

# Attention Mechanisms

---

Jacob Andreas / MIT 6.804-6.864 / Spring 2021

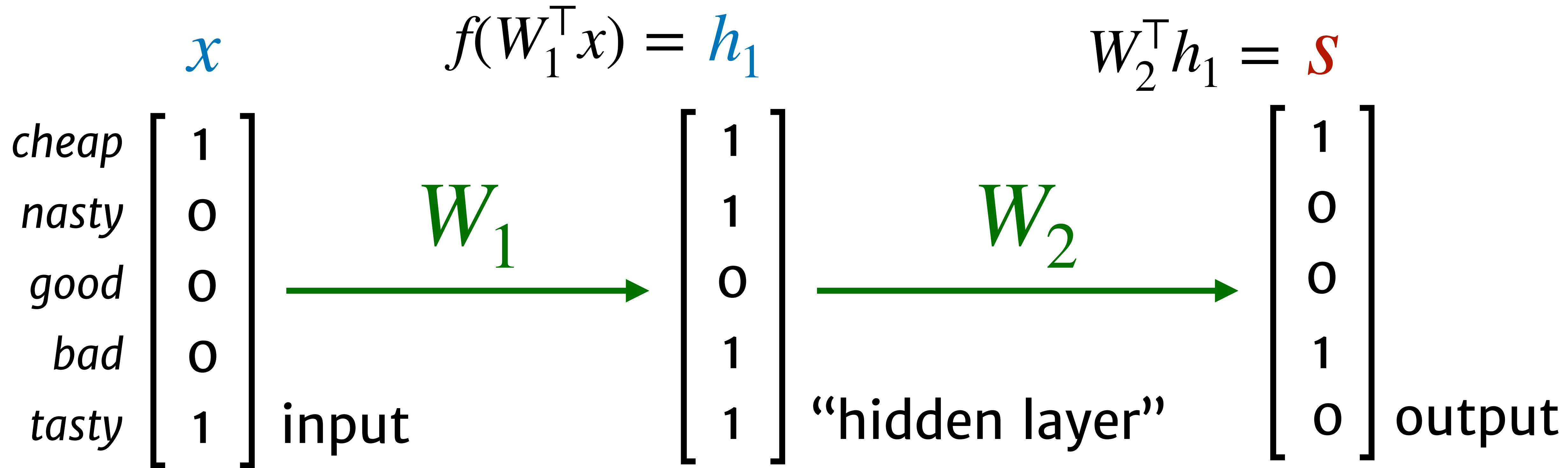
# Admin

# Recap: recurrent neural networks

# Neural networks

---

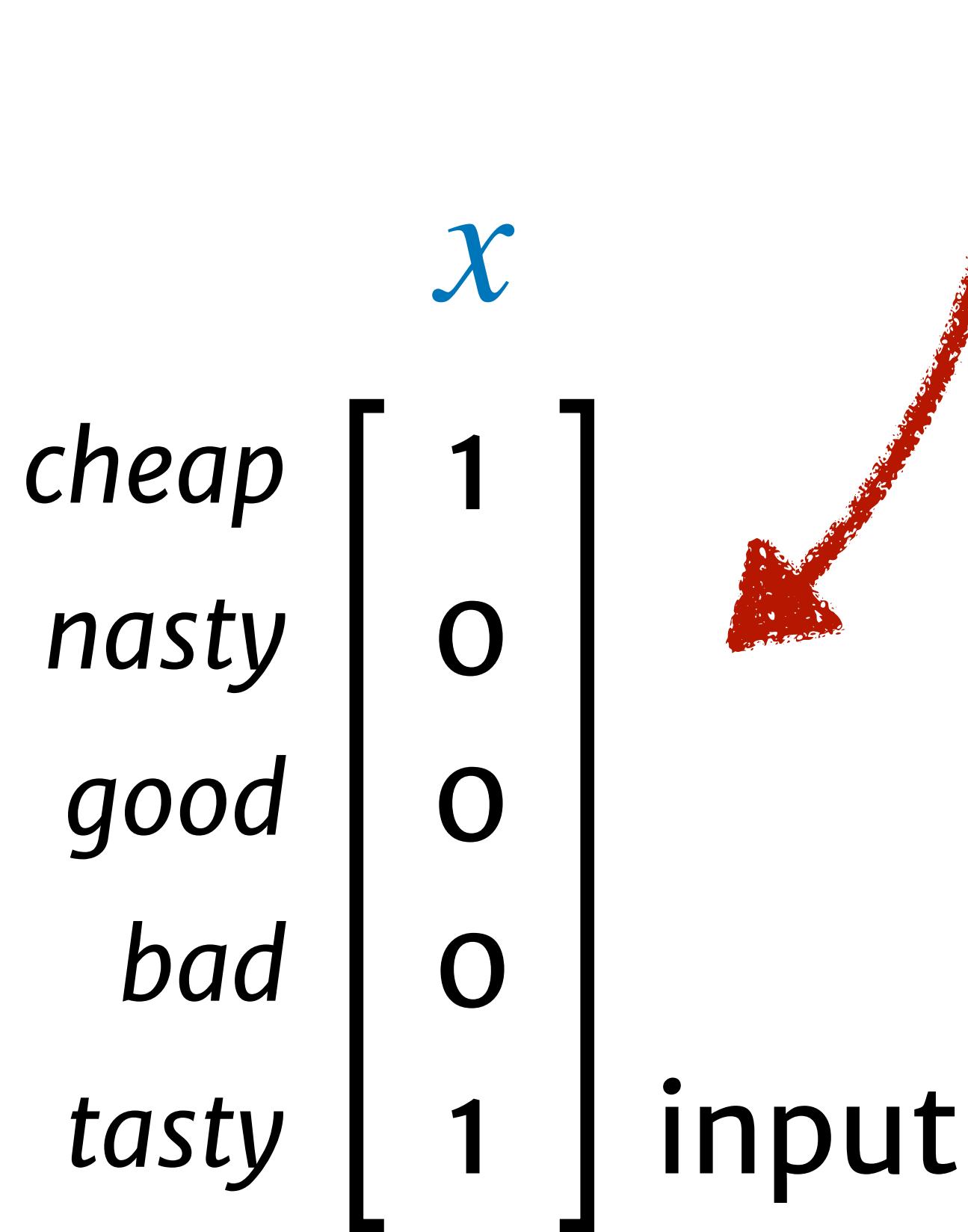
$$s = W_2^\top f(W_1^\top x)$$



# Variable-sized inputs

---

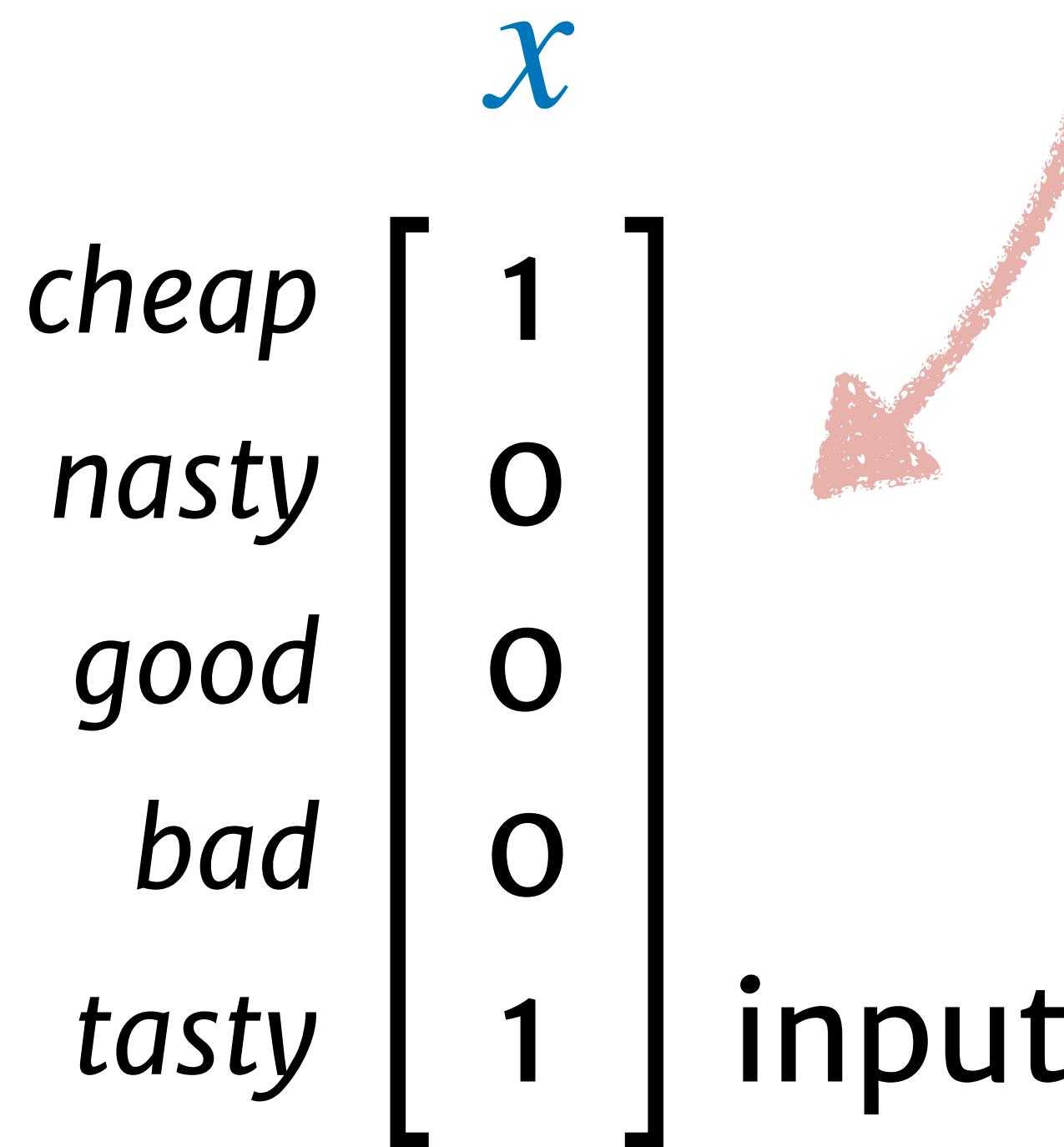
No ordering information!



# Variable-sized inputs

---

No ordering information!



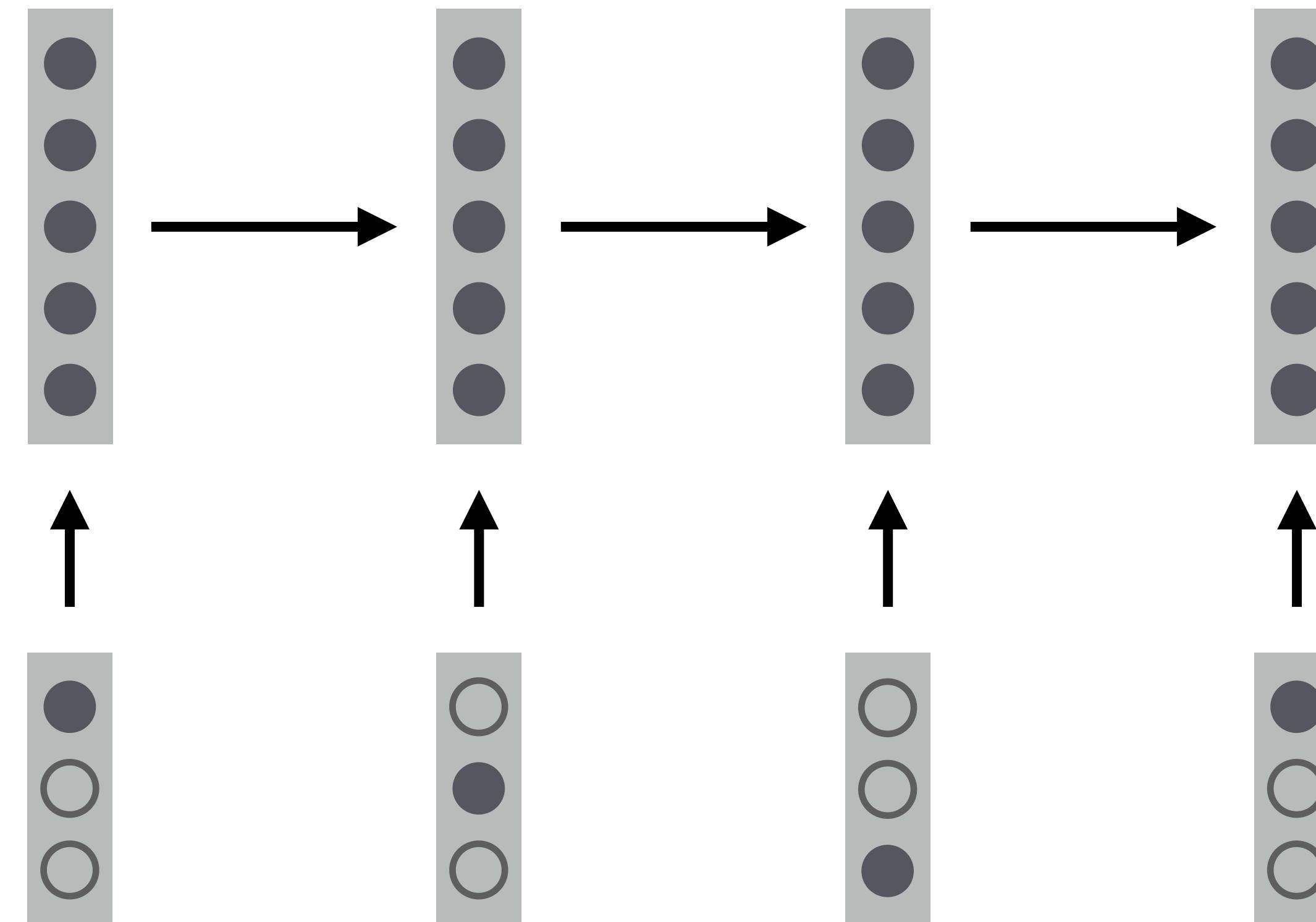
*cheap*    *and*    *very*    *tasty*

A diagram showing four vectors side-by-side, each representing a word. Above the vectors, the words *cheap*, *and*, *very*, and *tasty* are listed. Each vector is enclosed in large black brackets and has five entries. The first vector (*cheap*) has entries [1, 0, 0, 0, :]. The second vector (*and*) has entries [0, 0, 0, 0, :]. The third vector (*very*) has entries [0, 1, 0, 0, :]. The fourth vector (*tasty*) has entries [0, 0, 1, 0, :]. Ellipses (:) are used to indicate that there are more entries than shown.

No fixed input dimension!

