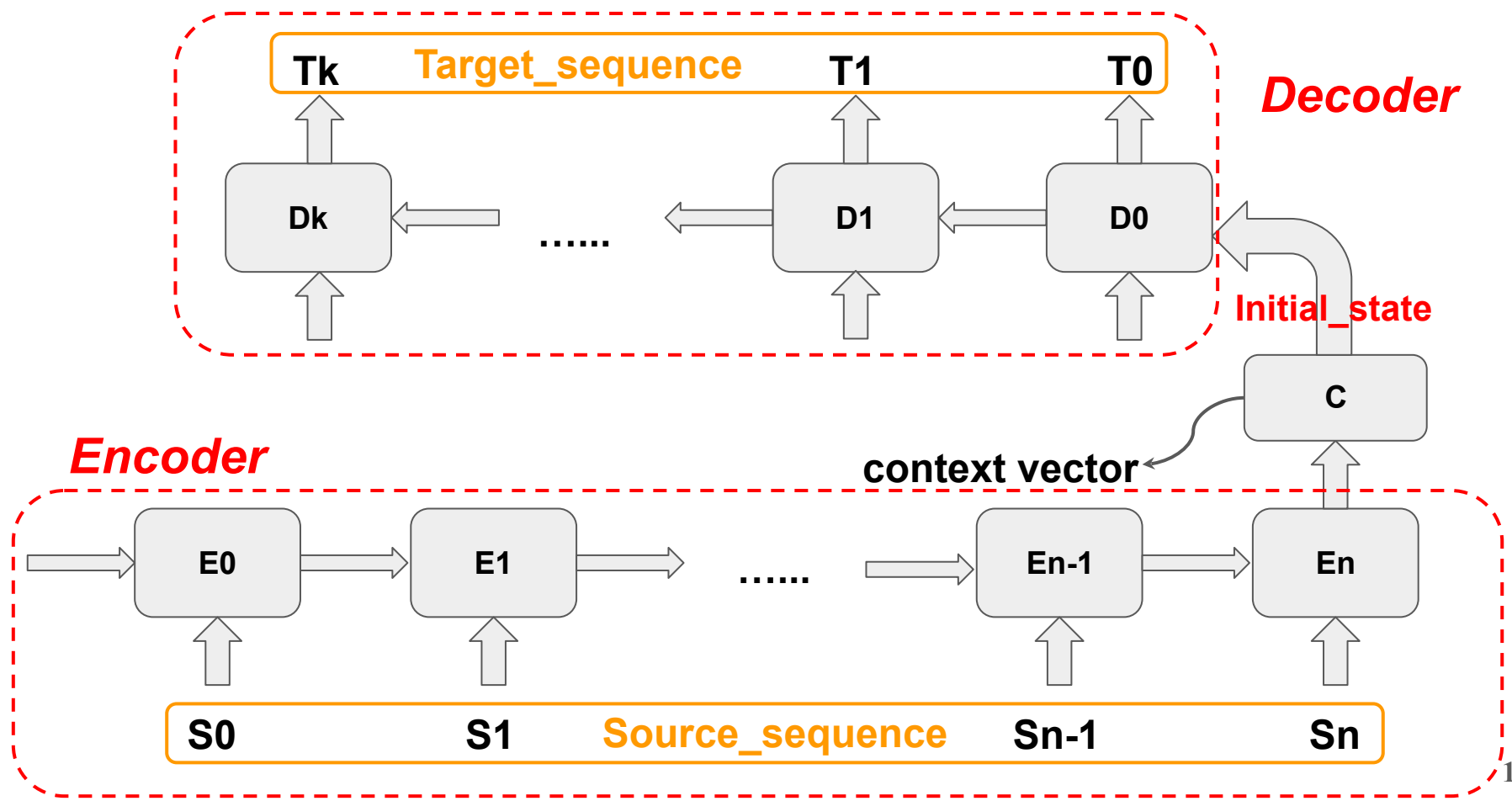


Seq2seq



Seq2seq

- Same weight for each target sentence

$$p(y_t \mid \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$

Decoder

Initial_state

c

context vector

Encoder

- Encode each source sentence into a **fix-length** context vector

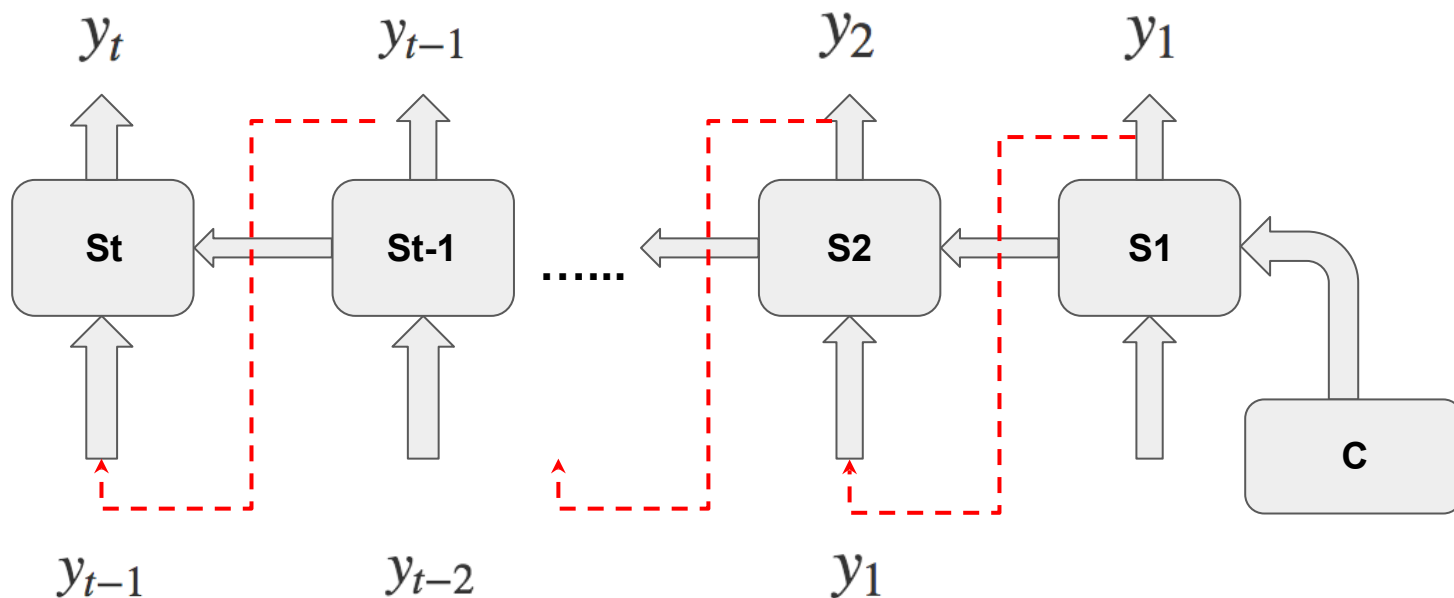
$$h_t = f(x_t, h_{t-1})$$
$$c = q(\{h_1, \dots, h_{T_x}\})$$

