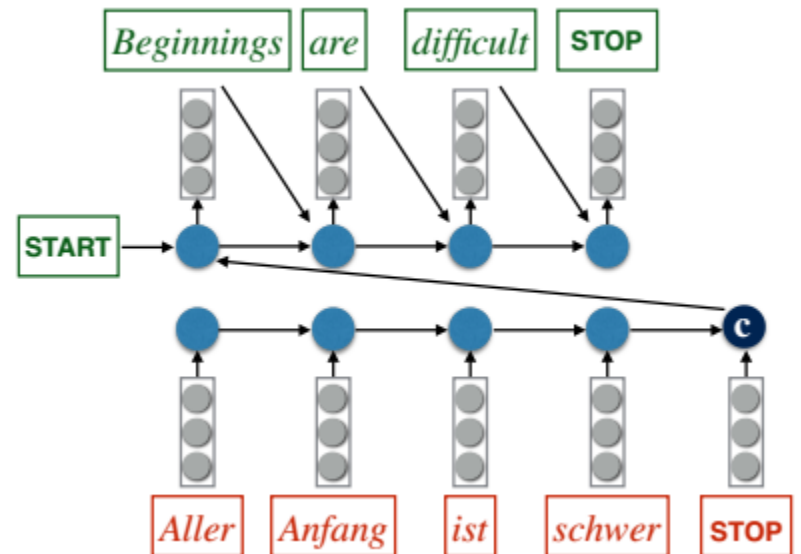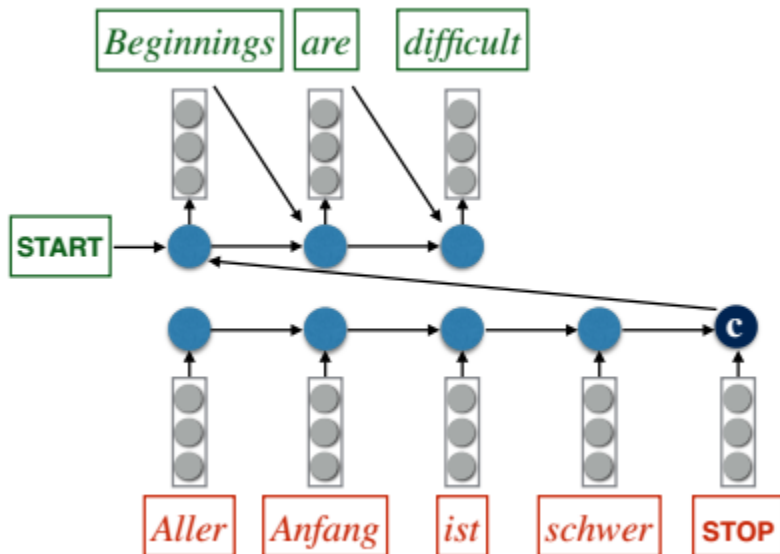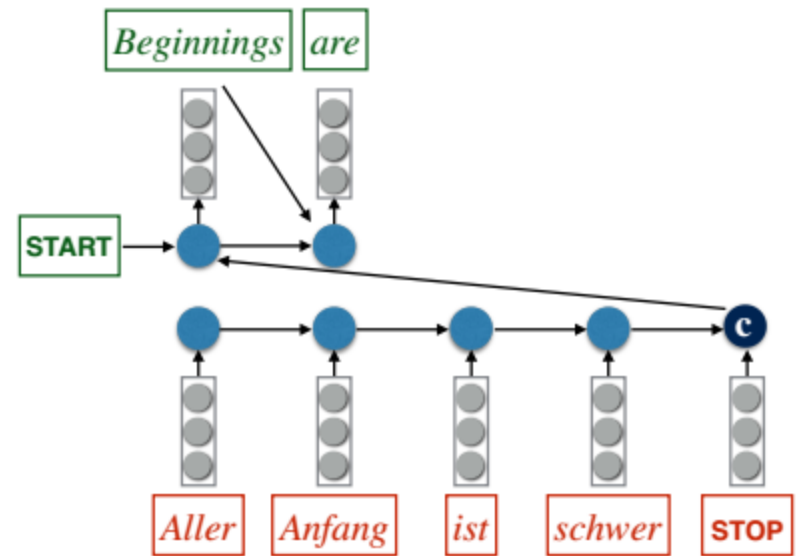# Machine translation models

