# Attention mechanisms

## CS 585, Fall 2019

Introduction to Natural Language Processing

## Mohit Iyyer

College of Information and Computer Sciences
University of Massachusetts Amherst
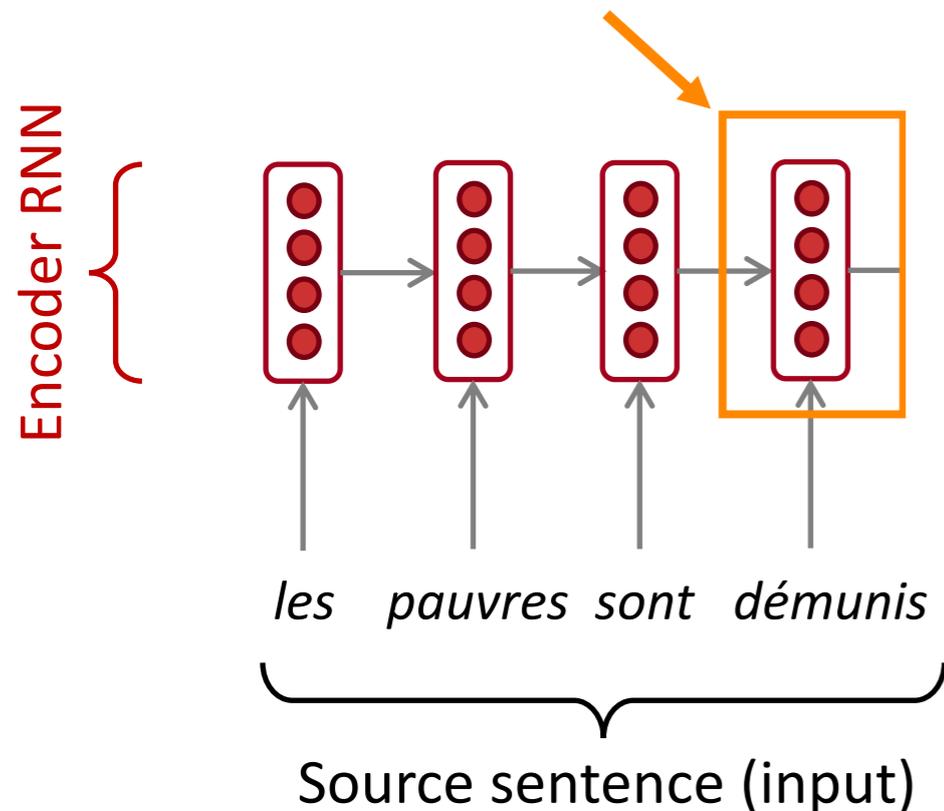
*some slides from Richard Socher*

# stuff from last time

- Colab issues :(

- HW1 time mixup, won't count anyone who submitted before 11:59pm as late

- Important dates:
  - Proposal due: Oct 4 (this Friday!!!)
  - Milestone 1 due: Oct 24
  - Midterm date: Oct 31
  - Milestone 2 due: Nov 21
  - HW 3 due: ???
  - Poster presentations: Dec 10/12
  - Final report due: Dec 19

- Can we spend a lot of time on attention? maybe

- Final exam instead of final project? NO!

# Neural Machine Translation (NMT)

The sequence-to-sequence model

Encoding of the source sentence.
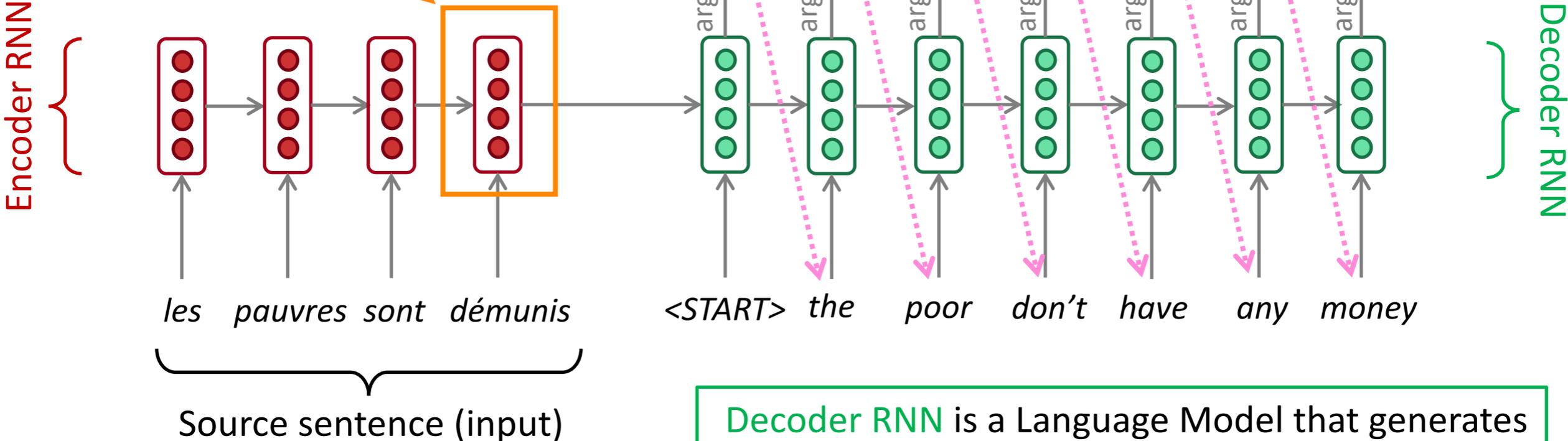Provides initial hidden state
for Decoder RNN.



Encoder RNN

*les    pauvres   sont    démunis*

Source sentence (input)

Encoder RNN produces
an encoding of the
source sentence.

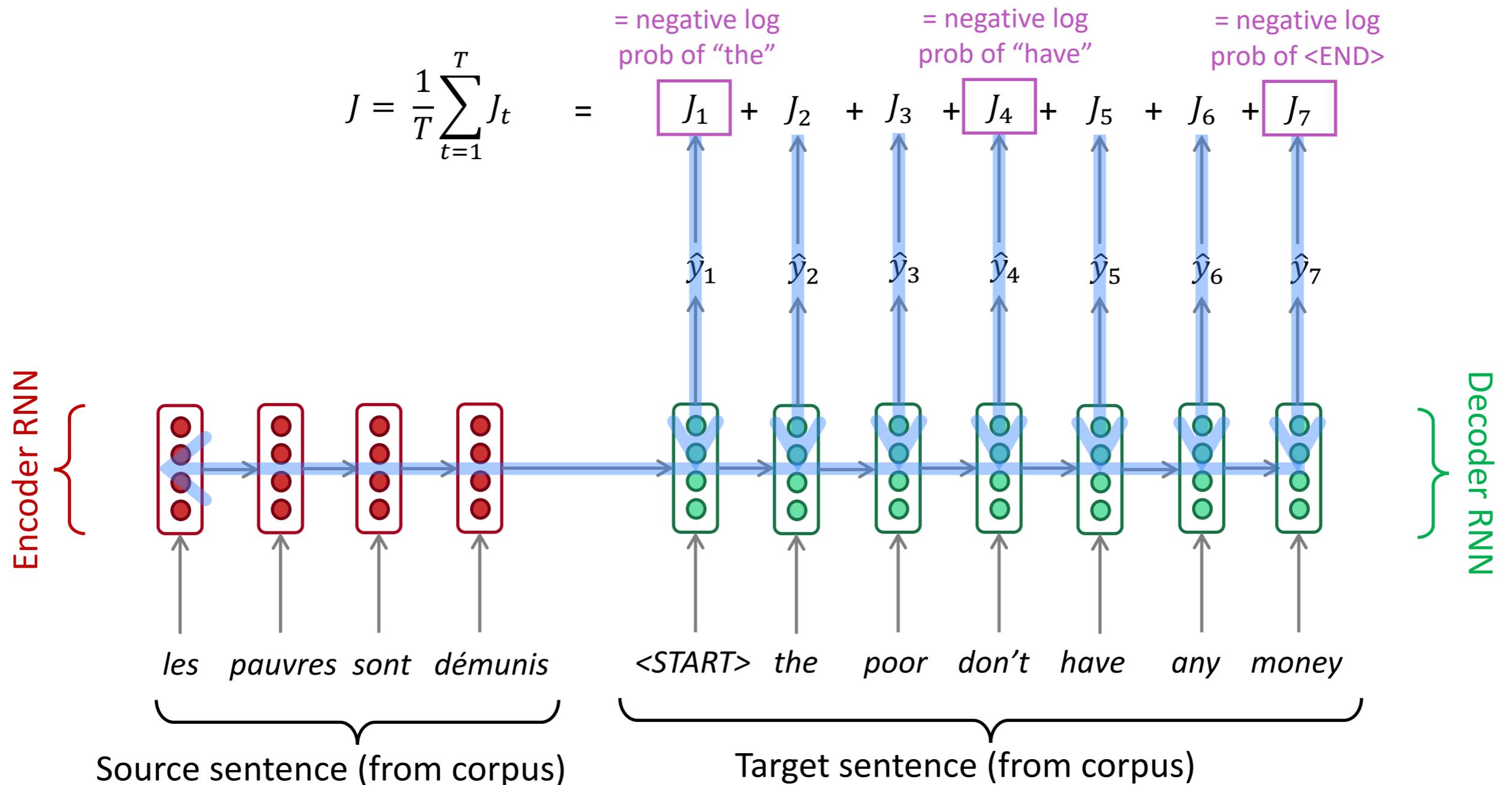# Neural Machine Translation (NMT)

The sequence-to-sequence model

Target sentence (output)

Encoding of the source sentence.
Provides initial hidden state
for Decoder RNN.

