

# CS 6501 Natural Language Processing

Machine Translation,  
Sequence-to-Sequence Models

---

Yangfeng Ji

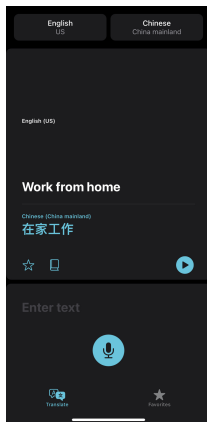
Information and Language Processing Lab  
Department of Computer Science  
University of Virginia



1. Noisy Channel Model
2. IBM Models
3. Seq2seq Models
4. Attention Mechanism

# Machine Translation Systems

A screenshot of my Apple Translate app



Commerical systems

- ▶ Google Translate
- ▶ Microsoft Translator
- ▶ Apple Translate
- ▶ Amazon Translate
- ▶ IBM Watson Language Translator
- ▶ ...

# Problem Formulation

- ▶ Goal: translate French to English
- ▶ Mathematical formulation

$$P(e | f) \tag{1}$$

where

- ▶  $f = (f_1, \dots, f_m)$  is a French sentence
- ▶  $e = (e_1, \dots, e_n)$  is an English translation

# Noisy Channel Model

---

# Noisy Channel Model: Definition

Consider a hypothetical communication channel that always adds some noise to the input clean signal, the task of decoding a noisy channel model is to decode the original clean signal  $x$  from the received noisy signal  $y$



- ▶ Input signal:  $x$
- ▶ Received signal:  $y$
- ▶ Task: Decode  $x$  from  $y$

# Noisy Channel Model: For translation

We can consider machine translation as a decoding task from a noisy channel

- ▶ decoding the clean signal from a noisy input, and
- ▶ decoding the English text from a foreign language



For example:

- ▶ Source language: French
- ▶ Target language: English
- ▶ Task: Translate French to English

# Problem Formulation

- Assume the probabilistic formulation of a noisy channel is  $P(f | e)$ , where  $f = (f_1, f_2, \dots, f_m)$  represents a French sentence and  $e = (e_1, e_2, \dots, e_n)$  represents the corresponding English sentence

Source Model  $P(e)$

Channel Model  $P(f | e)$





# Problem Formulation

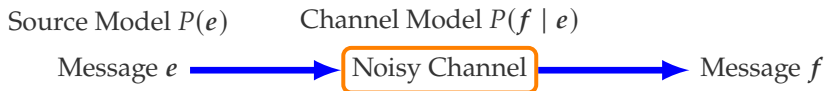
- Assume the probabilistic formulation of a noisy channel is  $P(f | e)$ , where  $f = (f_1, f_2, \dots, f_m)$  represents a French sentence and  $e = (e_1, e_2, \dots, e_n)$  represents the corresponding English sentence



- The task is to find out  $P(e | f)$ , which is the probabilistic formulation of the translation task

# Problem Formulation

- Assume the probabilistic formulation of a noisy channel is  $P(\mathbf{f} | \mathbf{e})$ , where  $\mathbf{f} = (f_1, f_2, \dots, f_m)$  represents a French sentence and  $\mathbf{e} = (e_1, e_2, \dots, e_n)$  represents the corresponding English sentence



- The task is to find out  $P(\mathbf{e} | \mathbf{f})$ , which is the probabilistic formulation of the translation task
- To solve the problem, we need the Bayes theorem

$$P(\mathbf{e} | \mathbf{f}) = \frac{P(\mathbf{f} | \mathbf{e}) \cdot P(\mathbf{e})}{P(\mathbf{f})} \quad (2)$$

where we need an extra component, the prior distribution of  $\mathbf{e}$ ,  $P(\mathbf{e})$

# Problem Formulation

- Assume the probabilistic formulation of a noisy channel is  $P(f | e)$ , where  $f = (f_1, f_2, \dots, f_m)$  represents a French sentence and  $e = (e_1, e_2, \dots, e_n)$  represents the corresponding English sentence



- The task is to find out  $P(e | f)$ , which is the probabilistic formulation of the translation task
- To solve the problem, we need the Bayes theorem

$$P(e | f) = \frac{P(f | e) \cdot P(e)}{P(f)} = \frac{P(f, e)}{\sum_{e'} P(f, e')} \quad (2)$$

where we need an extra component, the prior distribution of  $e$ ,  $P(e)$

# Two Components

For machine translation,

$$P(e | f) = \frac{P(e)P(f | e)}{P(f)} \quad (3)$$

Translating from  $f$  to  $e$  is essentially a decoding (prediction) problem, therefore we can ignore the

$$\hat{e} = \operatorname{argmax}_e P(e | f) = \operatorname{argmax}_e P(e)P(f | e) \quad (4)$$

