

# Conditional Language Modeling with Attention

Chris Dyer



# Review: Conditional LMs

A **conditional language model** assigns probabilities to sequences of words,  $\mathbf{w} = (w_1, w_2, \dots, w_\ell)$ , given some conditioning context,  $\mathbf{x}$ .

As with unconditional models, it is again helpful to use the chain rule to decompose this probability:

$$p(\mathbf{w} \mid \mathbf{x}) = \prod_{t=1}^{\ell} p(w_t \mid \mathbf{x}, w_1, w_2, \dots, w_{t-1})$$

*What is the probability of the next word, given the history of previously generated words **and** conditioning context  $\mathbf{x}$ ?*



明子は

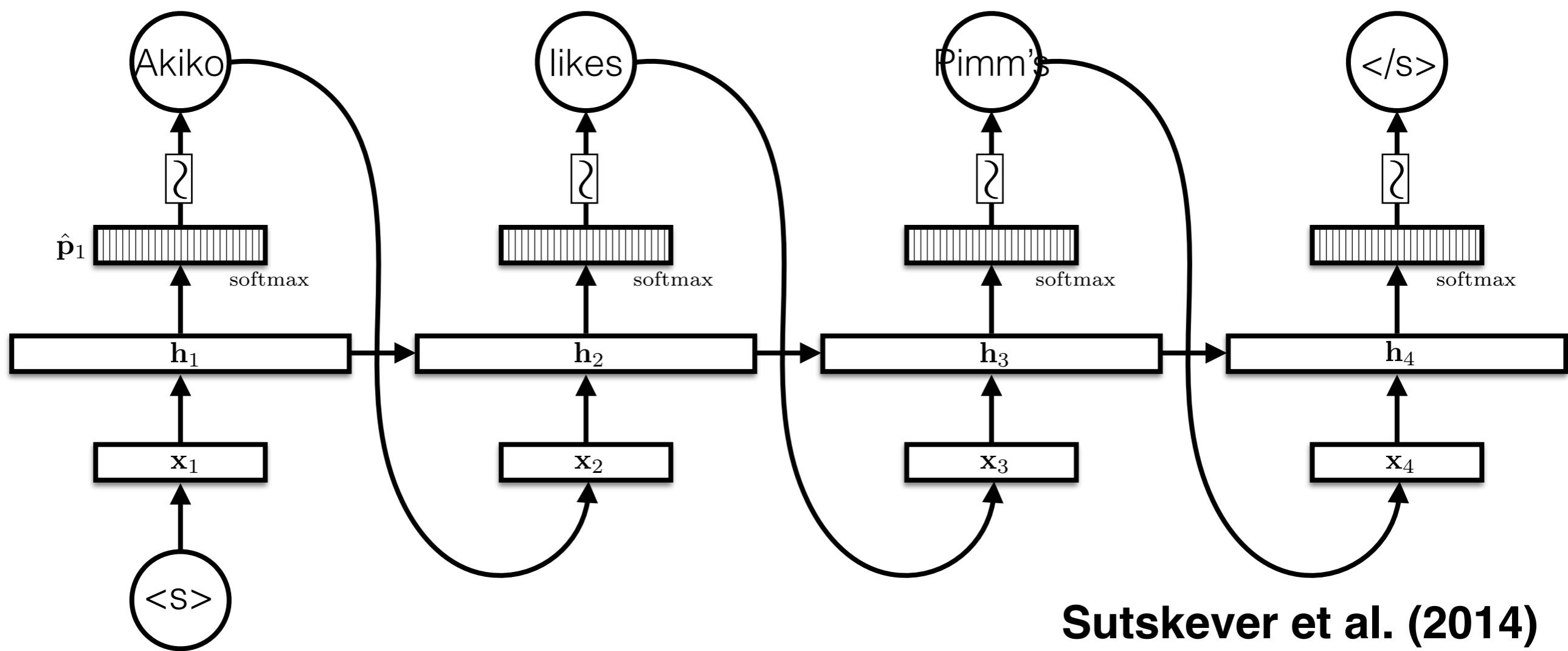
ピムスが

好きです

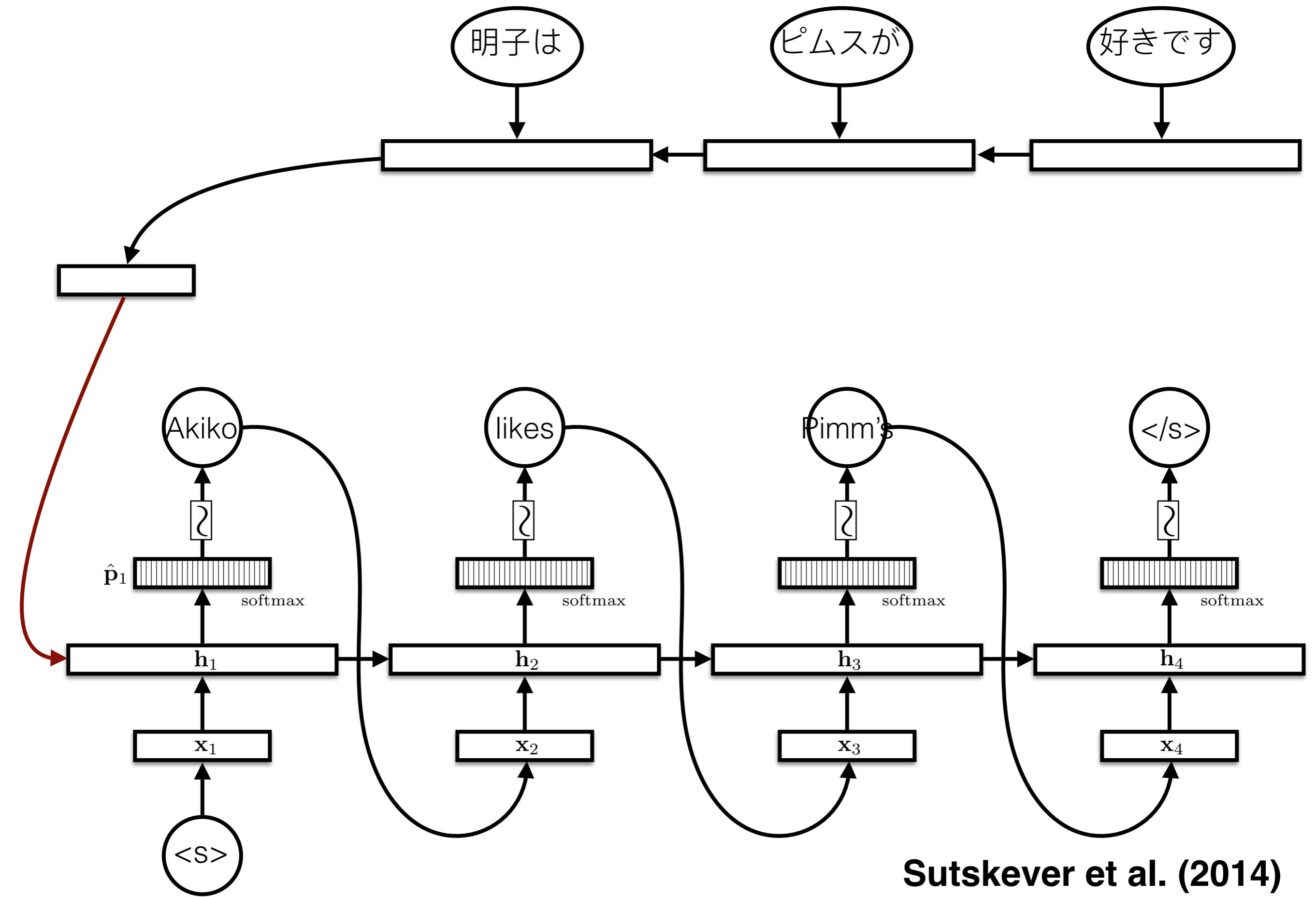
明子は

ピムスが

好きです



Sutskever et al. (2014)



# Conditioning with vectors

We are compressing a lot of information in a finite-sized vector.

# Conditioning with vectors

We are compressing a lot of information in a finite-sized vector.



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

Prof. Ray Mooney

# Conditioning with vectors

We are compressing a lot of information in a finite-sized vector.

Gradients have a long way to travel. Even LSTMs forget!

# Conditioning with vectors

We are compressing a lot of information in a finite-sized vector.

Gradients have a long way to travel. Even LSTMs forget!

**What is to be done?**

# Outline of Lecture

- Machine translation with attention
- Image caption generation with attention

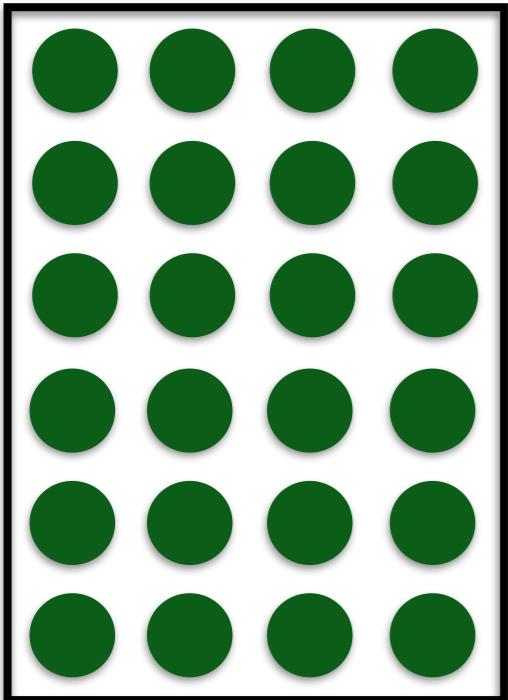
# Solving the Vector Problem in Translation

- Represent a source sentence as a matrix
- Generate a target sentence from a matrix
- This will
  - Solve the capacity problem
  - Solve the gradient flow problem

# Sentences as Matrices

- Problem with the fixed-size vector model
  - Sentences are of different sizes but vectors are of the same size
- Solution: use matrices instead
  - Fixed number of rows, but number of columns depends on the number of words
  - Usually  $|f| = \#cols$

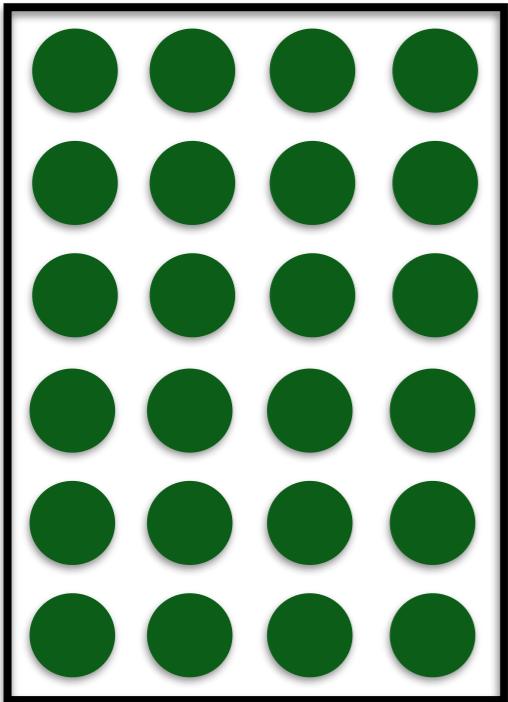
# Sentences as Matrices



*Ich möchte ein Bier*

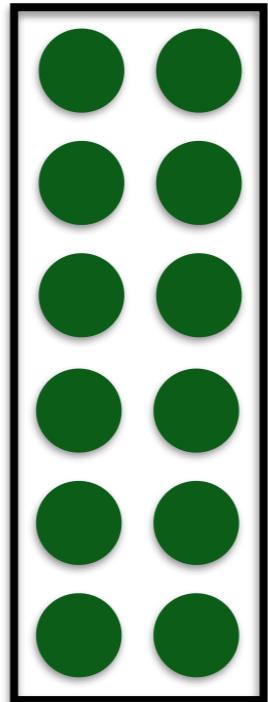
# Sentences as Matrices

动态的matrices

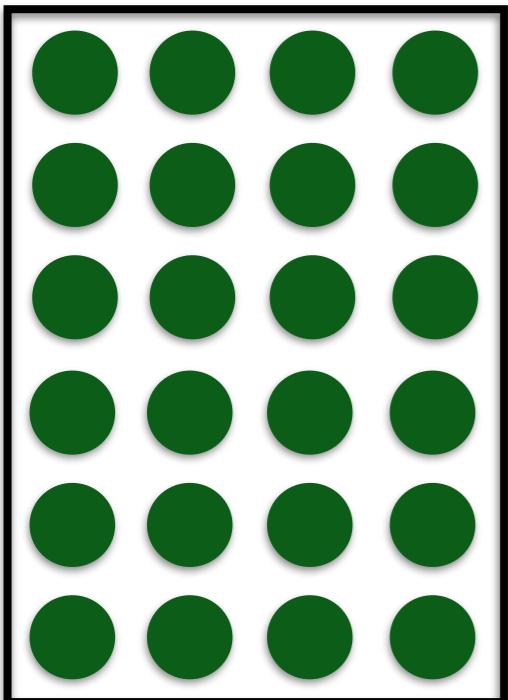


*Ich möchte ein Bier*

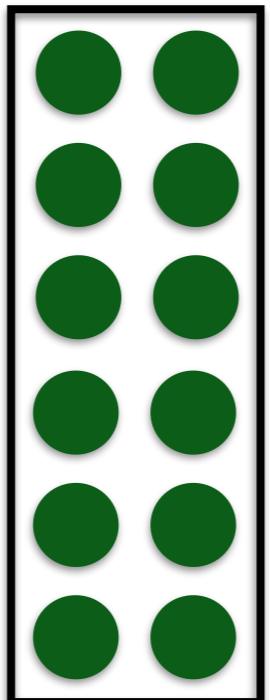
*Mach's gut*



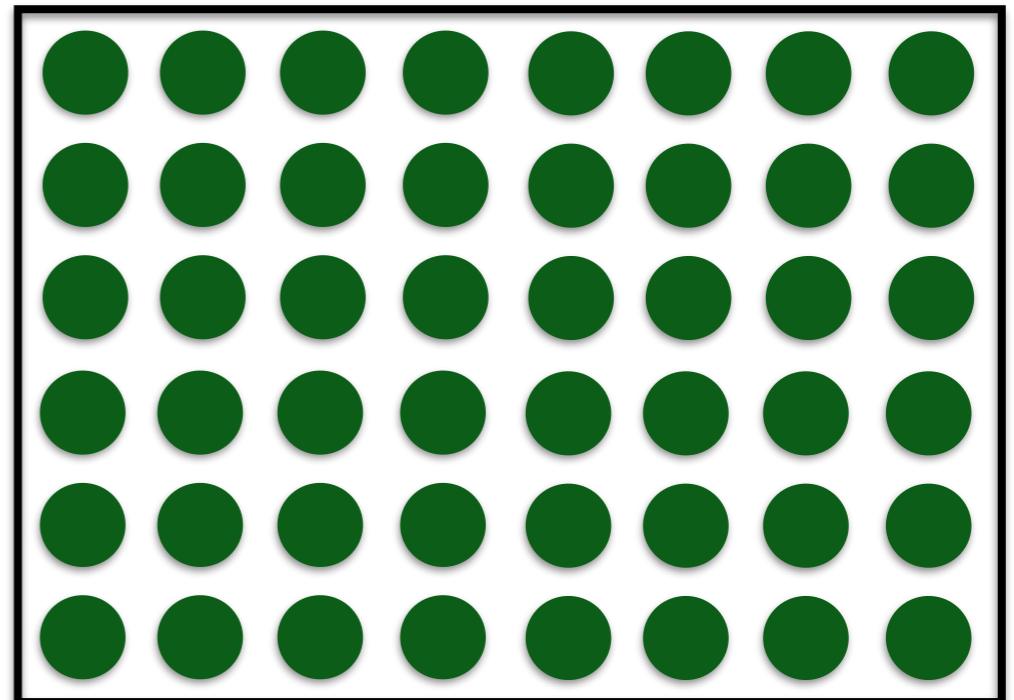
# Sentences as Matrices



*Ich möchte ein Bier*

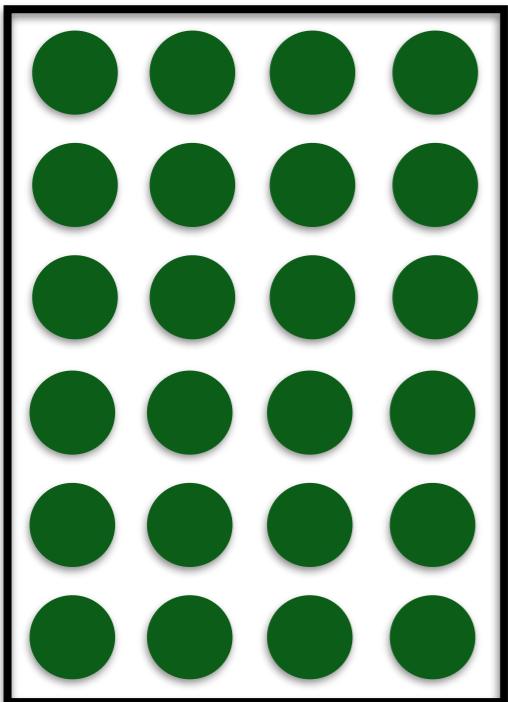


*Mach's gut*



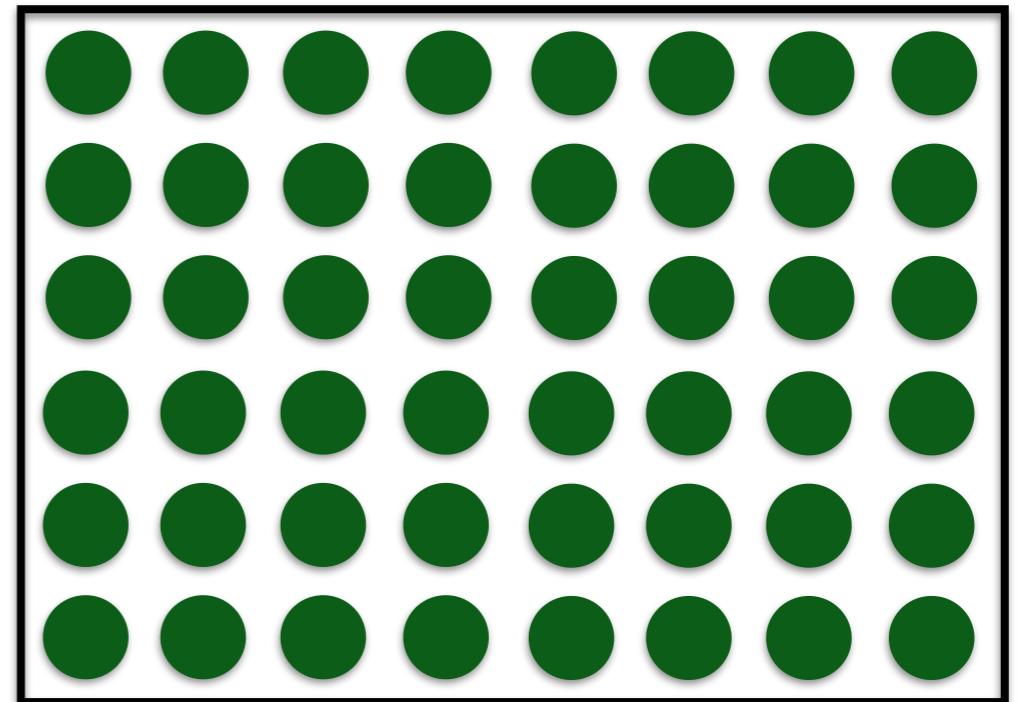
*Die Wahrheiten der Menschen sind die unwiderlegbaren Irrtümer*

# Sentences as Matrices



*Ich möchte ein Bier*

*Mach's gut*



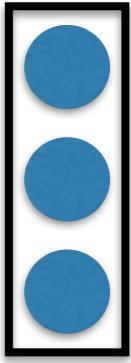
*Die Wahrheiten der Menschen sind die unwiderlegbaren Irrtümer*

**Question: How do we build these matrices?**

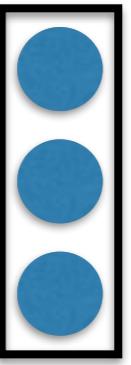
# With Concatenation

- Each word type is represented by an n-dimensional vector
- Take all of the vectors for the sentence and concatenate them into a matrix
- Simplest possible model
  - So simple, no one has bothered to publish how well/badly it works!