Encoder RNN

Decoder RNN

les  pauvres  sont  démunis

<START>  the  poor  don't  have  any  money

the  poor  don't  have  any  money  <END>

Source sentence (input)
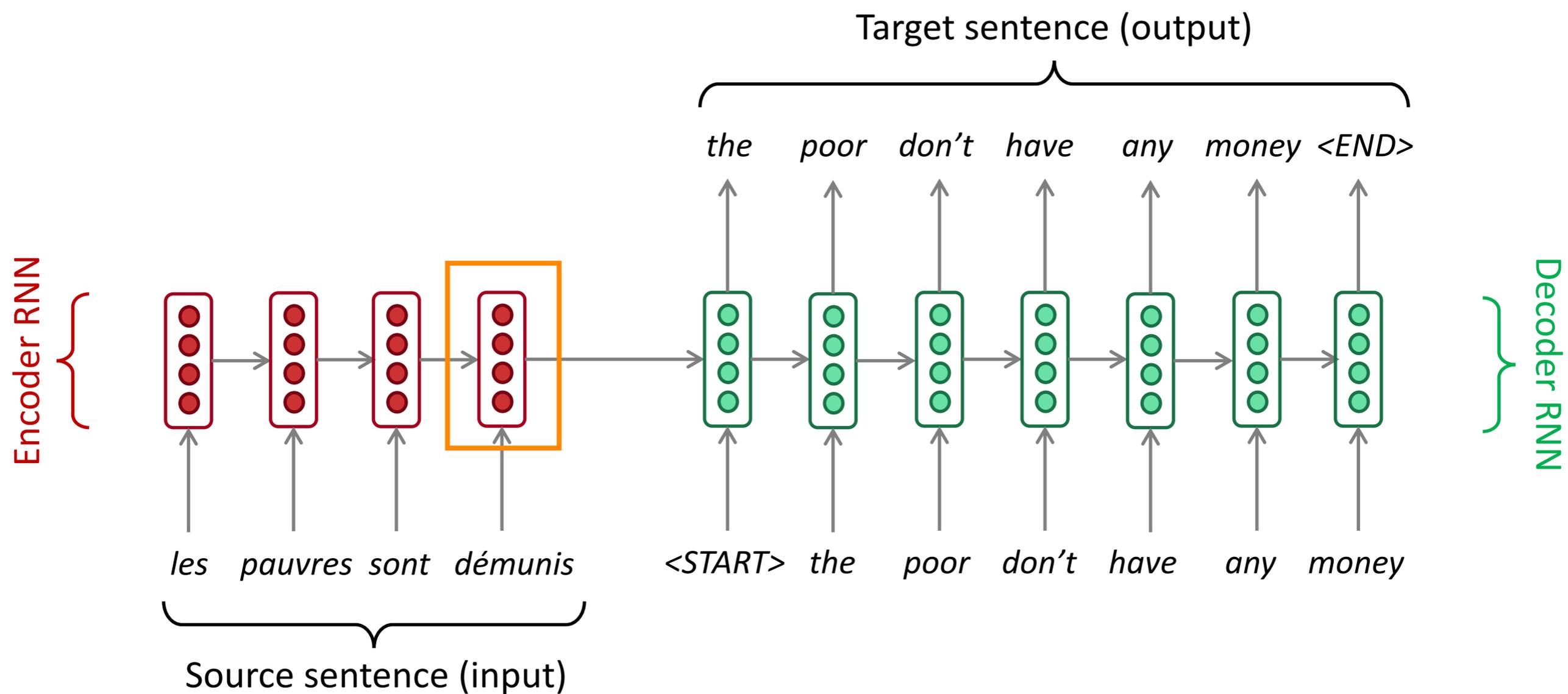
Encoder RNN produces
an encoding of the
source sentence.

Decoder RNN is a Language Model that generates
target sentence conditioned on encoding.

# Training a Neural Machine Translation system

= negative log prob of "the"

= negative log prob of "have"

= negative log prob of <END>

$$J = \frac{1}{T}\sum_{t=1}^{T} J_t \quad = \quad \boxed{J_1} + J_2 + J_3 + \boxed{J_4} + J_5 + J_6 + \boxed{J_7}$$

$\hat{y}_1$ $\hat{y}_2$ $\hat{y}_3$ $\hat{y}_4$ $\hat{y}_5$ $\hat{y}_6$ $\hat{y}_7$

Encoder RNN

Decoder RNN

*les   pauvres   sont   démunis*        *<START>   the   poor   don't   have   any   money*

Source sentence (from corpus)          Target sentence (from corpus)

what are the parameters of this model?

# Sequence-to-sequence: the bottleneck problem

# Sequence-to-sequence: the bottleneck problem

Encoding of the source sentence. This needs to capture *all information* about the source sentence. Information bottleneck!
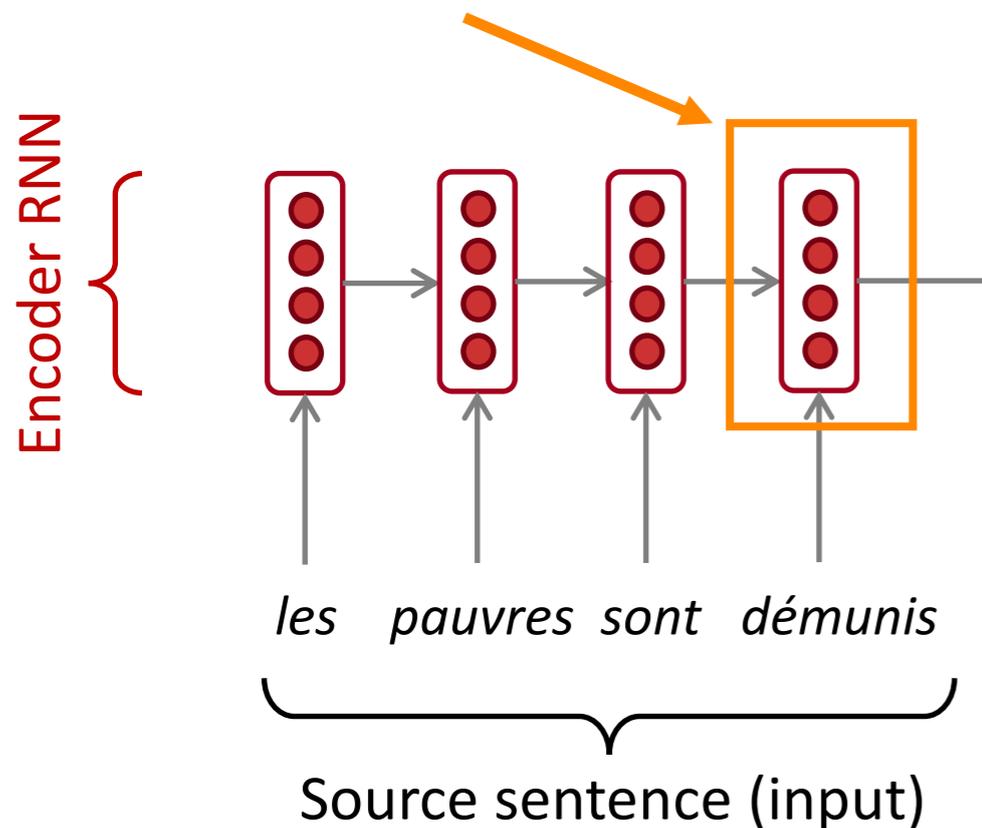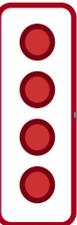
Target sentence (output)

the   poor   don't   have   any   money   <END>

Encoder RNN

Decoder RNN

les   pauvres   sont   démunis

<START>   the   poor   don't   have   any   money

Source sentence (input)

"you can't cram the meaning of a whole  %&@#&ing sentence into a single $*(&@ing vector!"

— Ray Mooney (NLP prof at UT Austin)

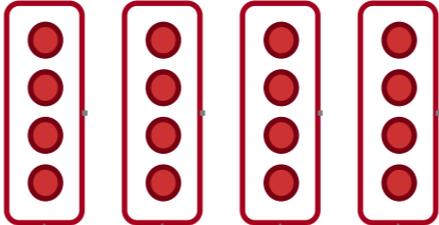# idea: what if we use multiple vectors?

Encoding of the
source sentence.
This needs to capture *all*
*information* about the
source sentence.
Information bottleneck!

Encoder RNN



*les*  *pauvres*  *sont*  *démunis*

Source sentence (input)

Instead of:

les pauvres sont démunis =

Let's try:

les pauvres sont démunis =

(all 4 hidden states!)

# The solution: **attention**

- **Attention mechanisms** (Bahdanau et al., 2015) allow the decoder to focus on a particular part of the source sequence at each time step
  - Conceptually similar to *word alignments*
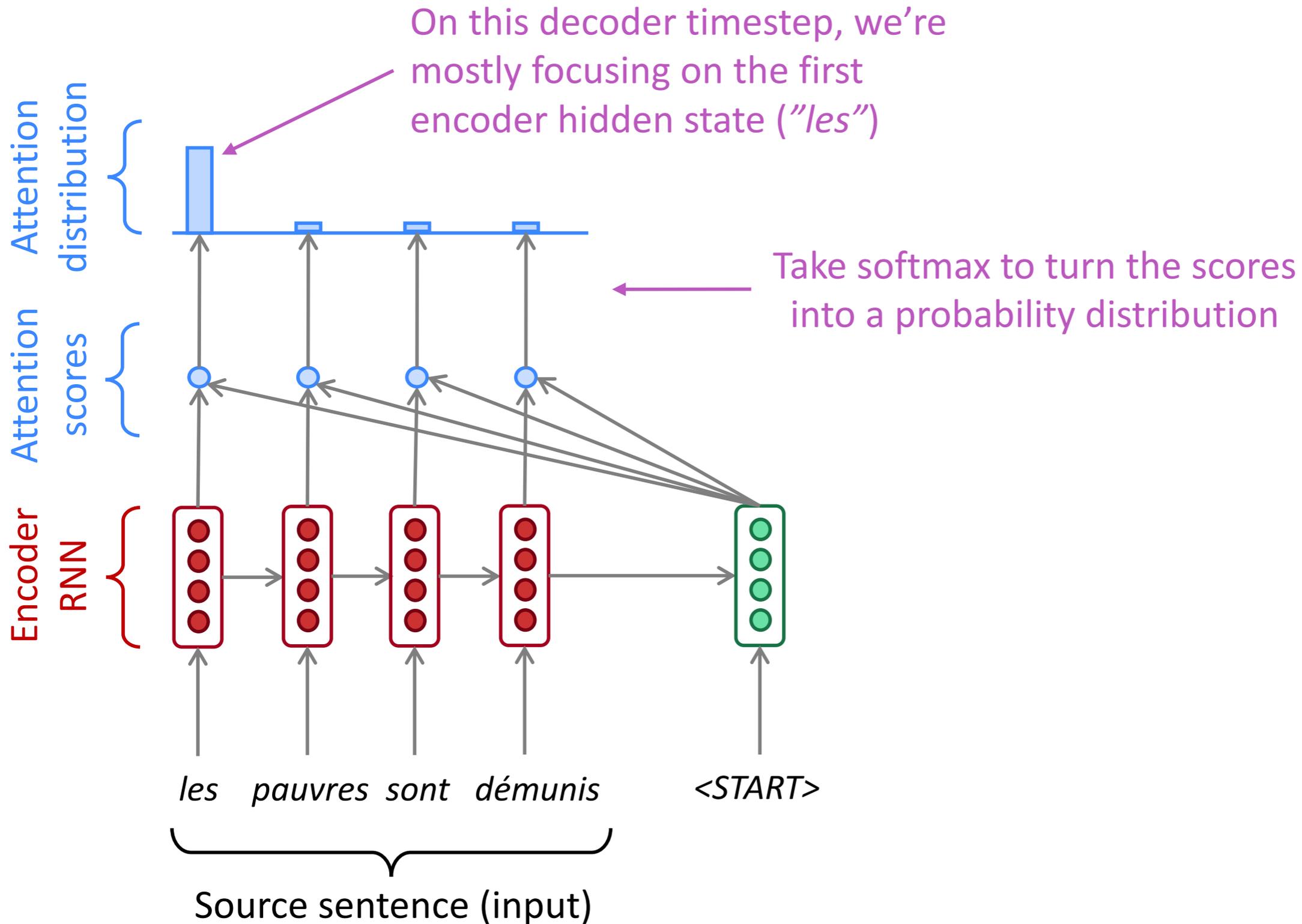
# How does it work?