# Recurrent neural networks

---



*cheap*

*and*

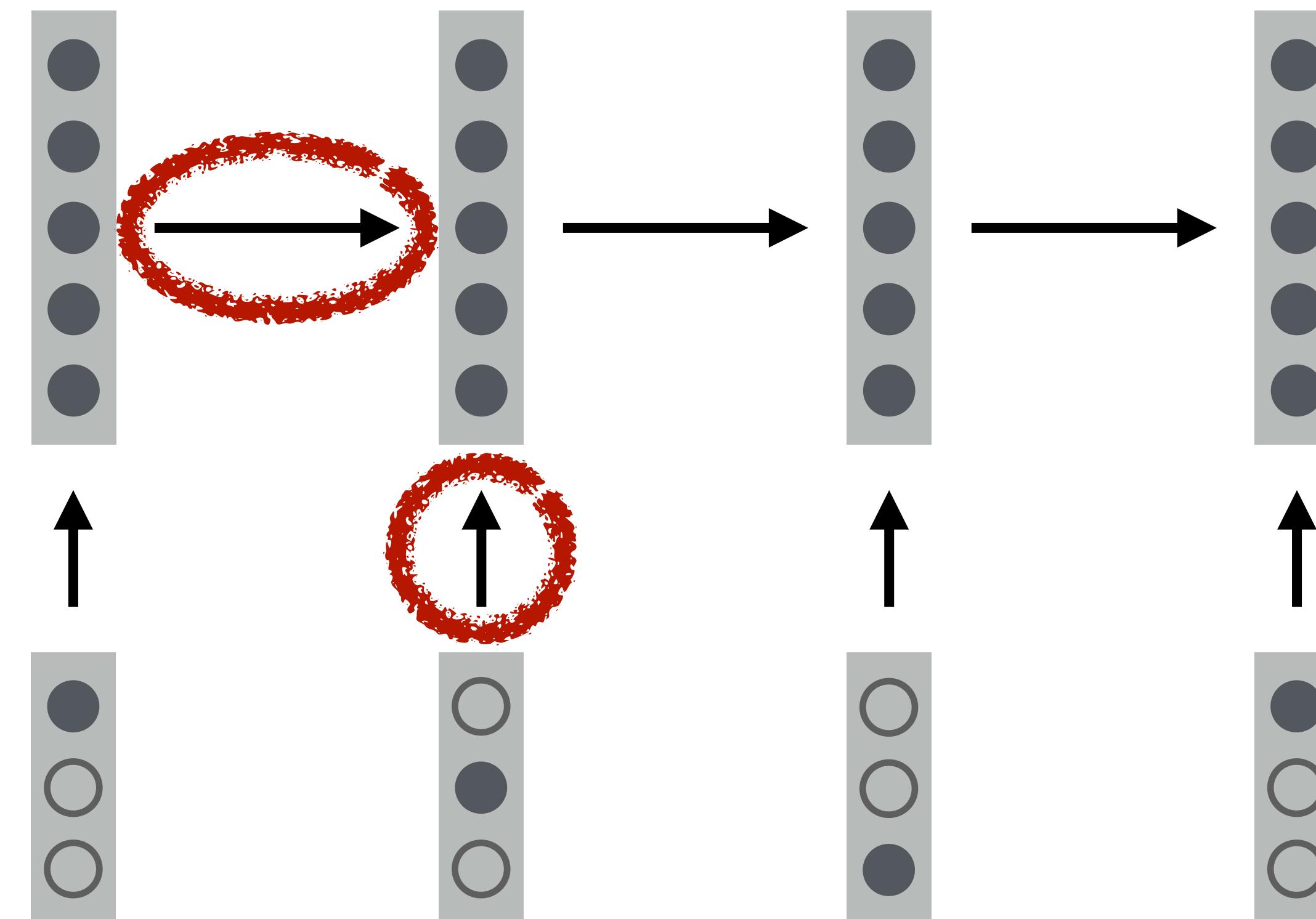
*very*

*tasty*

# Recurrent neural networks

---

Hidden states  
depend on an  
earlier state  
and an input



*cheap*

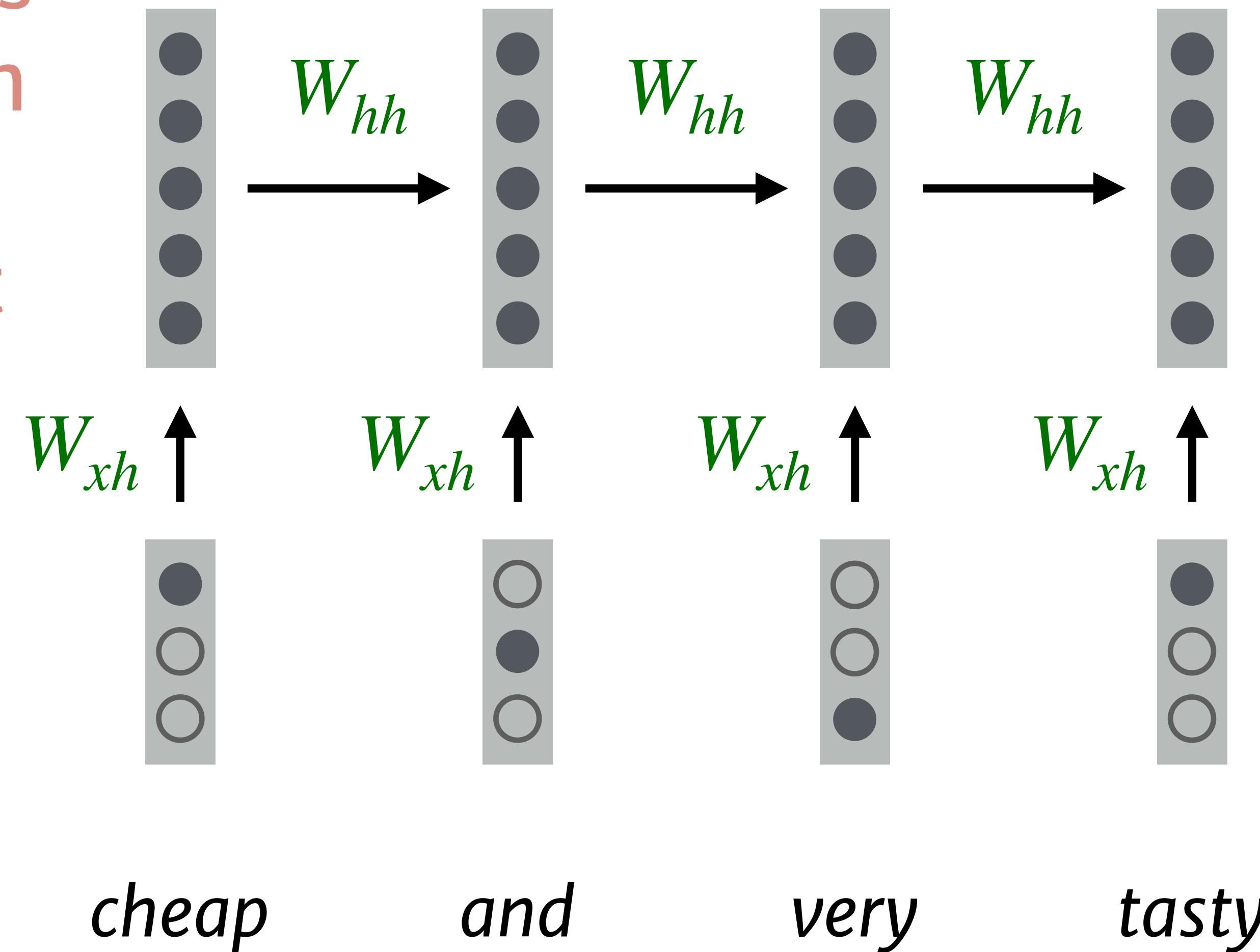
*and*

*very*

*tasty*

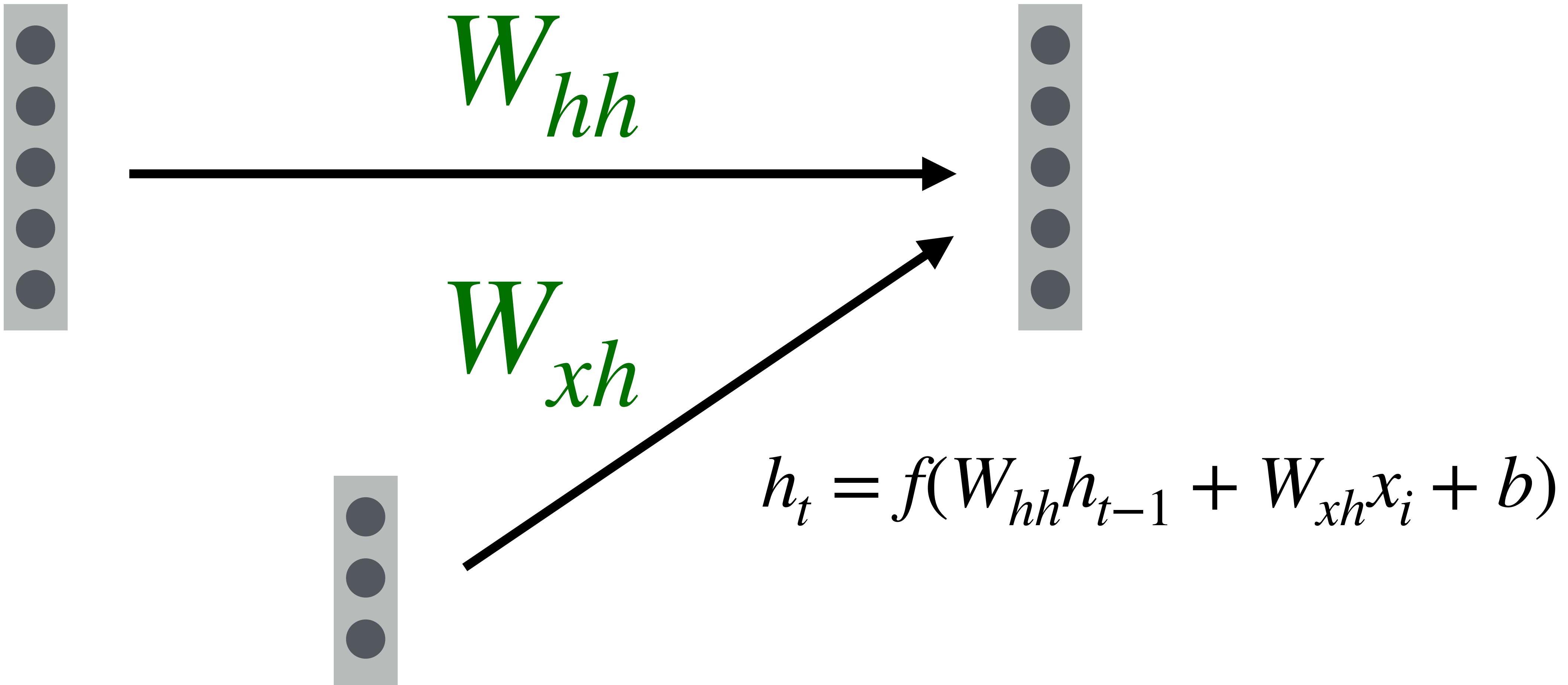
# Recurrent neural networks

Hidden states  
depend on an  
earlier state  
and an input



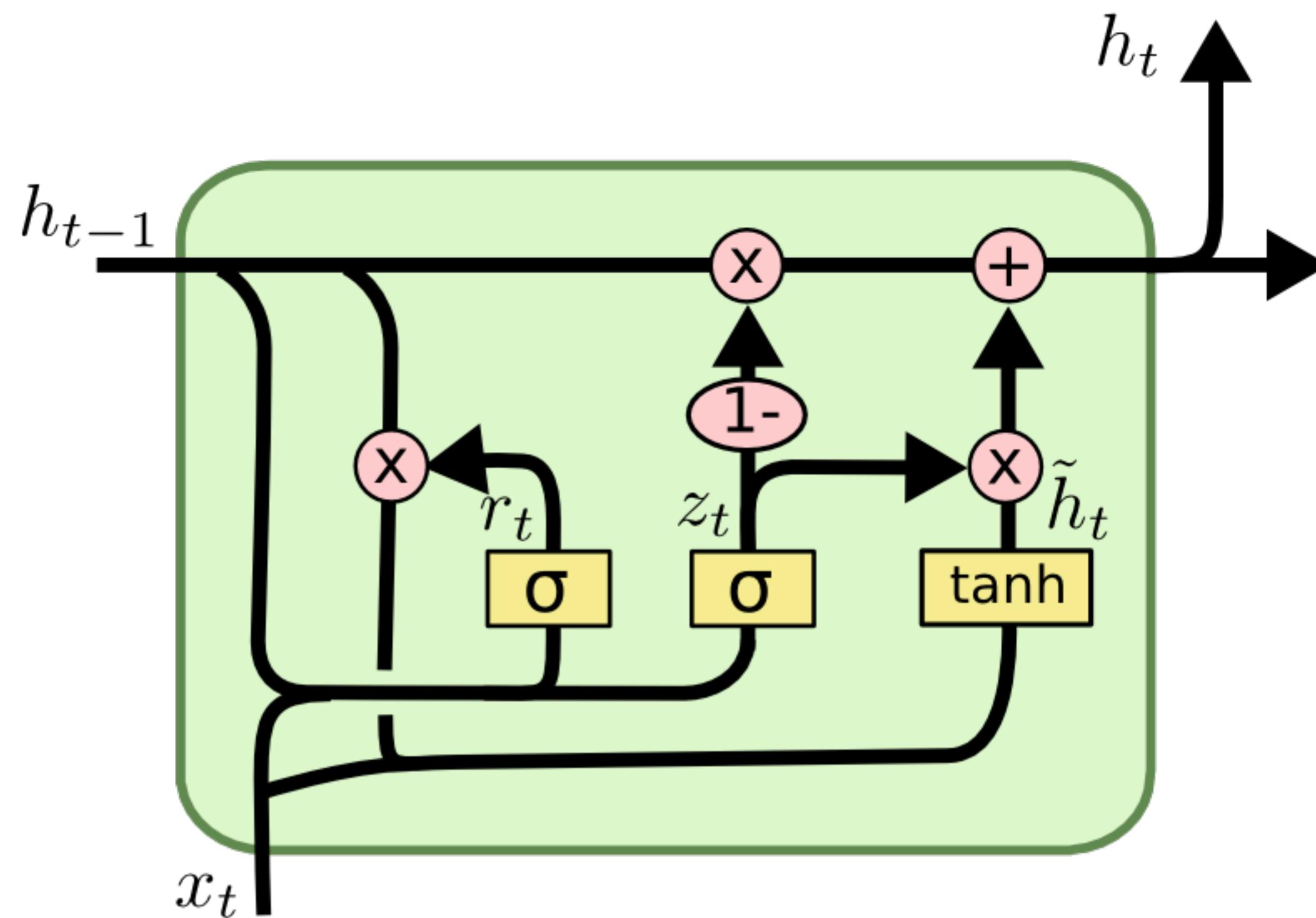
# “Vanilla” RNNs

---



# Gated Recurrent Units

---



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

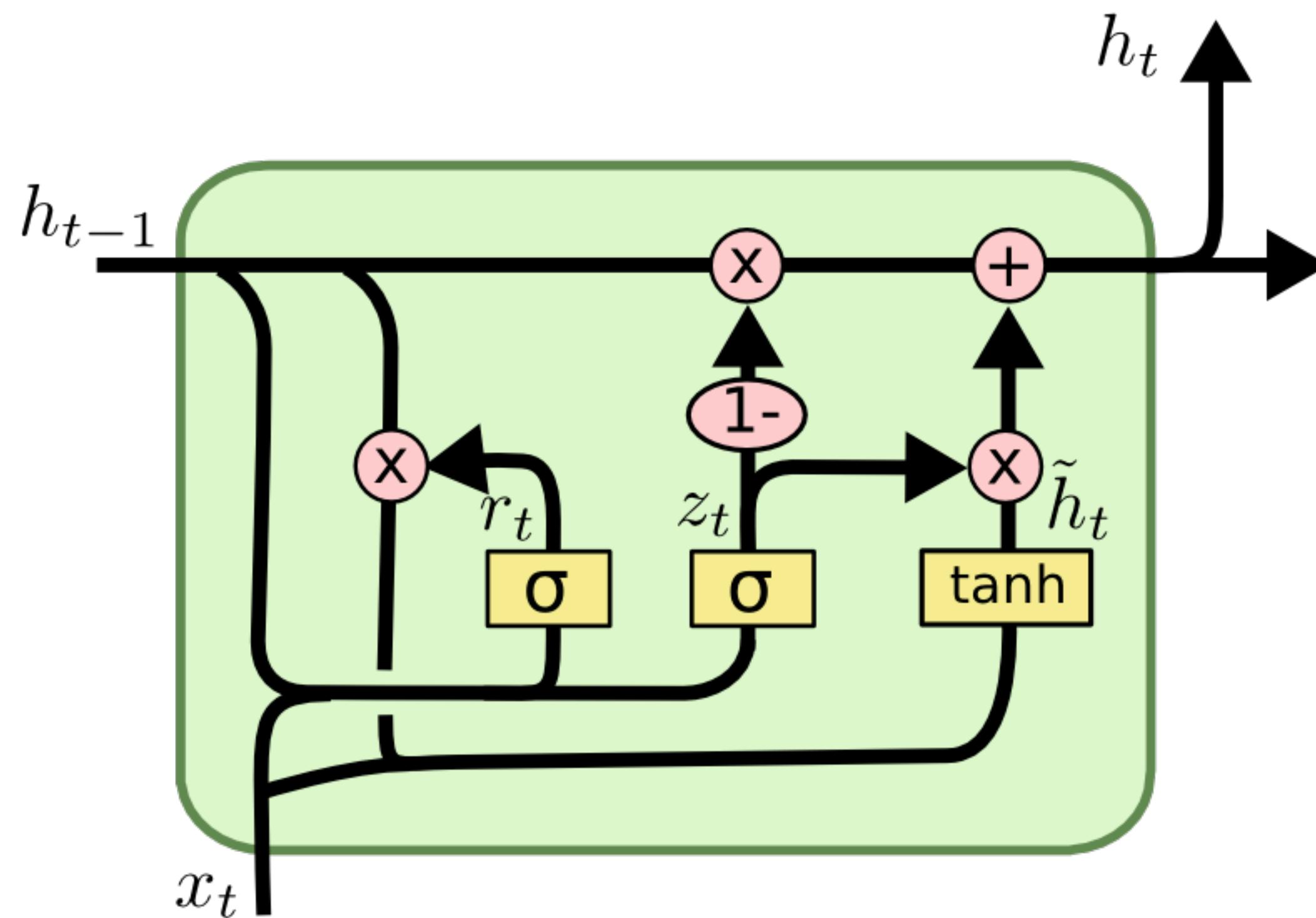
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

[Image: Cristopher Olah]

# Gated Recurrent Units



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

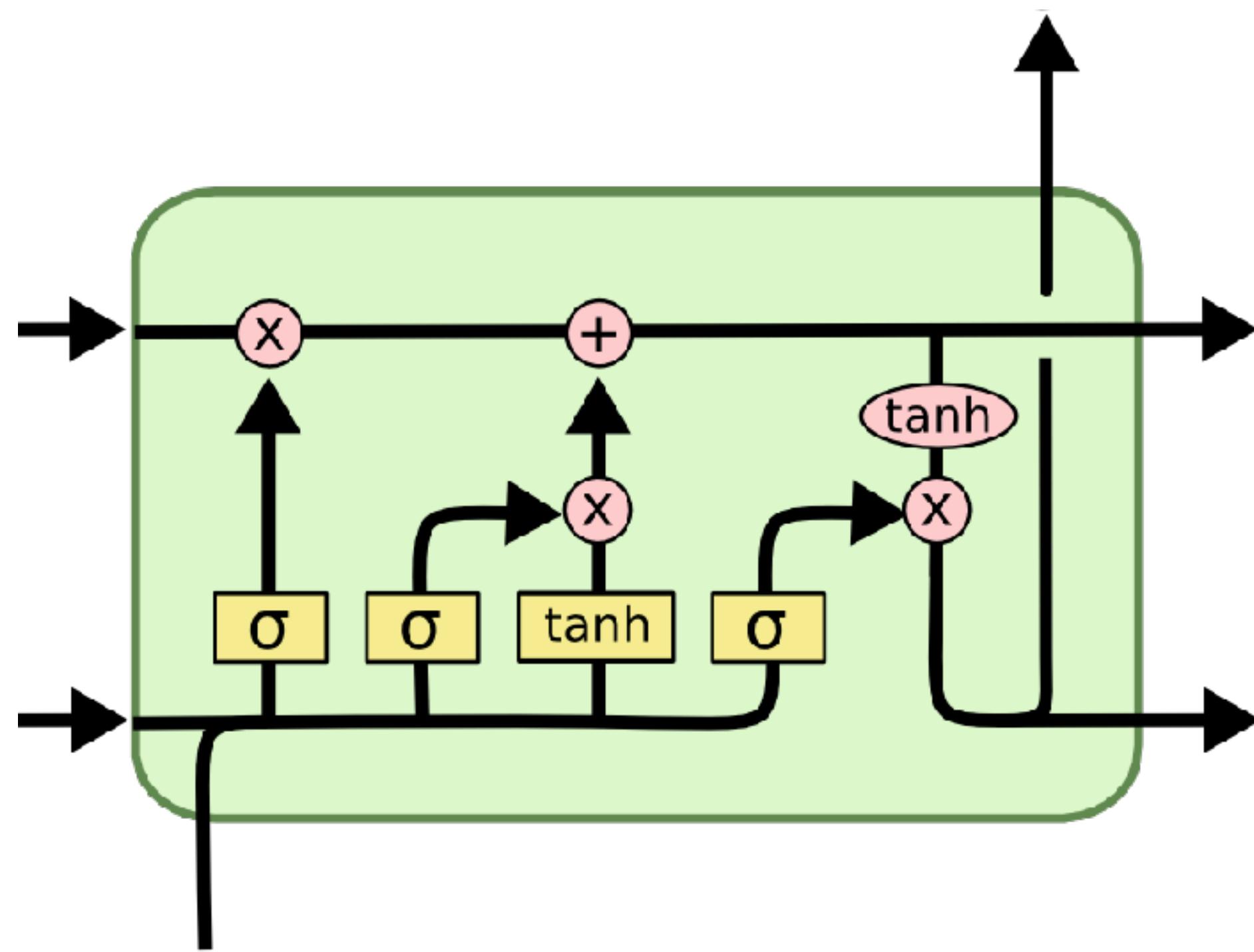
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = \boxed{(1 - z_t) * h_{t-1}} + z_t * \tilde{h}_t$$

[Image: Cristopher Olah]

# Long Short-Term Memory Units

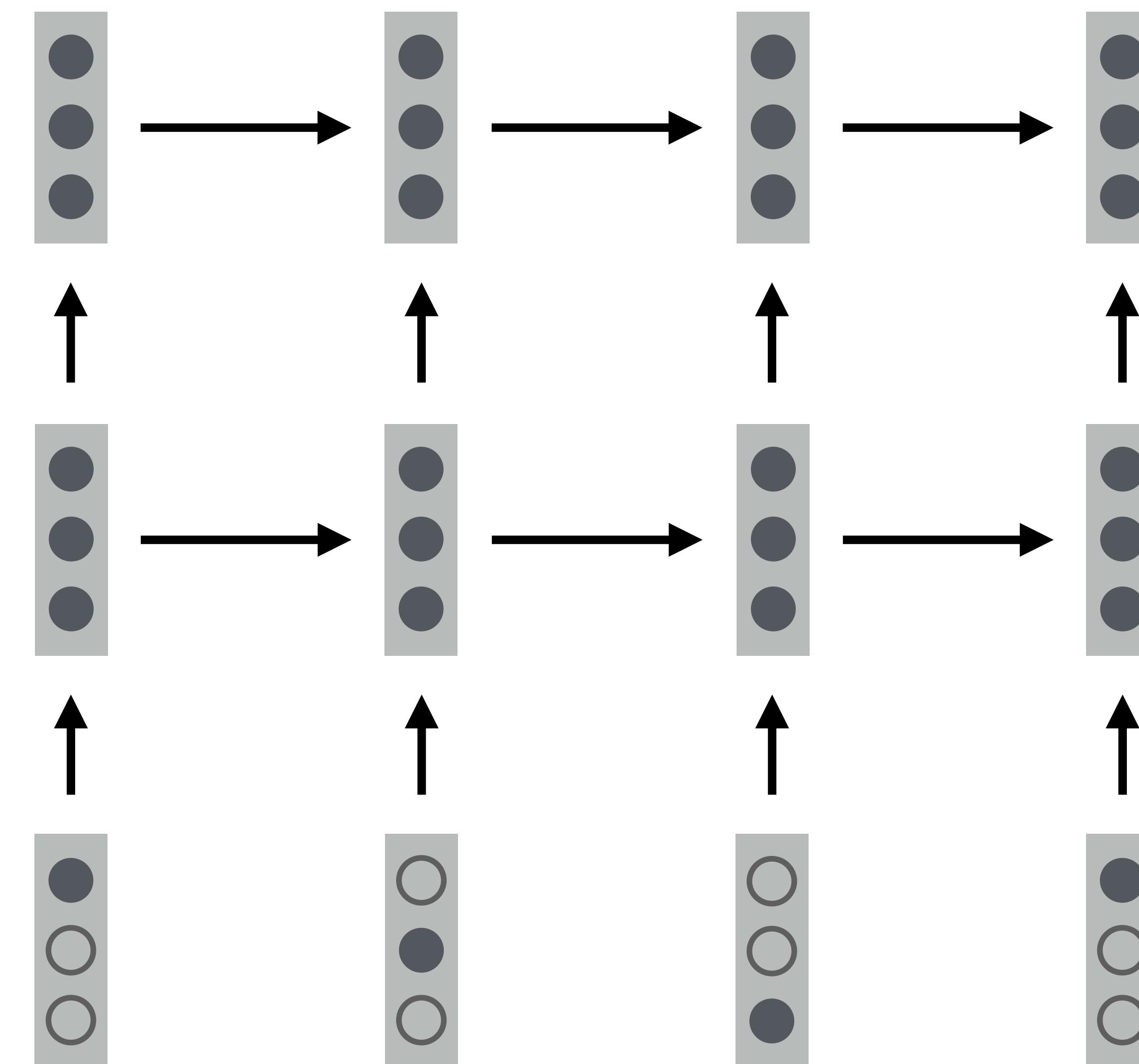
---



[Image: Christopher Olah]

