

Attention Mechanisms

Jacob Andreas / MIT 6.804-6.864 / Spring 2020

Admin

HW1 is done! Look out for survey.

HW2b will be released **tonight**.

6.864 students only!

Peer reviews will be assigned on OpenReview **tomorrow**.

Each student will get 2 papers to review.

Plan to spend ~15min / paper.

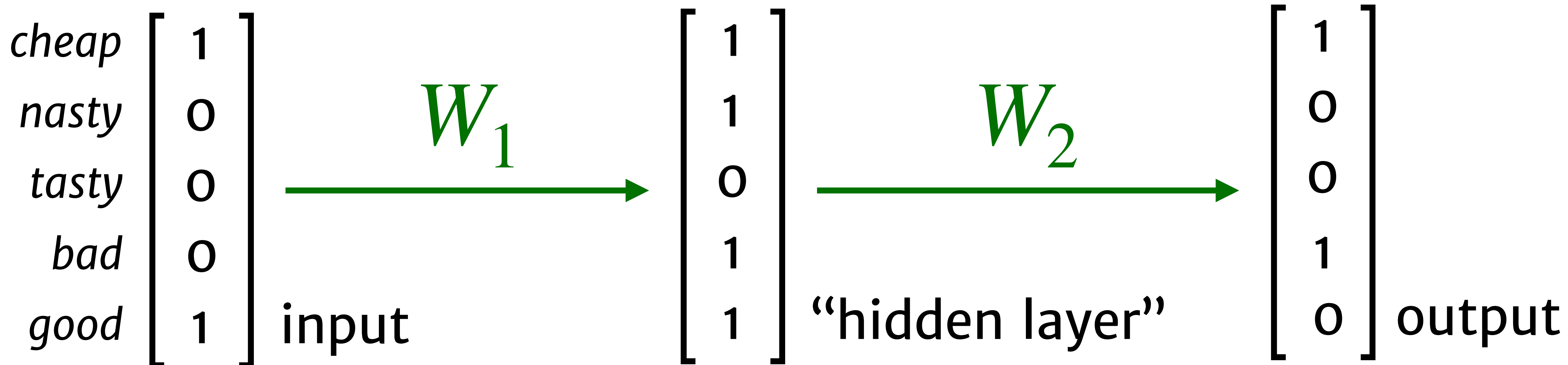
Recap: recurrent neural networks

Neural networks

$$s = W_2^T f(W_1^T x)$$

$$f(W_1^T x) = h_1$$

$$W_2^T h_1 = s$$



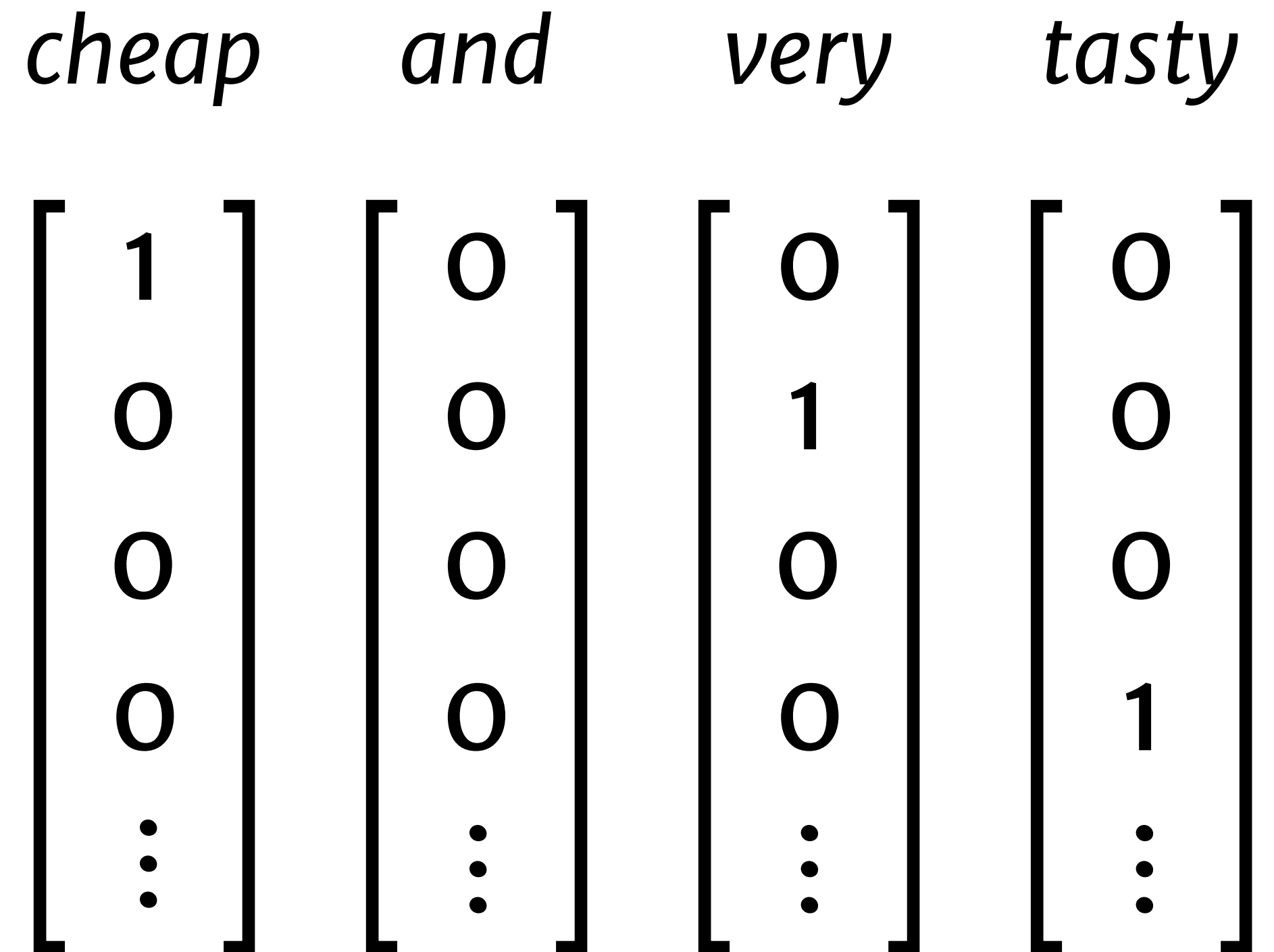
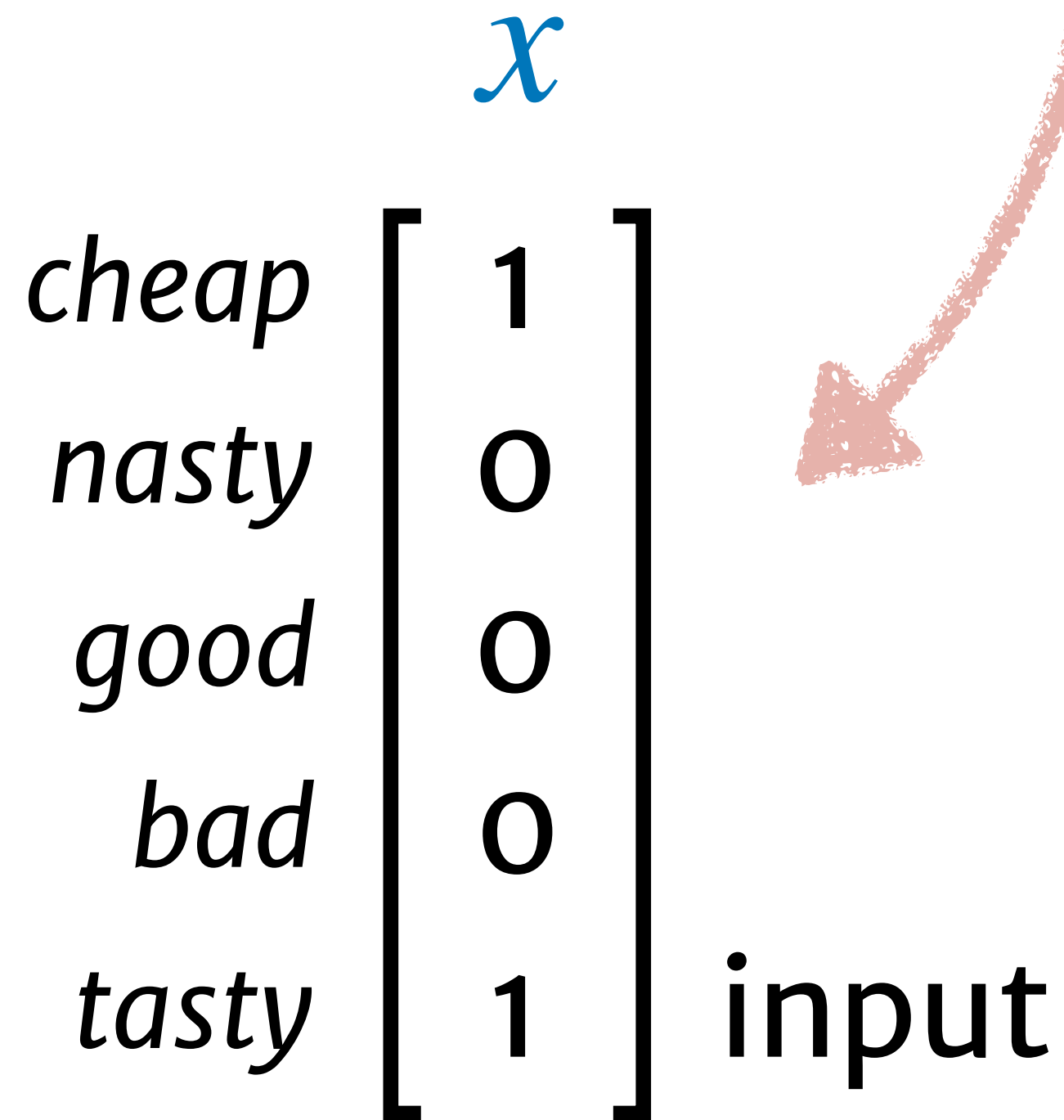
Variable-sized inputs

No ordering information!



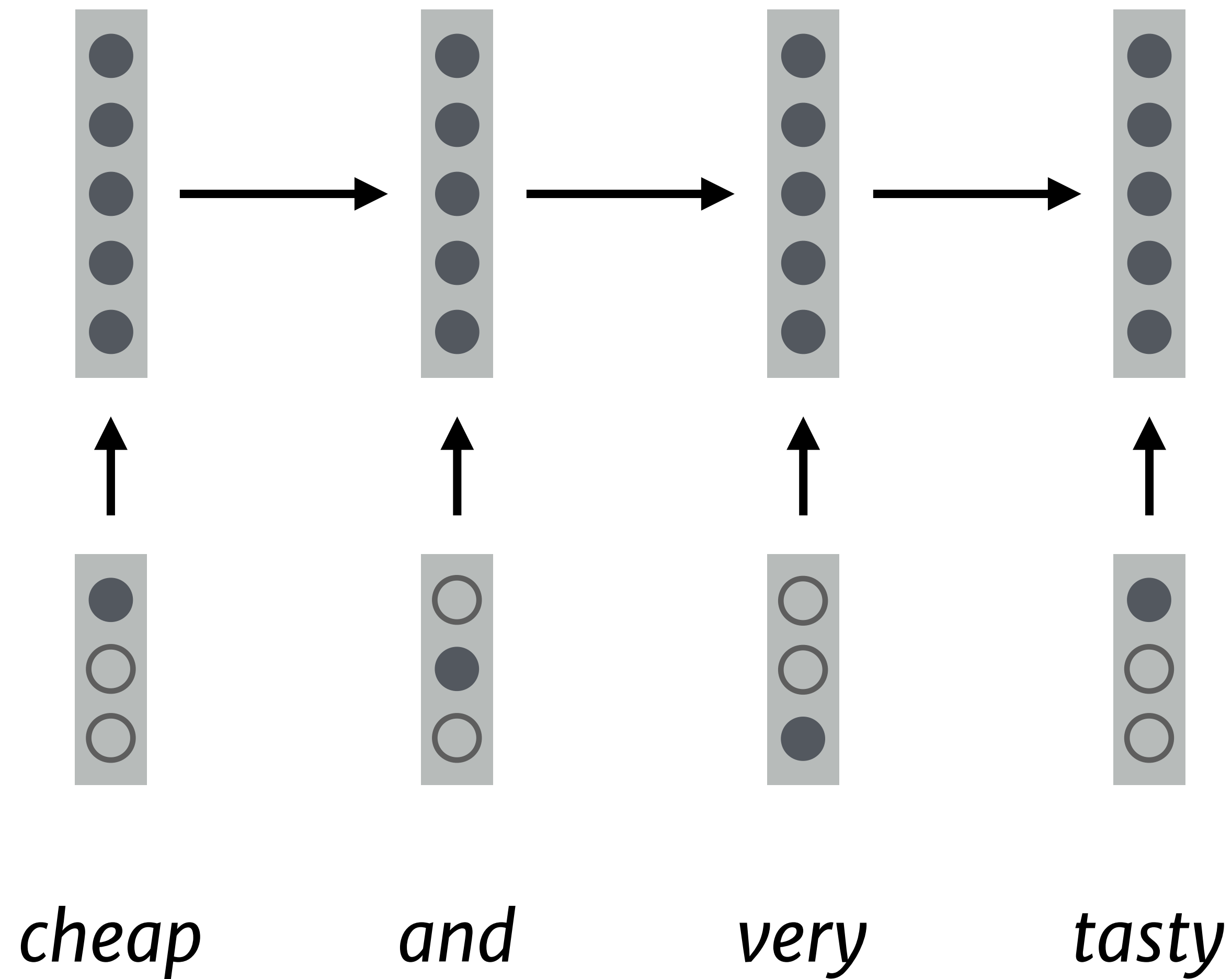
Variable-sized inputs

No ordering information!



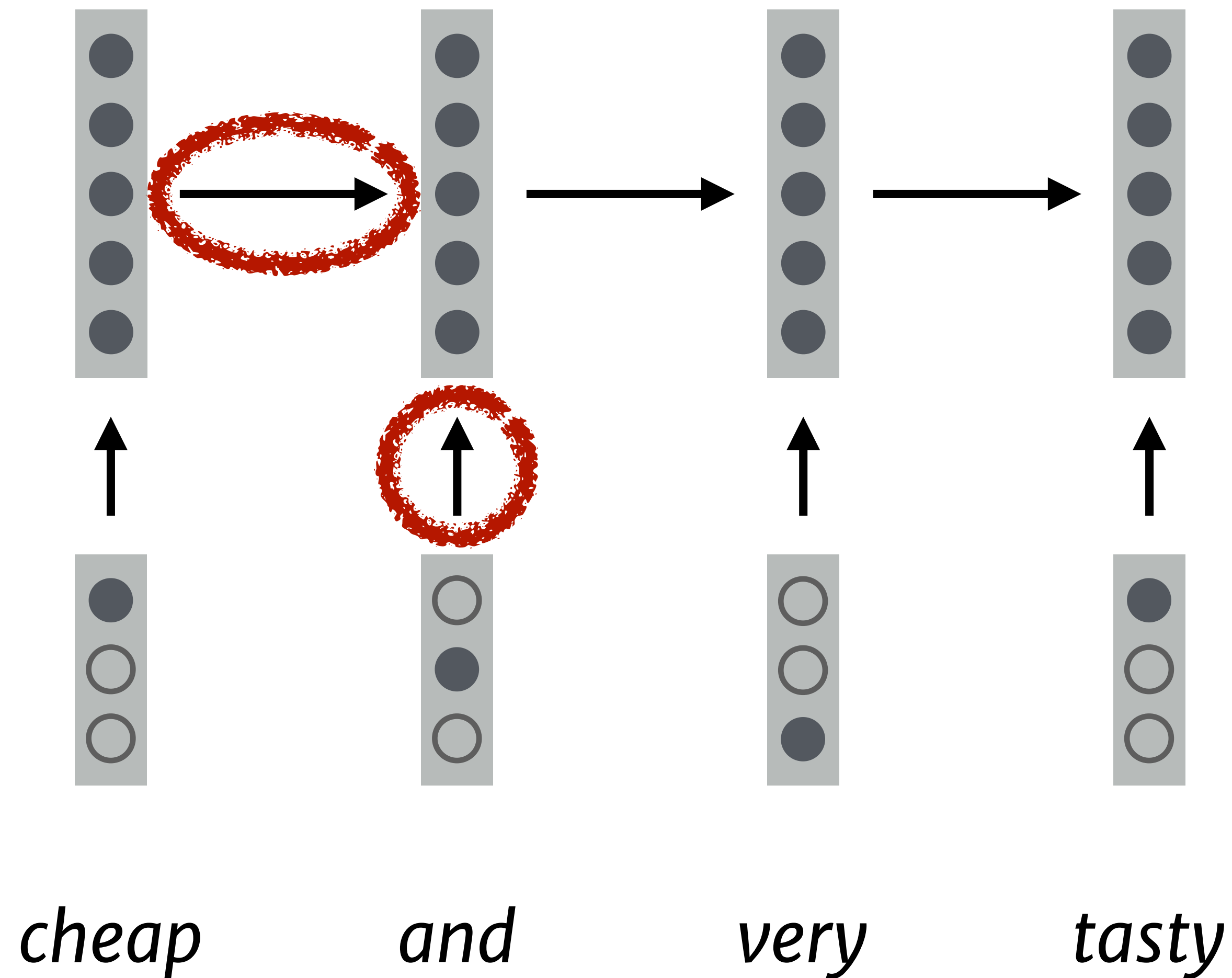
No fixed input dimension!

Recurrent neural networks



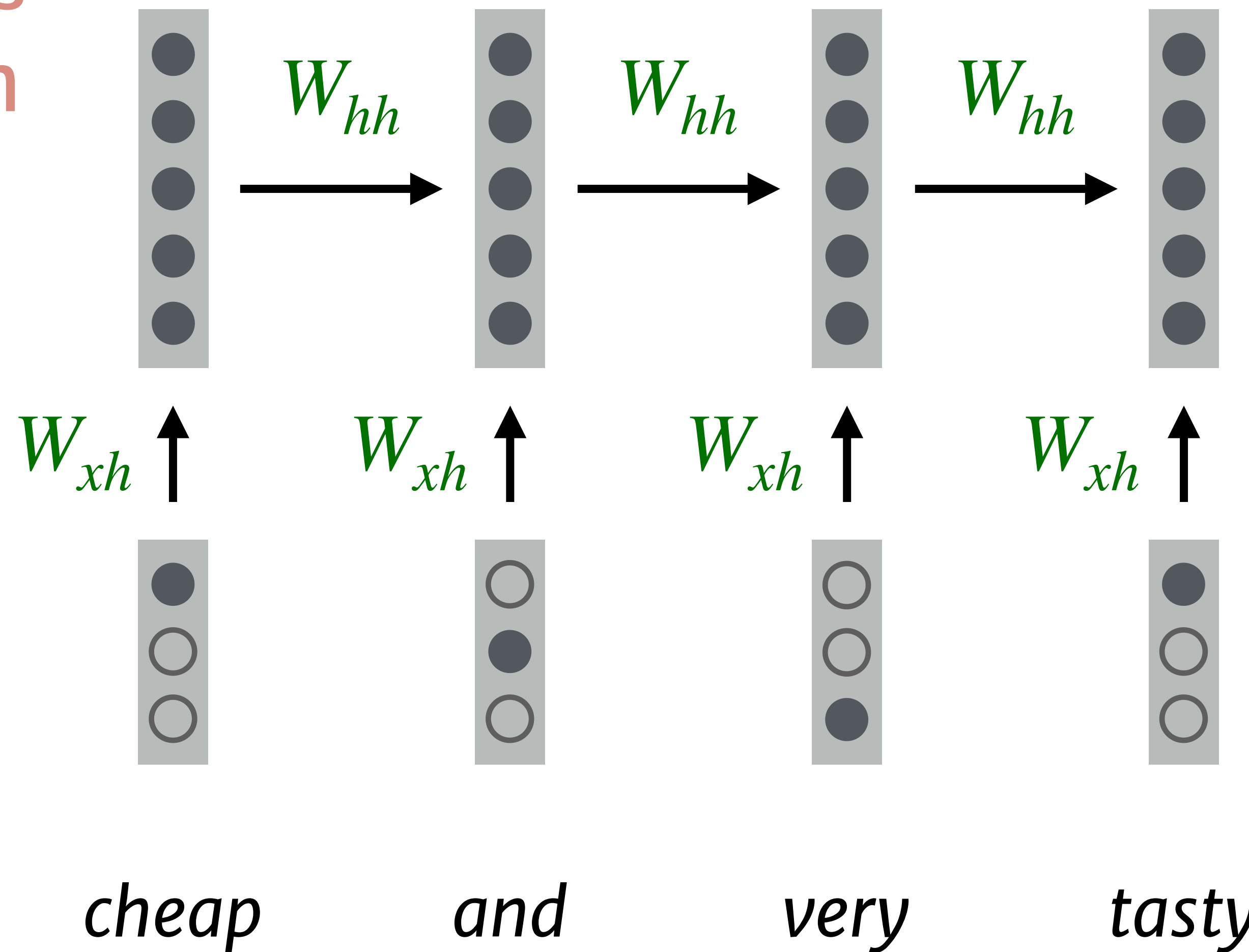
Recurrent neural networks

Hidden states depend on an earlier state and an input



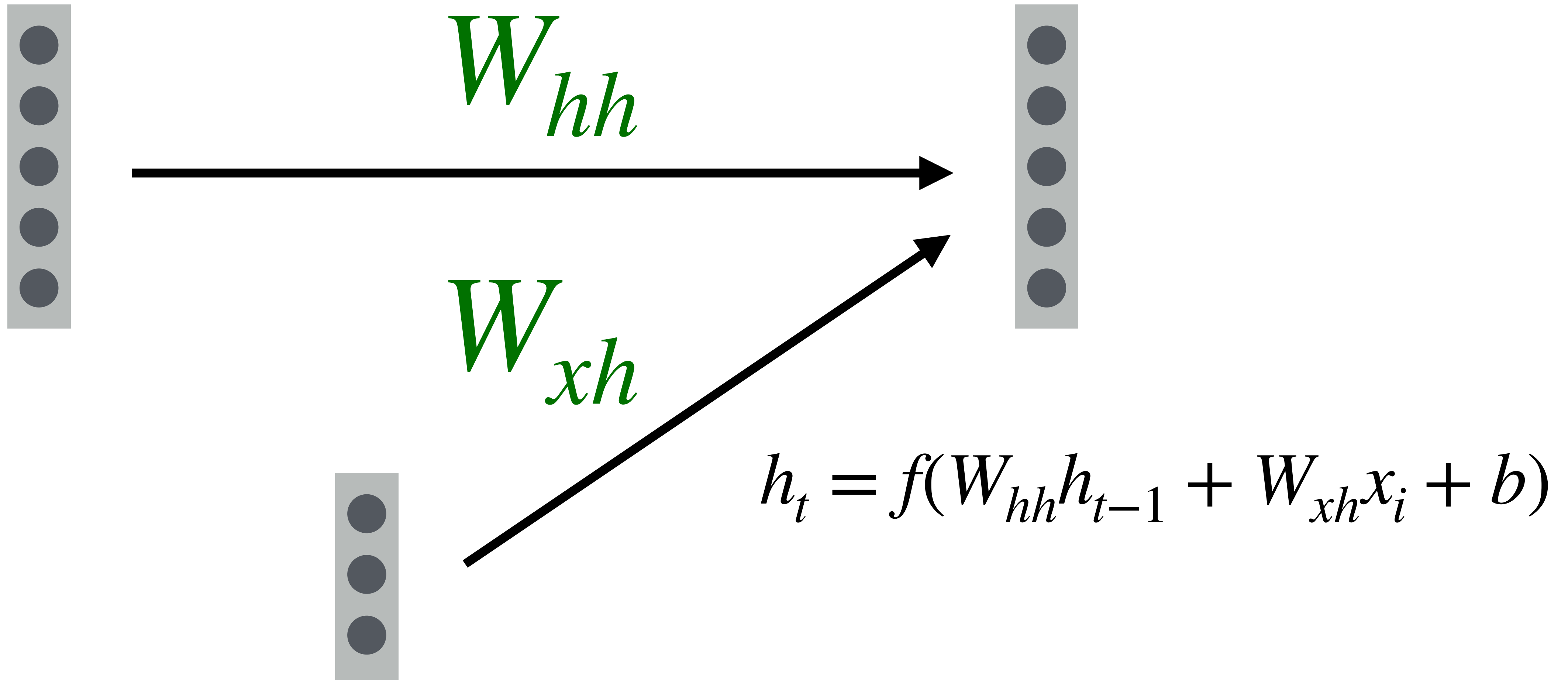
Recurrent neural networks

Hidden states depend on an earlier state and an input

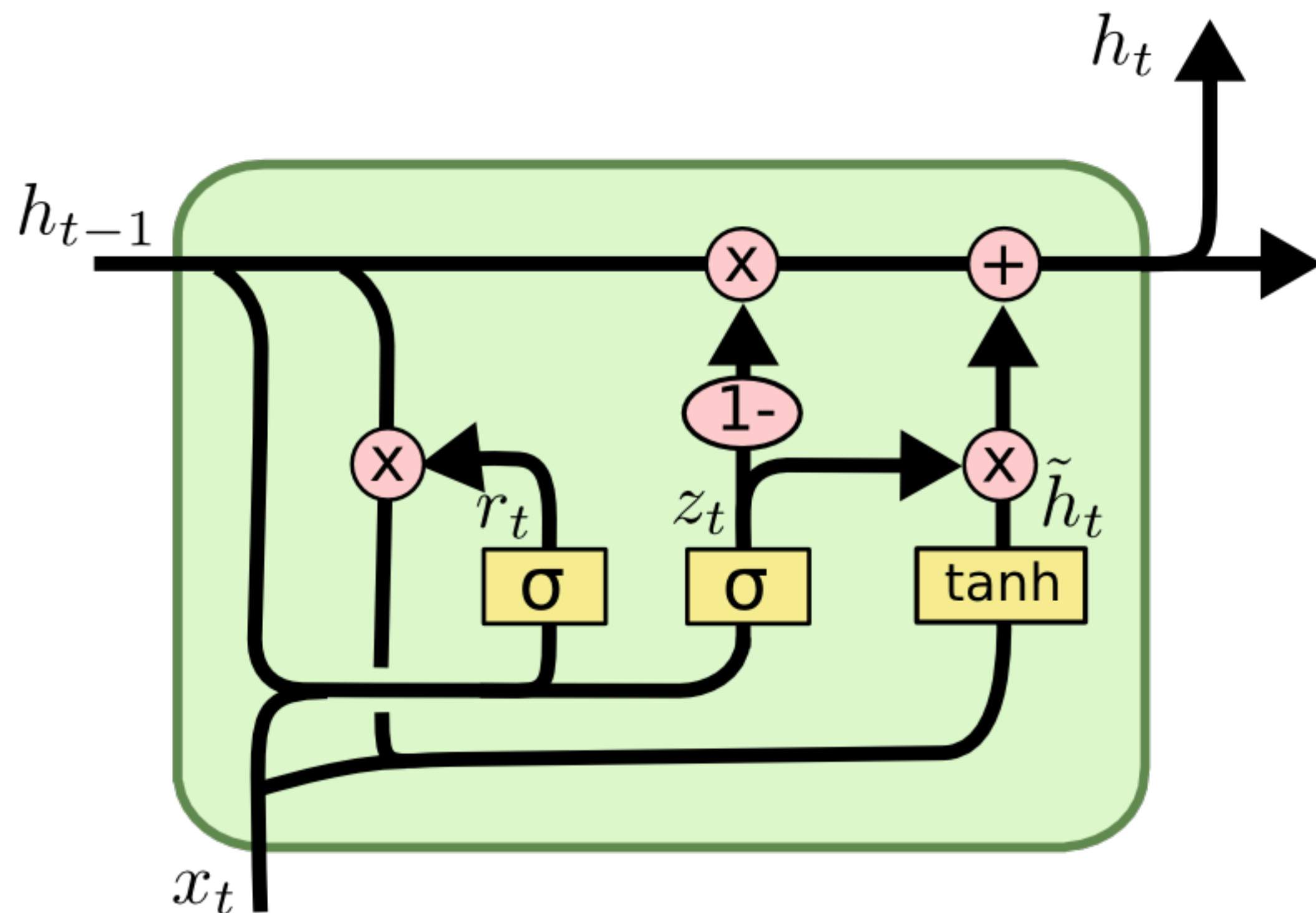


Same weights at every state!