- in general, we have a single *query* vector and multiple *key* vectors. We want to score each query-key pair

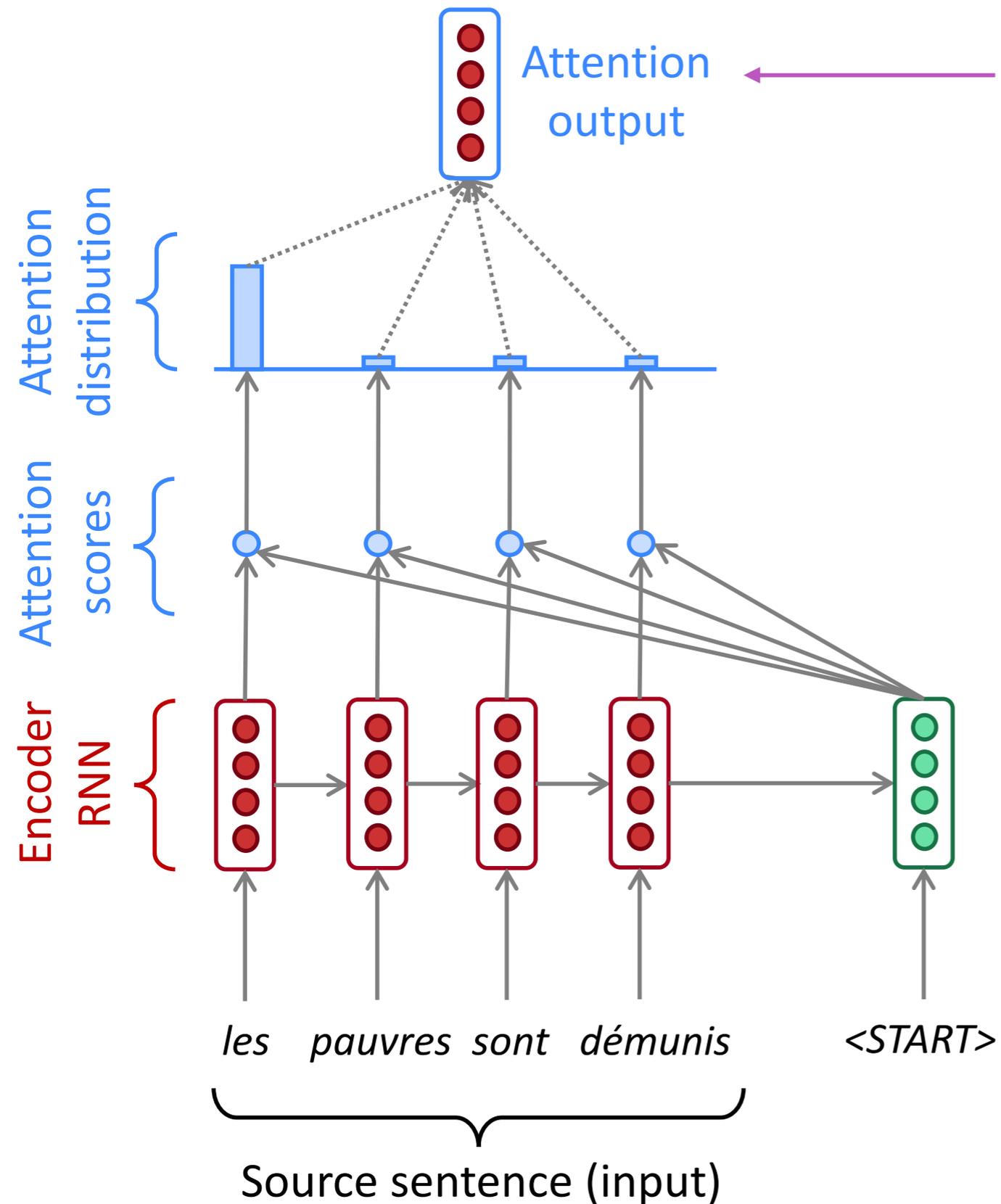in machine translation, what are the queries and keys?

# Sequence-to-sequence with attention



**Attention scores**

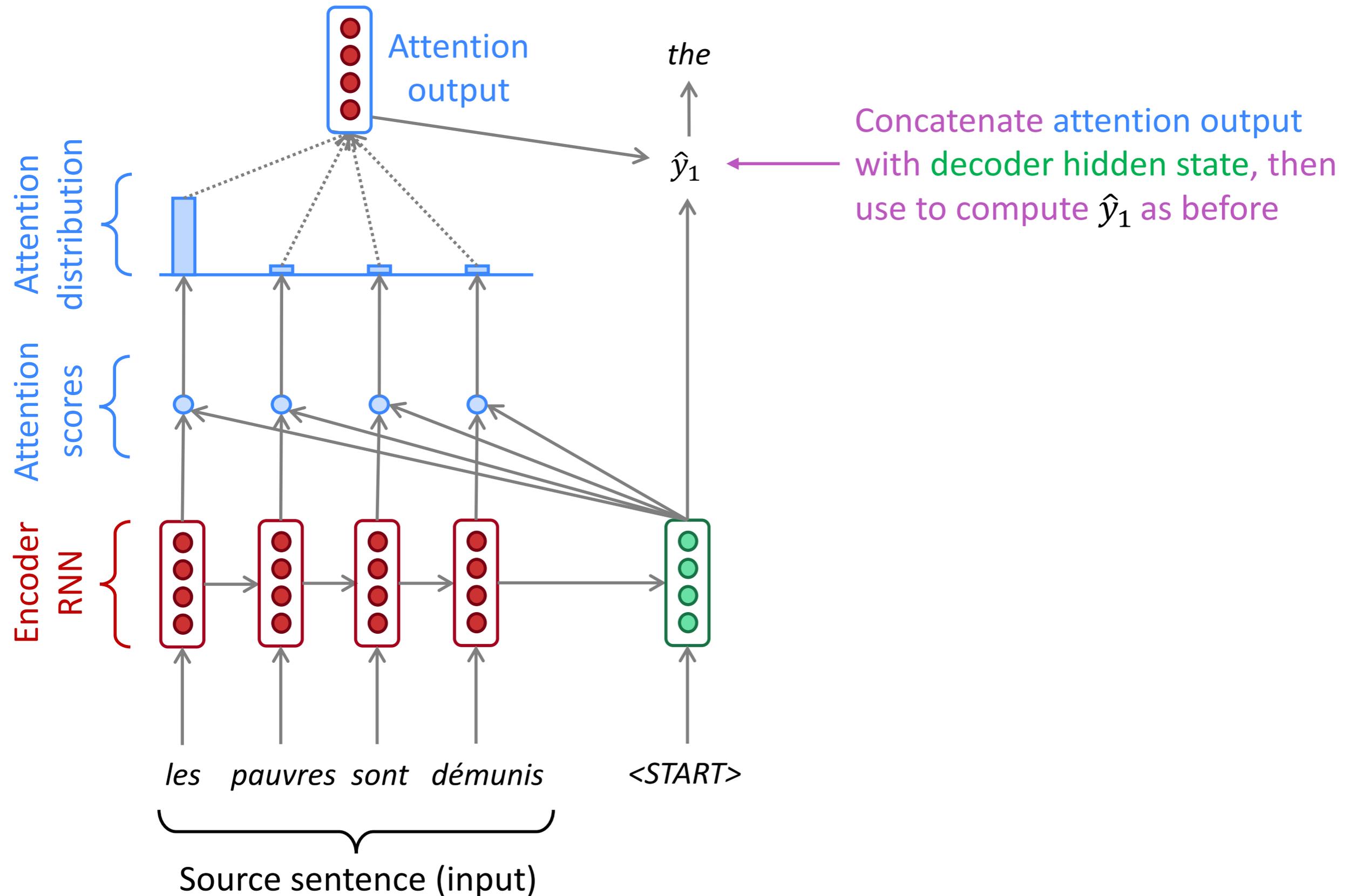**Encoder RNN**

dot product with *keys*
(encoder hidden states)

*Query 1:*
decoder, first time step

les  pauvres  sont  démunis  <START>

Source sentence (input)