# Deeper RNNs

---



*cheap*

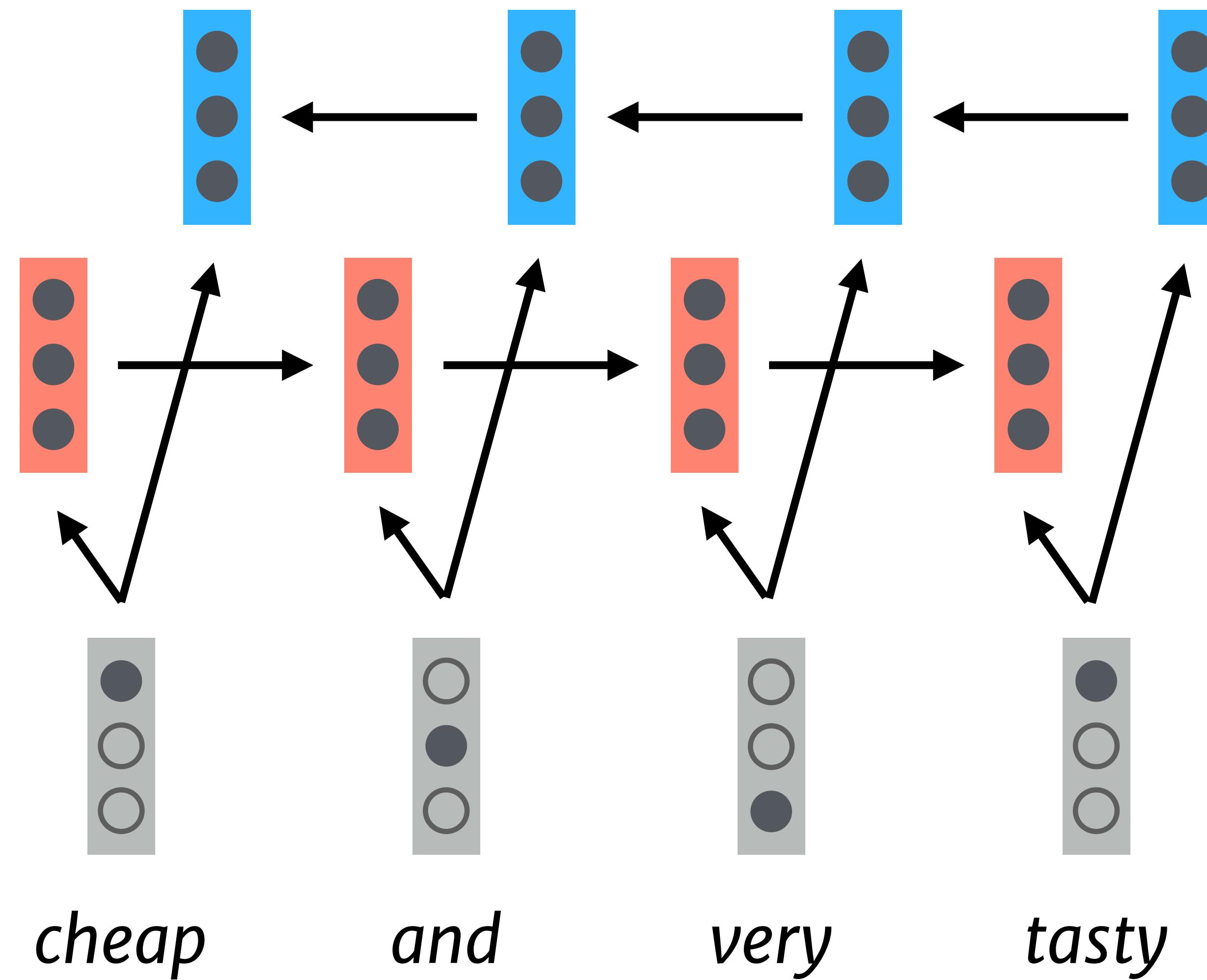
*and*

*very*

*tasty*

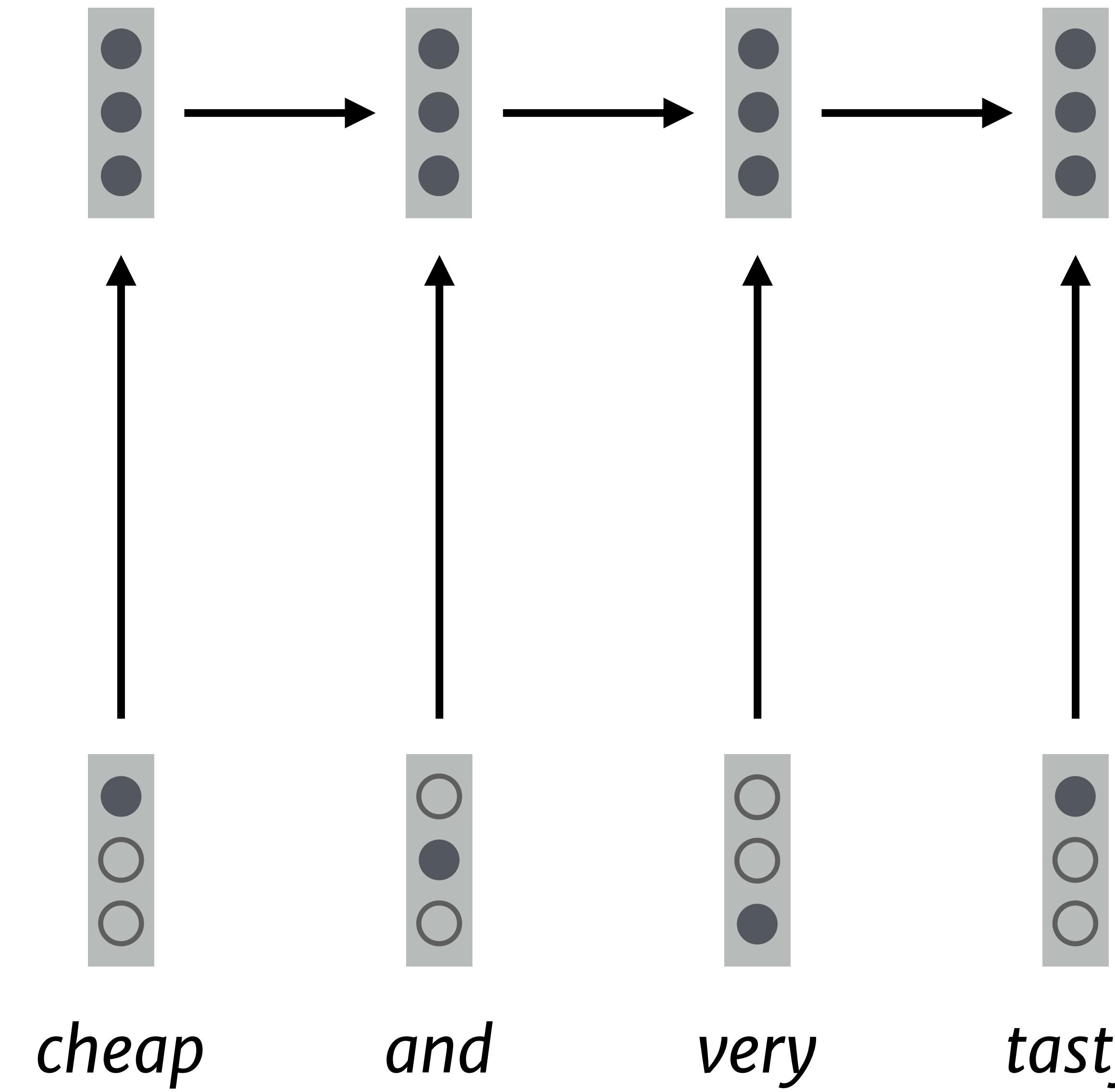
# Bidirectional RNNs

---



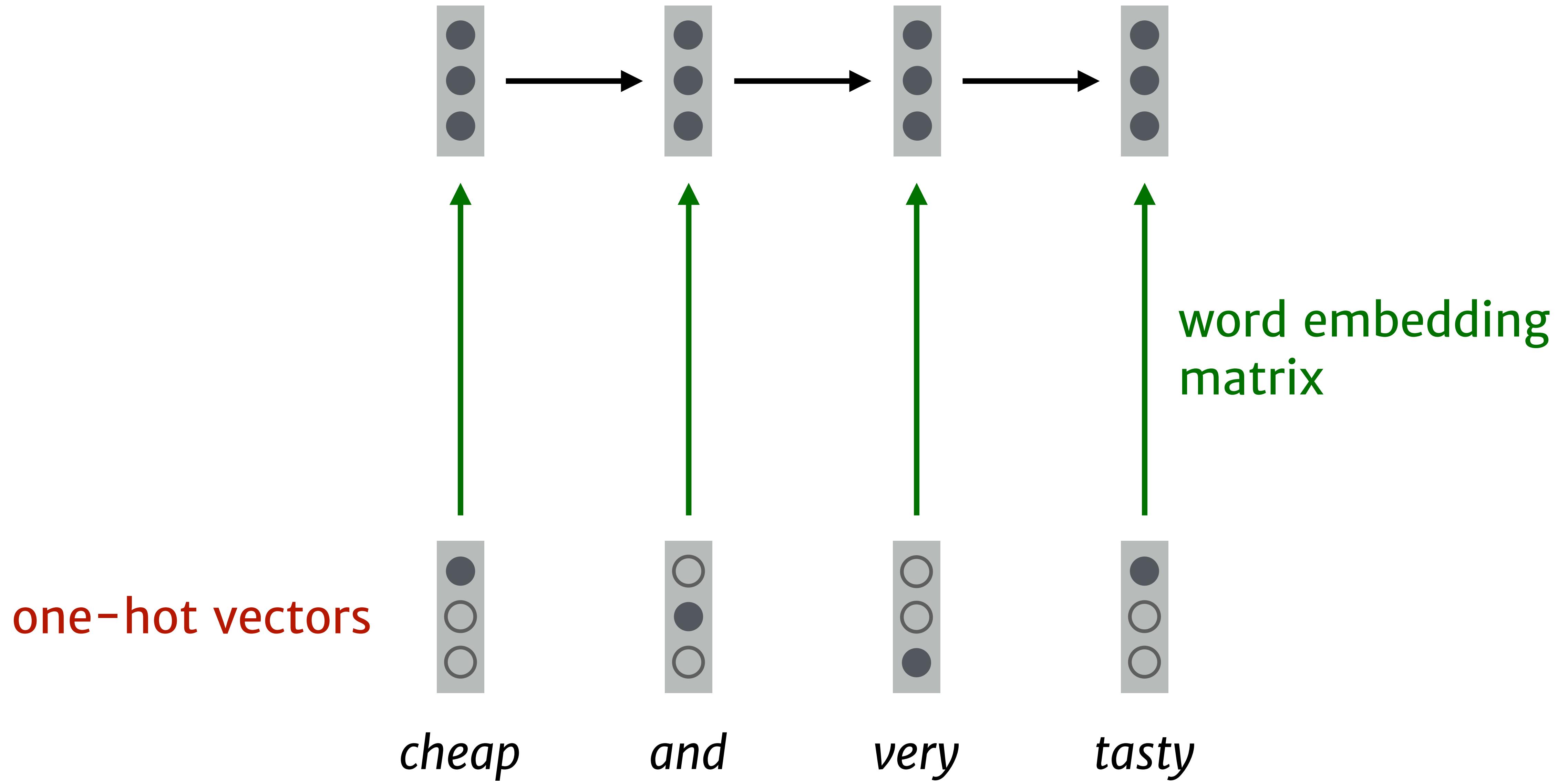
# RNNs and word embeddings

---



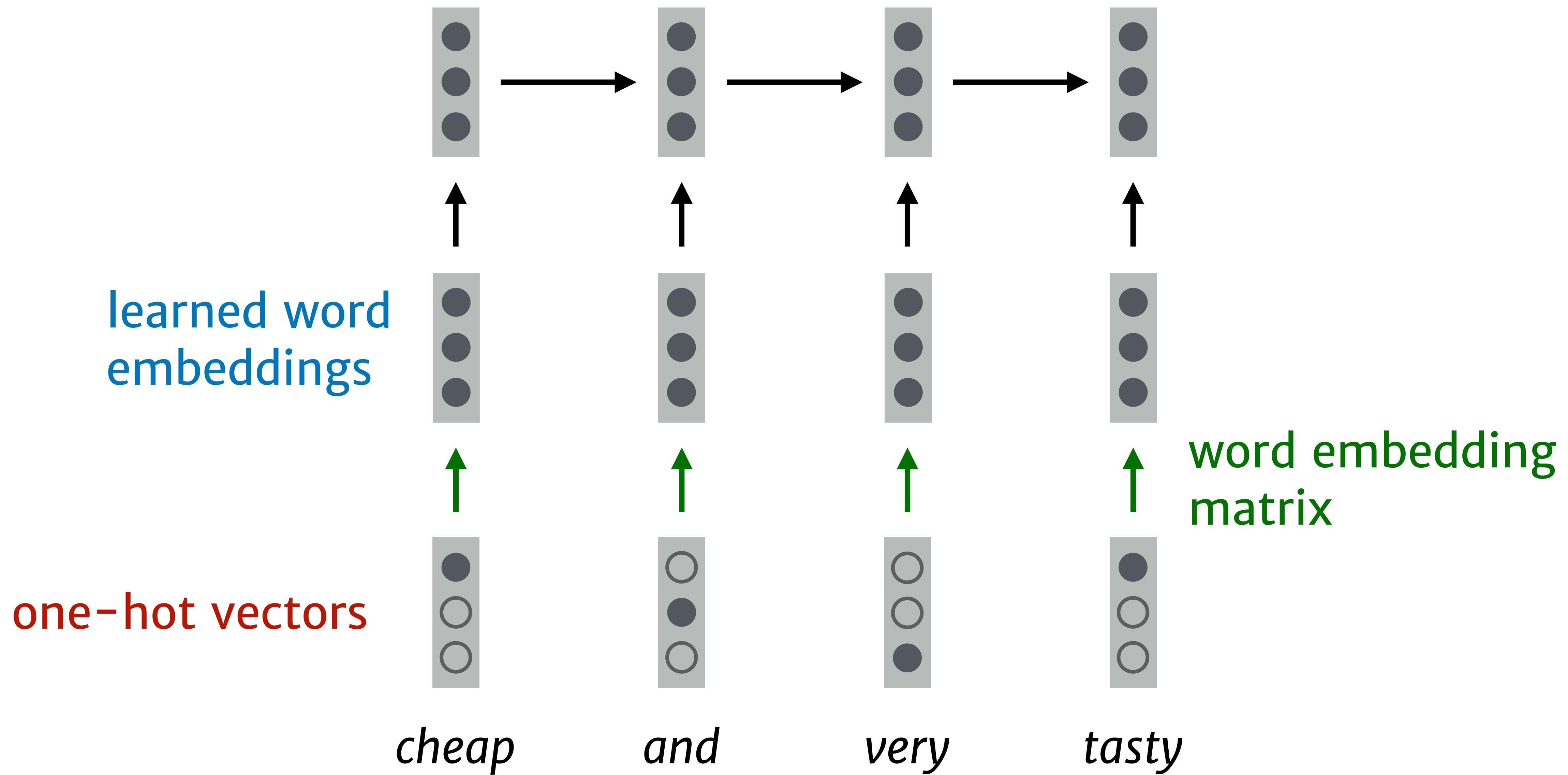
# RNNs and word embeddings

---



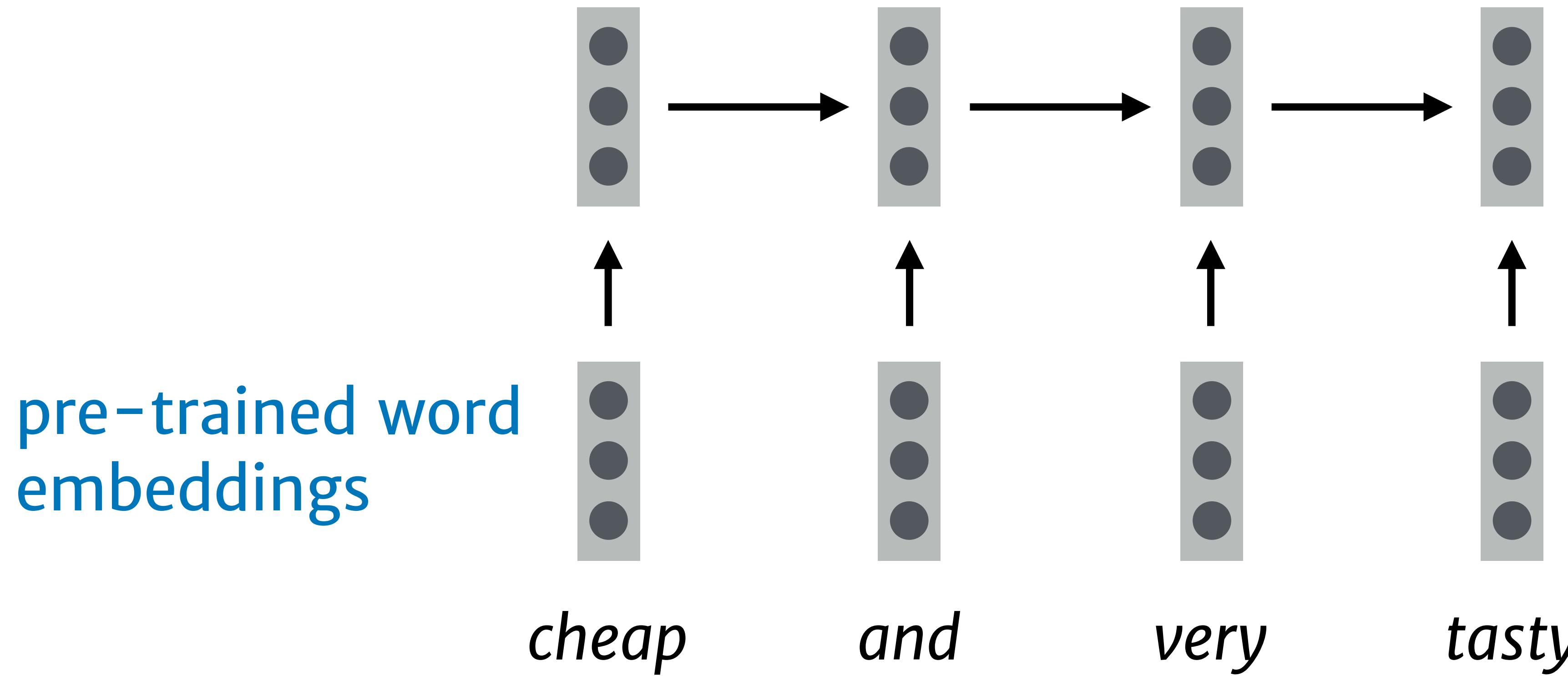
# RNNs and word embeddings

---



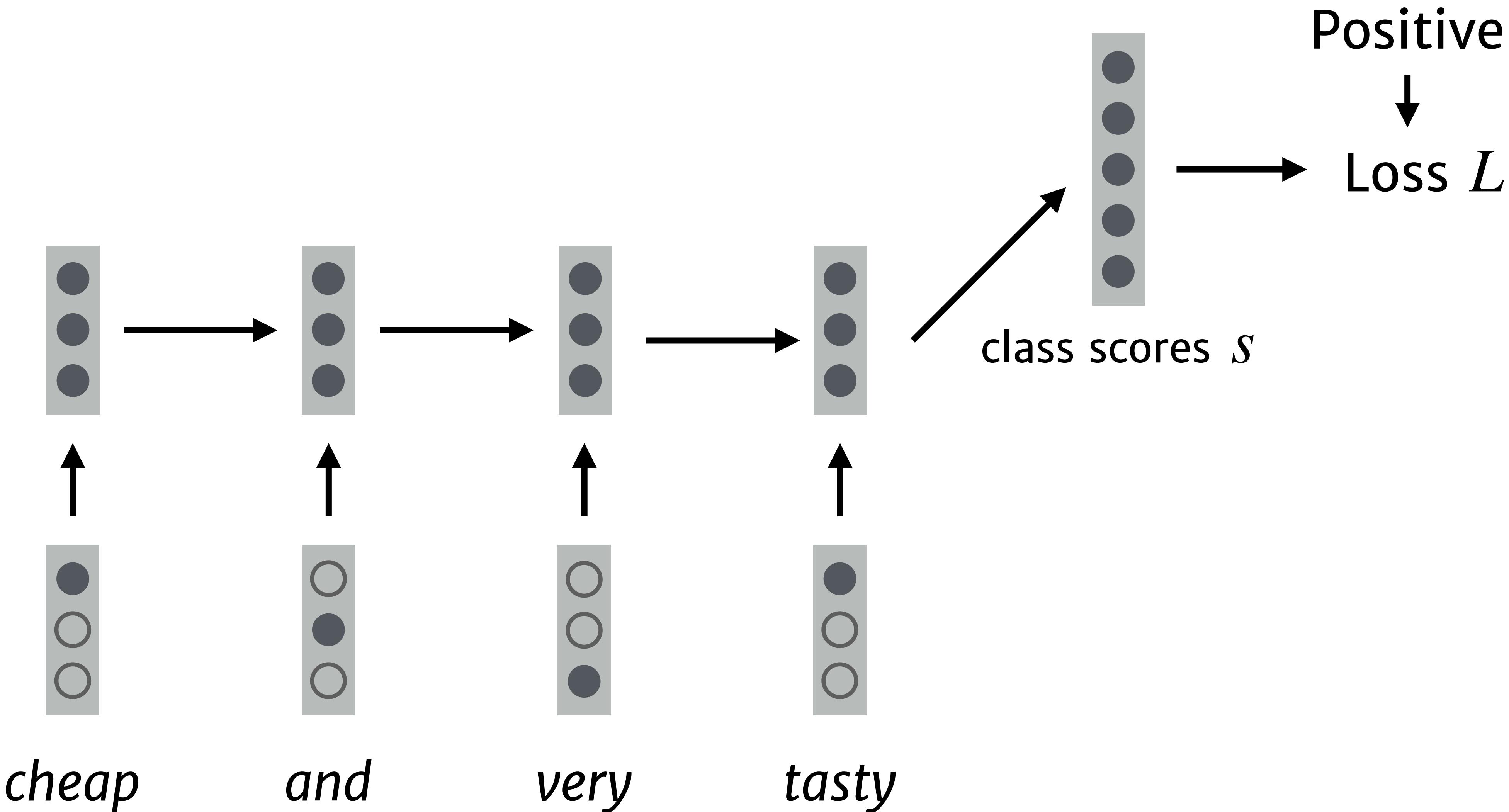
# RNNs and word embeddings

---

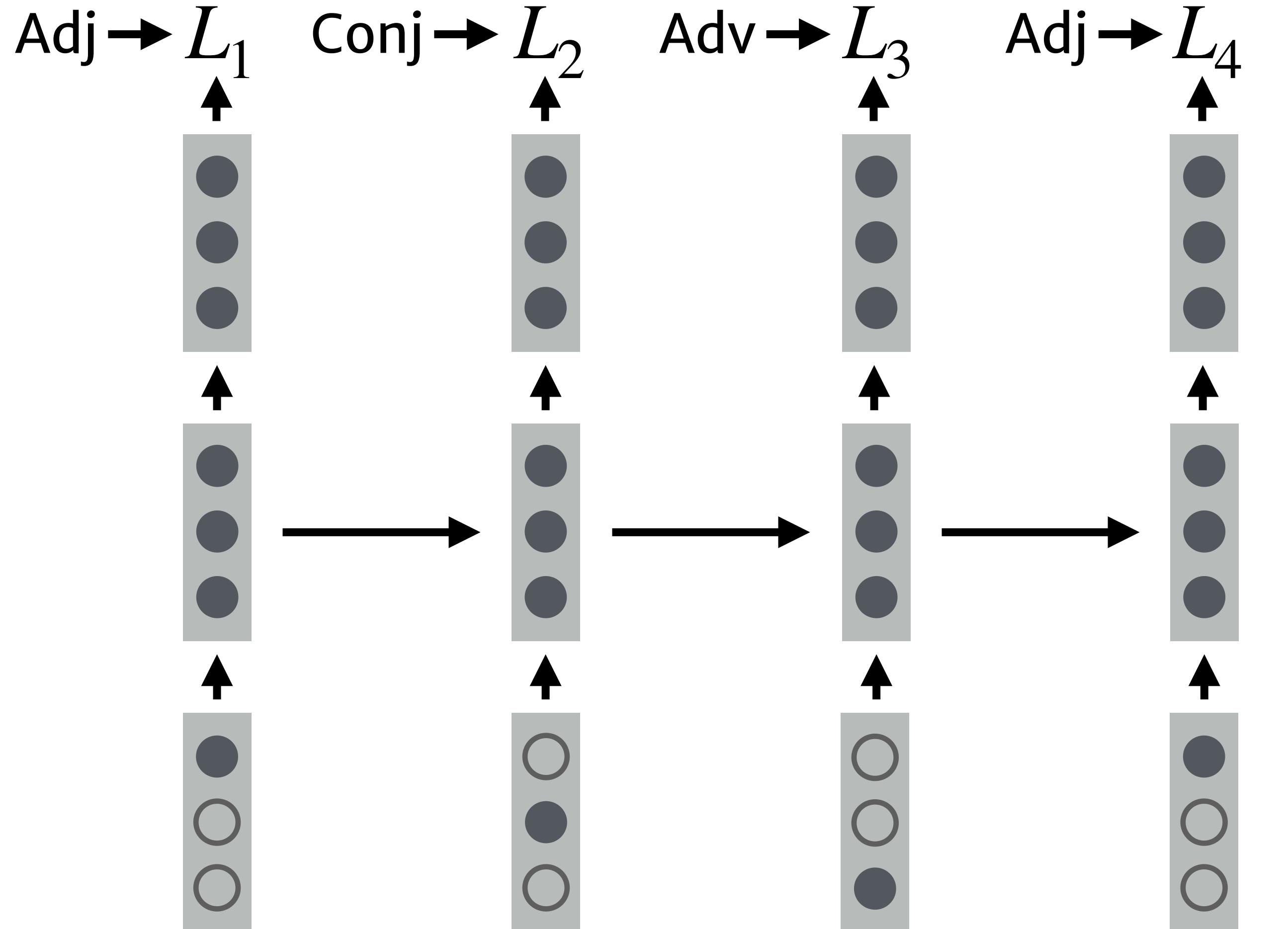


# Text classification

---



# Sequence labeling



*cheap*

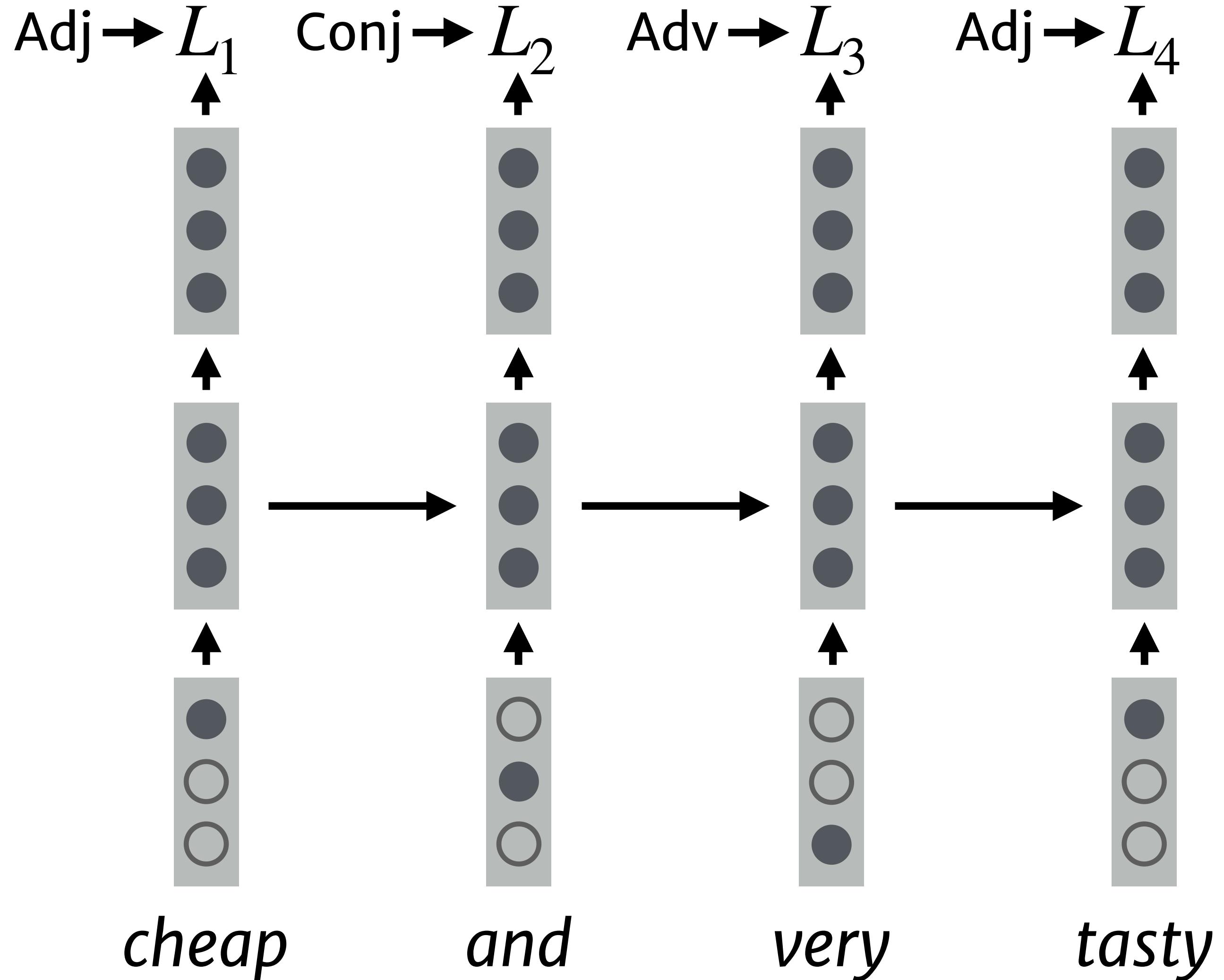
*and*

*very*

*tasty*

$$\begin{aligned} L &= \sum_t L_t \\ &= - \sum_t \log p(y_t | x_{:t}) \end{aligned}$$

# Sequence labeling



$$L = \sum_t L_t$$

e.g.

$$= - \sum_t \log p(y_t | x_{:t})$$

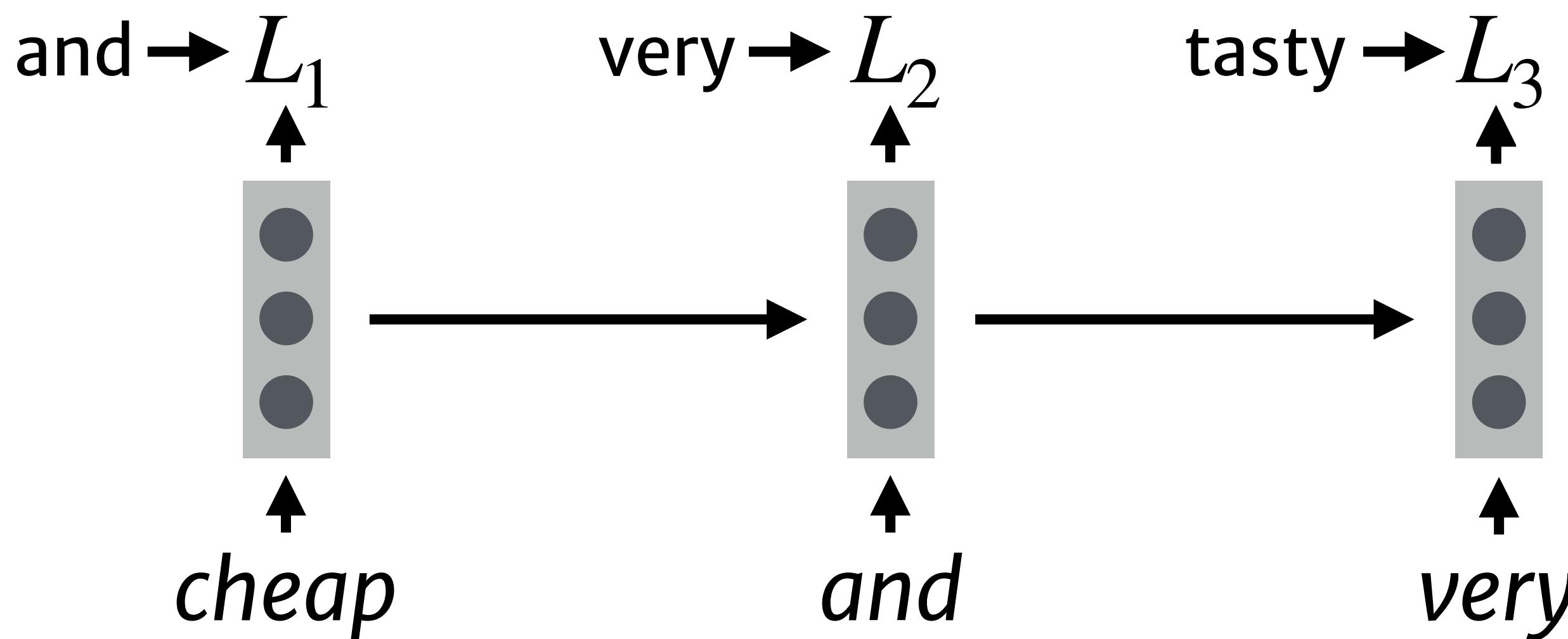
product of indep.  
conditionals!

# Langauge Modeling

---

A (unidirectional) RNN can compute  $p(y_t \mid x_{:t})$ .

Suppose for a sequence  $\mathcal{X}$  we set  $y_t = x_{t+1}$ .



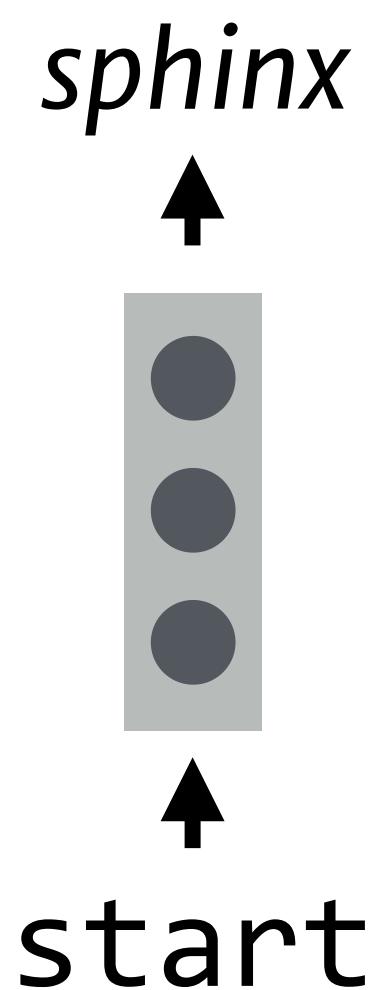
then 
$$\sum_t \log p(x_{t+1} \mid x_{:t}) = p(x)$$

# Language modeling: sampling

---

How do we sample from  $p(x)$ ?

$$x_1 \sim p(\cdot | \text{start})$$



# **SPHINX OF BLACK QUARTZ, JUDGE MY VOW**



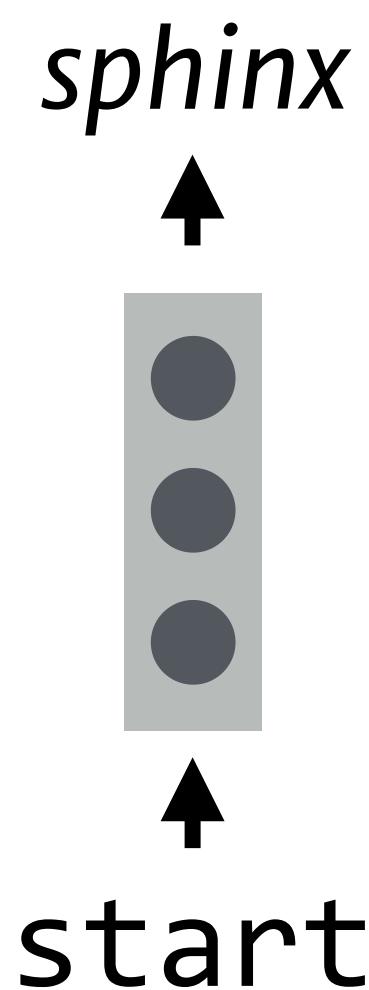
[Image: [egyptianmarketplace.com](http://egyptianmarketplace.com)]

# Language modeling: sampling

---

How do we sample from  $p(x)$ ?

$$x_1 \sim p(\cdot | \text{start})$$

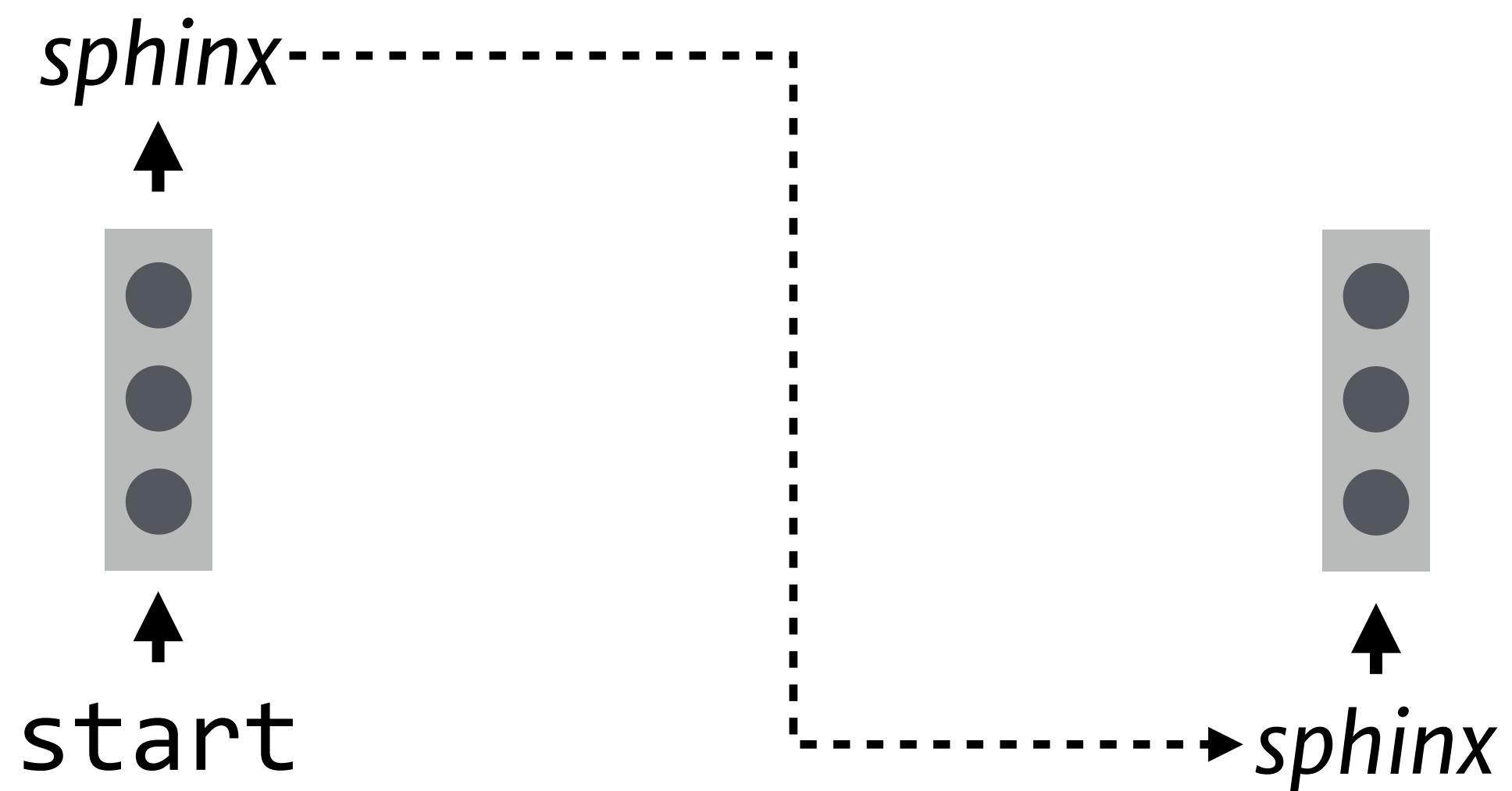


# Language modeling: sampling

---

How do we sample from  $p(x)$ ?

$$x_1 \sim p(\cdot | \text{start})$$



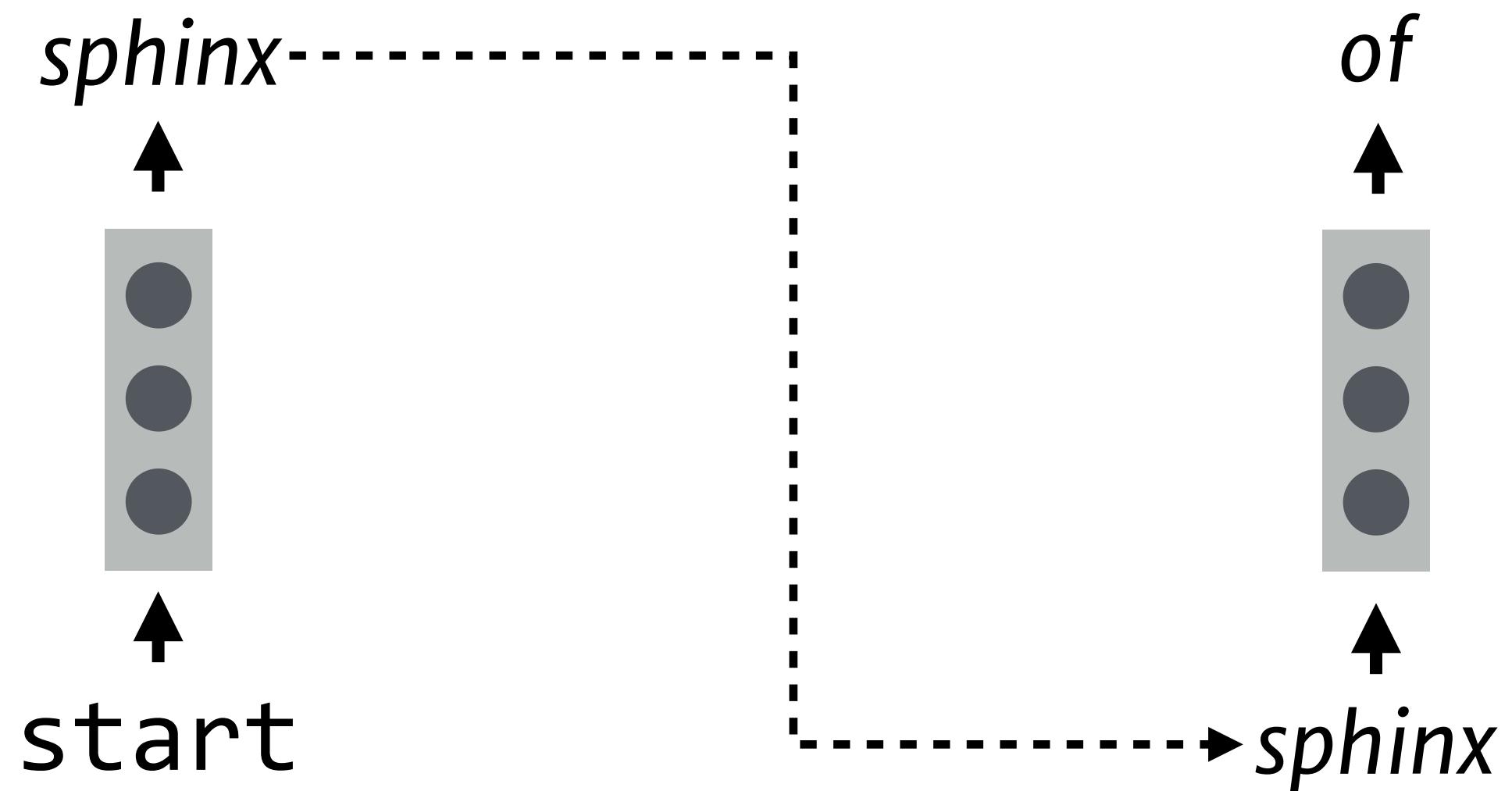
# Language modeling: sampling

---

How do we sample from  $p(x)$ ?

