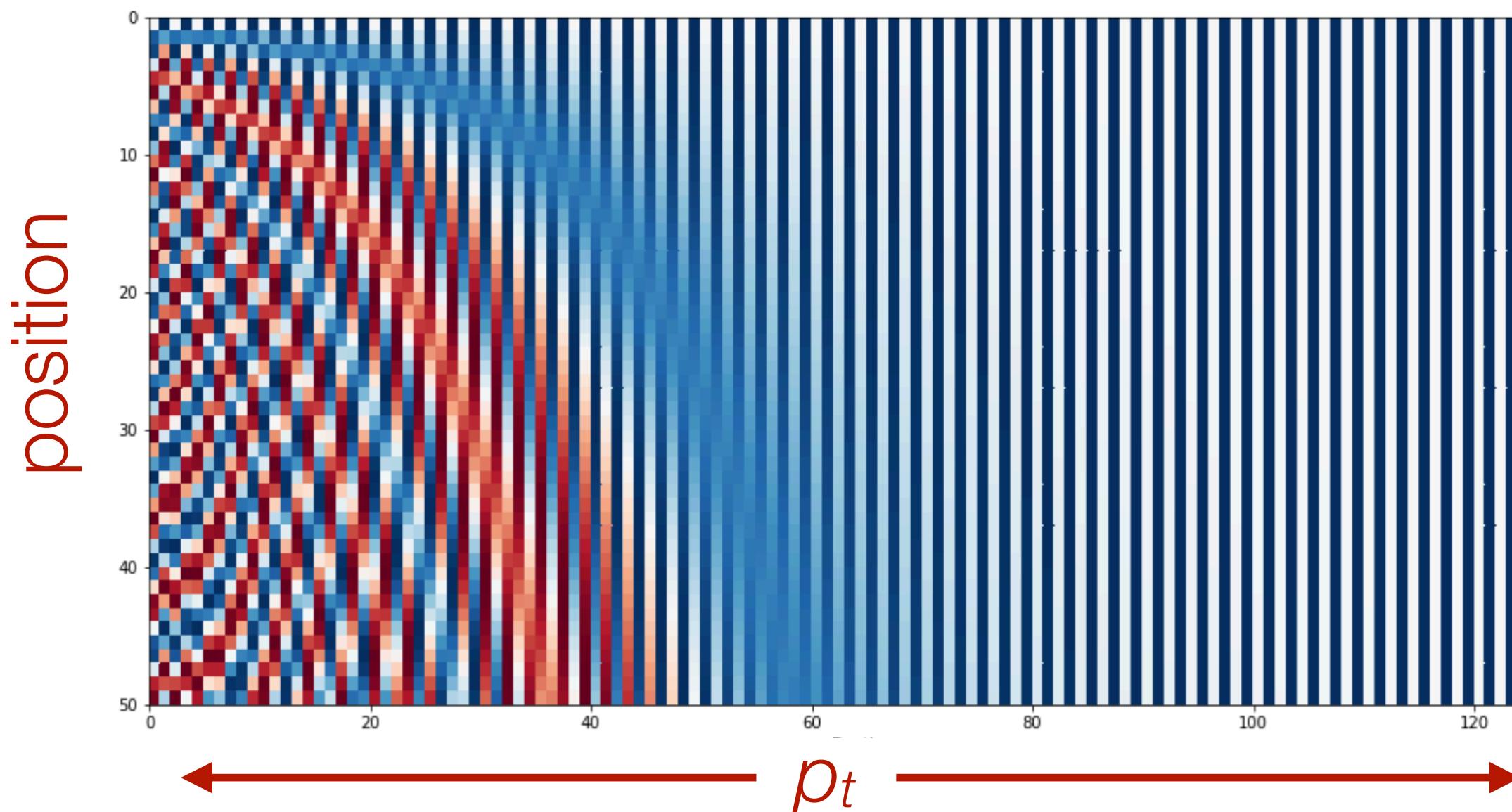
# Positional encoding

- **Positional embeddings** encode a token's position in the sentence

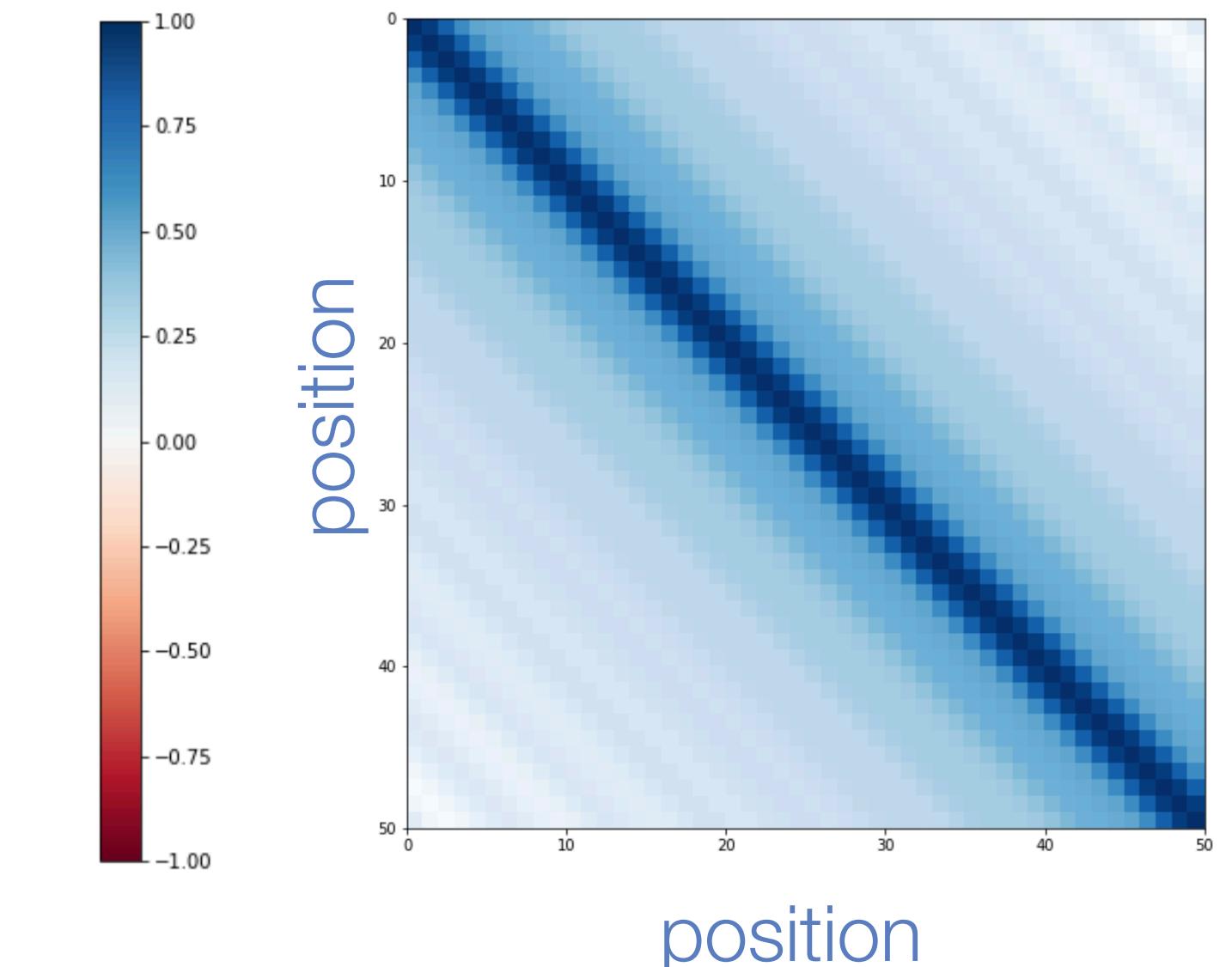


- **Sinusoidal position encodings** generalize to any sequence length.