# Sequence-to-sequence with attention
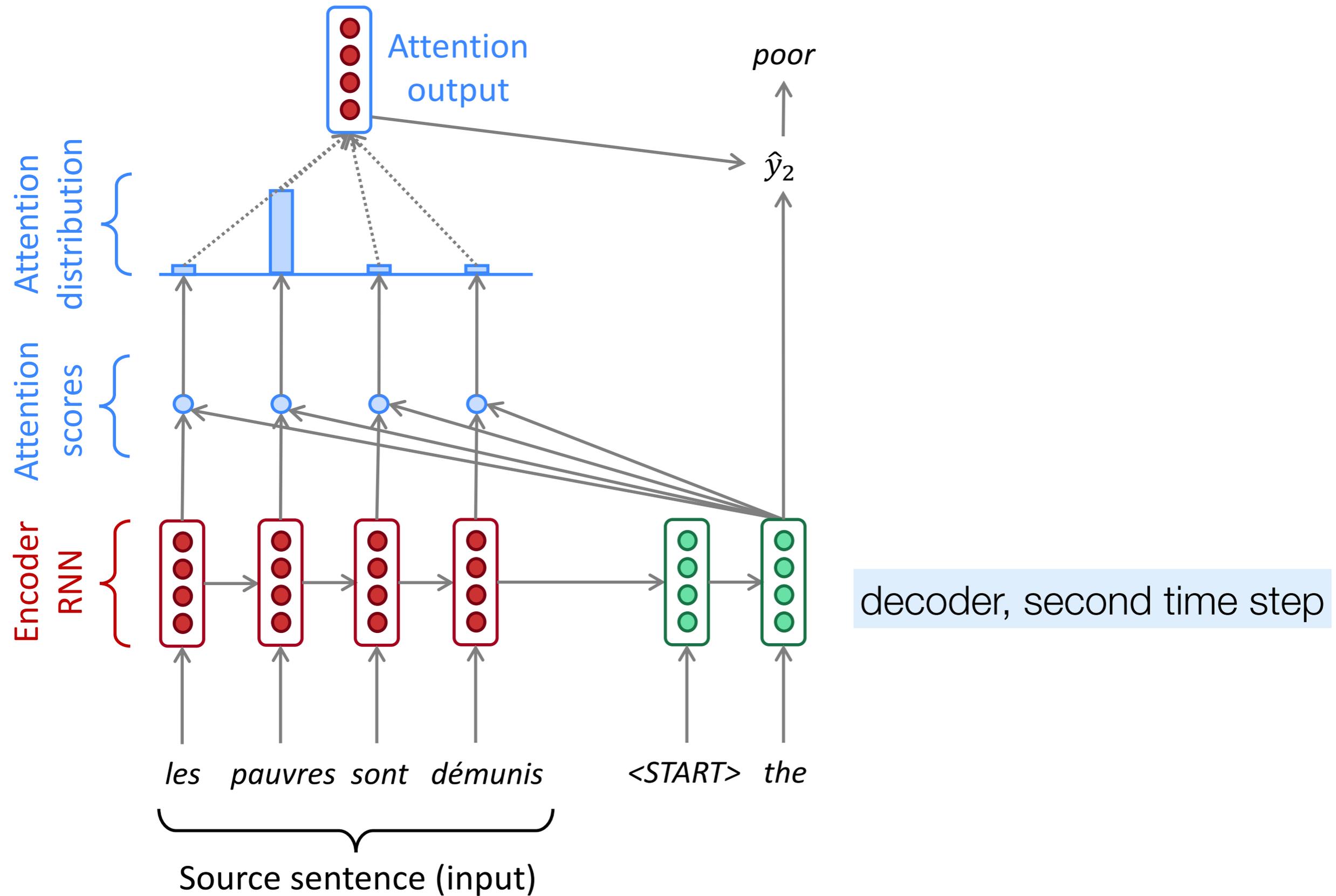
# Sequence-to-sequence with attention



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information the hidden states that received high attention.
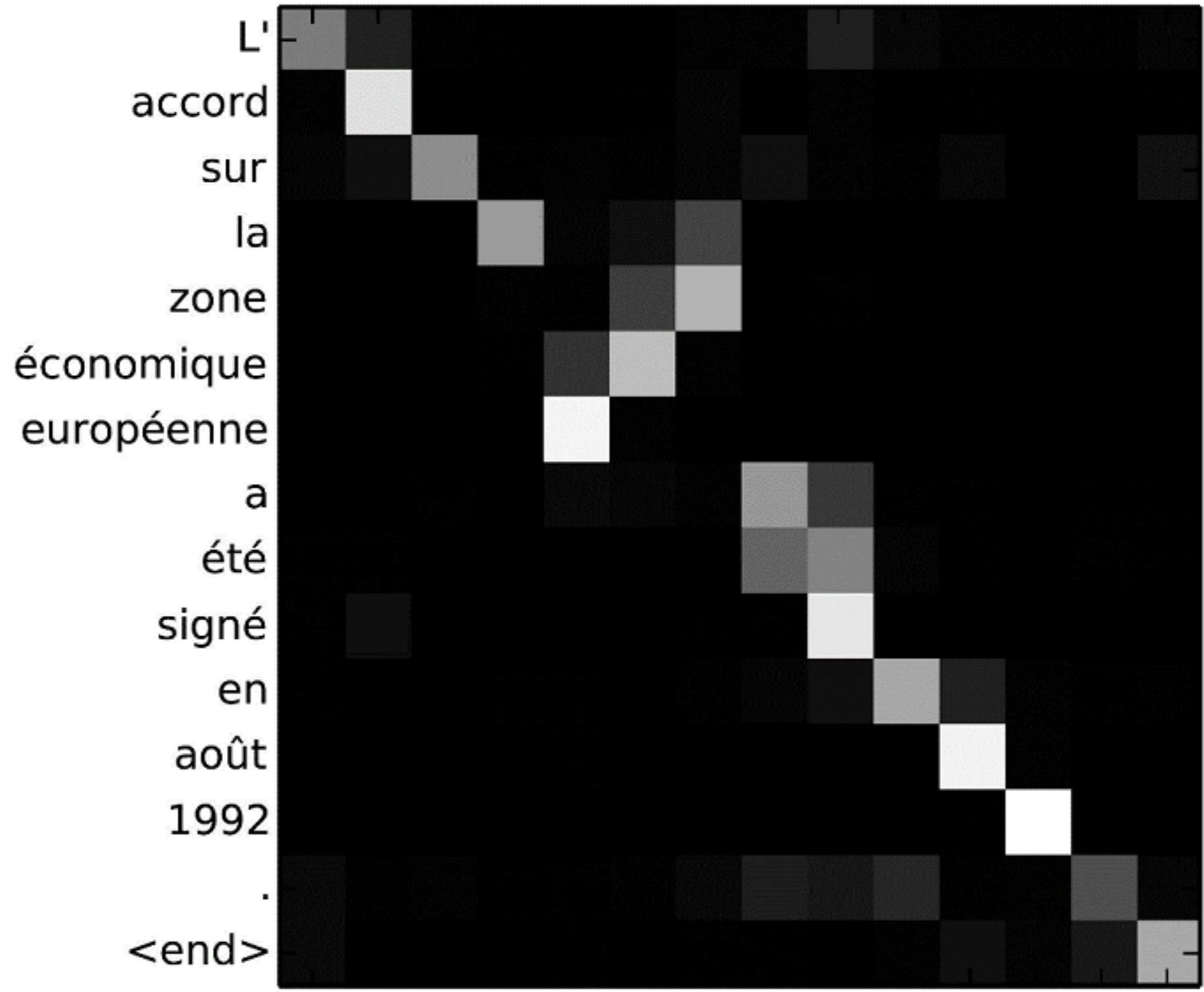
# Sequence-to-sequence with attention



Attention output

*the*

Concatenate attention output with decoder hidden state, then use to compute $\hat{y}_1$ as before

$\hat{y}_1$

Attention distribution

Attention scores

Encoder RNN

*les* *pauvres* *sont* *démunis*

*<START>*

Source sentence (input)

# Sequence-to-sequence with attention



**Attention output**

Attention distribution

Attention scores

Encoder RNN

$\hat{y}_2$

*poor*

decoder, second time step

*les*   *pauvres*   *sont*   *démunis*   <START>   *the*

Source sentence (input)

# Attention is great

- Attention significantly improves NMT performance
  - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
  - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
  - Provides shortcut to faraway states
- Attention provides some interpretability
  - By inspecting attention distribution, we can see what the decoder was focusing on
  - We get alignment for free!
  - This is cool because we never explicitly trained an alignment system
  - The network just learned alignment by itself

18

# Many variants of attention

- Original formulation:  $a(\mathbf{q}, \mathbf{k}) = w_2^T \tanh(W_1[\mathbf{q}; \mathbf{k}])$

- Bilinear product:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T W \mathbf{k}$       Luong et al., 2015
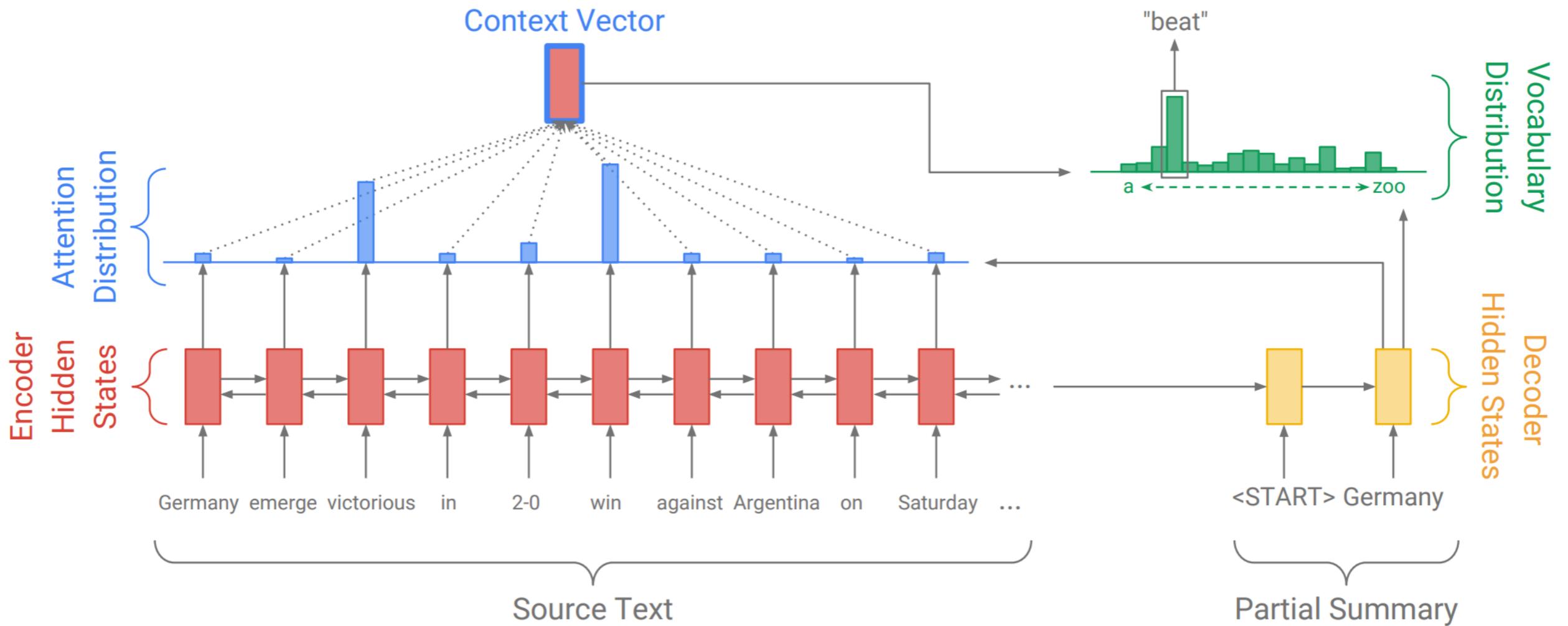
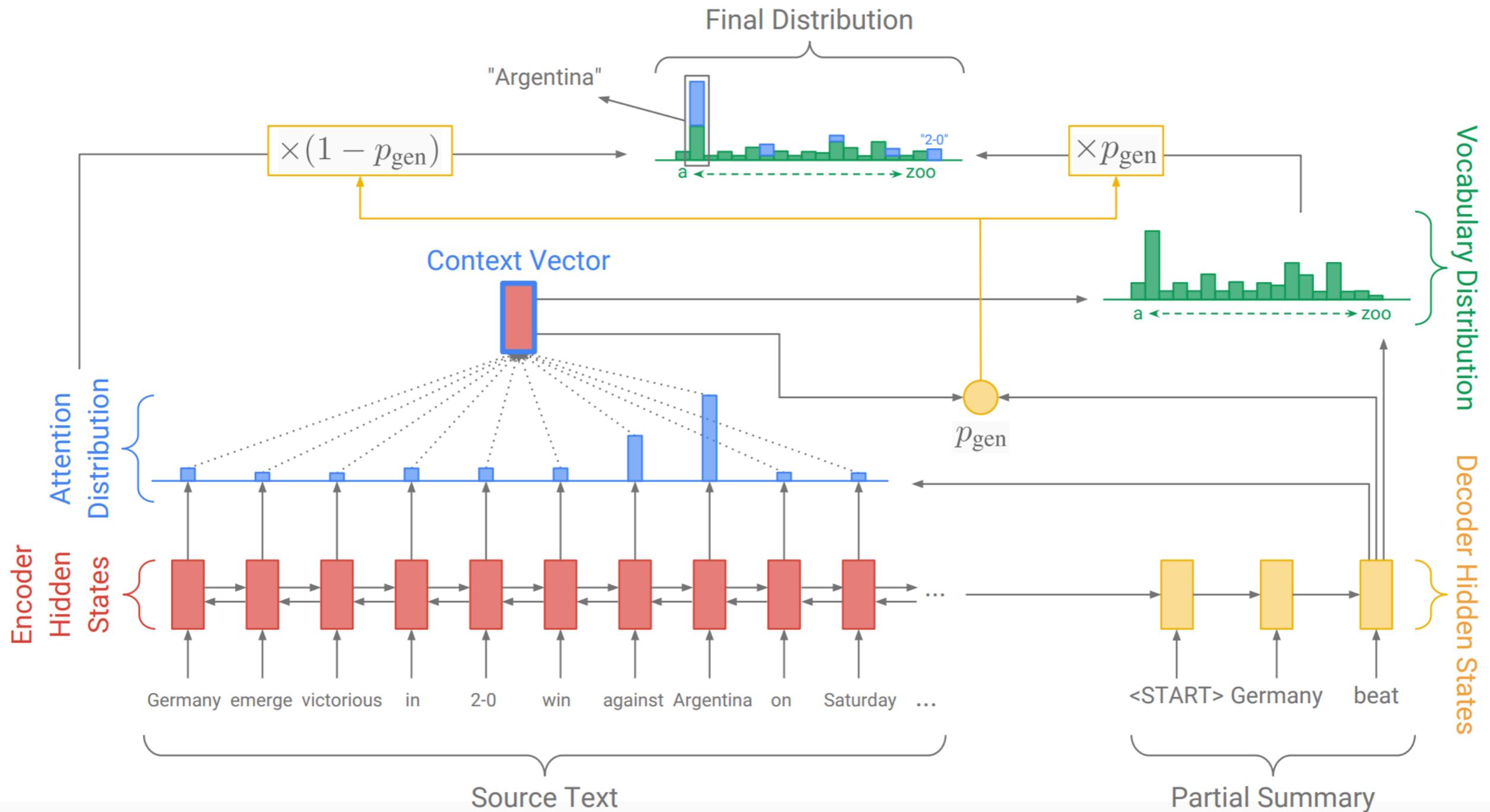- Dot product:  $a(\mathbf{q}, \mathbf{k}) = \mathbf{q}^T \mathbf{k}$       Luong et al., 2015