# Two Components

For machine translation,

$$P(e | f) = \frac{P(e)P(f | e)}{P(f)} \quad (3)$$

Translating from  $f$  to  $e$  is essentially a decoding (prediction) problem, therefore we can ignore the

$$\hat{e} = \underset{e}{\operatorname{argmax}} P(e | f) = \underset{e}{\operatorname{argmax}} P(e)P(f | e) \quad (4)$$

- ▶  $P(e)$ : the language model
  - ▶  $n$ -gram language models
  - ▶ recurrent neural network language models

# Two Components

For machine translation,

$$P(e | f) = \frac{P(e)P(f | e)}{P(f)} \quad (3)$$

Translating from  $f$  to  $e$  is essentially a decoding (prediction) problem, therefore we can ignore the

$$\hat{e} = \operatorname{argmax}_e P(e | f) = \operatorname{argmax}_e P(e)P(f | e) \quad (4)$$

- ▶  $P(e)$ : the language model
  - ▶  $n$ -gram language models
  - ▶ recurrent neural network language models
- ▶  $P(f | e)$ : the translation model
  - ▶ A probabilistic mapping from English to French

## Two Components (Cont.)

Divide one big problem into two subproblems:  $P(e)$  and  $P(f | e)$ , then solve them separately (with extra resources)



## Two Components (Cont.)

Divide one big problem into two subproblems:  $P(e)$  and  $P(f | e)$ , then solve them separately (with extra resources)



- ▶ For example,  $P(e)$  is essentially a language model, which can be trained with as many training example as we have (no annotation needed)



## Two Components (Cont.)

Divide one big problem into two subproblems:  $P(e)$  and  $P(f | e)$ , then solve them separately (with extra resources)



- ▶ For example,  $P(e)$  is essentially a language model, which can be trained with as many training examples as we have (no annotation needed)
- ▶ A key problem in statistical machine translation is to model  $P(f | e)$ . In this lecture we will discuss two methods
  - ▶ IBM Model 1
  - ▶ IBM Model 2
  - ▶ IBM Model ...

## IBM Models

---

# Challenge of Probability Definition

For modeling the probability distribution of a long sequence: a similar argument has also been used in the language modeling task.

Given

- ▶ an English sentence  $e$  with  $n$  words  $(e_1, \dots, e_n)$  and
- ▶ a French sentence  $f$  with  $m$  words  $(f_1, \dots, f_m)$ ,

directly modeling

$$P(f | e) = P(f_1, \dots, f_m | e_1, \dots, e_n) \quad (5)$$

is challenging.

Consider the specific example

$e$  = And the program has been implemented

$f$  = Le programme a ete mis en application

- ▶ To directly model the conditional probability of  $f$  given  $e$ ,  $P(f | e)$  defines a probability on a 13-dimensional space<sup>1</sup>

---

<sup>1</sup>13 is the total number of English and French words in the source and target sentences.

# Example

Consider the specific example

$e$  = And the program has been implemented

$f$  = Le programme a ete mis en application

- ▶ To directly model the conditional probability of  $f$  given  $e$ ,  $P(f | e)$  defines a probability on a 13-dimensional space<sup>1</sup>
- ▶ Each  $f_j$  depends on only part of  $e$ . In other words, (1) there are **alignments** between French and English words, and (2) word dependency only exists between source words and target words with alignments.

---

<sup>1</sup>13 is the total number of English and French words in the source and target sentences.

## Example (Cont.)

For the previous example, here are some example alignments between words in English and French

	1	2	3	4	5	6	7
<i>e</i>	And	the	program	has	been	implemented	
<i>f</i>	Le	programme	a	ete	mis	en	application

## Example (Cont.)

For the previous example, here are some example alignments between words in English and French

	1	2	3	4	5	6	7
<i>e</i>	And	the	program	has	been	implemented	
<i>f</i>	Le	programme	a	ete	mis	en	application

## Example (Cont.)

For the previous example, here are some example alignments between words in English and French

	1	2	3	4	5	6	7
<i>e</i>	And	the	program	has	been	implemented	
<i>f</i>	Le	programme	a	ete	mis	en	application