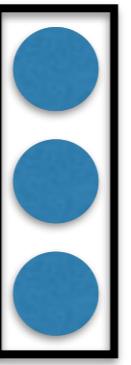
$x_1$



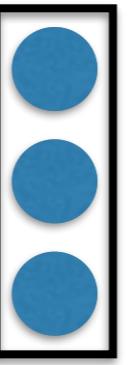
$x_2$



$x_3$



$x_4$



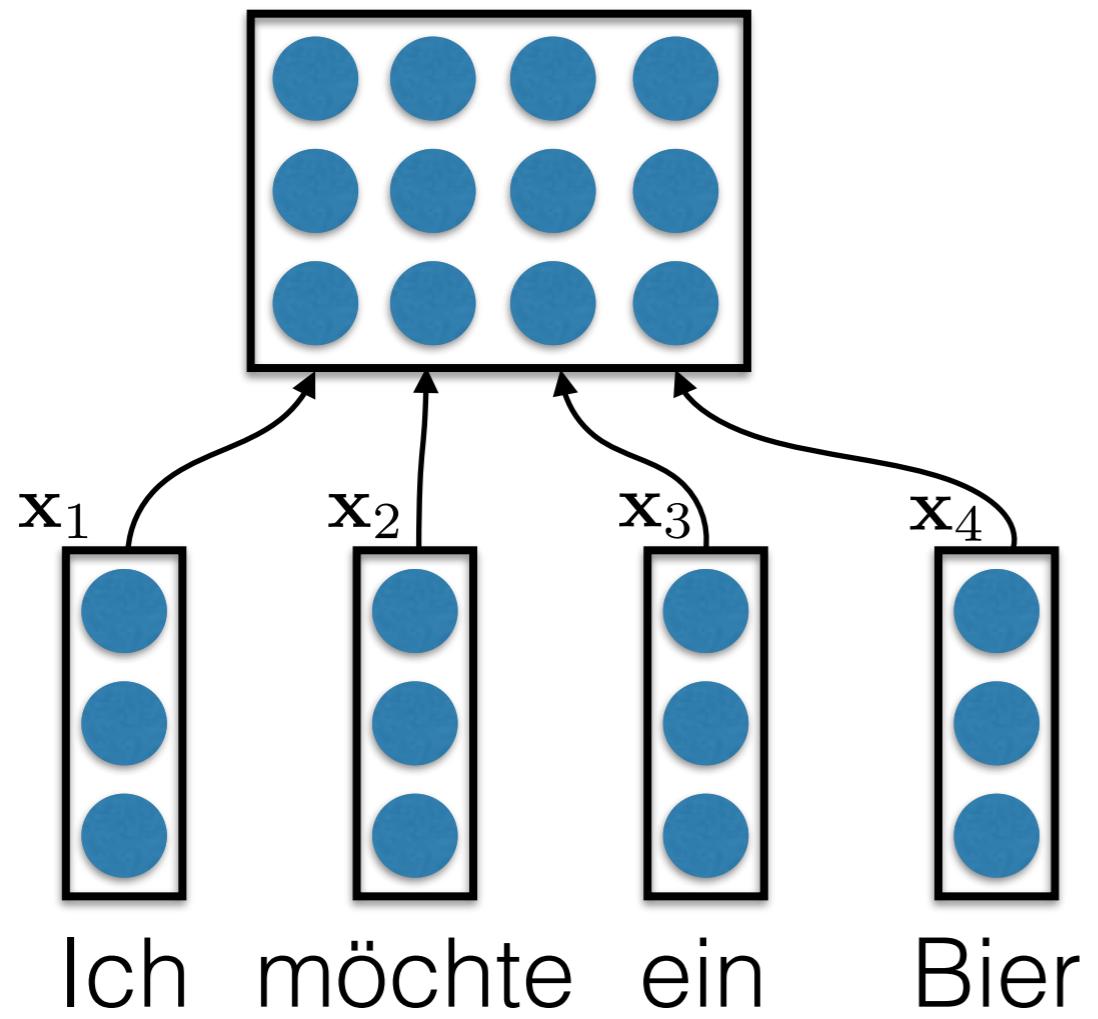
Ich

möchte

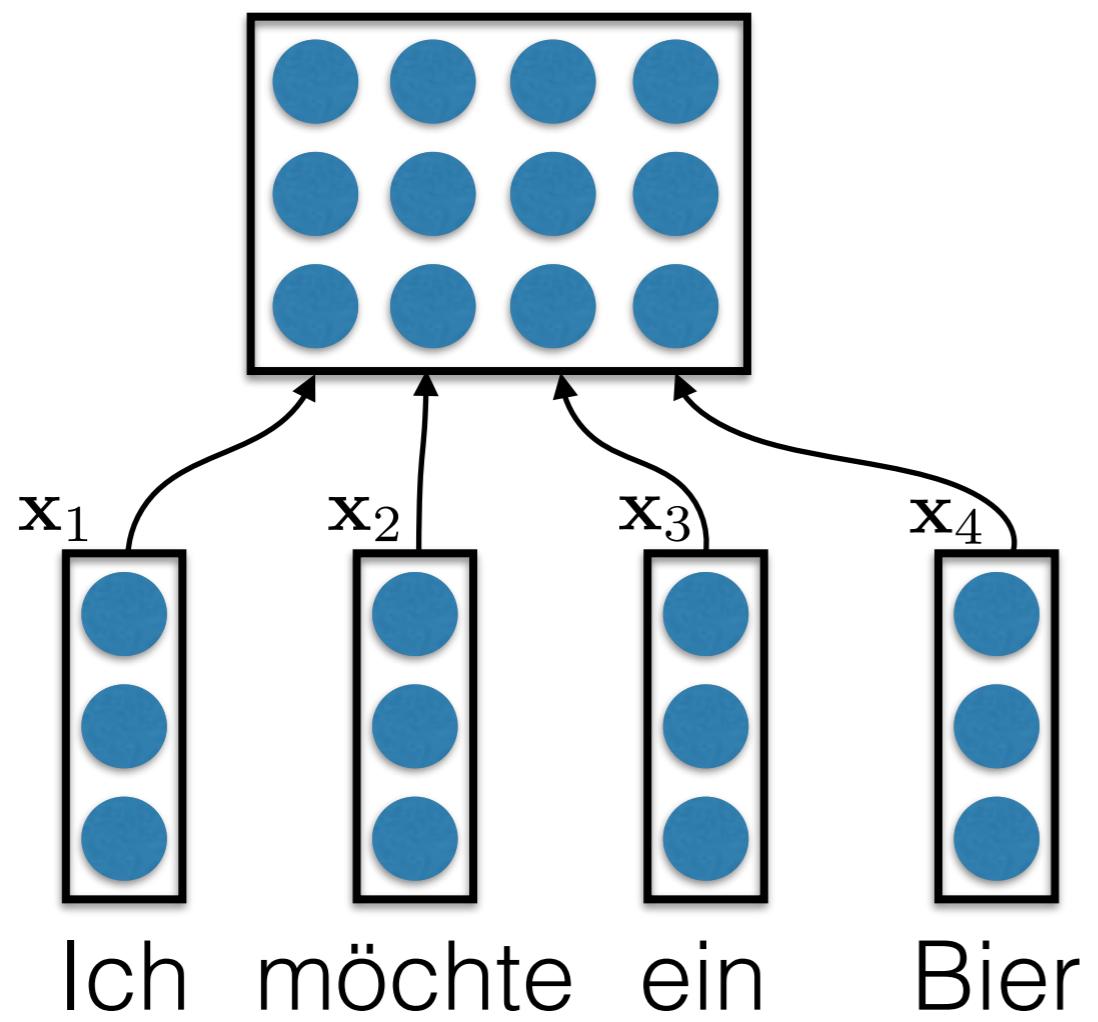
ein

Bier

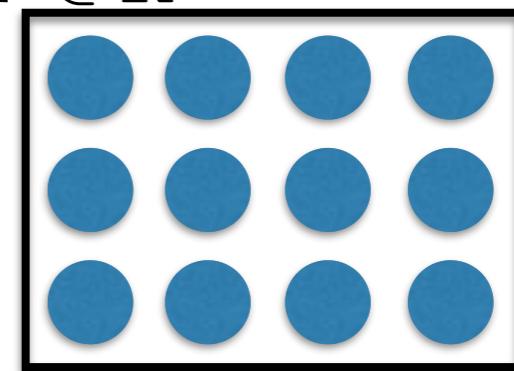
$$\mathbf{f}_i = \mathbf{x}_i$$



$$\mathbf{f}_i = \mathbf{x}_i$$



$$\mathbf{F} \in \mathbb{R}^{n \times |\mathcal{F}|}$$

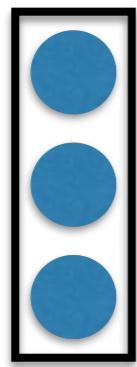


*Ich möchte ein Bier*

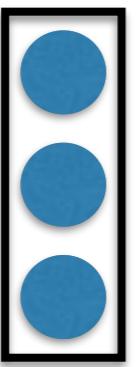
# With Convolutional Nets

- Apply convolutional networks to transform the naive concatenated matrix to obtain a context-dependent matrix
- Explored in a recent ICLR submission by Gehring et al., 2016 (from FAIR)
  - Closely related to the neural translation model proposed by Kalchbrenner and Blunsom, 2013
- Note: convnets usually have a “pooling” operation at the top level that results in a fixed-sized representation. For sentences, leave this out.

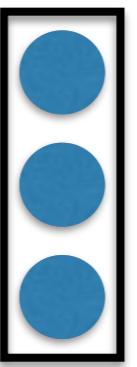
$x_1$



$x_2$



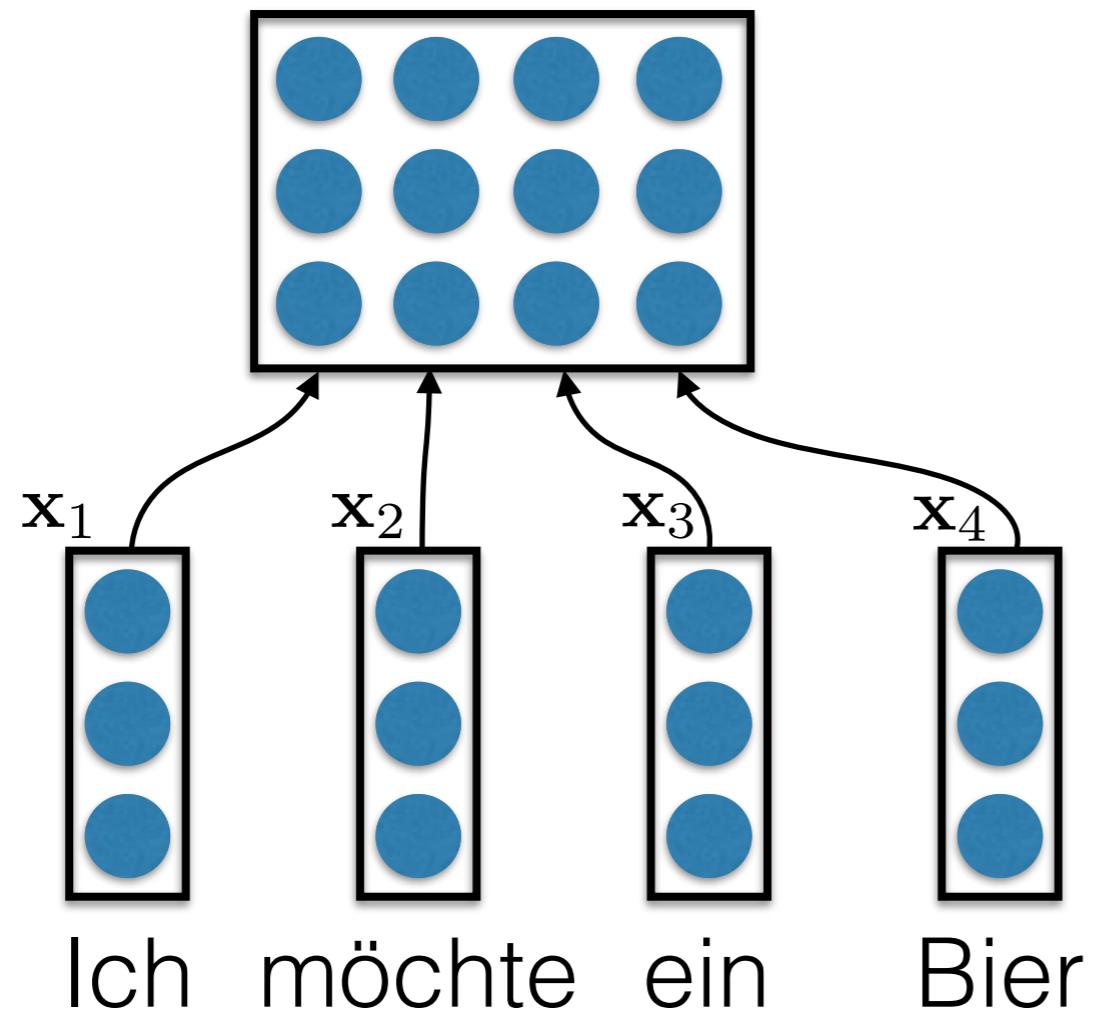
$x_3$

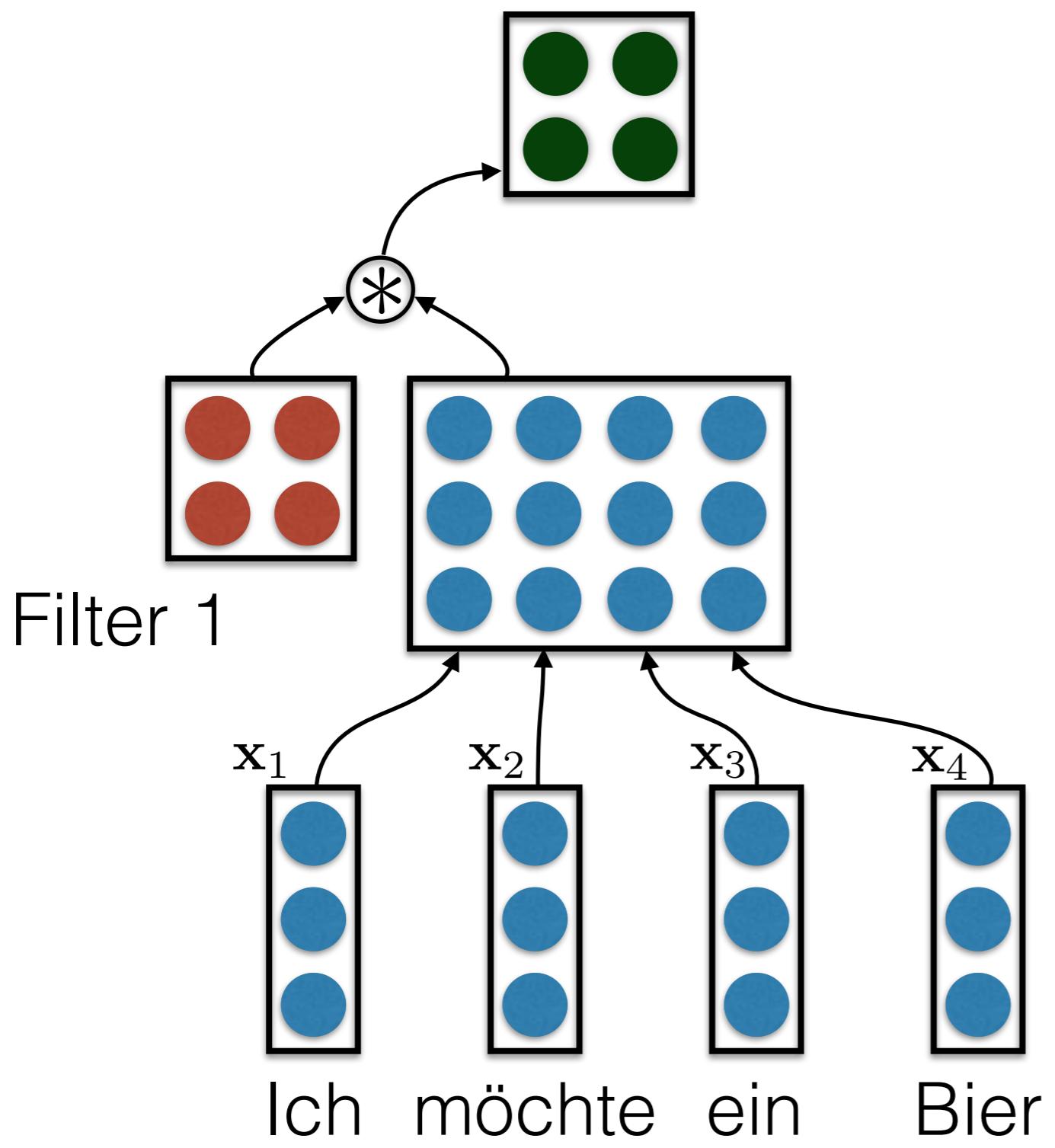


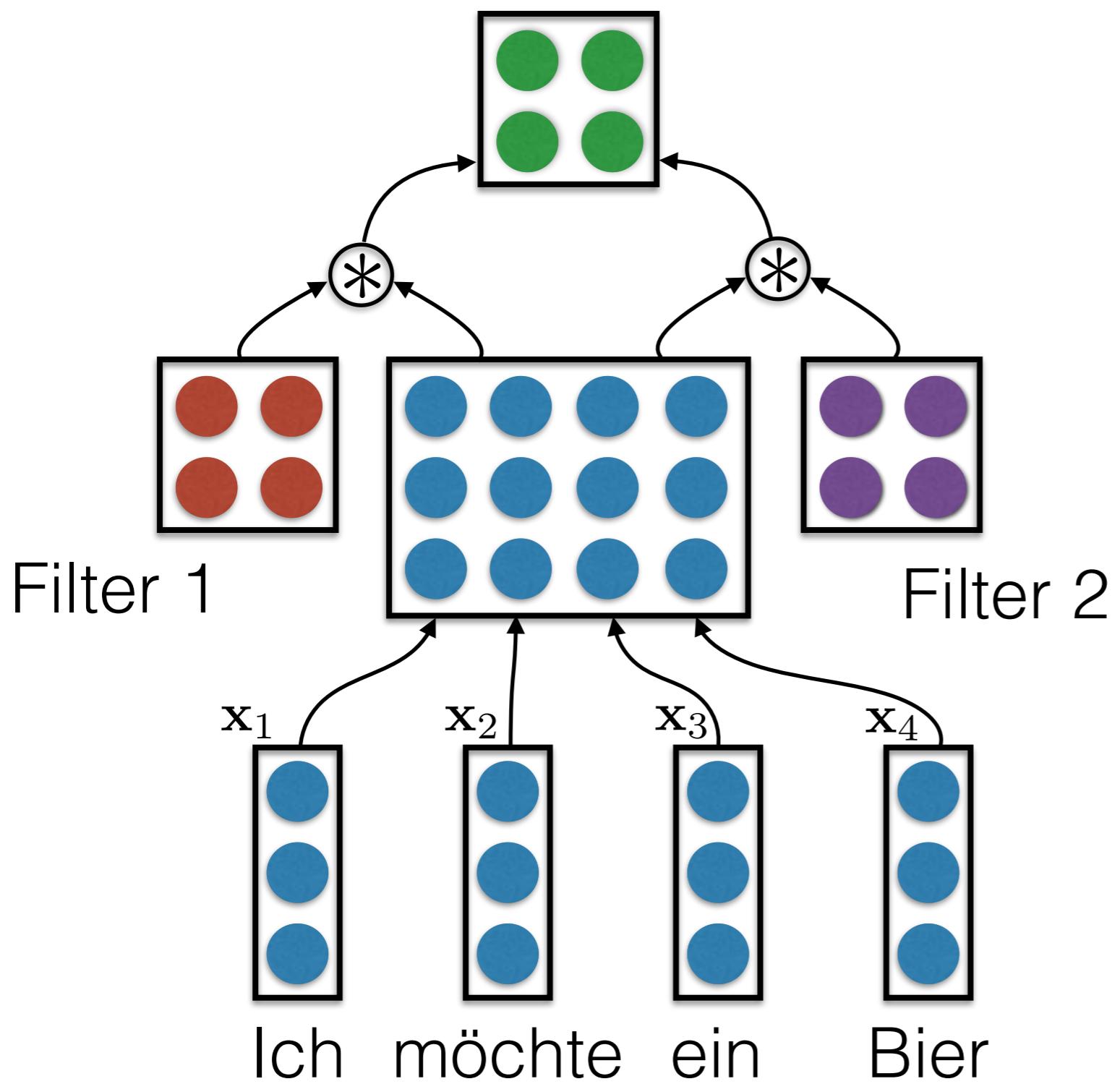
$x_4$

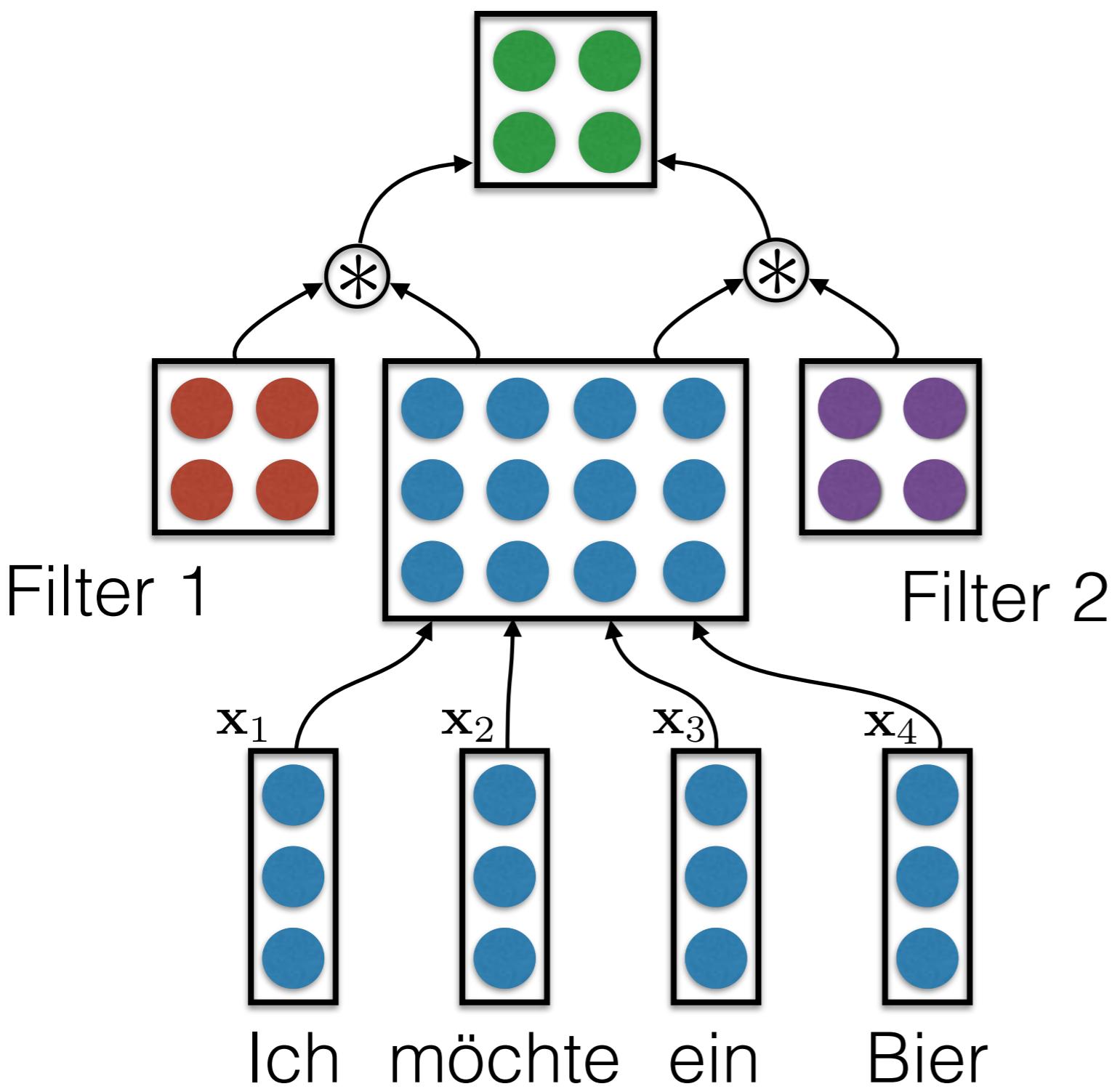


Ich möchte ein Bier

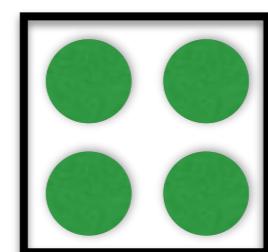








$$\mathbf{F} \in \mathbb{R}^{f(n) \times g(|\mathcal{F}|)}$$

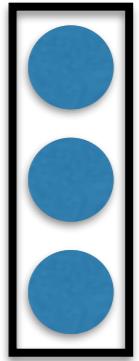


*Ich möchte ein Bier*

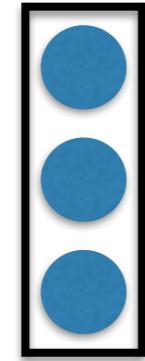
# With Bidirectional RNNs

- By far the most widely used matrix representation, due to Bahdanau et al (2015)
- One column per word
- Each column (word) has two halves concatenated together:
  - a “forward representation”, i.e., a word and its left context
  - a “reverse representation”, i.e., a word and its right context
- Implementation: bidirectional RNNs (GRUs or LSTMs) to read **f** from left to right and right to left, concatenate representations