- Scaled dot product:  $a(\mathbf{q}, \mathbf{k}) = \dfrac{\mathbf{q}^T \mathbf{k}}{\sqrt{|\mathbf{k}|}}$       Vaswani et al., 2017

# Attention is not just for MT!

Here we have a standard seq2seq
model for summarization

See et al., 2017

Here we have a seq2seq model with a **copy mechanism** for summarization

22

See et al., 2017

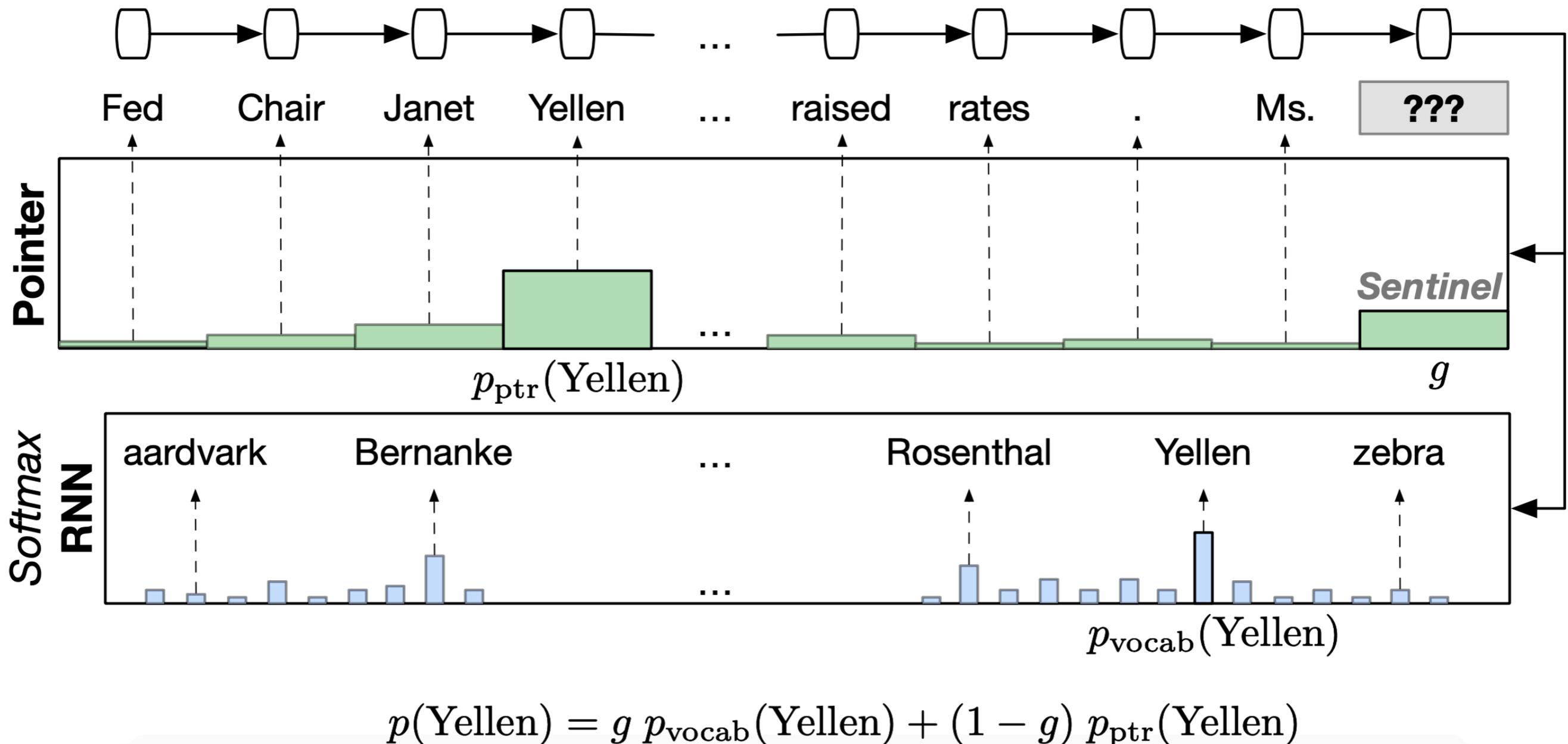# Target-side attention (in LMs or more complex MT models)



$$p(\text{Yellen}) = g\, p_{\text{vocab}}(\text{Yellen}) + (1 - g)\, p_{\text{ptr}}(\text{Yellen})$$

Merity et al., 2016

# Image Captioning with Attention



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

Xu et al., 2015

# visual attention

- Use the question representation *q* to determine where in the image to look



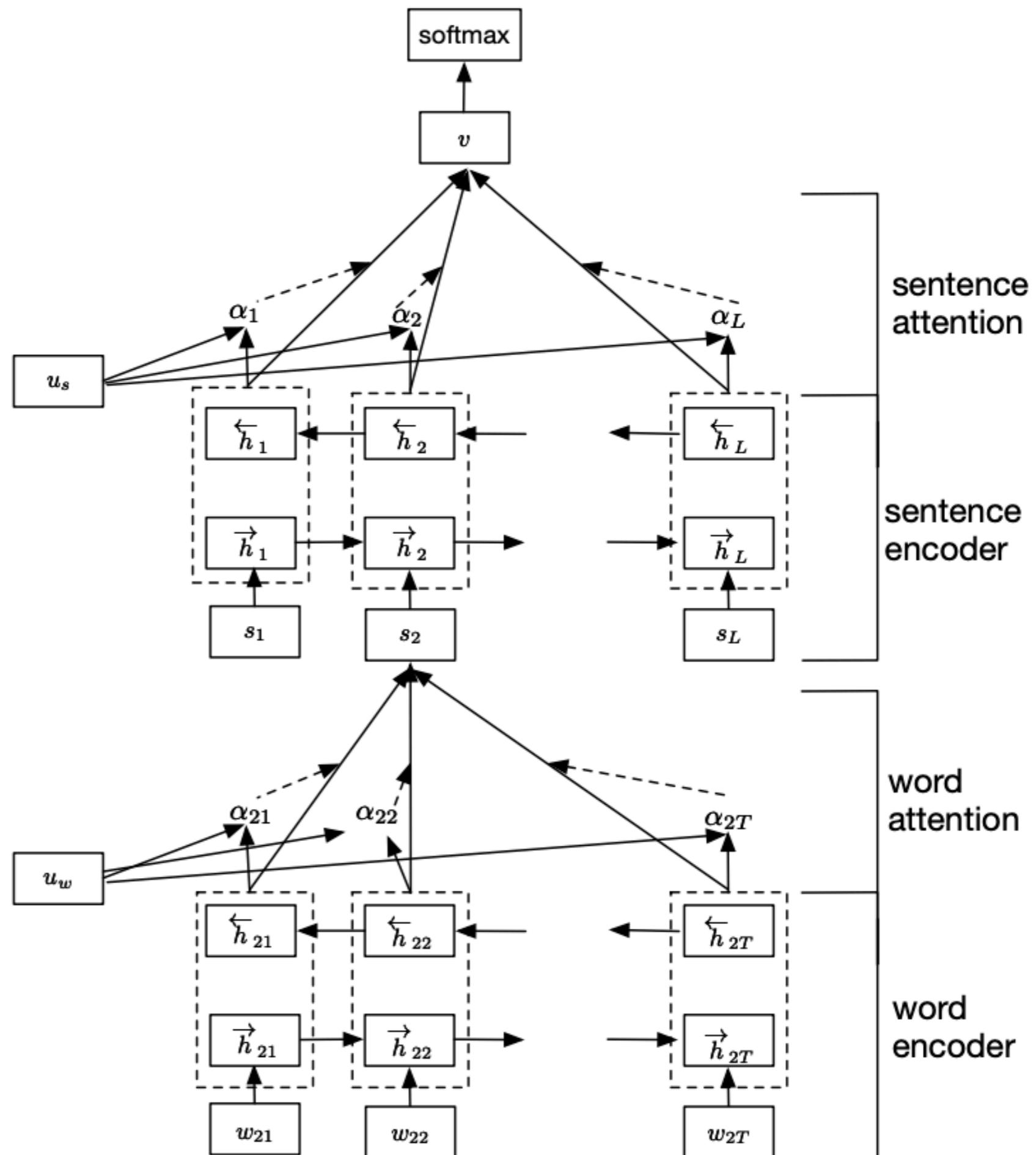How many benches are shown?