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases}$$



with  $\omega_k = \frac{1}{10000^{2k/d}}$



$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$

# Sub-word representations

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units.**

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units.**

**BPE** [Gage, 1994; Sennrich et al., 2016]

---

## Algorithm 1 Byte-pair encoding

---

```
1: Input: set of strings  $D$ , target vocab size  $k$ 
2: procedure BPE( $D, k$ )
3:    $V \leftarrow$  all unique characters in  $D$ 
4:   (about 4,000 in English Wikipedia)
5:   while  $|V| < k$  do            $\triangleright$  Merge tokens
6:      $t_L, t_R \leftarrow$  Most frequent bigram in  $D$ 
7:      $t_{\text{NEW}} \leftarrow t_L + t_R$      $\triangleright$  Make new token
8:      $V \leftarrow V + [t_{\text{NEW}}]$ 
9:     Replace each occurrence of  $t_L, t_R$  in
10:         $D$  with  $t_{\text{NEW}}$ 
11:   end while
12:   return  $V$ 
13: end procedure
```

---

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units.**

**BPE** [Gage, 1994; Sennrich et al., 2016]

---

## Algorithm 1 Byte-pair encoding

---

```
1: Input: set of strings  $D$ , target vocab size  $k$ 
2: procedure BPE( $D, k$ )
3:    $V \leftarrow$  all unique characters in  $D$ 
4:   (about 4,000 in English Wikipedia)
5:   while  $|V| < k$  do            $\triangleright$  Merge tokens
6:      $t_L, t_R \leftarrow$  Most frequent bigram in  $D$ 
7:      $t_{\text{NEW}} \leftarrow t_L + t_R$      $\triangleright$  Make new token
8:      $V \leftarrow V + [t_{\text{NEW}}]$ 
9:     Replace each occurrence of  $t_L, t_R$  in
10:         $D$  with  $t_{\text{NEW}}$ 
11:   end while
12:   return  $V$ 
13: end procedure
```

## unigram LM [Kudo, 2018]

---

### Algorithm 2 Unigram language modeling

---

```
1: Input: set of strings  $D$ , target vocab size  $k$ 
2: procedure UNIGRAMLM( $D, k$ )
3:    $V \leftarrow$  all substrings occurring more than
4:     once in  $D$  (not crossing words)
5:   while  $|V| > k$  do            $\triangleright$  Prune tokens
6:     Fit unigram LM  $\theta$  to  $D$ 
7:     for  $t \in V$  do       $\triangleright$  Estimate token ‘loss’
8:        $L_t \leftarrow p_\theta(D) - p_{\theta'}(D)$ 
9:       where  $\theta'$  is the LM without token  $t$ 
10:    end for
11:    Remove  $\min(|V| - k, \lfloor \alpha |V| \rfloor)$  of the
12:      tokens  $t$  with highest  $L_t$  from  $V$ ,
13:      where  $\alpha \in [0, 1]$  is a hyperparameter
14:   end while
15:   Fit final unigram LM  $\theta$  to  $D$ 
16:   return  $V, \theta$ 
17: end procedure
```

---

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units**.

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units**.

<b>Original:</b>	furiously
<b>BPE:</b>	_fur   iously
<b>Unigram LM:</b>	_fur   ious   ly

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units**.

Original:	furiously
BPE:	_fur   iously
Unigram LM:	_fur   ious   ly

Original:	tricycles
BPE:	_t   ric   y   cles
Unigram LM:	_tri   cycle   s

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units**.

**Original:** furiously  
**BPE:** \_fur | iously  
**Unigram LM:** \_fur | ious | ly

**Original:** tricycles  
**BPE:** \_t | ric | y | cles  
**Unigram LM:** \_tri | cycle | s

**Original:** Completely preposterous suggestions  
**BPE:** \_Comple | t | ely | \_prep | ost | erous | \_suggest | ions  
**Unigram LM:** \_Complete | ly | \_pre | post | er | ous | \_suggestion | s

# Sub-word representations

- Word embeddings suffer from UNK/OOV, but character embeddings are slow.
- Solution: **sub-word units**.
- BPE and unigram LM both implemented in the [SentencePiece library](#)

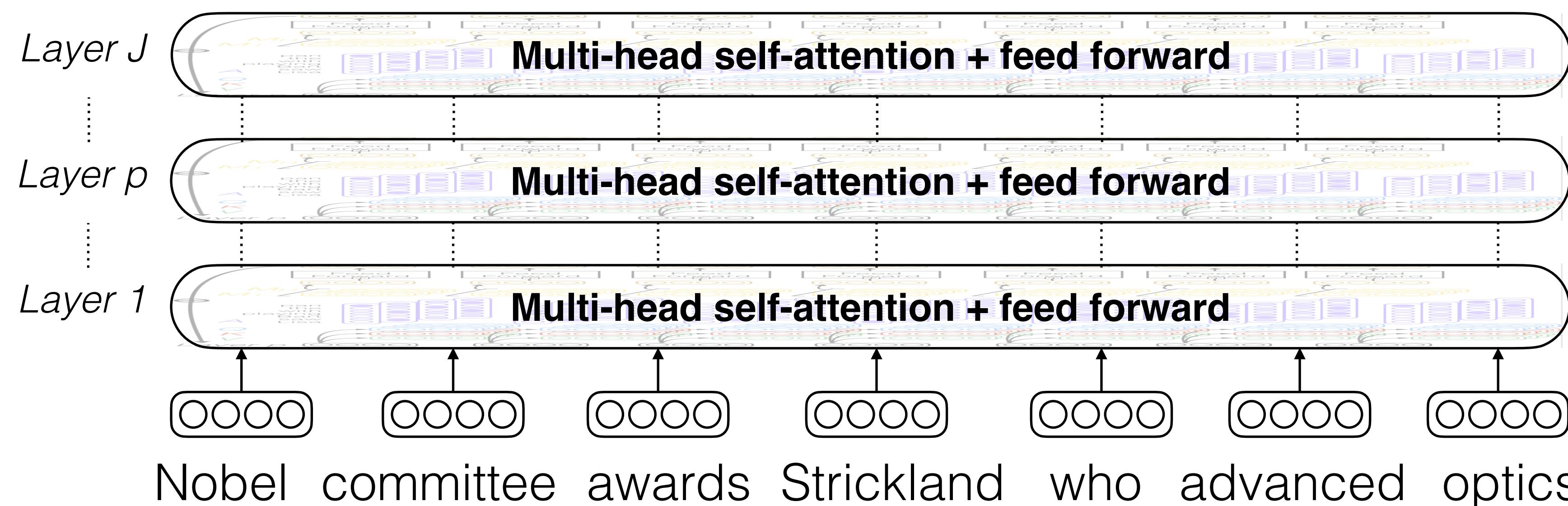
**Original:** furiously  
**BPE:** \_fur | iously  
**Unigram LM:** \_fur | ious | ly

**Original:** tricycles  
**BPE:** \_t | ric | y | cles  
**Unigram LM:** \_tri | cycle | s

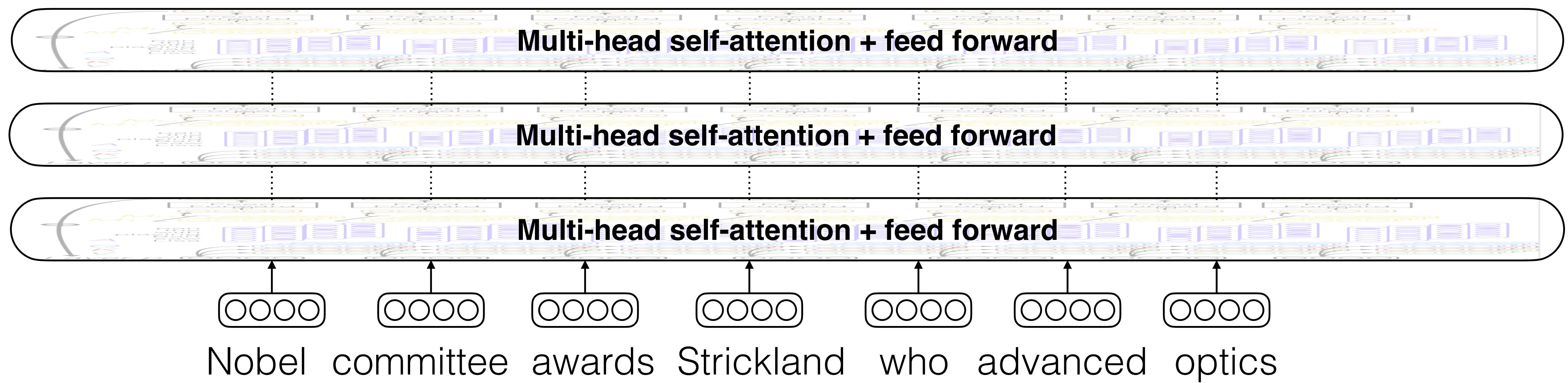
**Original:** Completely preposterous suggestions  
**BPE:** \_Comple | t | ely | \_prep | ost | erous | \_suggest | ions  
**Unigram LM:** \_Complete | ly | \_pre | post | er | ous | \_suggestion | s

# BERT

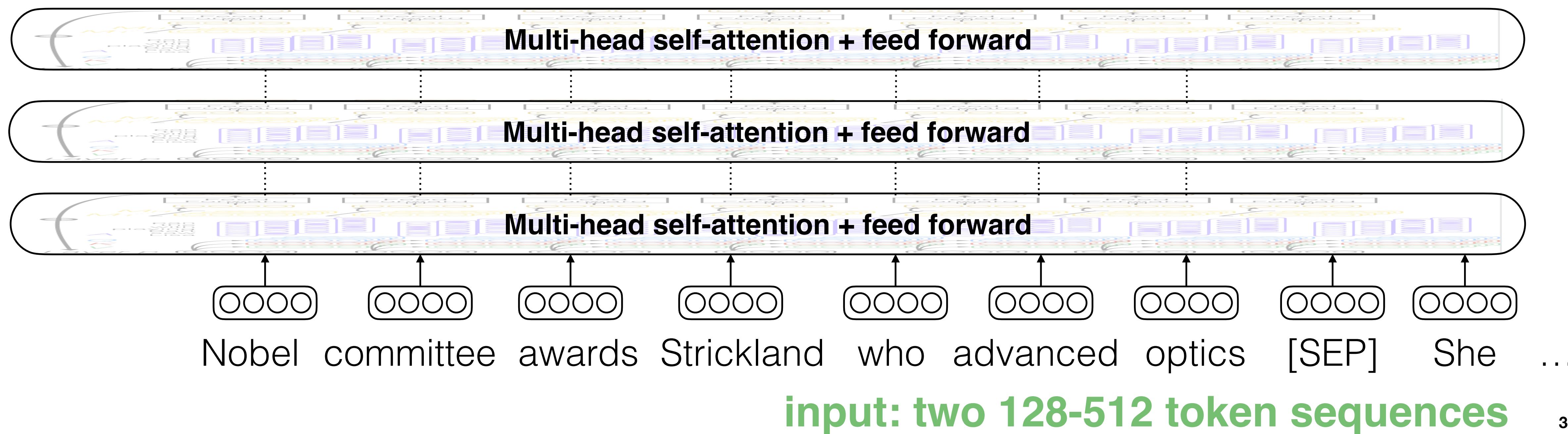
- Transformer language model trained with **non-autoregressive** masked language modeling (MLM) objective.