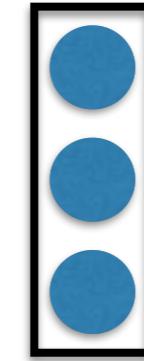
$x_1$



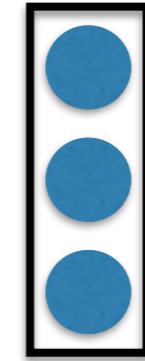
$x_2$



$x_3$

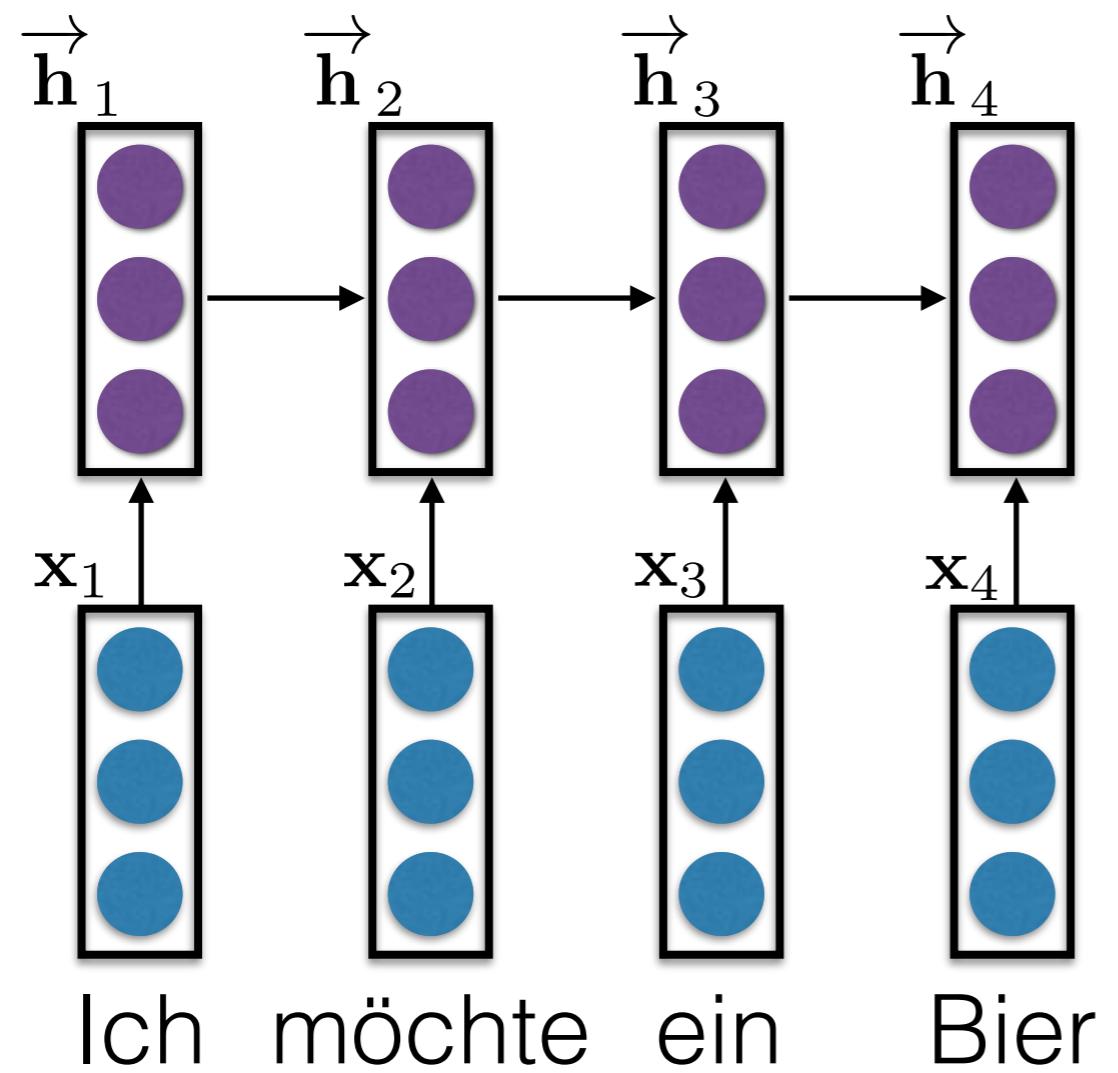


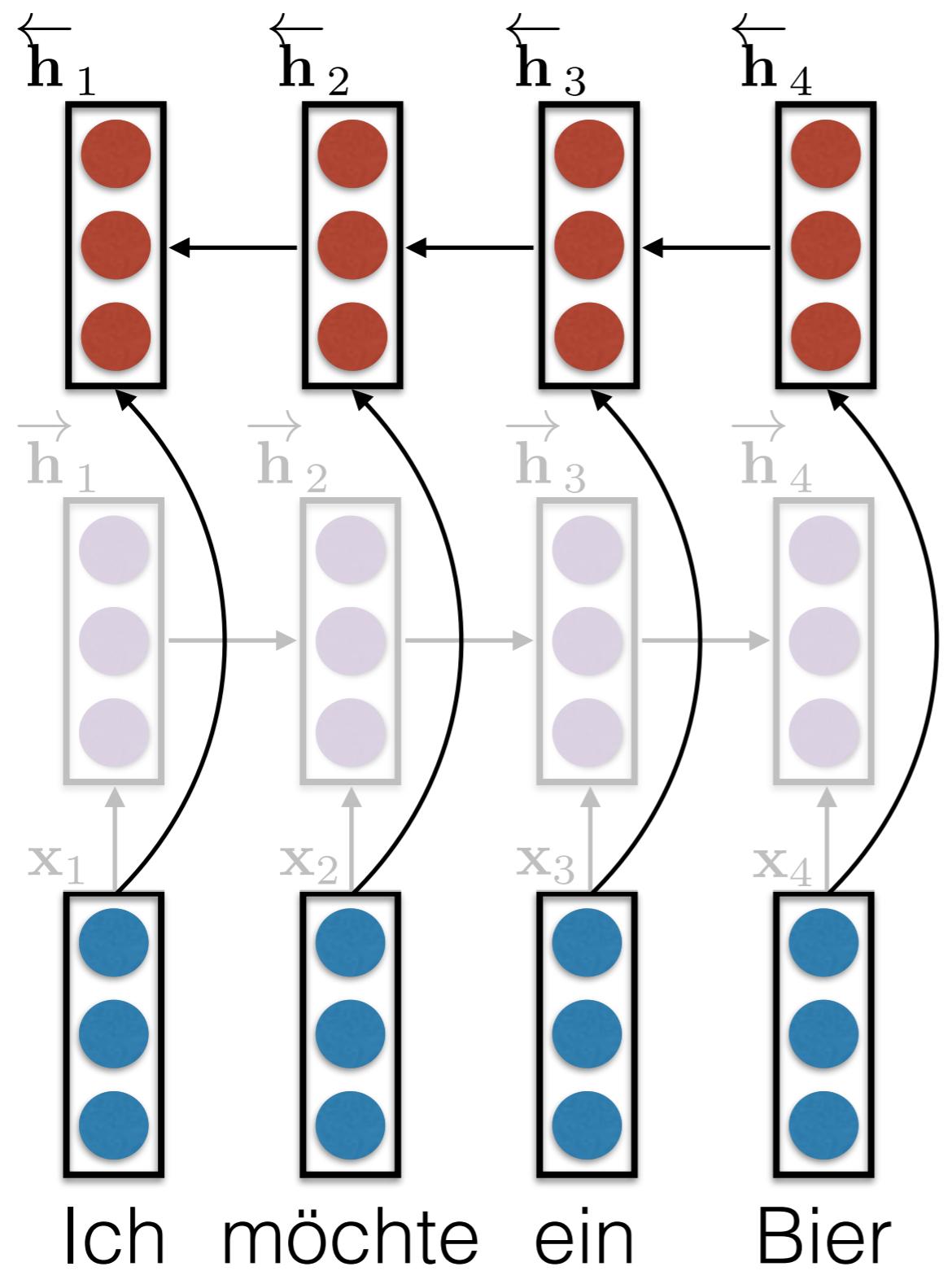
$x_4$



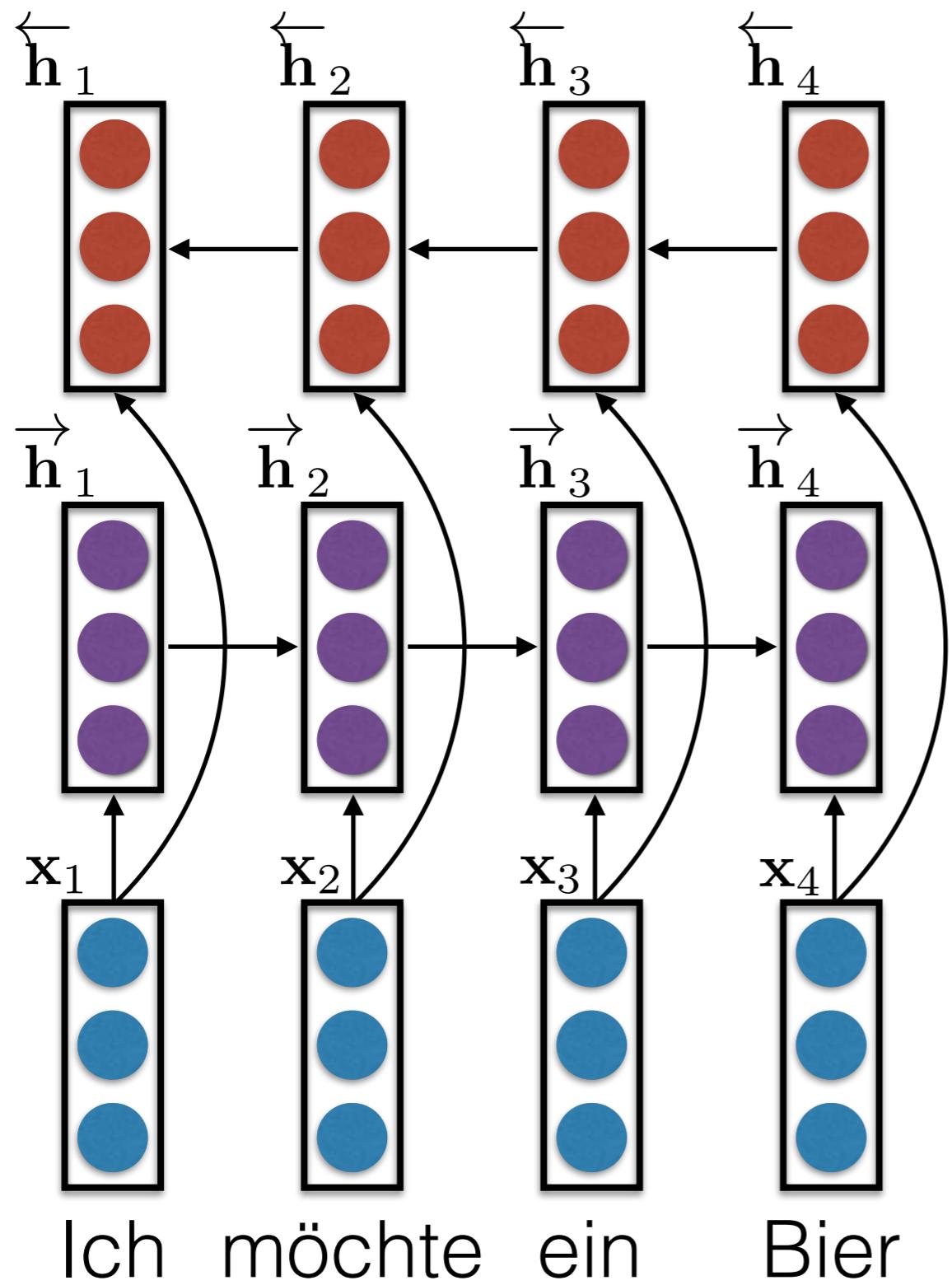
Ich möchte ein

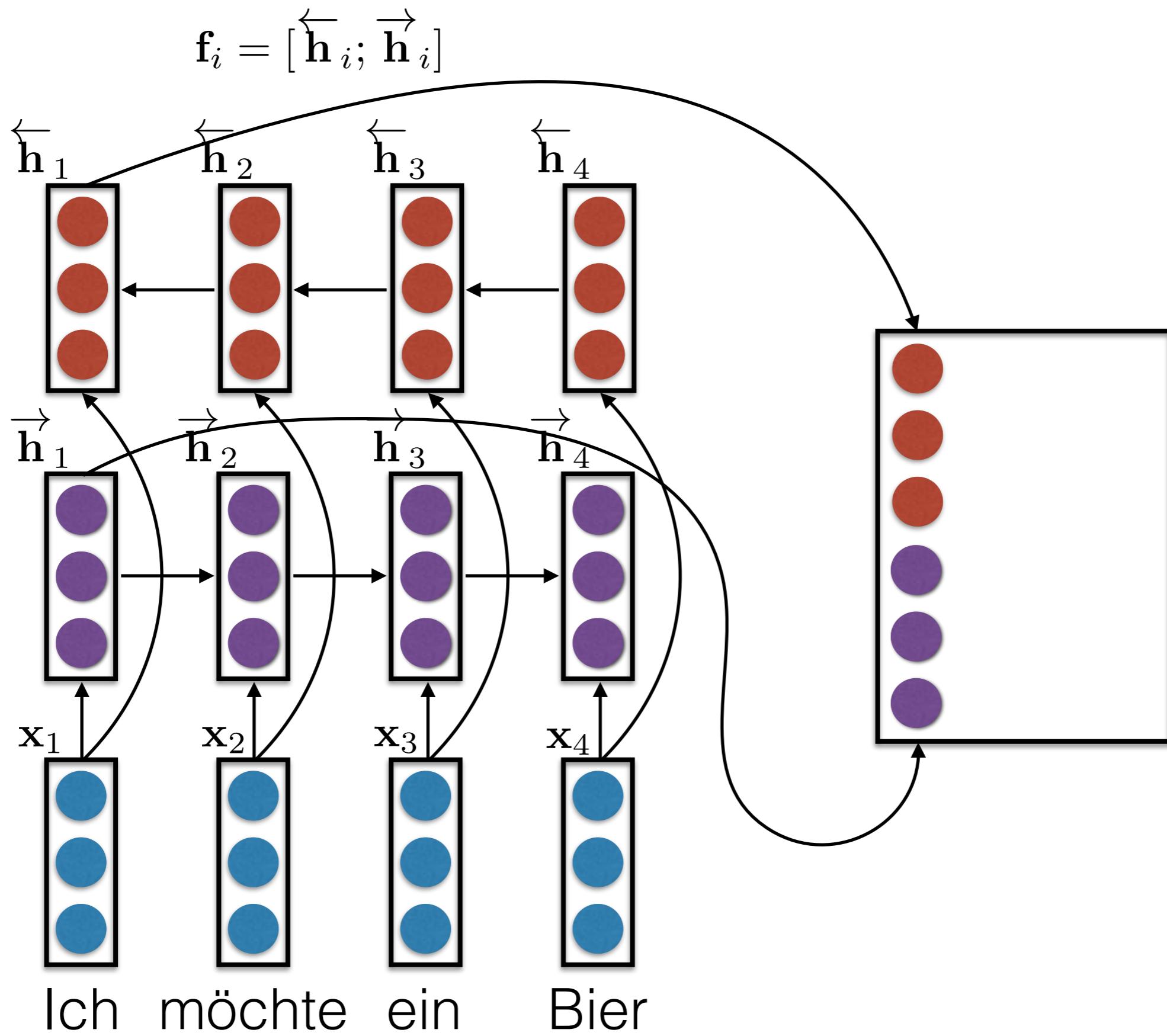
Bier

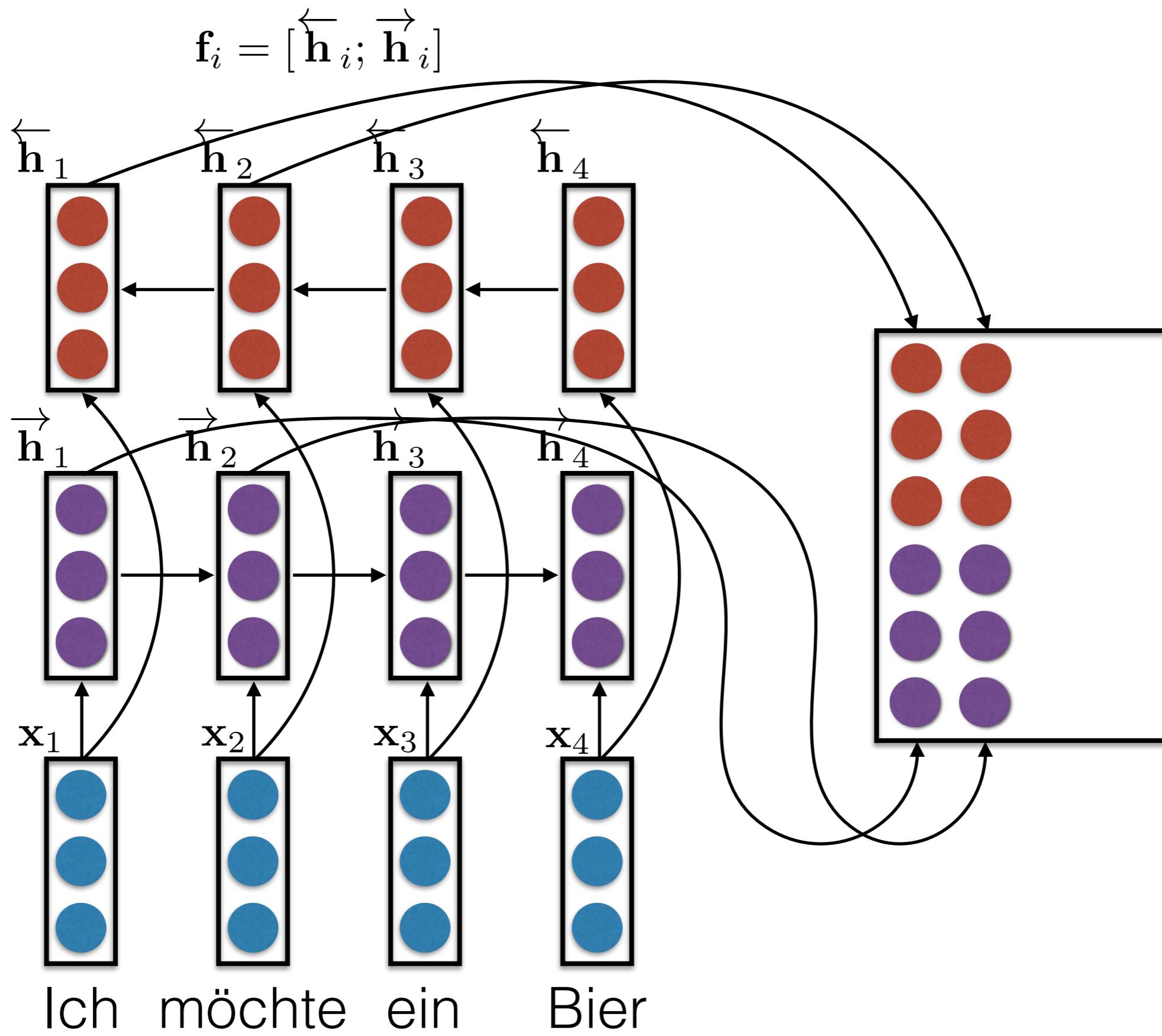


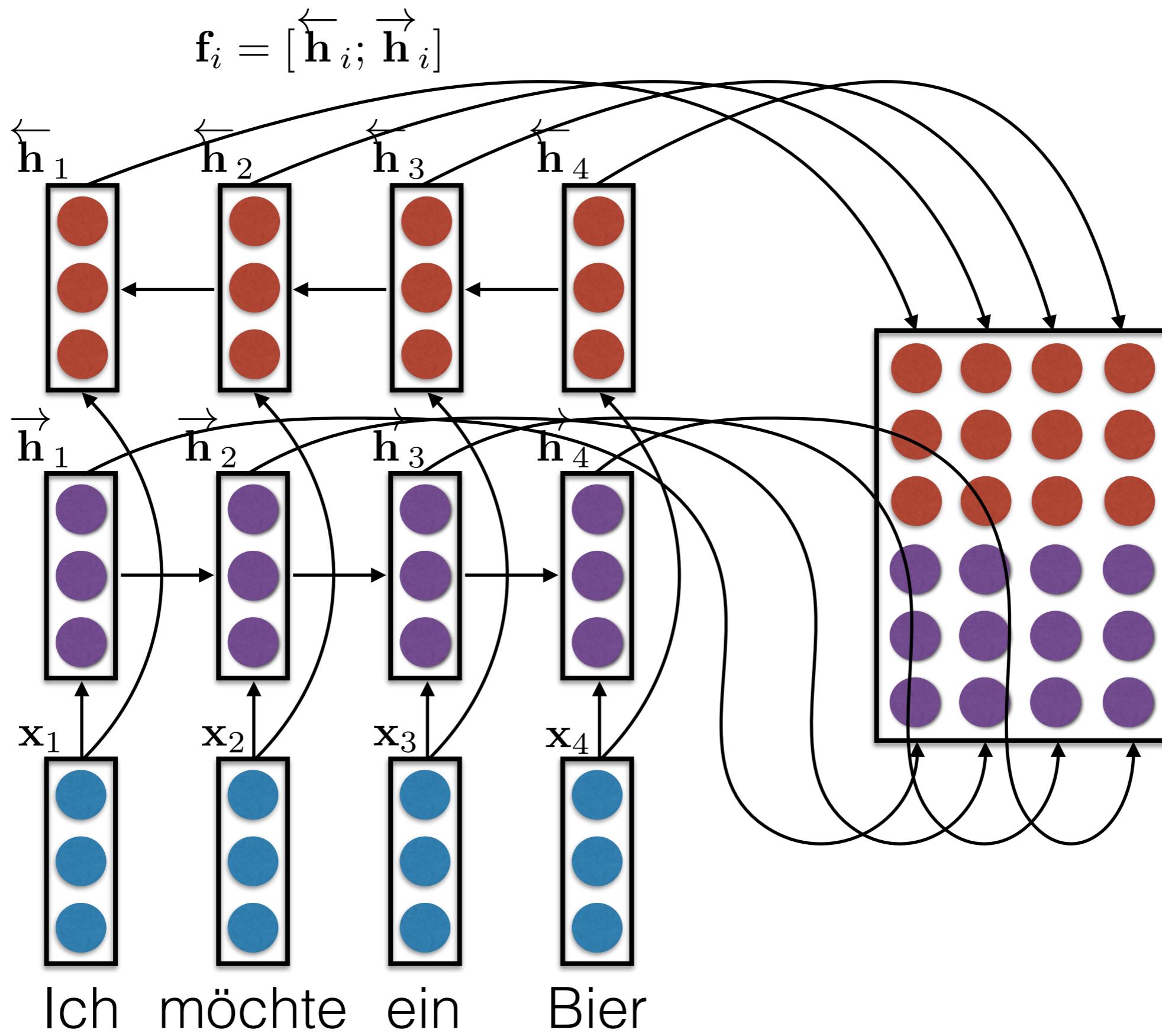


$$\mathbf{f}_i = [\overleftarrow{\mathbf{h}}_i; \overrightarrow{\mathbf{h}}_i]$$

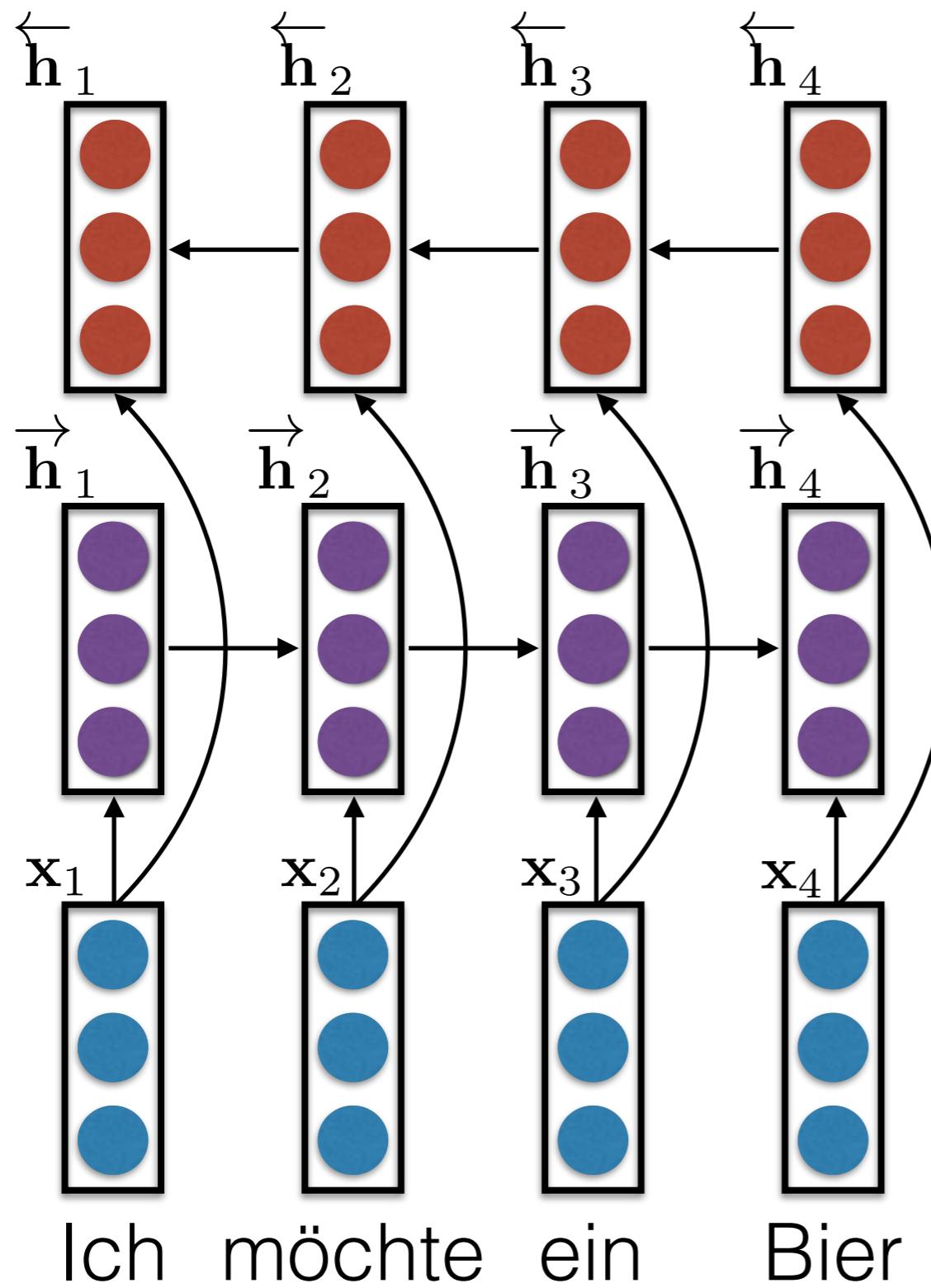




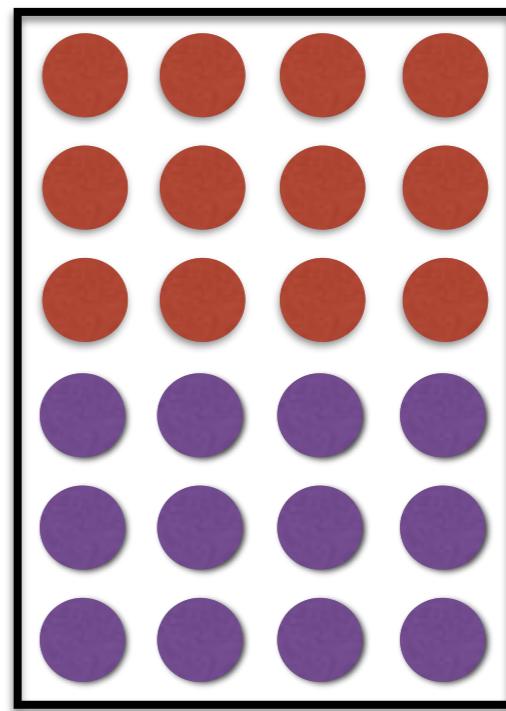




$$\mathbf{f}_i = [\overleftarrow{\mathbf{h}}_i; \overrightarrow{\mathbf{h}}_i]$$



$$\mathbf{F} \in \mathbb{R}^{2n \times |f|}$$



*Ich möchte ein Bier*

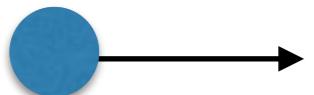
# Where are we in 2017?