Seq2seq

$$p(y_t \mid \{y_1, \dots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$

$$s_t = \tanh(Uy_{t-1} + Ws_{t-1} + b)$$

$$y_t = \text{sigmoid}(Vs_t)$$



NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau

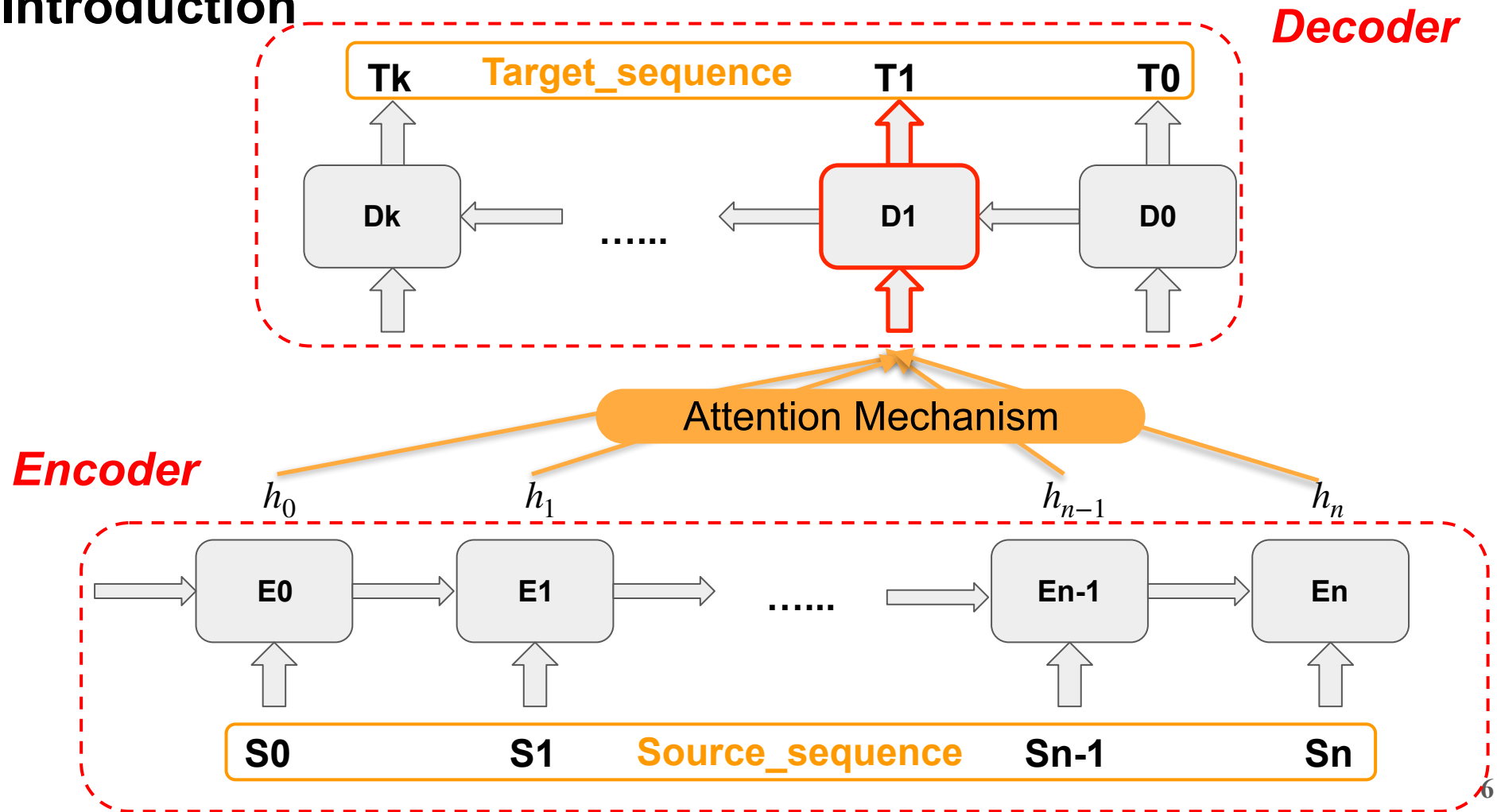
Jacobs University Bremen, Germany

KyungHyun Cho Yoshua Bengio*

Université de Montréal

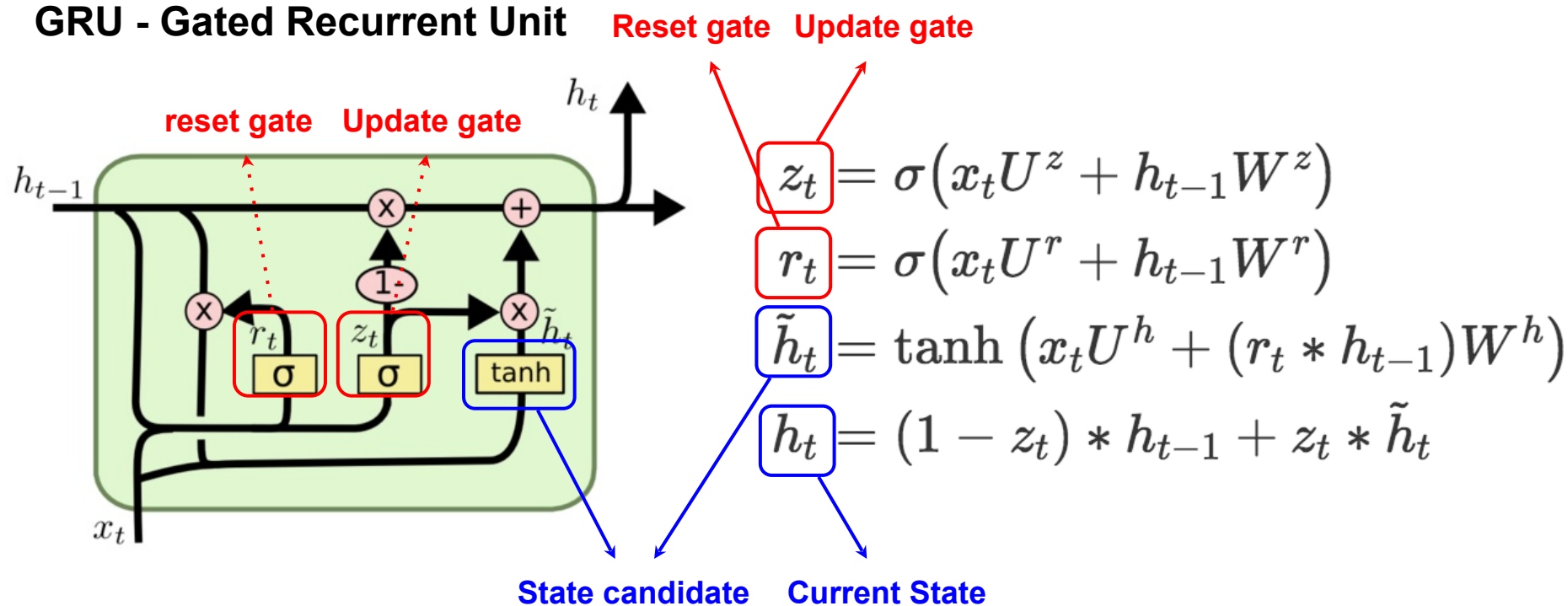
- Introduction
- Related work
- Attention mechanism
- Datasets
- Comparison
- Evaluation