- ▶ If we use  $a_j = i$  to represent that the  $j$ -th word in  $f$  is aligned with the  $i$ -th word in  $e$ . For the abovementioned example, we have  $a_1 = 2$ ,  $a_6 = 6$



## Example (Cont.)

For the previous example, here are some example alignments between words in English and French

	1	2	3	4	5	6	7
<i>e</i>	And	the	program	has	been	implemented	
<i>f</i>	Le	programme	a	ete	mis	en	application

- ▶ If we use  $a_j = i$  to represent that the  $j$ -th word in  $f$  is aligned with the  $i$ -th word in  $e$ . For the abovementioned example, we have  $a_1 = 2$ ,  $a_6 = 6$
- ▶ With the *explicit* alignments, we can simplify the conditional probability. For example,

$$P(f_1 | e, a_1) = P(f_1 | e_2, a_1) \quad (6)$$

# Alignments

Similarly, we introduce new alignment variables  $\mathbf{a} = (a_1, \dots, a_m)$

$$P(f_1, \dots, f_m, a_1, \dots, a_m \mid e_1, \dots, e_n) \quad (7)$$

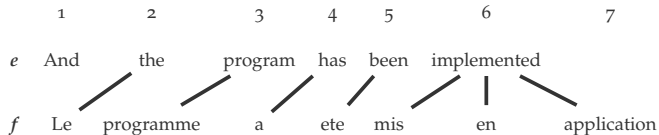
where  $a_j \in \{0, 1, \dots, n\}$

# Alignments

Similarly, we introduce new alignment variables  $\mathbf{a} = (a_1, \dots, a_m)$

$$P(f_1, \dots, f_m, a_1, \dots, a_m \mid e_1, \dots, e_n) \quad (7)$$

where  $a_j \in \{0, 1, \dots, n\}$



$$a_1 = 2 \quad a_2 = 3$$

$$a_3 = 4 \quad a_4 = 5$$

$$a_5 = 6 \quad a_6 = 6$$

$$a_7 = 6$$

Further break down  $P(f, a | e)$  into two parts:

$$P(f, a | e) = P(a | e)P(f | a, e) \quad (8)$$

- ▶ Alignment  $P(a | e)$
- ▶ Translation with a given alignment  $a, P(f | a, e)$
- ▶ Both  $P(a | e)$  and  $P(f | a, e)$  can be further factorized

## IBM Model 1: $P(\mathbf{a} \mid \mathbf{e})$

In **IBM Model 1**, all alignments are equally likely

$$P(\mathbf{a} \mid \mathbf{e}) = \prod_{j=1}^m q(a_j = i \mid j, n, m) = \frac{1}{(n+1)^m} \quad (9)$$

where  $m$  and  $n$  are the lengths of  $\mathbf{f}$  and  $\mathbf{e}$  respectively

# IBM Model 1: $P(\mathbf{a} \mid \mathbf{e})$

In **IBM Model 1**, all alignments are equally likely

$$P(\mathbf{a} \mid \mathbf{e}) = \prod_{j=1}^m q(a_j = i \mid j, n, m) = \frac{1}{(n+1)^m} \quad (9)$$

where  $m$  and  $n$  are the lengths of  $\mathbf{f}$  and  $\mathbf{e}$  respectively

- ▶ Uniform distribution: major simplification, great starting point

$$q(a_j = i \mid j, n, m) = \frac{1}{(n+1)} \quad (10)$$

# IBM Model 1: $P(a | e)$

In **IBM Model 1**, all alignments are equally likely

$$P(a | e) = \prod_{j=1}^m q(a_j = i | j, n, m) = \frac{1}{(n + 1)^m} \quad (9)$$

where  $m$  and  $n$  are the lengths of  $f$  and  $e$  respectively

- ▶ Uniform distribution: major simplification, great starting point

$$q(a_j = i | j, n, m) = \frac{1}{(n + 1)} \quad (10)$$

- ▶ Why  $n + 1$ ? A: some words  $f$  cannot find an alignment with any word  $e$ , so we need a dummy token.