- There are lots of ways to construct **F**
  - Very little systematic work comparing them
  - There are many more undiscovered things out there
    - convolutions are particularly interesting and under-explored
    - syntactic information can help (Sennrich & Haddow, 2016; Nadejde et al., 2017), but many more integration strategies are possible
    - try something with phrase types instead of word types?

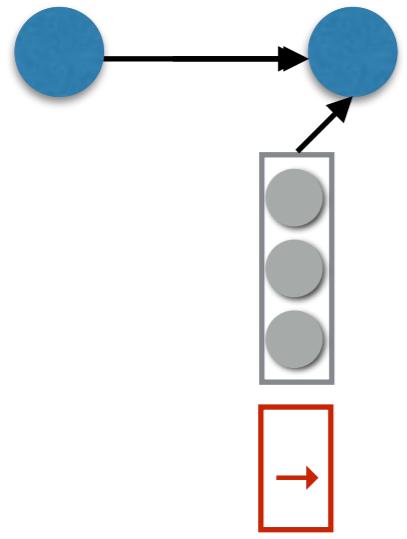
**Multi-word expressions are a pain in the neck .**

# Generation from Matrices

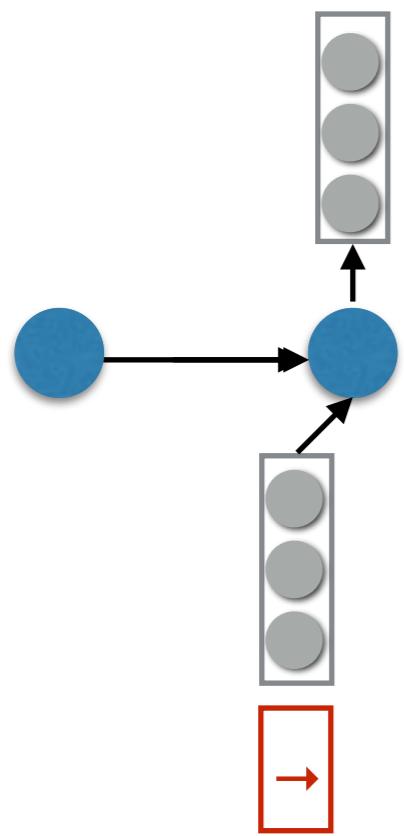
- We have a matrix  $\mathbf{F}$  representing the input, now we need to generate from it
- Bahdanau et al. (2015) were the first to propose using **attention** for translating from matrix-encoded sentences
- High-level idea
  - Generate the output sentence word by word using an RNN
  - At each output position  $t$ , the RNN receives **two** inputs (in addition to any recurrent inputs)
    - a fixed-size vector embedding of the previously generated output symbol  $e_{t-1}$
    - a fixed-size vector encoding a “view” of the input matrix
  - How do we get a fixed-size vector from a matrix that changes over time?
    - Bahdanau et al: do a weighted sum of the columns of  $\mathbf{F}$  (i.e., words) based on how important they are *at the current time step*. (i.e., just a matrix-vector product  $\mathbf{F}\mathbf{a}_t$ )
    - The weighting of the input columns at each time-step ( $\mathbf{a}_t$ ) is called **attention**



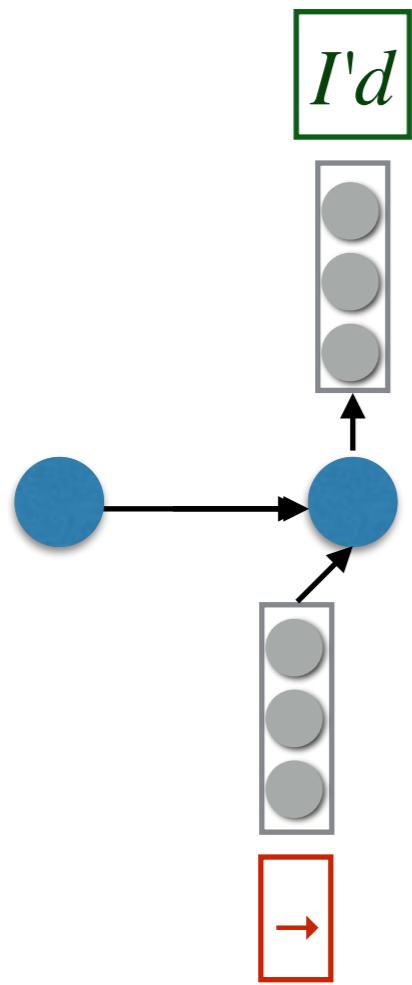
**Recall RNNs...**



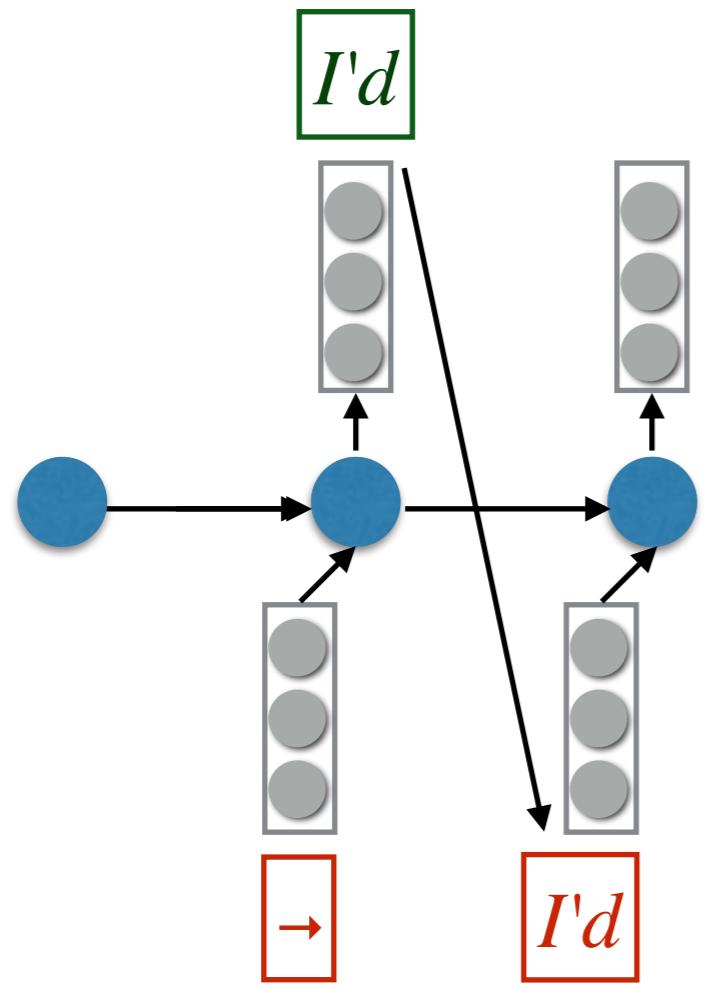
**Recall RNNs...**



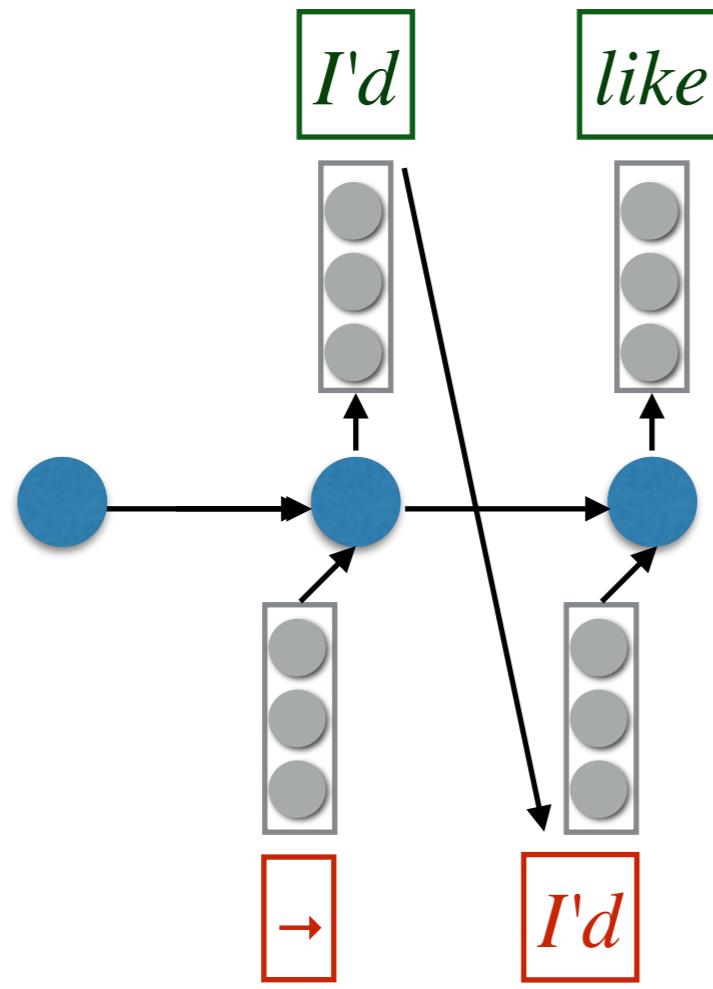
# Recall RNNs...



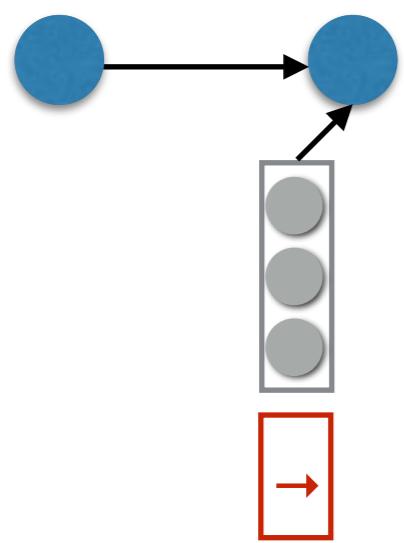
# Recall RNNs...

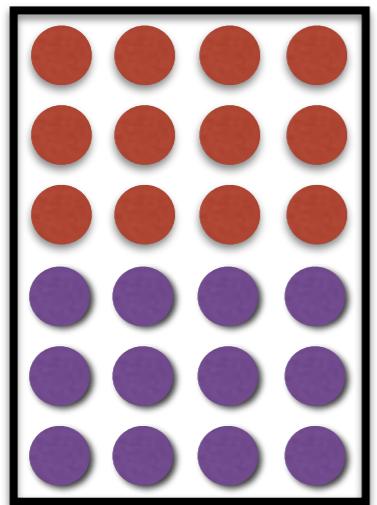
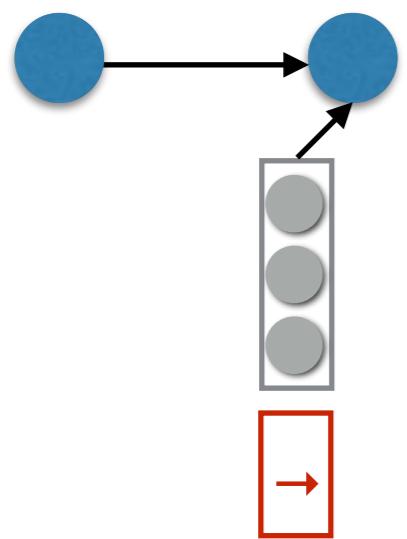


# Recall RNNs...

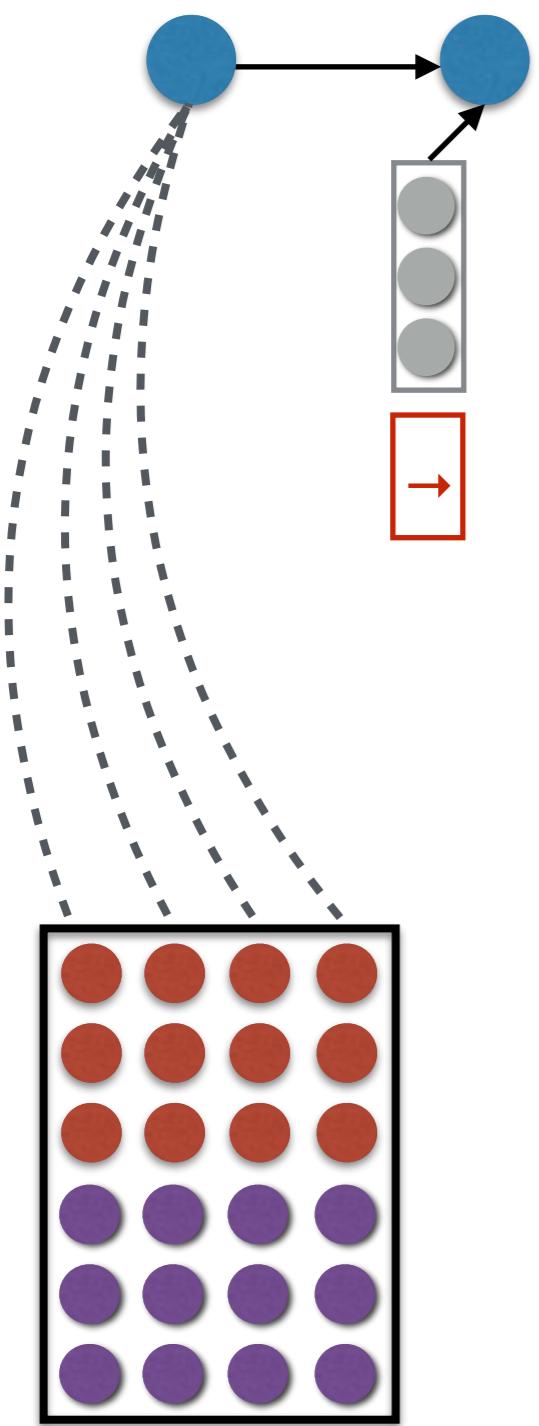


# Recall RNNs...

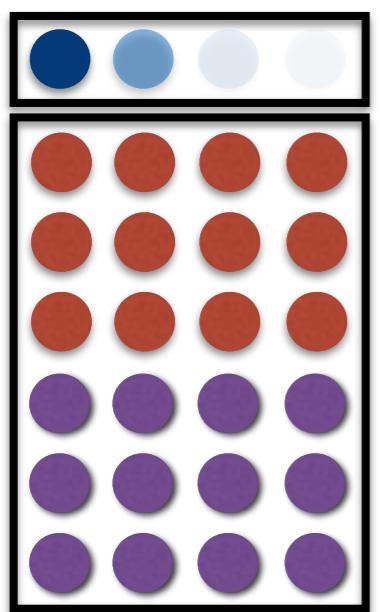
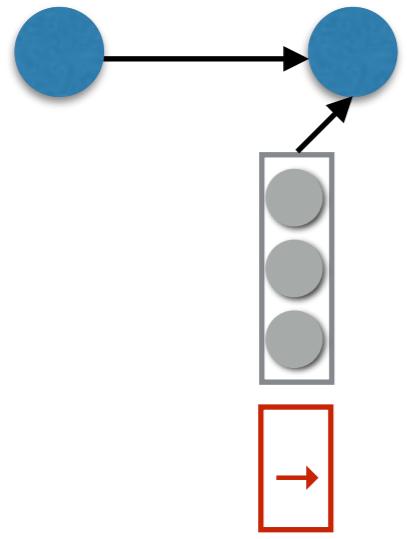




*Ich möchte ein Bier*



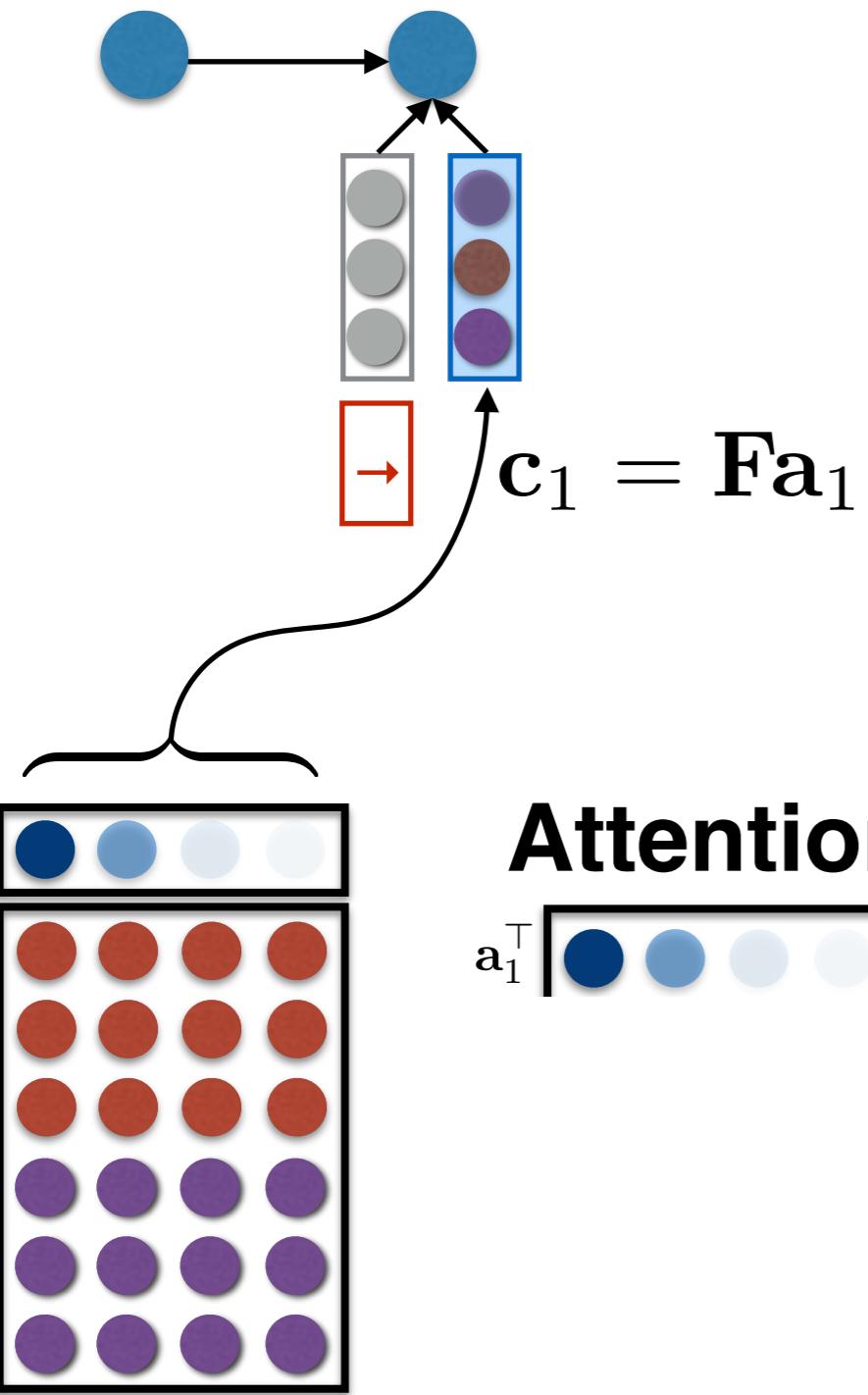
*Ich möchte ein Bier*



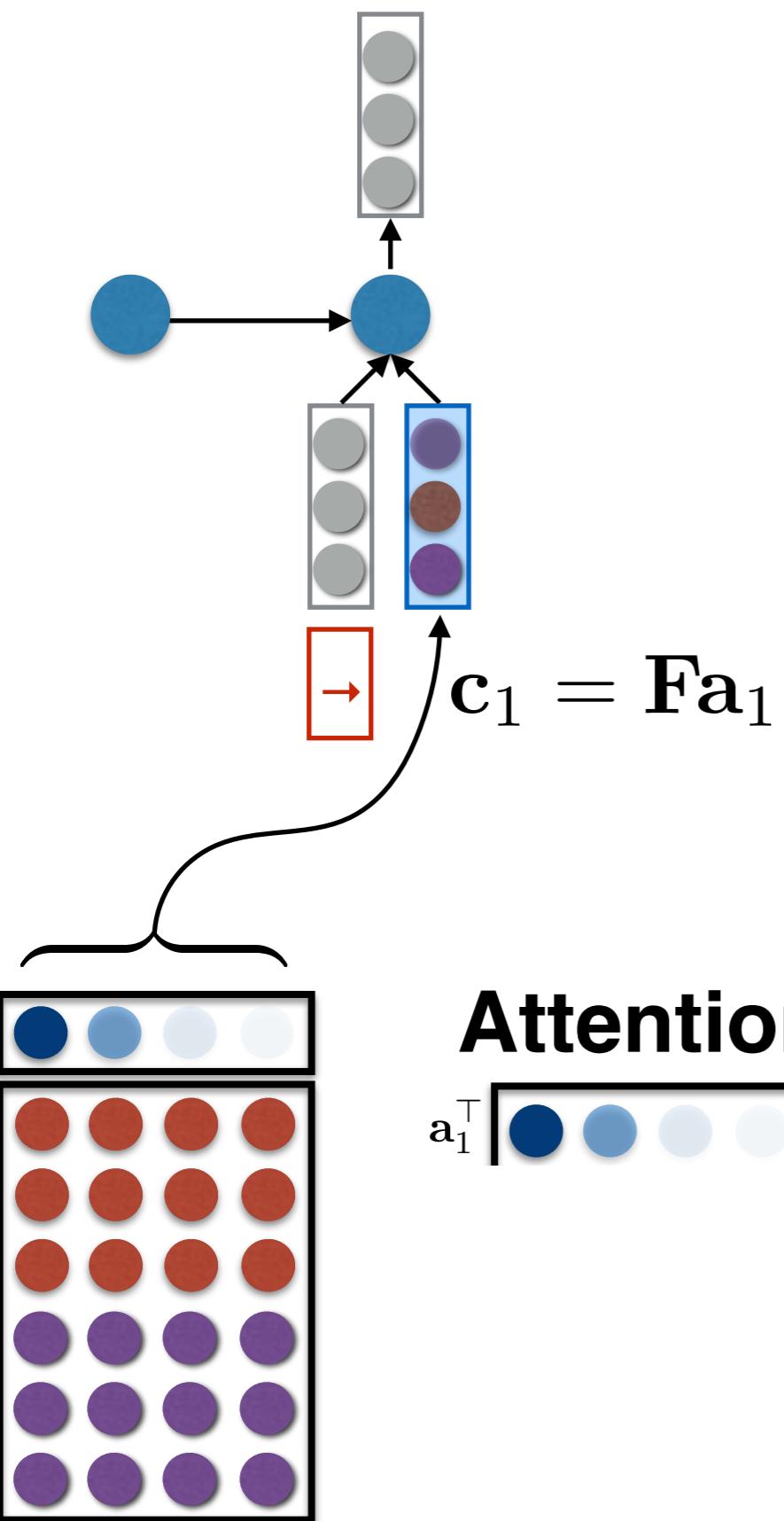
## Attention history:

$$a_1^T \boxed{\text{dark blue} \quad \text{medium blue} \quad \text{light blue} \quad \text{white}}$$