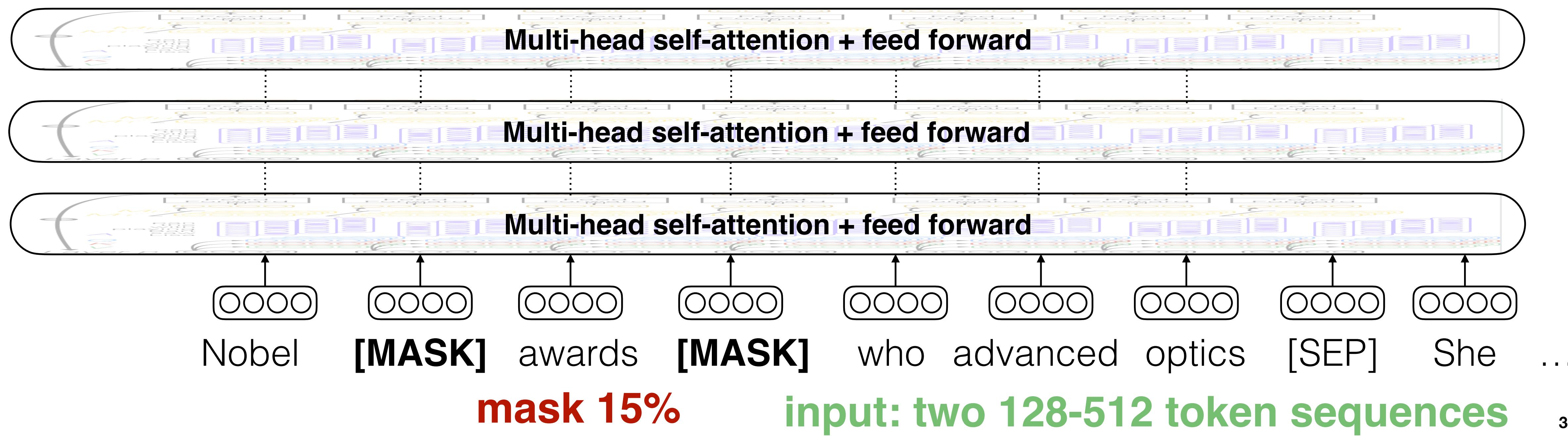
# BERT



# BERT

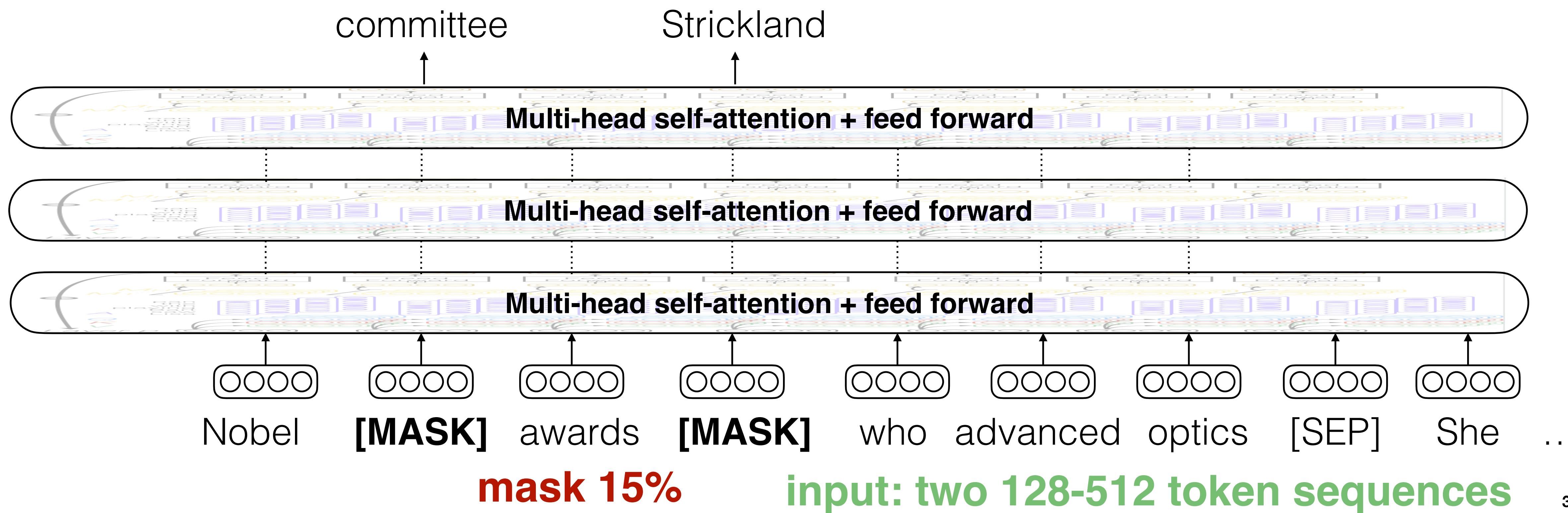


# BERT



# BERT

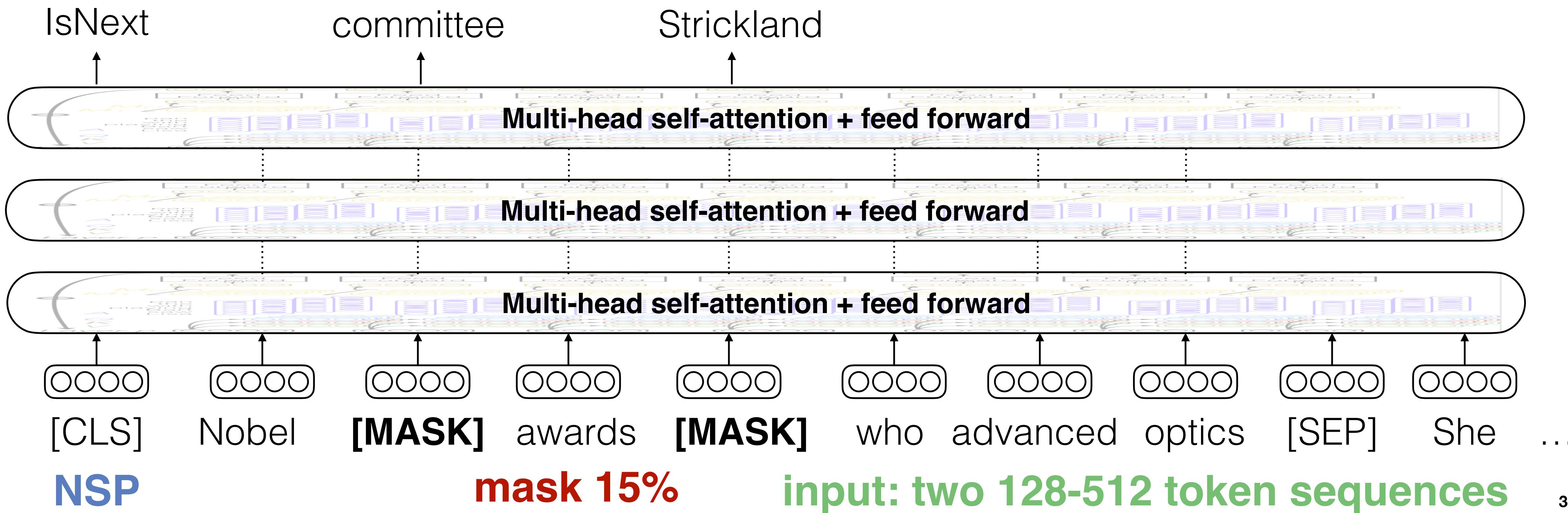
**predict masked tokens**



# BERT

Next sentence prediction:  
does the second sequence follow the first?