softmax:
predict answer

attention over final convolutional
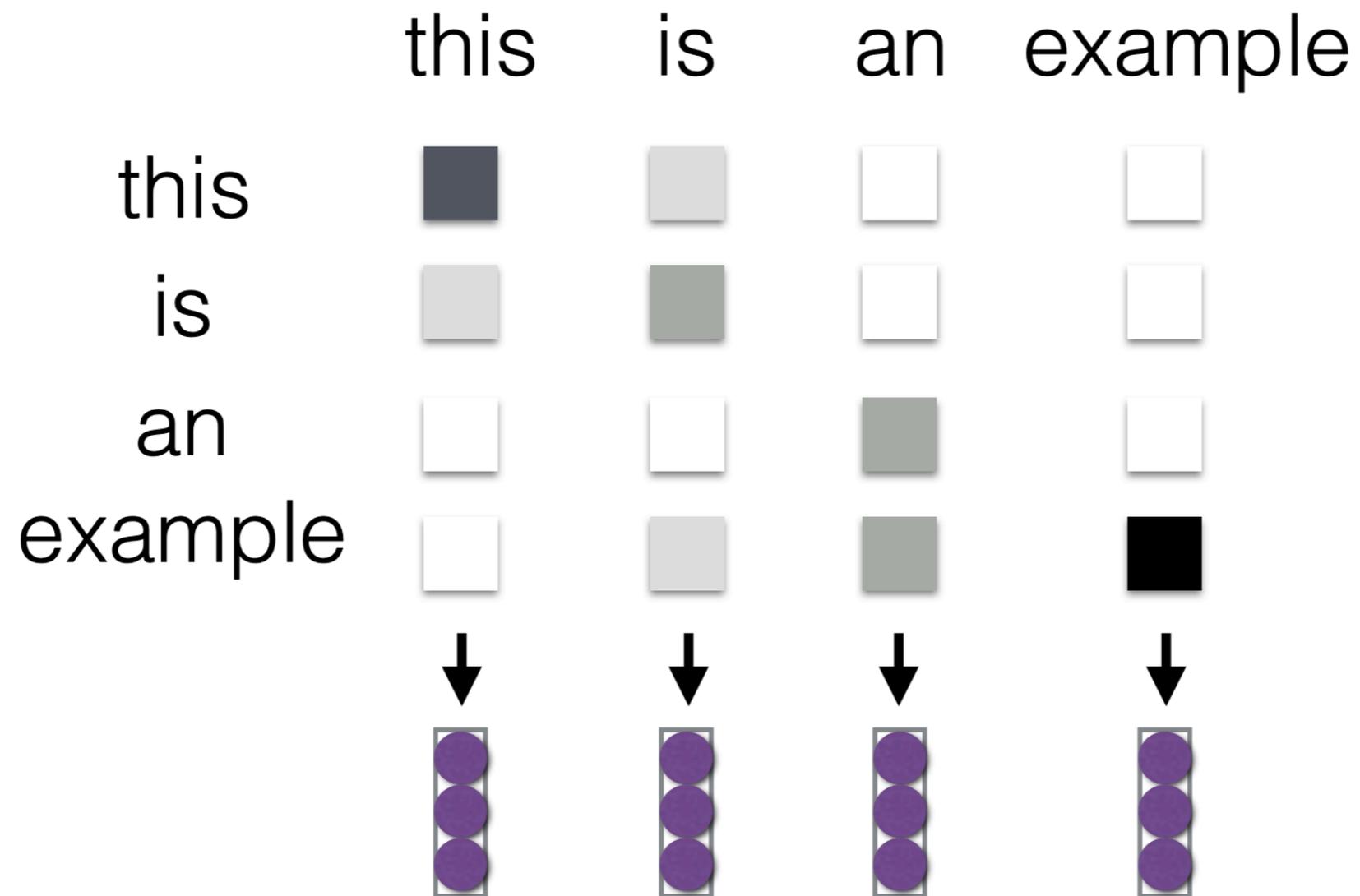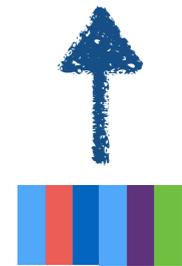layer in network: 196 boxes, captures
color and positional information

| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.05 | 0.05 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.2 | 0.1 | 0.05 | 0.0 | 0.0 | 0.0 |
| 0.3 | 0.2 | 0.05 | 0.0 | 0.0 | 0.0 |

How many benches are shown?

# Hierarchical attention



softmax

$v$

sentence attention

$\alpha_1$  $\alpha_2$  $\alpha_L$

$u_s$

$\overleftarrow{h}_1$  $\overleftarrow{h}_2$  $\overleftarrow{h}_L$

sentence encoder

$\overrightarrow{h}_1$  $\overrightarrow{h}_2$  $\overrightarrow{h}_L$

$s_1$  $s_2$  $s_L$

word attention

$\alpha_{21}$  $\alpha_{22}$  $\alpha_{2T}$

$u_w$

$\overleftarrow{h}_{21}$  $\overleftarrow{h}_{22}$  $\overleftarrow{h}_{2T}$

word encoder

$\overrightarrow{h}_{21}$  $\overrightarrow{h}_{22}$  $\overrightarrow{h}_{2T}$

$w_{21}$  $w_{22}$  $w_{2T}$

Yang et al., 2016

# Self-attention as an encoder!
## (core component of Transformer)

# Attention variants

# hard attention

attention over final convolutional layer in network: 196 boxes, captures color and positional information

softmax: predict answer



| | | | | | |
|---|---|---|---|---|---|
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1.0 | 0.0 | 0.05 | 0.0 | 0.0 | 0.0 |

we can use *reinforcement learning* to focus on just one box

How many benches are shown?

Xu et al., 2015

# Multi-headed attention

- Intuition: $k$ different attentions, each of which is computed independently and focuses on different parts of the sentence

- Transformers = stacked layers of multi-headed self-attention
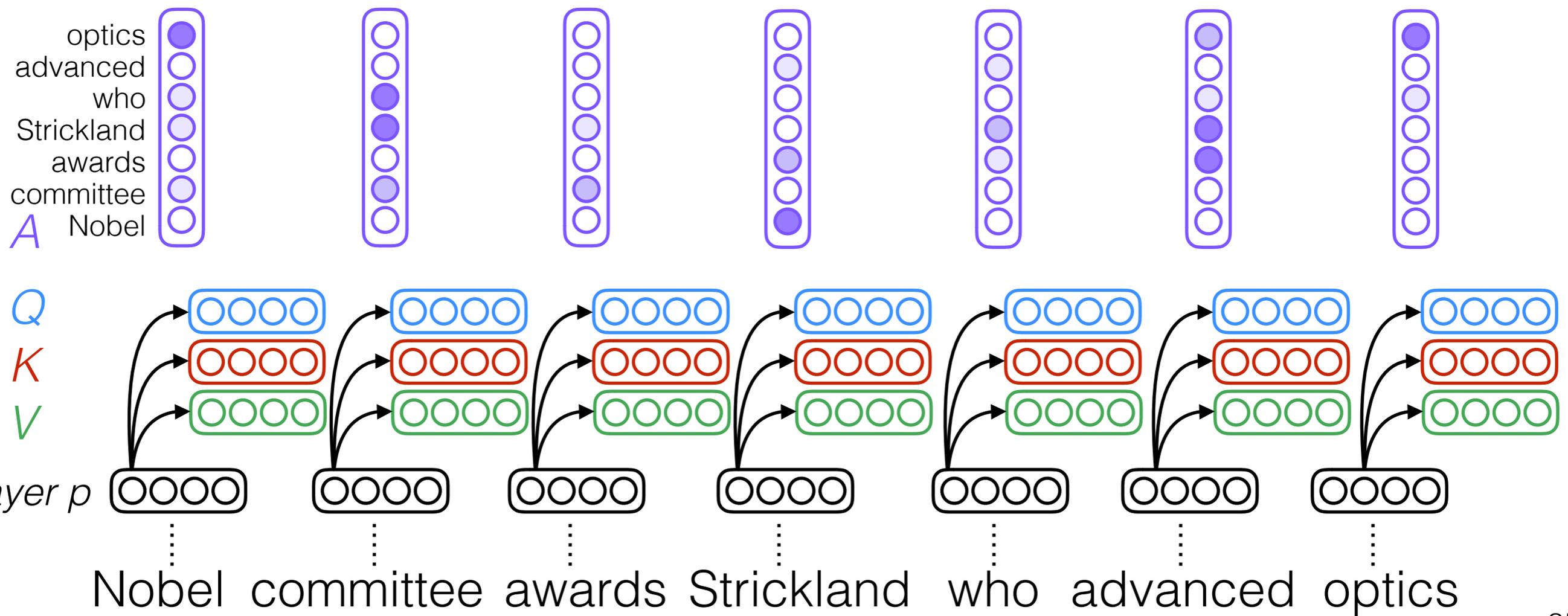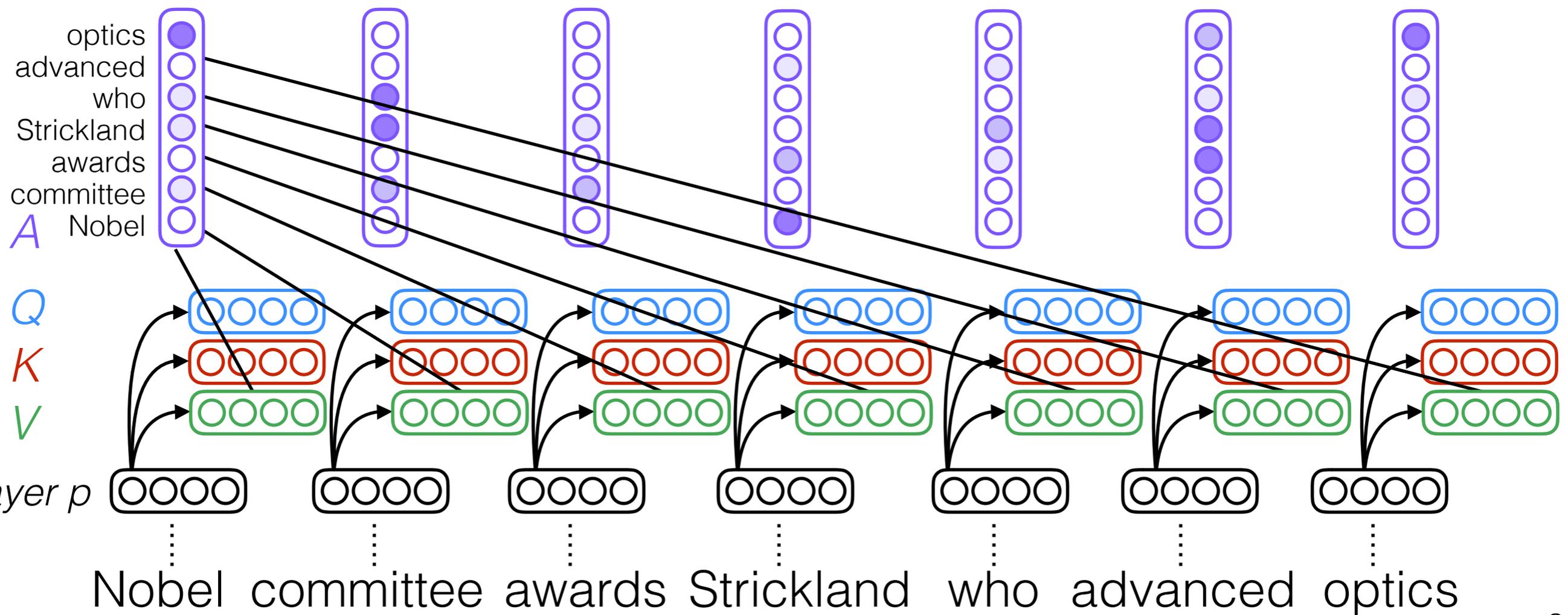
# Self-attention