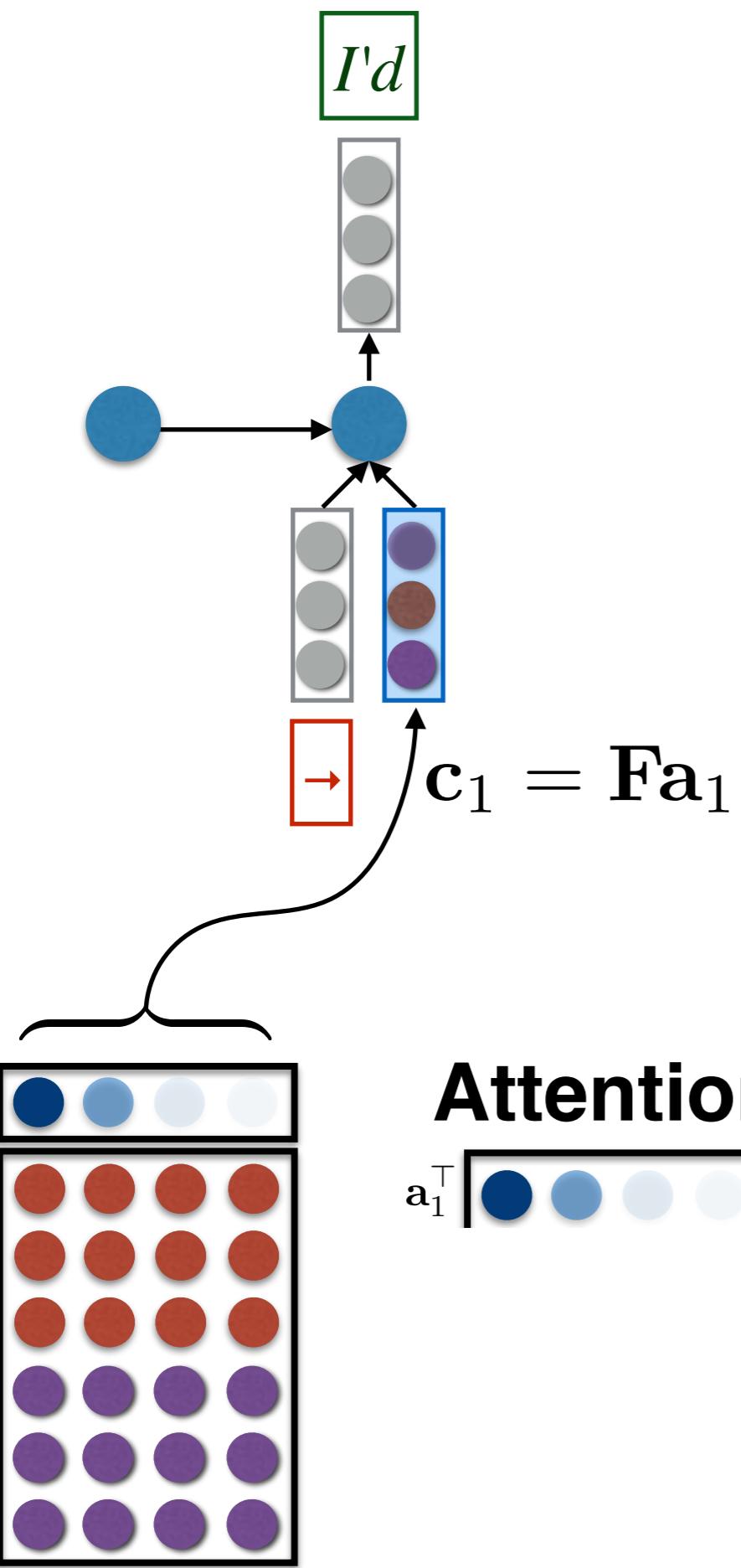
*Ich möchte ein Bier*



*Ich möchte ein Bier*



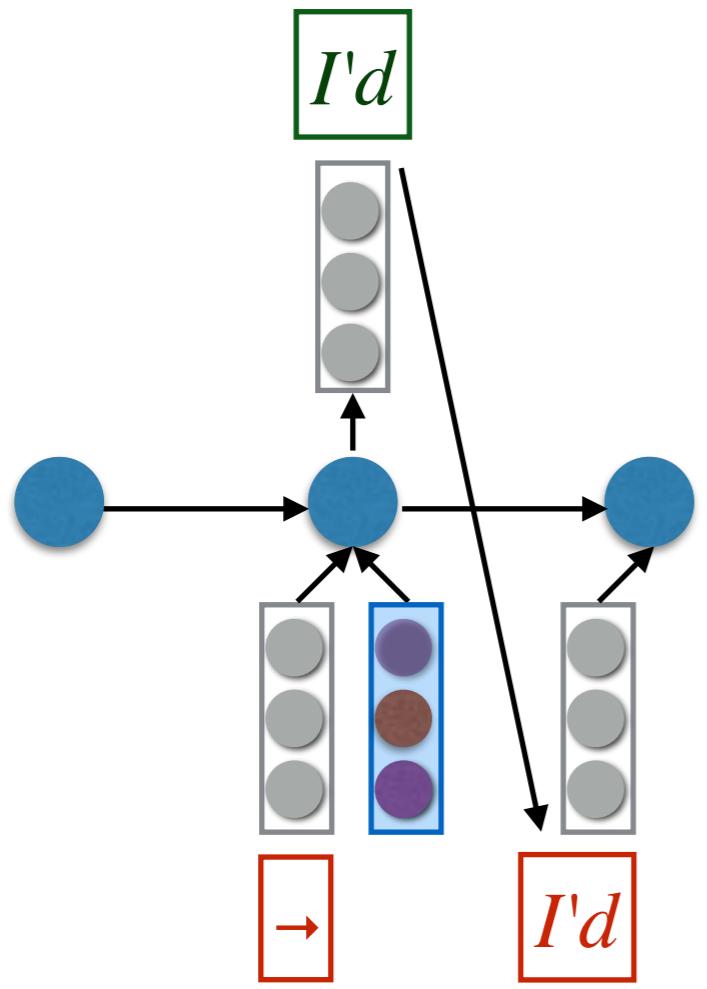
*Ich möchte ein Bier*



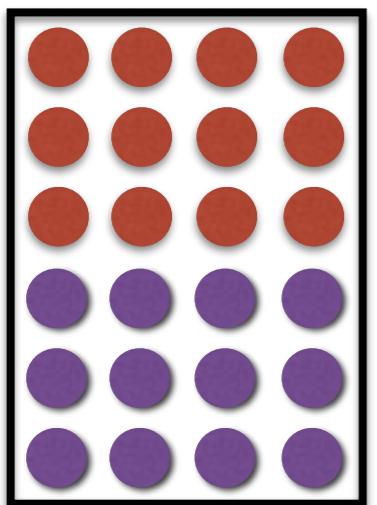
**Attention history:**

$$a_1^T \quad \boxed{\text{blue} \quad \text{red} \quad \text{red} \quad \text{purple}}$$

*Ich möchte ein Bier*

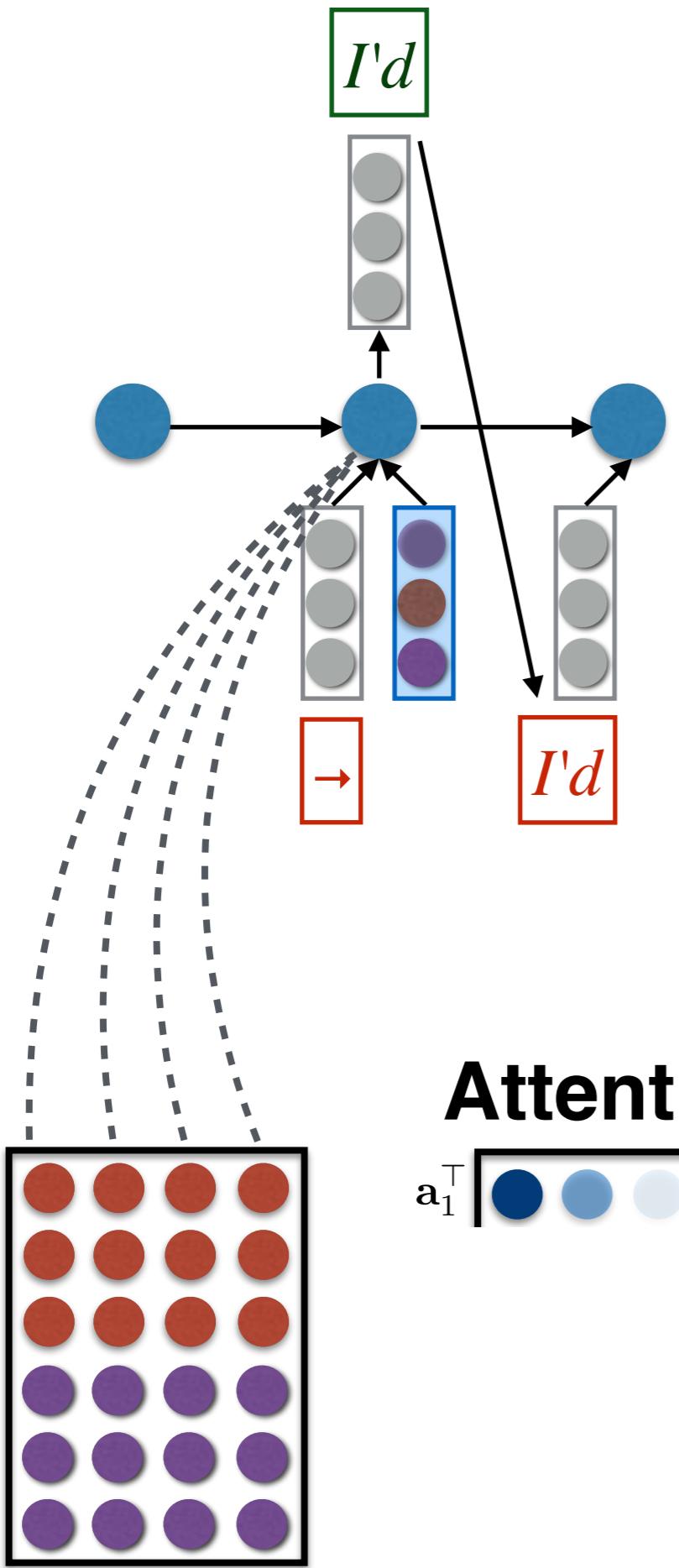


## Attention history:



$$a_1^\top [ \text{blue} \text{ } \text{blue} \text{ } \text{light blue} \text{ } \text{light blue} ]$$

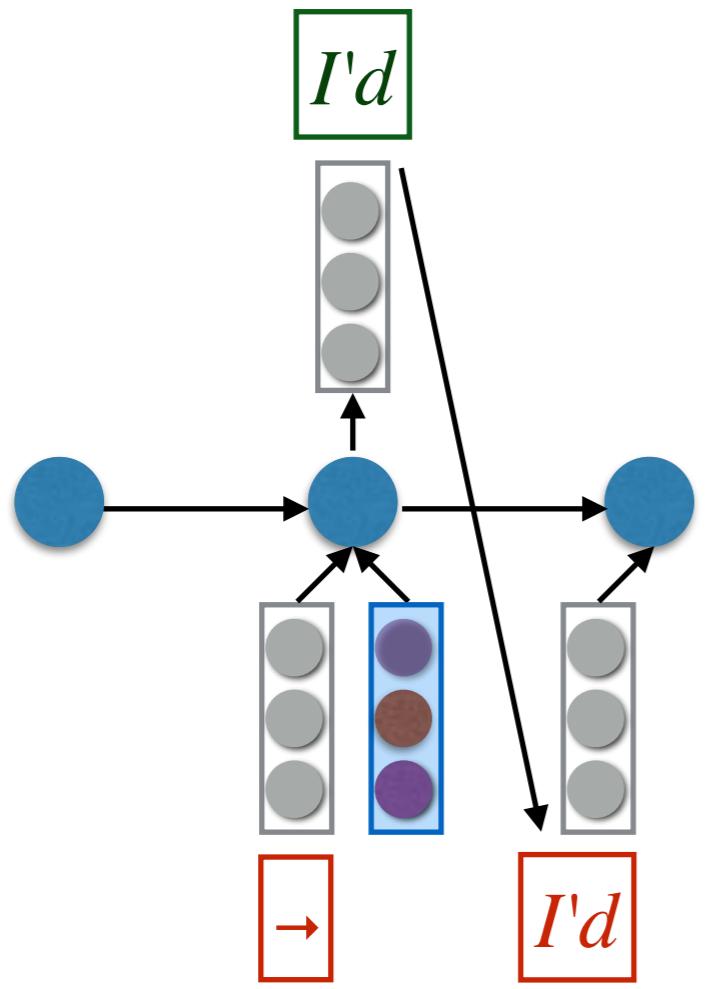
*Ich möchte ein Bier*



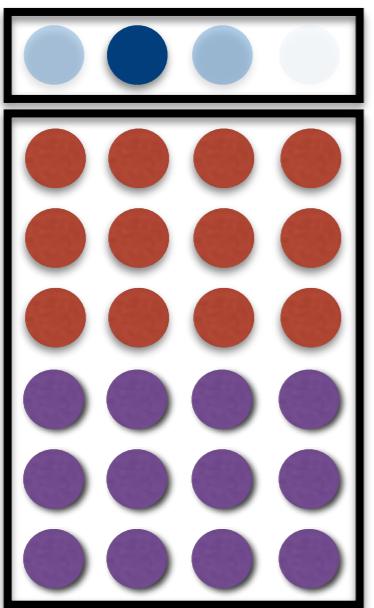
**Attention history:**

$$a_1^\top [ \text{blue dot} \quad \text{blue dot} \quad \text{light blue dot} \quad \text{light blue dot} ]$$

*Ich möchte ein Bier*



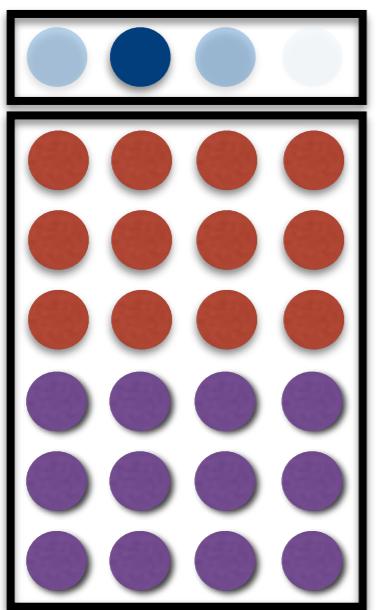
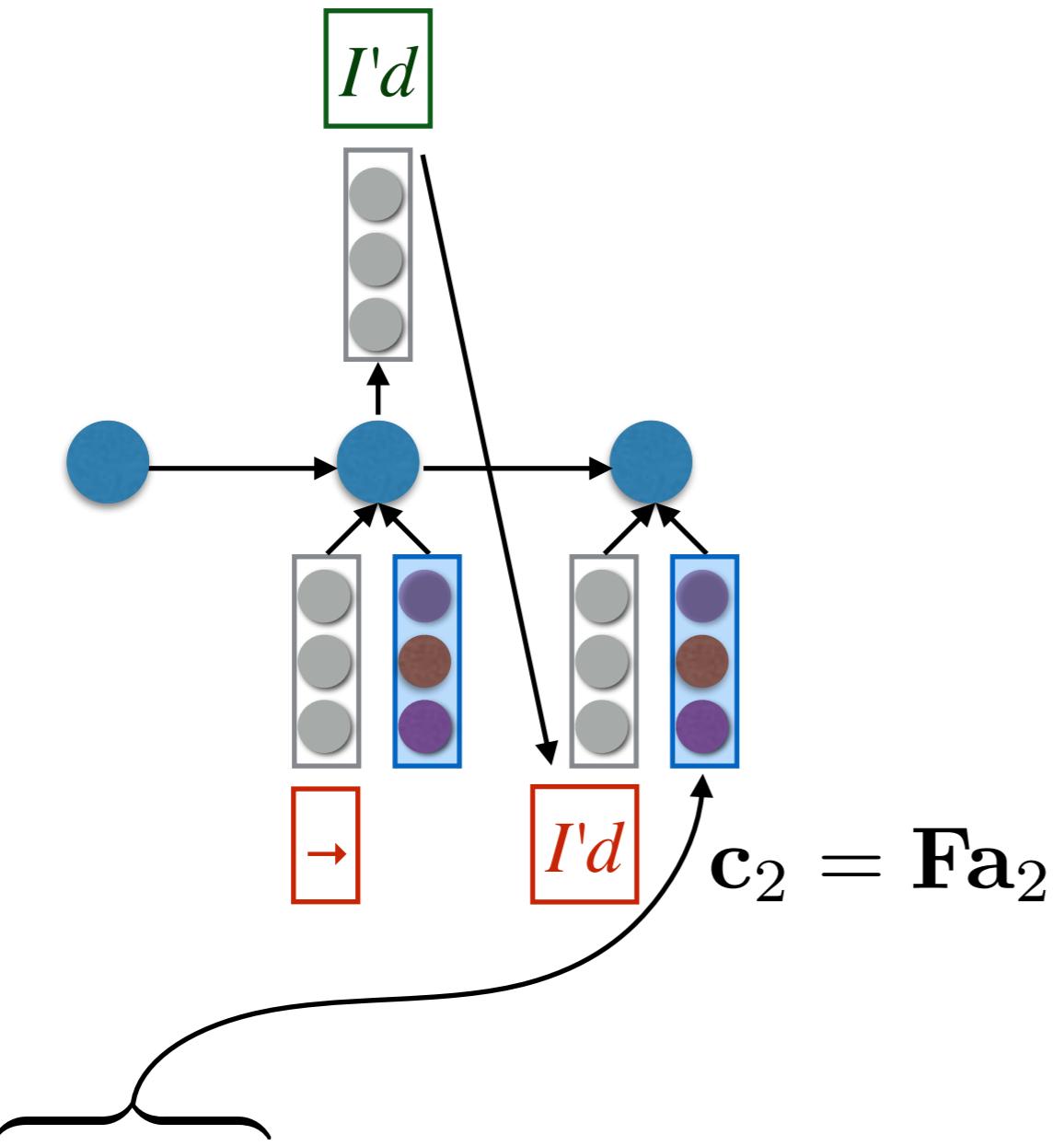
## Attention history:



$$a_1^\top \begin{bmatrix} \text{dark blue} \\ \text{light blue} \\ \text{white} \\ \text{white} \end{bmatrix}$$

$$a_2^\top \begin{bmatrix} \text{light blue} \\ \text{dark blue} \\ \text{light blue} \\ \text{white} \end{bmatrix}$$

*Ich möchte ein Bier*

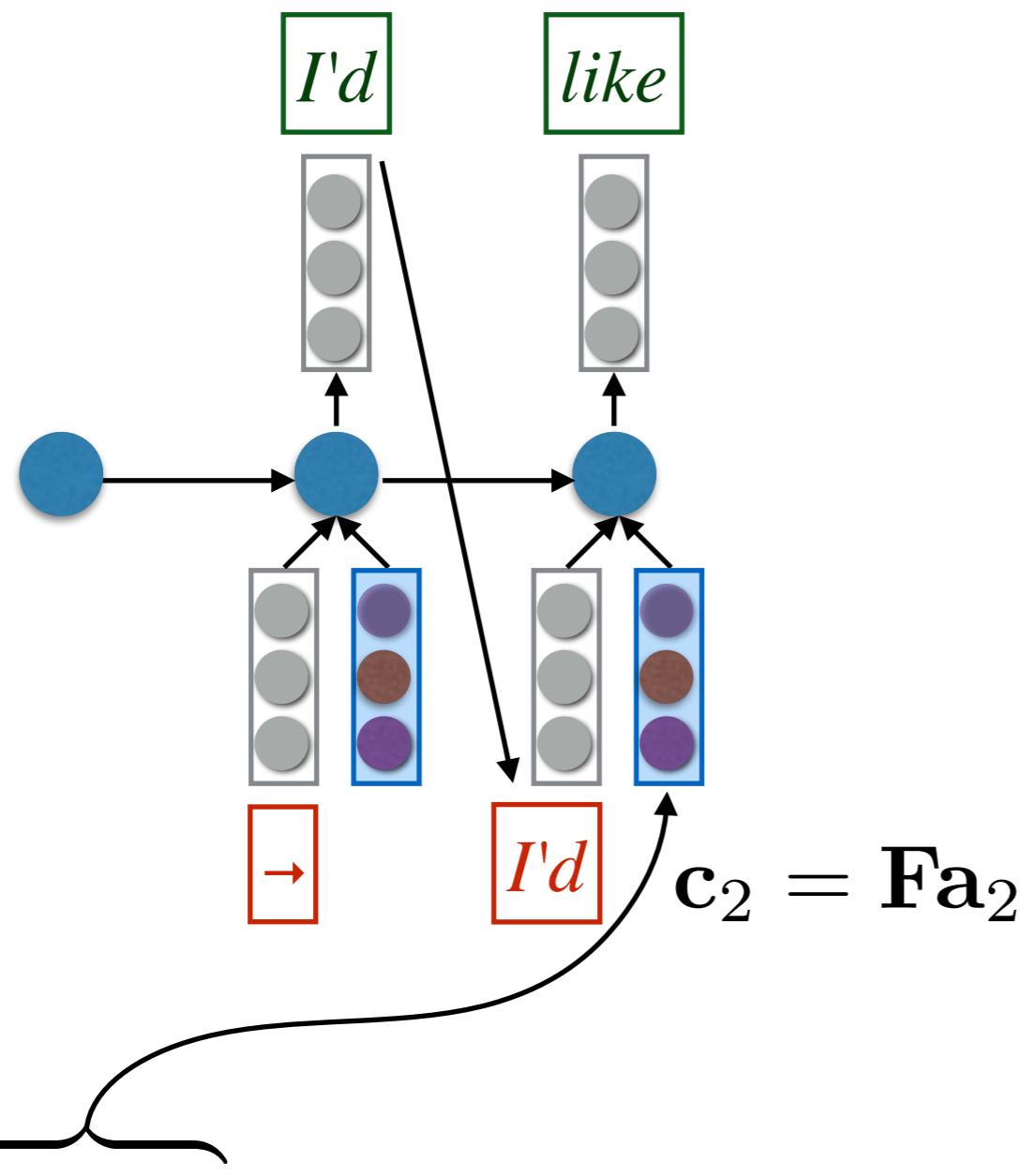


**Attention history:**

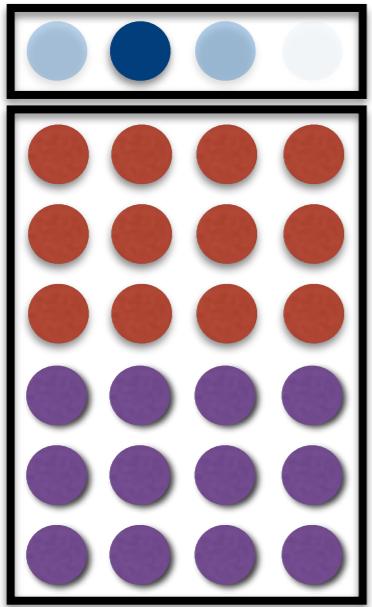
$$a_1^\top \begin{bmatrix} \text{dark blue} \\ \text{light blue} \\ \text{light blue} \\ \text{white} \end{bmatrix}$$

$$a_2^\top \begin{bmatrix} \text{light blue} \\ \text{dark blue} \\ \text{light blue} \\ \text{white} \end{bmatrix}$$