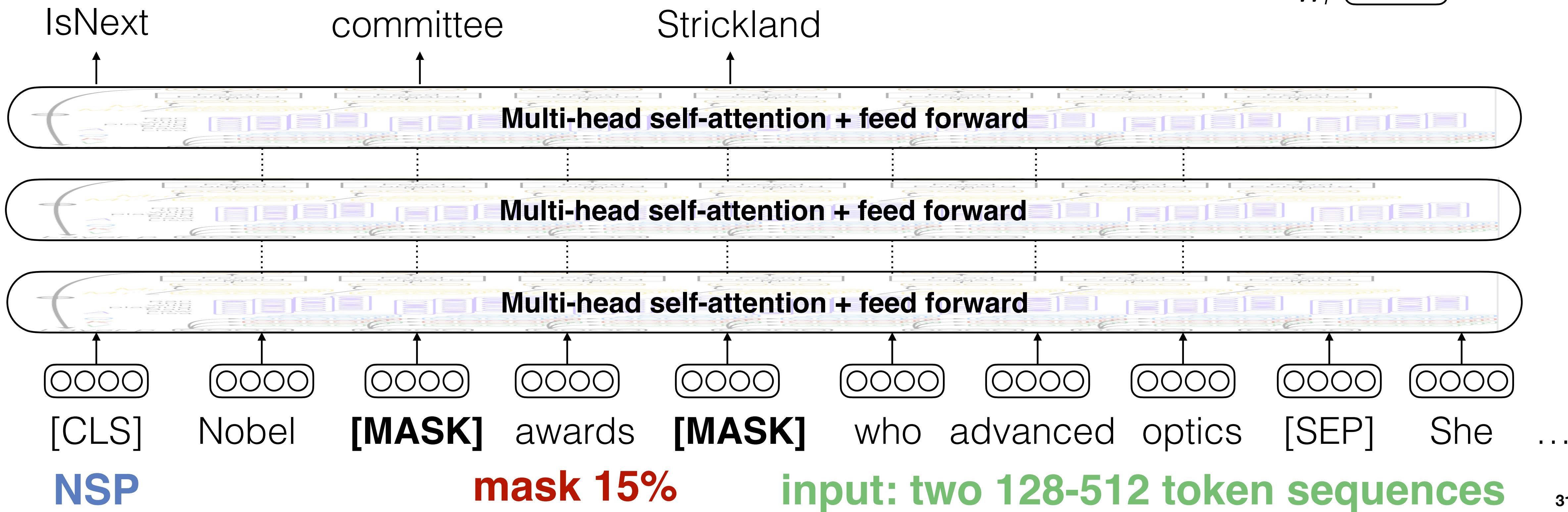
**predict masked tokens**



# BERT

Next sentence prediction:  
does the second sequence follow the first?

**predict masked tokens**

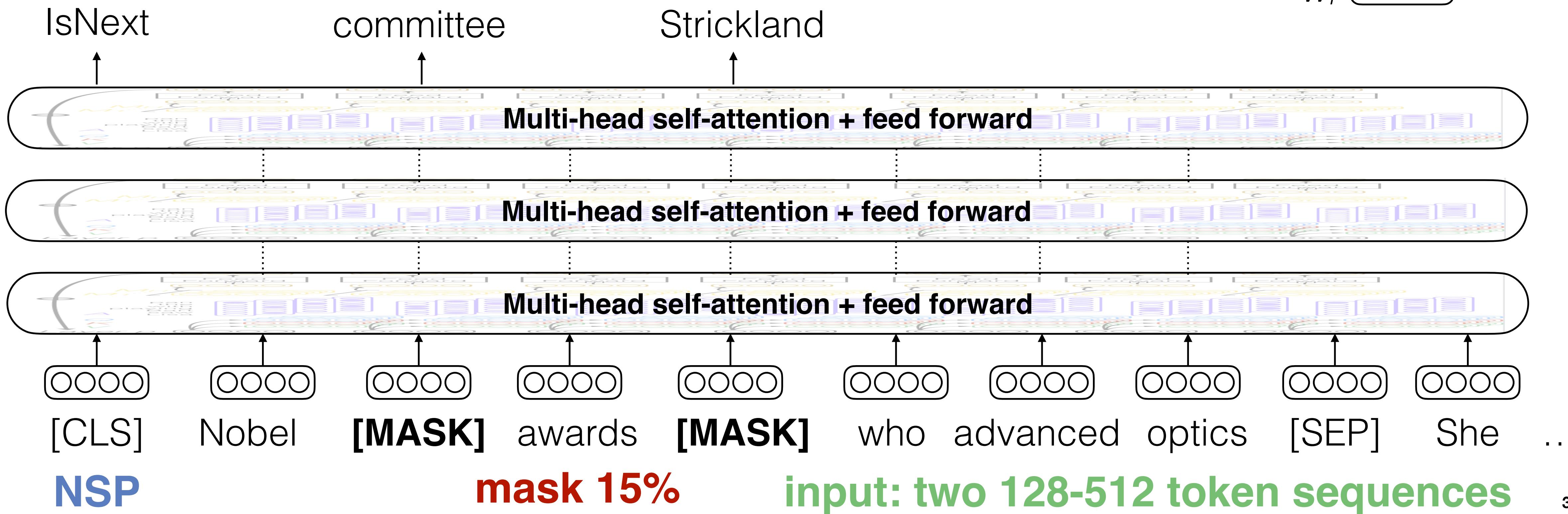


# BERT

$s_i$ : which sequence  
is this token from?

Next sentence prediction:  
does the second sequence follow the first?

**predict masked tokens**



$$x_i \quad \boxed{\text{oooo}} = p_i \quad \boxed{\text{oooo}} + w_i \quad \boxed{\text{oooo}}$$

# Evaluating pretrained LMs

# Evaluating pretrained LMs

- Intrinsic (perplexity) and extrinsic (fine-tuning on downstream task) evaluations.

# Evaluating pretrained LMs

- Intrinsic (perplexity) and extrinsic (fine-tuning on downstream task) evaluations.
- Standard extrinsic benchmarks for LM pretraining: **(Super)GLUE** and **SQuAD**

# Evaluating pretrained LMs

- Intrinsic (perplexity) and extrinsic (fine-tuning on downstream task) evaluations.
- Standard extrinsic benchmarks for LM pretraining: **(Super)GLUE** and **SQuAD**

**GLUE:** <https://gluebenchmark.com/>

## General Language Understanding Evaluation

Corpus	Train	Test	Task	Metrics	Domain
Single-Sentence Tasks					
CoLA	8.5k	<b>1k</b>	acceptability	Matthews corr.	misc.
SST-2	67k	1.8k	sentiment	acc.	movie reviews
Similarity and Paraphrase Tasks					
MRPC	3.7k	1.7k	paraphrase	acc./F1	news
STS-B	7k	1.4k	sentence similarity	Pearson/Spearman corr.	misc.
QQP	364k	<b>391k</b>	paraphrase	acc./F1	social QA questions
Inference Tasks					
MNLI	393k	<b>20k</b>	NLI	matched acc./mismatched acc.	misc.
QNLI	105k	5.4k	QA/NLI	acc.	Wikipedia
RTE	2.5k	3k	NLI	acc.	news, Wikipedia
WNLI	634	<b>146</b>	coreference/NLI	acc.	fiction books

# Evaluating pretrained LMs

- Intrinsic (perplexity) and extrinsic (fine-tuning on downstream task) evaluations.
- Standard extrinsic benchmarks for LM pretraining: **(Super)GLUE** and **SQuAD**

**GLUE:** <https://gluebenchmark.com/>

## General Language Understanding Evaluation

Corpus	Train	Test	Task	Metrics	Domain
Single-Sentence Tasks					
CoLA	8.5k	<b>1k</b>	acceptability	Matthews corr.	misc.
SST-2	67k	1.8k	sentiment	acc.	movie reviews
Similarity and Paraphrase Tasks					
MRPC	3.7k	1.7k	paraphrase	acc./F1	news
STS-B	7k	1.4k	sentence similarity	Pearson/Spearman corr.	misc.
QQP	364k	<b>391k</b>	paraphrase	acc./F1	social QA questions
Inference Tasks					
MNLI	393k	<b>20k</b>	NLI	matched acc./mismatched acc.	misc.
QNLI	105k	5.4k	QA/NLI	acc.	Wikipedia
RTE	2.5k	3k	NLI	acc.	news, Wikipedia
WNLI	634	<b>146</b>	coreference/NLI	acc.	fiction books

**SQuAD:** <https://rajpurkar.github.io/SQuAD-explorer/>

## Stanford Question Answering Dataset

The Yuan dynasty (Chinese: 元朝; pinyin: Yuán Cháo), officially the Great Yuan (Chinese: 大元; pinyin: Dà Yuán; Mongolian: Yehe Yuan Ulus[a]), was the empire or ruling dynasty of China established by Kublai Khan, leader of the Mongolian Borjigin clan. Although the Mongols had ruled territories including today's North China for decades, it was not until 1271 that Kublai Khan officially proclaimed the dynasty in the traditional Chinese style. His realm was, by this point, isolated from the other khanates and controlled most of present-day China and its surrounding areas, including modern Mongolia and Korea. It was the first foreign dynasty to rule all of China and lasted until 1368, after which its Genghisid rulers returned to their Mongolian homeland and continued to rule the Northern Yuan dynasty. Some of the Mongolian Emperors of the Yuan mastered the Chinese language, while others only used their native language (i.e. Mongolian) and the 'Phags-pa script.

What is the Chinese name for the Yuan dynasty?