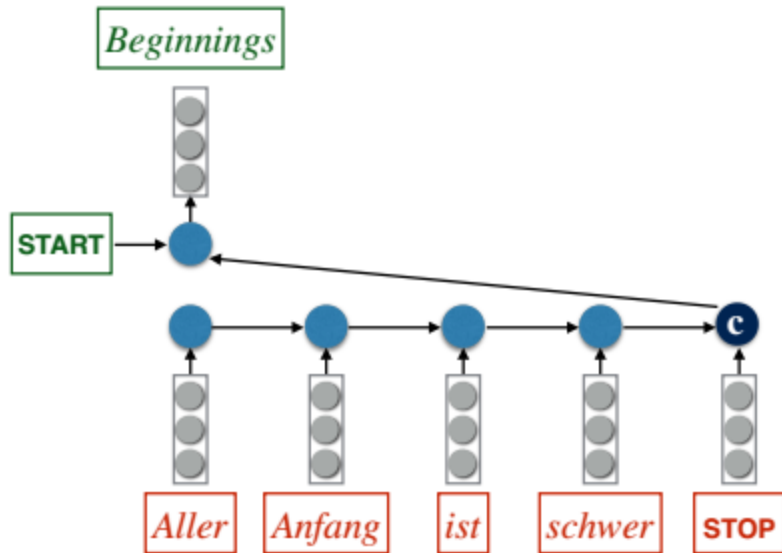# Seq2Seq model1: Embed, Encode

# Sentence encoder: fixed size vector

# Decode from encoded vector

# Pros & Cons

- **Good**

  - RNNs deal naturally with sequences of various lengths

  - LSTMs in principle can propagate gradients a long distance

  - Very simple architecture!

- **Bad**

  - The hidden state has to remember a lot of information!

# Seq2Seq model2: Embed, Encode & Attention

# Sentence encoding as fixed size vector

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?

# Sentence encoding as matrix

- Represent a source sentence as a matrix
- Generate a target sentence from a matrix



*Ich möchte ein Bier*    *Mach's gut*    *Die Wahrheiten der Menschen sind die unwiderlegbaren Irrtümer*

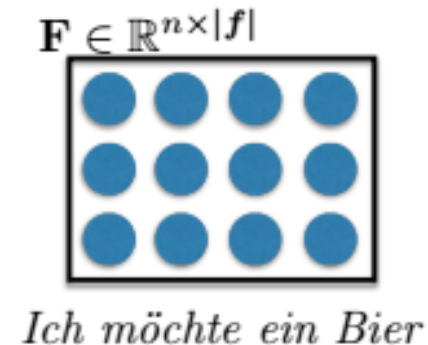*Each matrix has fixed number of rows, but number of columns depends on the number of words