*Ich möchte ein Bier*



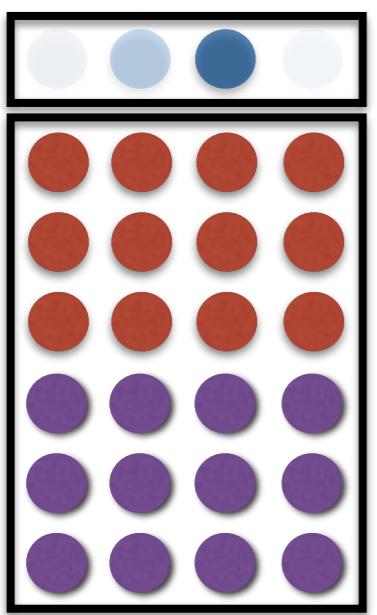
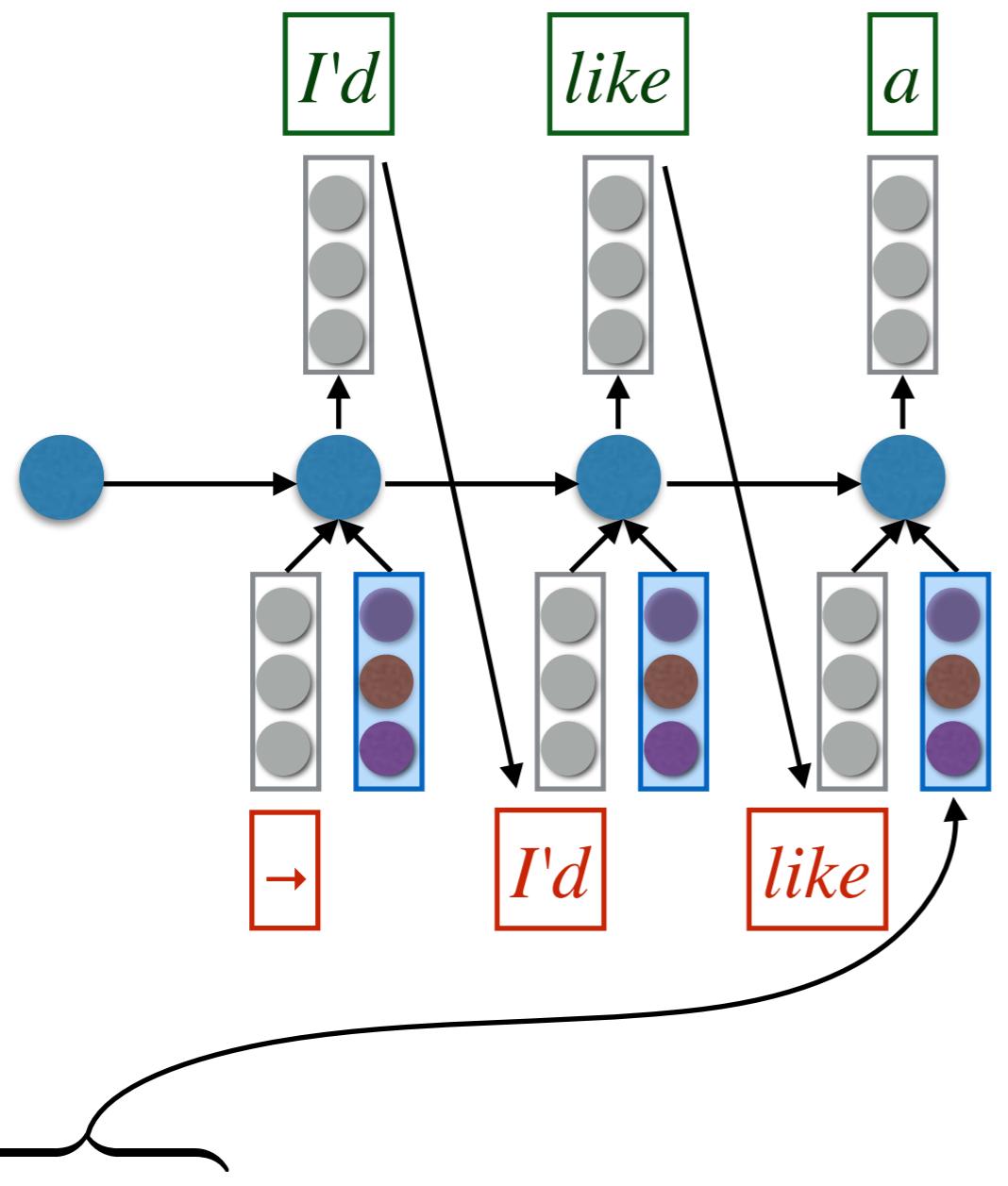
**Attention history:**



$$a_1^\top \begin{bmatrix} \text{Light Blue} \\ \text{Dark Blue} \\ \text{Light Blue} \\ \text{White} \end{bmatrix}$$

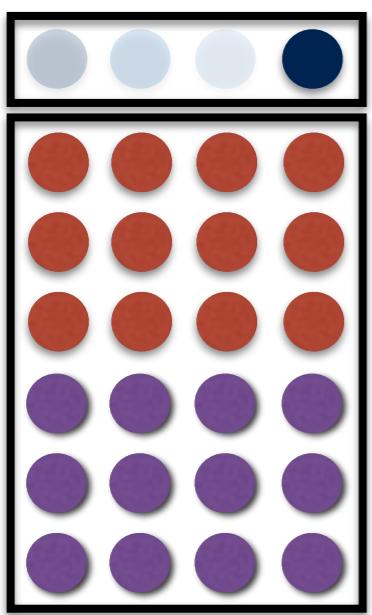
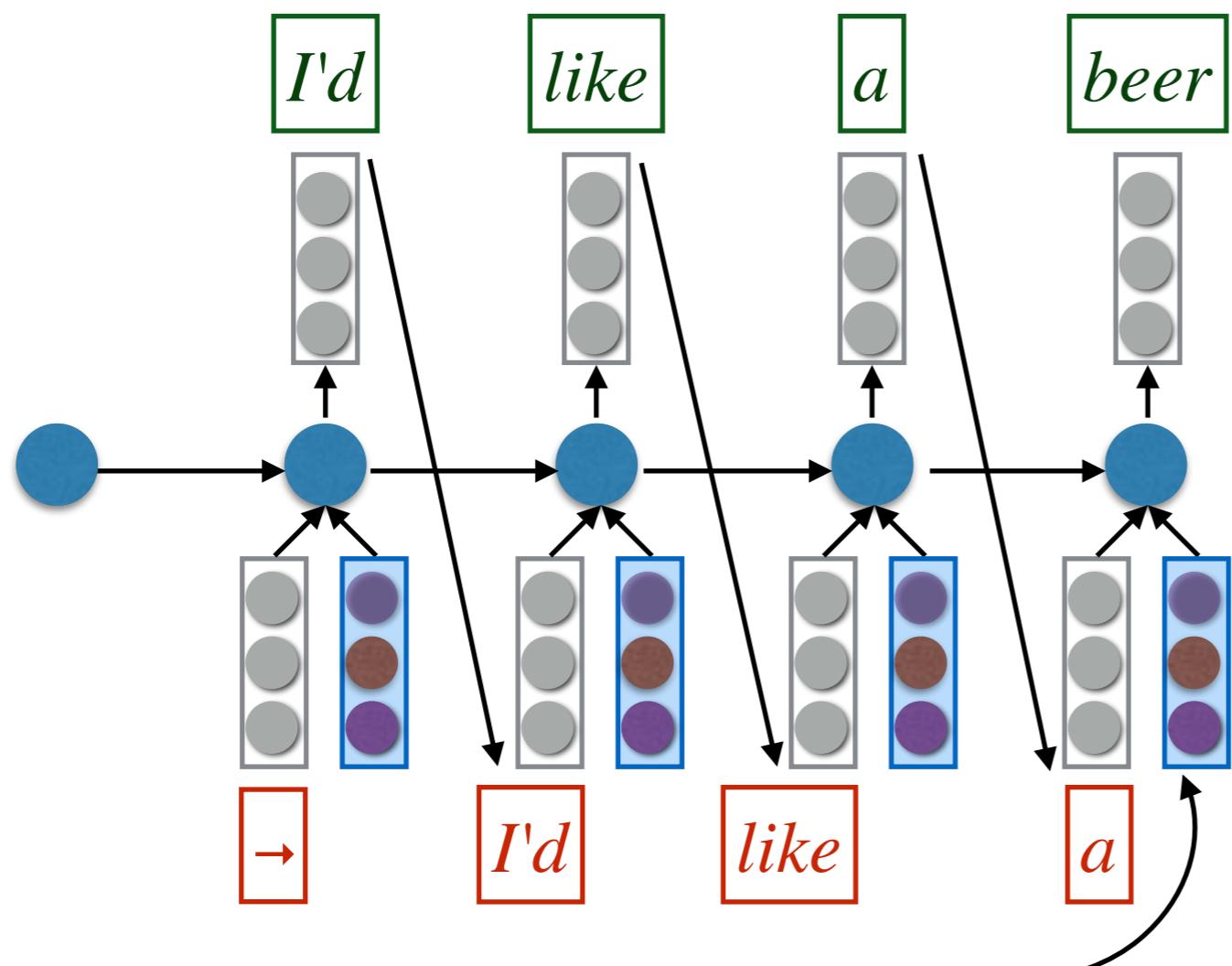
$$a_2^\top \begin{bmatrix} \text{Light Blue} \\ \text{Dark Blue} \\ \text{Light Blue} \\ \text{White} \end{bmatrix}$$

*Ich möchte ein Bier*



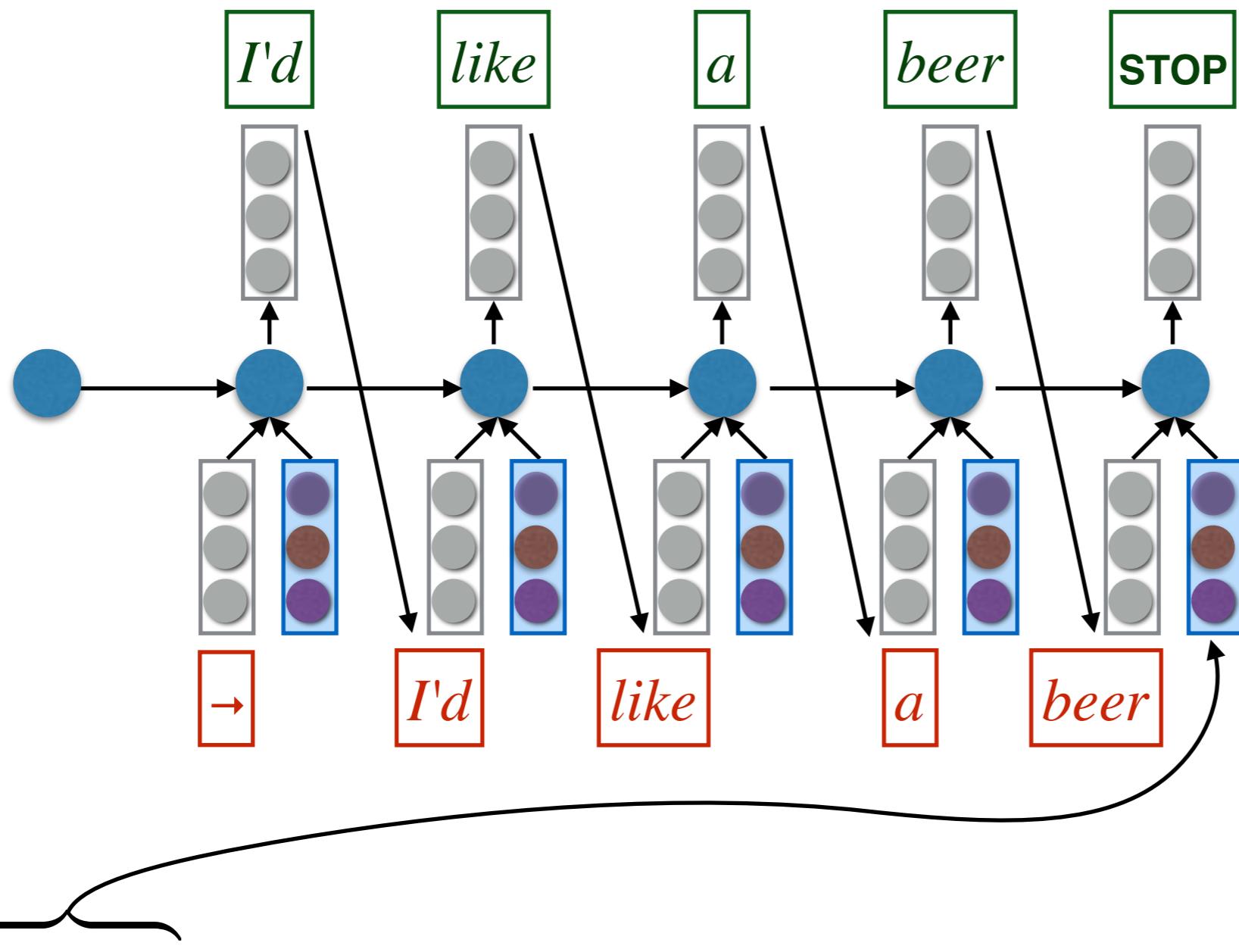
## Attention history:

*Ich möchte ein Bier*

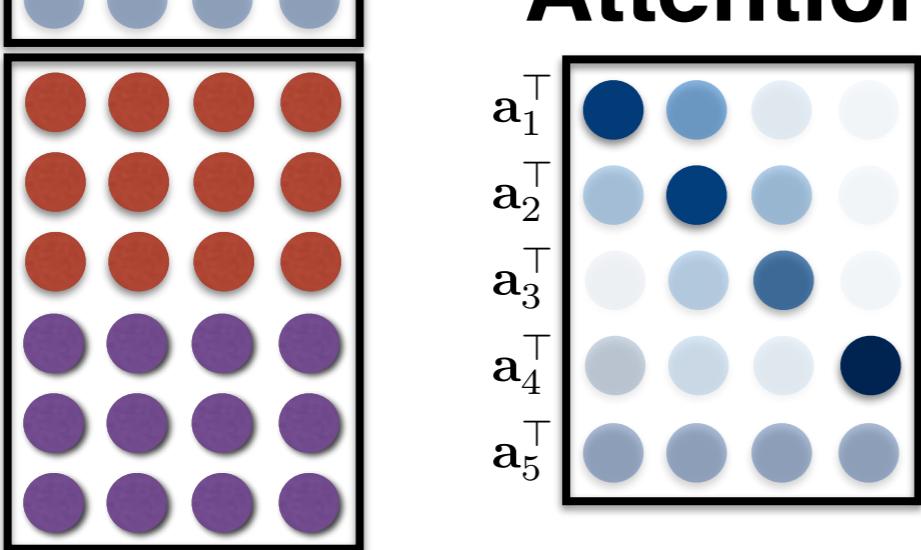


## Attention history:

*Ich möchte ein Bier*



## Attention history:



*Ich möchte ein Bier*

# Attention

- How do we know what to attend to at each time-step?
- That is, how do we compute  $\mathbf{a}_t$ ?

# Computing Attention

每个time step，都会有个 $|f|$ 长的向量，决定每个词的权重

- At each time step (one time step = one output word), we want to be able to “attend” to different words in the source sentence
  - We need a weight for every column: this is an  $|f|$ -length vector  $\mathbf{a}_t$
  - Here is a simplified version of Bahdanau et al.’s solution
    - Use an RNN to predict model output, call the hidden states  $\mathbf{s}_t$  ( $\mathbf{s}_t$  has a fixed dimensionality, call it  $m$ )

# Computing Attention

- At each time step (one time step = one output word), we want to be able to “attend” to different words in the source sentence
  - We need a weight for every column: this is an  $|\mathbf{f}|$ -length vector  $\mathbf{a}_t$
  - Here is a simplified version of Bahdanau et al.’s solution
    - Use an RNN to predict model output, call the hidden states  $\mathbf{s}_t$   
( $\mathbf{s}_t$  has a fixed dimensionality, call it  $m$ )
    - At time  $t$  compute the ***expected input embedding***  $\mathbf{r}_t = \mathbf{V}\mathbf{s}_{t-1}$   
( $\mathbf{V}$  is a learned parameter)

# Computing Attention

- At each time step (one time step = one output word), we want to be able to “attend” to different words in the source sentence
  - We need a weight for every column: this is an  $|\mathbf{f}|$ -length vector  $\mathbf{a}_t$
  - Here is a simplified version of Bahdanau et al.’s solution
    - Use an RNN to predict model output, call the hidden states  $\mathbf{s}_t$  ( $\mathbf{s}_t$  has a fixed dimensionality, call it  $m$ )
    - At time  $t$  compute the ***expected input embedding***  $\mathbf{r}_t = \mathbf{V}\mathbf{s}_{t-1}$  ( $\mathbf{V}$  is a learned parameter)
    - Take the dot product with every column in the source matrix to compute the ***attention energy***.  $\mathbf{u}_t = \mathbf{F}^\top \mathbf{r}_t$  (called  $\mathbf{e}_t$  in the paper)  
(Since  $\mathbf{F}$  has  $|\mathbf{f}|$  columns,  $\mathbf{u}_t$  has  $|\mathbf{f}|$  rows)

# Computing Attention

- At each time step (one time step = one output word), we want to be able to “attend” to different words in the source sentence
  - We need a weight for every column: this is an  $|\mathbf{f}|$ -length vector  $\mathbf{a}_t$
  - Here is a simplified version of Bahdanau et al.’s solution
    - Use an RNN to predict model output, call the hidden states  $\mathbf{s}_t$  ( $\mathbf{s}_t$  has a fixed dimensionality, call it  $m$ )
    - At time  $t$  compute the **expected input embedding**  $\mathbf{r}_t = \mathbf{V}\mathbf{s}_{t-1}$  ( $\mathbf{V}$  is a learned parameter)
    - Take the dot product with every column in the source matrix to compute the **attention energy**.  $\mathbf{u}_t = \mathbf{F}^\top \mathbf{r}_t$  (called  $\mathbf{e}_t$  in the paper)  
(Since  $\mathbf{F}$  has  $|\mathbf{f}|$  columns,  $\mathbf{u}_t$  has  $|\mathbf{f}|$  rows)
    - Exponentiate and normalize to 1:  $\mathbf{a}_t = \text{softmax}(\mathbf{u}_t)$   
(called  $\alpha_t$  in the paper)

# Computing Attention

- At each time step (one time step = one output word), we want to be able to “attend” to different words in the source sentence
  - We need a weight for every column: this is an  $|\mathbf{f}|$ -length vector  $\mathbf{a}_t$
  - Here is a simplified version of Bahdanau et al.’s solution
    - Use an RNN to predict model output, call the hidden states  $\mathbf{s}_t$  ( $\mathbf{s}_t$  has a fixed dimensionality, call it  $m$ )
    - At time  $t$  compute the **expected input embedding**  $\mathbf{r}_t = \mathbf{V}\mathbf{s}_{t-1}$  ( $\mathbf{V}$  is a learned parameter)
    - Take the dot product with every column in the source matrix to compute the **attention energy**.  $\mathbf{u}_t = \mathbf{F}^\top \mathbf{r}_t$  (called  $\mathbf{e}_t$  in the paper)  
(Since  $\mathbf{F}$  has  $|\mathbf{f}|$  columns,  $\mathbf{u}_t$  has  $|\mathbf{f}|$  rows)
    - Exponentiate and normalize to 1:  $\mathbf{a}_t = \text{softmax}(\mathbf{u}_t)$   
(called  $\alpha_t$  in the paper)
    - Finally, the **input source vector** for time  $t$  is  $\mathbf{c}_t = \mathbf{F}\mathbf{a}_t$

# Nonlinear Attention-Energy Model

- In the actual model, Bahdanau et al. replace the dot product between the columns of  $\mathbf{F}$  and  $\mathbf{r}_t$  with an MLP:

$$\mathbf{u}_t = \mathbf{F}^\top \mathbf{r}_t \quad (\text{simple model})$$

# Nonlinear Attention-Energy Model

- In the actual model, Bahdanau et al. replace the dot product between the columns of  $\mathbf{F}$  and  $\mathbf{r}_t$  with an MLP:

$$\mathbf{u}_t = \mathbf{F}^\top \mathbf{r}_t \quad (\text{simple model})$$

$$\mathbf{u}_t = \mathbf{v}^\top \tanh(\mathbf{W}\mathbf{F} + \mathbf{r}_t) \quad (\text{Bahdanau et al})$$

# Nonlinear Attention-Energy Model

- In the actual model, Bahdanau et al. replace the dot product between the columns of  $\mathbf{F}$  and  $\mathbf{r}_t$  with an MLP:

$$\mathbf{u}_t = \mathbf{F}^\top \mathbf{r}_t \quad (\text{simple model})$$

$$\mathbf{u}_t = \mathbf{v}^\top \tanh(\mathbf{W}\mathbf{F} + \mathbf{r}_t) \quad (\text{Bahdanau et al})$$

- Here,  $\mathbf{W}$  and  $\mathbf{v}$  are learned parameters of appropriate dimension and + “broadcasts” over the  $|\mathbf{f}|$  columns in  $\mathbf{WF}$
- This can learn more complex interactions
  - It is unclear if the added complexity is necessary for good performance

# Putting it all together

$\mathbf{F} = \text{EncodeAsMatrix}(f)$  (Part 1 of lecture)

$e_0 = \langle s \rangle$

$s_0 = w$  (Learned initial state; Bahdanau uses  $U^{\leftarrow} h_1$ )

$t = 0$

**while**  $e_t \neq \langle /s \rangle$  :

$t = t + 1$

$\mathbf{r}_t = \mathbf{V}s_{t-1}$

$\mathbf{u}_t = \mathbf{v}^\top \tanh(\mathbf{W}\mathbf{F} + \mathbf{r}_t)$

$\mathbf{a}_t = \text{softmax}(\mathbf{u}_t)$

$\mathbf{c}_t = \mathbf{F}\mathbf{a}_t$

$\mathbf{s}_t = \text{RNN}(\mathbf{s}_{t-1}, [\mathbf{e}_{t-1}; \mathbf{c}_t])$  ( $\mathbf{e}_{t-1}$  is a learned embedding of  $e_t$ )

$\mathbf{y}_t = \text{softmax}(\mathbf{P}\mathbf{s}_t + \mathbf{b})$  ( $\mathbf{P}$  and  $\mathbf{b}$  are learned parameters)

$e_t \mid e_{<t} \sim \text{Categorical}(\mathbf{y}_t)$

}

(Compute attention; part 2 of lecture)

# Putting it all together

$\mathbf{F} = \text{EncodeAsMatrix}(f)$  (Part 1 of lecture)

$e_0 = \langle s \rangle$

$s_0 = w$  (Learned initial state; Bahdanau uses  $U^{\leftarrow} h_1$ )

$t = 0$

**while**  $e_t \neq \langle /s \rangle$  :

$t = t + 1$

$r_t = \mathbf{V}s_{t-1}$

$\mathbf{u}_t = \mathbf{v}^\top \tanh(\mathbf{W}\mathbf{F} + r_t)$

$\mathbf{a}_t = \text{softmax}(\mathbf{u}_t)$

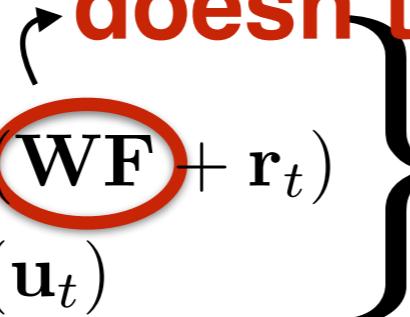
$\mathbf{c}_t = \mathbf{F}\mathbf{a}_t$

$\mathbf{s}_t = \text{RNN}(\mathbf{s}_{t-1}, [\mathbf{e}_{t-1}; \mathbf{c}_t])$  ( $\mathbf{e}_{t-1}$  is a learned embedding of  $e_t$ )

$\mathbf{y}_t = \text{softmax}(\mathbf{P}\mathbf{s}_t + \mathbf{b})$  ( $\mathbf{P}$  and  $\mathbf{b}$  are learned parameters)

$e_t | e_{<t} \sim \text{Categorical}(\mathbf{y}_t)$

**doesn't depend on output decisions**



# Putting it all together

$\mathbf{F} = \text{EncodeAsMatrix}(f)$  (Part 1 of lecture)