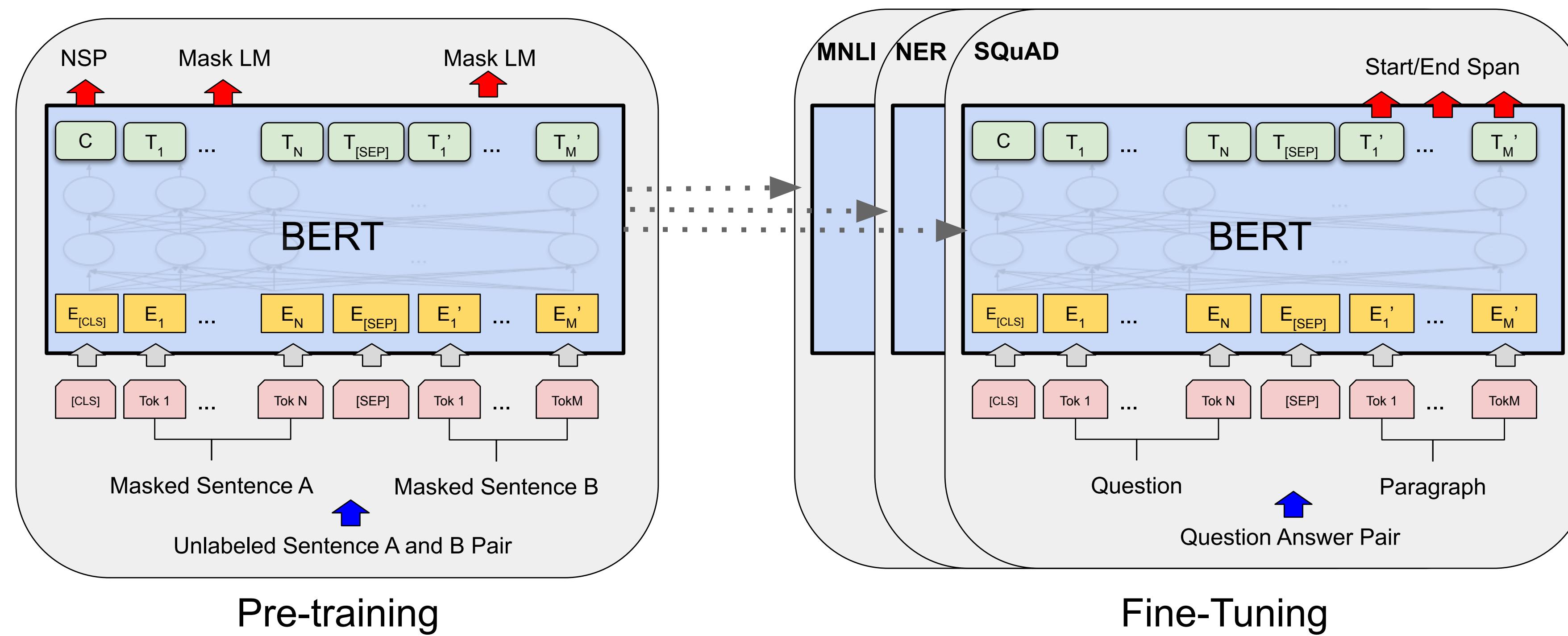
Ground Truth Answers: Yuán Cháo Yuán Cháo 元朝

When did Khan formally declare the Yuan dynasty?

Ground Truth Answers: 1271 1271 1271

# Fine-tuning pretrained LMs

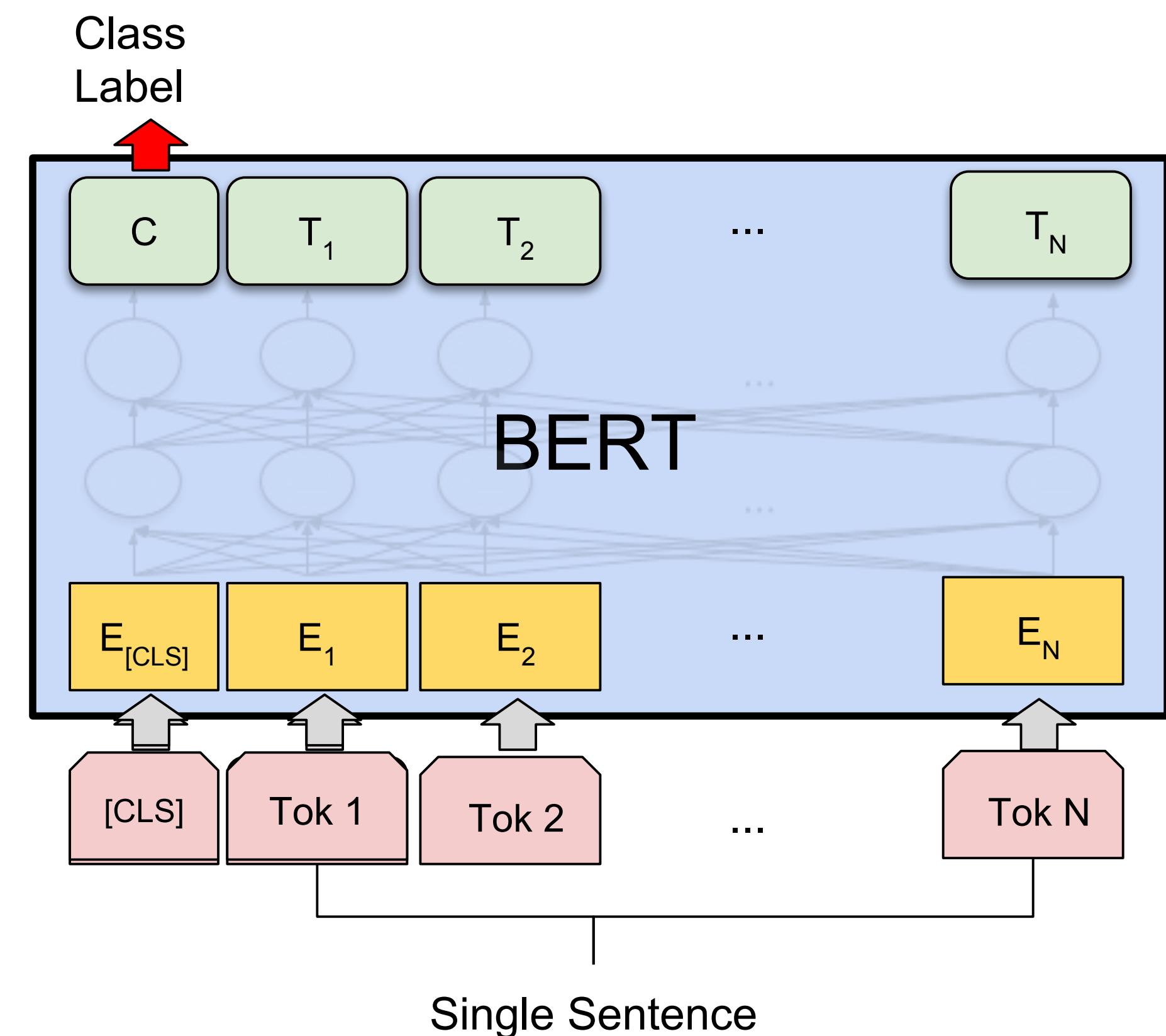
- Unlike in ELMo, all model parameters are fine-tuned (receive gradient updates w.r.t. supervised downstream task.)