$$x_1 \sim p(\cdot | \text{start})$$

$$x_2 \sim p(\cdot | \text{sphinx})$$



# Language modeling: sampling

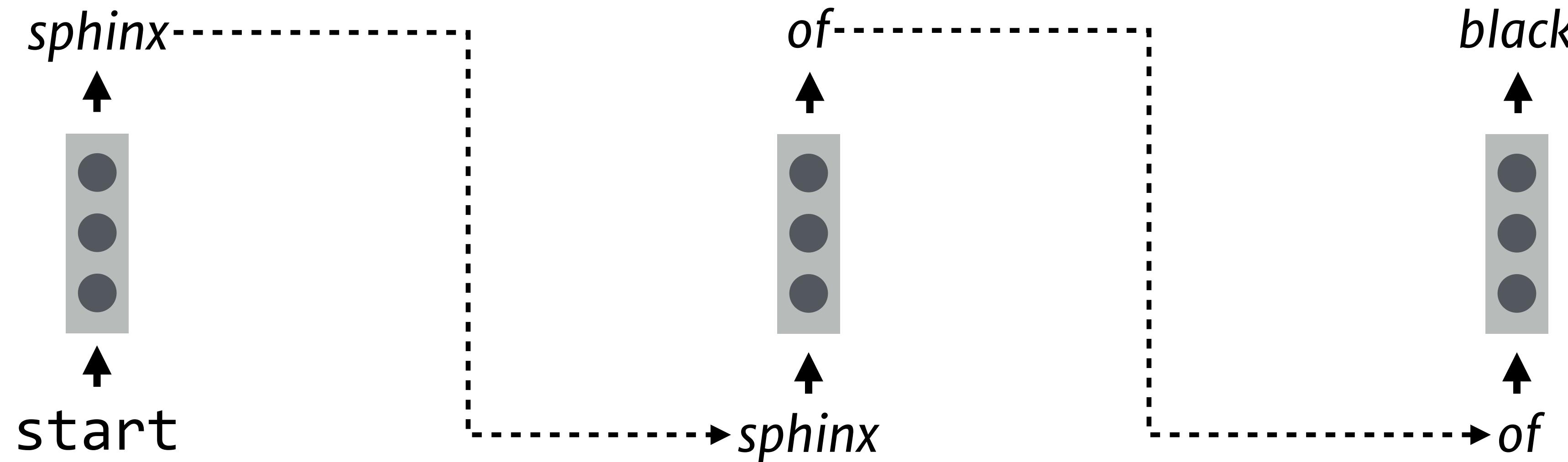
---

How do we sample from  $p(x)$ ?

$$x_1 \sim p(\cdot | \text{start})$$

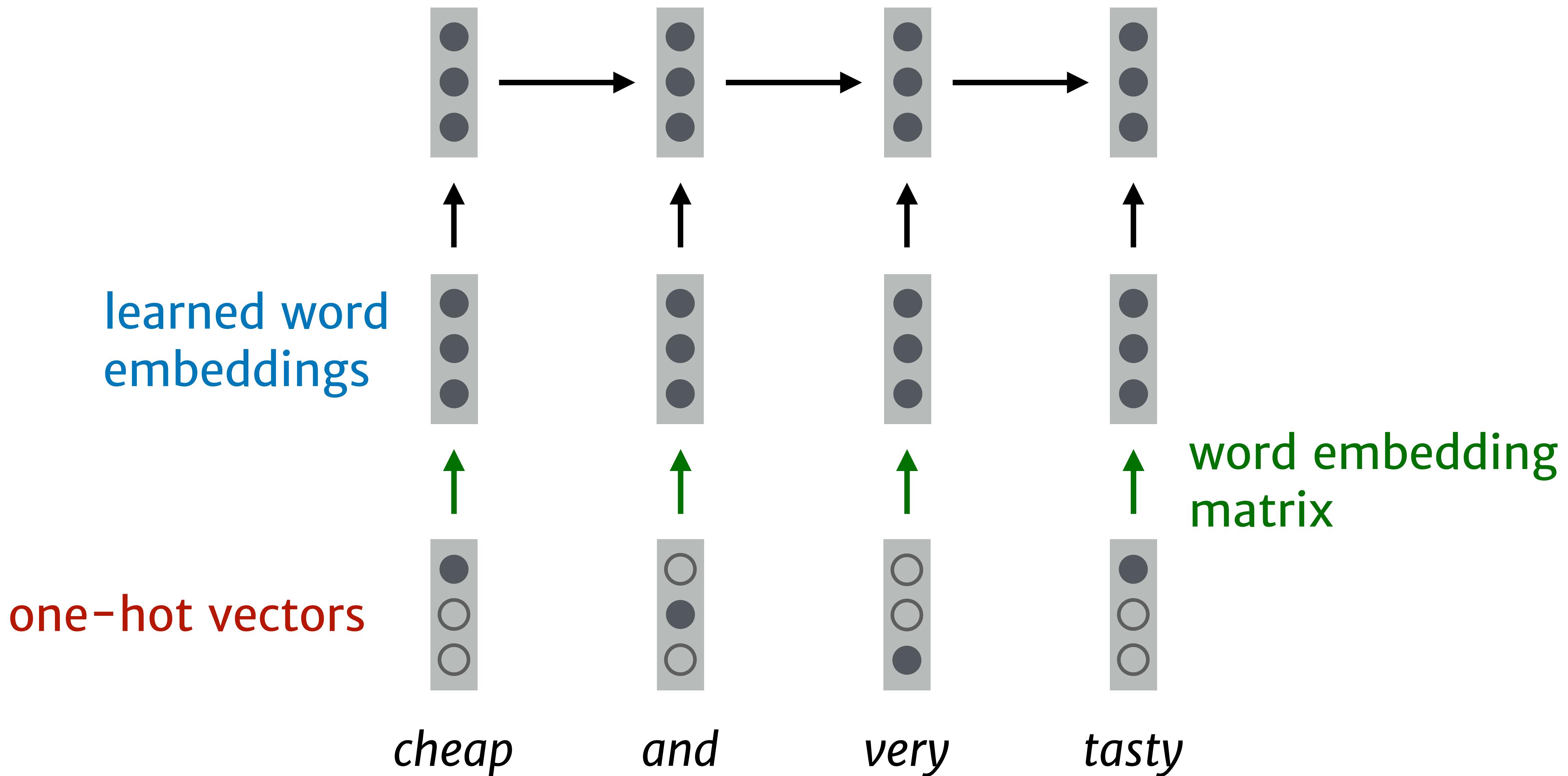
$$x_2 \sim p(\cdot | \text{sphinx})$$

$$x_3 \sim p(\cdot | \text{sphinx of})$$



# Language modeling as representation learning

---

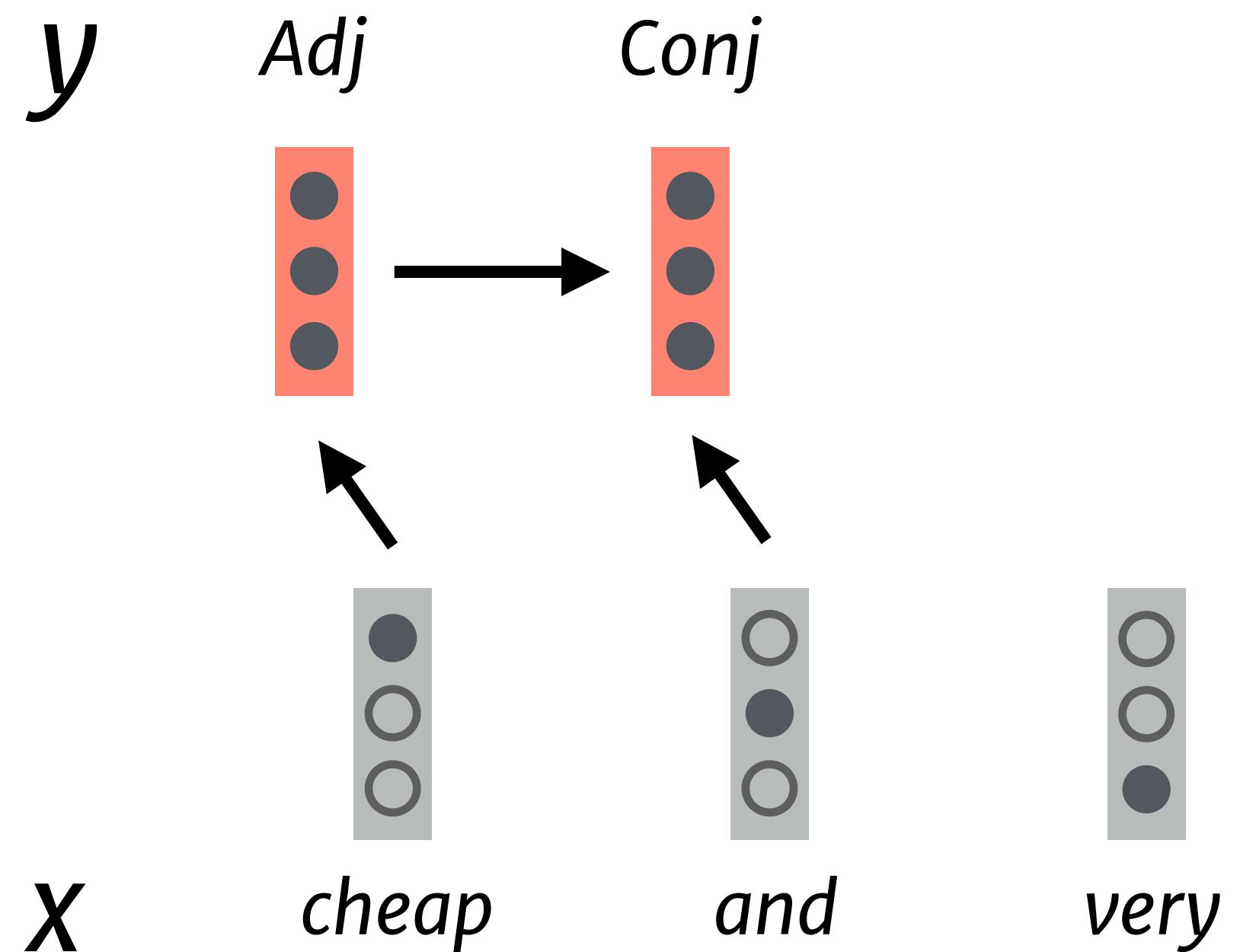


# RNNs as Markov models

---

I can train this network to predict:

$$\log p(y_t | x_{:t})$$

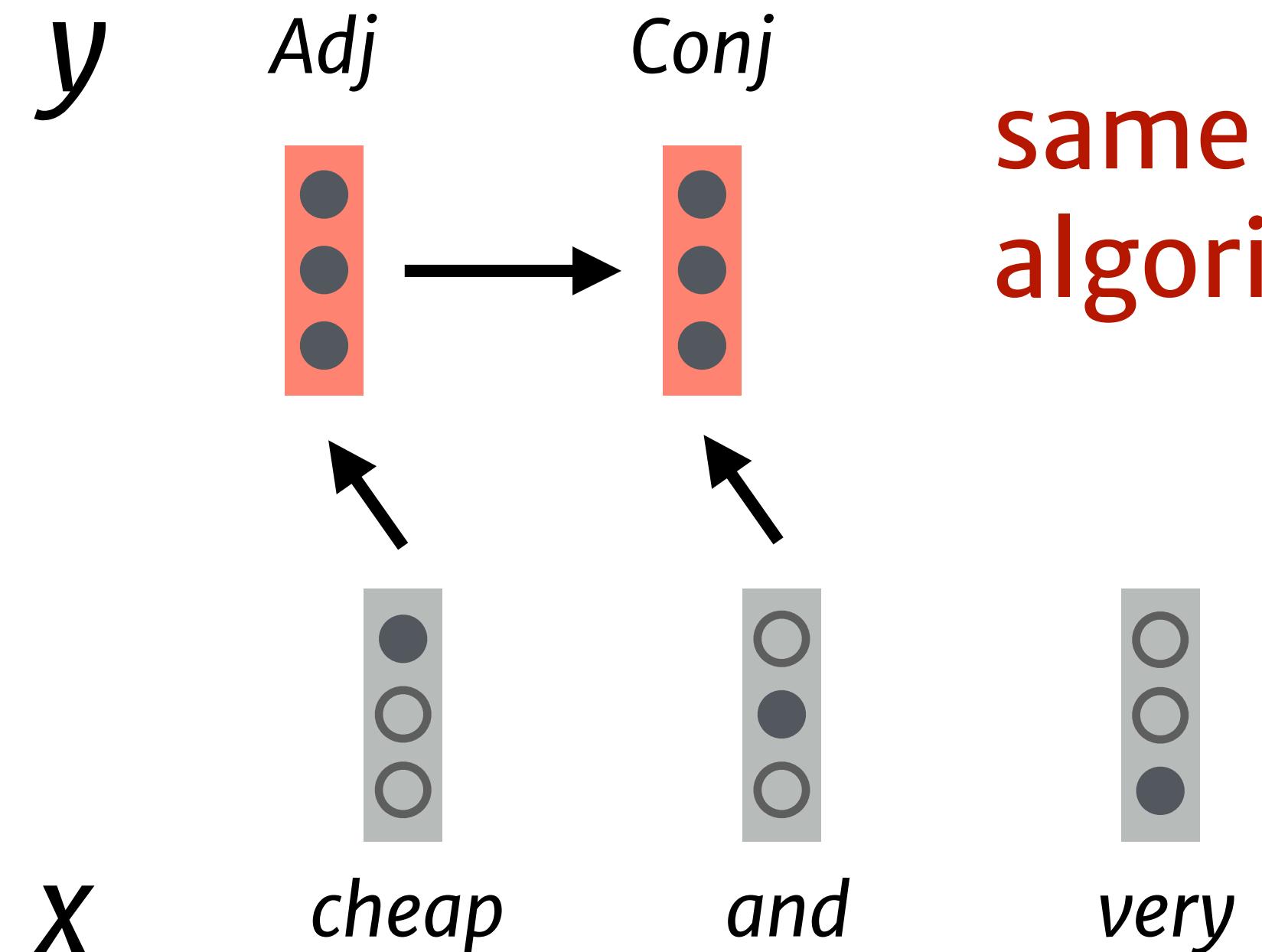


# RNNs as Markov models

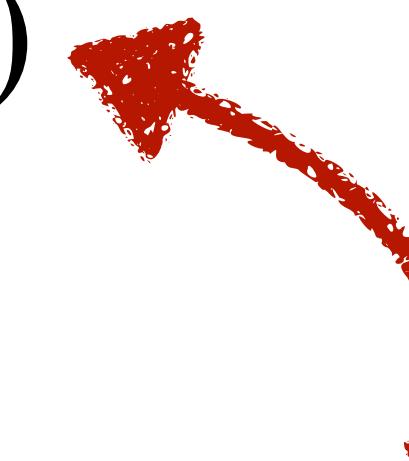
---

I can train this network to predict:

$$\log p(y_t | x_{:t}) = p(q_t | O_{:t})$$



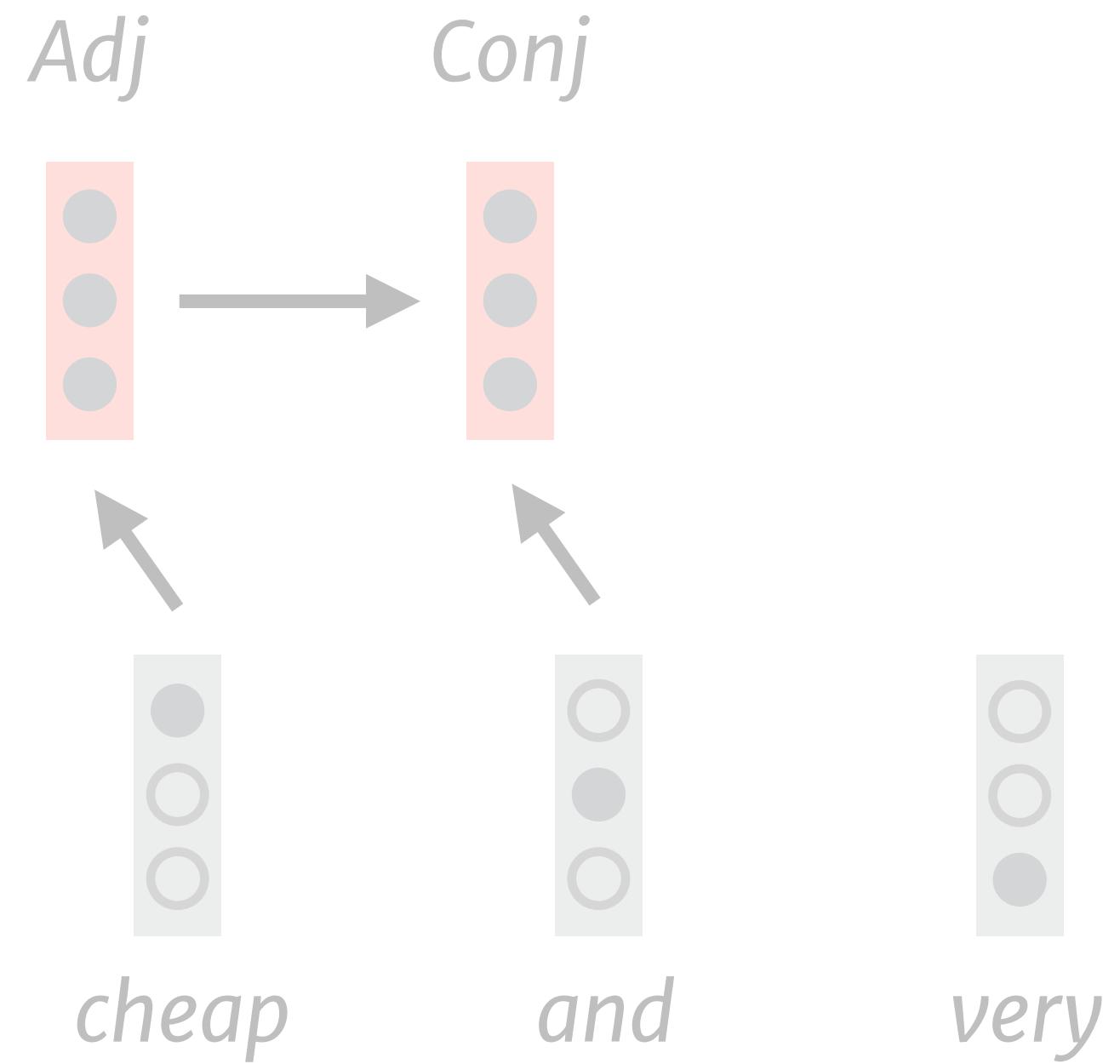
same as forward  
algorithm!



# RNNs as Markov models

I can train this network to predict:

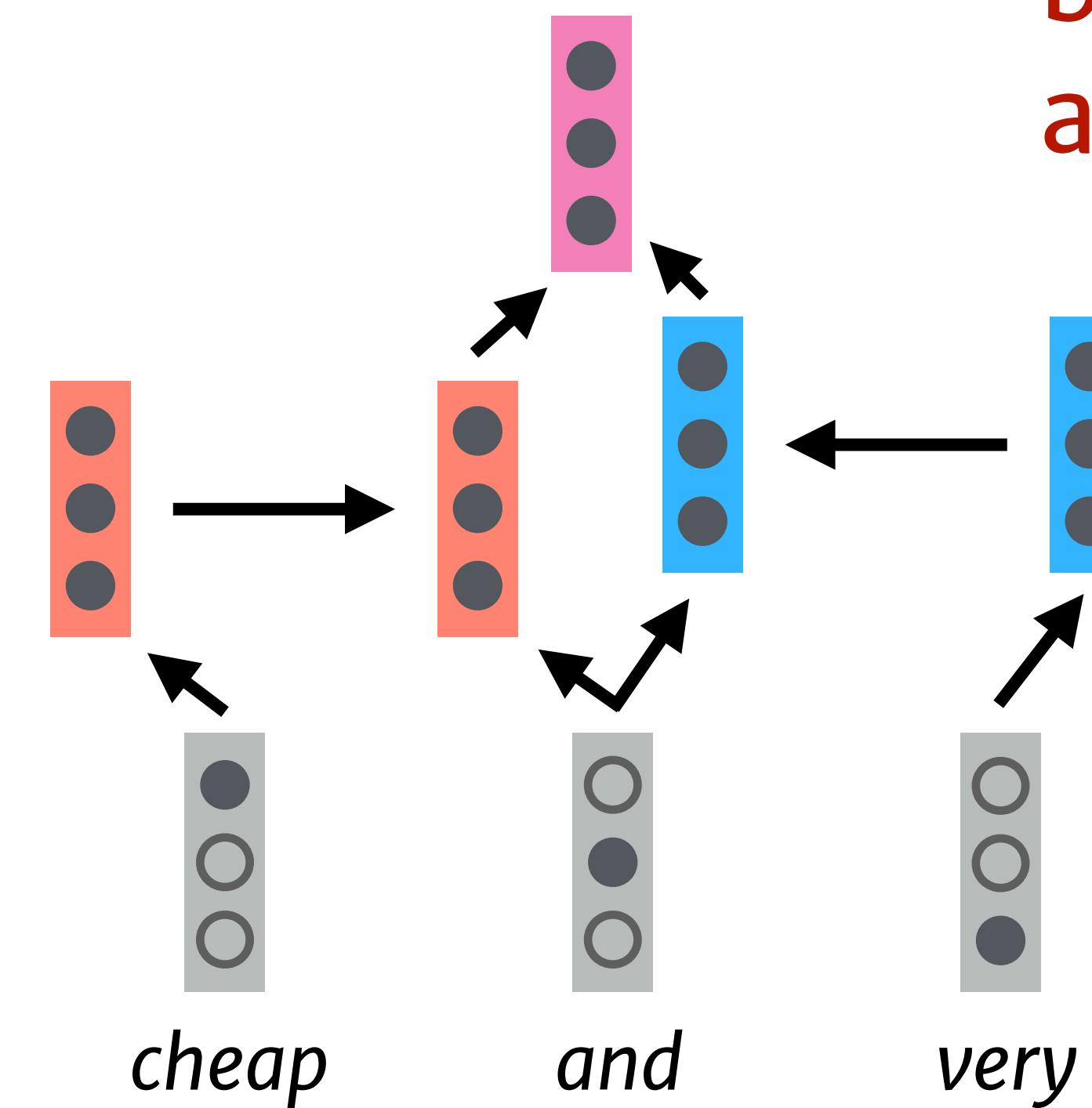
$$\log p(y_t | x_{:t}) = p(q_t | O_{:t})$$



I can train this network to predict:

$$\log p(y_t | x) = p(q_t | O)$$

**forward-backward algo!**



# Sequence-to-sequence models

# A dataset of math problems

---

*One plus one equals two.*

*Two times two equals four.*

*Seven is prime.*

*One plus two times three equals seven.*

# A dataset of math problems

---

*One plus one equals two.*

*Two times two equals four.*

*Seven is prime.*

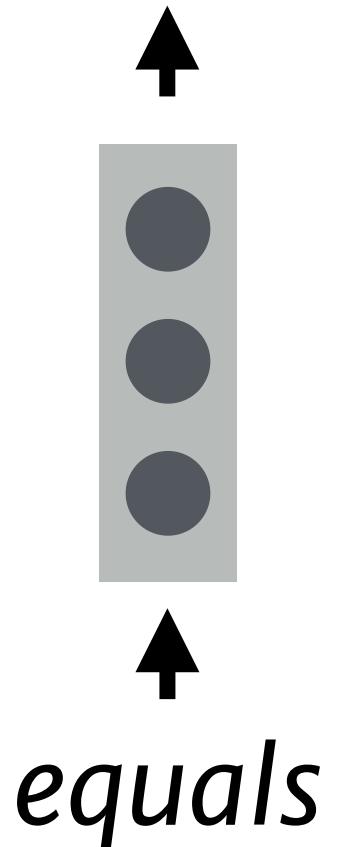
*One plus two times three equals seven.*

*Two times three times three equals ???*

# Answering math problems with LMs

---

$x_1 \sim p(x_1 | \dots \text{times three equals})$



# Answering math problems with LMs

---

$x_1 \sim p(x_1 | \dots \text{times three equals})$

*twenty*



*equals*

# Answering math problems with LMs

---

$x_1 \sim p(x_1 | \dots \text{times three equals})$

*twenty*  
-----  
↑  
  
↑  
*equals*

$x_2 \sim p(x_2 | \dots \text{equals twenty})$

*seven*  
-----  
↑  
  
↑  
----->*twenty*

(don't try this at home)

# A dataset of translated sentences

---

*Caecilius est in horto. [SEP] Caecilius is in the garden.*

*Caecilius in horto sedet. [SEP] Caecilius sits in the garden.*

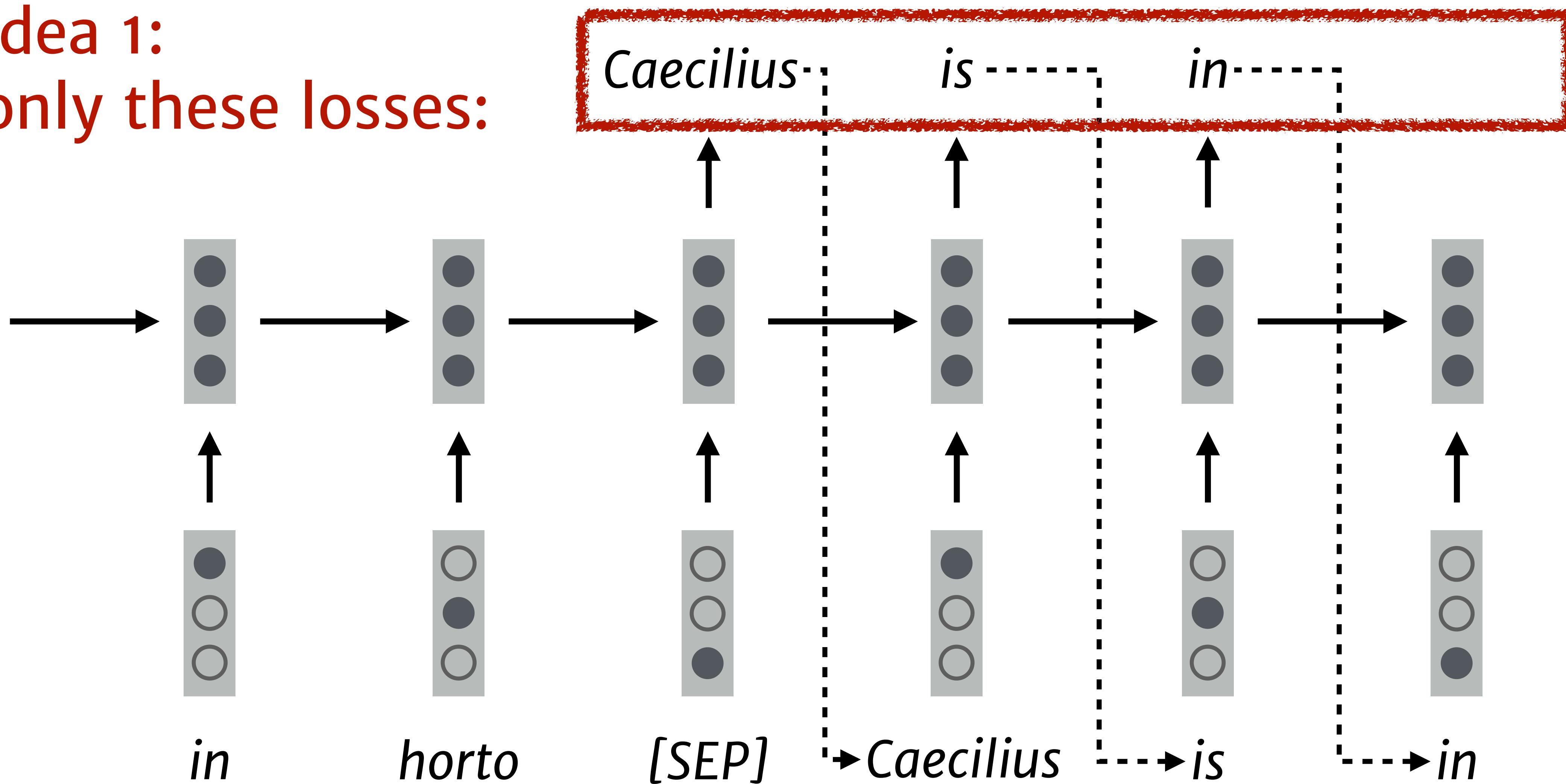
*Grumio est in atrio. [SEP] Grumio is in the atrium.*

*Grumio in atrio laborat. [SEP] ???*

(try this at home!)