“Vanilla” RNNs



Gated Recurrent Units



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

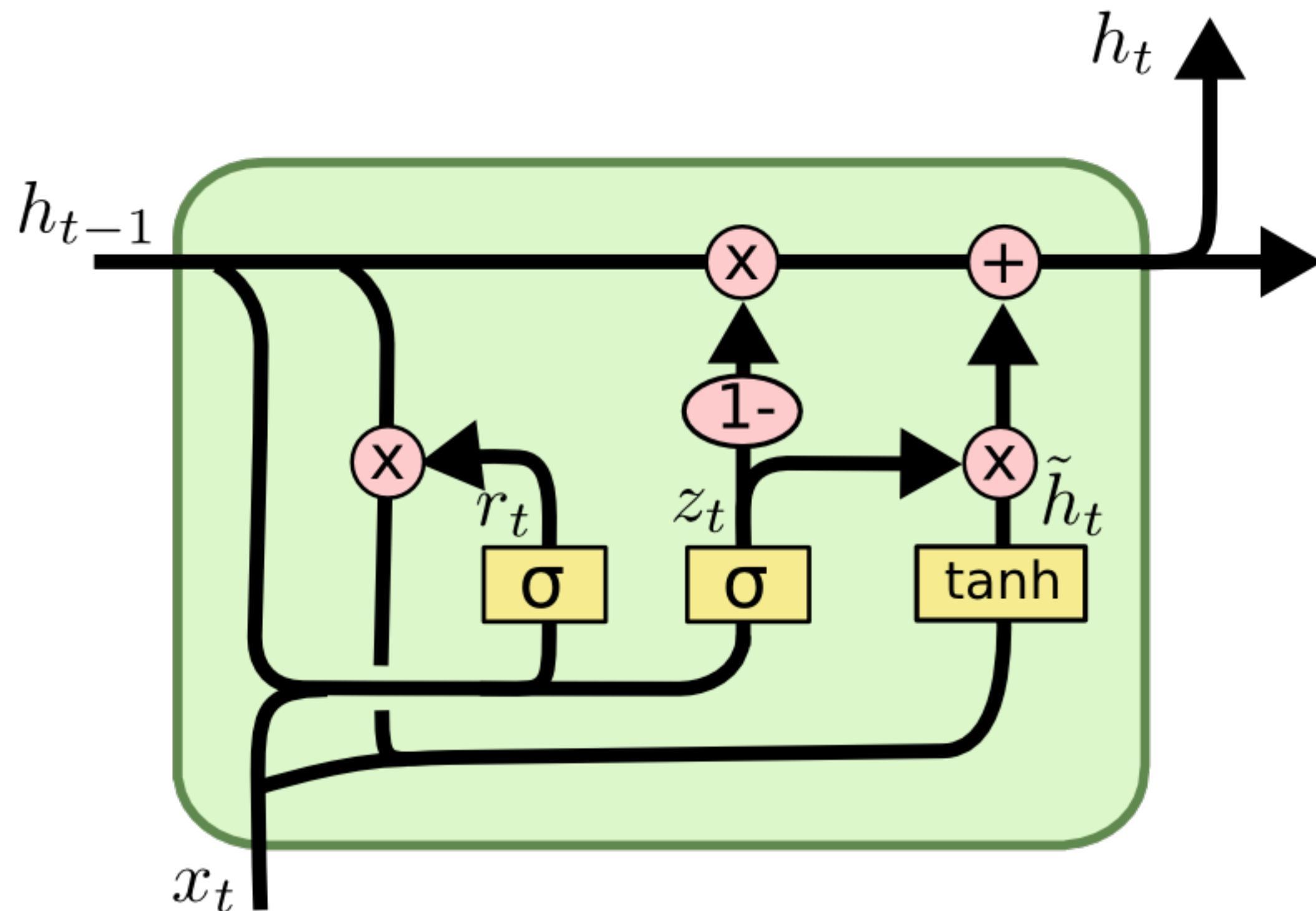
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

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

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

[Image: Cristopher Olah]

Gated Recurrent Units



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

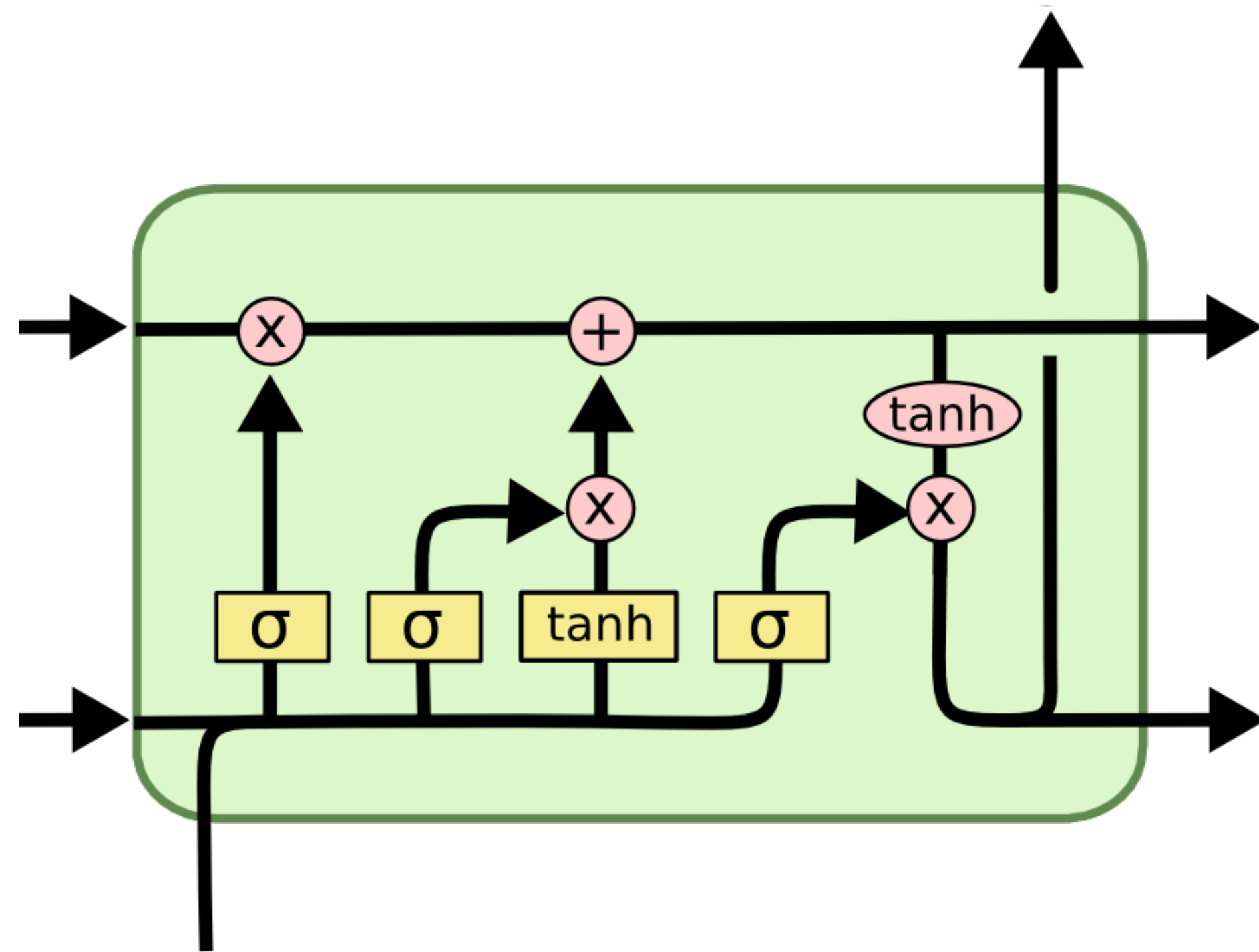
$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

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

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

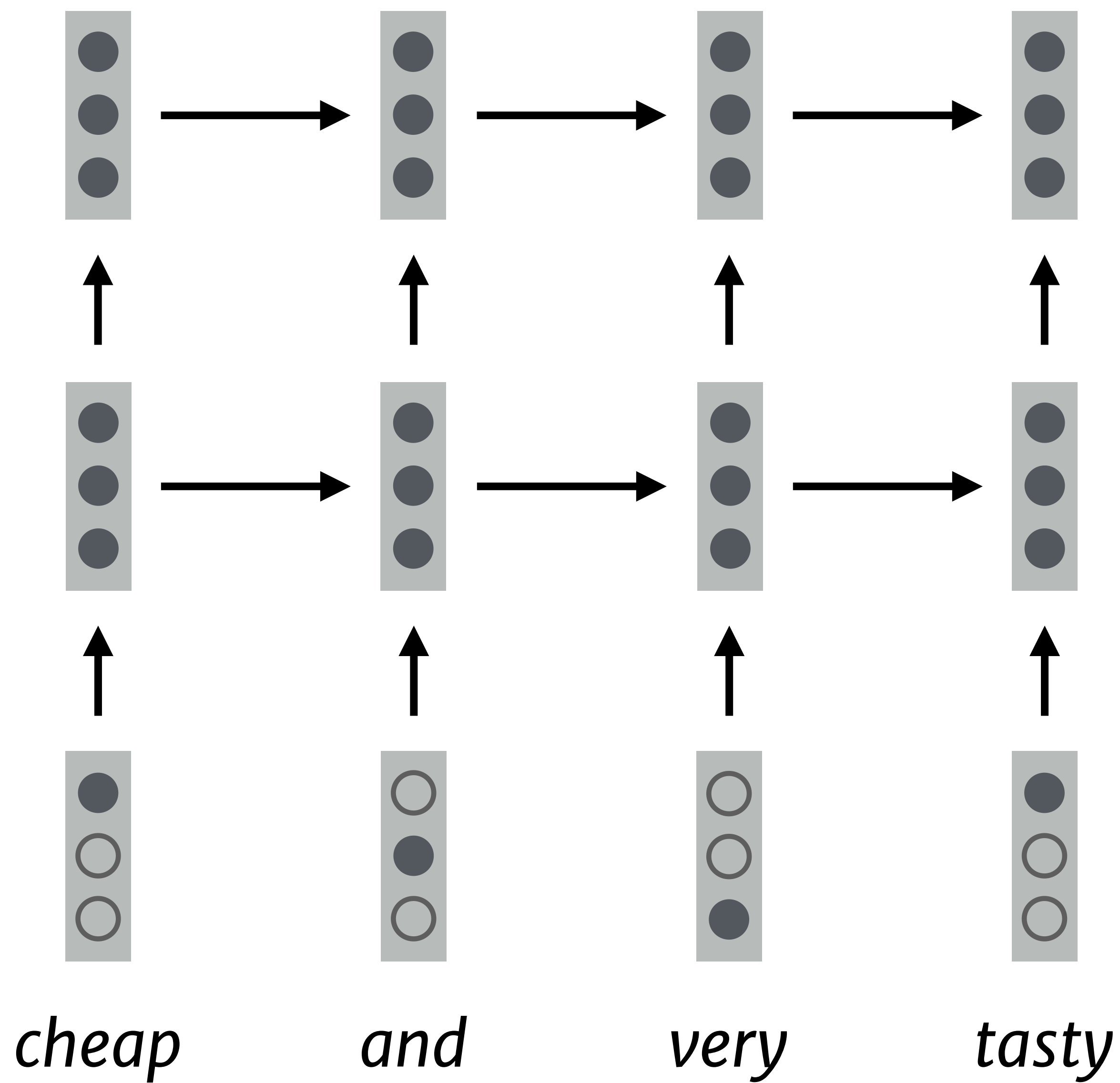
[Image: Cristopher Olah]

Long Short-Term Memory Units

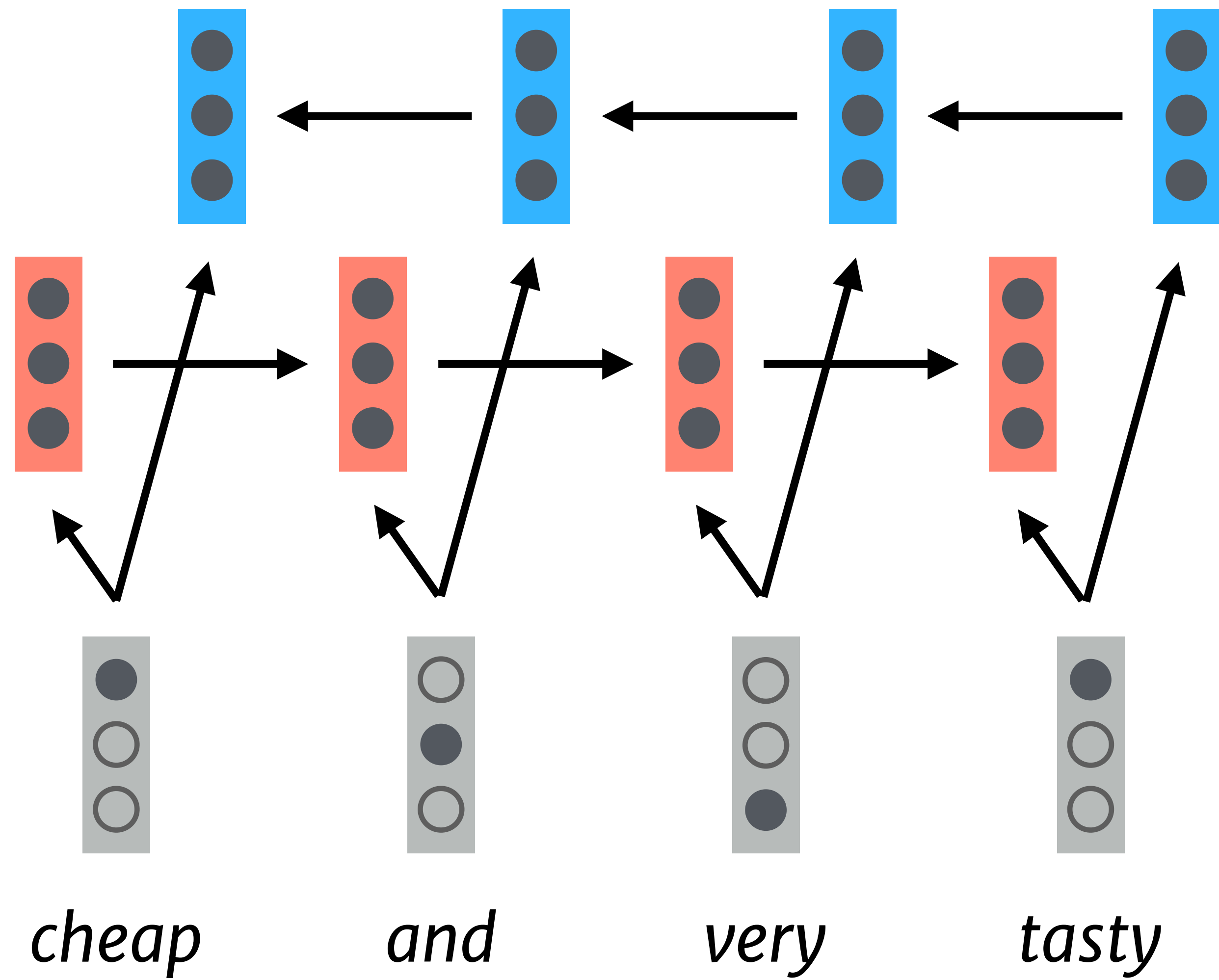


[Image: Cristopher Olah]