# Sentence encoding as matrix

I.   Using word vectors alone

II.  Using CNN

III. Using RNN/Bi-directional RNN

# I. Sentence matrix with word vectors alone

- Each word type is represented by an n-dimensional vector

- Take all of the vectors for the sentence and concatenate them into a matrix

# I. Sentence matrix with word vectors alone

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

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

$\mathbf{x}_1$     $\mathbf{x}_2$     $\mathbf{x}_3$     $\mathbf{x}_4$

Ich   möchte   ein   Bier

*Ich möchte ein Bier*

# II. Sentence matrix with CNN

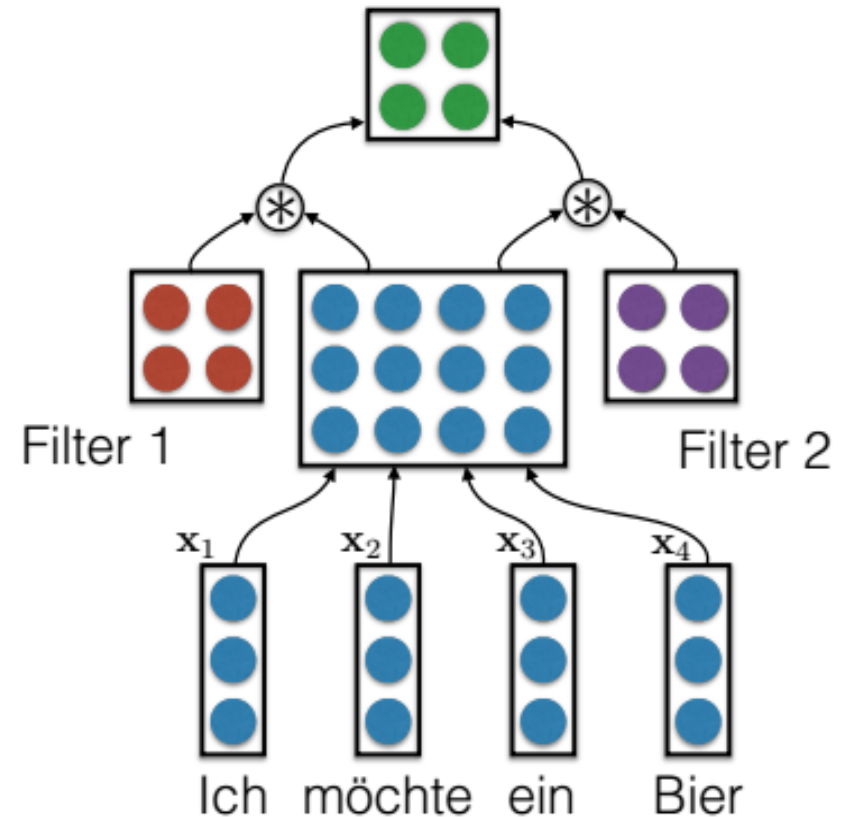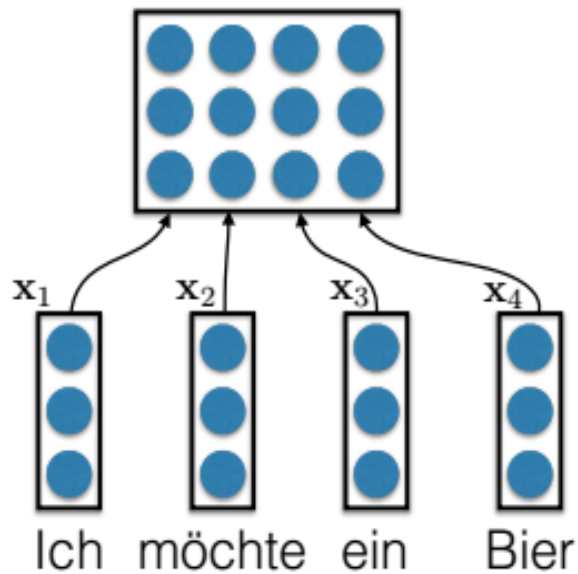- Apply convolutional networks to transform the naïve concatenated matrix to obtain a context-dependent matrix