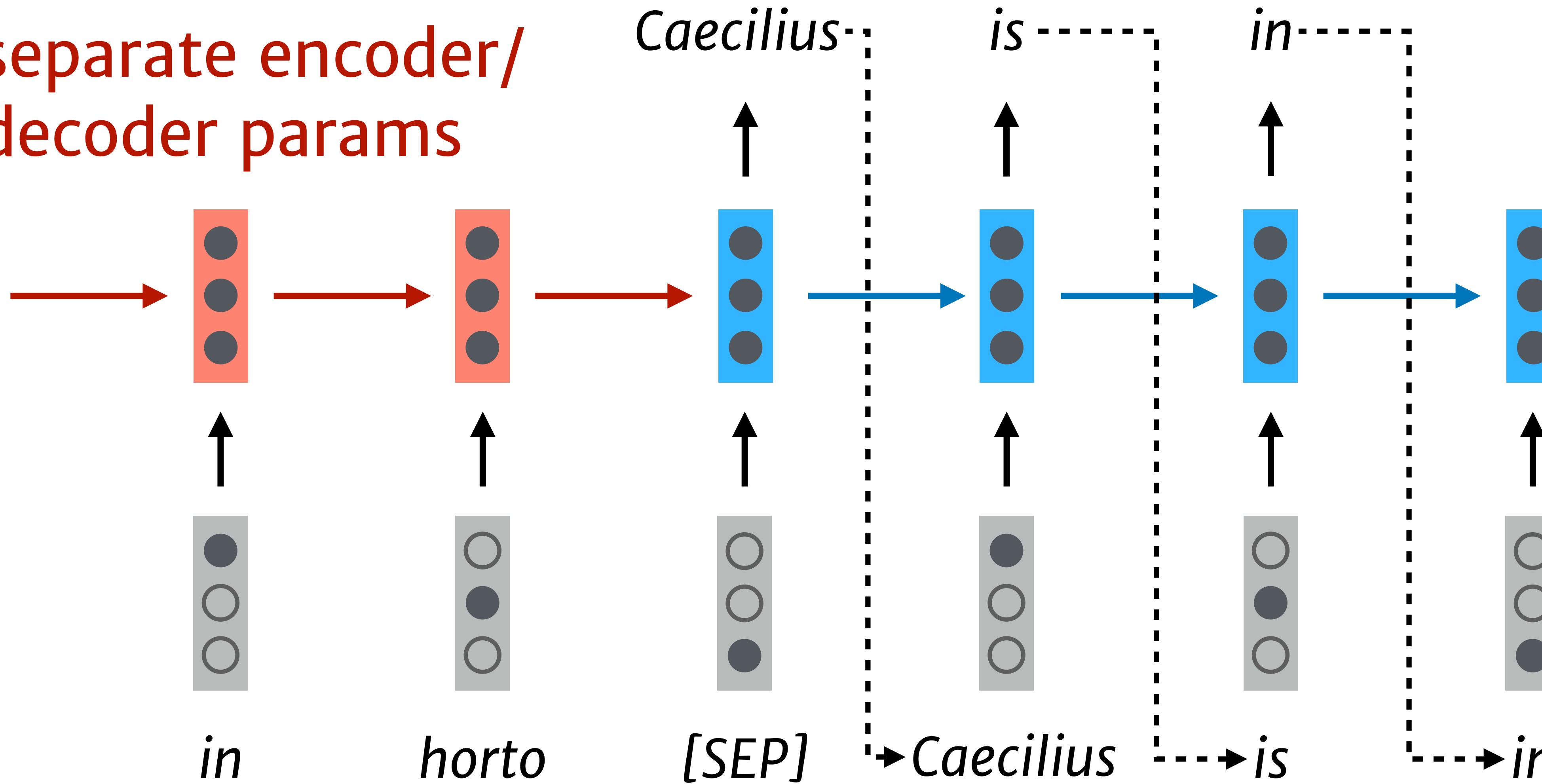
# Sequence-to-sequence models

Idea 1:  
only these losses:



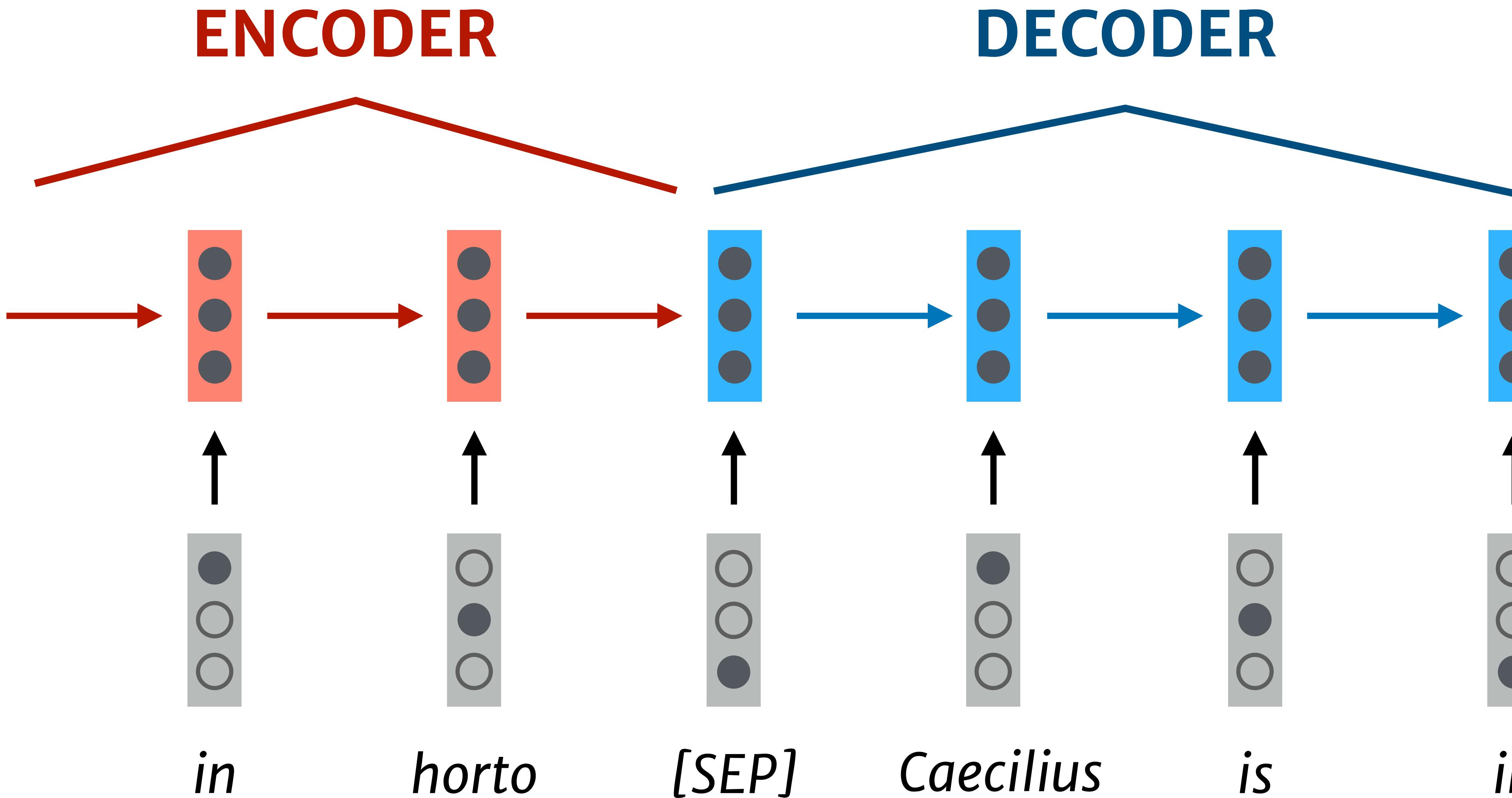
# Sequence-to-sequence models

Idea 2:  
separate encoder/  
decoder params



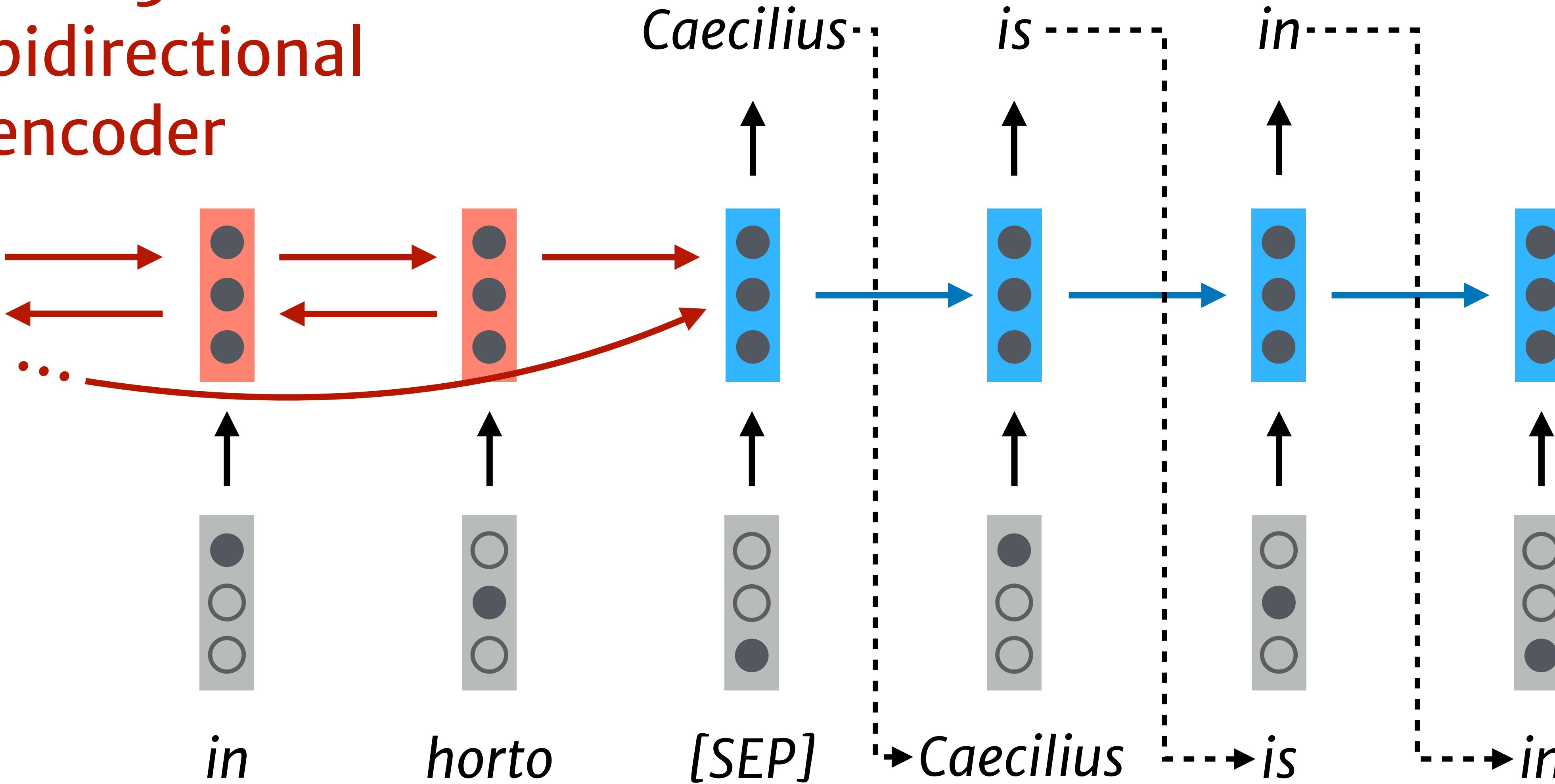
# Sequence-to-sequence models

---

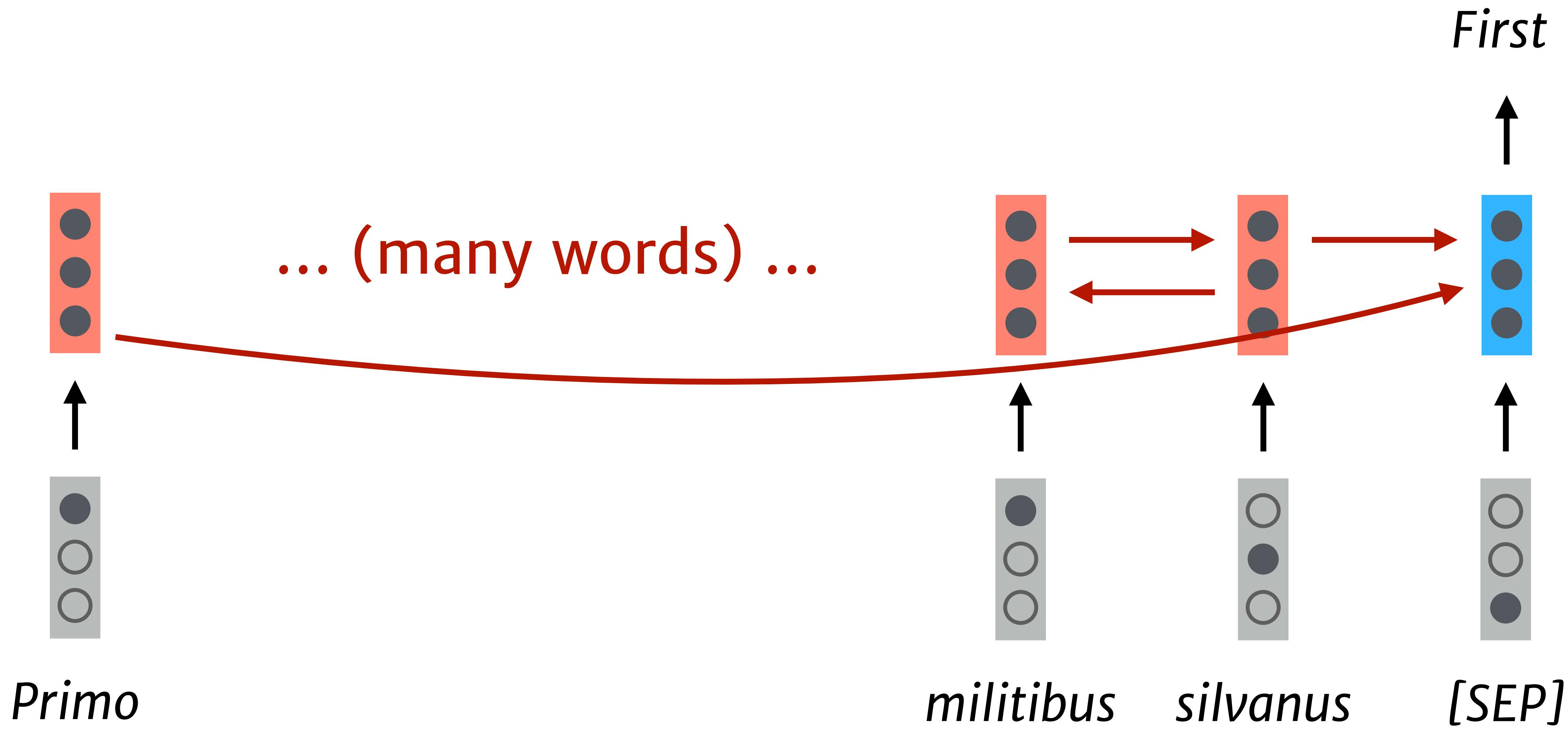


# Sequence-to-sequence models

Idea 3:  
bidirectional  
encoder

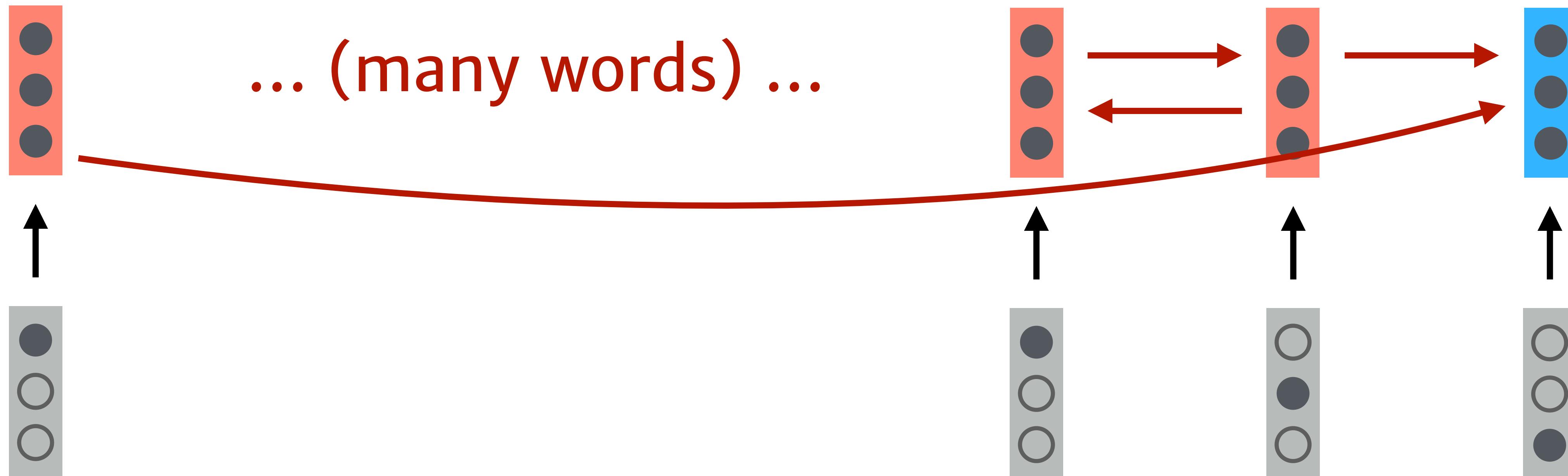


# Revenge of the vanishing gradients



# Revenge of the vanishing gradients

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_{hh}} = \sum_{i=1}^t (\mathbf{W}_{hh}^T)^{t-i} \mathbf{h}_i \quad \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_{xh}} = \sum_{i=1}^t (\mathbf{W}_{hh}^T)^{t-i} \mathbf{v}_i$$



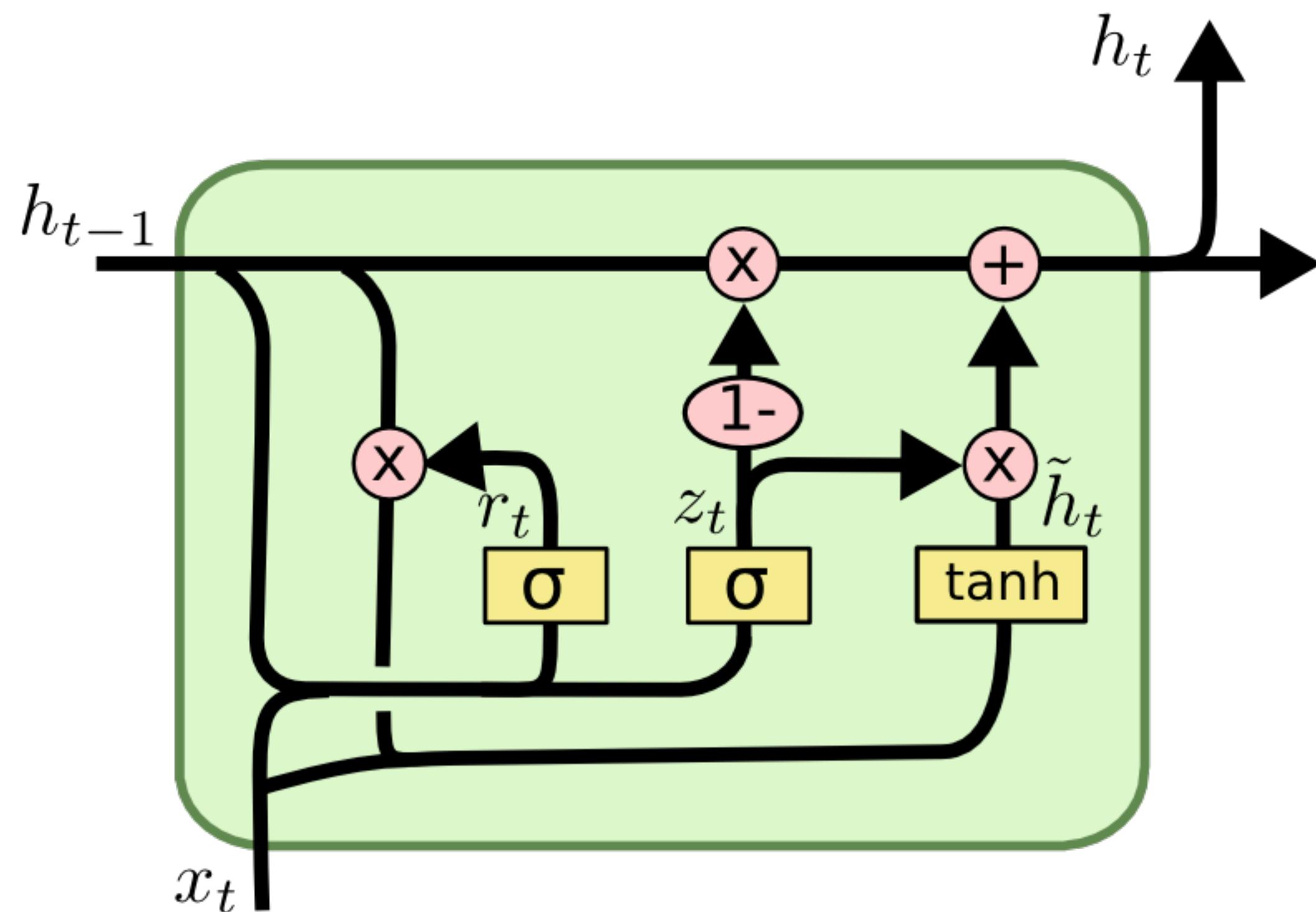
*Primo*

*militibus silvanus*

*[SEP]*

# Attention mechanisms

# Gated Recurrent Units



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

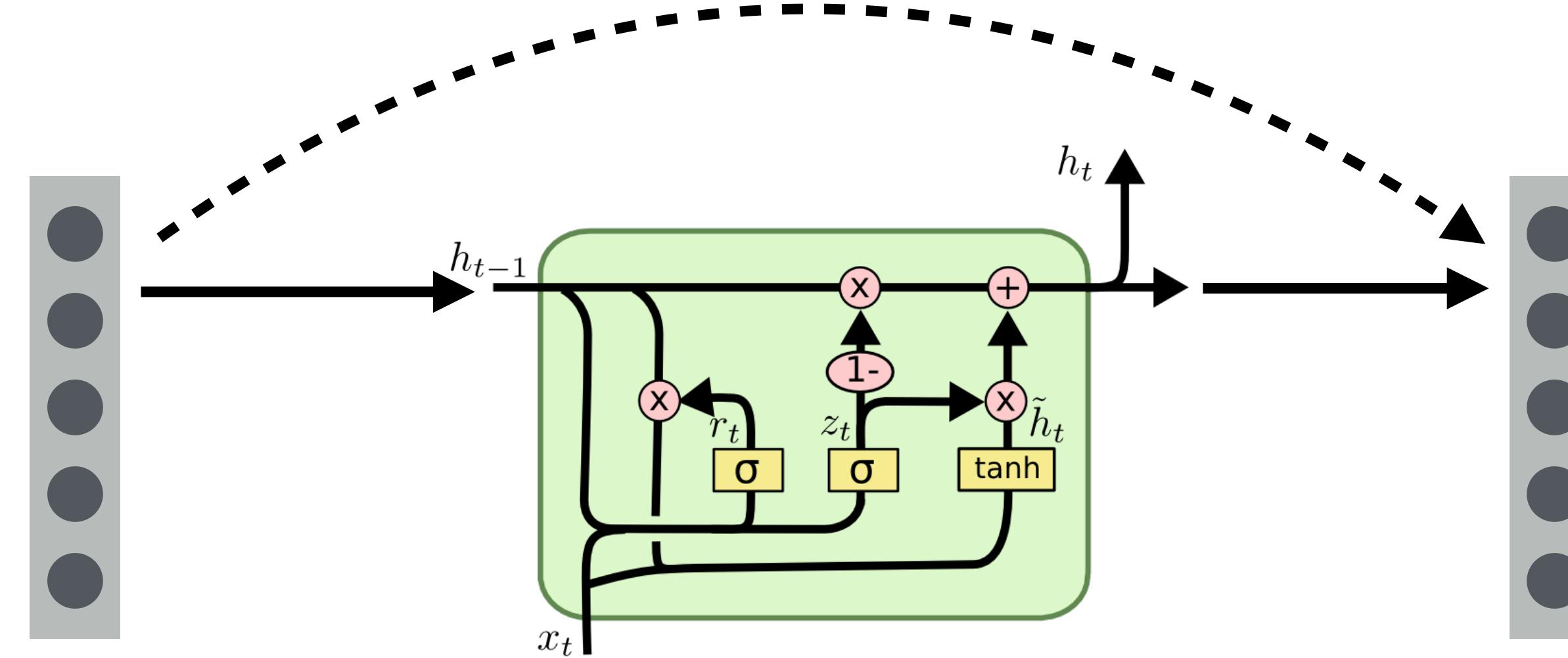
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = \boxed{(1 - z_t) * h_{t-1}} + z_t * \tilde{h}_t$$

[Image: Cristopher Olah]

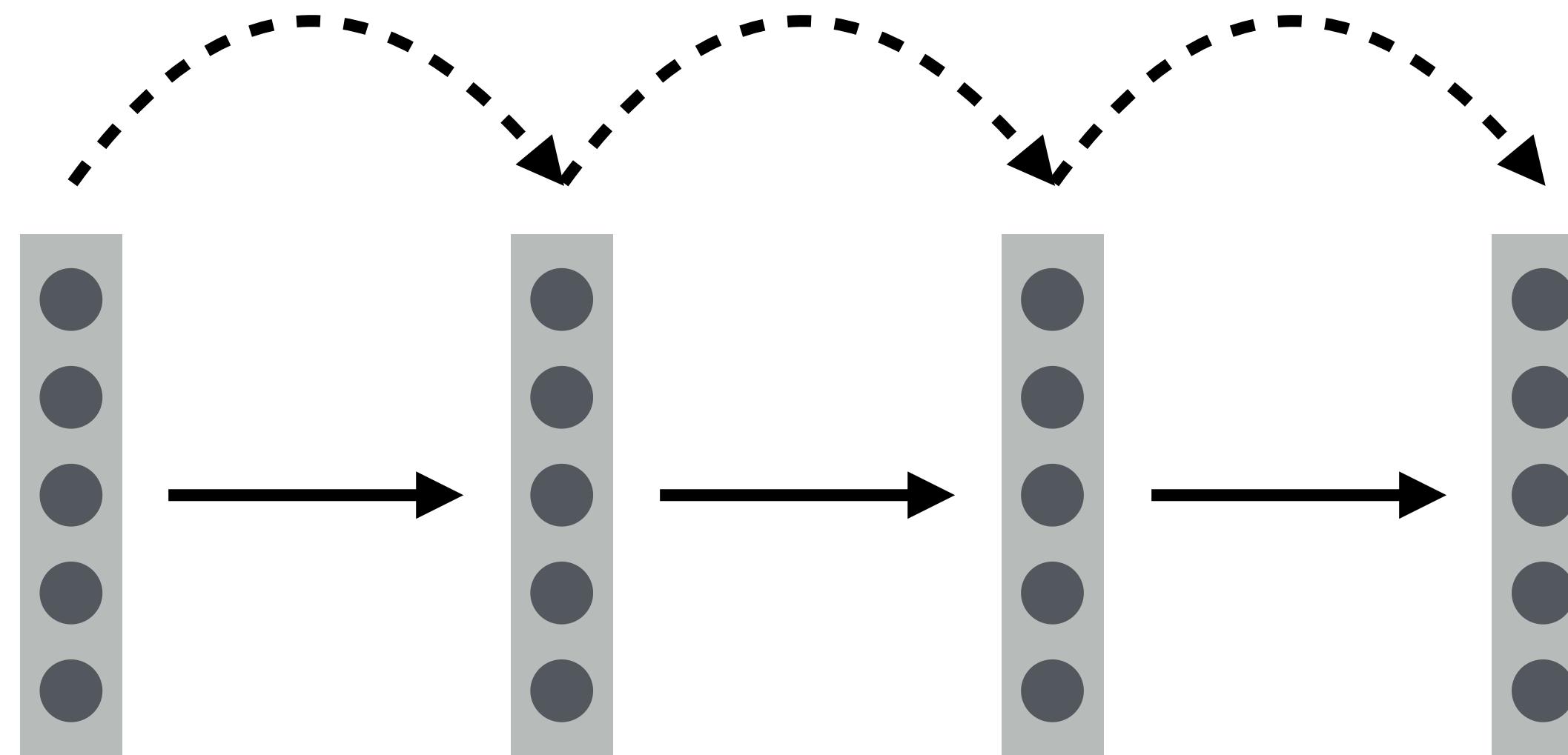
# Shortcuts

---



# Shortcuts

---

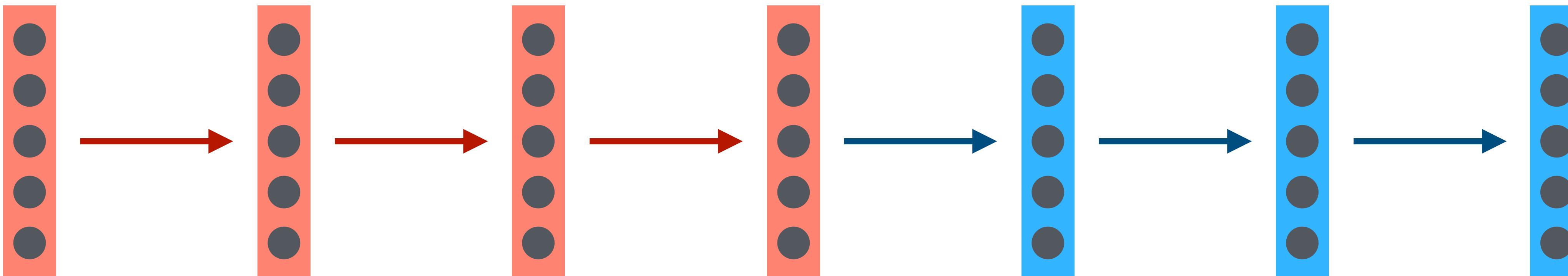


Direct "copying" between hidden states makes it easy to propagate information.