Nobel  committee  awards  Strickland  who  advanced  optics

# Self-attention

$Q$
$K$
$V$

*Layer p*

Nobel  committee  awards  Strickland  who  advanced  optics

# Self-attention



optics
advanced
who
Strickland
awards
committee
Nobel

$Q$
$K$
$V$

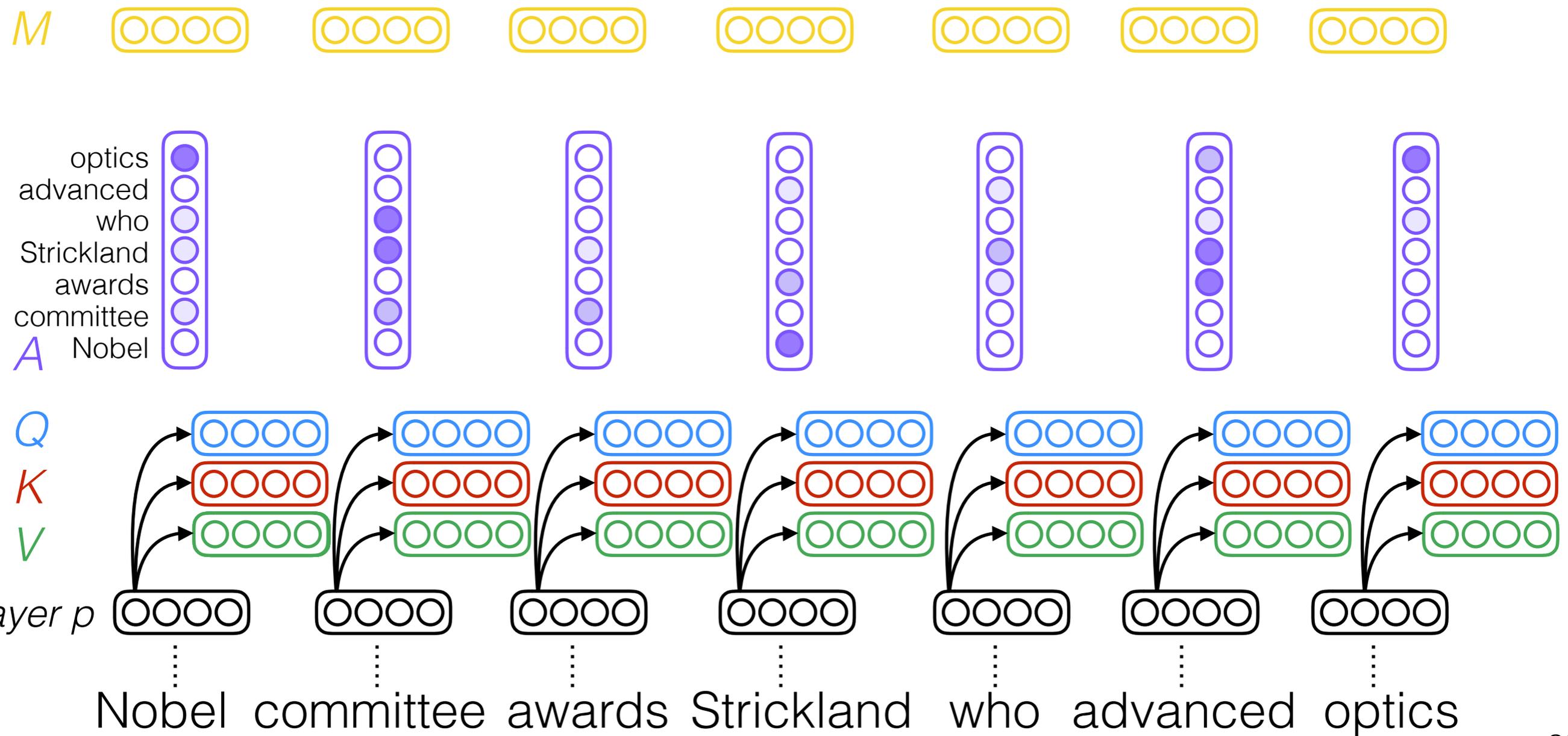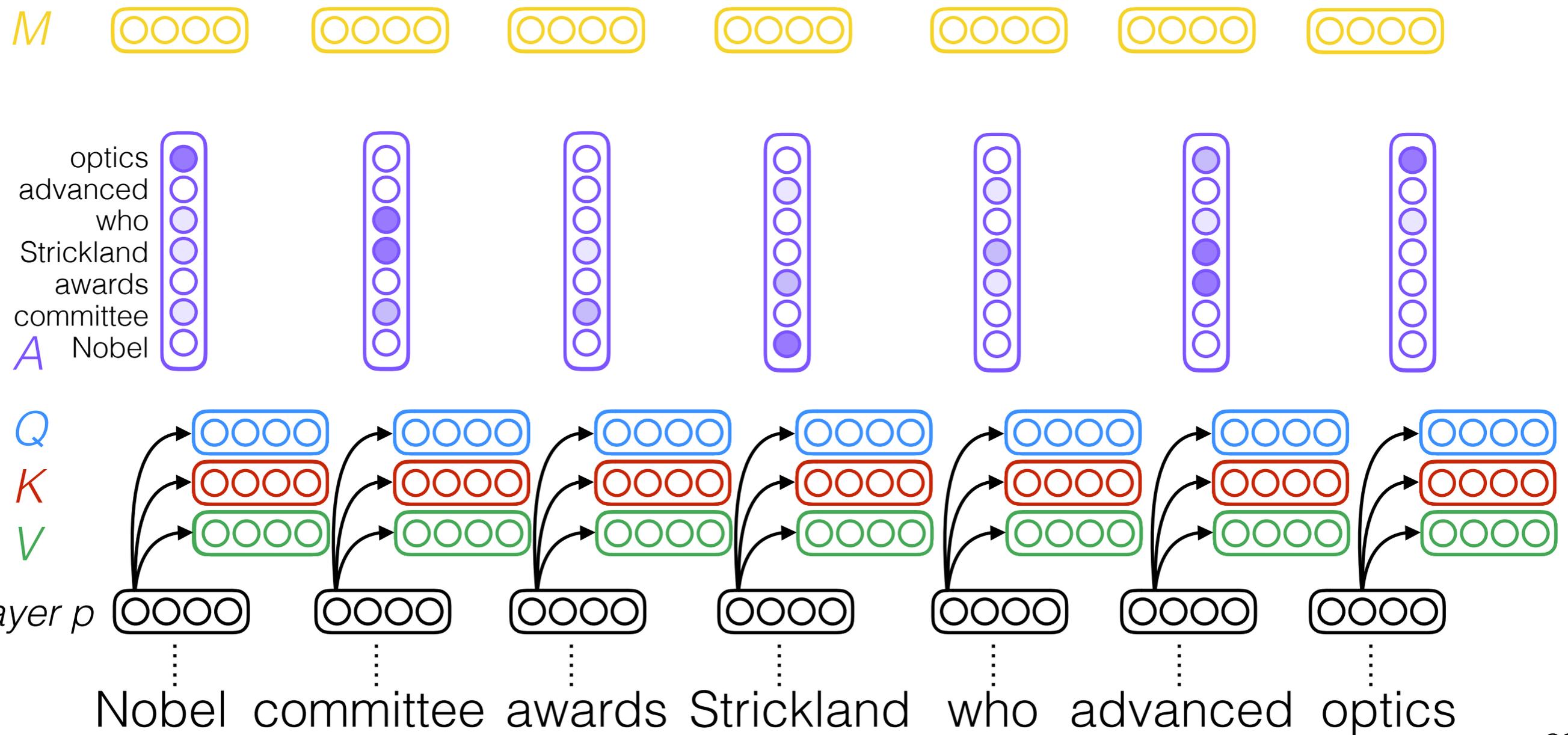*Layer p*