# Fine-tuning pretrained LMs

## ■ Sentence classification:

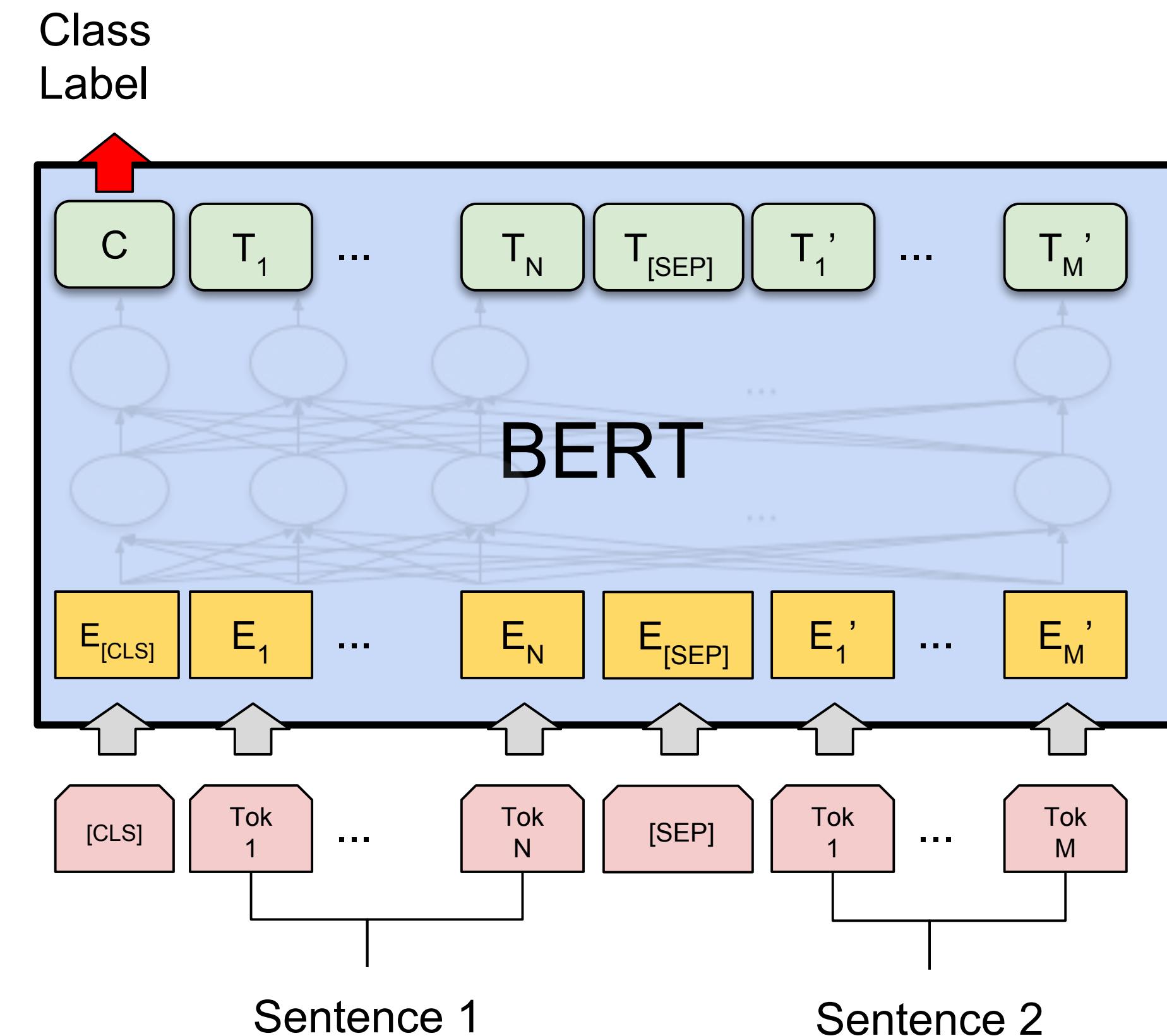
- Provide single sentence as input
- Use final representation at CLS position as features for a linear classifier (e.g. logistic regression).



# Fine-tuning pretrained LMs

## ■ Sentence pair classification:

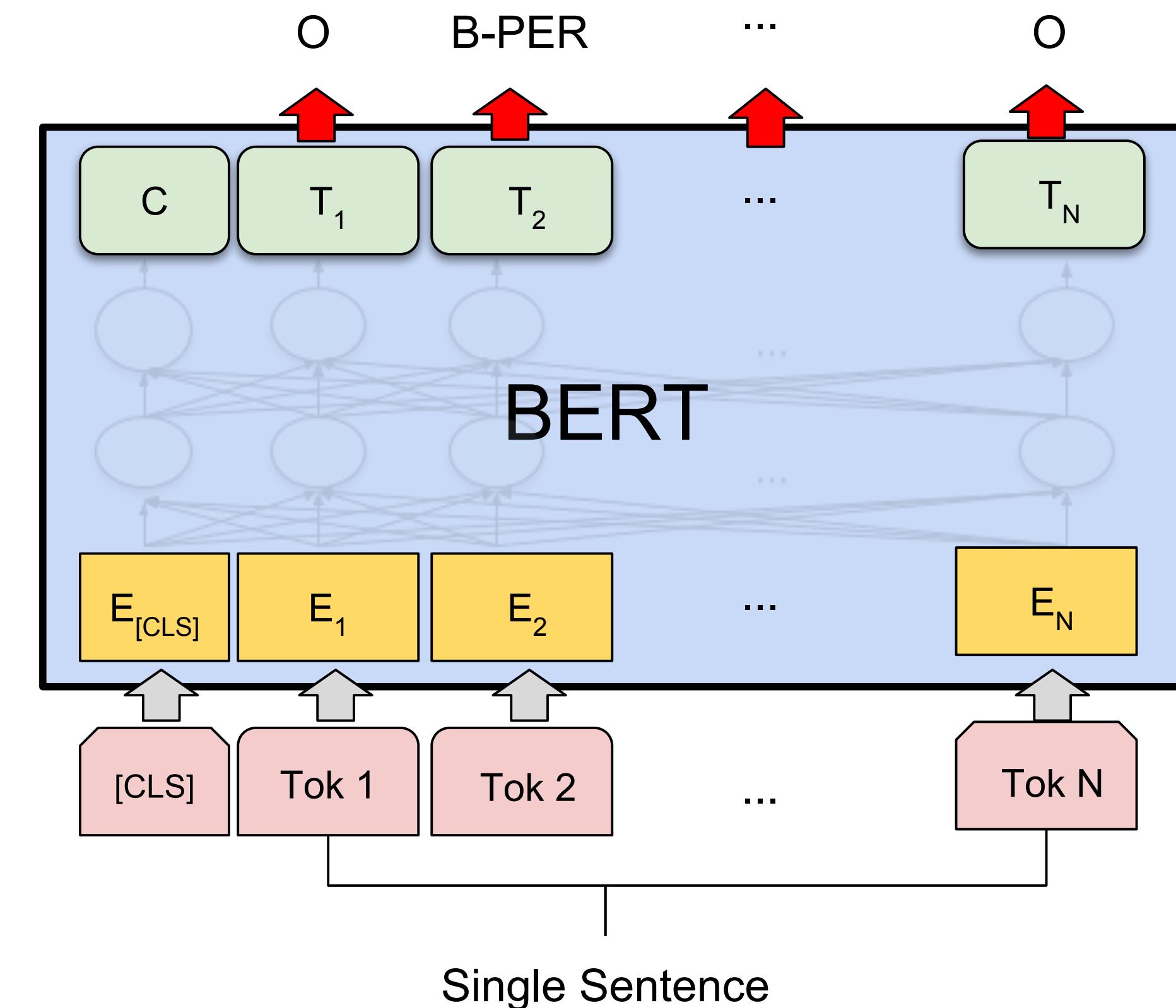
- Provide both sentences as input, separated by SEP token.
- Input embeddings indicate sentence identity.
- Use final representation at CLS position as features for a linear classifier (e.g. logistic regression).



# Fine-tuning pretrained LMs

## ■ Sequence labeling:

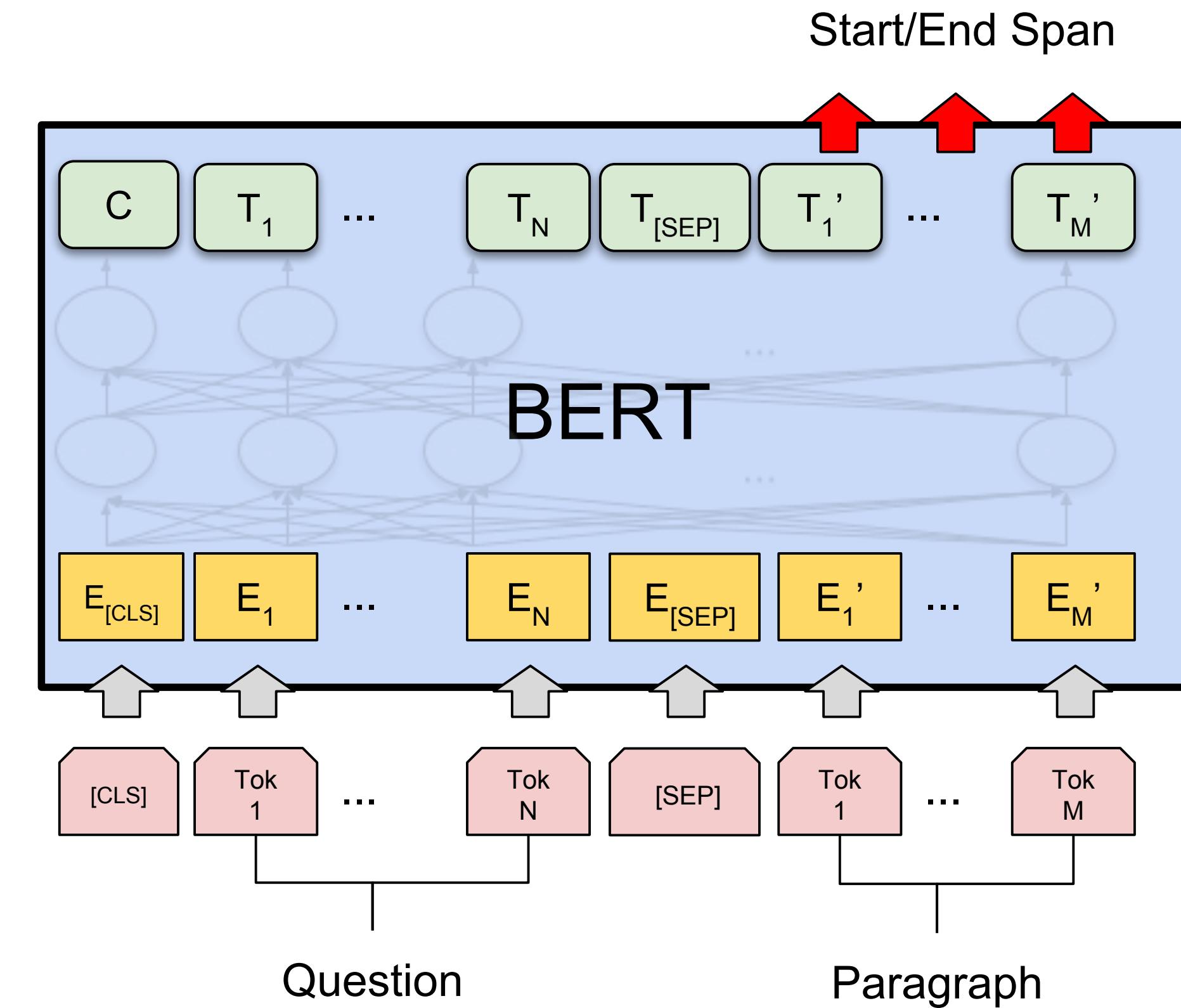
- Use final representations at each timestep as features for a linear classifier (e.g. logistic regression).
- Need to deal with sub-tokens.