$e_0 = \langle s \rangle$

$s_0 = w$  (Learned initial state; Bahdanau uses  $U^{\leftarrow} h_1$ )

$t = 0$

$X = WF$

**while**  $e_t \neq \langle /s \rangle$  :

$t = t + 1$

$r_t = Vs_{t-1}$

$u_t = v^\top \tanh(WF + r_t)$

$a_t = \text{softmax}(u_t)$

$c_t = Fa_t$

$s_t = \text{RNN}(s_{t-1}, [e_{t-1}; c_t])$

$y_t = \text{softmax}(Ps_t + b)$

$e_t | e_{<t} \sim \text{Categorical}(y_t)$

} (Compute attention; part 2 of lecture)

( $e_{t-1}$  is a learned embedding of  $e_t$ )

( $P$  and  $b$  are learned parameters)

# Putting it all together

$\mathbf{F} = \text{EncodeAsMatrix}(f)$  (Part 1 of lecture)

$e_0 = \langle s \rangle$

$s_0 = w$  (Learned initial state; Bahdanau uses  $U^{\leftarrow} h_1$ )

$t = 0$

$\mathbf{X} = \mathbf{WF}$

**while**  $e_t \neq \langle /s \rangle$  :

$t = t + 1$

$\mathbf{r}_t = \mathbf{Vs}_{t-1}$

$\mathbf{u}_t = \mathbf{v}^\top \tanh(\mathbf{X} + \mathbf{r}_t)$

$\mathbf{a}_t = \text{softmax}(\mathbf{u}_t)$

$\mathbf{c}_t = \mathbf{Fa}_t$

$\mathbf{s}_t = \text{RNN}(\mathbf{s}_{t-1}, [\mathbf{e}_{t-1}; \mathbf{c}_t])$

$\mathbf{y}_t = \text{softmax}(\mathbf{Ps}_t + \mathbf{b})$

$e_t | e_{<t} \sim \text{Categorical}(\mathbf{y}_t)$

}

(Compute attention; part 2 of lecture)

( $\mathbf{e}_{t-1}$  is a learned embedding of  $e_t$ )

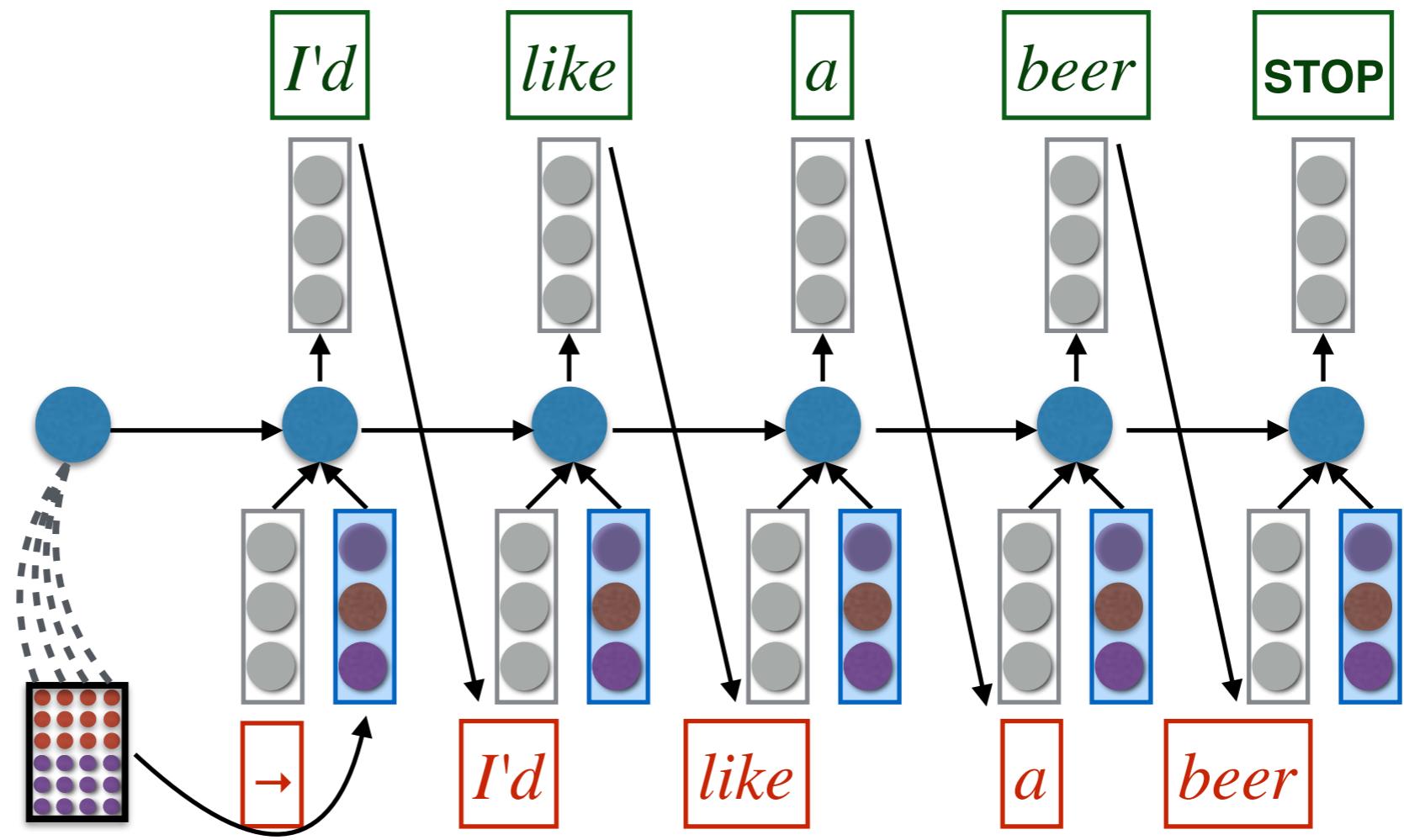
( $\mathbf{P}$  and  $\mathbf{b}$  are learned parameters)

# Attention in MT

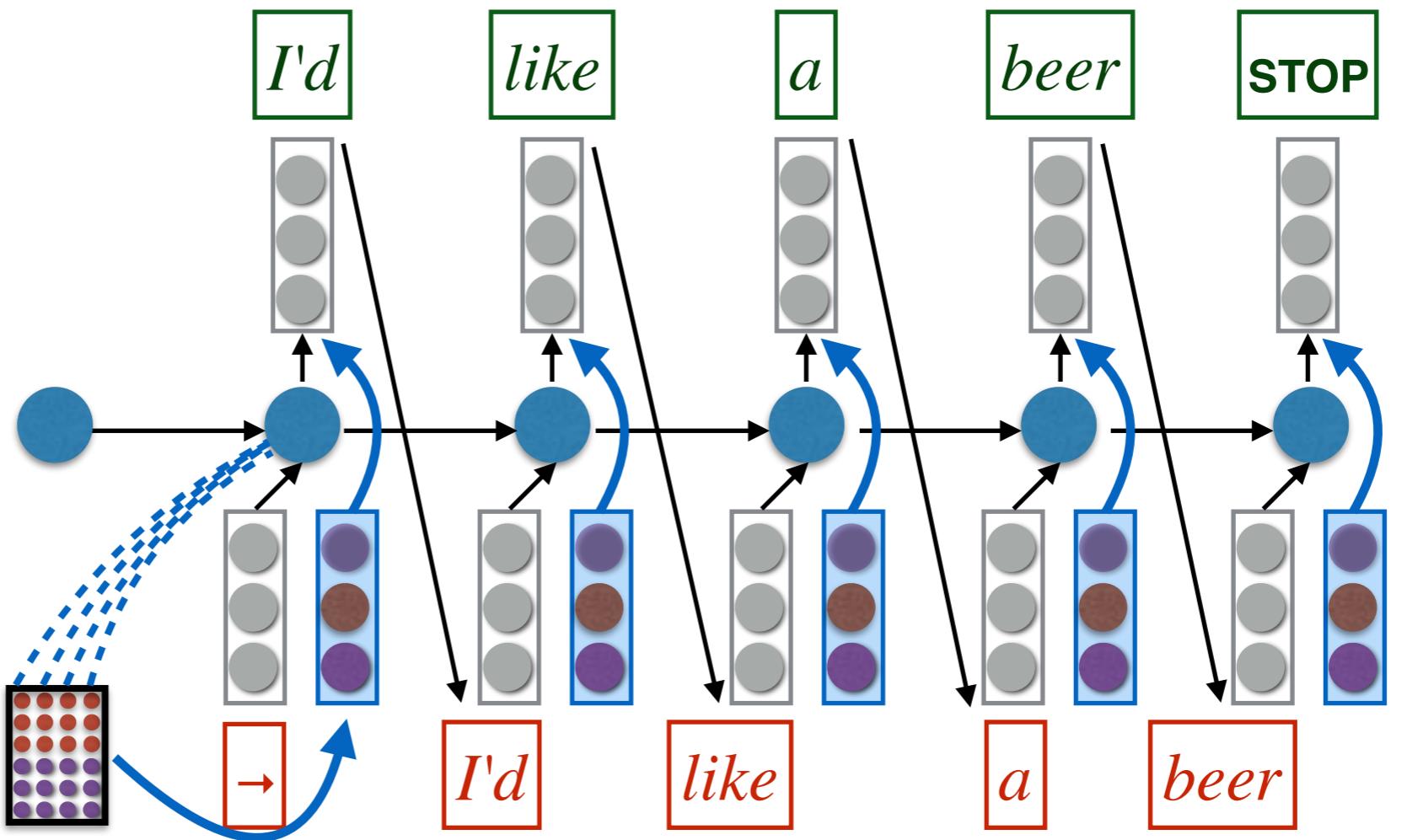
Add attention to seq2seq translation: **+11 BLEU**

# Model Variant

“Early binding”



“Late binding”



# Model Variant

- What are the relative advantages of early binding versus late binding?

# Summary

- Attention is closely related to “pooling” operations in convnets (and other architectures)
- Bahdanau’s attention model seems to only cares about “content”
  - No obvious bias in favor of diagonals, short jumps, fertility, etc.
  - Some work has begun to add other “structural” biases (Luong et al., 2015; Cohn et al., 2016), but there are lots more opportunities
- Attention weights provide interpretation you can look at

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

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

<end>

(a)

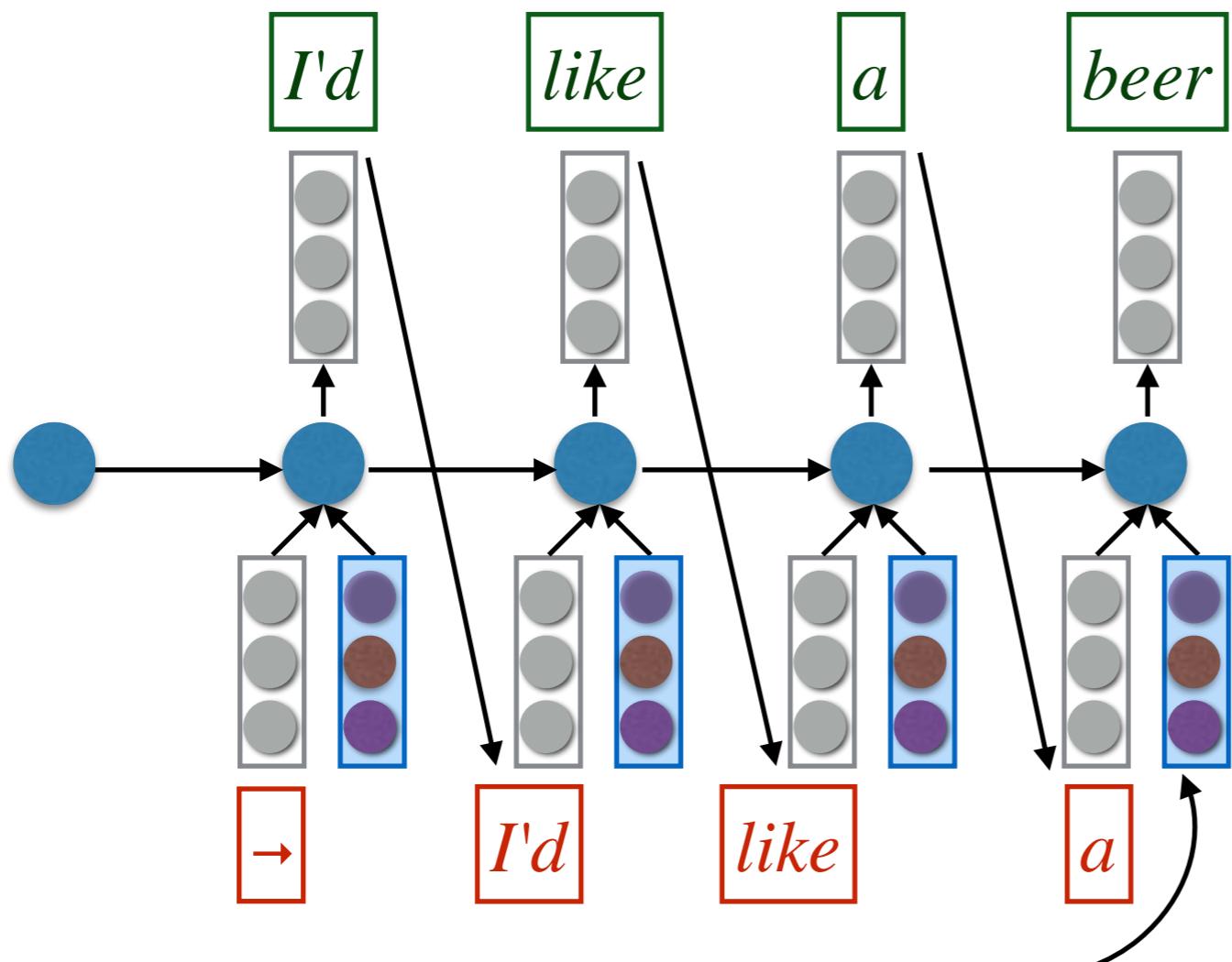
Il convient de noter que l'environnement marin est le moins connu de l'environnement.

It should be noted that the marine environment is the least known of environments.

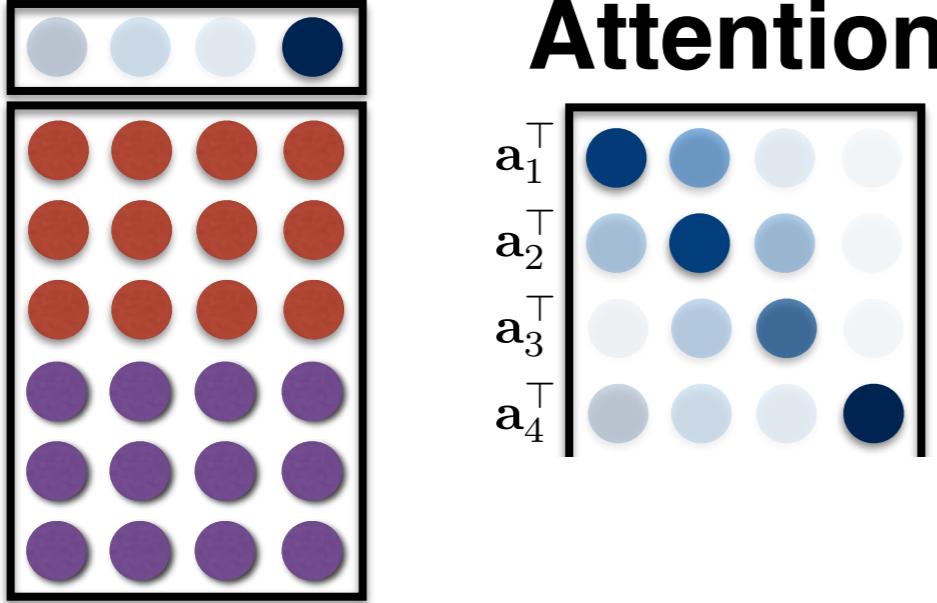
<end>

(b)

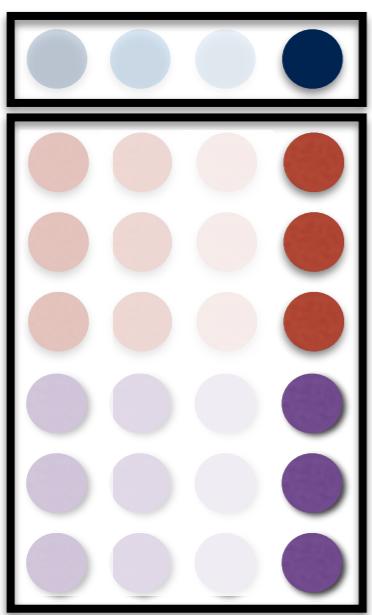
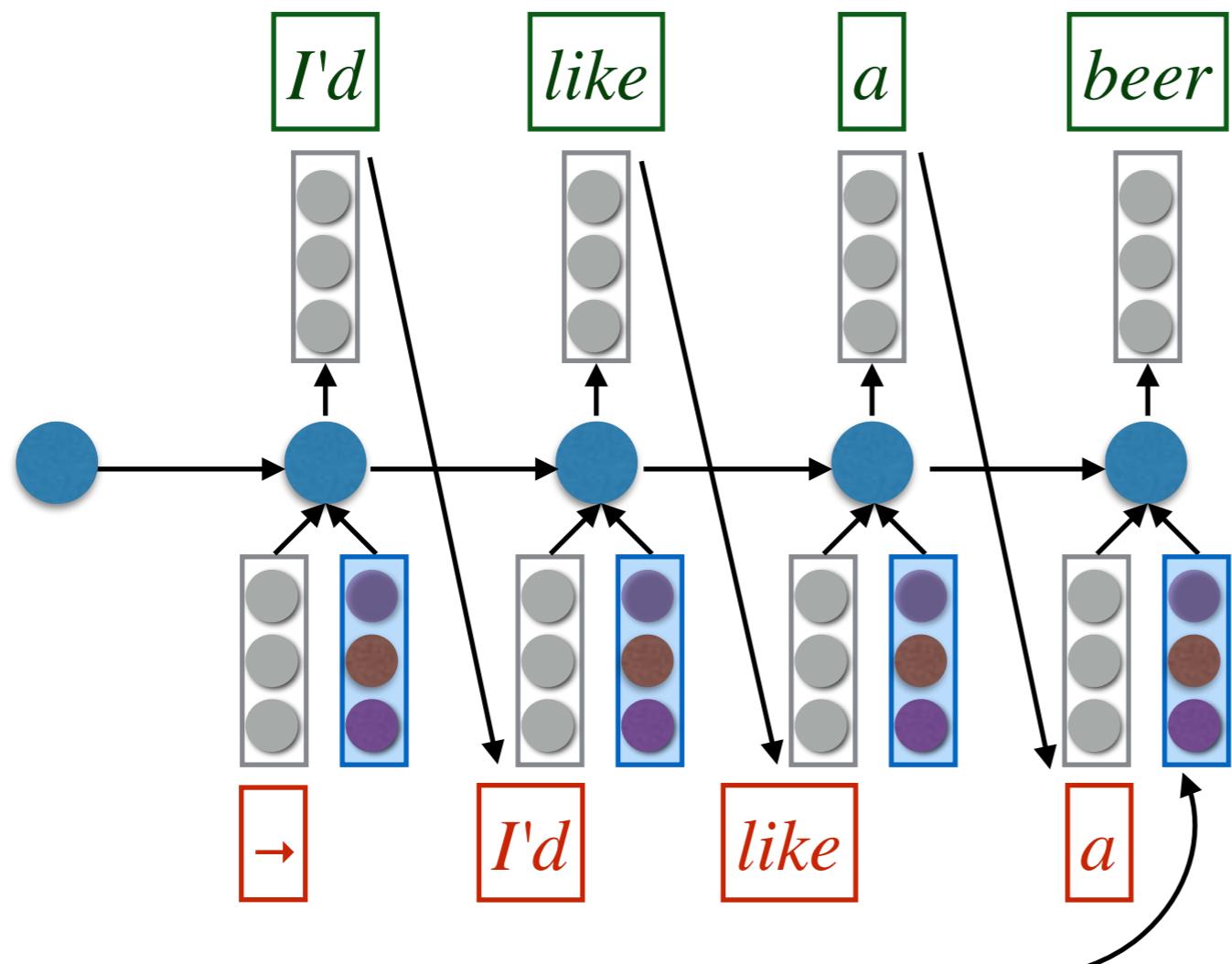
A word about gradients



**Attention history:**



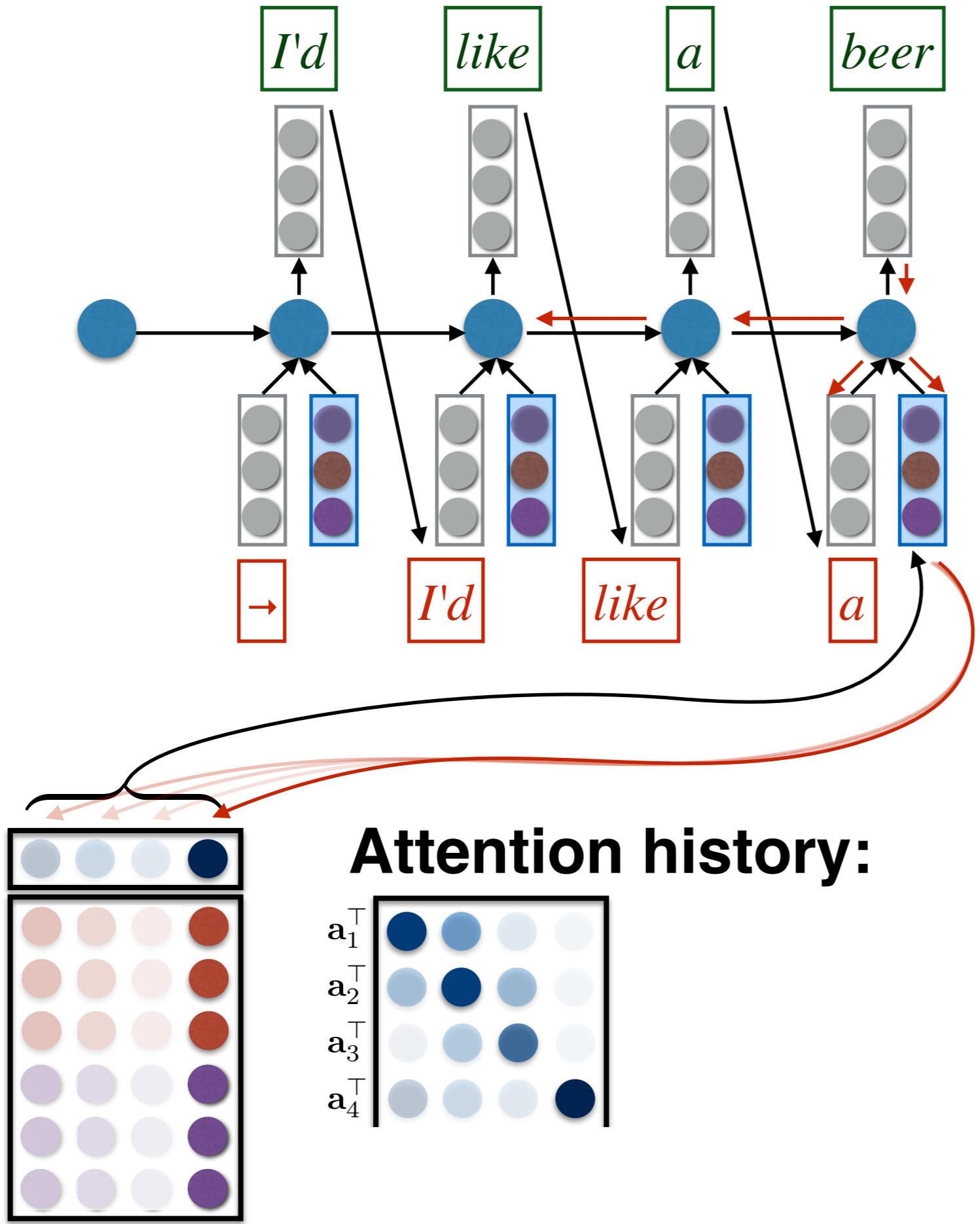
*Ich möchte ein Bier*



**Attention history:**

$$\begin{matrix}
 & a_1^T & a_2^T & a_3^T & a_4^T \\
 a_1 & \text{dark blue} & \text{light blue} & \text{light blue} & \text{light blue} \\
 a_2 & \text{light blue} & \text{dark blue} & \text{light blue} & \text{light blue} \\
 a_3 & \text{light blue} & \text{light blue} & \text{dark blue} & \text{light blue} \\
 a_4 & \text{light blue} & \text{light blue} & \text{light blue} & \text{dark blue}
 \end{matrix}$$