Introduction

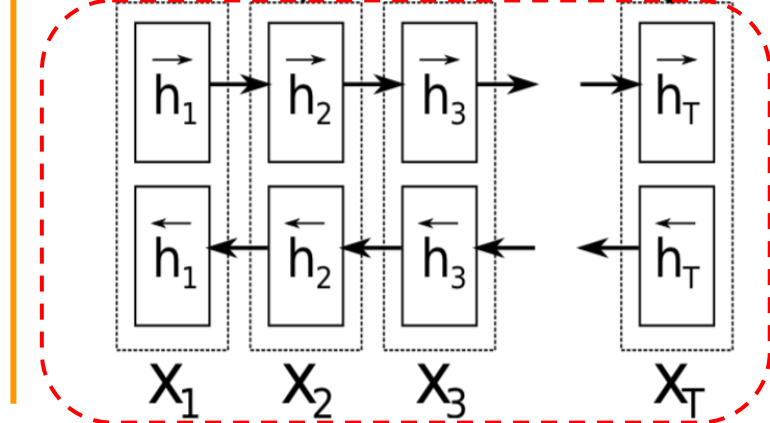
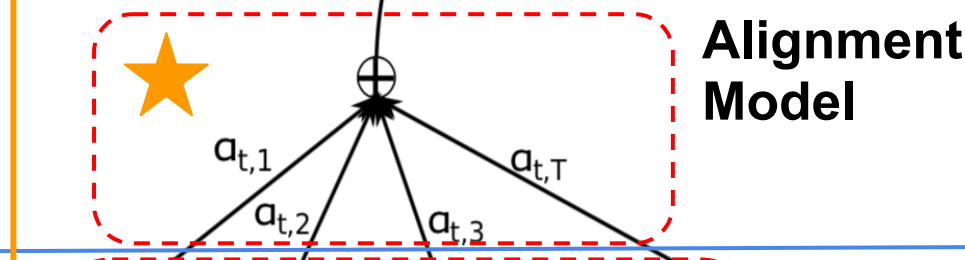
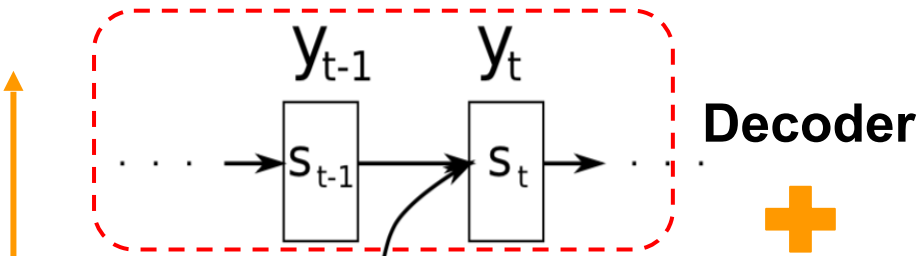


Related work

GRU - Gated Recurrent Unit



Attention Mechanism



$$p(y_i|y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, \underline{c_i})$$

$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

$$\mathbf{x} = (x_1, \dots, x_{T_x}), x_i \in \mathbb{R}^{K_x}$$

- 1. $e_{ij} = a(s_{i-1}, h_j)$
- 2. $\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$
- 3. $c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$

Encoder

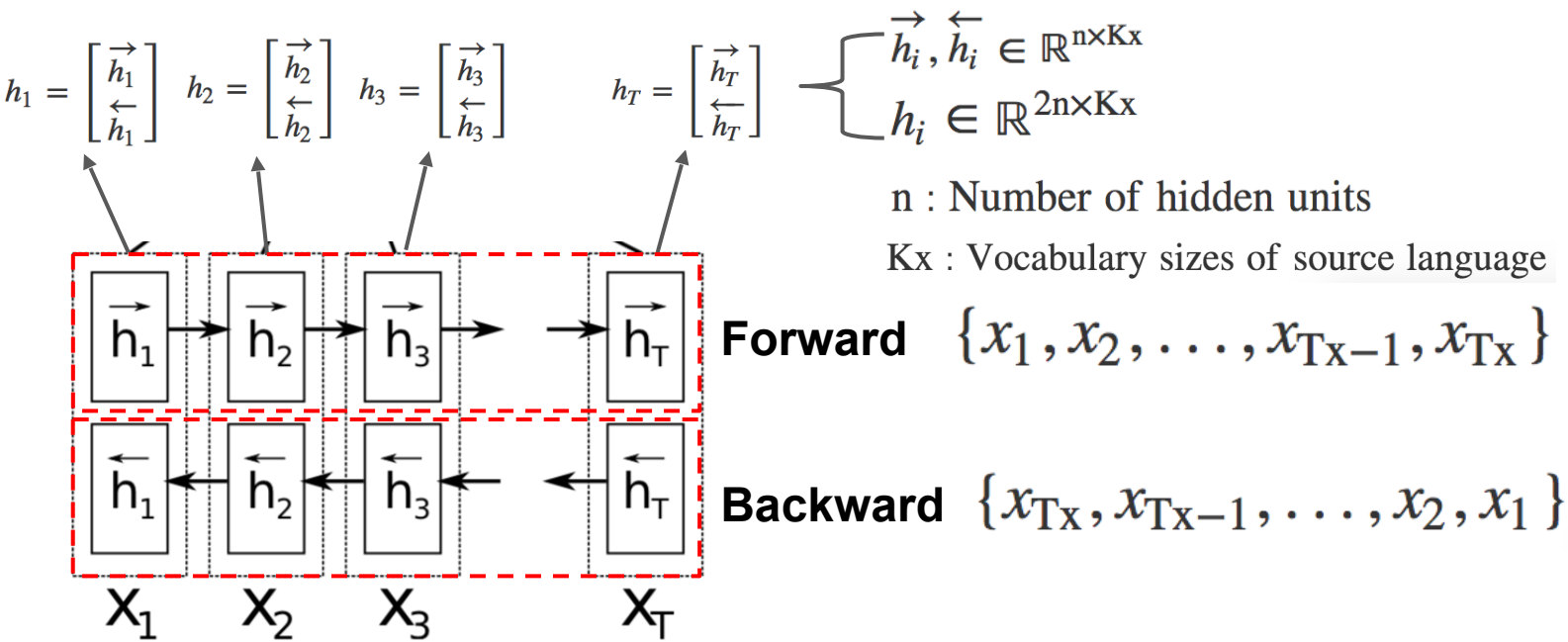
$$\vec{h}_i = \begin{cases} (1 - \vec{z}_i) \circ \vec{h}_{i-1} + \vec{z}_i \circ \underline{h}_i & , \text{ if } i > 0 \\ 0 & , \text{ if } i = 0 \end{cases}$$