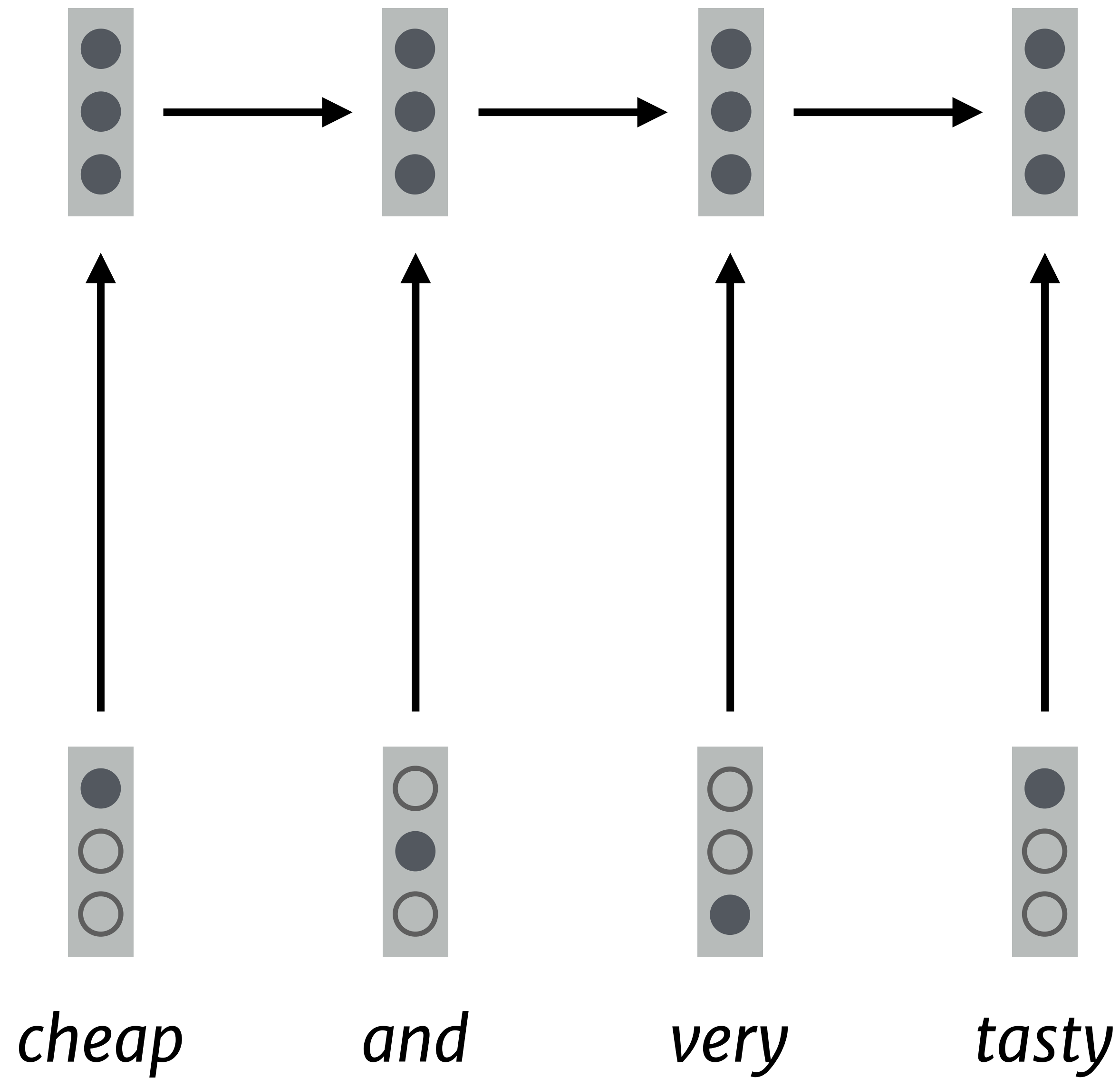
Deeper RNNs



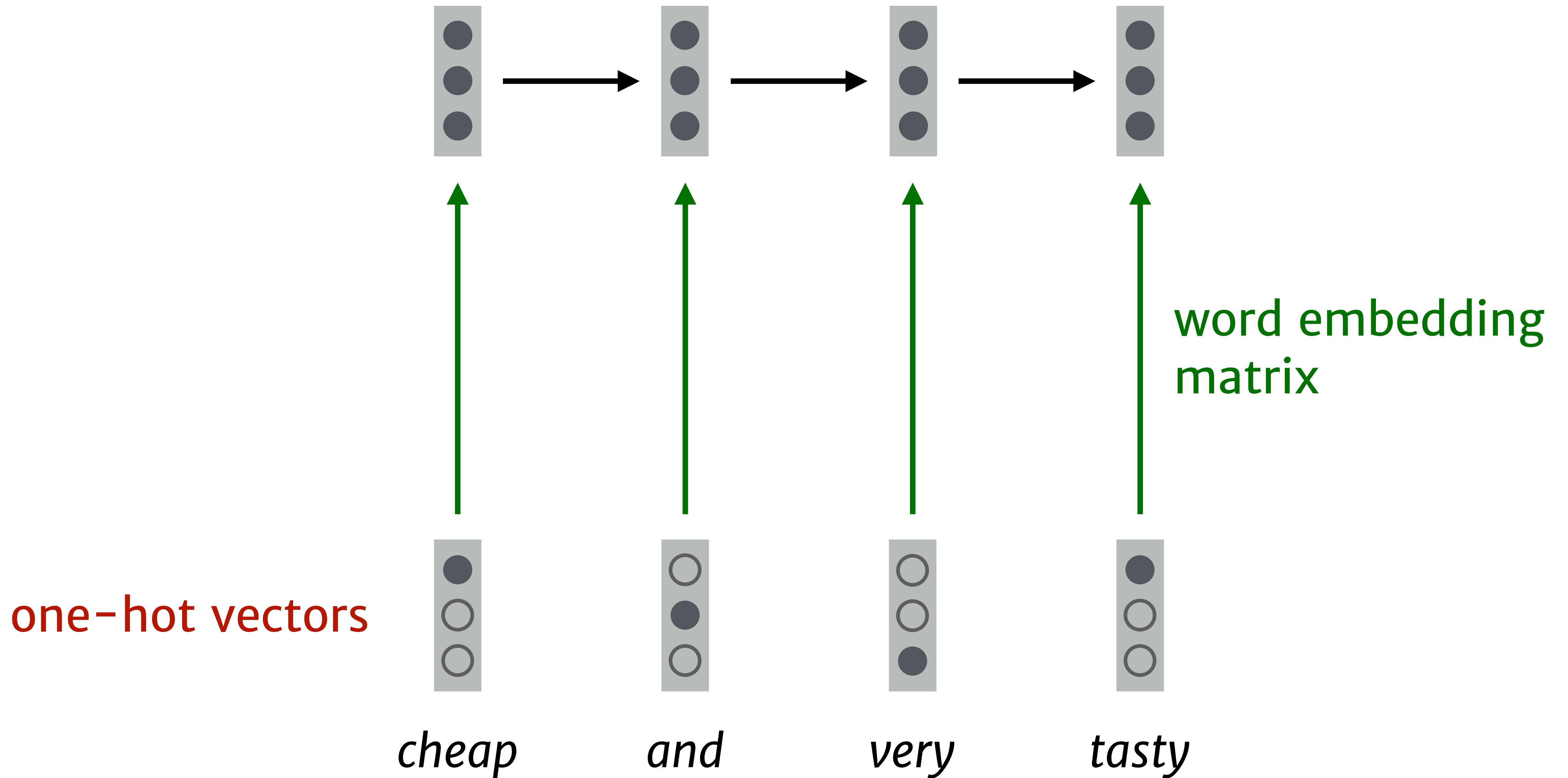
Bidirectional RNNs



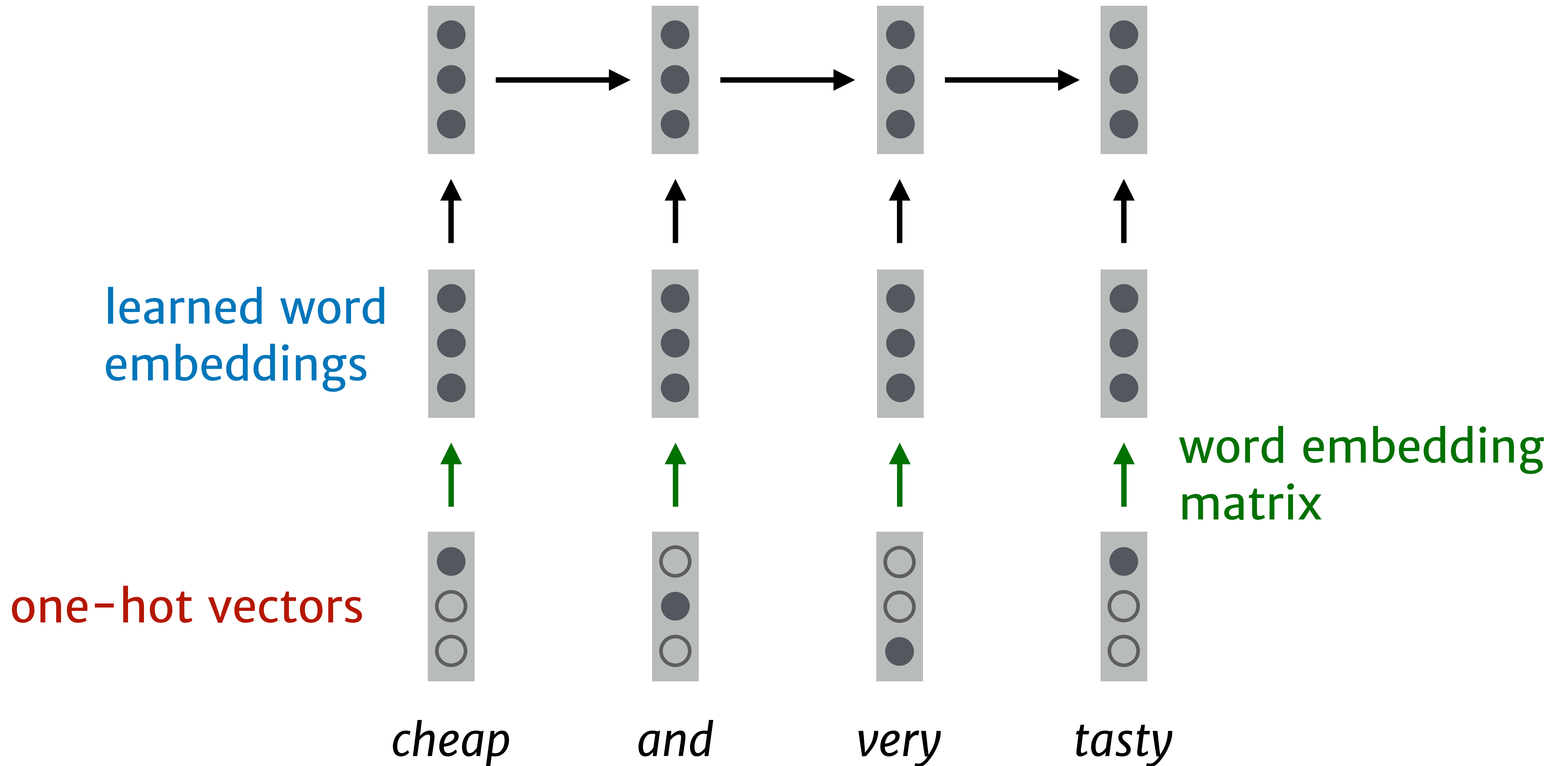
RNNs and word embeddings



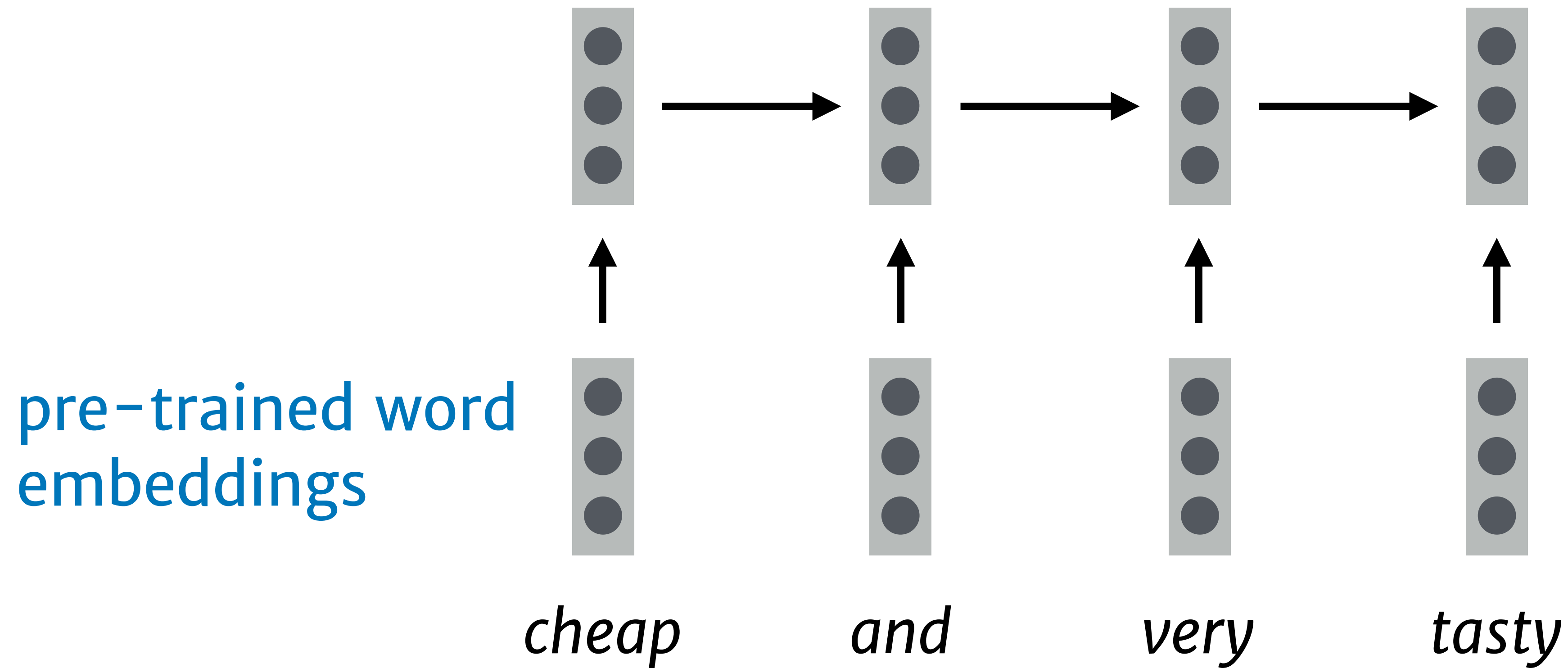
RNNs and word embeddings



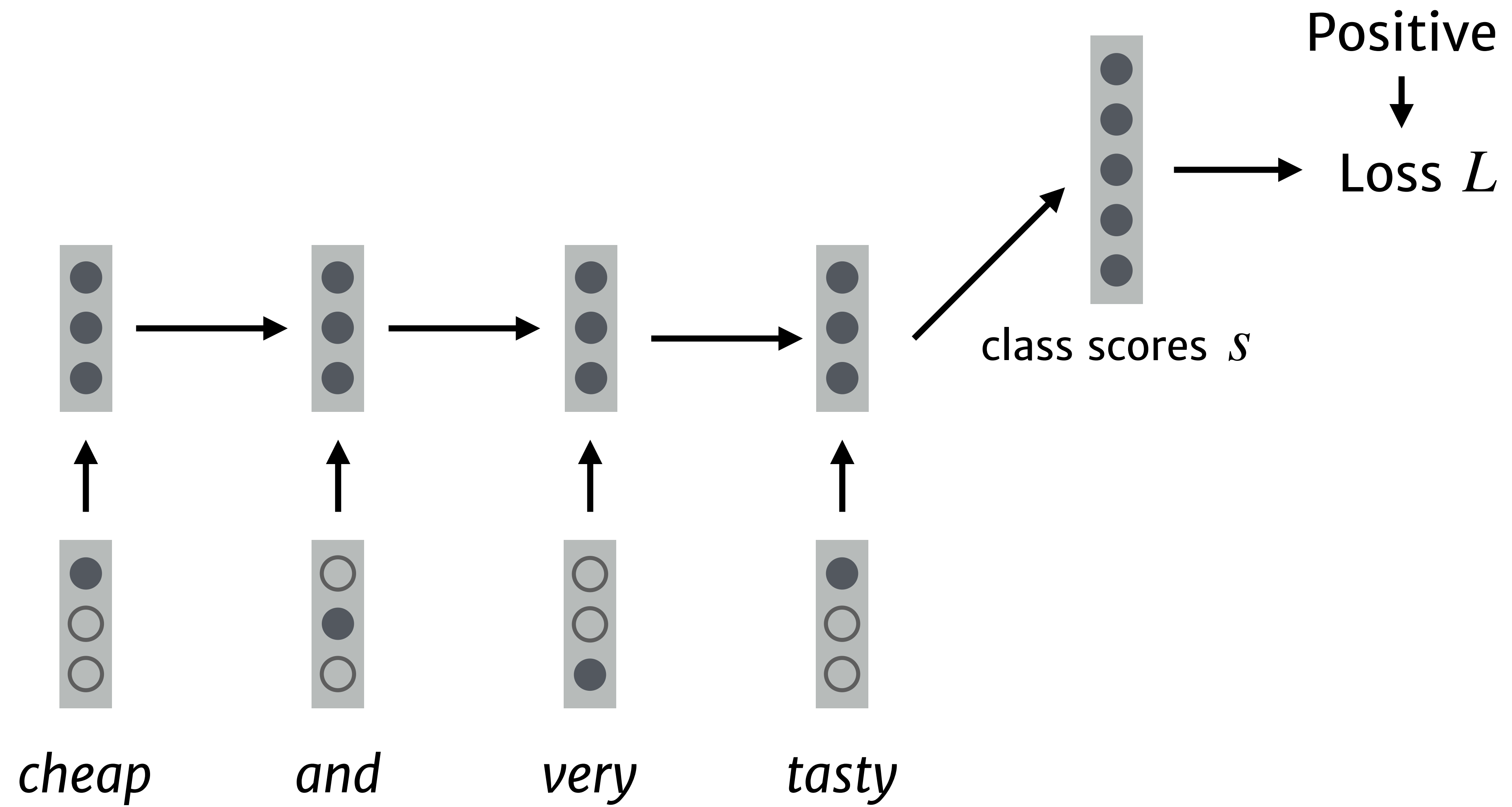
RNNs and word embeddings



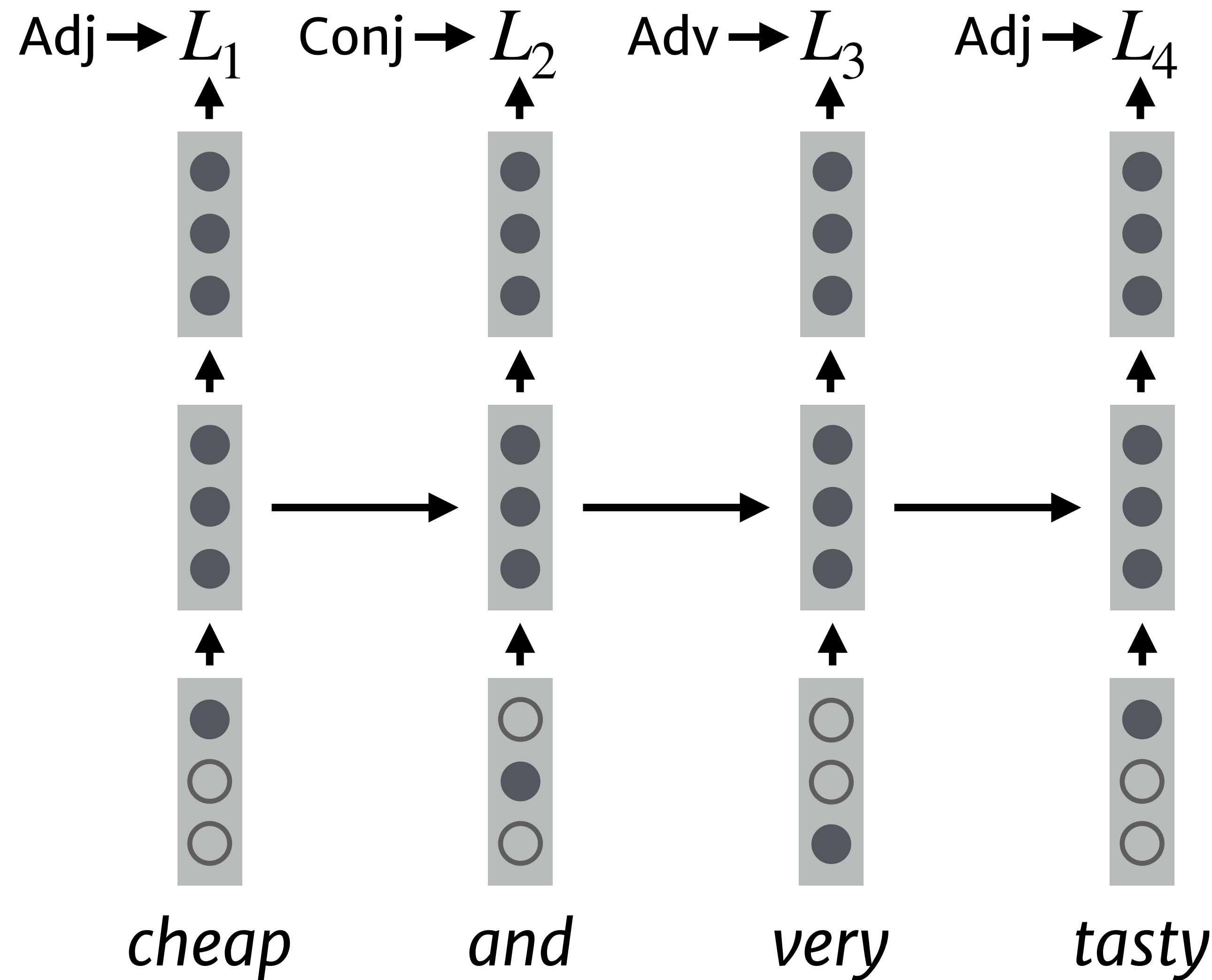
RNNs and word embeddings



Text classification



Sequence labeling

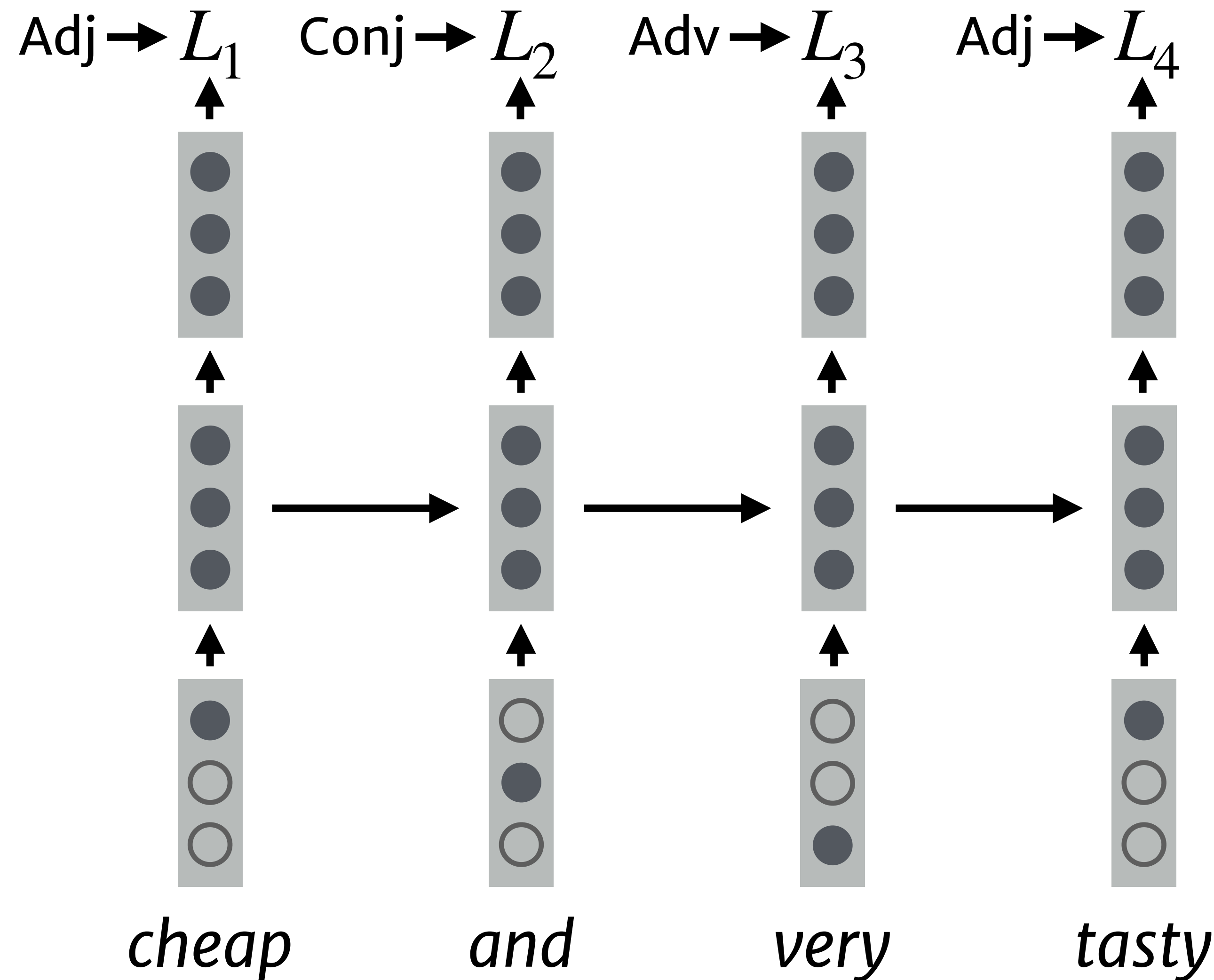


$$L = \sum_t L_t$$

e.g.

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

Sequence labeling



$$L = \sum_t L_t$$

e.g.

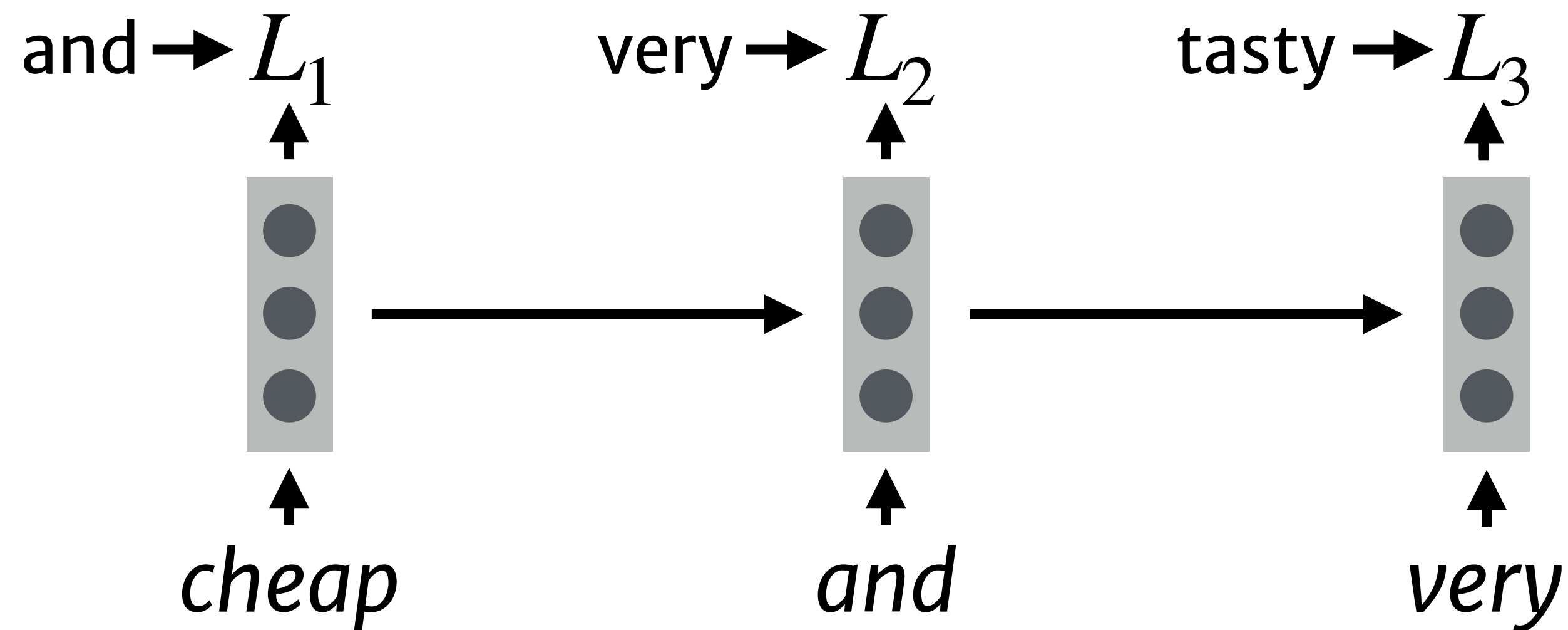
$$= - \sum_t \log p(y_t | x_{:t})$$

product of indep.
conditionals!

Language Modeling

A (unidirectional) can compute $p(y_t | x_{:t})$.