The translation probability with the alignment  $\mathbf{a}$  on condition

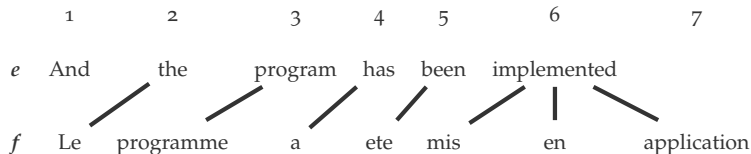
$$P(\mathbf{f} \mid \mathbf{a}, \mathbf{e}) = \prod_{j=1}^m t(f_j \mid e_{a_j}) \quad (11)$$

- ▶  $f_j$ : the  $j$ -th French word
- ▶  $a_j = i$ : the alignment of the  $j$ -th French word
- ▶  $e_{a_j}$ : the aligned English word of the  $j$ -th French word
- ▶  $t(f_j \mid e_{a_j})$ : the translation probability from the  $a_j (= i)$ -th English word to the  $j$ -th French word



## Example (Cont.)

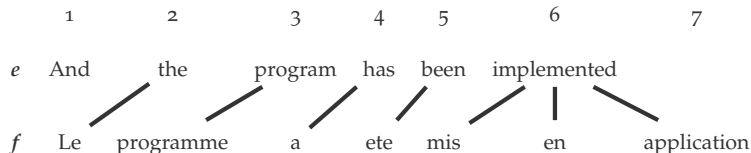
Given the previous example



we have the translation probability with alignment  $P(f | a, e)$

## Example (Cont.)

Given the previous example



we have the translation probability with alignment  $P(f | a, e)$

$$\begin{aligned} P(f | a, e) = & t(\text{Le} | \text{the}) \cdot t(\text{programme} | \text{program}) \cdot \\ & t(\text{a} | \text{has}) \cdot t(\text{ete} | \text{been}) \cdot \\ & t(\text{mis} | \text{implemented}) \cdot t(\text{en} | \text{implemented}) \cdot \\ & t(\text{application} | \text{implemented}) \end{aligned}$$

The final probability  $P(\mathbf{f} | \mathbf{e})$  after we have each pieces of information

$$\begin{aligned} P(\mathbf{f} | \mathbf{e}) &= \sum_{\mathbf{a}} P(\mathbf{f}, \mathbf{a} | \mathbf{e}) \\ &= \sum_{\mathbf{a}} P(\mathbf{a} | \mathbf{e}) P(\mathbf{f} | \mathbf{a}, \mathbf{e}) \\ &= \sum_{\mathbf{a}} \frac{1}{(n+1)^m} \prod_{j=1}^m t(f_j | e_{a_j}) \end{aligned} \tag{12}$$

*Basic idea:* Break a big conditional probability into small pieces on word pairs

## Seq2seq Models

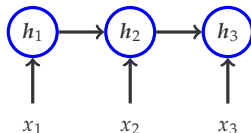
---

# Recurrent Neural Networks (RNNs)

A simple RNN is defined by the following recursive function

$$h_t = f(x_t, h_{t-1}) \quad (13)$$

and depicted as



where

- ▶  $h_{t-1}$ : hidden state at time step  $t - 1$
- ▶  $x_t$ : input at time step  $t$
- ▶  $h_t$ : hidden state at time step  $t$

For a machine translation problem

- ▶ Input:  $f = (f_1, f_2, \dots, f_m)$
- ▶ Output:  $e = (e_1, e_2, \dots, e_n)$

Neural machine translation usually model the conditional probability  $P(e | f)$  directly as

$$P(e | f) = P(e_1 | f) \cdot P(e_2 | e_1, f) \quad (14)$$

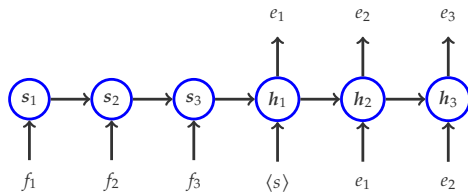
$$\cdot P(e_3 | e_{1:2}, f) \cdots \quad (15)$$

$$\cdot P(e_n | e_{1:n-1}, f) \quad (16)$$

$$= \prod_{i=1}^n P(e_i | e_{1:i-1}, f) \quad (17)$$

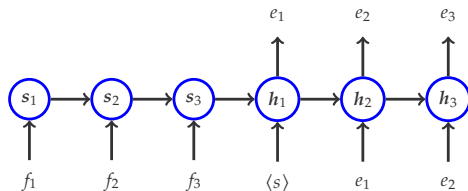
# Neural Sequence-to-sequence Models

*Ideally*, a sequence-to-sequence model offers a natural framework of mapping a sequence to another sequence *step-by-step*



# Neural Sequence-to-sequence Models

*Ideally*, a sequence-to-sequence model offers a natural framework of mapping a sequence to another sequence *step-by-step*