\vec{z}_i : Update gate \underline{h}_i : State candidate
 \vec{h}_i : Current state

Attention Mechanism

Bi-directional GRU (Forward)

$$\vec{h}_i = \begin{cases} (1 - \vec{z}_i) \circ \vec{h}_{i-1} + \vec{z}_i \circ \vec{h}_i & , \text{ if } i > 0 \\ 0 & , \text{ if } i = 0 \end{cases} \quad \bullet : \text{Element-wise Multiplication}$$

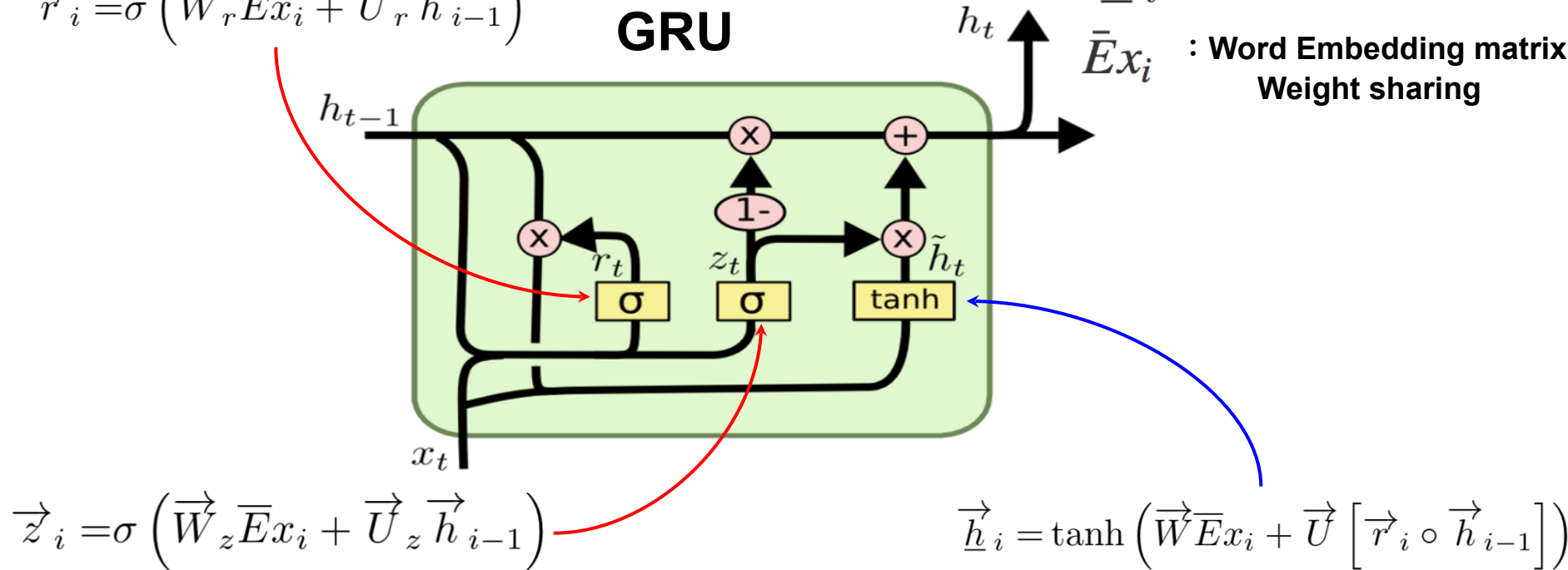


Attention Mechanism

$$\vec{h}_i = \begin{cases} (1 - \vec{z}_i) \circ \vec{h}_{i-1} + \vec{z}_i \circ \vec{h}_i & , \text{ if } i > 0 \\ 0 & , \text{ if } i = 0 \end{cases}$$

$$\vec{r}_i = \sigma \left(\vec{W}_r \bar{E}x_i + \vec{U}_r \vec{h}_{i-1} \right)$$