# II. Sentence matrix with CNN

# III. Sentence matrix with Bi-RNN

- 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

# III. Sentence matrix with Bi-RNN

$$\overleftarrow{\mathbf{h}}_1 \qquad \overleftarrow{\mathbf{h}}_2 \qquad \overleftarrow{\mathbf{h}}_3 \qquad \overleftarrow{\mathbf{h}}_4$$

$$\overrightarrow{\mathbf{h}}_1 \qquad \overrightarrow{\mathbf{h}}_2 \qquad \overrightarrow{\mathbf{h}}_3 \qquad \overrightarrow{\mathbf{h}}_4$$

$$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3 \qquad \mathbf{x}_4$$

Ich   möchte   ein   Bier

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



$\overleftarrow{\mathbf{h}}_1$ $\overleftarrow{\mathbf{h}}_2$ $\overleftarrow{\mathbf{h}}_3$ $\overleftarrow{\mathbf{h}}_4$

$\overrightarrow{\mathbf{h}}_1$ $\overrightarrow{\mathbf{h}}_2$ $\overrightarrow{\mathbf{h}}_3$ $\overrightarrow{\mathbf{h}}_4$

$\mathbf{x}_1$ $\mathbf{x}_2$ $\mathbf{x}_3$ $\mathbf{x}_4$

Ich  möchte  ein  Bier

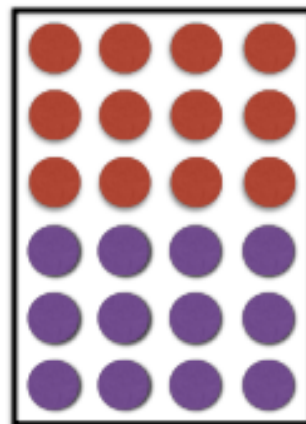$$\mathbf{f}_i = [\overleftarrow{\mathbf{h}}_i; \overrightarrow{\mathbf{h}}_i]$$
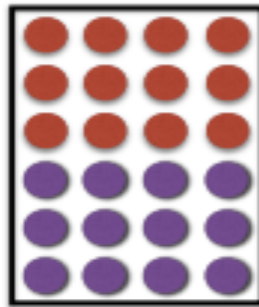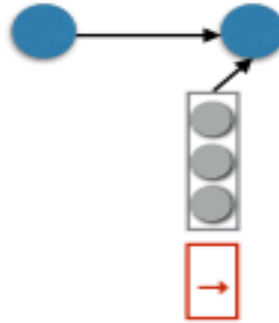


$$\mathbf{F} \in \mathbb{R}^{2n \times |f|}$$

*Ich möchte ein Bier*

# Decoder with Attention



Ich möchte ein Bier

Attention network

*Ich möchte ein Bier*

*I'd*

Attention network

*Ich möchte ein Bier*

Attention network

Attention network

| | I'd | like | | | | I'd | like | a |
|---|---|---|---|---|---|---|---|---|

Attention network

Attention network

Left diagram labels (top to bottom):

I'd    like    a

→    I'd    like    a

Attention network

Right diagram labels (top to bottom):

I'd    like    a    beer

→    I'd    like    a

Attention network

Left diagram labels: I'd, like, a, beer — → , I'd, like, a, beer — Attention network

Right diagram labels: I'd, like, a, beer, STOP — → , I'd, like, a, beer — Attention network

# Decoding with seq2seq model

# Decoding: Greedy search

Simply take the best prediction of the previous timestep as the input to the current timestep.

$$w_1^* = \arg\max_{w_1} p(w_1 \mid \boldsymbol{x})$$

$$w_2^* = \arg\max_{w_2} p(w_2 \mid \boldsymbol{x}, w_1^*)$$

$$\vdots$$

$$w_t^* = \arg\max_{w_2} p(w_t \mid \boldsymbol{x}, \boldsymbol{w}_{<t}^*)$$

# Decoding: Greedy search

- The problem with greedy search, as its name suggested, is that it may not be optimal. More concretely, the joint probability of the output sequence may not be highest possible.

- To do better than greedy search, one can do a full search. We use the recurrent networks to compute the joint probability of every possible output sequence. This can guarantee to find the best possible sequence with a large computation cost.