# Can we go farther?

---

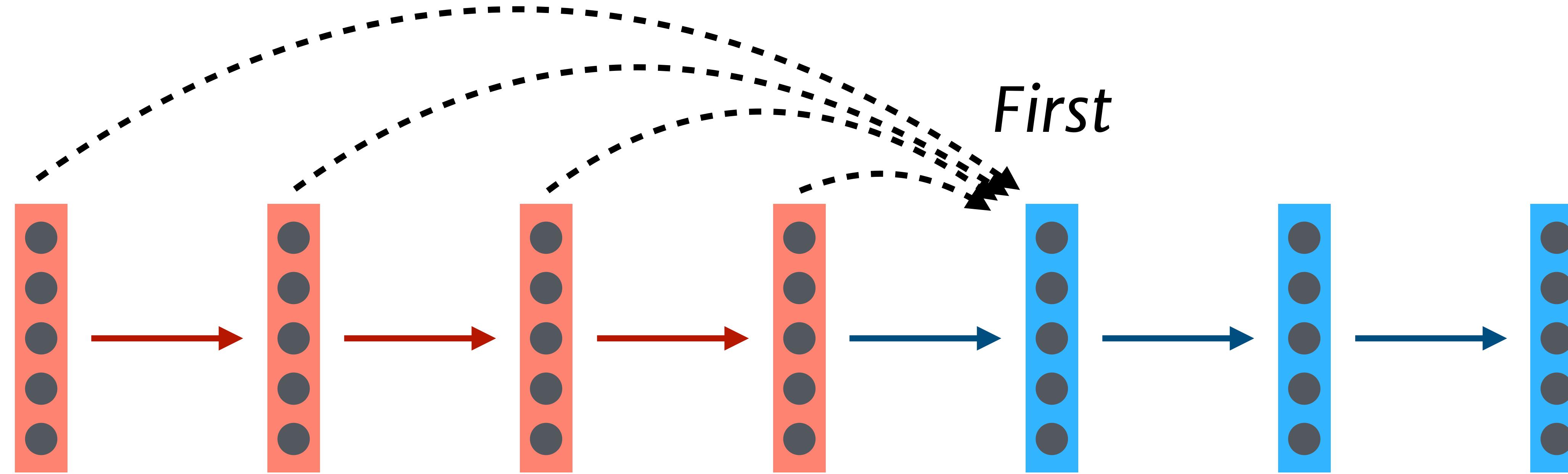
*First*



*Primo*

# Can we go farther?

---

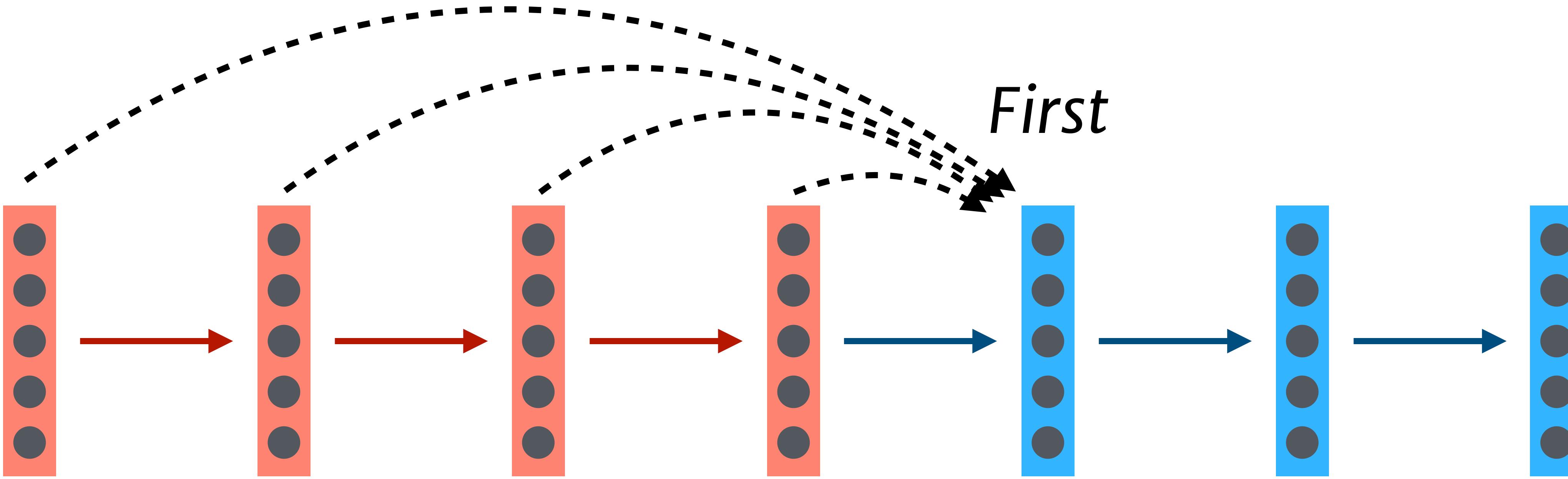


*Primo*

*First*

# Can we go farther?

---

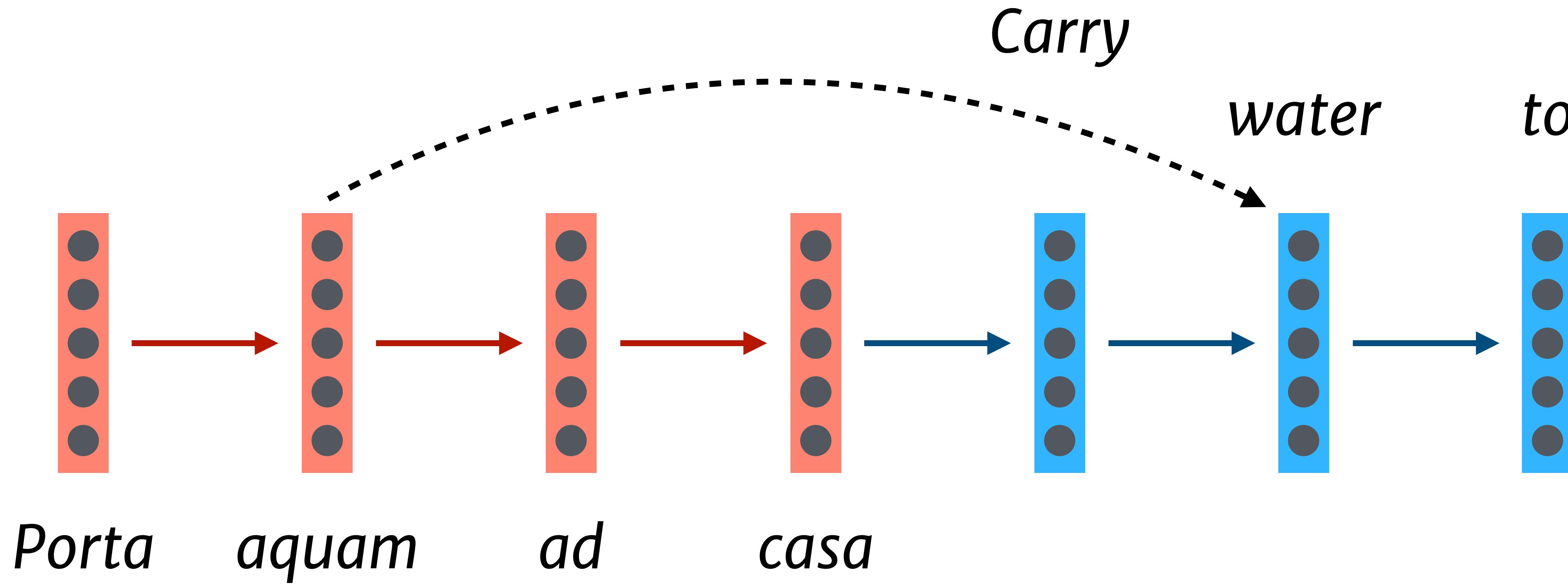


*Primo*

Not super useful: no selectivity for the relevant word  
(since we don't know which word is relevant when we add connections)

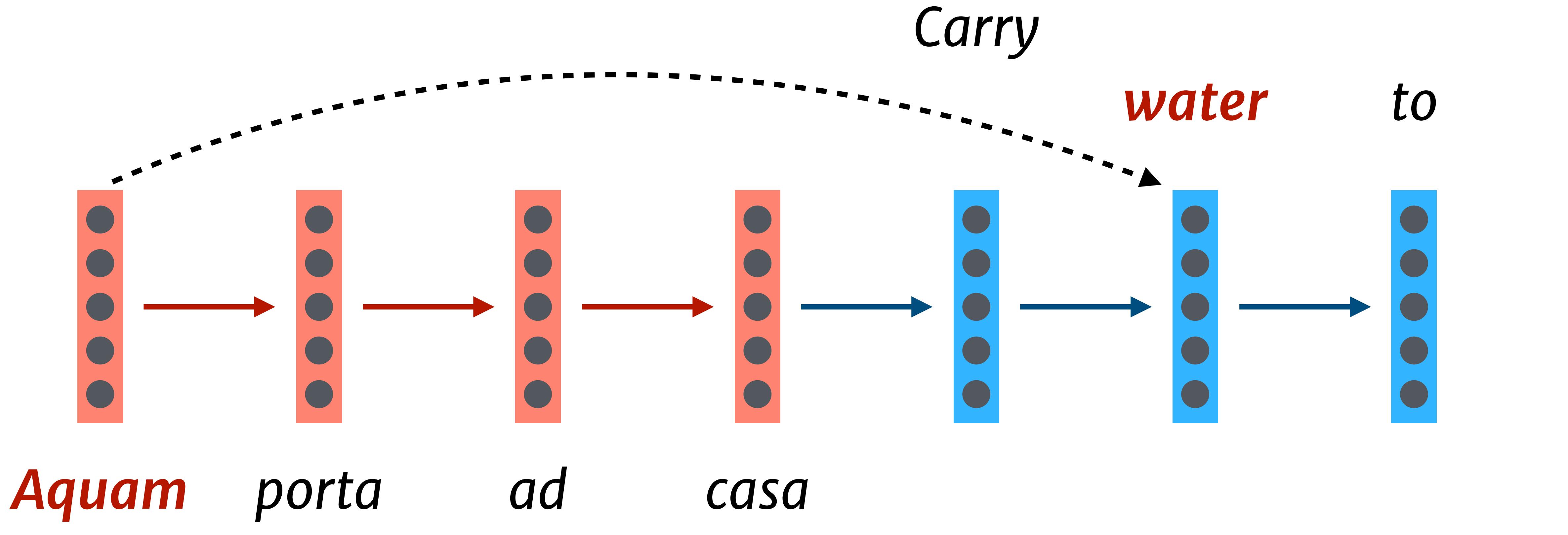
# Can we hard-code connections?

---



# Can we hard-code connections?

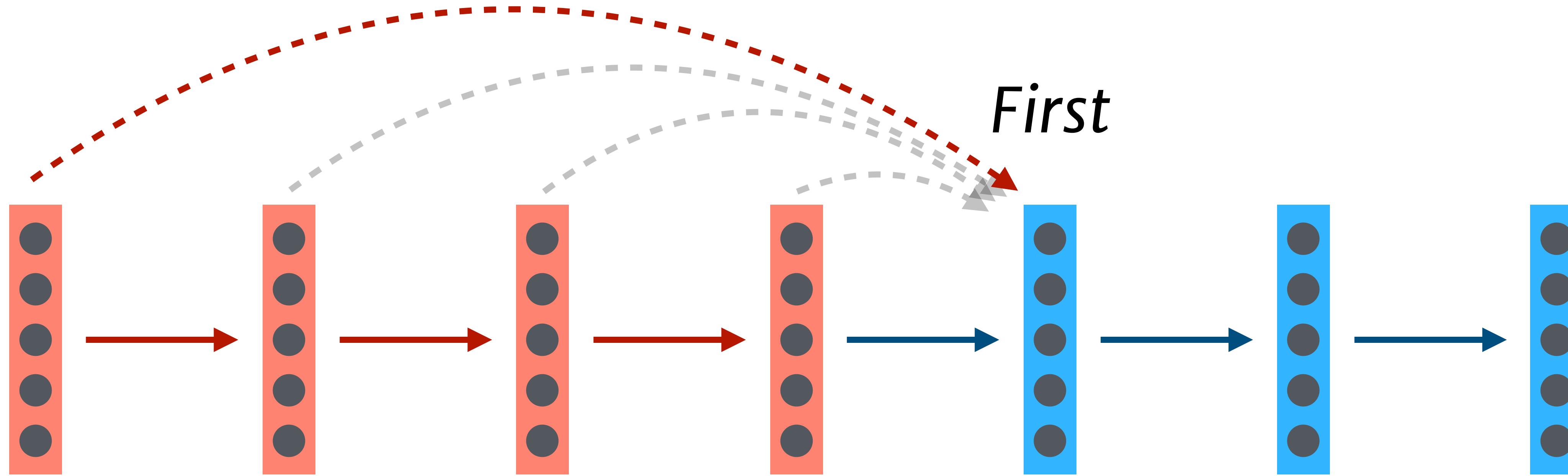
---



Words aren't one-to-one (and order can change!)

# Can we learn connections?

---

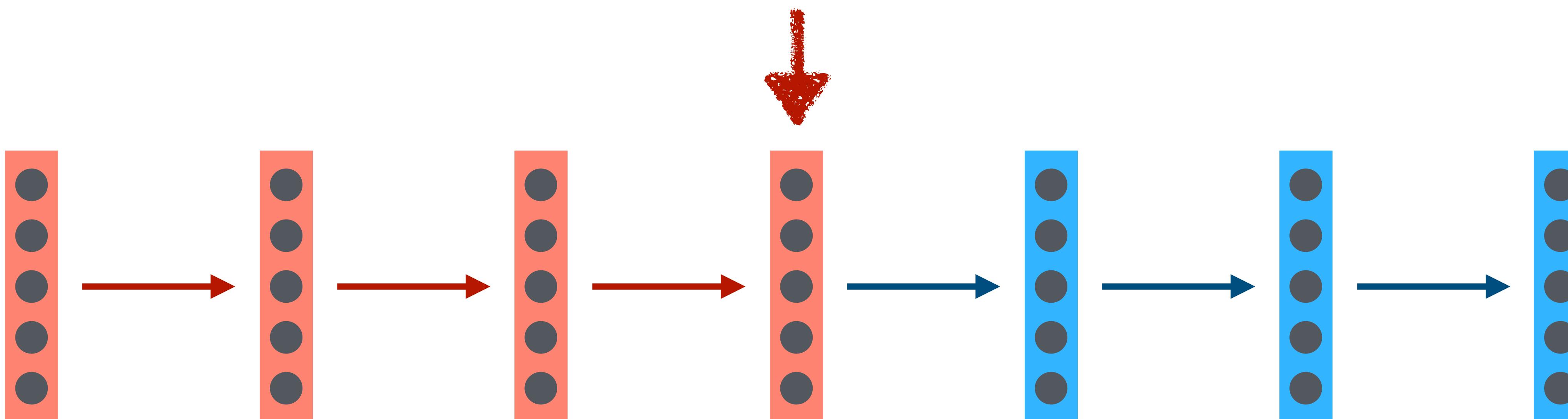


*Primo*

*First*

# Sentence representations

*This vector represents the whole sentence!*



*Aquam Aquam Aquam Aquam*  
*porta porta ad ad casa*

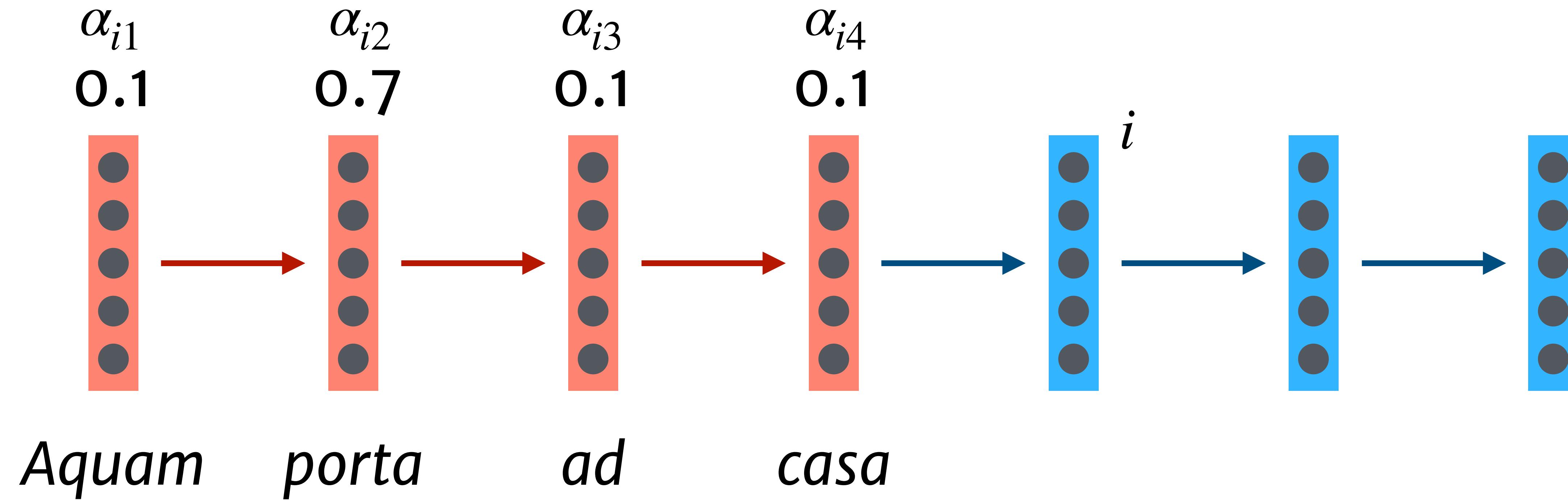


*You can't cram the  
meaning of a whole  
sentence into a  
single vector!*

[Ray Mooney, ca. 2014]

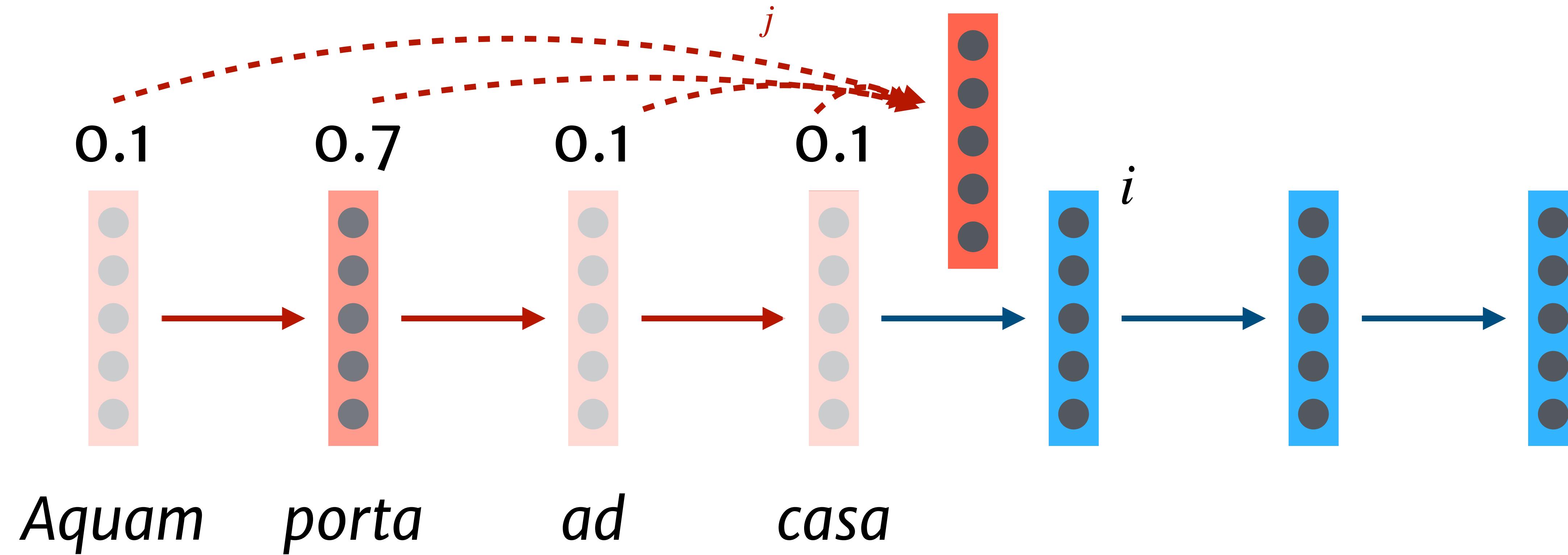
# Attention mechanisms

1. When predicting output  $i$ , assign a weight  $\alpha_{ij}$  to each encoder state  $h_j$   
(weights sum to 1)



# Attention mechanisms

1. When predicting output  $i$ , assign a weight  $\alpha_{ij}$  to each encoder state  $h_j$
2. Compute a pooled input  $c_i = \sum \alpha_{ij} h_j$



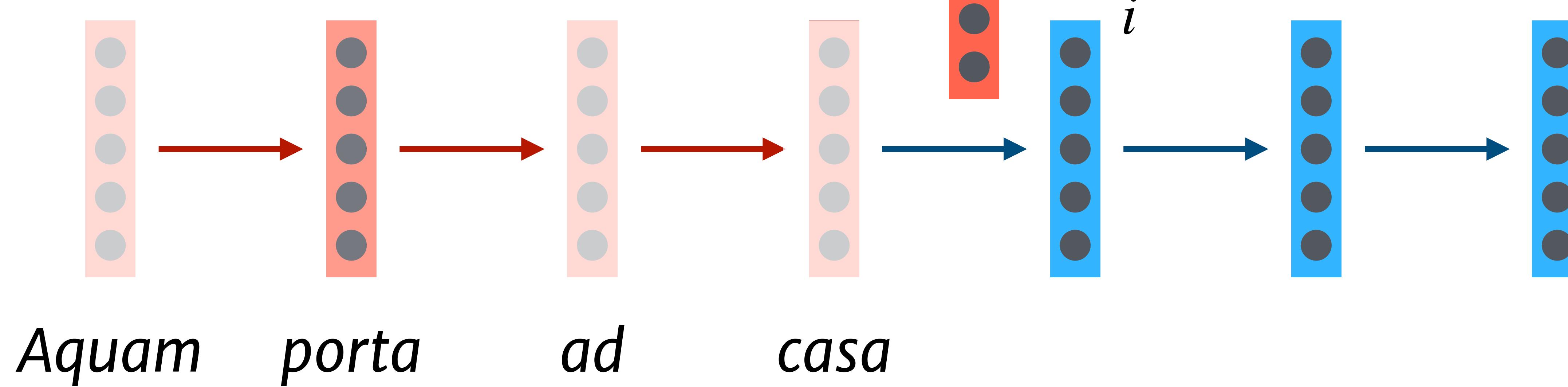
# Attention mechanisms

---

1. When predicting output  $i$ , assign a weight  $\alpha_{ij}$  to each encoder state  $h_j$

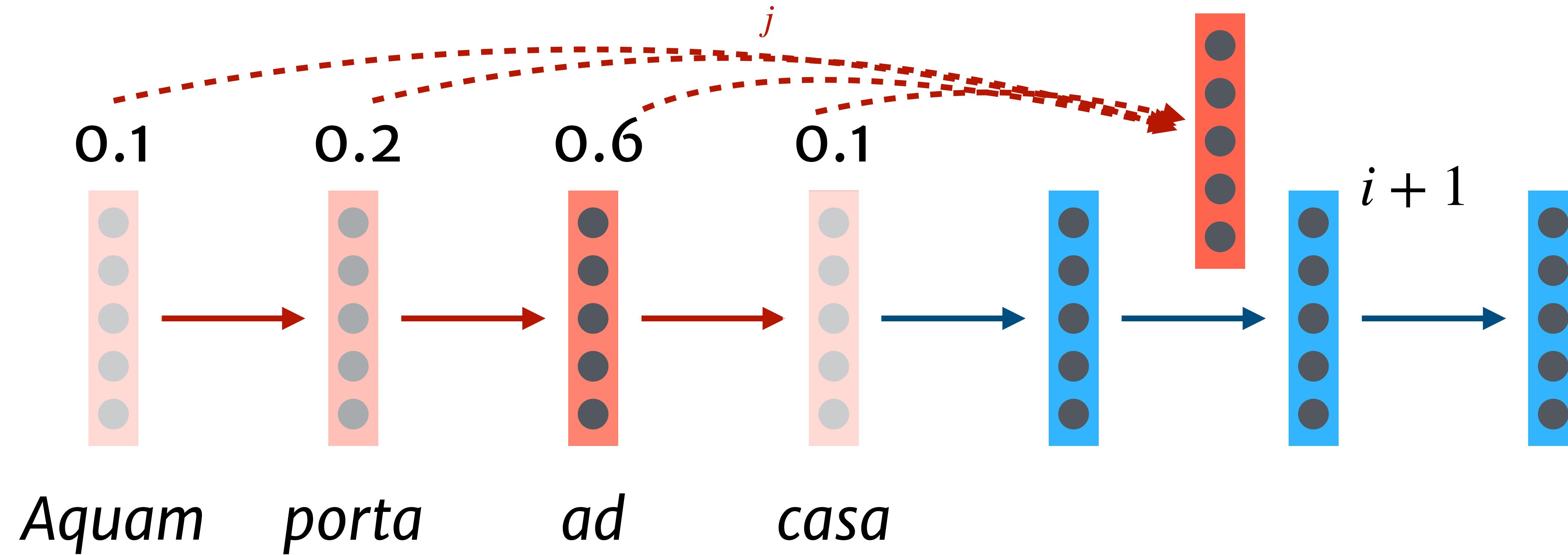
2. Compute a pooled input  $c_i = \sum_j \alpha_{ij} h_j$