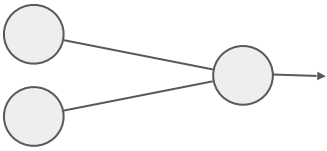
- \vec{z}_i : Update gate
- \vec{r}_i : Reset gate
- σ : Sigmoid
- \vec{h}_i : State candidate
- $\bar{E}x_i$: Word Embedding matrix Weight sharing



Attention Mechanism

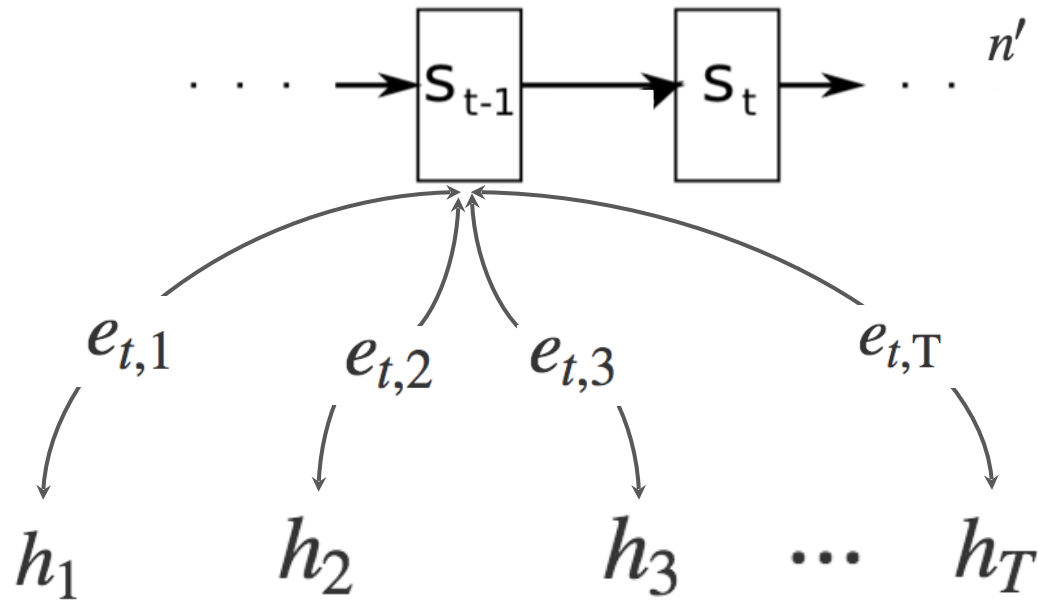
1. $e_{ij} = a(s_{i-1}, h_j)$ **Soft-Alignment**

$$e_{ij} = v_a^\top \tanh(W_a s_{i-1} + U_a h_j) \quad \begin{cases} W_a \in \mathbb{R}^{n' \times n} \\ U_a \in \mathbb{R}^{n' \times 2n} \\ v_a \in \mathbb{R}^{n'} \end{cases}$$