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

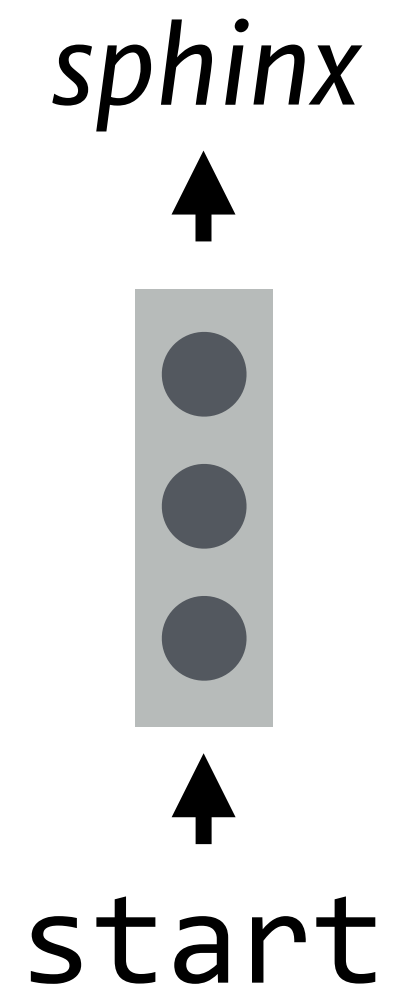


then $\sum_t \log p(x_{t+1} | x_{:t}) = p(x)$

Language modeling: sampling

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

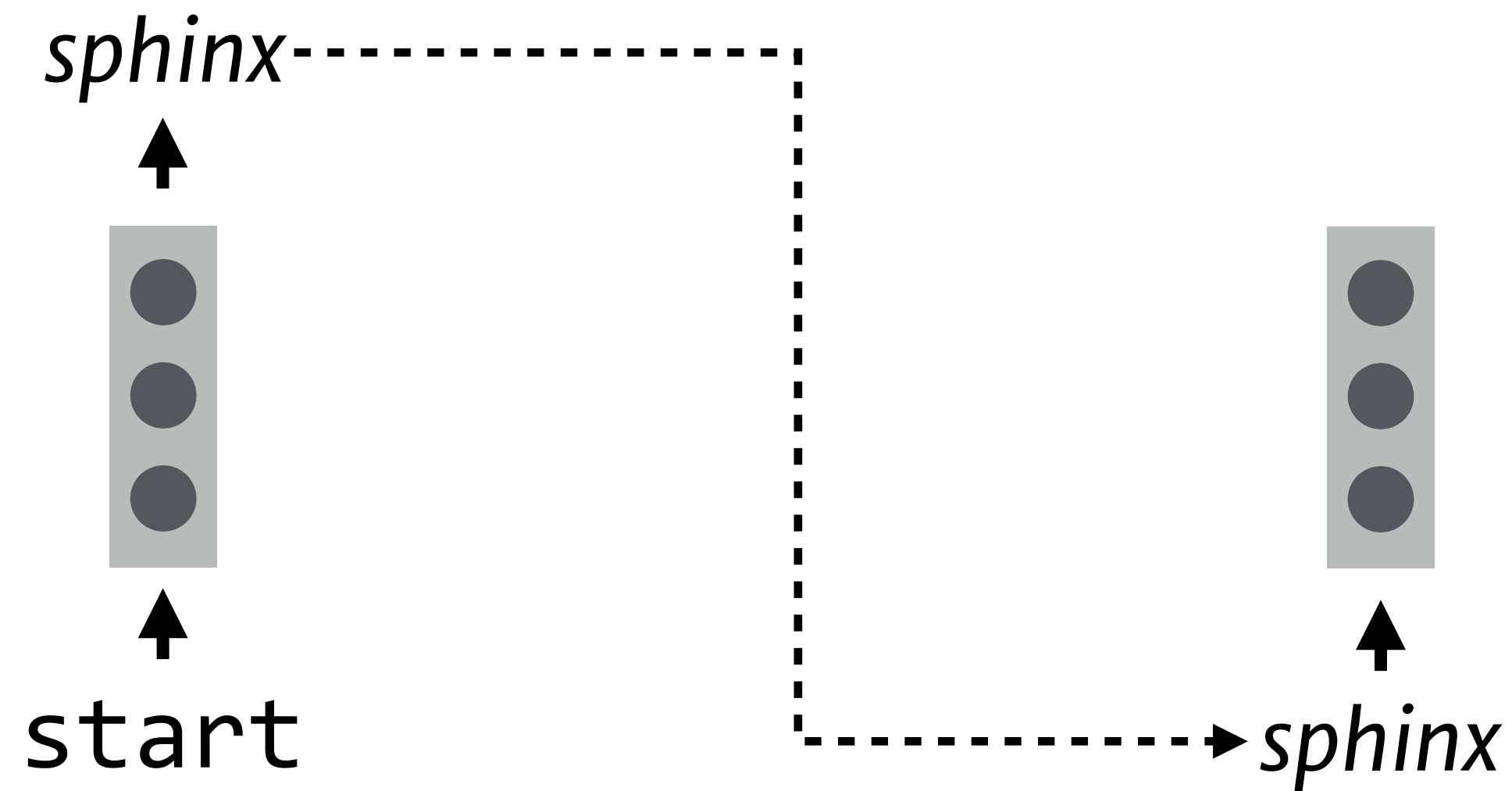
$$x_1 \sim p(x_1 | \text{start})$$



Language modeling: sampling

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

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

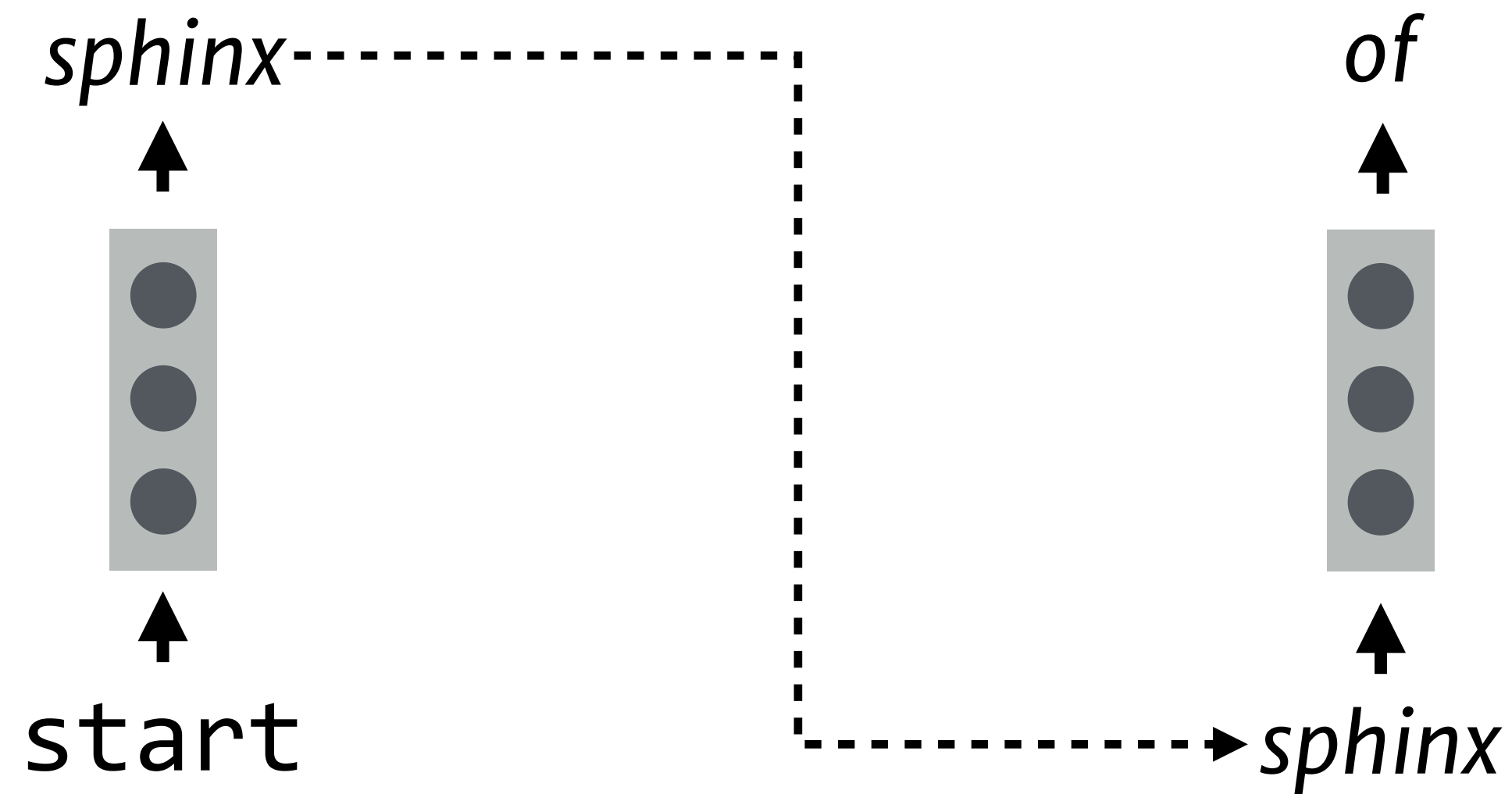


Language modeling: sampling

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

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

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



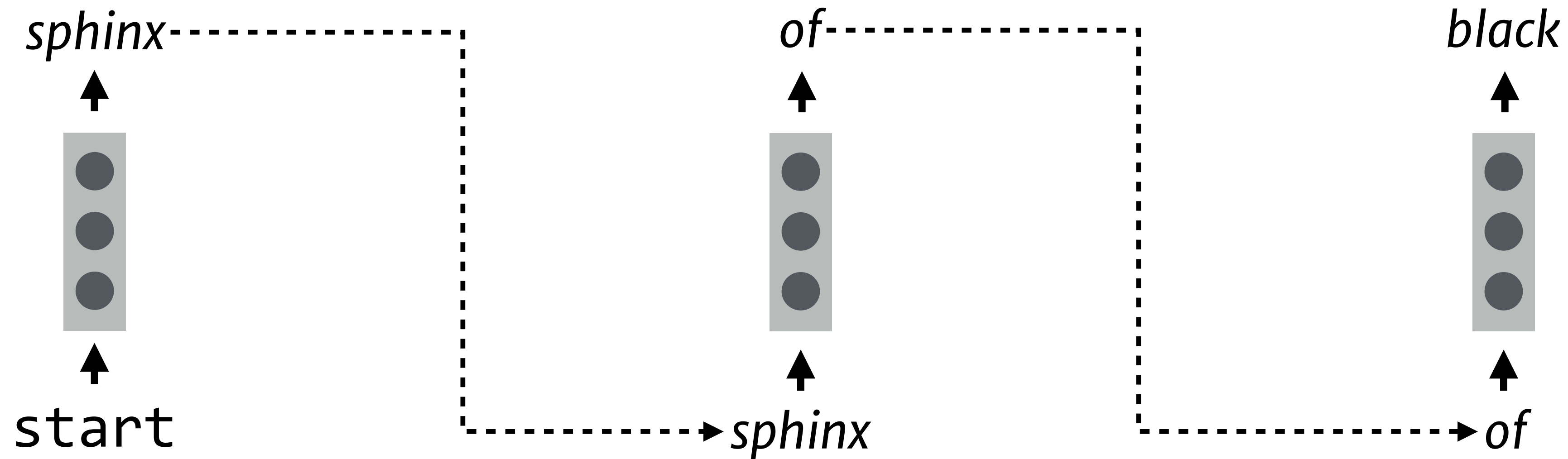
Language modeling: sampling

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

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

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

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

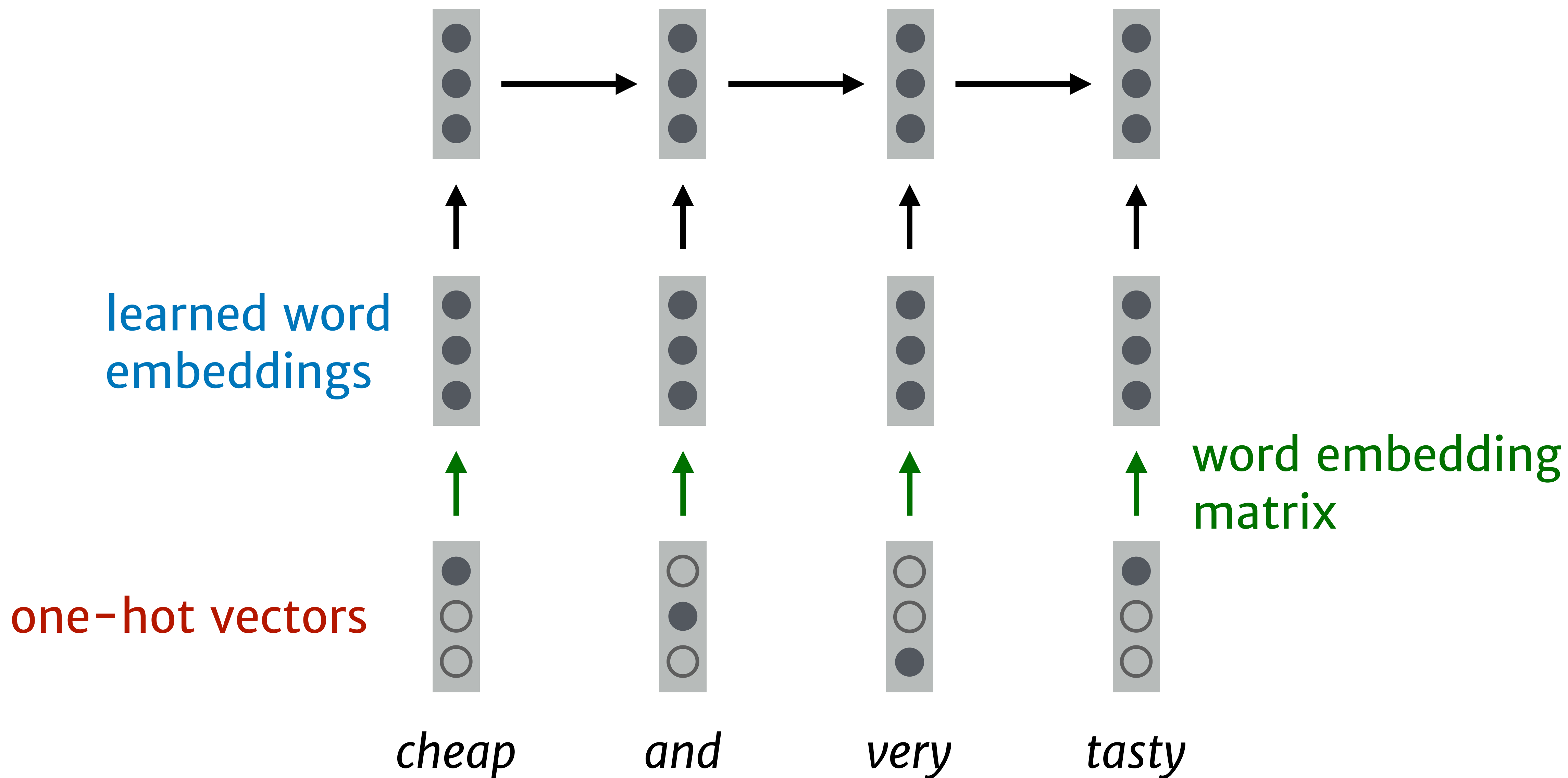


SPHINX OF BLACK QUARTZ, JUDGE MY VOW



[Image: egyptianmarketplace.com]

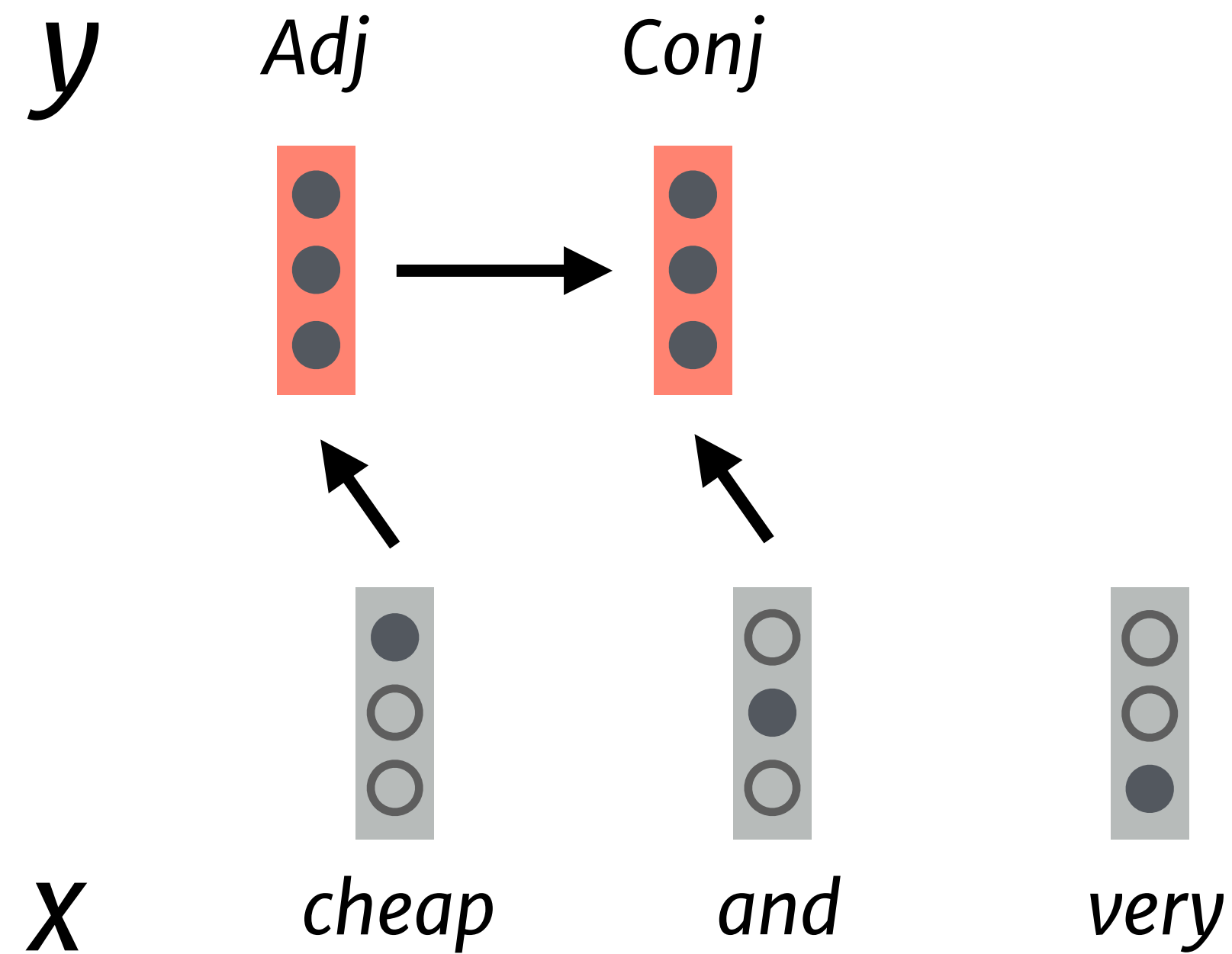
Language modeling as representation learning



RNNs as Markov chains

I can train this network to predict:

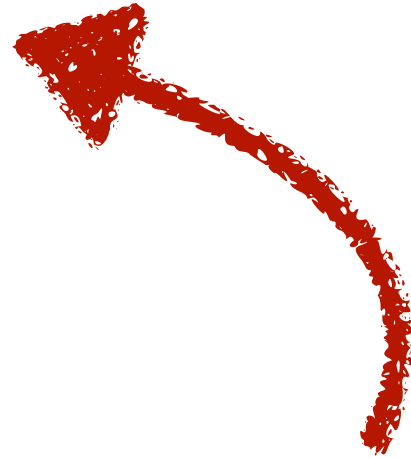
$$\log p(y_t | x_{:t})$$



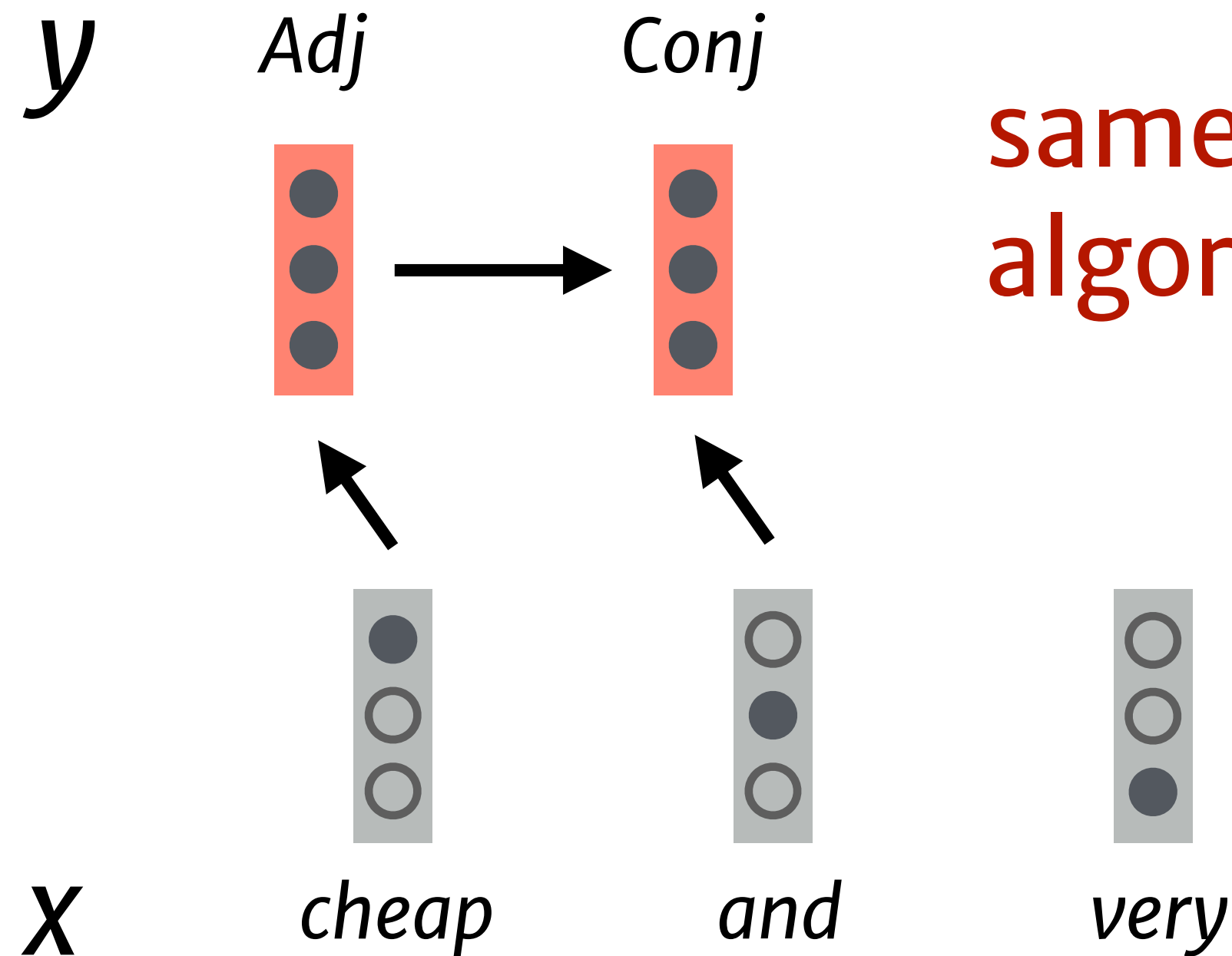
RNNs as Markov chains

I can train this network to predict:

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



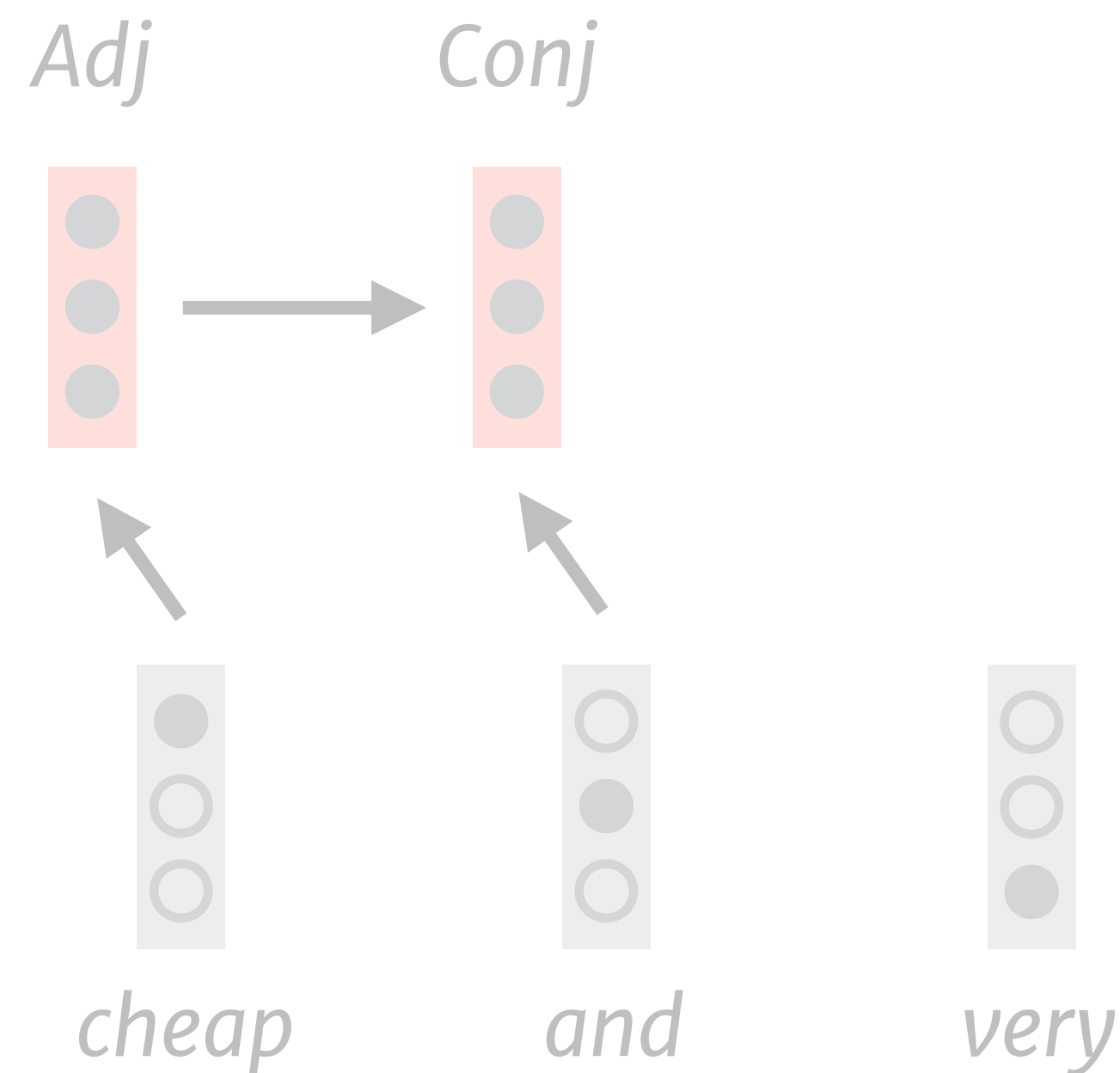
same as forward algorithm!



RNNs as Markov chains

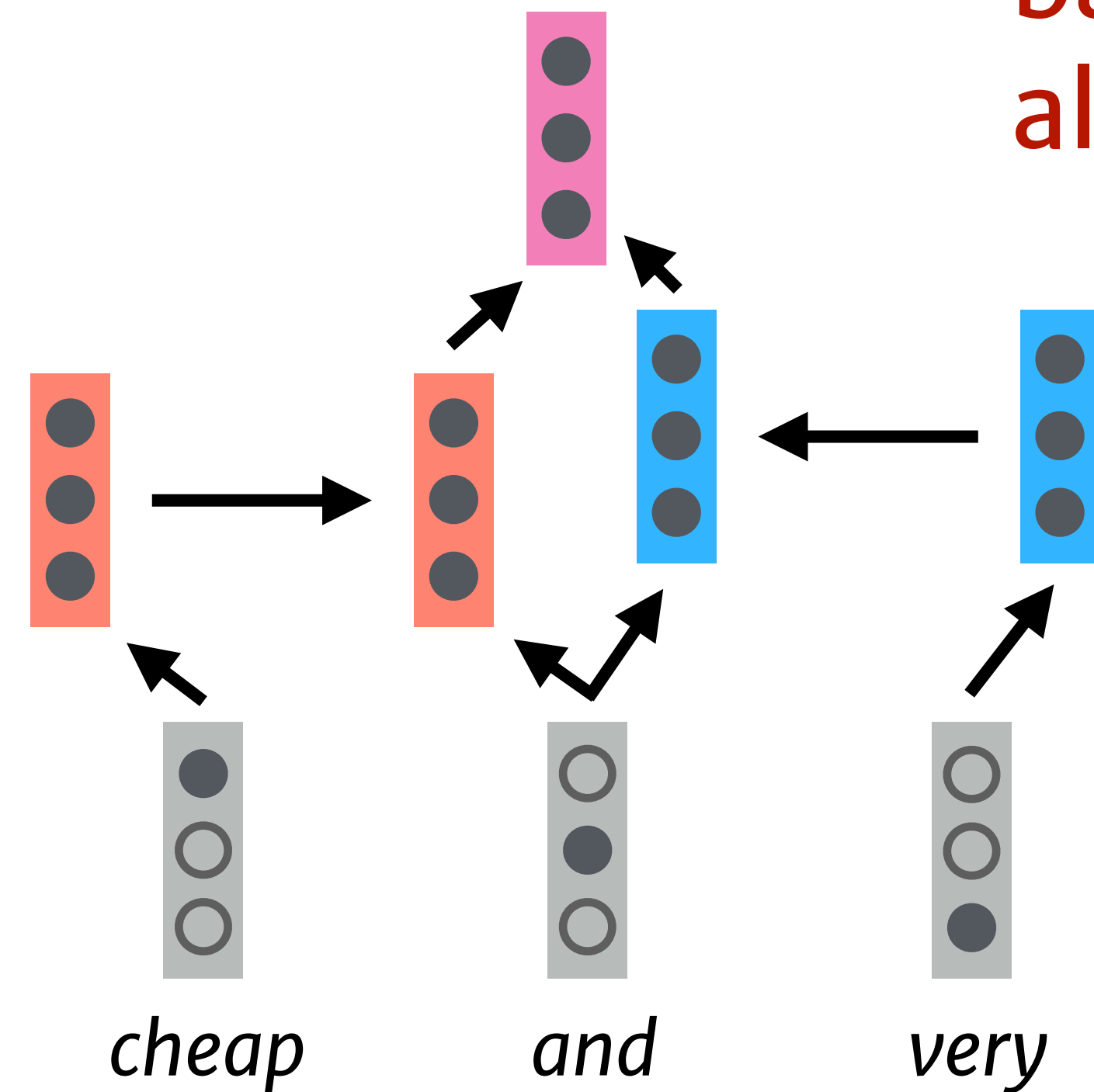
I can train this network to predict:

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



I can train this network to predict:

$$\log p(y_t | x) = p(q_t | O) \text{ forward-backward algo!}$$



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.

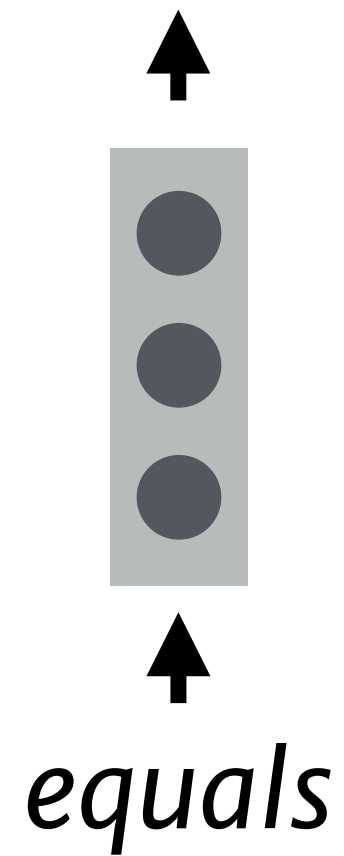
Seven is prime.

One plus two times three equals seven.

Two times three times three equals ???

Answering math problems with LMs

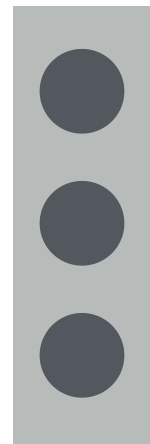
$x_1 \sim p(x_1 | \dots \text{times three equals})$



Answering math problems with LMs

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

twenty

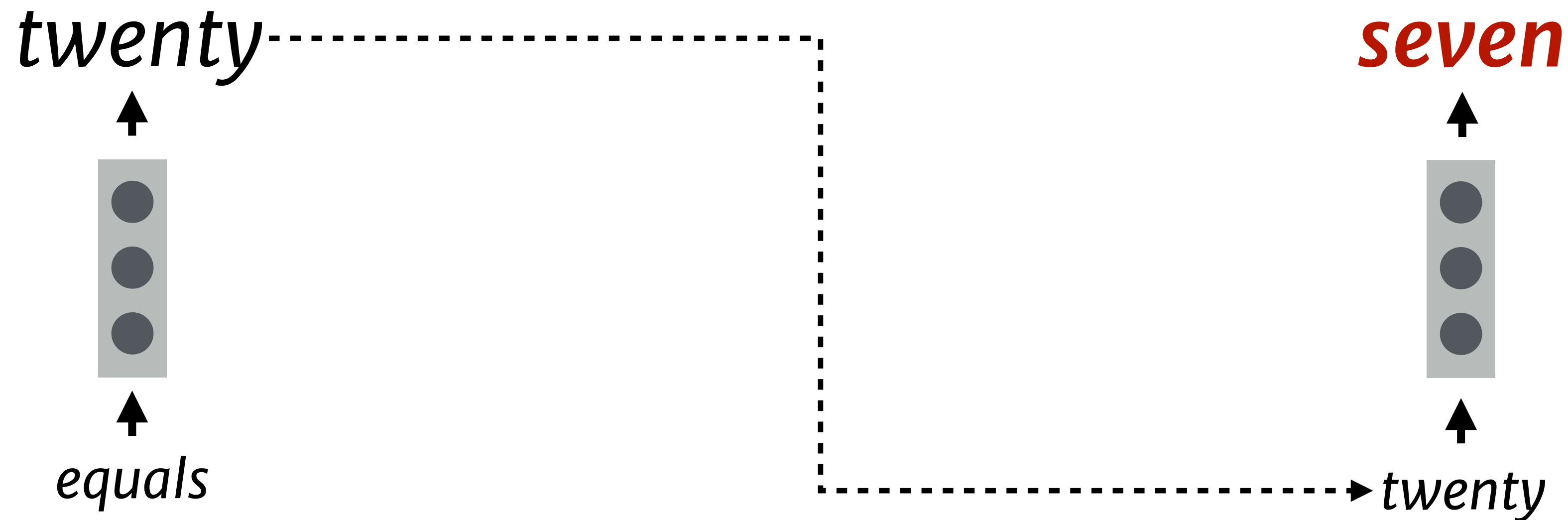


equals

Answering math problems with LMs

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

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



(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.