3. Use  $c_i$  to update the decoder



# Attention mechanisms

1. When predicting output  $i$ , assign a weight  $\alpha_{ij}$  to each encoder state  $h_j$
2. Compute a pooled input  $c_i = \sum \alpha_{ij} h_j$



# Design decision: how to compute $\alpha_{ij}$ ?

1. When predicting output  $i$ , assign a weight  $\alpha_{ij}$  to each encoder state  $h_j$

$$e_{ij} = w^\top \tanh(W^\top [h_i, h_j])$$

[Bahdanau 2014]

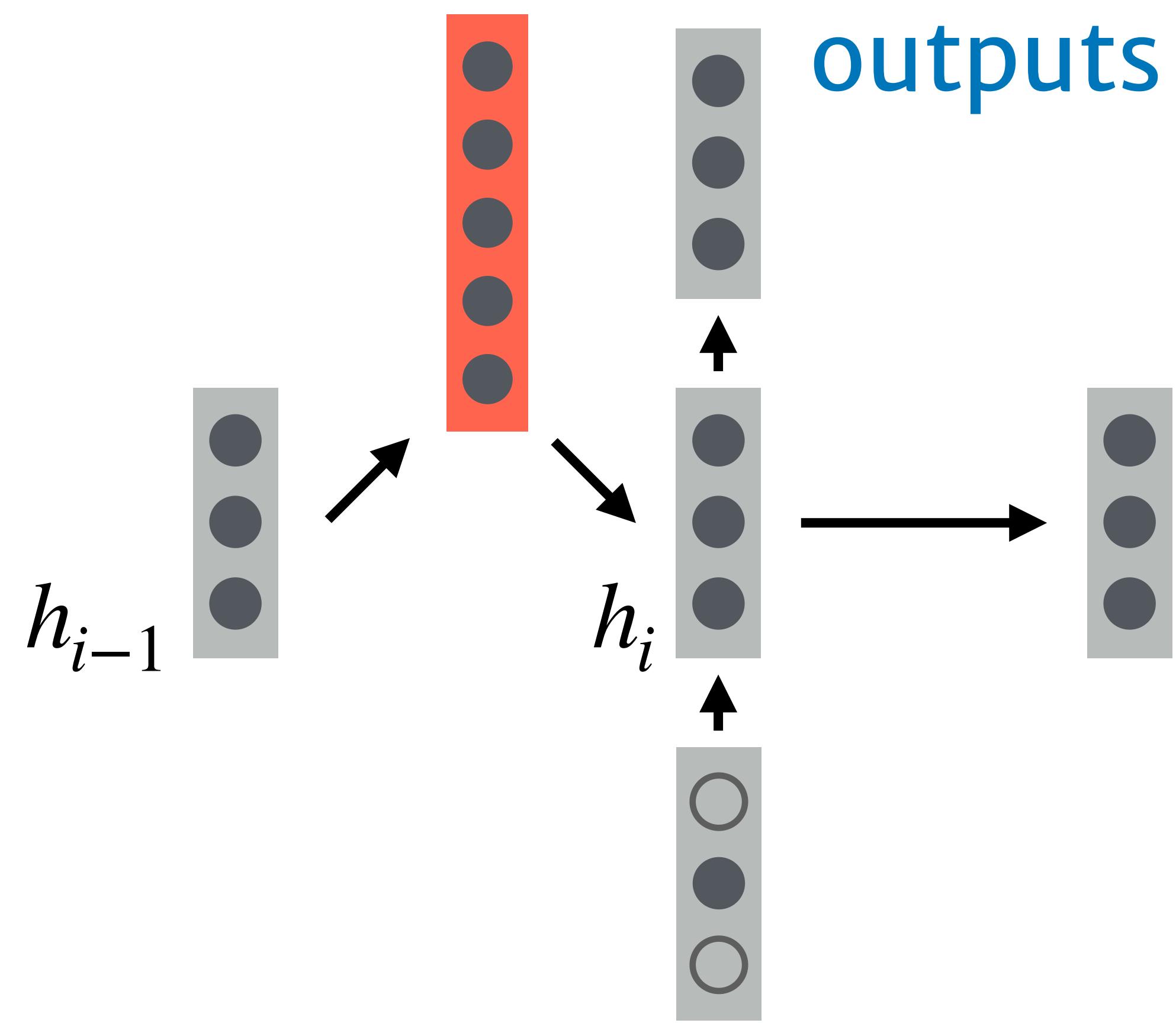
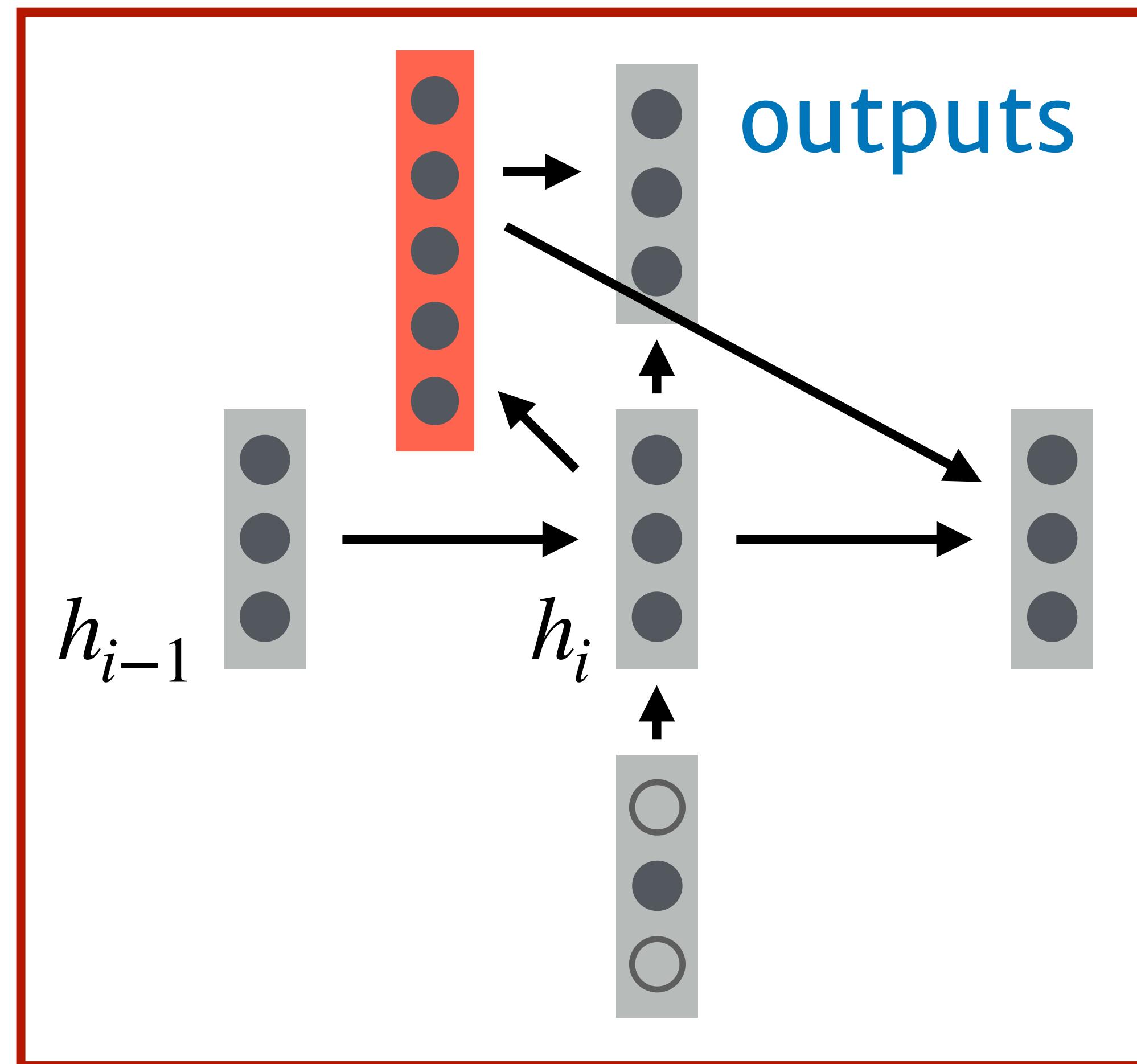
$$e_{ij} = h_i^\top W h_j$$

[Luong 2015]

$$\alpha_{i,:} = \text{softmax}(e_{i,:})$$

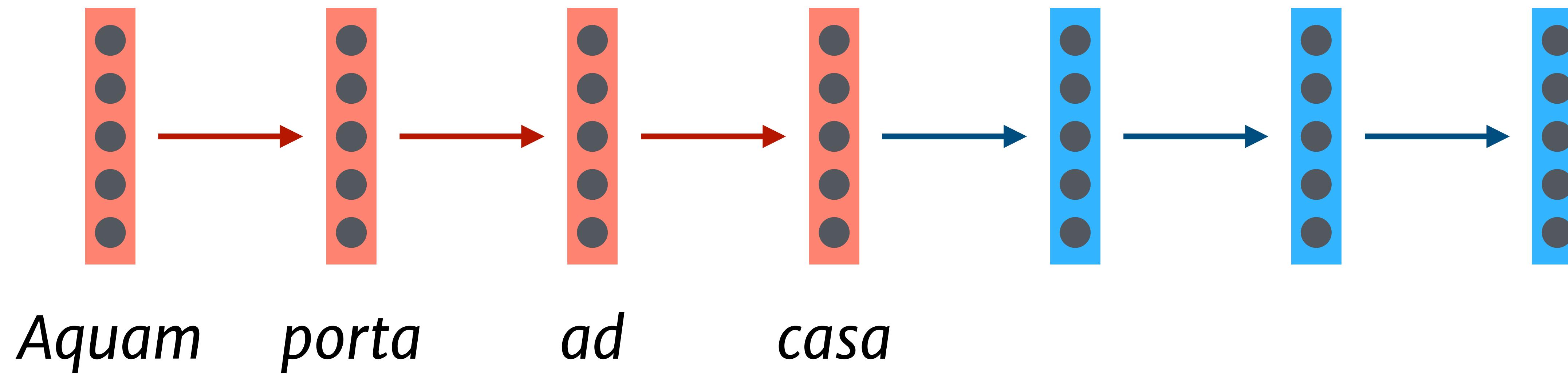
# Design decision: how to use $c_i$ ?

## 3. Use $c_i$ to update the decoder



# Why does this work?

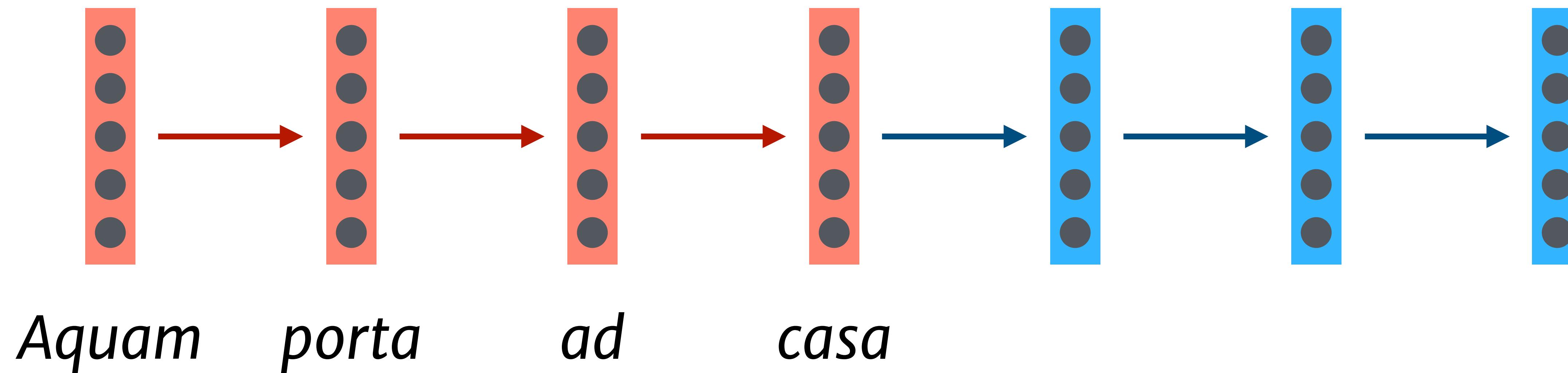
---



# Why does this work?

---

MAIN VERB  
IMPERATIVE

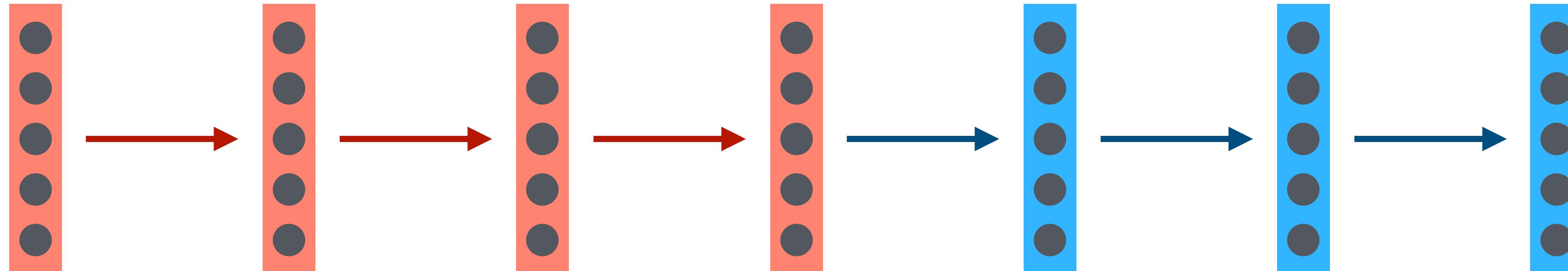


# Why does this work?

---

MAIN VERB  
IMPERATIVE

SUBJECT?  
IMP. VERB?



*Aquam*

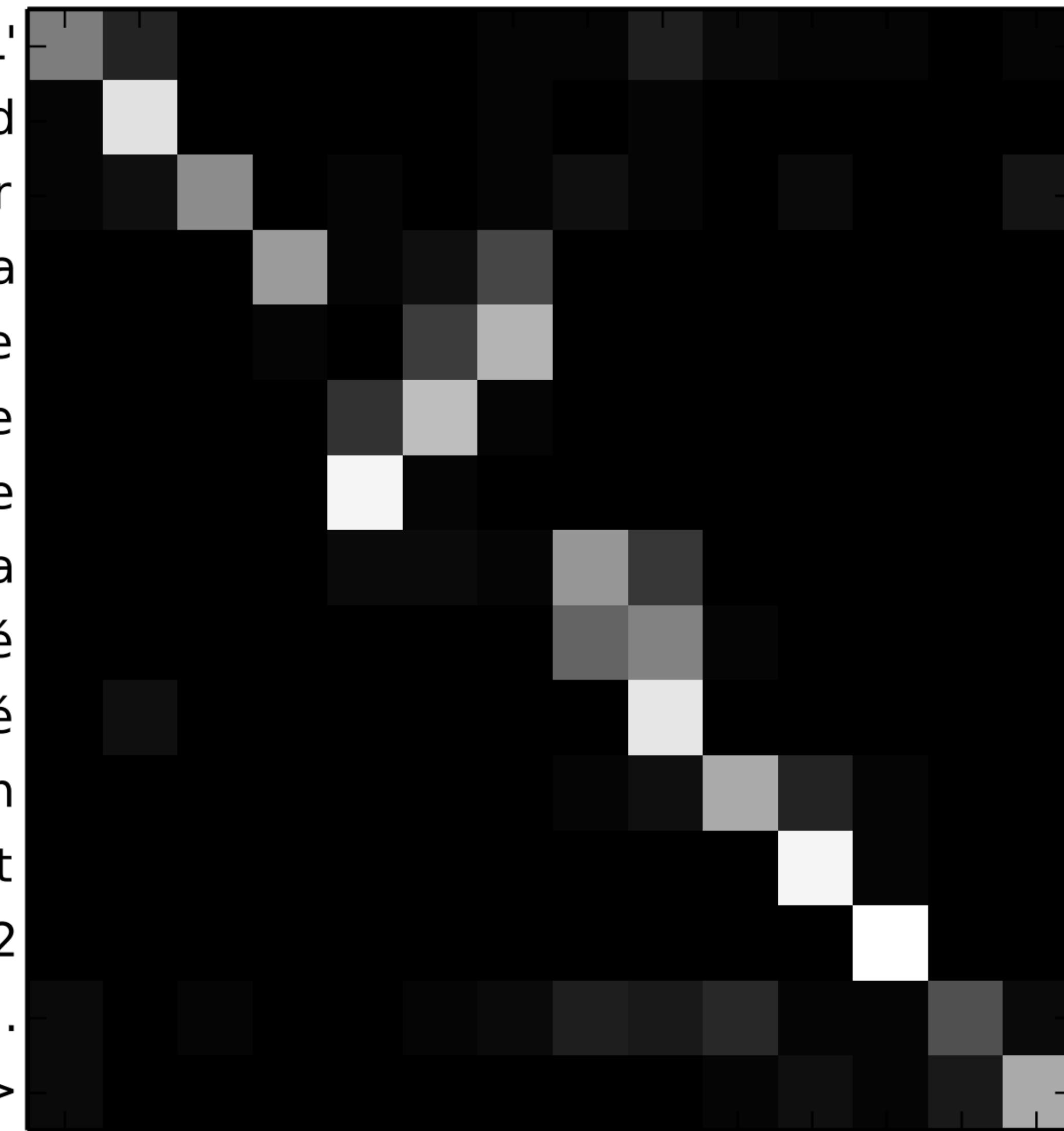
*porta*

*ad*

*casa*

The  
agreement  
on  
the  
European  
Economic  
Area  
was  
signed  
in  
August  
1992  
. <end>

L'  
accord  
sur  
la  
zone  
économique  
européenne  
a  
été  
signé  
en  
août  
1992  
. <end>



[Example from Greg Durrett]

# Multi-headed attention

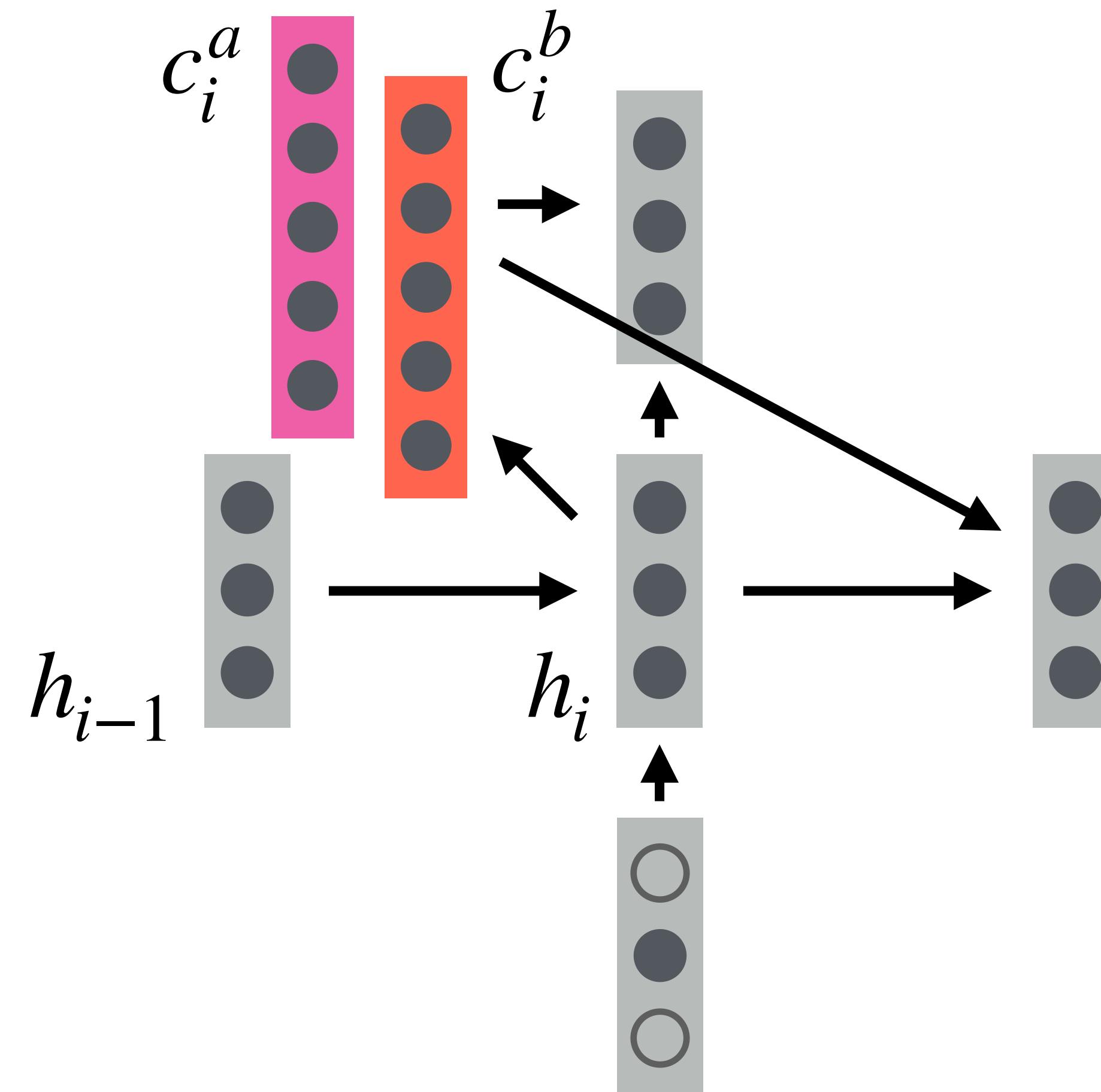
---

Look two places at once!

$$e_{ij}^a = h_i^\top W_a h_j$$

$$e_{ij}^b = h_i^\top W_b h_j$$

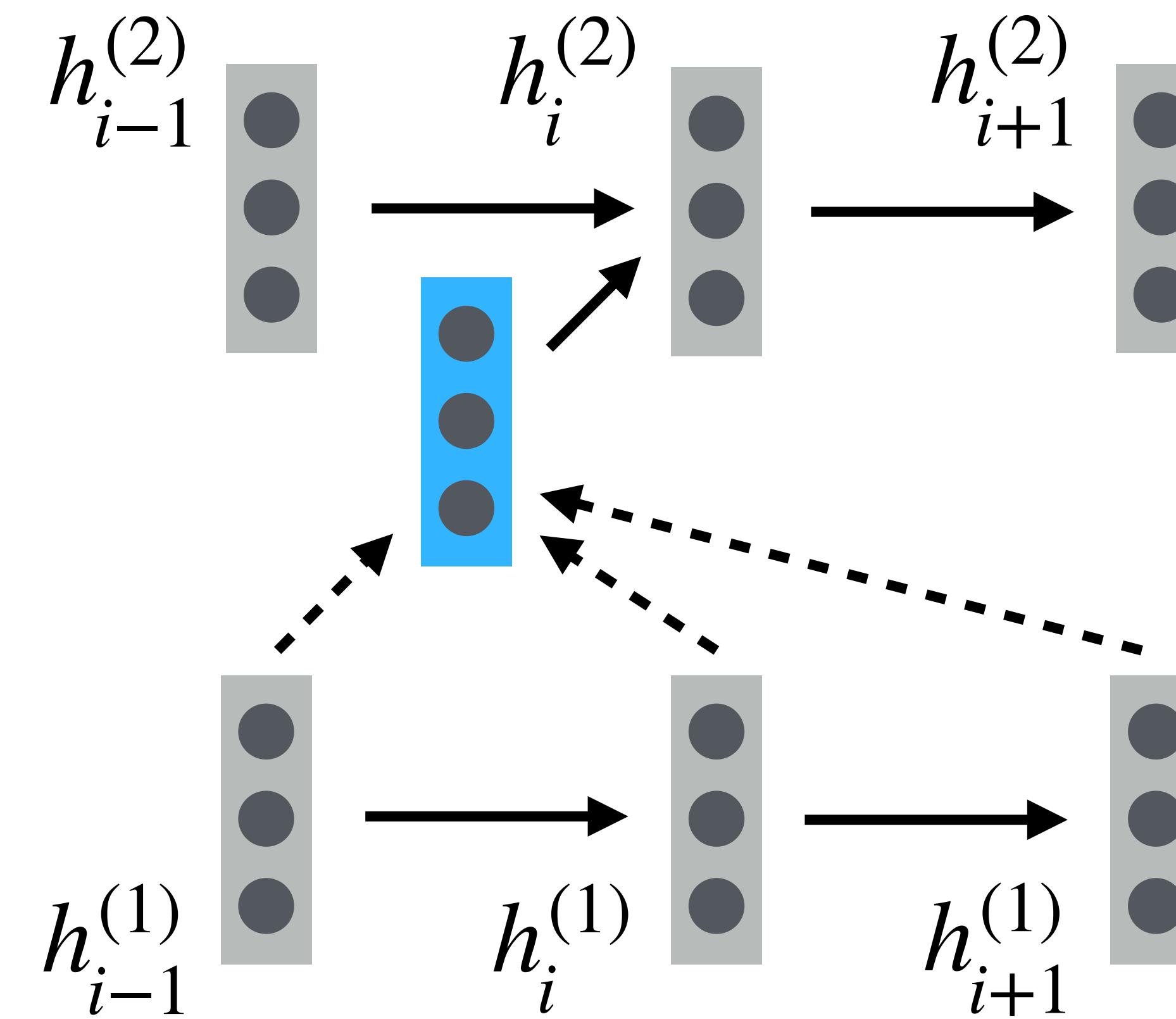
etc.