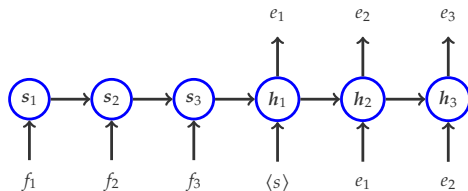


- ▶ The first RNN encoding the input sequence  $f$  is called the **encoder**
- ▶ The second RNN decoding the output sequence  $e$  step-by-step is called the **decoder**



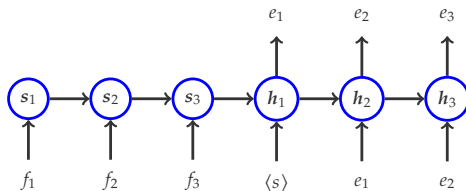
# Neural Sequence-to-sequence Models

*Ideally*, a sequence-to-sequence model offers a natural framework of mapping a sequence to another sequence *step-by-step*



- ▶ The first RNN encoding the input sequence  $f$  is called the **encoder**
- ▶ The second RNN decoding the output sequence  $e$  step-by-step is called the **decoder**
- ▶ In a broader sense, this model is also called the encoder-decoder model

# Neural Sequence-to-sequence Models (Cont.)



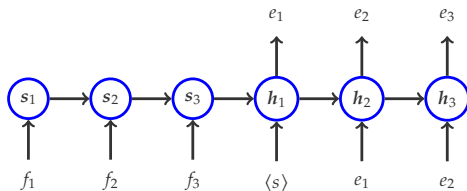
Specifically, for each decoding step, we can write the conditional probability as

$$P(e_i | e_{1:i-1}, f) = \text{softmax}(\mathbf{W}_o \mathbf{h}_i) \quad (18)$$

where  $\mathbf{h}_i$  is the hidden state on step  $i$

## Tricks (I)

A basic seq2seq model is nothing more than connecting two RNNs (LSTMs) together



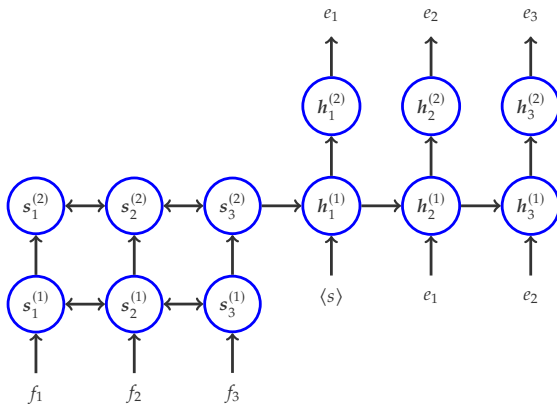
There are some additional tricks that to make it work

- ▶ two different LSTMs: one for input sequence and the other for output sequence
- ▶ it can greatly improve the performance on the output side, by reversing the order of input sequence

[Sutskever et al., 2014]

## Tricks (II)

- ▶ stacked (or deep) LSTM with multiple layers on both encoder and decoder, which has much more potential than single-layer LSTMs
- ▶ bi-directional LSTM for the encoder



# Seq2seq Models for Machine Translation

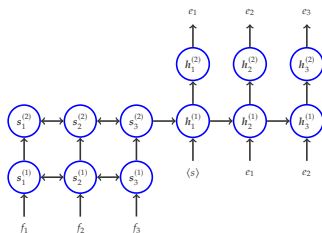
Early ways of using seq2seq models for machine translation

- ▶ Use seq2seq model evaluation scores to **re-rank** the  $k$ -best list [Sutskever et al., 2014]

# Seq2seq Models for Machine Translation

Early ways of using seq2seq models for machine translation

- ▶ Use seq2seq model evaluation scores to **re-rank** the  $k$ -best list [Sutskever et al., 2014]
- ▶ Use seq2seq models evaluation scores as **additional features** in the translation model [Cho et al., 2014]



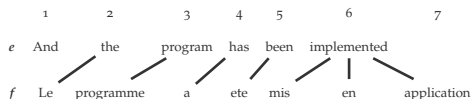
Major limitation: the input signal is probably too weak on the decoder side

# Attention Mechanism

---

# Alignment in Statistical MT

Another important module in statistical machine translation is the alignments offered by  $P(f | a, e)$ , which essentially is a **word-level mapping** between the input and output