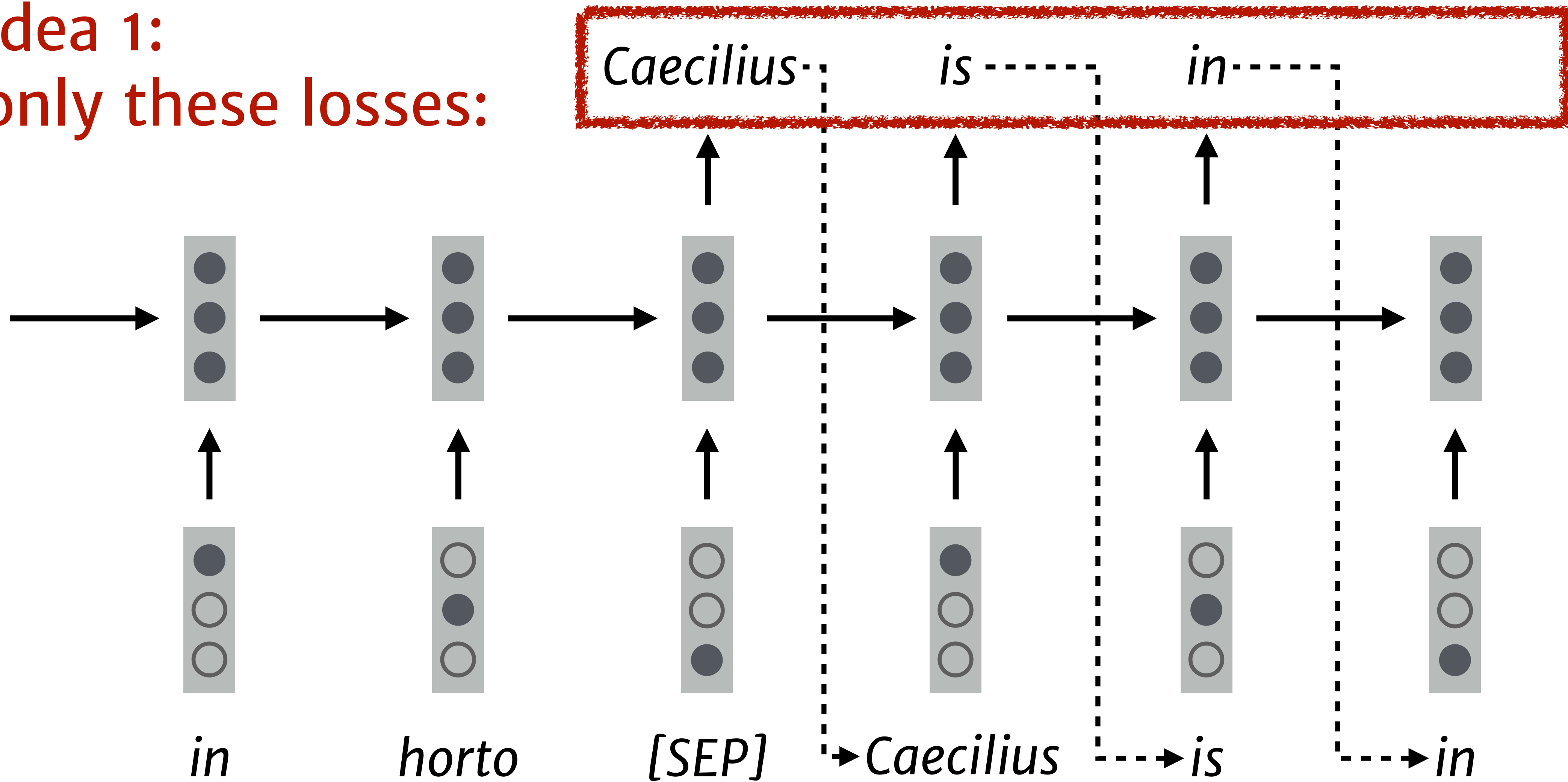
Grumio est in atrio. [SEP] Grumio is in the atrium.

Grumio in atrio laborat. [SEP] ???

(try this at home!)

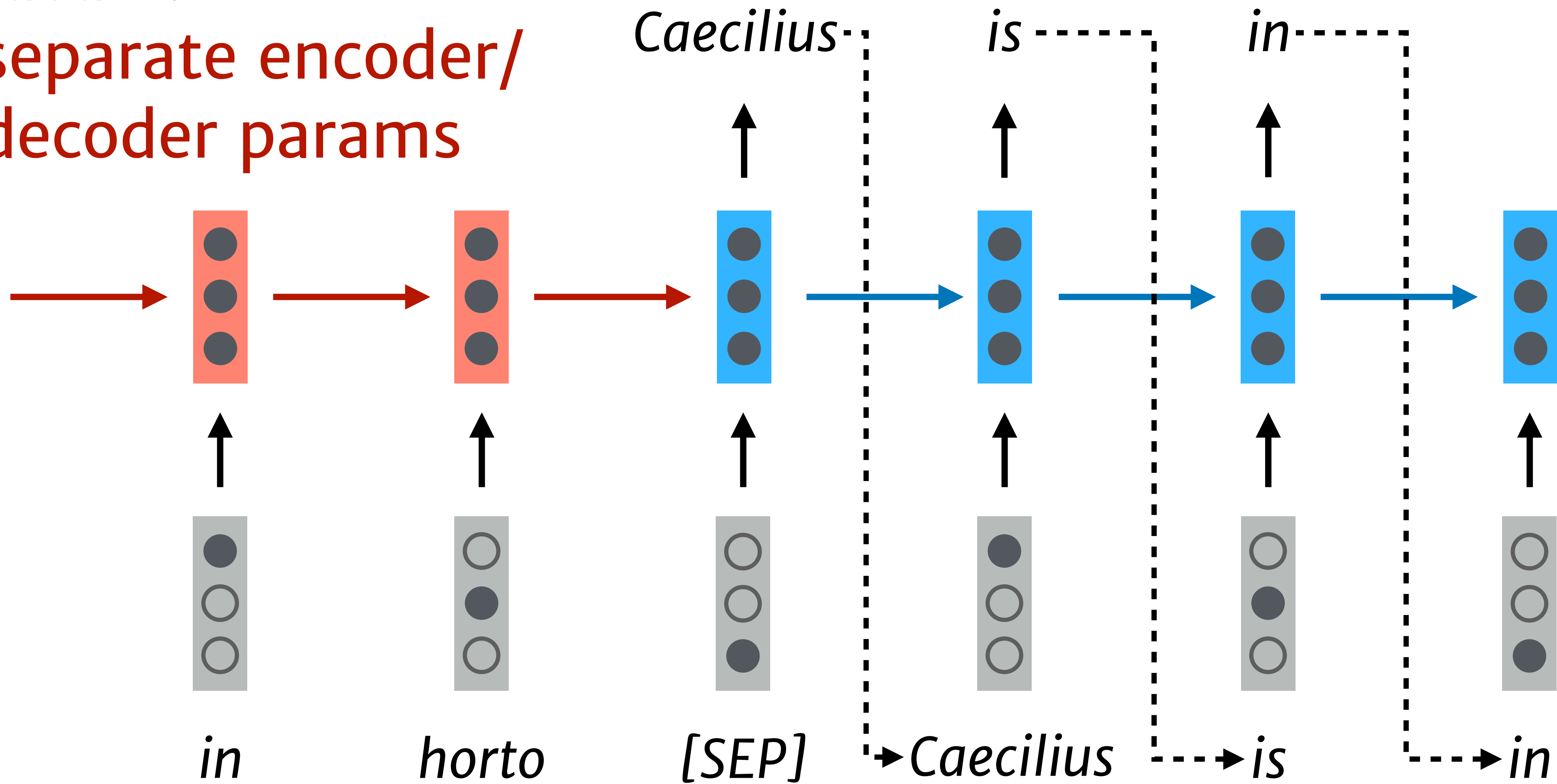
Sequence-to-sequence models

Idea 1:
only these losses:

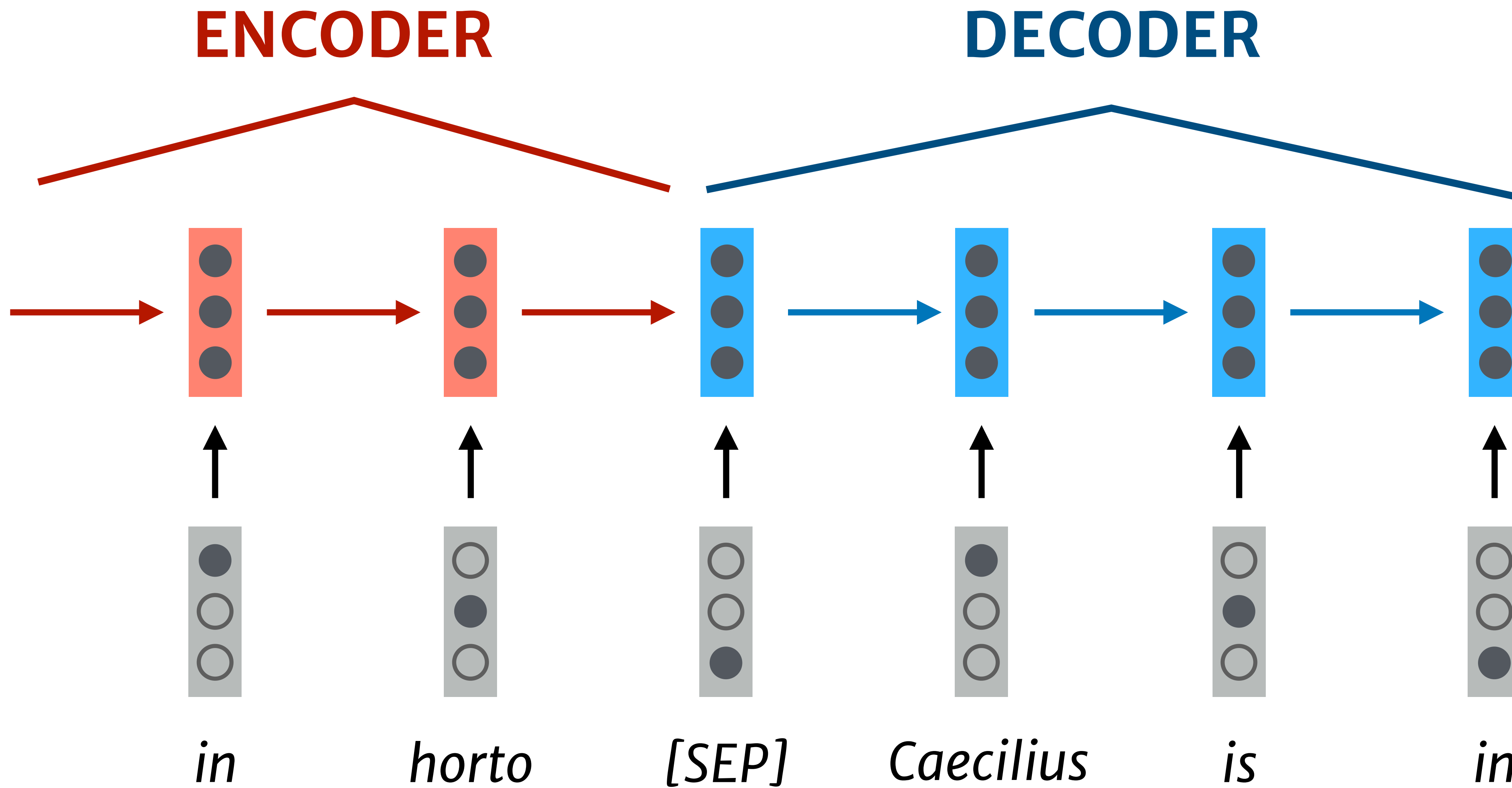


Sequence-to-sequence models

Idea 2:
separate encoder/
decoder params

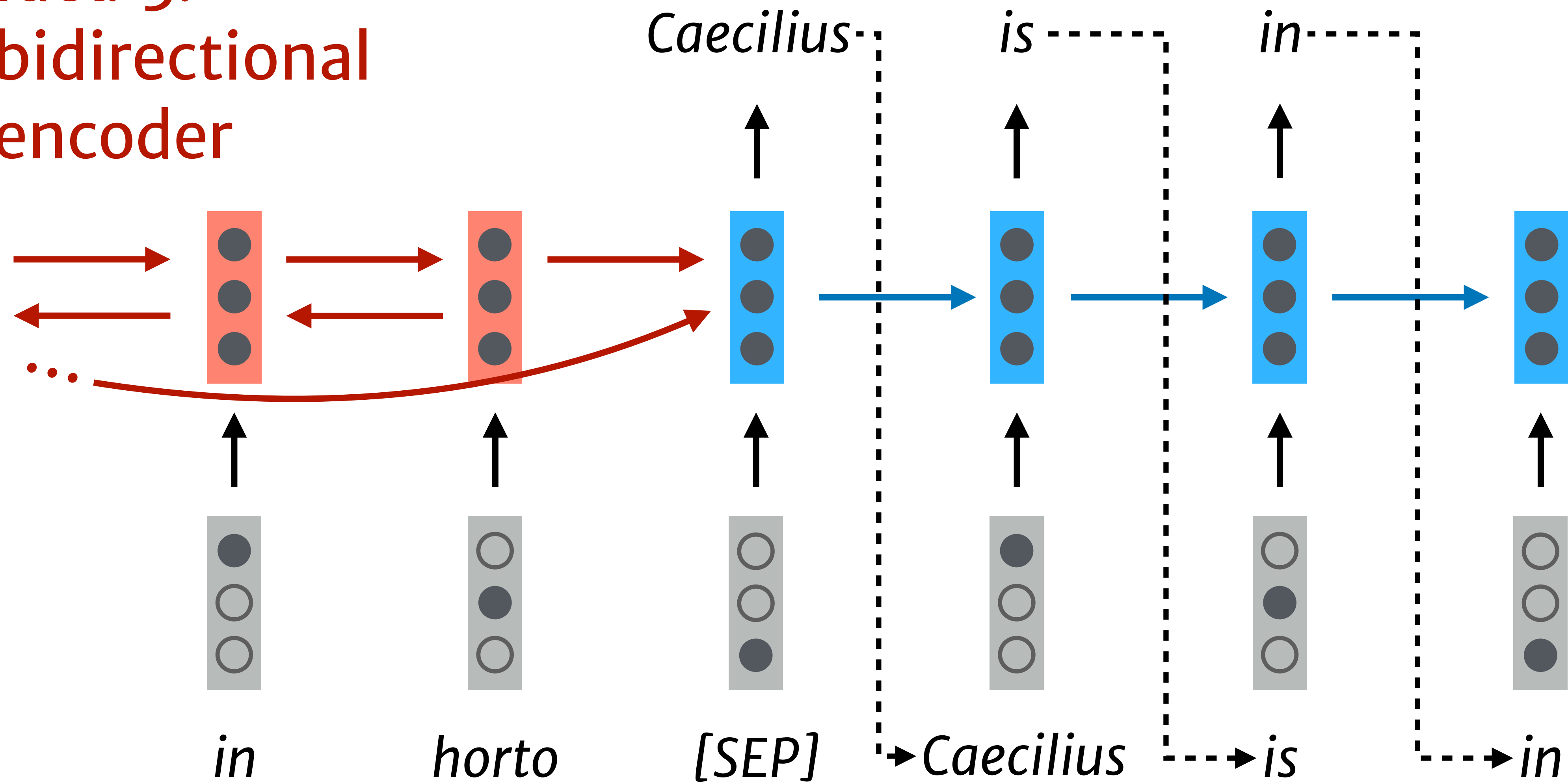


Sequence-to-sequence models



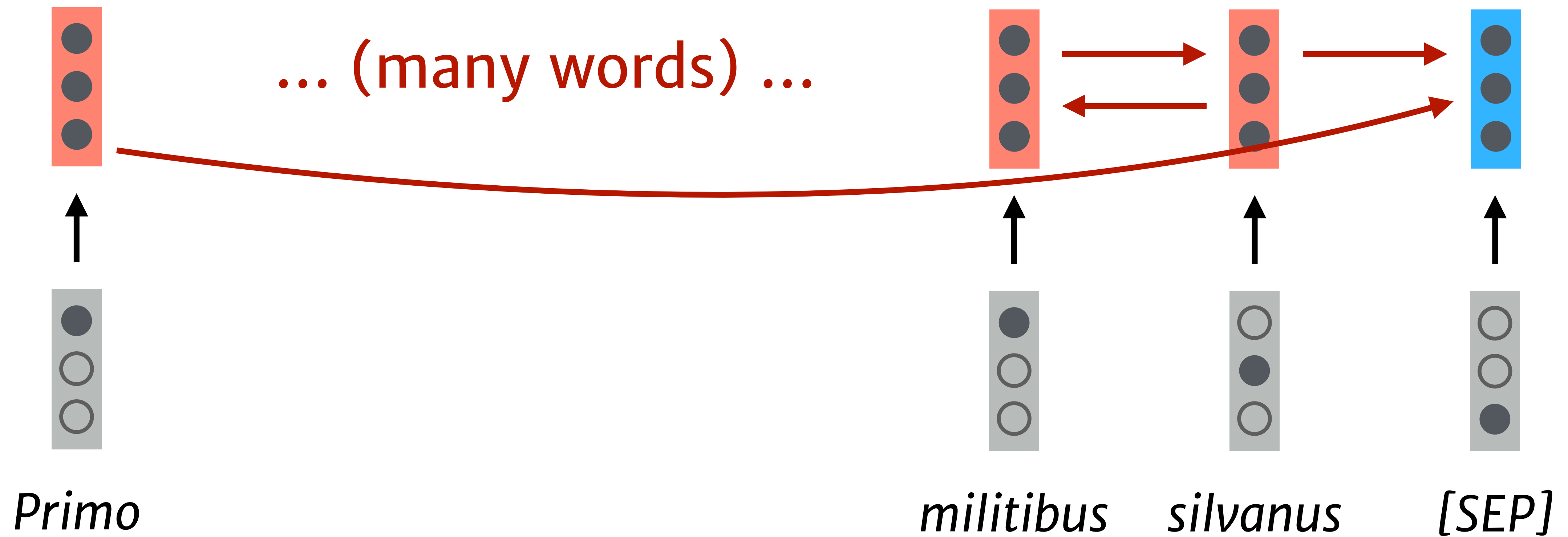
Sequence-to-sequence models

Idea 3:
bidirectional
encoder



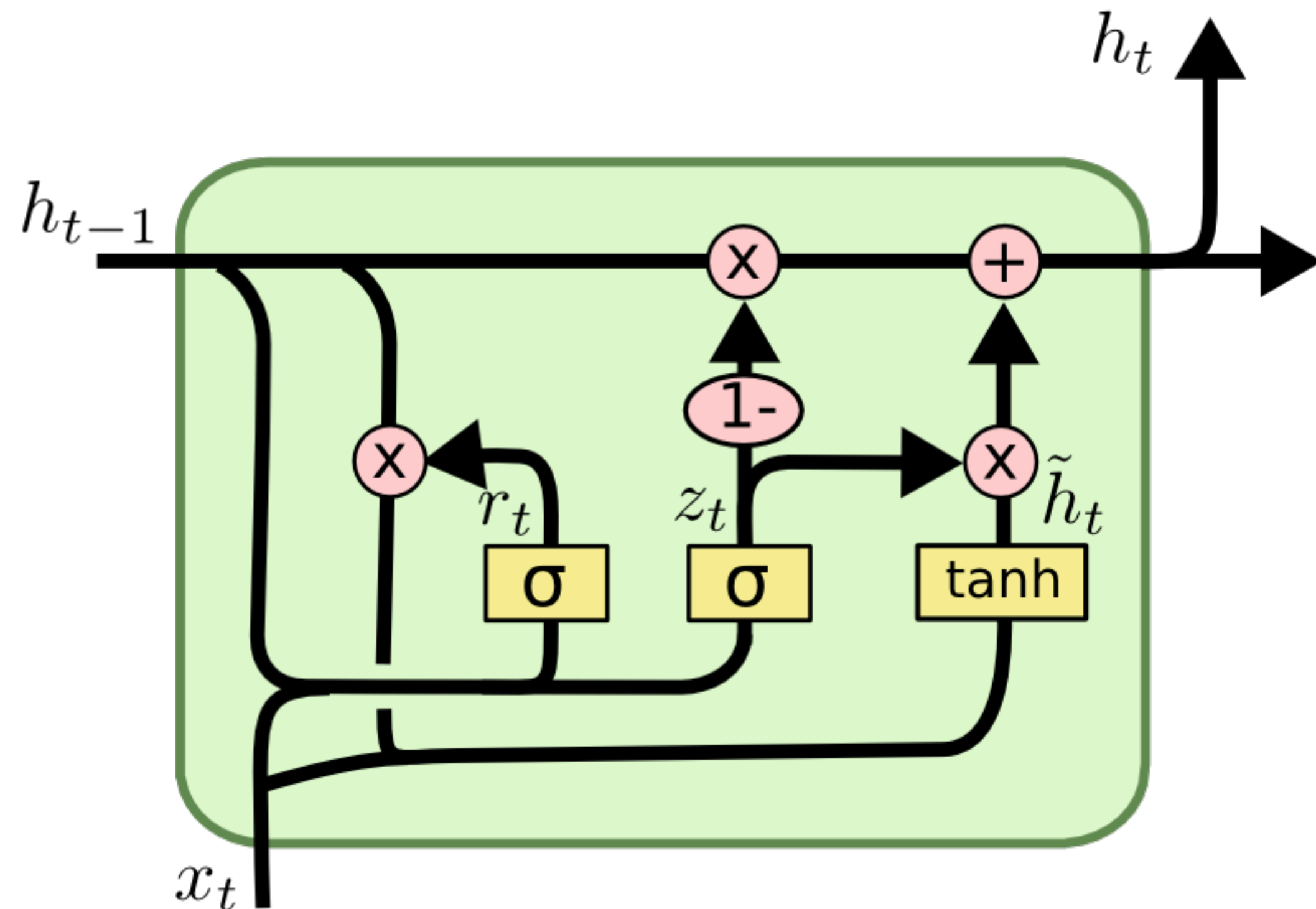
Revenge of the vanishing gradients

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



Attention mechanisms

Gated Recurrent Units



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

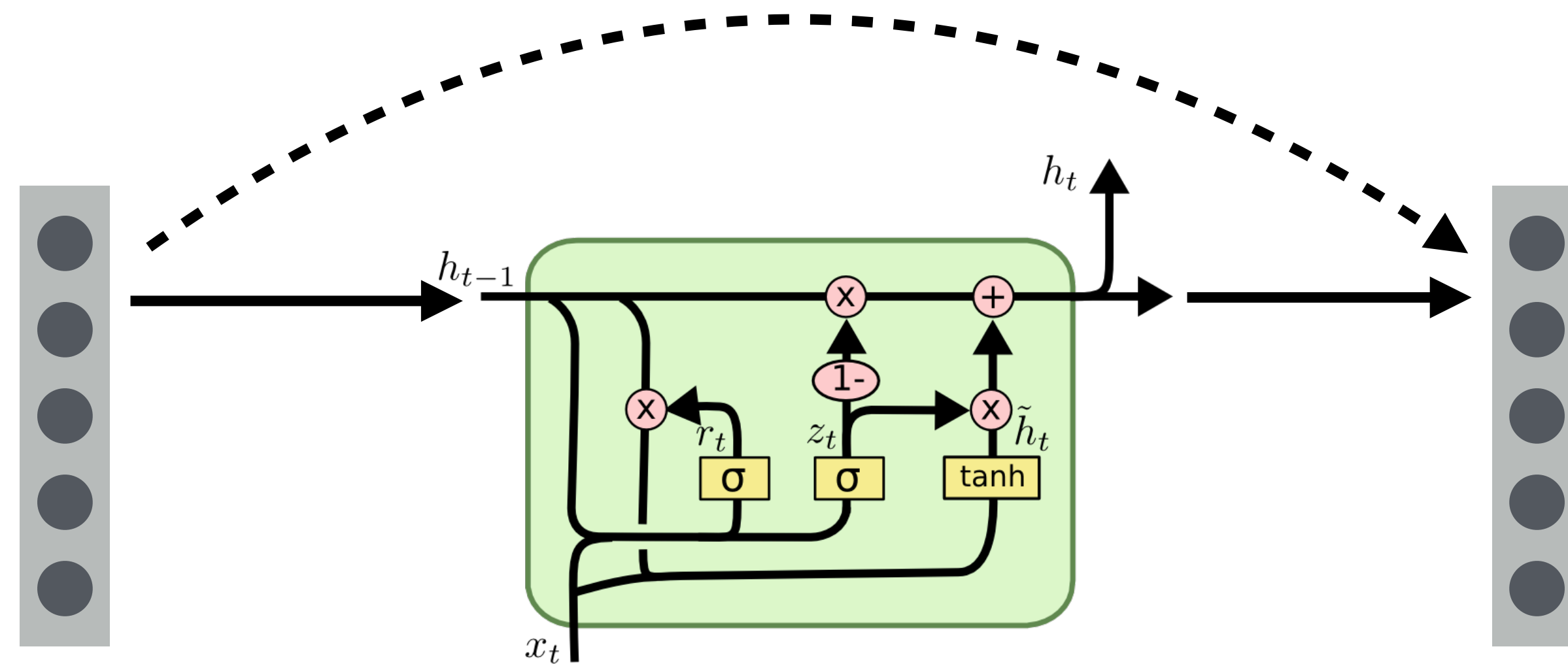
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