# Self-attention

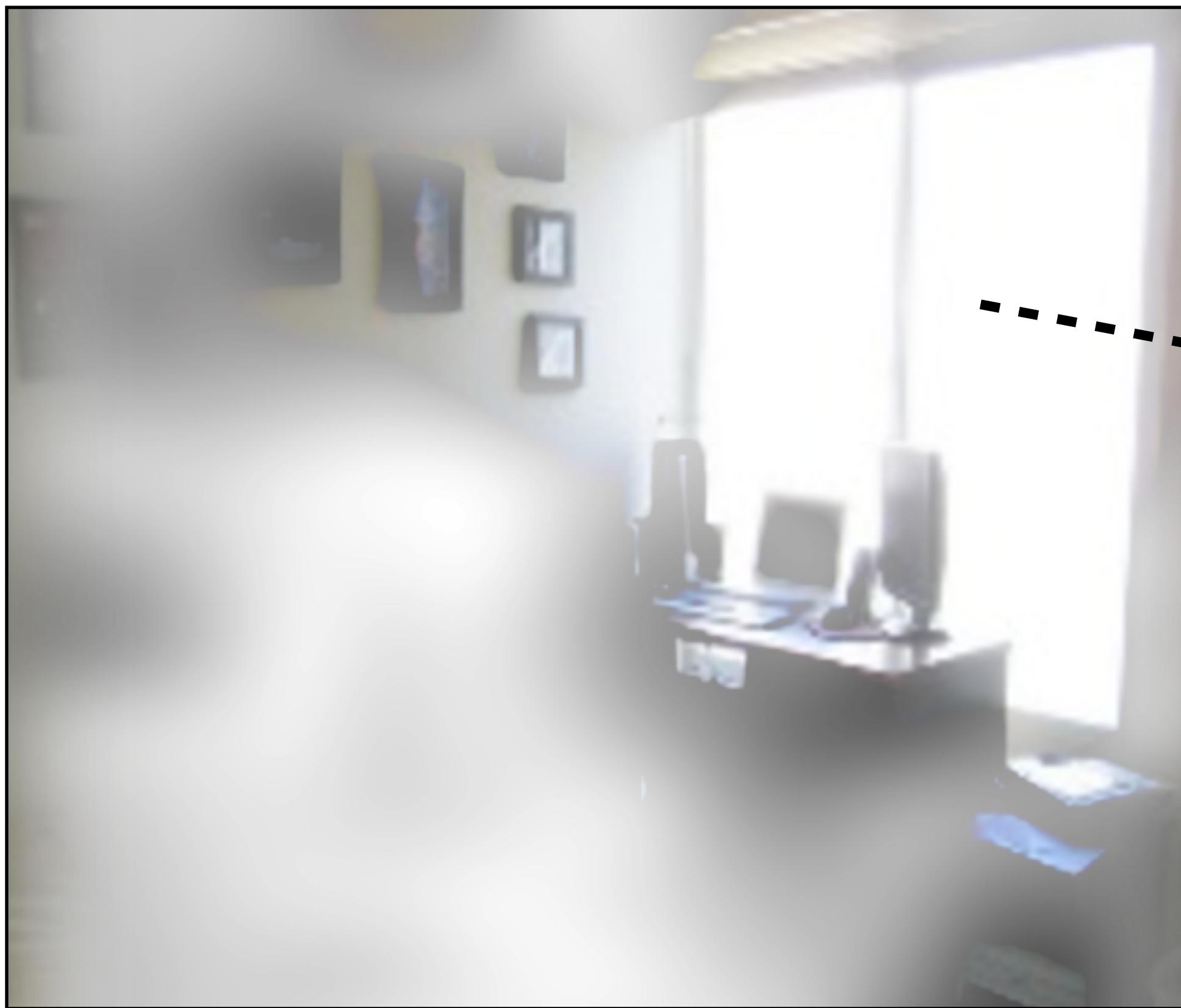
---

Attention to lower RNN layers (instead of decoder → encoder)

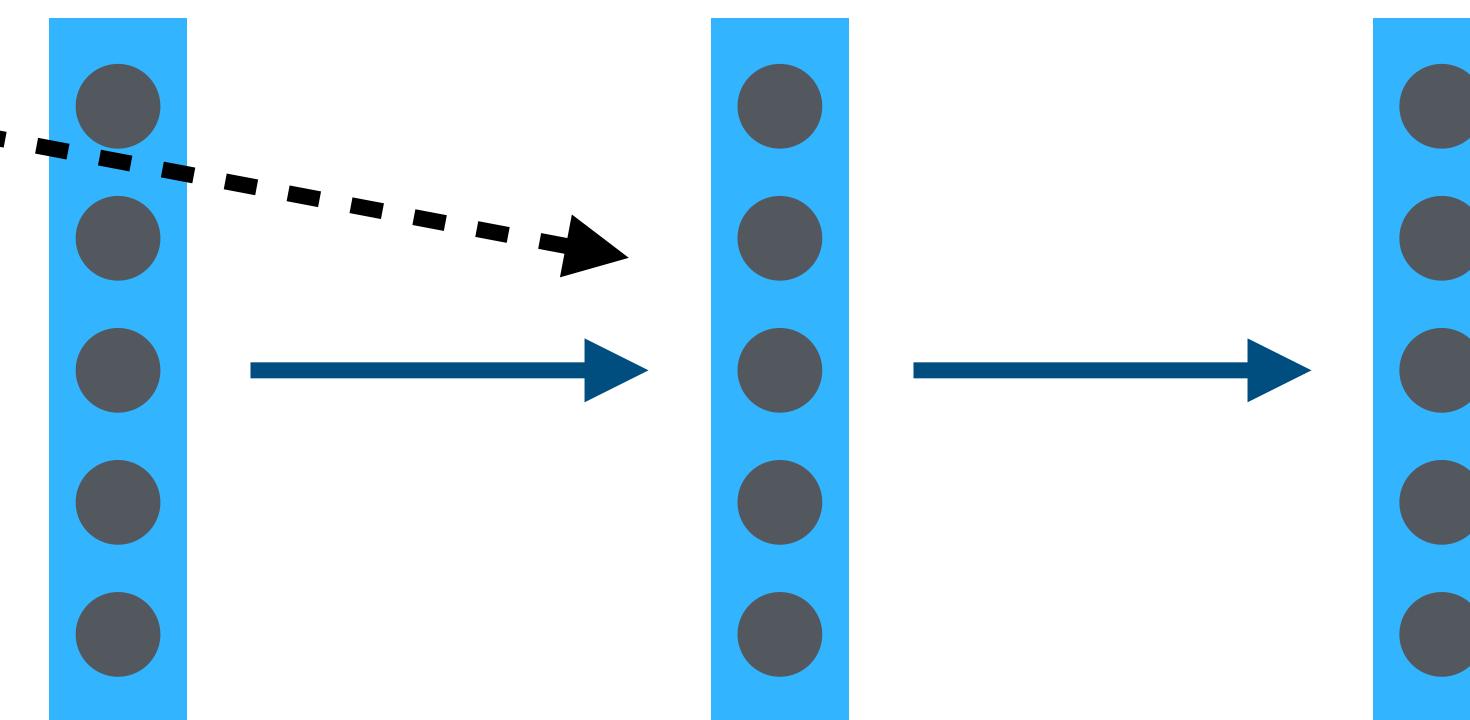


# Non-textual attention

---

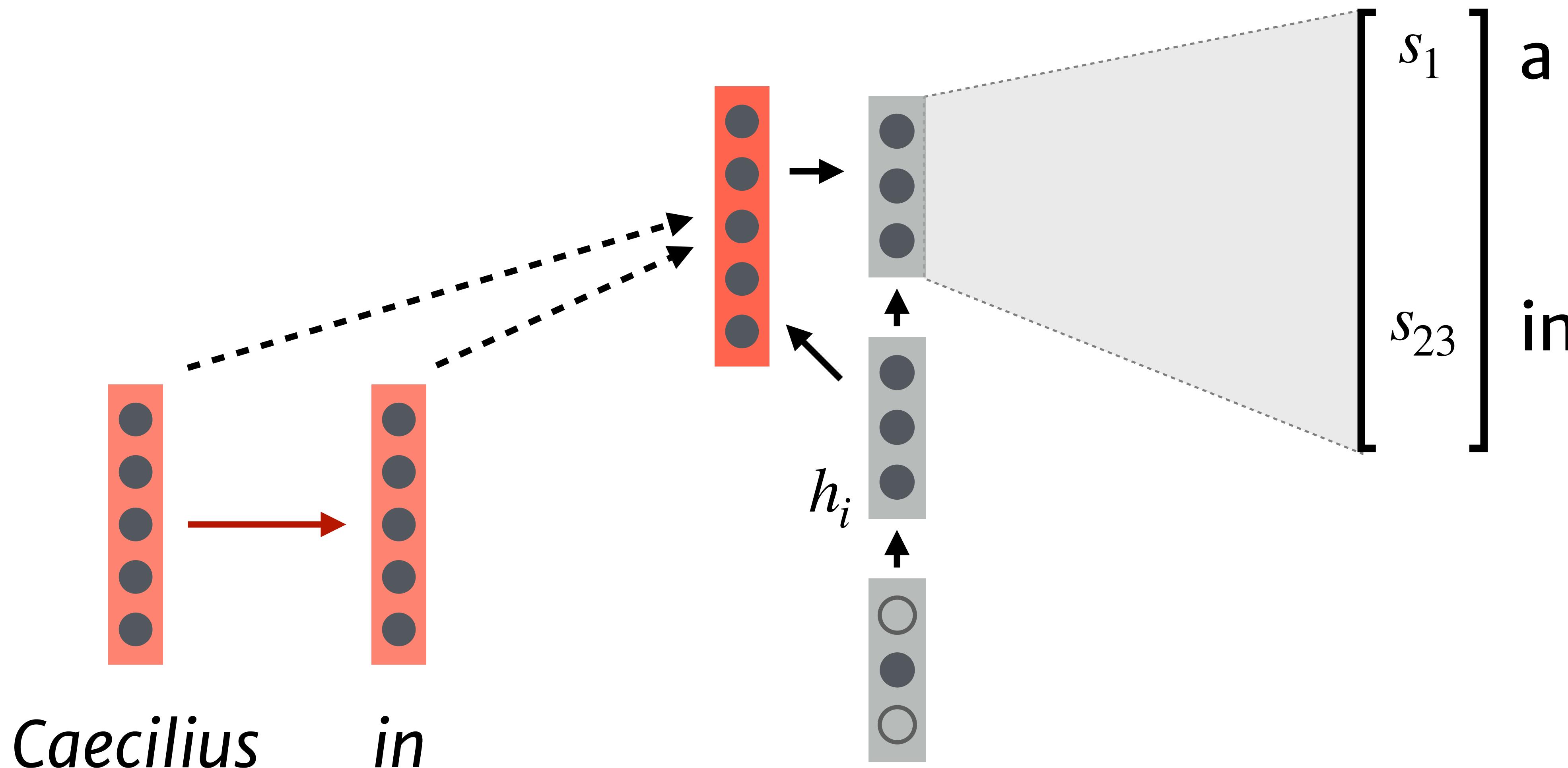


*a      desk      behind*



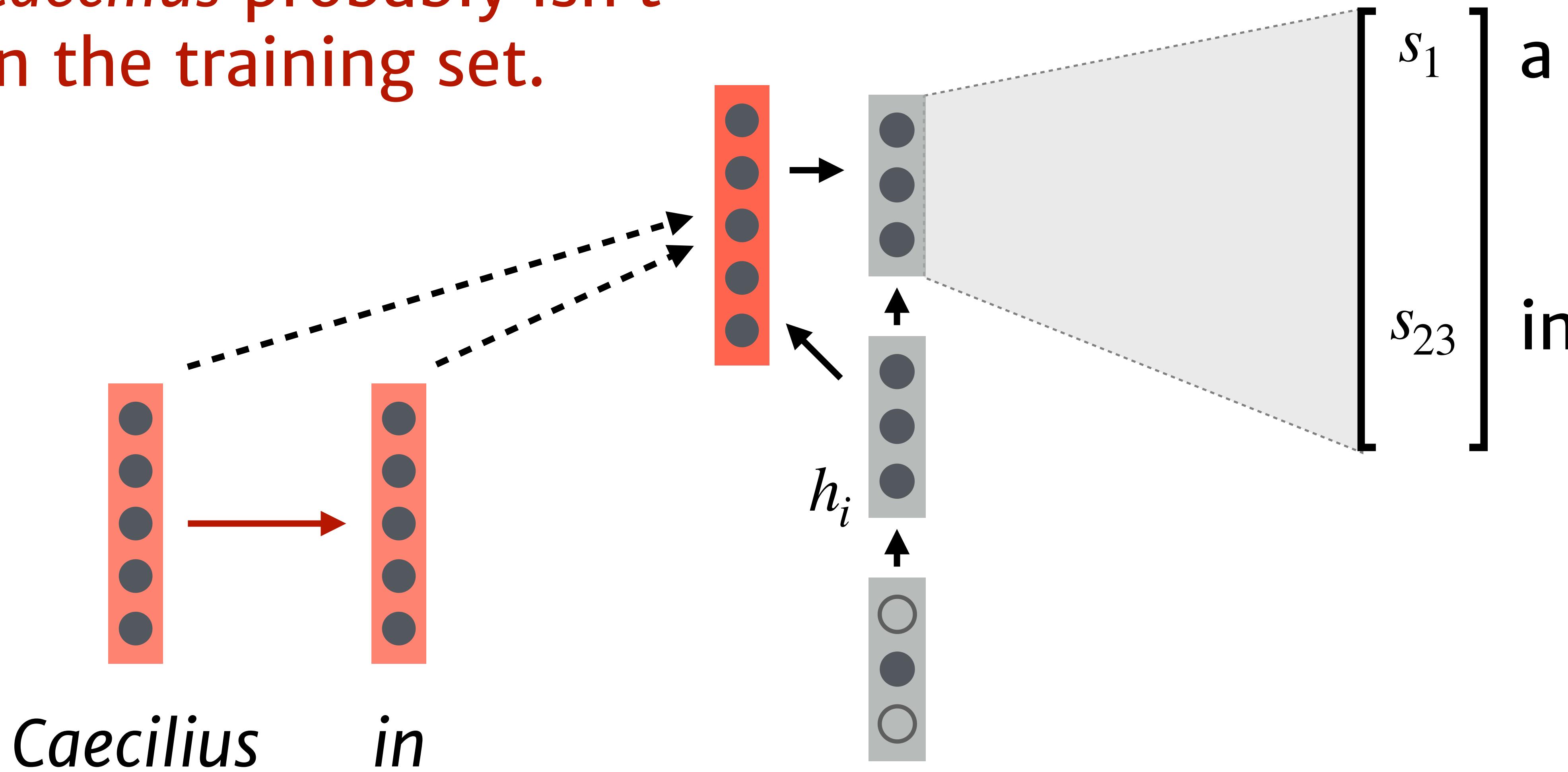
# Copying

---



# Copying

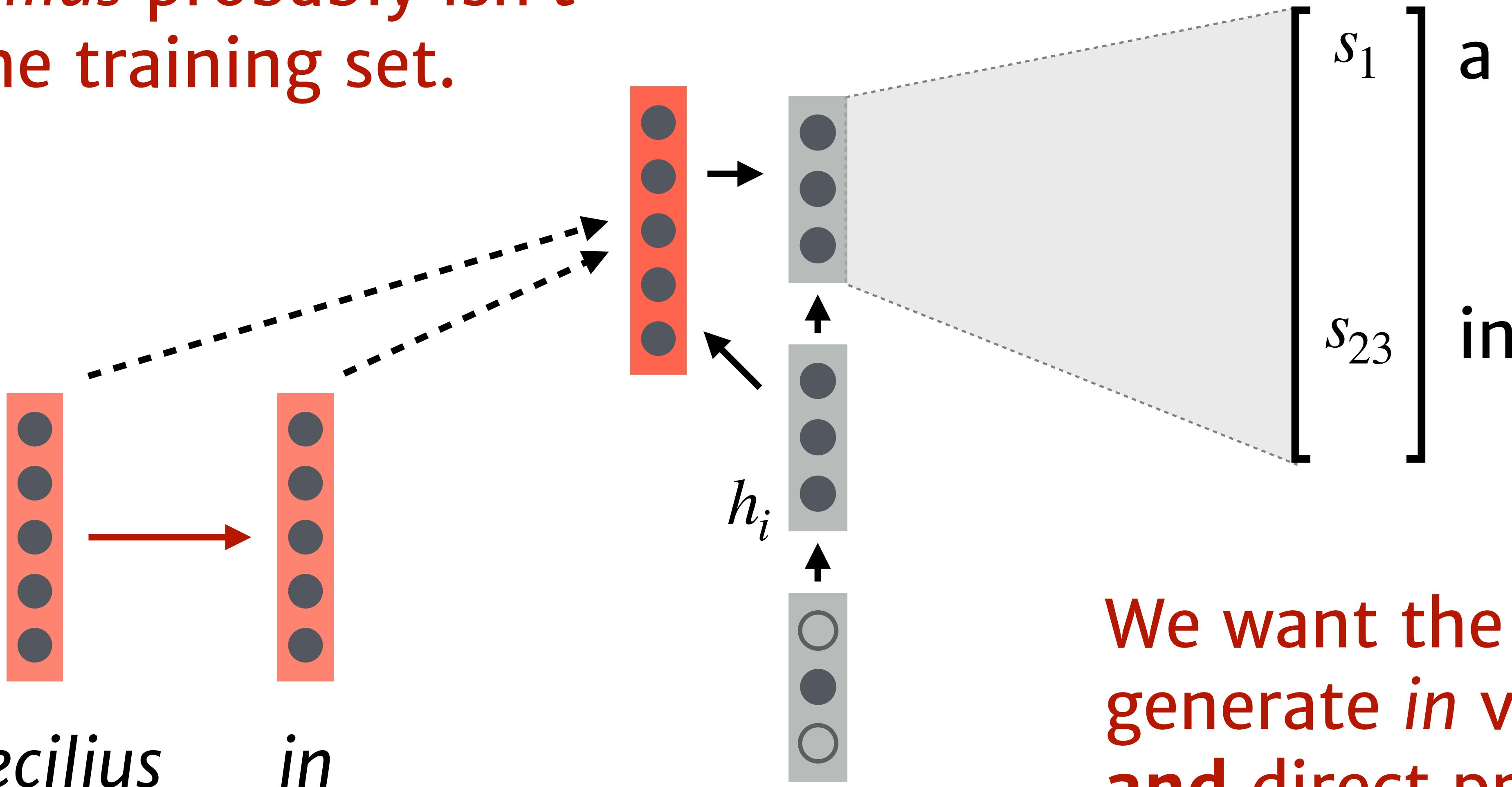
*Caecilius* probably isn't  
in the training set.



# Copying

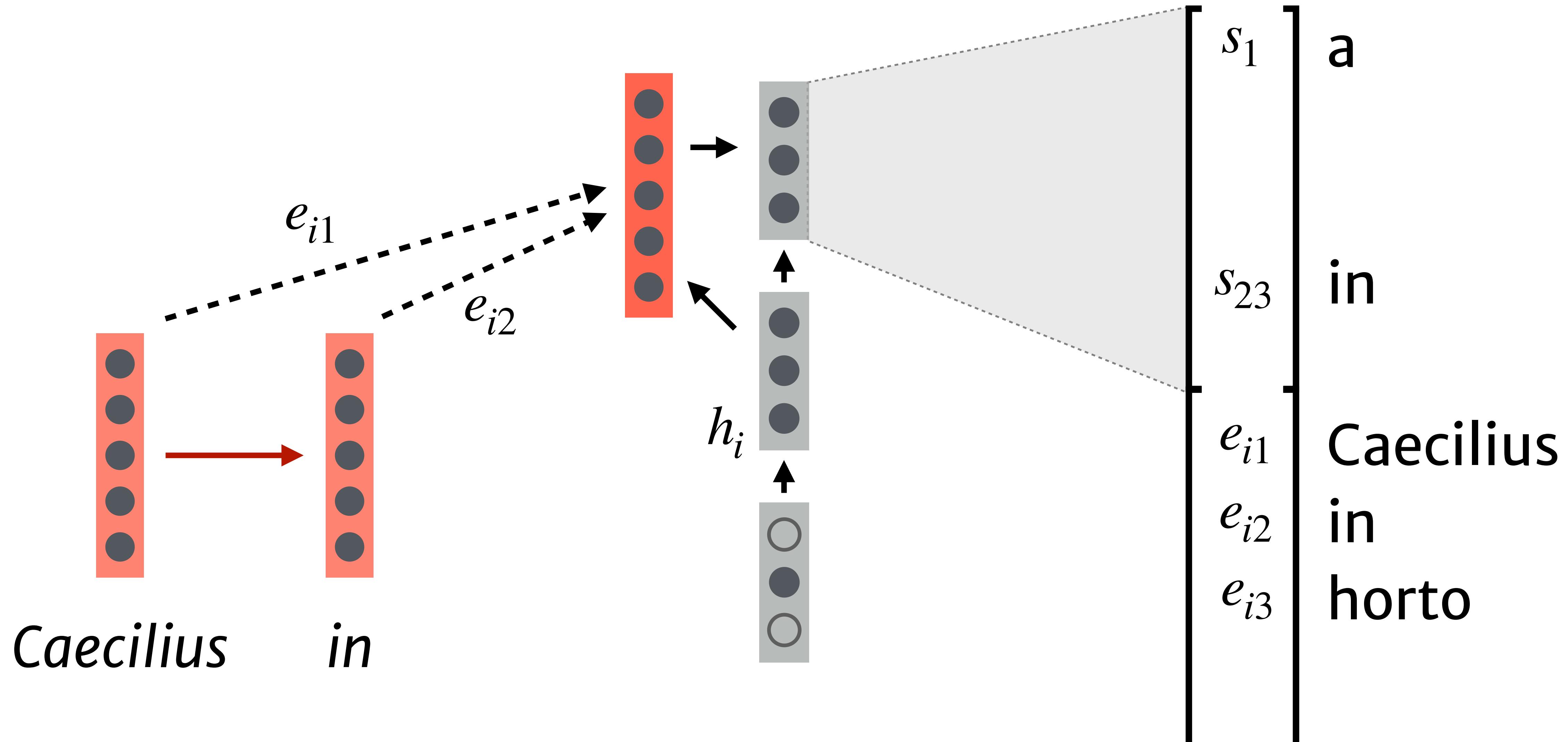
*Caecilius* probably isn't  
in the training set.

*Caecilius*      *in*



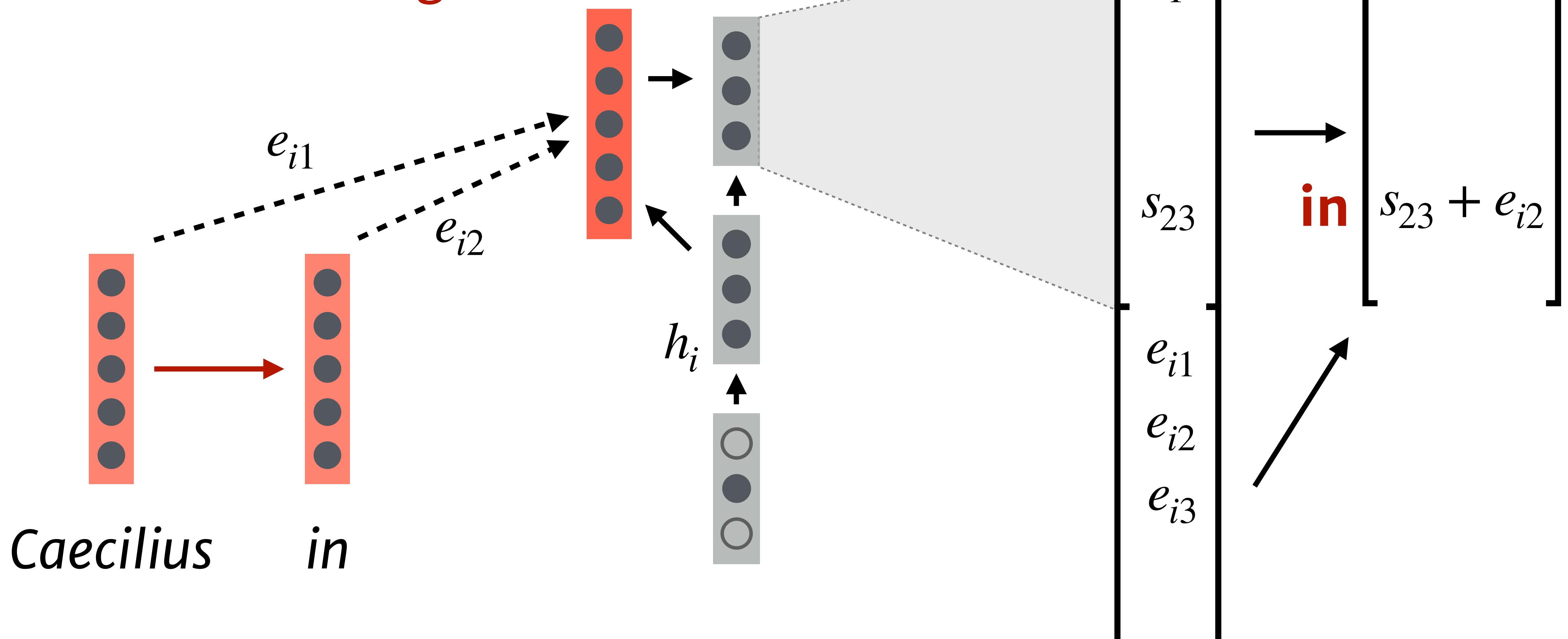
We want the ability to  
generate *in* via copying  
and direct prediction.

# Copying



# Copying

*In* is double-counted:  
just add scores together



# Hard attention

---

When predicting output  $i$ , assign a weight  $\alpha_{ij}$  to each encoder state  $h_j$

$$e_{ij} = w^\top \tanh(W^\top [h_i, h_j])$$

[Bahdanau 2014]

$$e_{ij} = h_i^\top W h_j$$

[Luong 2015]

attention       $\alpha_i = \text{argmax}(e_{i:})$

context repr     $c_i = h_{\alpha_i}$

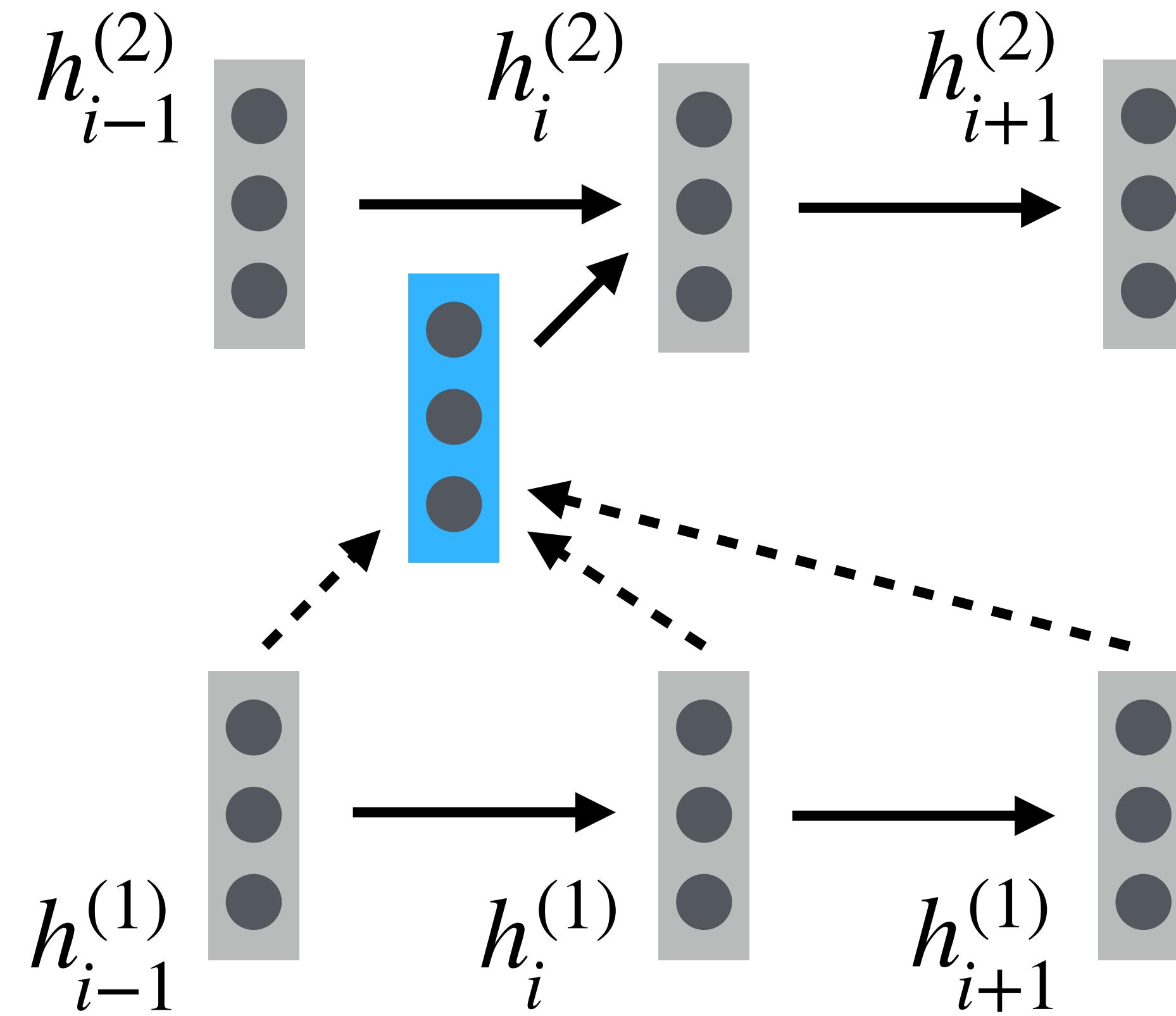
**nondifferentiable!**

**but sometimes better generalization**

now you know how to  
build [almost] anything

# Self-attention revisited

---



# Next class: transformers