# Decoding: Full search

In general, we want to find the most probable (MAP) output given the input, i.e.

$$\boldsymbol{w}^* = \arg\max_{\boldsymbol{w}} p(\boldsymbol{w} \mid \boldsymbol{x})$$

$$= \arg\max_{\boldsymbol{w}} \sum_{t=1}^{|\boldsymbol{w}|} \log p(w_t \mid \boldsymbol{x}, \boldsymbol{w}_{<t})$$

This is, for general RNNs, a hard problem.

# Decoding: Beam search

- Keep a list of b possible sequences sorted by the joint probability (which is computed by multiplying the output probability of each prediction in the sequence produced so far).

- During the decode procedure, any sequence that does not belong to the top-b highest joint probability will be removed from our list of candidates.

- Even though it's not guaranteed for the beam search to achieve the optimal sequence, in practice, the generated sequences in the top-b list are generally very good.

# Decoding: Beam search

A slightly better approximation is to use a **beam search** with beam size $b$. Key idea: keep track of top b hypothesis.

E.g., for $b=2$:

$$\boldsymbol{x} = Bier\ trinke\ ich$$

| **beer** | **drink** | **I** |

```
⟨s⟩
logprob=0
```

$w_0$ $\qquad\qquad$ $w_1$ $\qquad\qquad$ $w_2$ $\qquad\qquad$ $w_3$

E.g., for $b=2$:

$x = Bier\ trinke\ ich$
       **beer**    **drink**    **I**

| $\langle s \rangle$ |
| logprob=0 |

| $beer$ |
| logprob=-1.82 |

| $I$ |
| logprob=-2.11 |

$w_0$           $w_1$           $w_2$           $w_3$

E.g., for $b$=2:

$\boldsymbol{x} = Bier \ trinke \ ich$
beer     drink     I



$$\langle s \rangle \quad logprob=0$$

beer
logprob=-1.82

I
logprob=-2.11

drink
logprob=-6.93

I
logprob=-5.80

beer
logprob=-8.66

drink
logprob=-2.87

$w_0$          $w_1$          $w_2$          $w_3$

E.g., for $b=2$:

$\boldsymbol{x} = Bier \ trinke \ ich$

**beer**     **drink**     **I**



```
               drink
               logprob=-6.93

      beer
      logprob=-1.82        I
                           logprob=-5.80
⟨s⟩
logprob=0
      I                    beer
      logprob=-2.11        logprob=-8.66

                           drink
                           logprob=-2.87

$w_0$          $w_1$          $w_2$          $w_3$
```

E.g., for $b$=2:

$x = Bier\ trinke\ ich$
beer      drink      I



| | | | |
|---|---|---|---|
| ⟨s⟩<br>logprob=0 | beer<br>logprob=-1.82 | drink<br>logprob=-6.93 | drink<br>logprob=-6.28 |
| | | I<br>logprob=-5.80 | like<br>logprob=-7.31 |
| | I<br>logprob=-2.11 | beer<br>logprob=-8.66 | beer<br>logprob=-3.04 |
| | | drink<br>logprob=-2.87 | wine<br>logprob=-5.12 |

$w_0$        $w_1$        $w_2$        $w_3$

E.g., for $b$=2:

$\boldsymbol{x} = Bier\ trinke\ ich$

beer      drink      I



$w_0$            $w_1$            $w_2$            $w_3$

Diagram nodes:

- ⟨s⟩ logprob=0
- beer logprob=-1.82
- I logprob=-2.11
- drink logprob=-6.93
- I logprob=-5.80
- beer logprob=-8.66
- drink logprob=-2.87
- drink logprob=-6.28
- like logprob=-7.31
- beer logprob=-3.04
- wine logprob=-5.12

A slightly better approximation is to use a **beam search** with beam size $b$. Key idea: keep track of top b hypothesis.

E.g., for $b=2$:



$\mathbf{x} = Bier\ trinke\ ich$