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

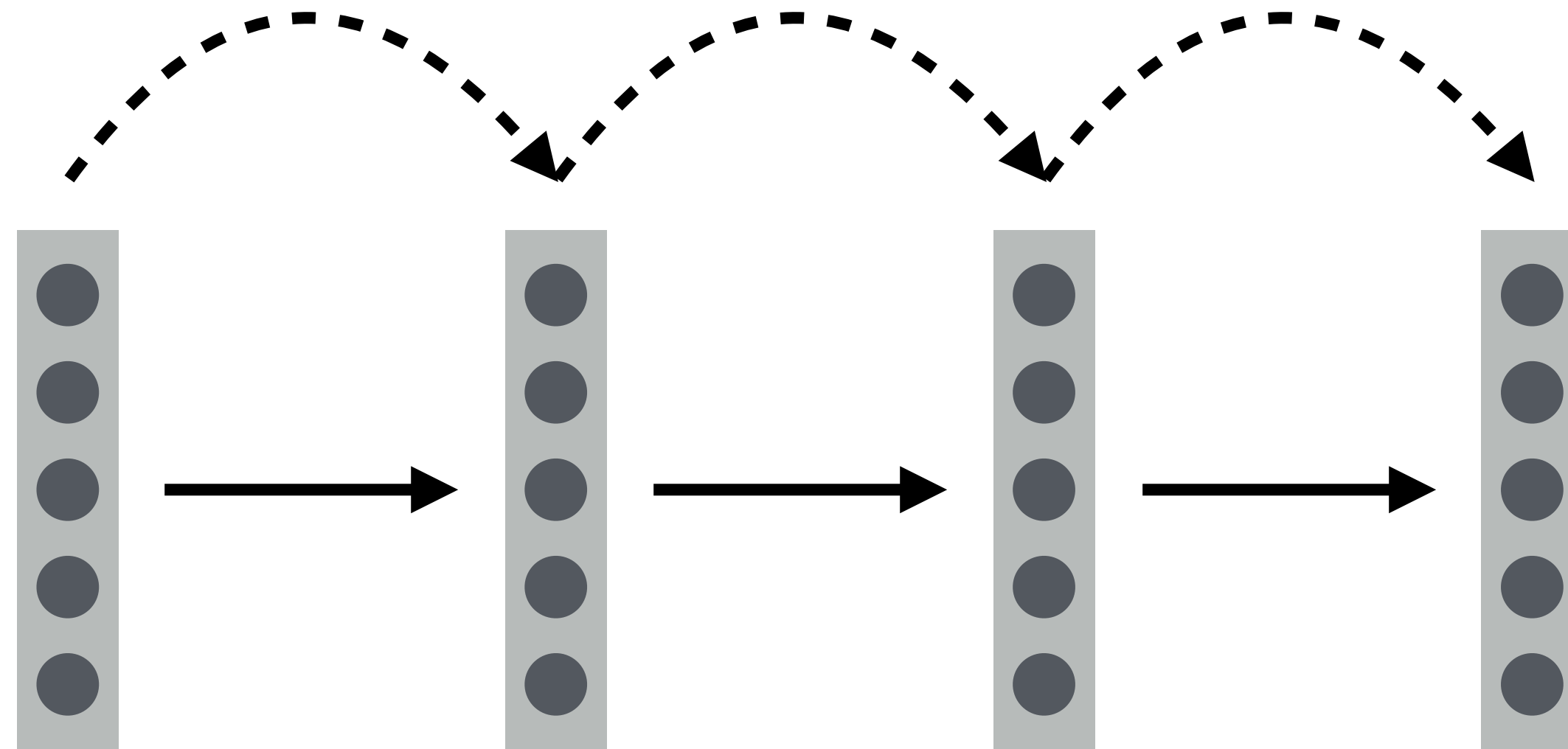
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

[Image: Cristopher Olah]

Shortcuts

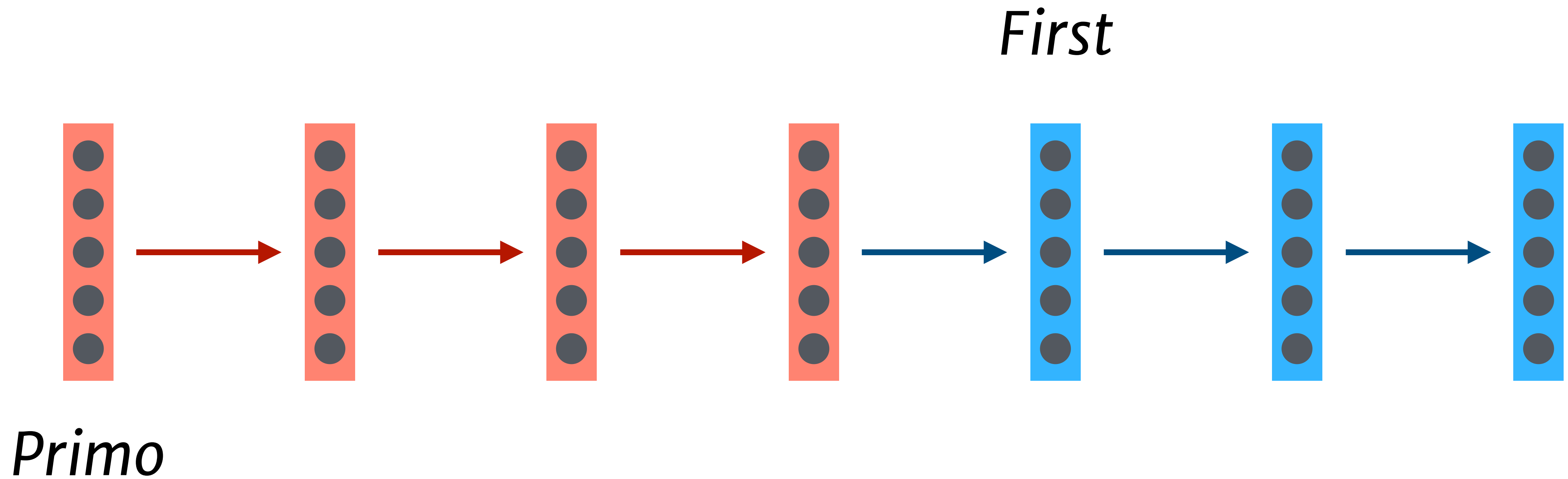


Shortcuts

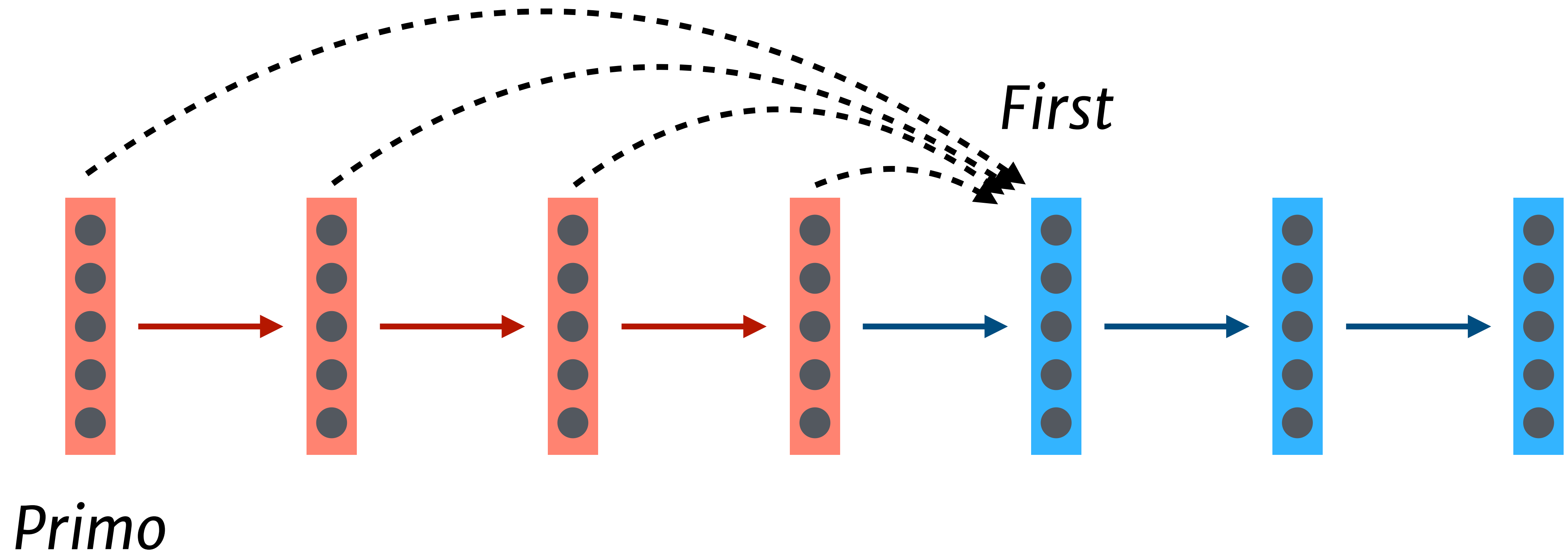


Direct "copying" between hidden states makes it easy to propagate information.

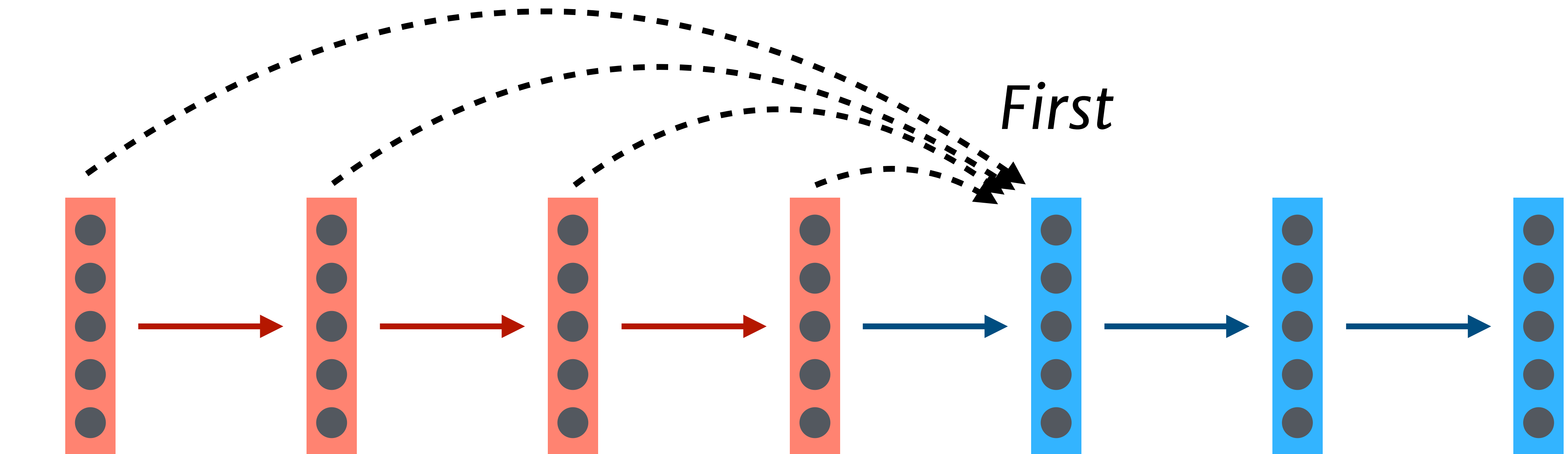
Can we go farther?



Can we go farther?



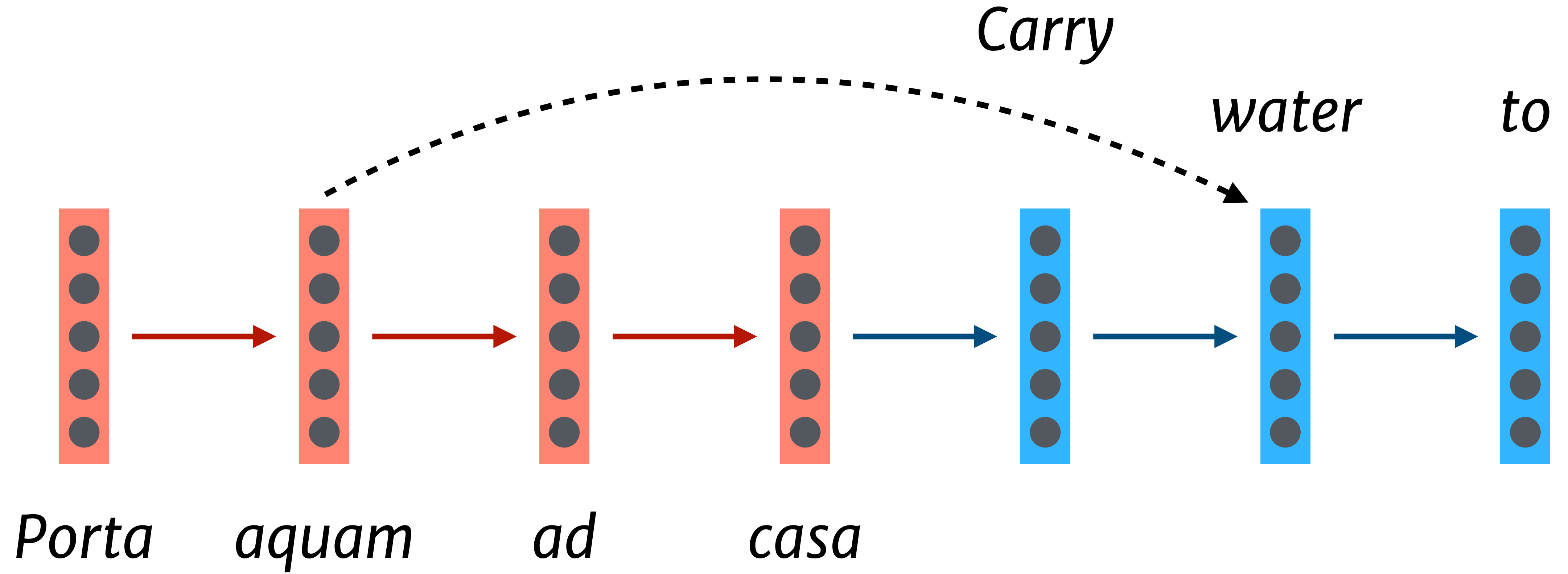
Can we go farther?



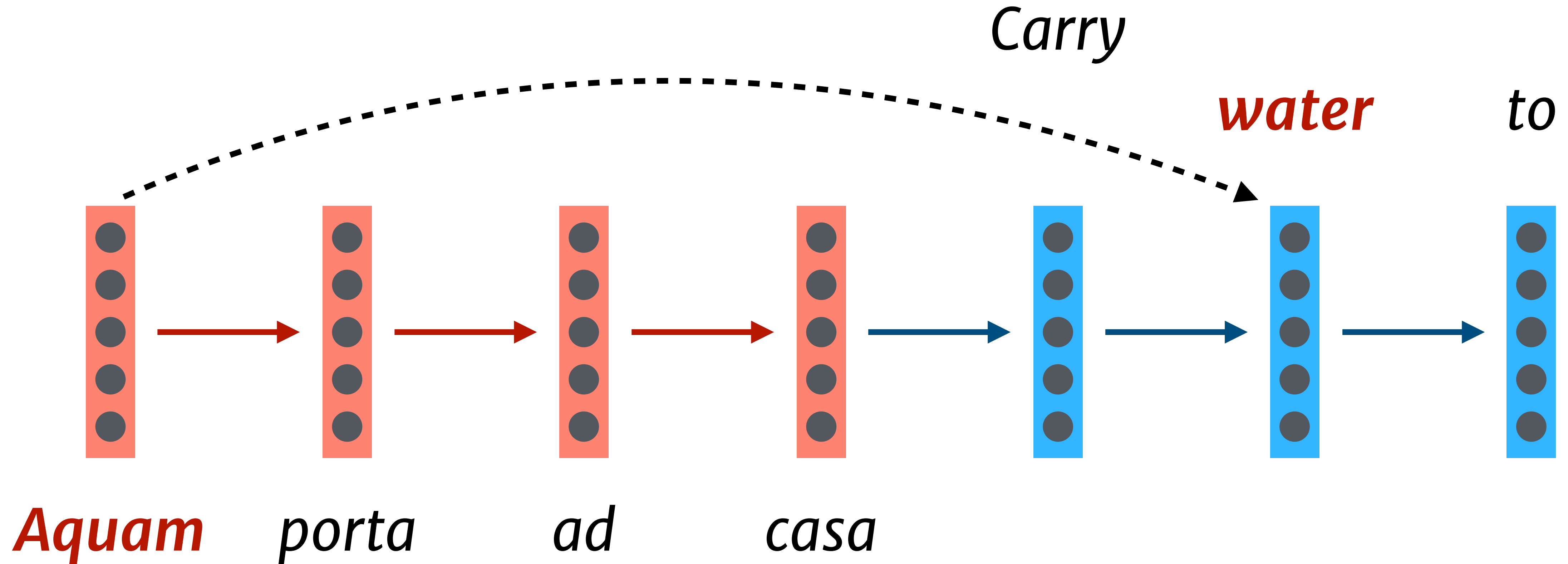
Primo

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?

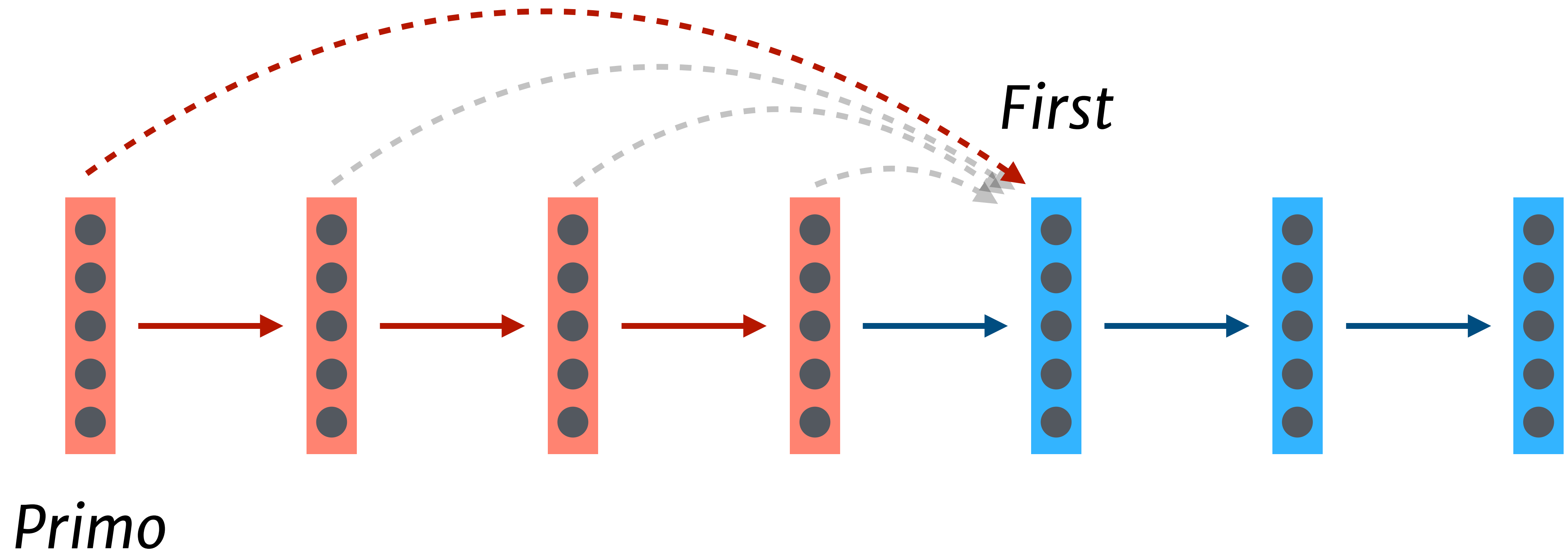


Can we hard-code connections?



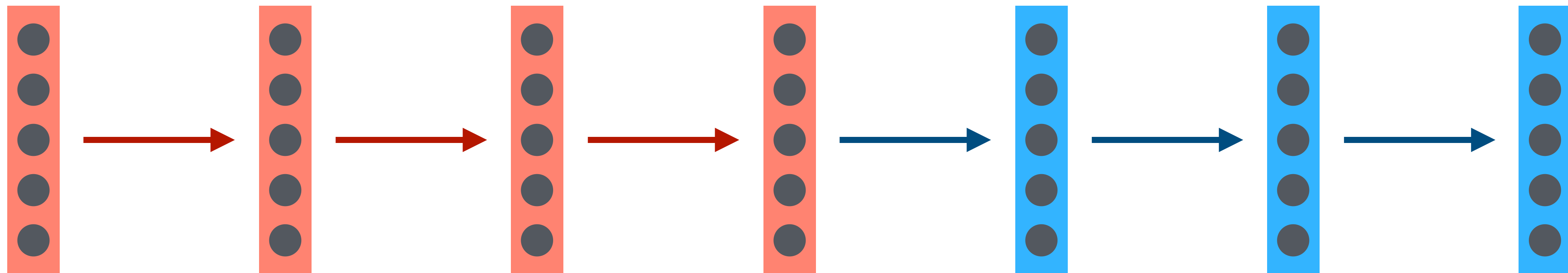
Words aren't one-to-one (and order can change!)

Can we **learn** connections?



Sentence representations

This vector represents the whole sentence!



Aquam *Aquam*
porta

Aquam *Aquam*
porta
ad

Aquam
porta
ad
casa



*You can't cram the
meaning of a whole
%&!\$# sentence into a
single \$&!#* vector!*

[Ray Mooney, ca. 2014]



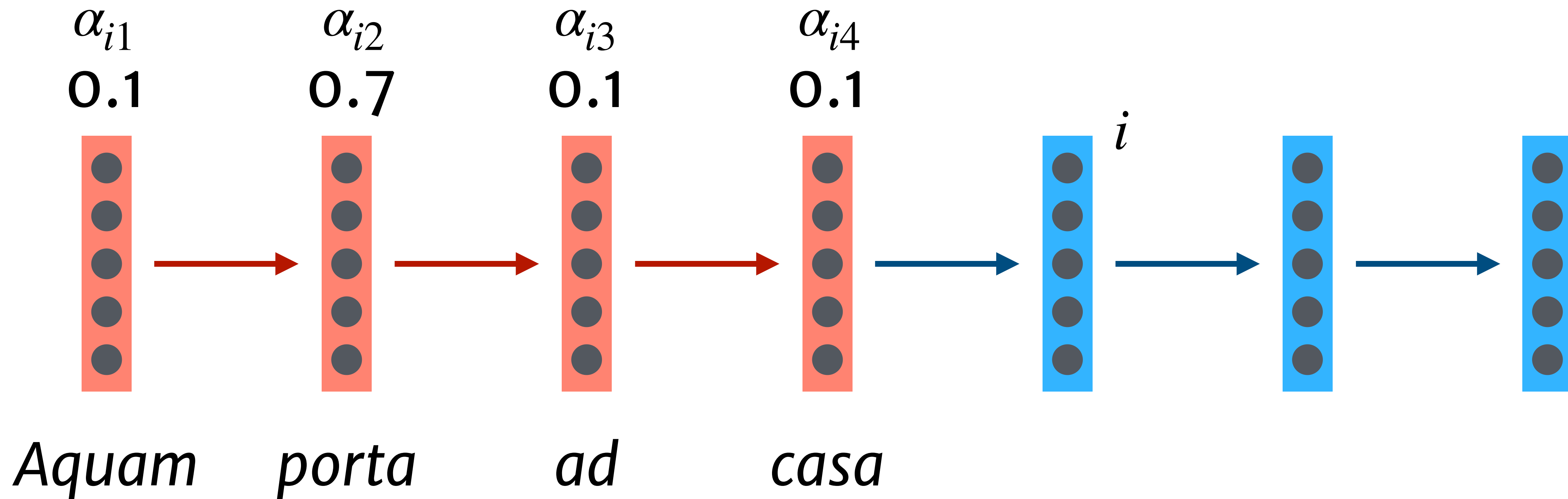
*You can't cram the
meaning of a whole
%&!\$# sentence into a
single \$&!#* vector!*

**Actually you can!
(But you usually shouldn't.)**

[Ray Mooney, ca. 2014]

Attention mechanisms

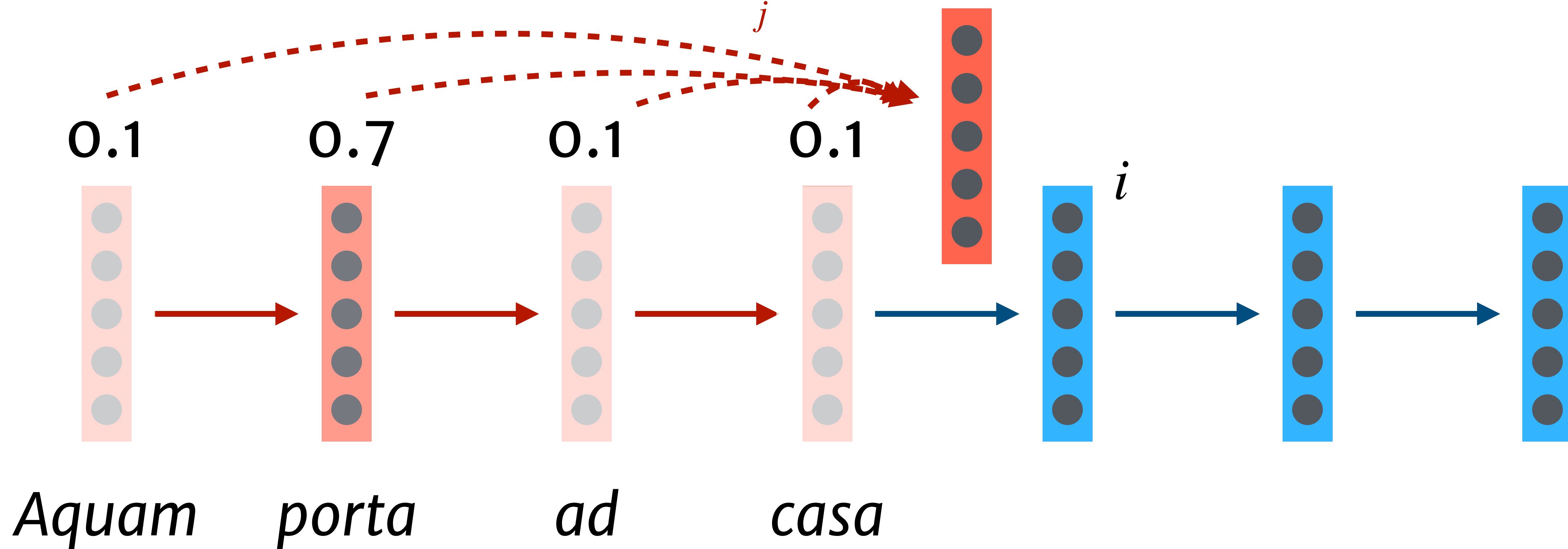
1. When predicting output i , assign a weight α_{ij} to each encoder state h_j (weights sum to 1)



Attention mechanisms

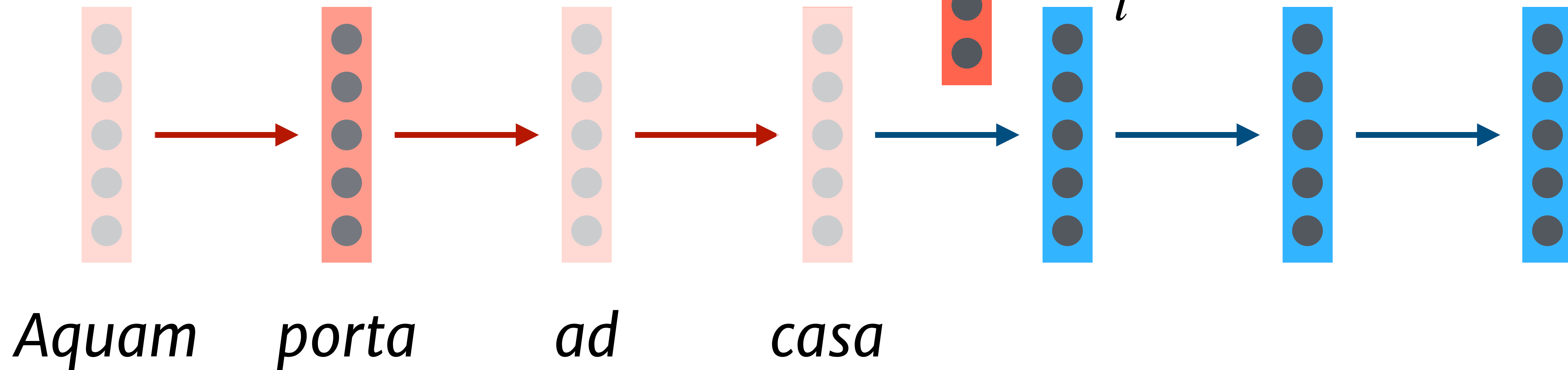
1. When predicting output i , assign a weight α_{ij} to each encoder state h_j

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



Attention mechanisms

1. When predicting output i , assign a weight α_{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



Design decision: how to compute α_{ij} ?