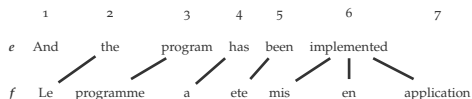


$$\begin{aligned} P(f | a, e) = & t(\text{Le} | \text{the}) \cdot t(\text{programme} | \text{program}) \cdot \\ & t(\text{a} | \text{has}) \cdot t(\text{ete} | \text{been}) \cdot \\ & t(\text{mis} | \text{implemented}) \cdot t(\text{en} | \text{implemented}) \cdot \\ & t(\text{application} | \text{implemented}) \end{aligned}$$



# Alignment in Statistical MT

Another important module in statistical machine translation is the alignments offered by  $P(f | a, e)$ , which essentially is a **word-level mapping** between the input and output

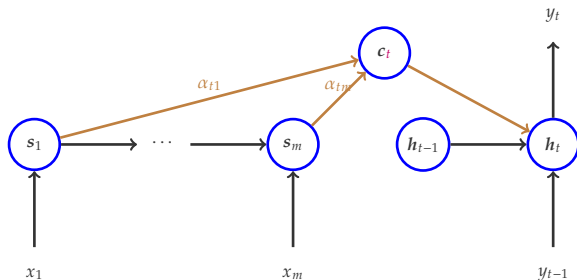


$$\begin{aligned} P(f | a, e) = & t(\text{Le} | \text{the}) \cdot t(\text{programme} | \text{program}) \cdot \\ & t(\text{a} | \text{has}) \cdot t(\text{ete} | \text{been}) \cdot \\ & t(\text{mis} | \text{implemented}) \cdot t(\text{en} | \text{implemented}) \cdot \\ & t(\text{application} | \text{implemented}) \end{aligned}$$

This is also something missed from the seq2seq models we discussed in the previous section.

# Attention Mechanism

With the attention mechanism [Bahdanau et al., 2015], the hidden states in the decoder are computed as



$$c_t = \sum_{j=1}^m \alpha_{tj} s_j \quad h_t = f(h_{t-1}, y_{t-1}, c_t) \quad (19)$$

where  $c_t$  is dynamically changing over time

# Attention Weights

In [Bahdanau et al., 2015], the attention weights are computed as

$$\alpha_{tj} = \frac{\exp(\phi(\mathbf{h}_{t-1}, \mathbf{s}_j))}{\sum_{j'=1}^m \exp(\phi(\mathbf{h}_{t-1}, \mathbf{s}_{j'}))}. \quad (20)$$

where  $\phi(\mathbf{h}_{t-1}, \mathbf{s}_j)$  is specifically defined as

$$\phi(\mathbf{h}_{t-1}, \mathbf{s}_j) = \mathbf{v}_a^\top \tanh(\mathbf{W}_{a0} \mathbf{h}_{t-1} + \mathbf{W}_{ai} \mathbf{s}_j) \quad (21)$$

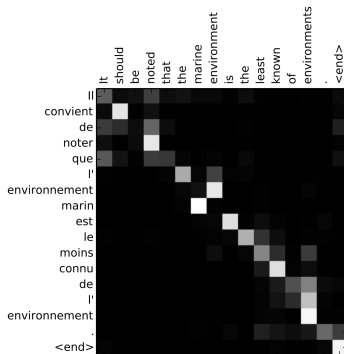
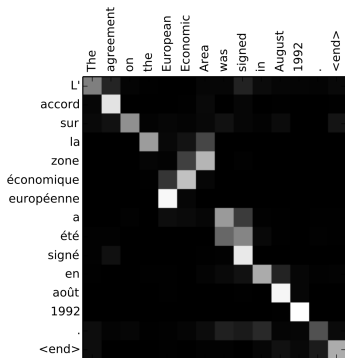
with parameters  $\mathbf{W}_{a0}$ ,  $\mathbf{W}_{ai}$  and  $\mathbf{v}_a$ .

The softmax function in Equation 20 implies

$$\sum_{j=1}^m \alpha_{tj} = 1 \quad (22)$$

# Example: Attention Weights in Machine Translation

By visualizing the attention weights, we can see the (soft) alignments between the source sentence and the target translation<sup>2</sup>



<sup>2</sup>The examples are from [Bahdanau et al., 2015], in which the translation is from English to French.



Bahdanau, D., Cho, K., and Bengio, Y. (2015).  
Neural machine translation by jointly learning to align and translate.  
In *ICLR*.



Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014).  
Learning phrase representations using rnn encoder-decoder for statistical machine translation.  
*arXiv preprint arXiv:1406.1078*.



Sutskever, I., Vinyals, O., and Le, Q. V. (2014).  
Sequence to sequence learning with neural networks.  
In *Advances in neural information processing systems*, pages 3104–3112.