*Ich möchte ein Bier*



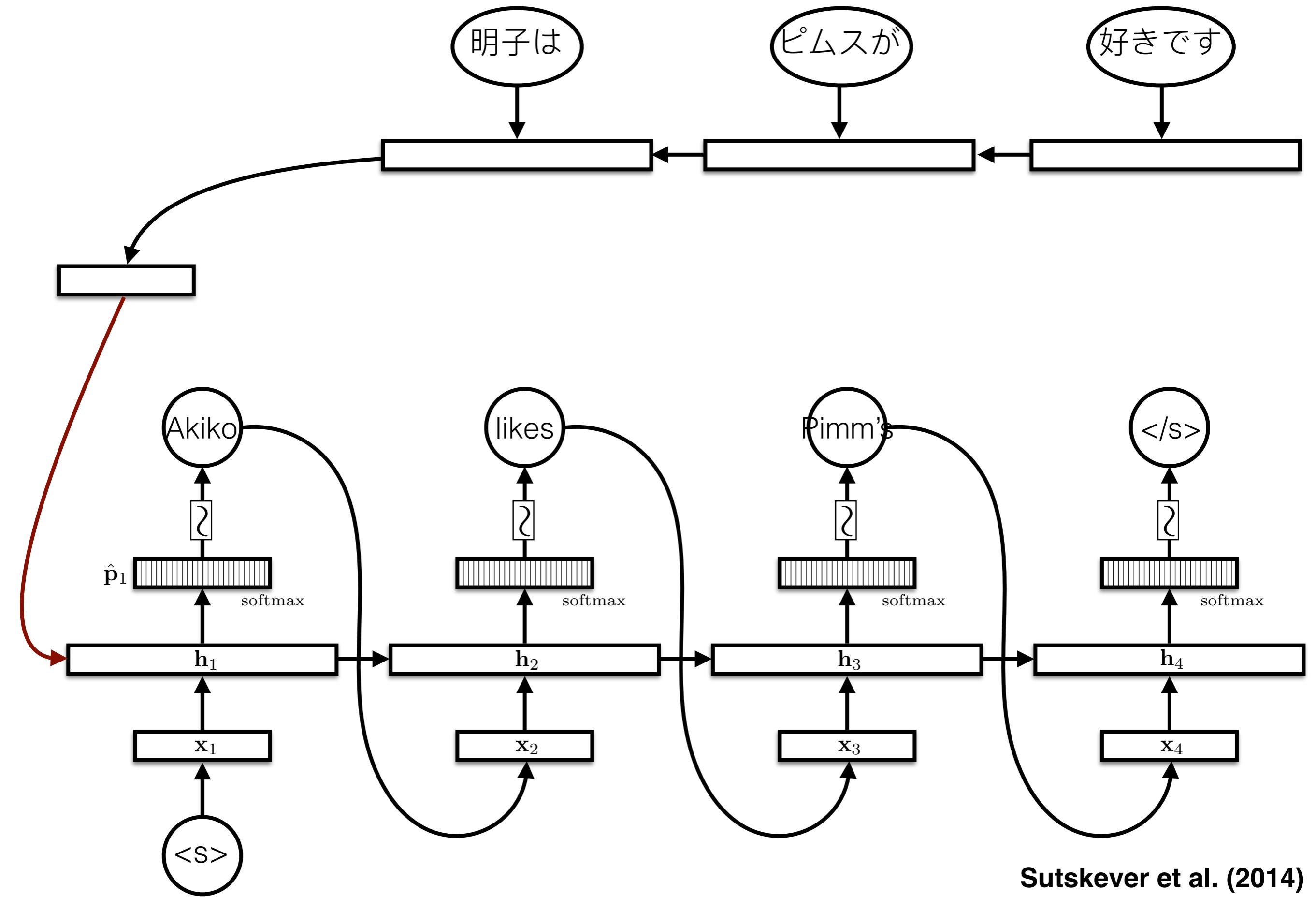
*Ich möchte ein Bier*

# Attention and Translation

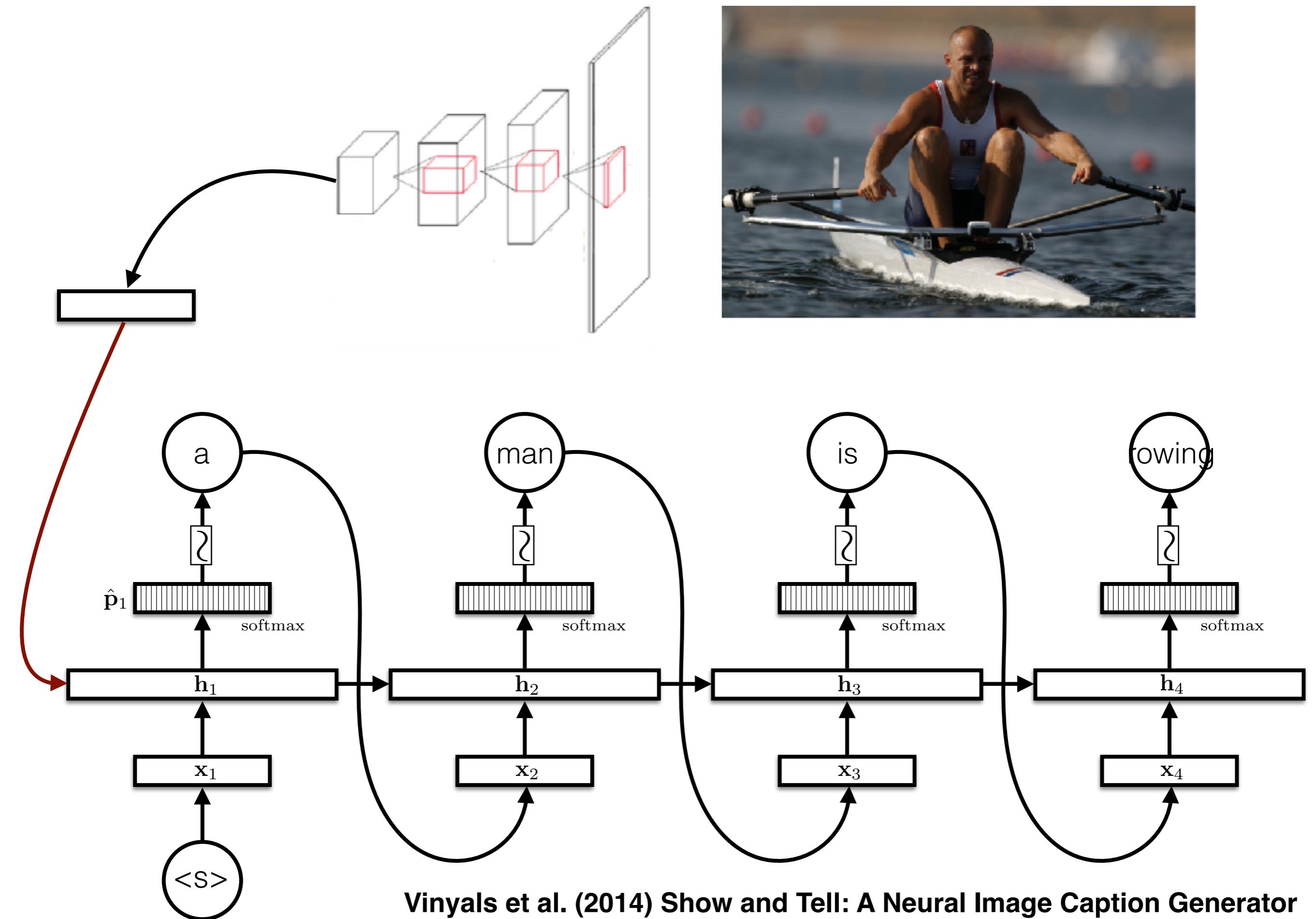
- Cho's question: does a translator read and memorize the input sentence/document and then generate the output?
  - Compressing the entire input sentence into a vector basically says “memorize the sentence”
  - Common sense experience says translators refer back and forth to the input. (also backed up by eye-tracking studies)
- Should humans be a model for machines?

# Outline of Lecture

- Machine translation with attention
- Image caption generation with attention



Sutskever et al. (2014)



# Image Caption Generation

- Can attention help caption modeling?

---

## Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

---

**Kelvin Xu**

**Jimmy Lei Ba**

**Ryan Kiros**

**Kyunghyun Cho**

**Aaron Courville**

**Ruslan Salakhutdinov**

**Richard S. Zemel**

**Yoshua Bengio**

KELVIN.XU@UMONTREAL.CA

JIMMY@PSI.UTORONTO.CA

RKIROS@CS.TORONTO.EDU

KYUNGHYUN.CHO@UMONTREAL.CA

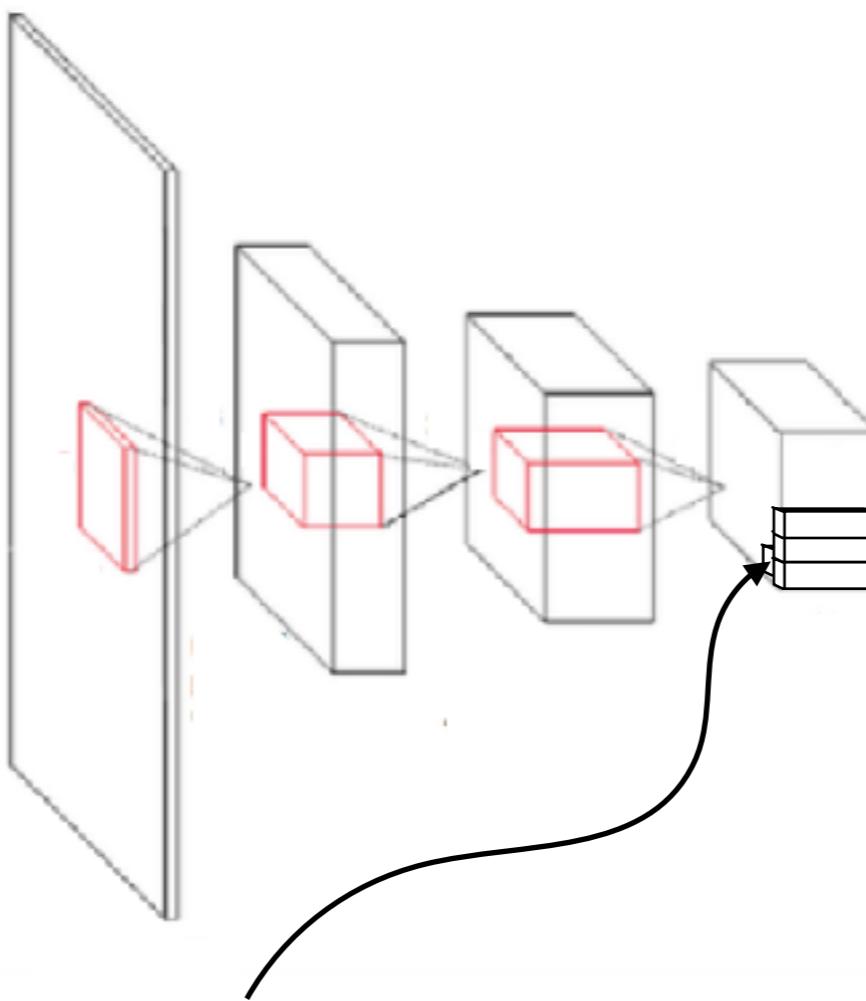
AARON.COURVILLE@UMONTREAL.CA

RSALAKHU@CS.TORONTO.EDU

ZEMEL@CS.TORONTO.EDU

FIND-ME@THE.WEB

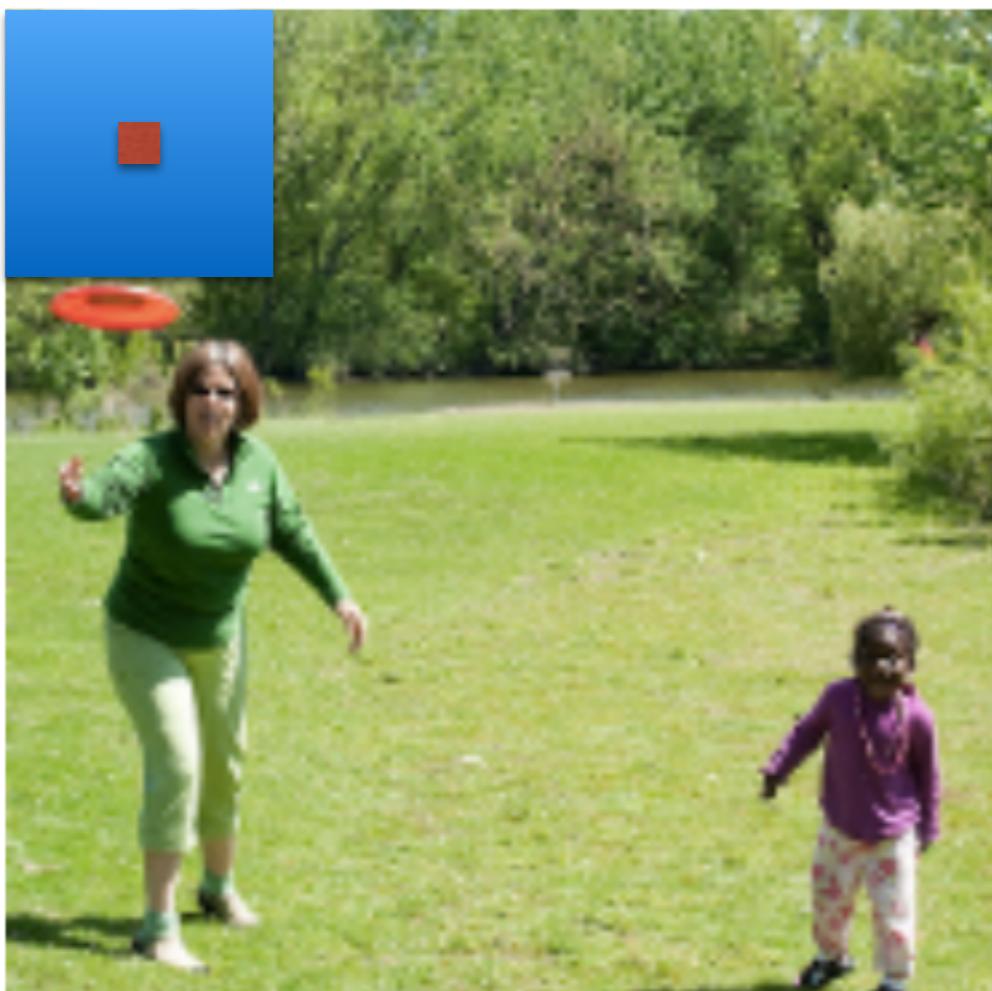
# Regions in ConvNets



Each point in a “higher” level of a convnet defines spatially localised feature vectors(/matrices).

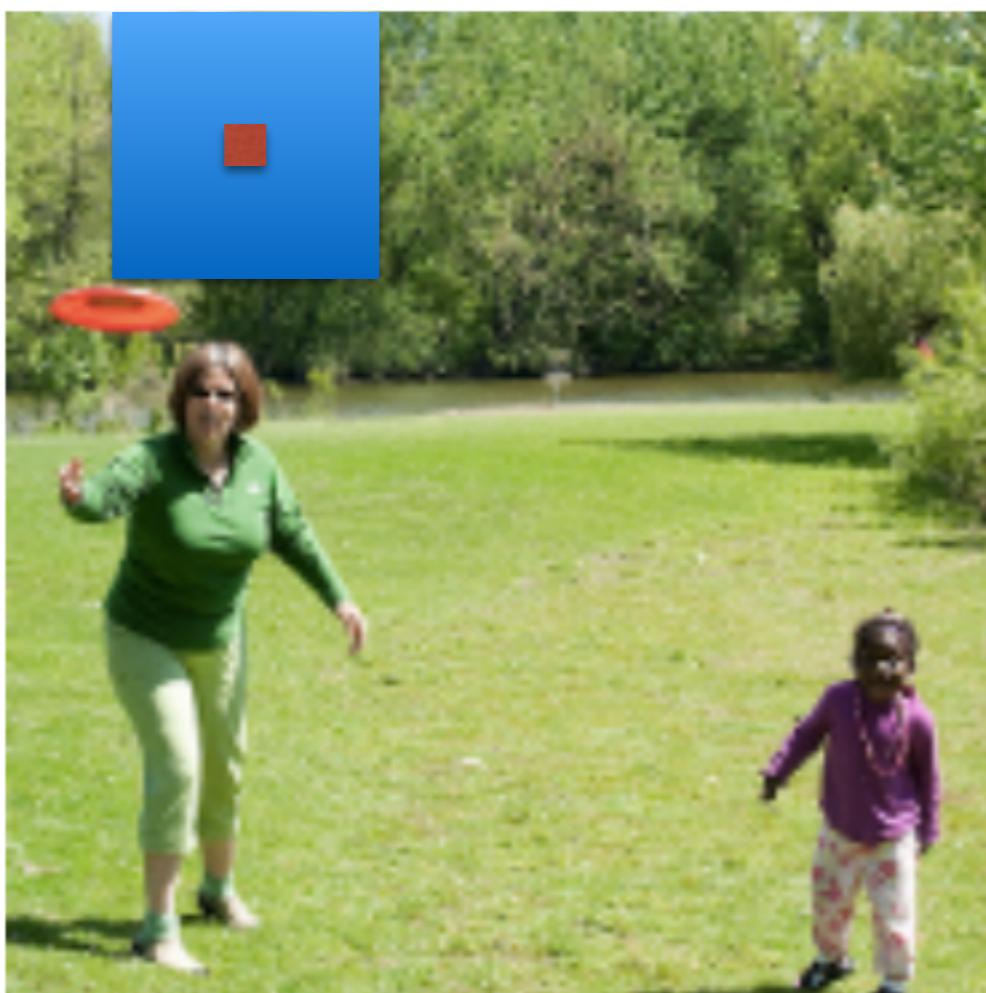
Xu et al. calls these “annotation vectors”,  $\mathbf{a}_i$ ,  $i \in \{1, \dots, L\}$

$\mathbf{a}_1$



$$\mathbf{F} = \begin{bmatrix} | \\ \mathbf{a}_1 \\ | \end{bmatrix}$$

$\mathbf{a}_2$



$$\mathbf{F} = \begin{bmatrix} | & | \\ \mathbf{a}_1 & \mathbf{a}_2 \\ | & | \end{bmatrix}$$

**a<sub>3</sub>**



$$\mathbf{F} = \begin{bmatrix} | & | & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 & \cdots \\ | & | & | \end{bmatrix}$$

# Attention

- Attention “weights” ( $a_t$ ) are computed using exactly the same technique as discussed above

# Attention

- Attention “weights” ( $\mathbf{a}_t$ ) are computed using exactly the same technique as discussed above
- Deterministic soft attention (Bahdanau et al., 2014)

$$\mathbf{c}_t = \mathbf{F}\mathbf{a}_t \quad (\text{weighted average})$$

# Attention

- Attention “weights” ( $\mathbf{a}_t$ ) are computed using exactly the same technique as discussed above
- Deterministic soft attention (Bahdanau et al., 2014)

$$\mathbf{c}_t = \mathbf{F}\mathbf{a}_t \quad (\text{weighted average})$$

- Stochastic hard attention (Xu et al., 2015)

$$s_t \sim \text{Categorical}(\mathbf{a}_t)$$

$$\mathbf{c}_t = \mathbf{F}_{:,s_t} \quad (\text{sample a column})$$

- What are the benefits of this model?
- What are the challenges of learning the parameters of this model?

# Learning Hard Attention

$$\begin{aligned}\mathcal{L} &= -\log p(\mathbf{w} \mid \mathbf{x}) \\ &= -\log \sum_{\mathbf{s}} p(\mathbf{w}, \mathbf{s} \mid \mathbf{x}) \\ &= -\log \sum_{\mathbf{s}} p(\mathbf{s} \mid \mathbf{x}) p(\mathbf{w} \mid \mathbf{x}, \mathbf{s})\end{aligned}$$

# Learning Hard Attention

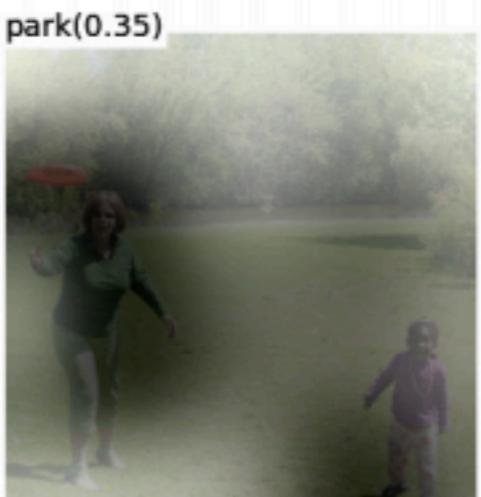
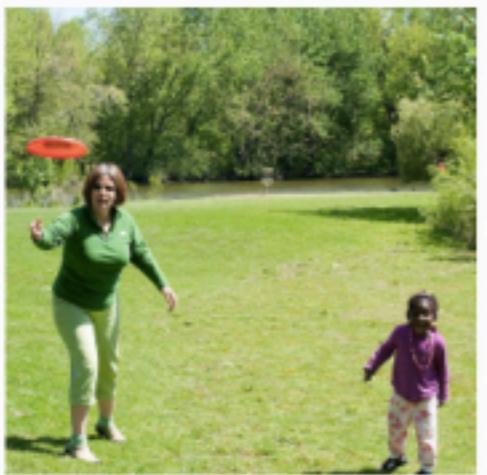
$$\begin{aligned}\mathcal{L} &= -\log p(\mathbf{w} \mid \mathbf{x}) \\ &= -\log \sum_{\mathbf{s}} p(\mathbf{w}, \mathbf{s} \mid \mathbf{x}) \\ &= -\log \sum_{\mathbf{s}} p(\mathbf{s} \mid \mathbf{x}) p(\mathbf{w} \mid \mathbf{x}, \mathbf{s}) \\ &\leq -\sum_{\mathbf{s}} p(\mathbf{s} \mid \mathbf{x}) \log p(\mathbf{w} \mid \mathbf{x}, \mathbf{s}) \quad (\text{Jensen's inequality})\end{aligned}$$

# Learning Hard Attention

$$\begin{aligned}\mathcal{L} &= -\log p(\mathbf{w} \mid \mathbf{x}) \\&= -\log \sum_{\mathbf{s}} p(\mathbf{w}, \mathbf{s} \mid \mathbf{x}) \\&= -\log \sum_{\mathbf{s}} p(\mathbf{s} \mid \mathbf{x}) p(\mathbf{w} \mid \mathbf{x}, \mathbf{s}) \\&\leq -\sum_{\mathbf{s}} p(\mathbf{s} \mid \mathbf{x}) \log p(\mathbf{w} \mid \mathbf{x}, \mathbf{s}) \quad (\text{Jensen's inequality}) \\&\stackrel{\text{MC}}{\approx} -\frac{1}{N} \sum_{i=1}^N p(\mathbf{s}^{(i)} \mid \mathbf{x}) \log p(\mathbf{w} \mid \mathbf{x}, \mathbf{s})\end{aligned}$$

# Learning Hard Attention

- Sample  $N$  sequences of attention decisions from the model
- The gradient is the probability of the gradient of the probability of this sequence scaled by the log probability of generating the target words using that sequence of attention decisions
- This is equivalent to using the REINFORCE algorithm (Williams, 1992) using the log probability of the observed words as a “reward function”. REINFORCE a policy gradient algorithm used for reinforcement learning.





A woman holding a clock in her hand.



A large white bird standing in a forest.

# Attention in Captioning

Add soft attention to image captioning: **+2 BLEU**

Add hard attention to image captioning: **+4 BLEU**

# Summary

- Significant performance improvements
  - Better performance over vector-based encodings
  - Better performance with smaller training data sets
- Model interpretability
- Better gradient flow
- Better capacity (especially obvious for translation)

# Questions?