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

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

[Bahdanau 2014]

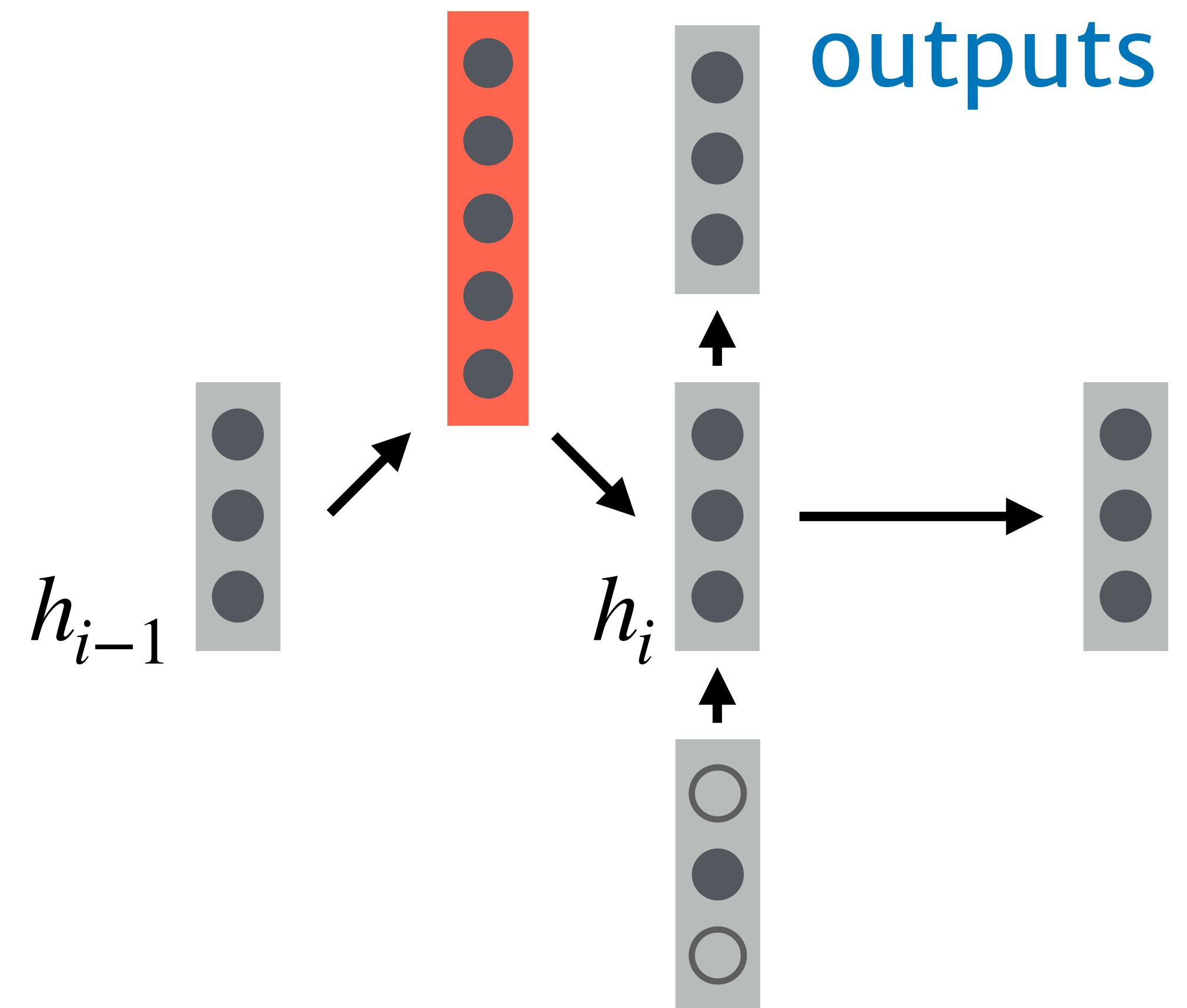
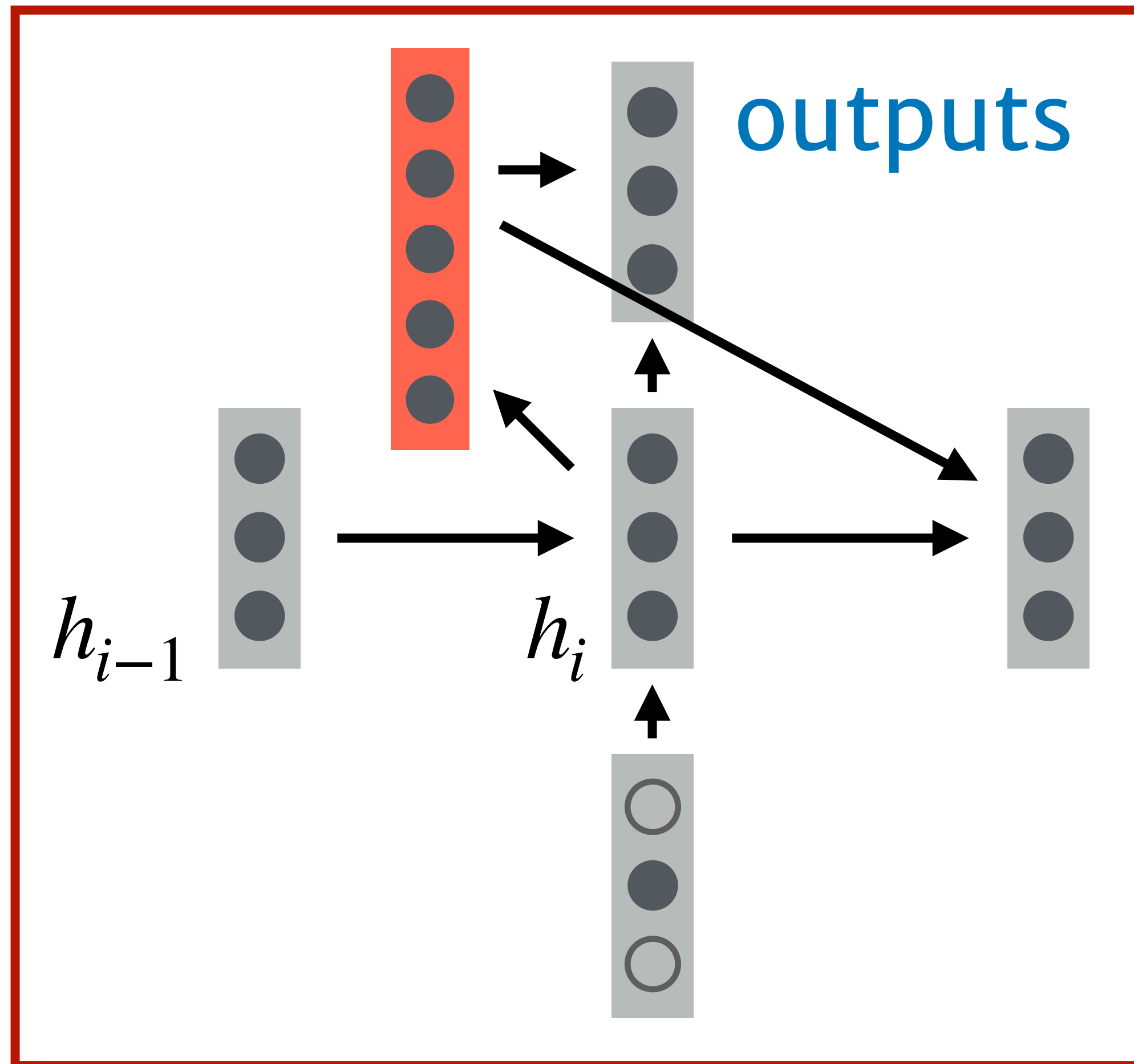
$$e_{ij} = h_i^\top W h_j$$

[Luong 2015]

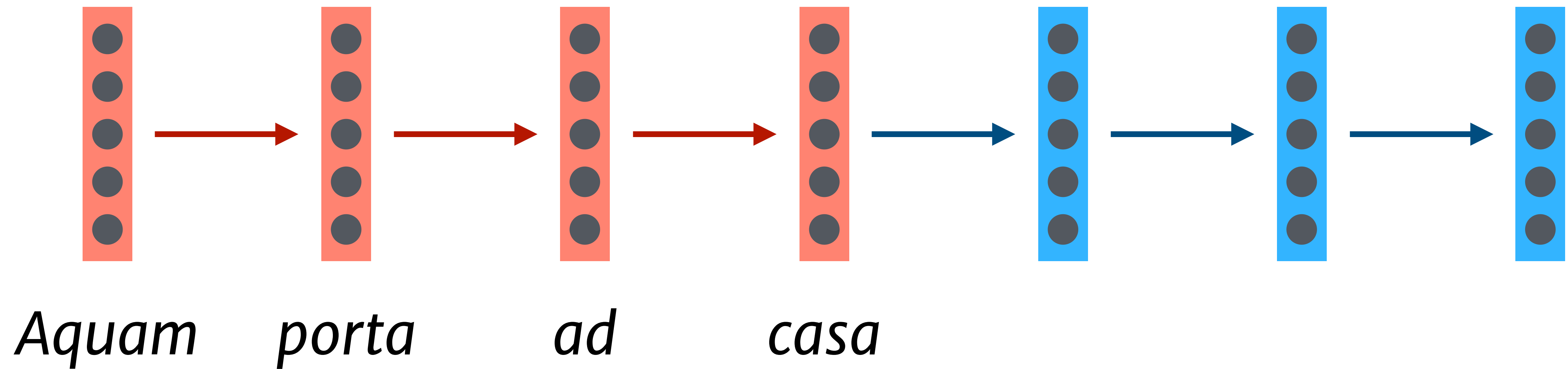
$$\alpha_{i\cdot} = \text{softmax}(e_{i\cdot})$$

Design decision: how to use c_i ?

3. Use c_i to update the decoder

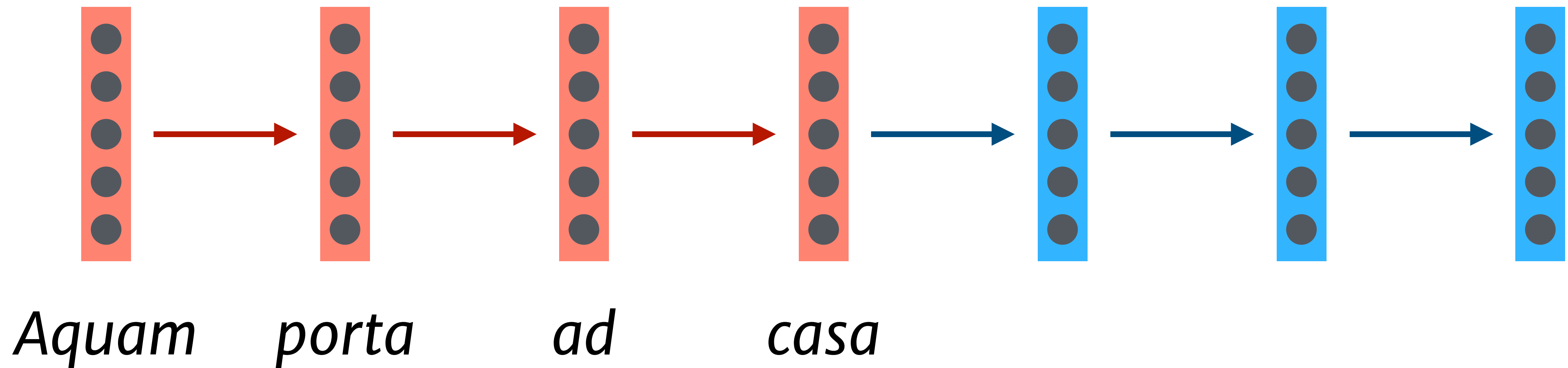


Why does this work?



Why does this work?

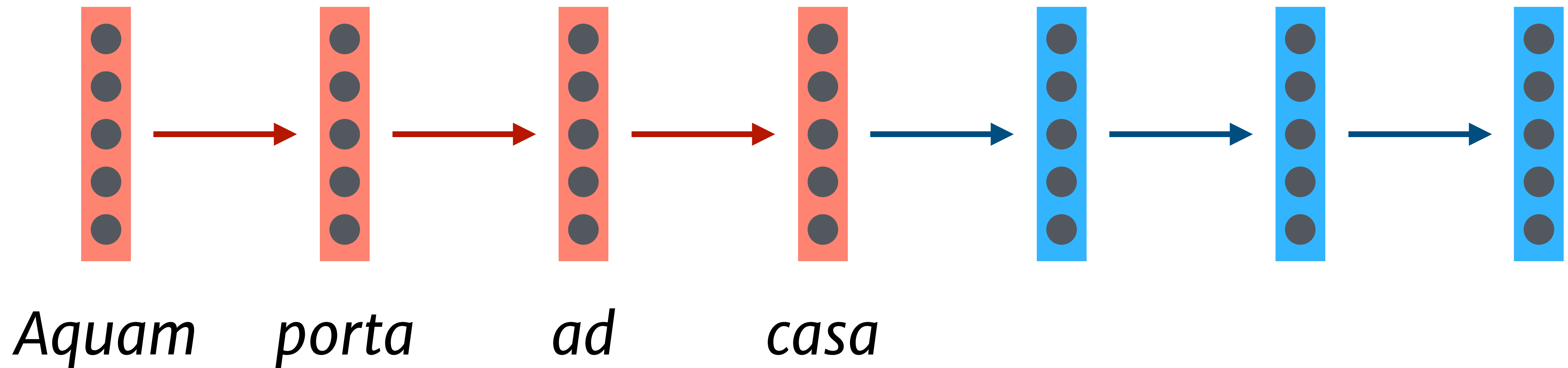
MAIN VERB
INDEX 2
IMPERATIVE

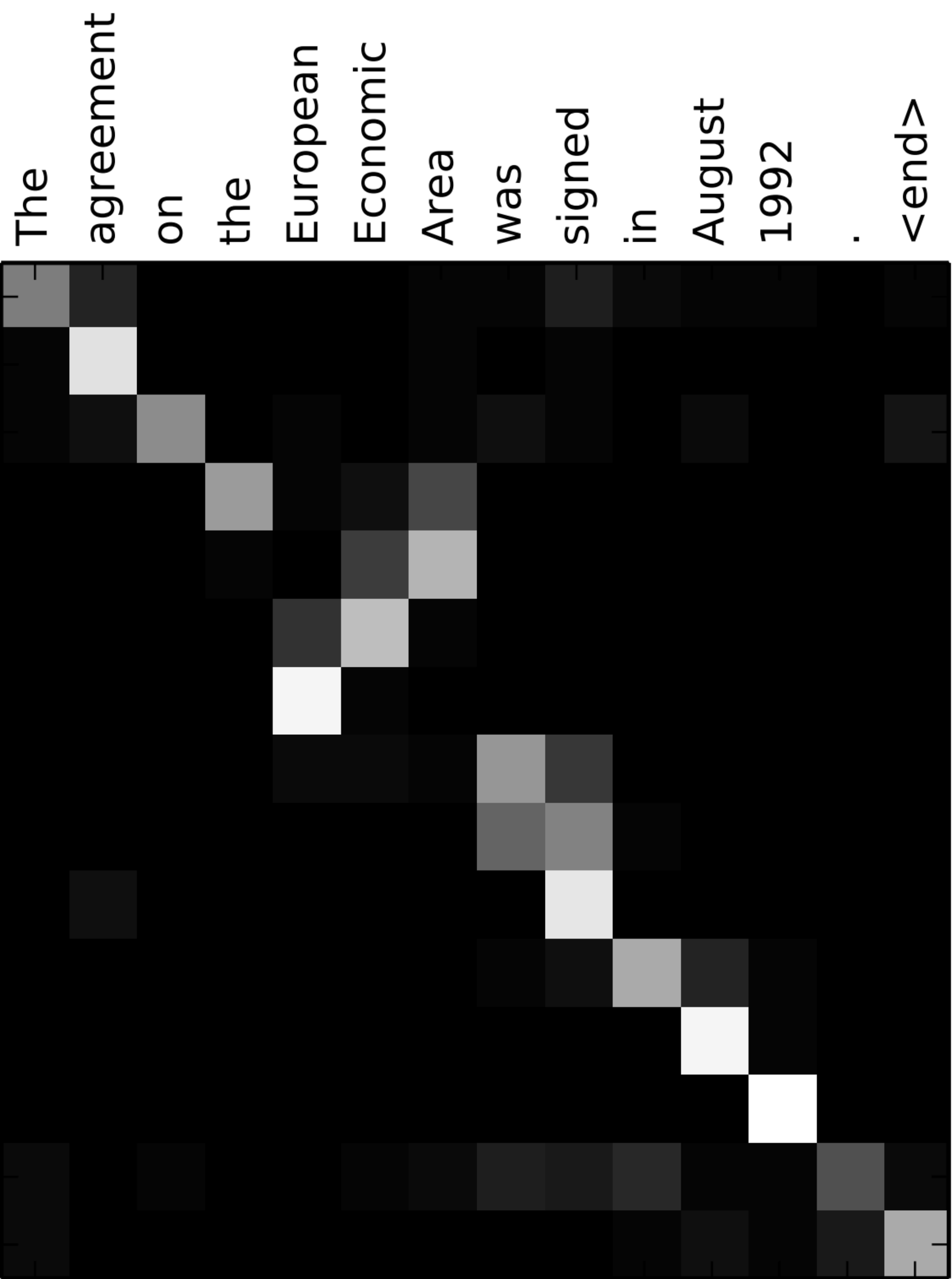


Why does this work?

MAIN VERB
INDEX 2
IMPERATIVE

SUBJECT?
IMP. VERB?





[Example from Greg Durrett]

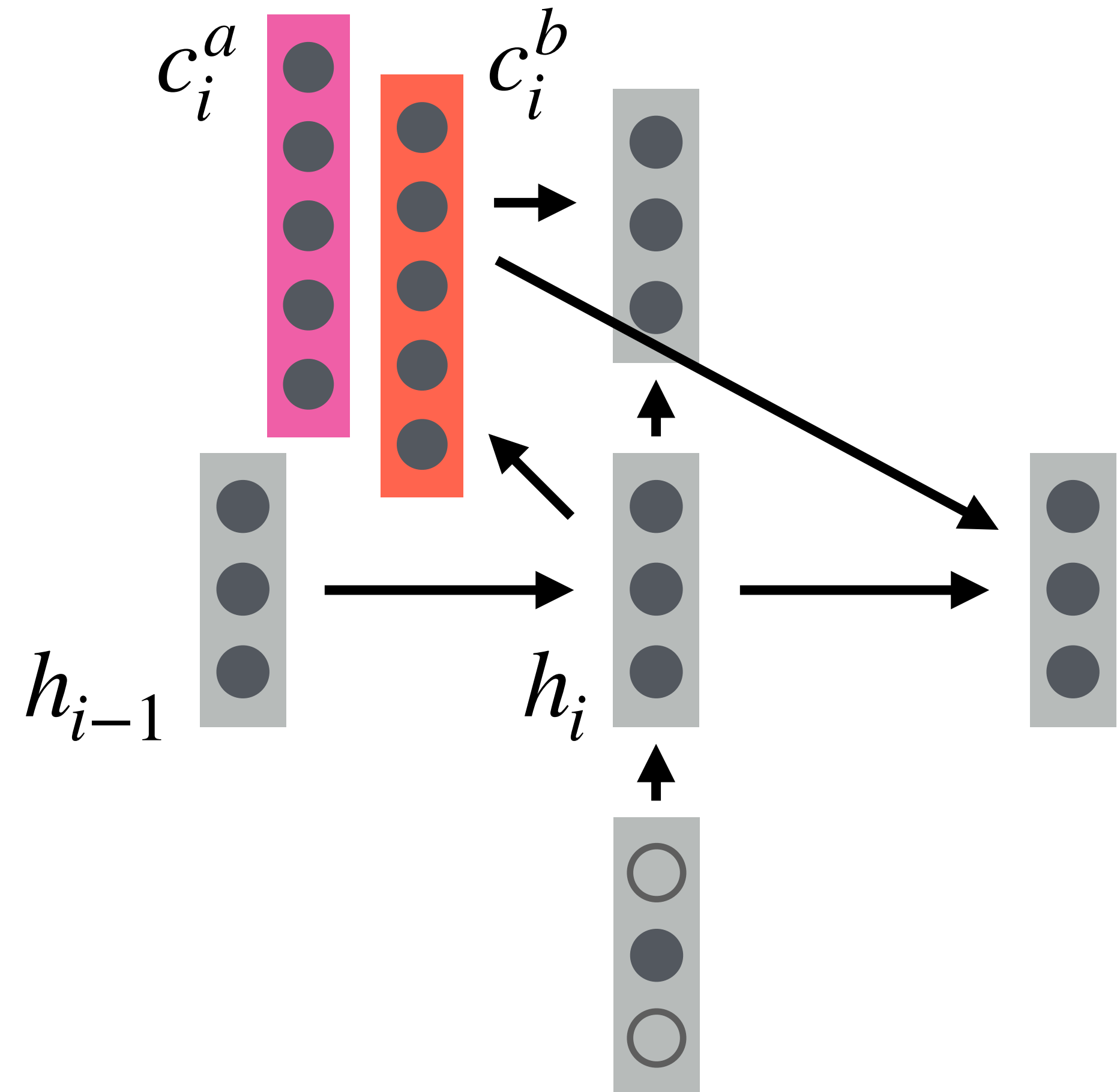
Multi-headed attention

Look two places at once!

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

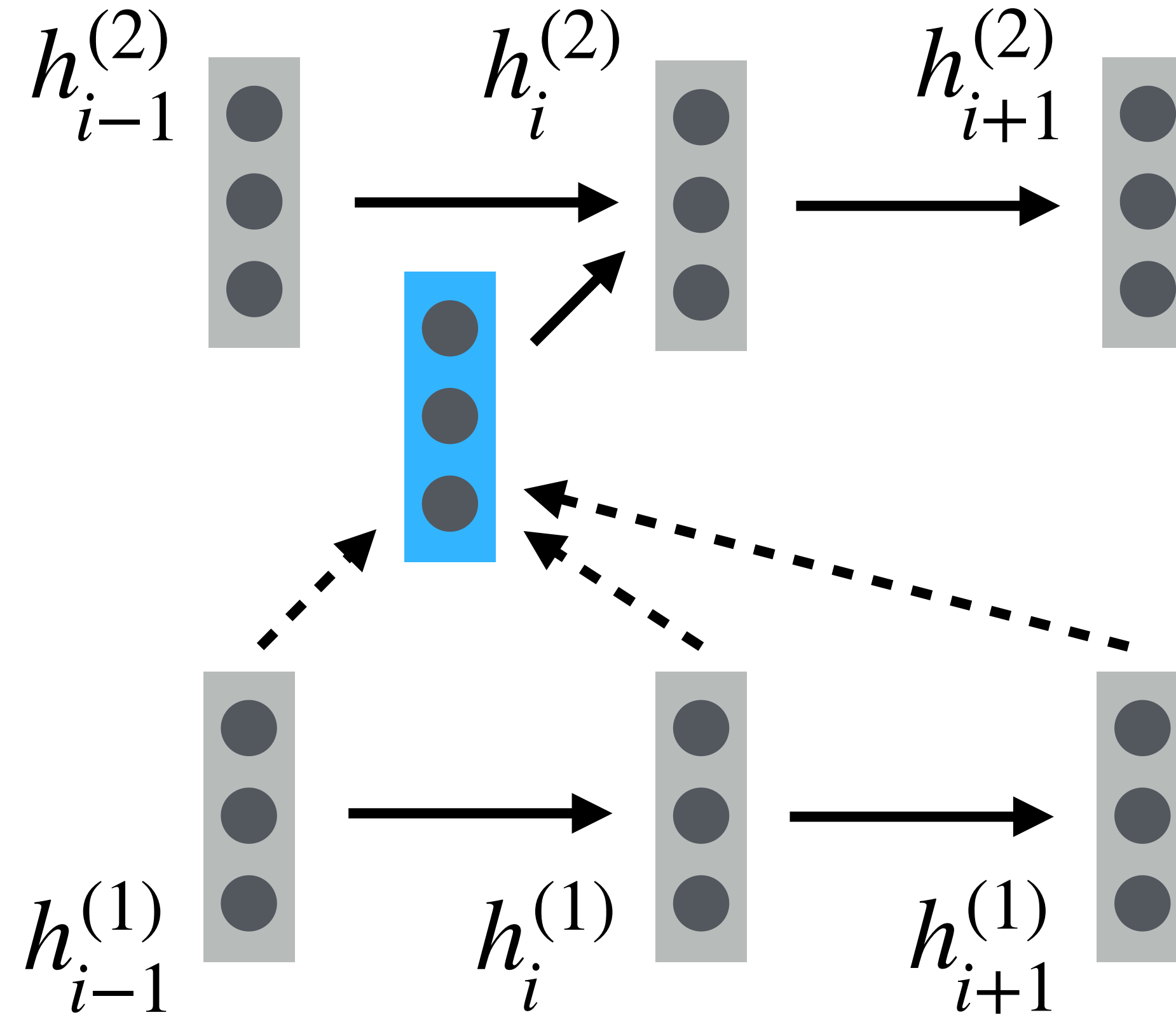
$$e_{ij}^b = h_i^\top W_b h_j$$

etc.