Nobel committee awards Strickland who advanced optics

# Self-attention

# Self-attention

# Self-attention

# Self-attention
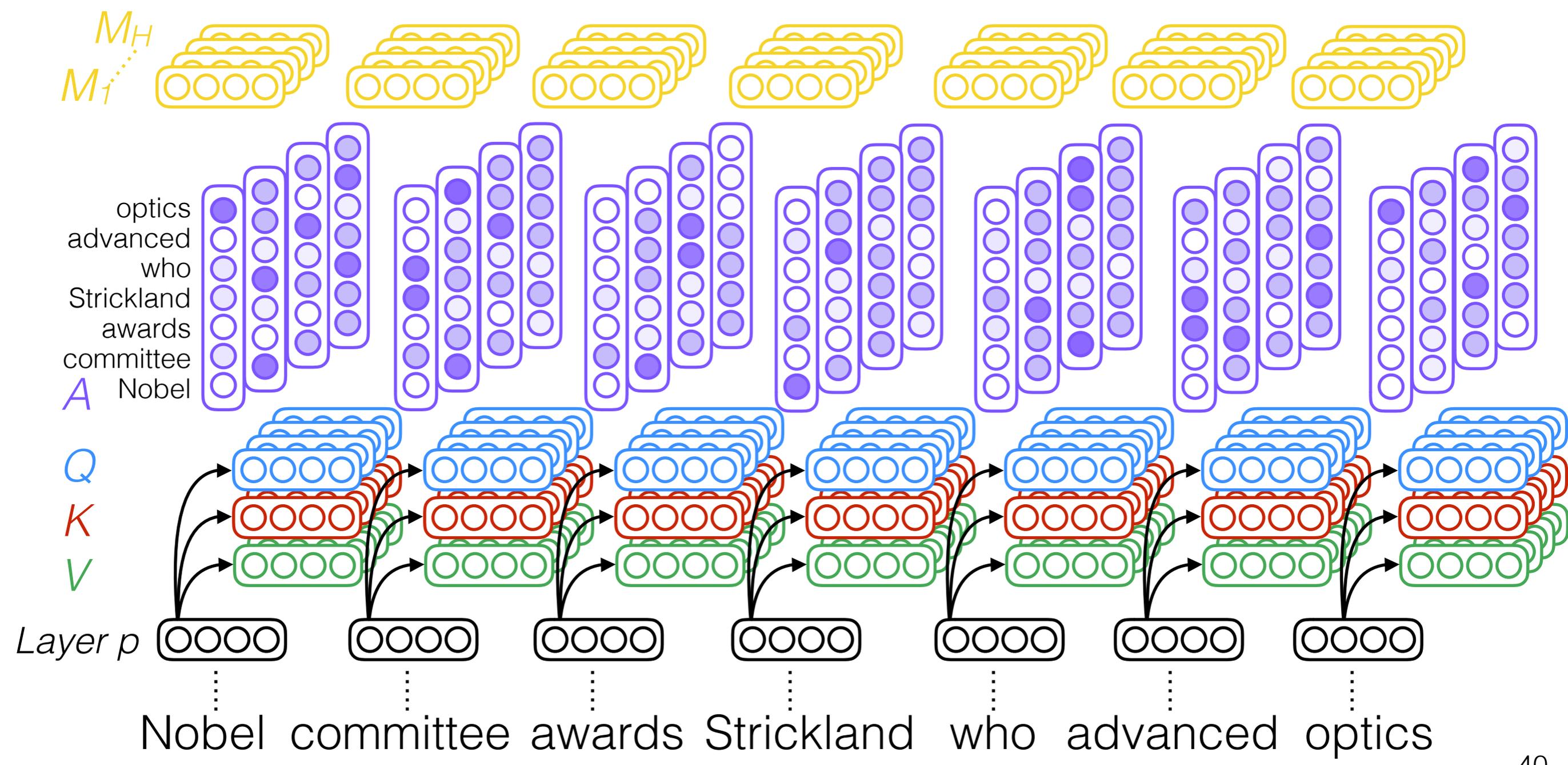


$M$

optics
advanced
who
Strickland
awards
committee
$A$  Nobel

$Q$

$K$

$V$

*Layer p*

Nobel  committee  awards  Strickland  who  advanced  optics
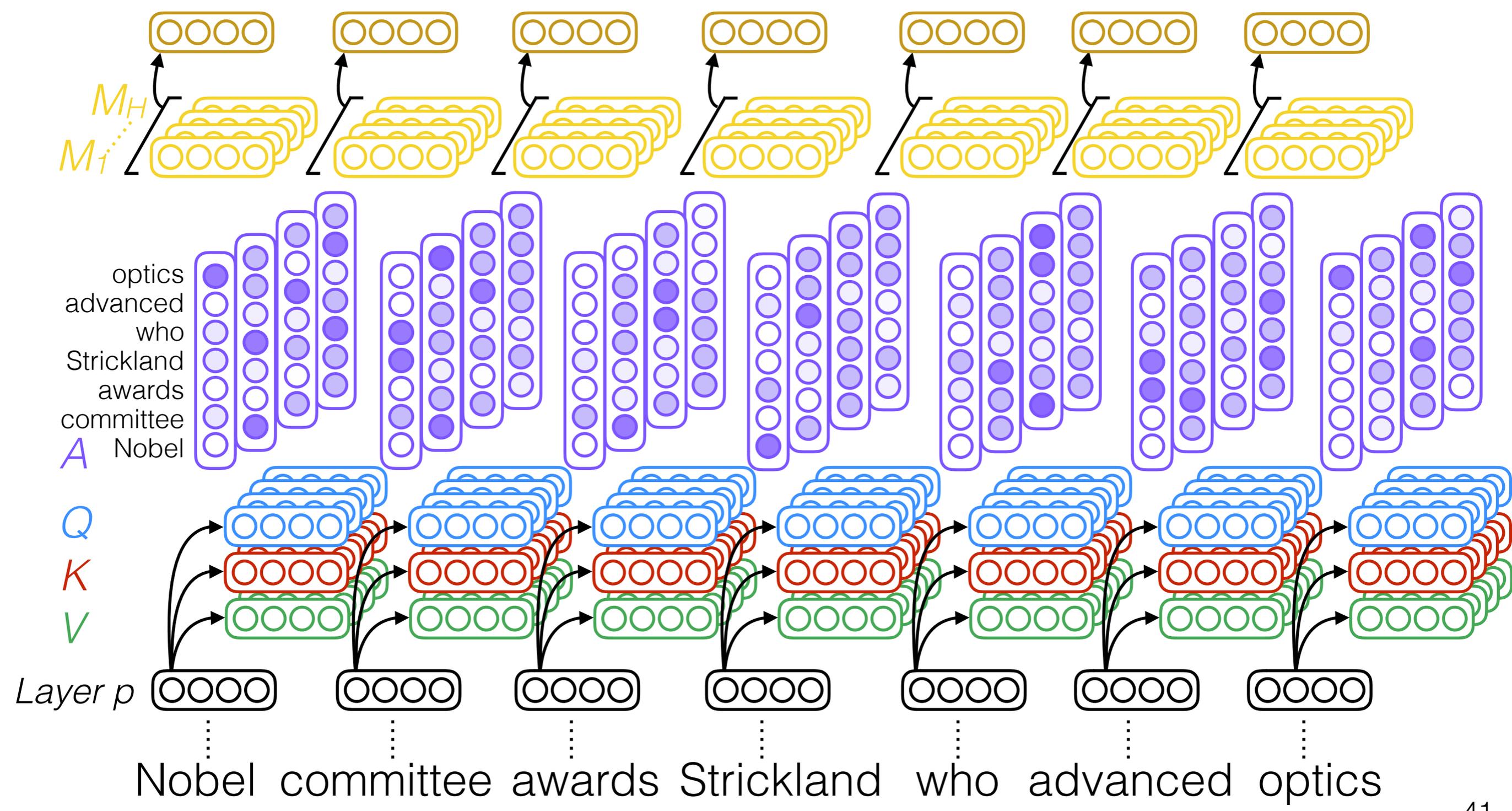
# Self-attention

39

# Multi-head self-attention

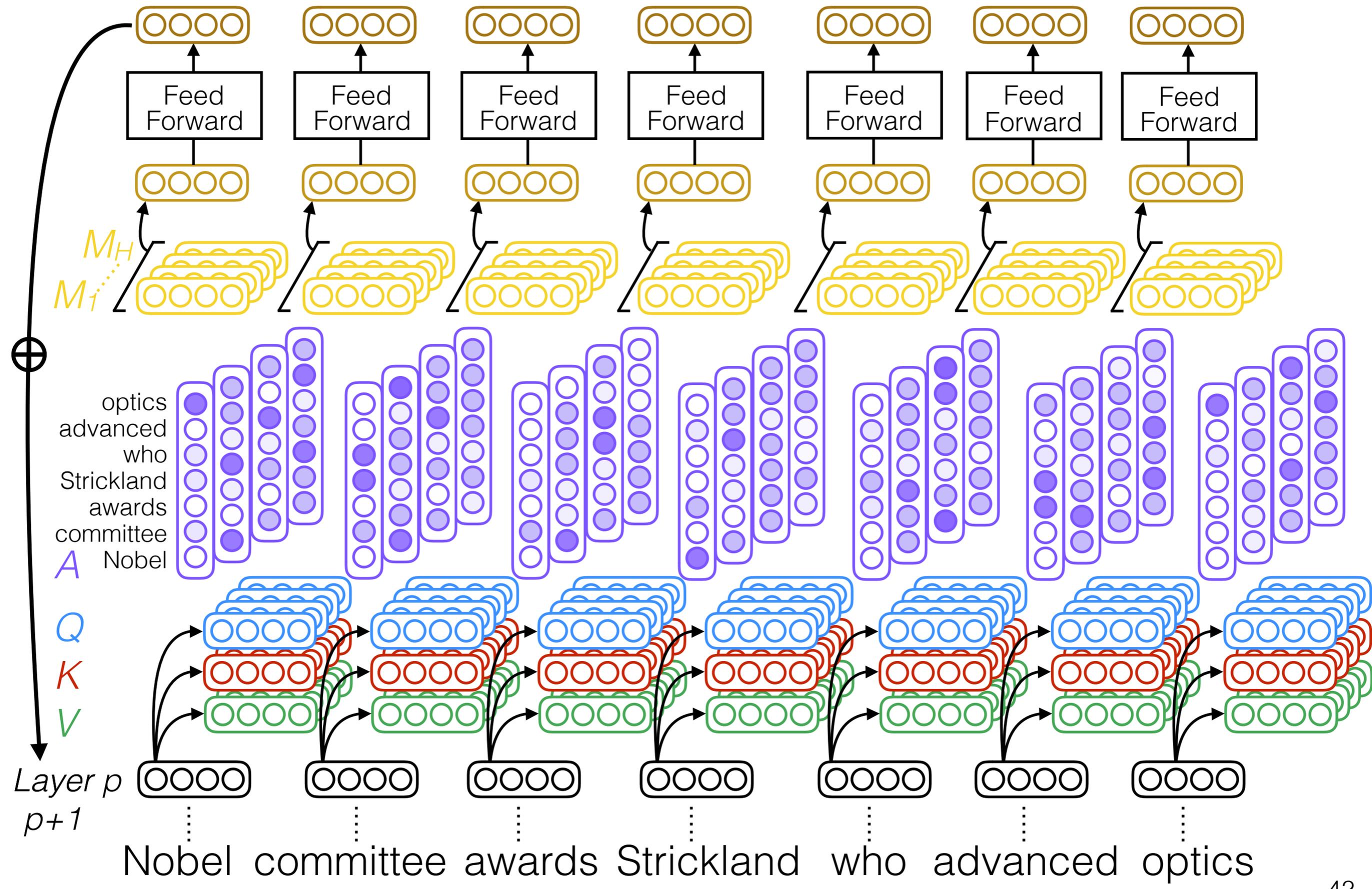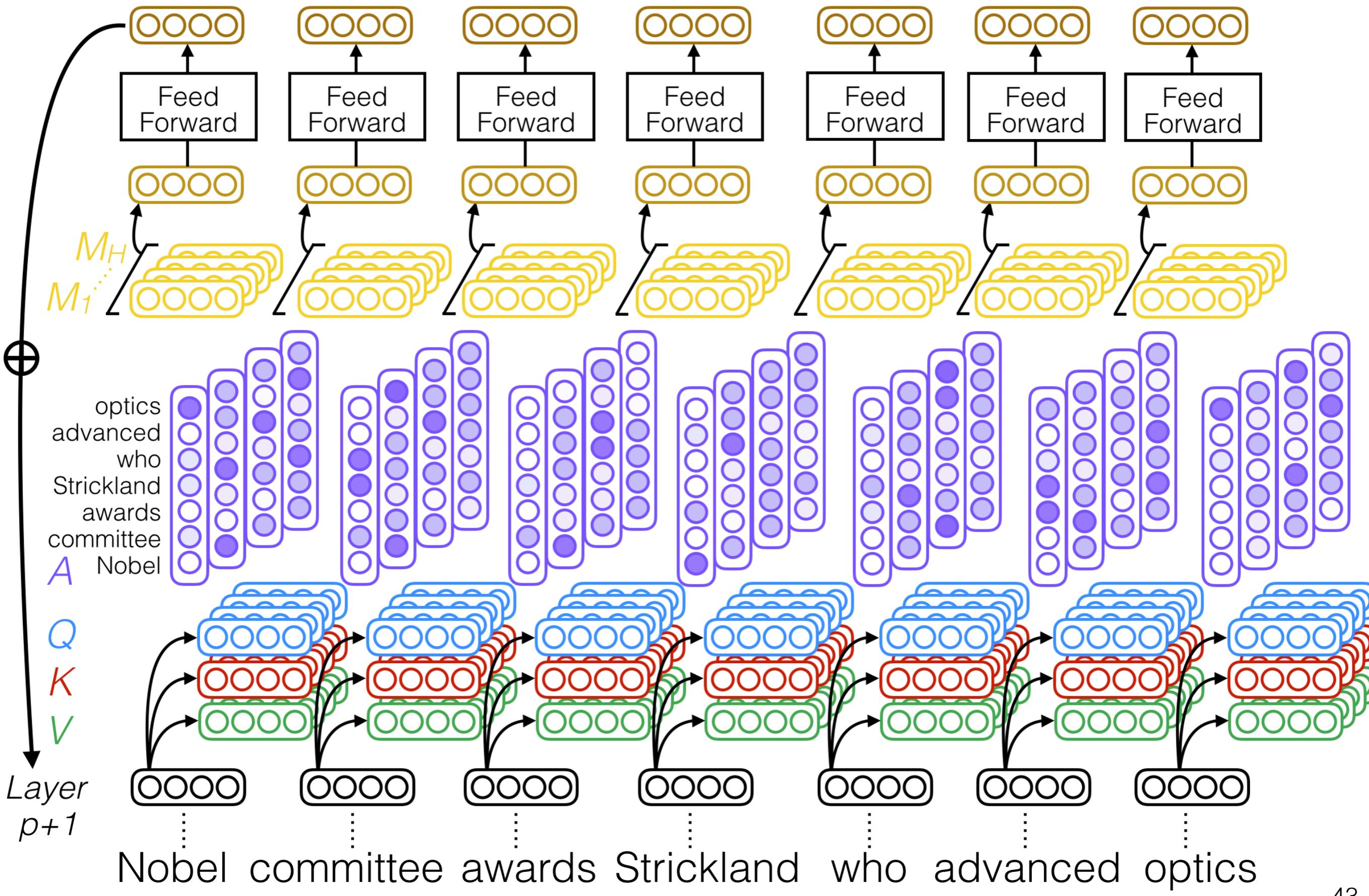# Multi-head self-attention

# Multi-head self-attention

42

# Multi-head self-attention

# Multi-head self-attention