n : Number of hidden units

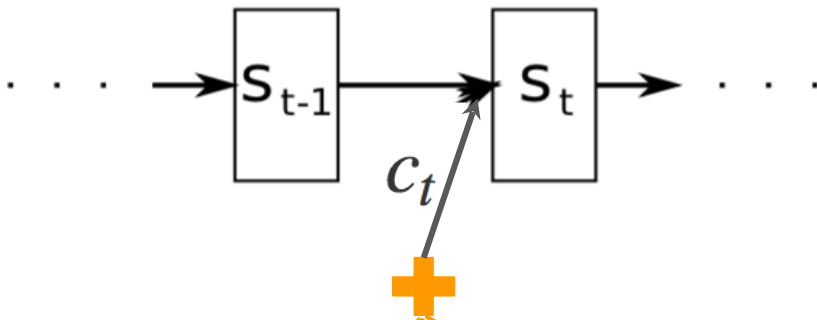
n' : Number of Alignment hidden size



How Match between **1-th** source word and **t-th** target word

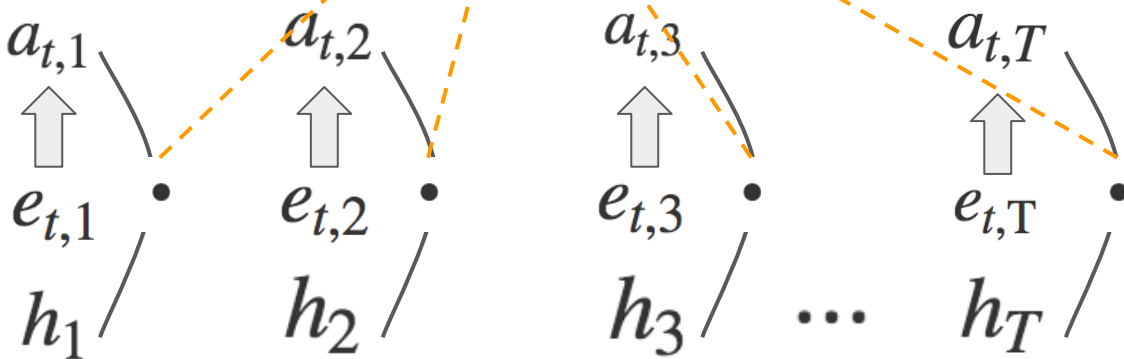
Attention Mechanism

3. $c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$



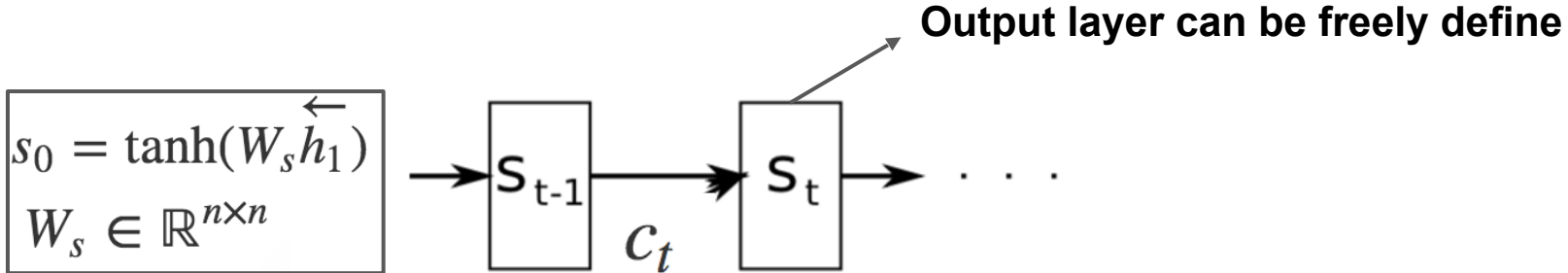
Regard α_{ij} as weight for each h

2. $\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}$



• : Dot product

Attention Mechanism



$$s_i = f(s_{i-1}, y_{i-1}, c_i)$$

$$= (1 - z_i) \circ s_{i-1} + z_i \circ \tilde{s}_i$$

$$\tilde{s}_i = \tanh(W E y_{i-1} + U [r_i \circ s_{i-1}] + C c_i)$$

$$z_i = \sigma(W_z E y_{i-1} + U_z s_{i-1} + C_z c_i)$$

$$r_i = \sigma(W_r E y_{i-1} + U_r s_{i-1} + C_r c_i)$$

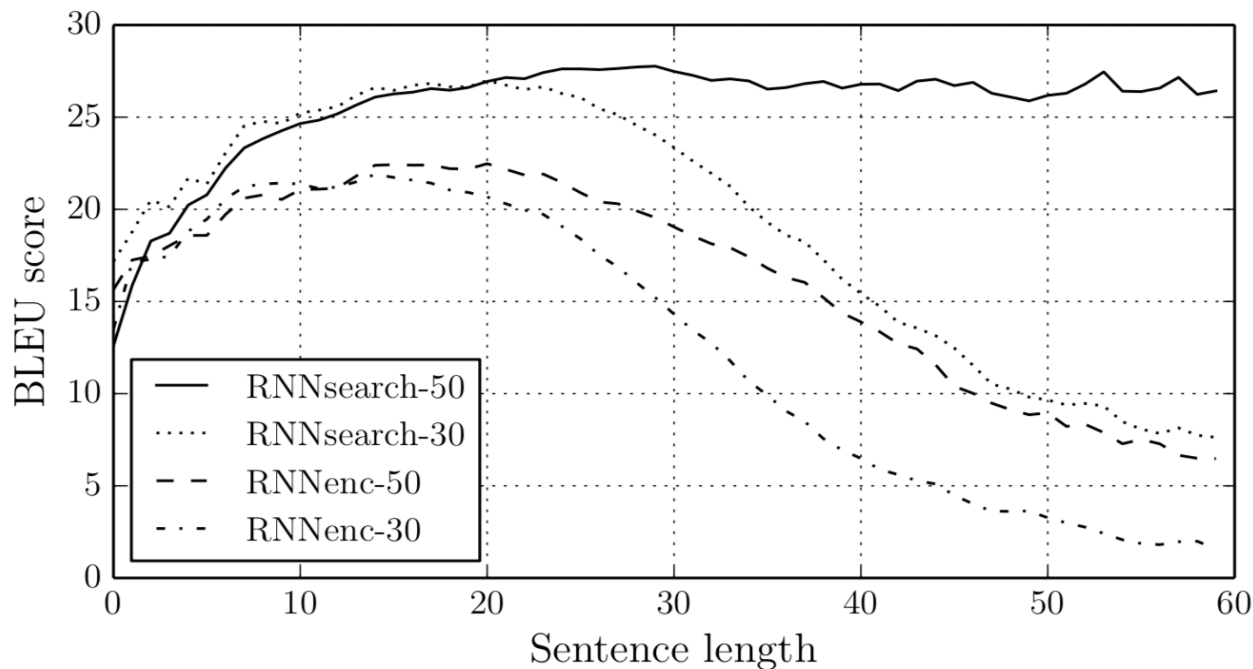
Be a little different from GRU

$U, U_z, U_r \in \mathbb{R}^{n \times n}$
 $C, C_z, C_r \in \mathbb{R}^{n \times 2n}$
 $W, W_z, W_r \in \mathbb{R}^{n \times m}$