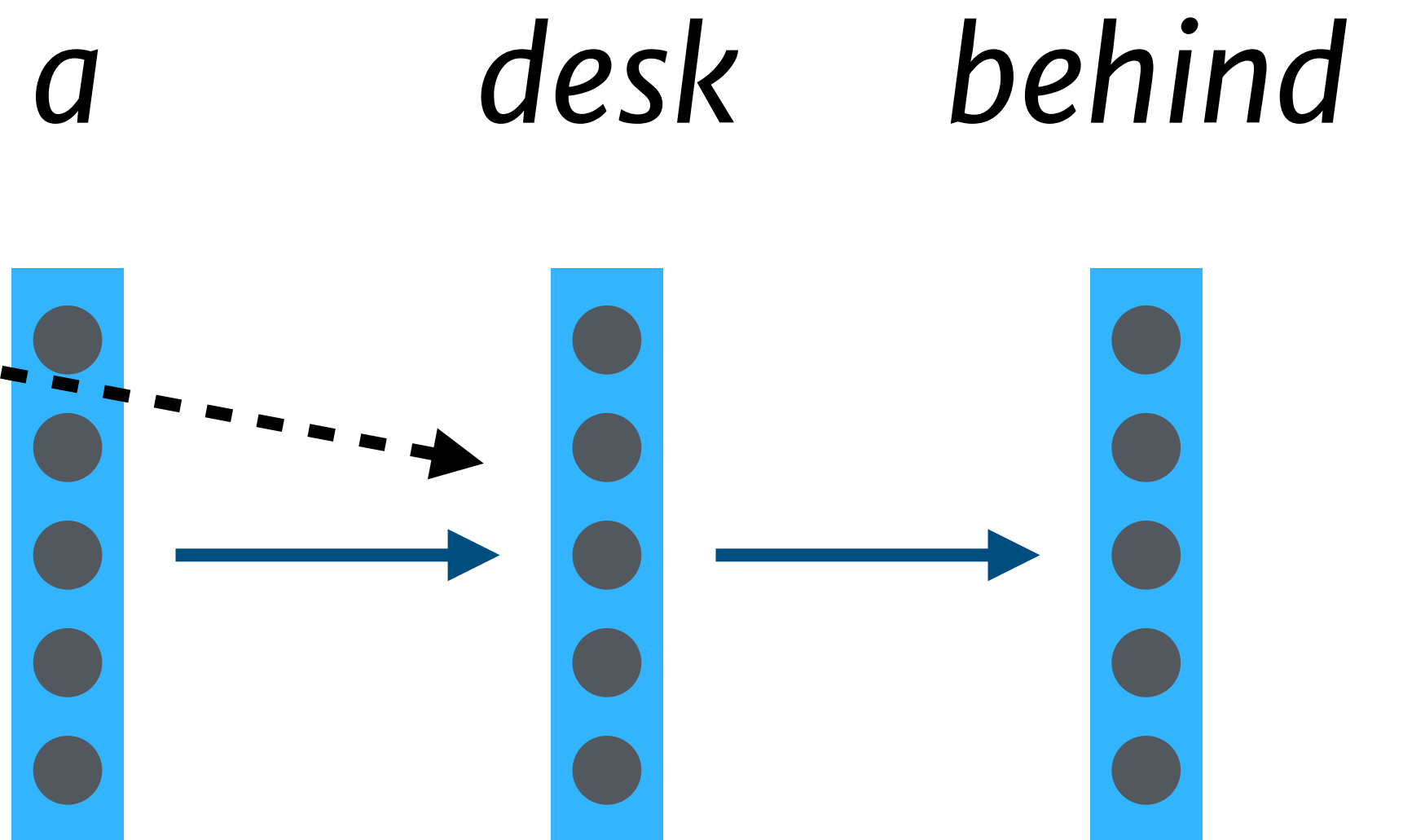
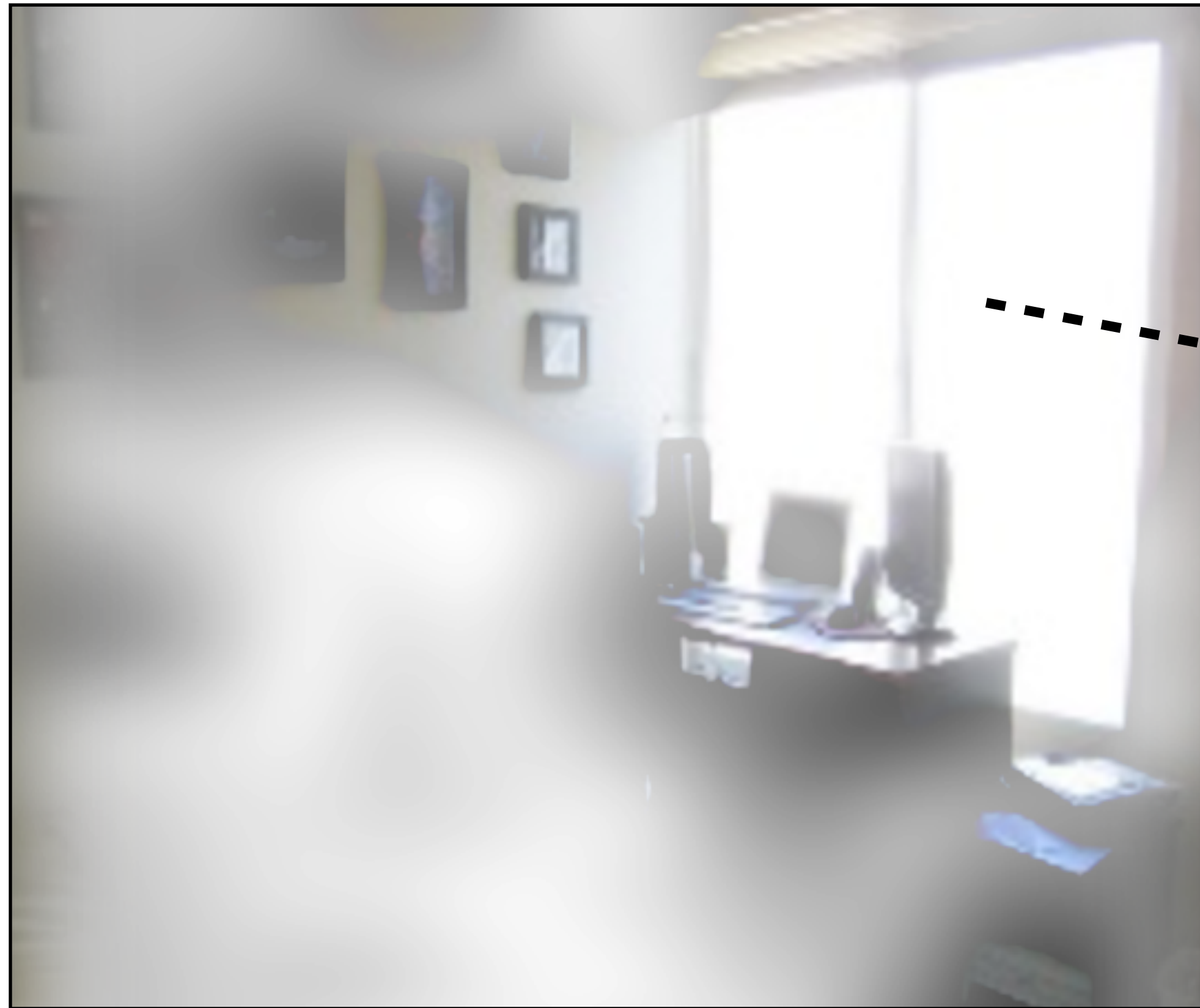


Self-attention

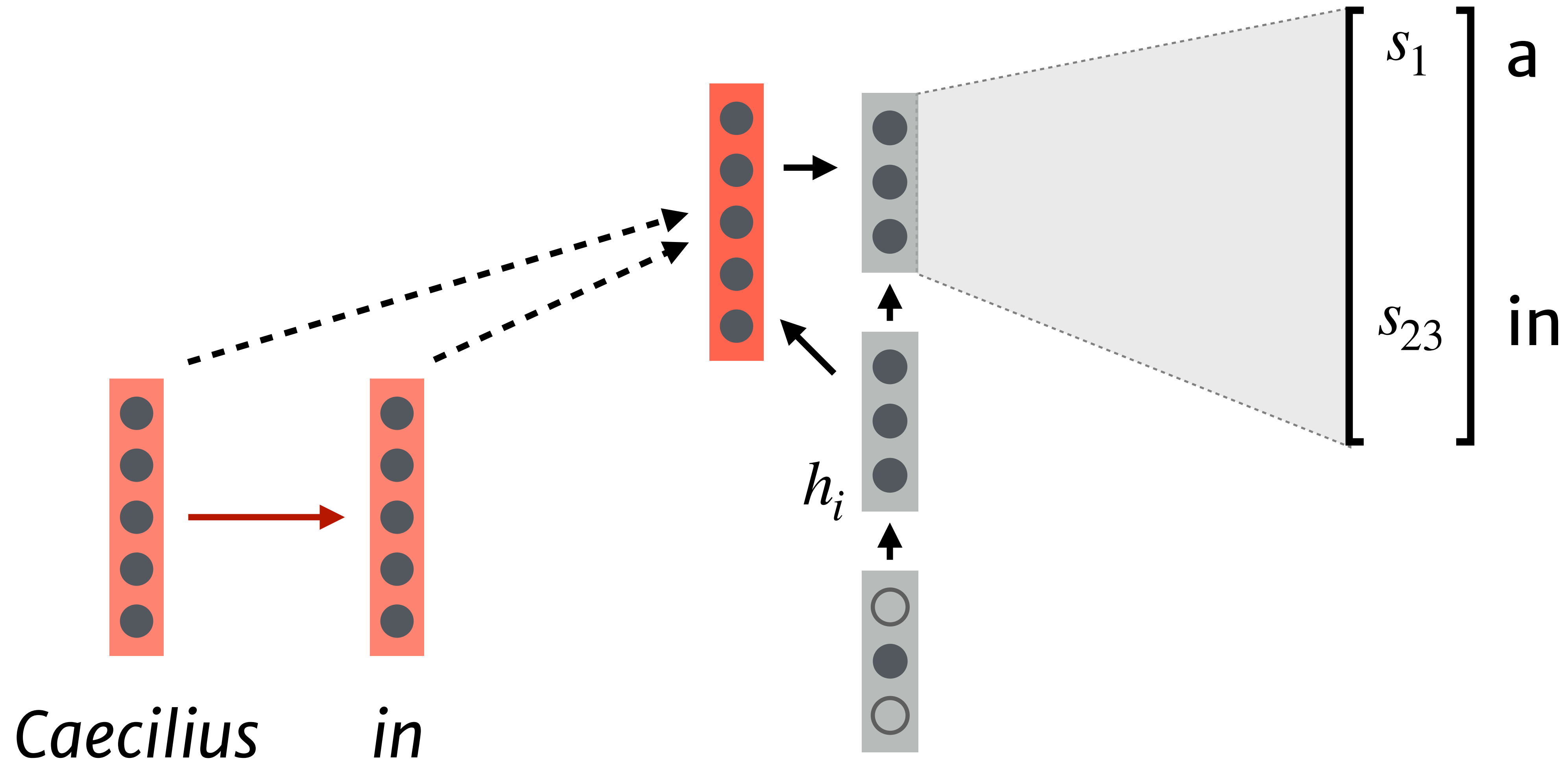
Attention to lower RNN layers (instead of decoder → encoder)



Non-textual attention

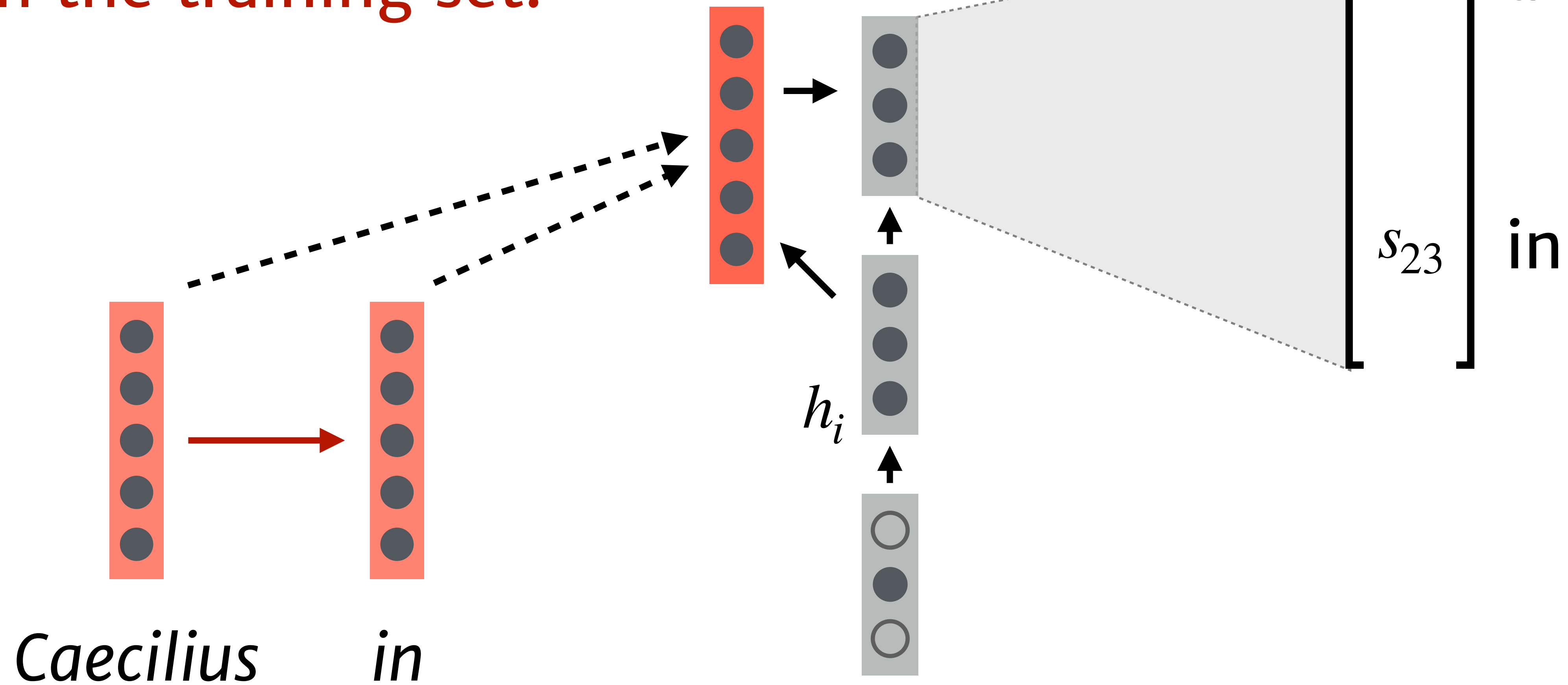


Copying



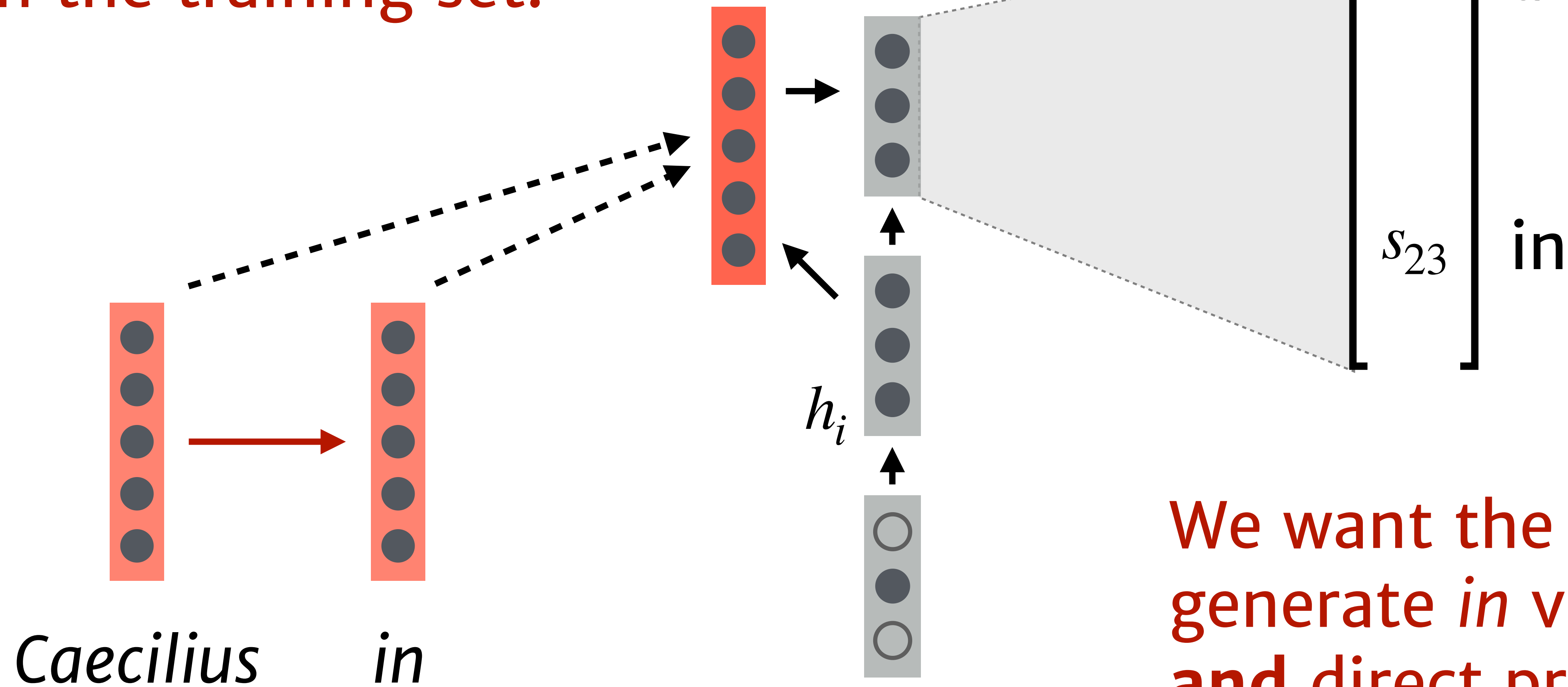
Copying

Caecilius probably isn't in the training set.



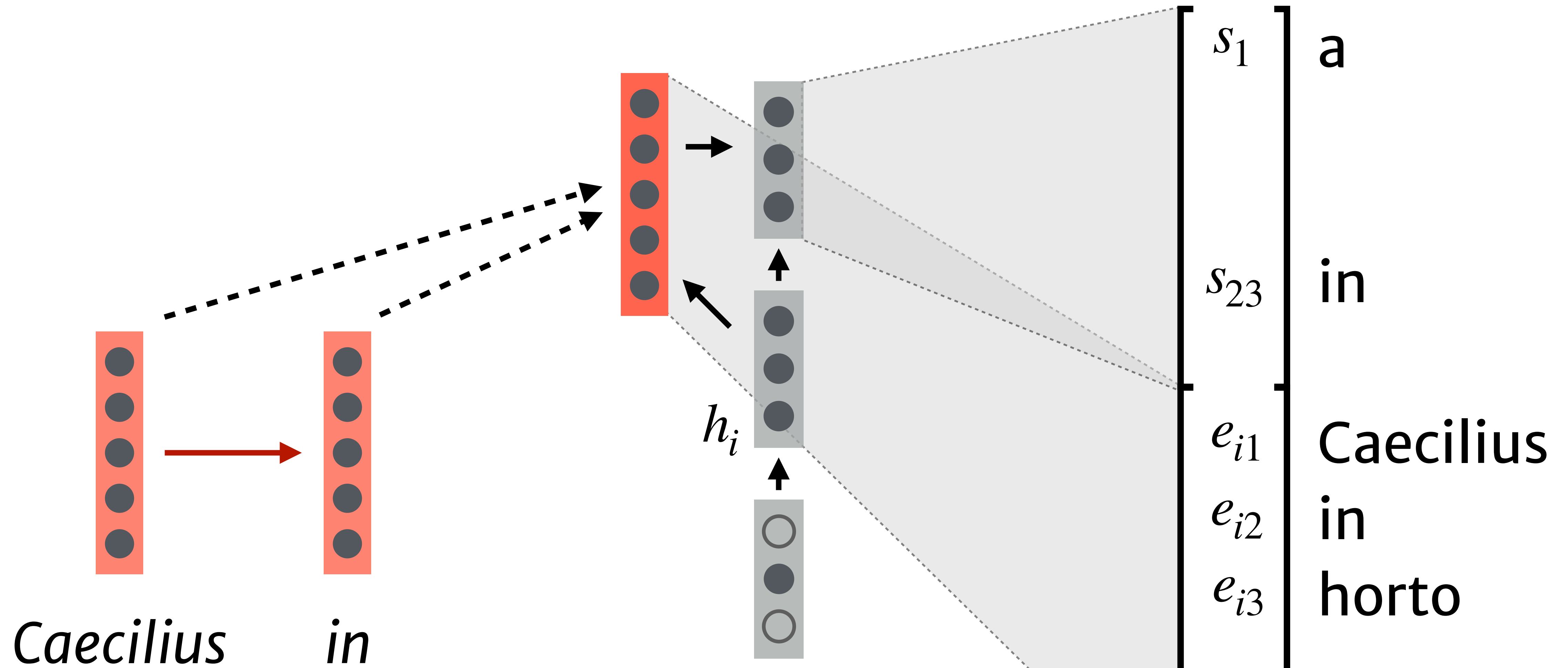
Copying

Caecilius probably isn't in the training set.



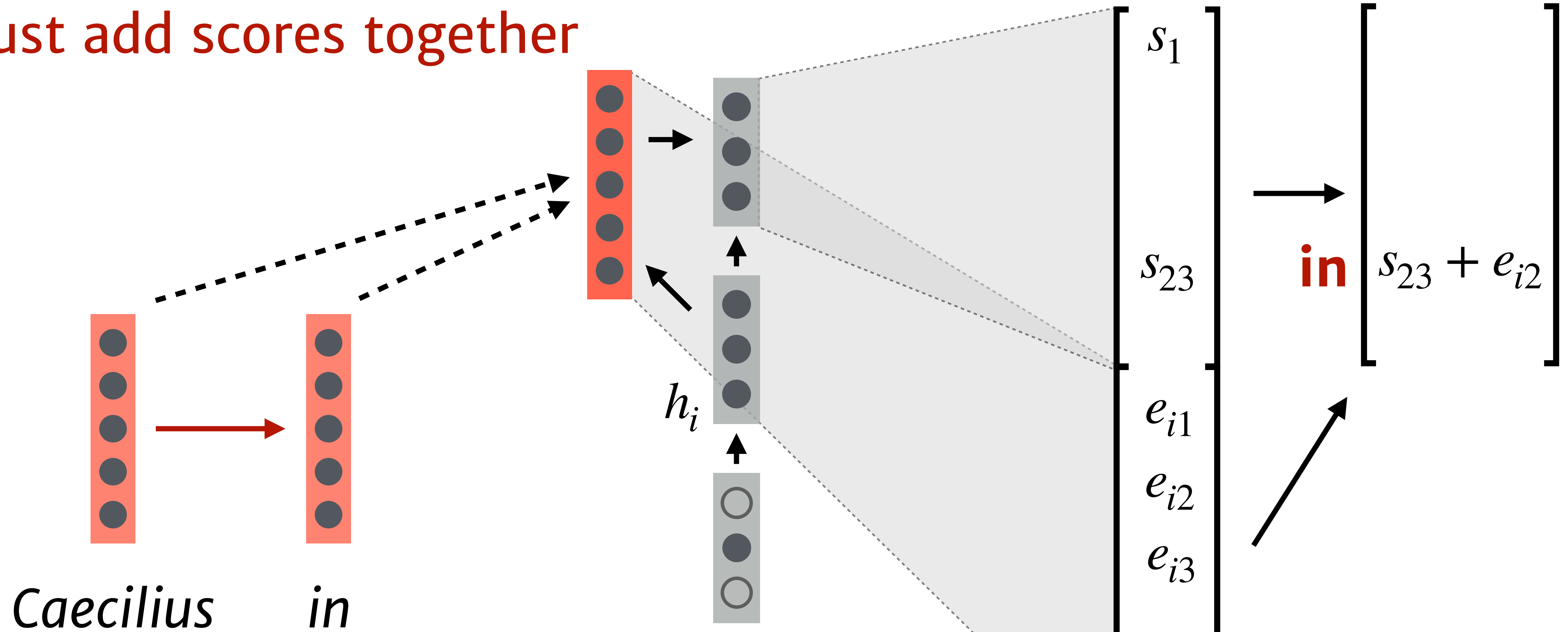
We want the ability to generate *in* via copying and direct prediction.

Copying



Copying

In is double-counted:
just add scores together



Hard attention

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

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

[Bahdanau 2014]

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

[Luong 2015]

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

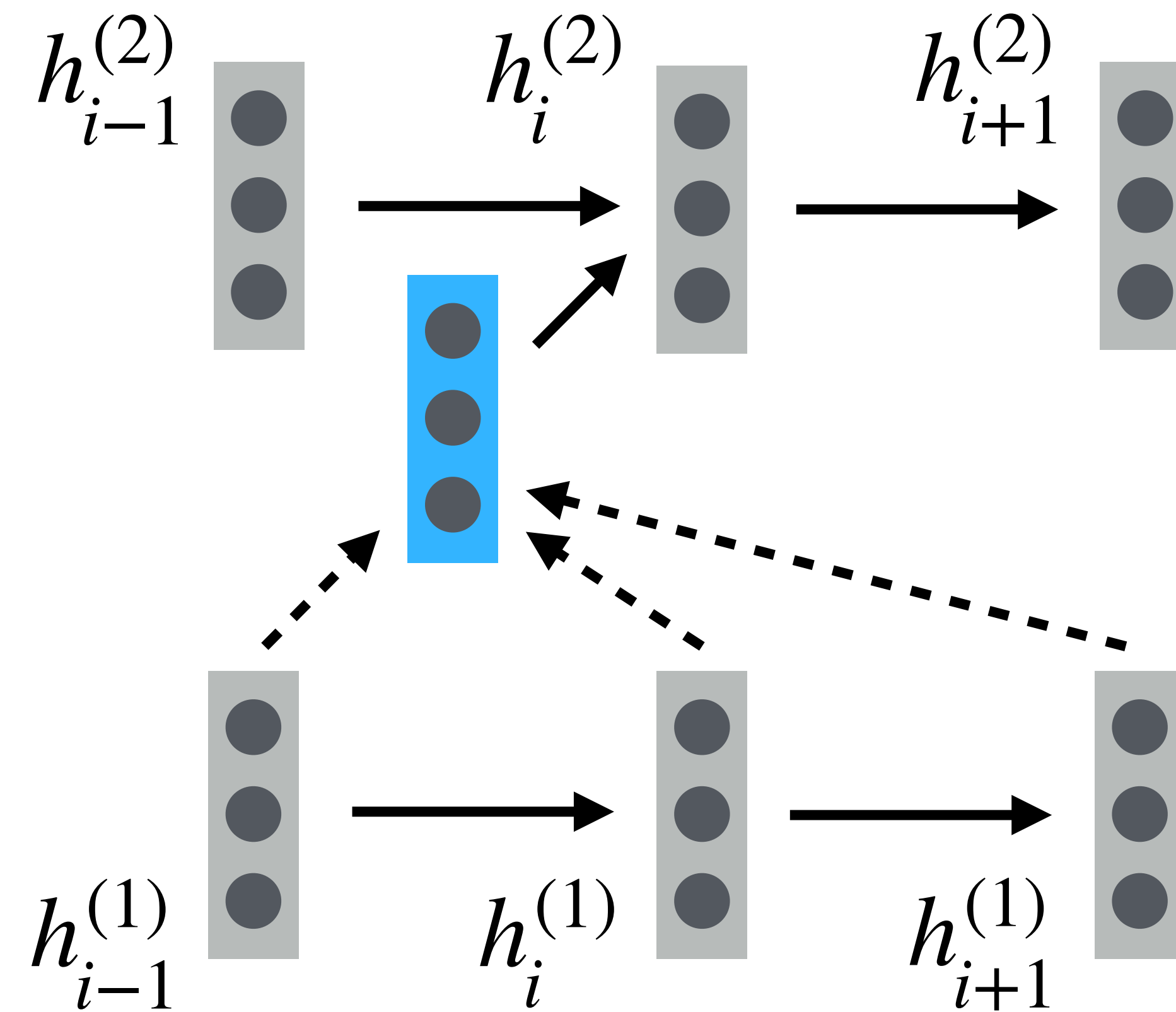
context repr $c_i = h_{\alpha_i}$

nondifferentiable!

but sometimes better generalization

(now you know how to
build anything)

Self-attention revisited



Next class: transformers