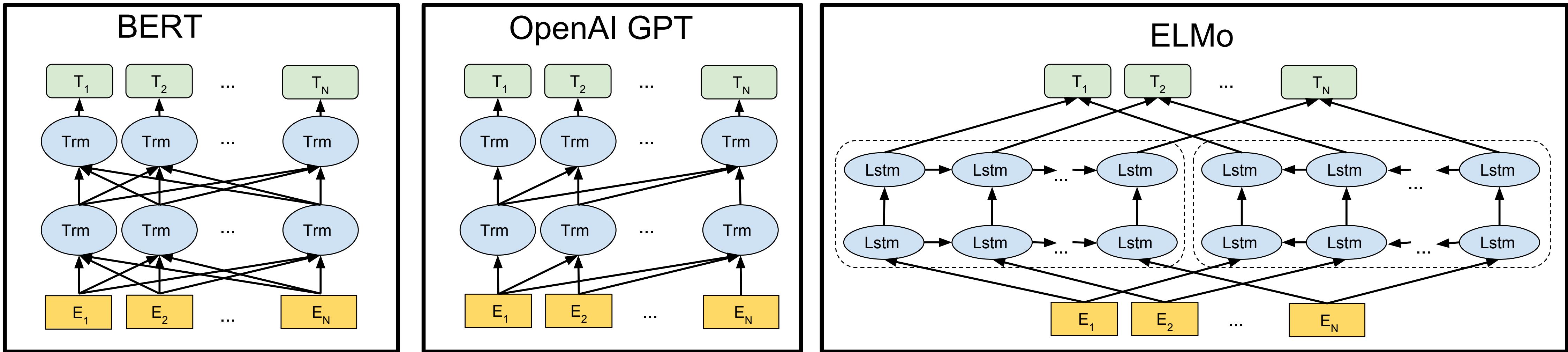
# Fine-tuning pretrained LMs

## ■ Question answering:

- Use final representations at each timestep as features to score every token as start or end of a span. Take the highest-scoring pair.

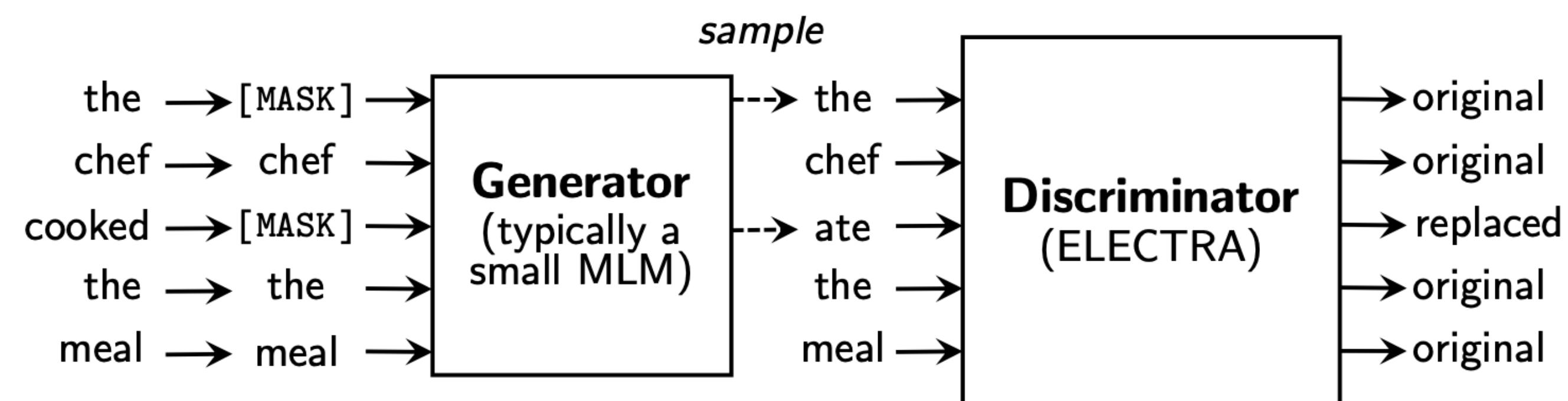


# Other popular Transformer LMs



- GPT/XLNet
- RoBERTa
- ELECTRA

ELECTRA [Clark et al. 2020]



# How does performance scale with data?



**Colin Raffel**  
@colinraffel

Hot take: The most surprising thing about BERT isn't how well it worked when it was proposed, but how much better it would have worked if they had just pre-trained for longer on a more diverse dataset.

2:28 PM · Apr 16, 2020 · [Twitter Web App](#)

**21** Retweets   **204** Likes



# How does performance scale with data?

RoBERTa  
([Liu et al. 2019](#))

Training data				
Dataset	16GB	160GB	$\Delta$	err reduction
SQuAD	87.3	89.4	2.1	16.5%
MNLI-m	89.0	90.2	1.2	10.9%
SST-2	95.3	96.4	0.9	19.1%

XLNet  
([Yang et al. 2019](#))

~3B	~3B	~3B	~3B	~3B
Dataset	13GB	126GB	$\Delta$	err reduction
SQuAD	87.8	88.8	1.0	8.2%
MNLI-m	88.4	89.8	1.4	12.1%
SST-2	94.4	95.6	1.2	21.4%

T5  
([Raffel et al. 2020](#))

~8M	~8M	~8M	~8M	~8M
Dataset	$2^{23}$	$2^{25}$	$2^{27}$	$2^{29}$
SQuAD	70.92	76.27	79.78	80.97
GLUE	76.34	79.55	82.62	82.87
SGLUE	59.29	64.76	69.97	72.03
				$2^{35}$
				80.88
				83.28
				71.36

# Data source effects performance

Dataset	Size	GLUE	CNNDM	SQuAD	SGLUE	English		
						news	wikipedia	fiction*
★ C4	745GB	83.28	<b>19.24</b>	80.88	71.36	<b>26.98</b>	<b>39.82</b>	<b>27.65</b>
C4, unfiltered	6.1TB	81.46	19.14	78.78	68.04	26.55	39.34	27.21
RealNews-like	35GB	<b>83.83</b>	<b>19.23</b>	80.39	72.38	<b>26.75</b>	<b>39.90</b>	<b>27.48</b>
WebText-like	17GB	<b>84.03</b>	<b>19.31</b>	<b>81.42</b>	71.40	<b>26.80</b>	<b>39.74</b>	<b>27.59</b>
Wikipedia	16GB	81.85	<b>19.31</b>	81.29	68.01	<b>26.94</b>	39.69	<b>27.67</b>
Wikipedia + TBC	20GB	83.65	<b>19.28</b>	<b>82.08</b>	<b>73.24</b>	<b>26.77</b>	39.63	<b>27.57</b>

**Table 8:** Performance resulting from pre-training on different datasets. The first four variants are based on our new C4 dataset.