Comparison

Training : Sentence of length up to 30 words (RNNencdec-30, RNNsearch-30)
Sentence of length up to 50 words (RNNencdec-50, RNNsearch-50)

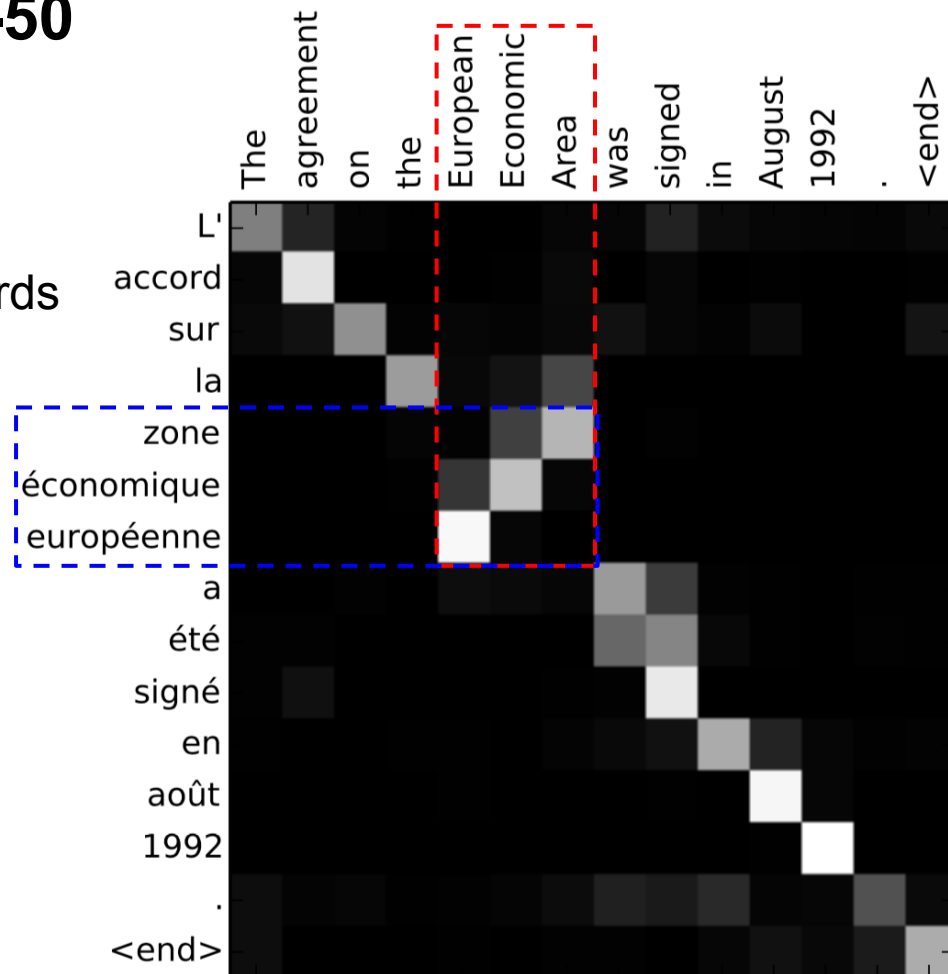
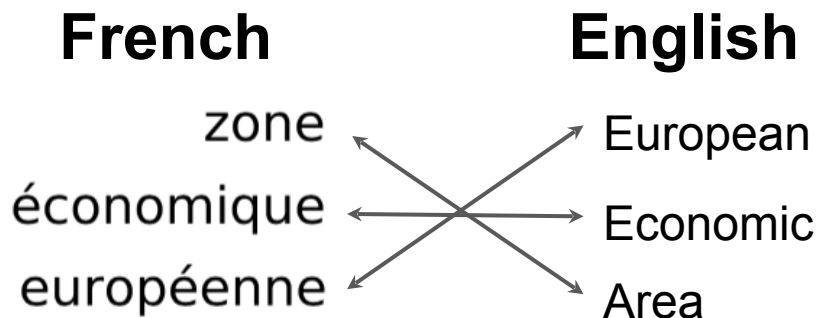
Testing : On the full test set which includes sentences having unknown words



Performance on RNNsearch-50

Cond : Arbitrary sentence on test set

Includes sentences having unknown words



Performance

BLEU scores of the models computed on the test set

All : All the test set

No UNK : Without any unknown words on test set

RNNsearch-50* : Train longer until the performance on the valid stopped improving

Model	All	No UNK ^o
RNNencdec-30	13.93	24.19
RNNsearch-30	21.50	31.44
RNNencdec-50	17.82	26.71
RNNsearch-50	26.75	34.16
RNNsearch-50*	28.45	36.15
Moses	33.30	35.63

Next...

Effective Approaches to Attention-based Neural Machine Translation

Minh-Thang Luong Hieu Pham Christopher D. Manning
Computer Science Department Stanford University, Stanford, CA, 94305
`{lmthang, hyhieu, manning}@stanford.edu`

New alignment model

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \textit{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \textit{general} \\ \mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s] & \textit{concat} \end{cases}$$