T5 [Raffel et al. 2020]

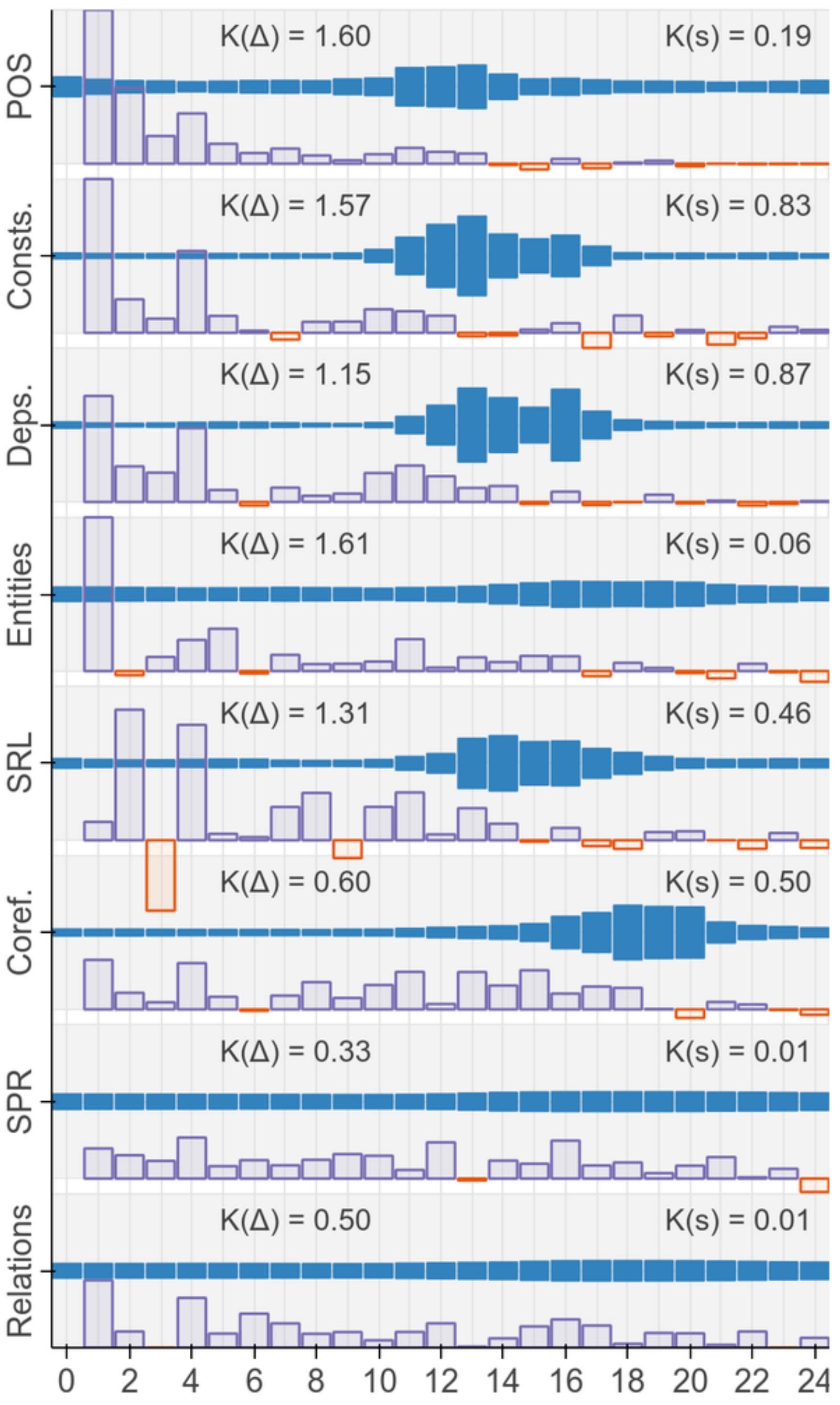
# Data source effects performance

Field	Task	Dataset	SOTA	BERT-Base		SciBERT	
				Frozen	Finetune	Frozen	Finetune
Bio	NER	BC5CDR (Li et al., 2016)	88.85 <sup>7</sup>	85.08	86.72	88.73	<b>90.01</b>
		JNLPBA (Collier and Kim, 2004)	<b>78.58</b>	74.05	76.09	75.77	77.28
		NCBI-disease (Dogan et al., 2014)	<b>89.36</b>	84.06	86.88	86.39	88.57
	PICO	EBM-NLP (Nye et al., 2018)	66.30	61.44	71.53	68.30	<b>72.28</b>
	DEP	GENIA (Kim et al., 2003) - LAS	<b>91.92</b>	90.22	90.33	90.36	90.43
		GENIA (Kim et al., 2003) - UAS	<b>92.84</b>	91.84	91.89	92.00	91.99
CS	REL	ChemProt (Kringelum et al., 2016)	76.68	68.21	79.14	75.03	<b>83.64</b>
	NER	SciERC (Luan et al., 2018)	64.20	63.58	65.24	65.77	<b>67.57</b>
	REL	SciERC (Luan et al., 2018)	n/a	72.74	78.71	75.25	<b>79.97</b>
	CLS	ACL-ARC (Jurgens et al., 2018)	67.9	62.04	63.91	60.74	<b>70.98</b>
Multi	CLS	Paper Field	n/a	63.64	65.37	64.38	<b>65.71</b>
		SciCite (Cohan et al., 2019)	84.0	84.31	84.85	<b>85.42</b>	<b>85.49</b>
Average				73.58	77.16	76.01	79.27

SciBERT [Beltagy et al. 2020]

pre-training LM on in-domain data improves end-task performance

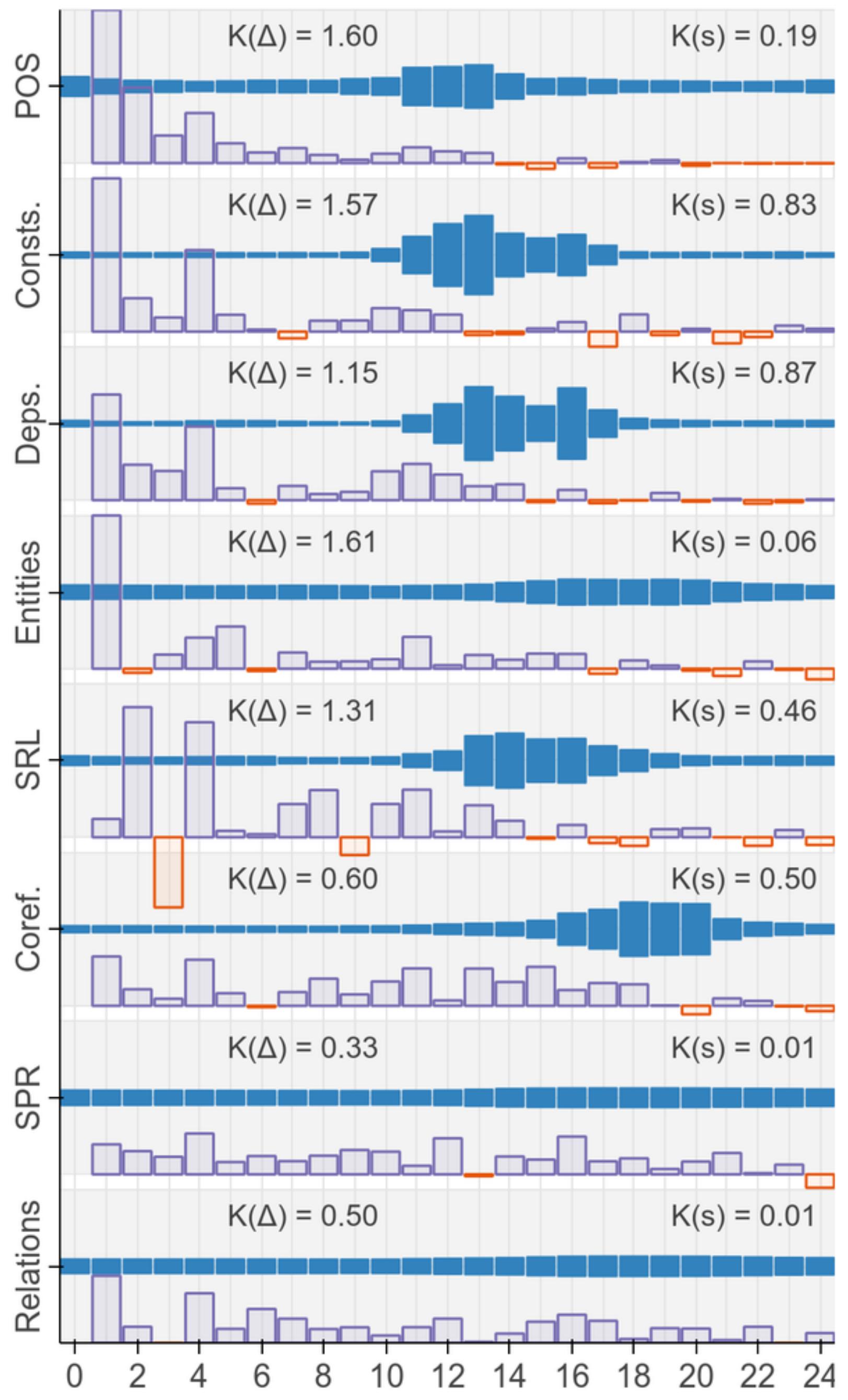
# What does BERT know?



# What does BERT know?

Context	BERT <sub>LARGE</sub> predictions
<i>A robin is a _____</i>	<i>bird, robin, person, hunter, pigeon</i>
<i>A daisy is a _____</i>	<i>daisy, rose, flower, berry, tree</i>
<i>A hammer is a _____</i>	<i>hammer, tool, weapon, nail, device</i>
<i>A hammer is an _____</i>	<i>object, instrument, axe, implement, explosive</i>
<i>A robin is not a _____</i>	<i>robin, bird, penguin, man, fly</i>
<i>A daisy is not a _____</i>	<i>daisy, rose, flower, lily, cherry</i>
<i>A hammer is not a _____</i>	<i>hammer, weapon, tool, gun, rock</i>
<i>A hammer is not an _____</i>	<i>object, instrument, axe, animal, artifact</i>

Table 13: NEG-88-SIMP predictions by BERT<sub>LARGE</sub>



# Announcements

- Share your writeup from P1 (bake-off design and paper summary) on Piazza!
- Friday's recitation will be an overview of P2.