

CS5242 : Neural Networks and Deep Learning

Lecture 11: Recurrent Neural Networks Applications

Semester 1 2021/22

Xavier Bresson

<https://twitter.com/xbresson>

Department of Computer Science
National University of Singapore (NUS)



Outline

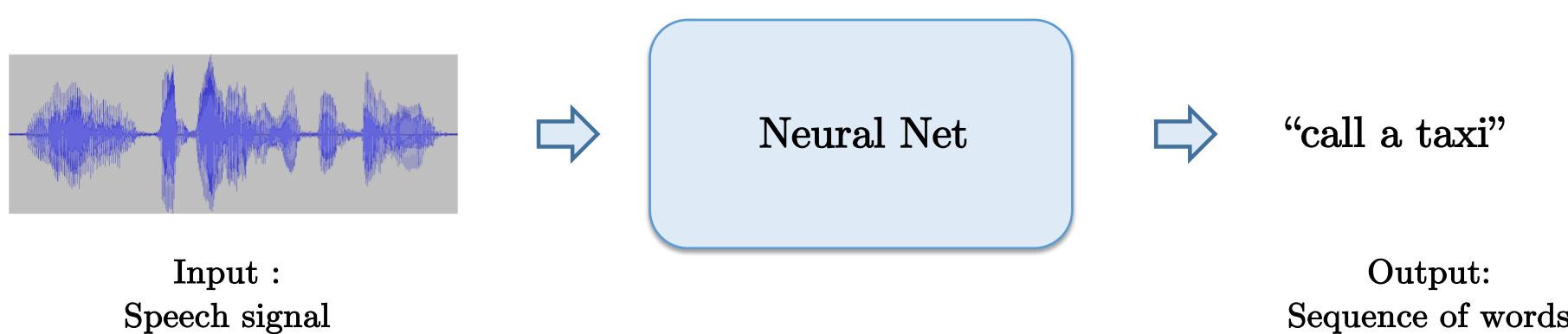
- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - Inference
 - Training

Outline

- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - Inference
 - Training

Speech recognition

- **Speech-to-text** is a fundamental task :
 - **Virtual assistants** : Apple Siri, Amazon Alexa, Microsoft Cortana, Google Assistant require to change the speech signal into a sequence of words.
 - Beyond **keyboard** : Speech recognition in-place of keyboard (when error is less than 1%).
 - **Current performances:** 5.1% error by Microsoft Research, as good as human transcript.

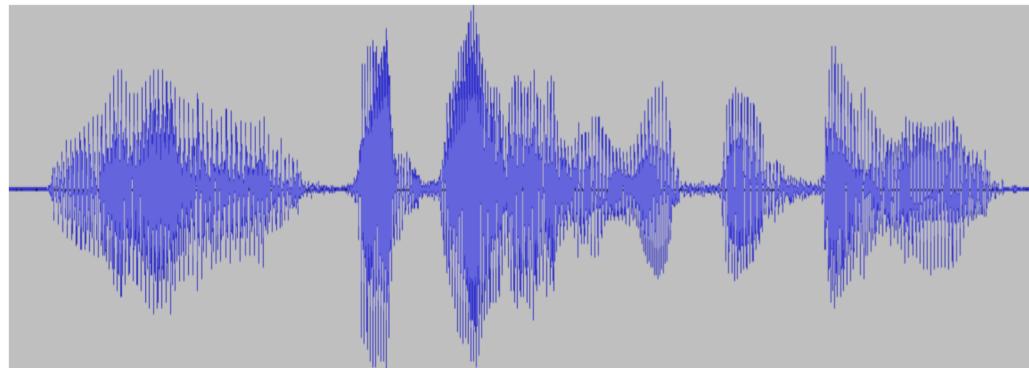


Outline

- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - Inference
 - Training

Spectrogram

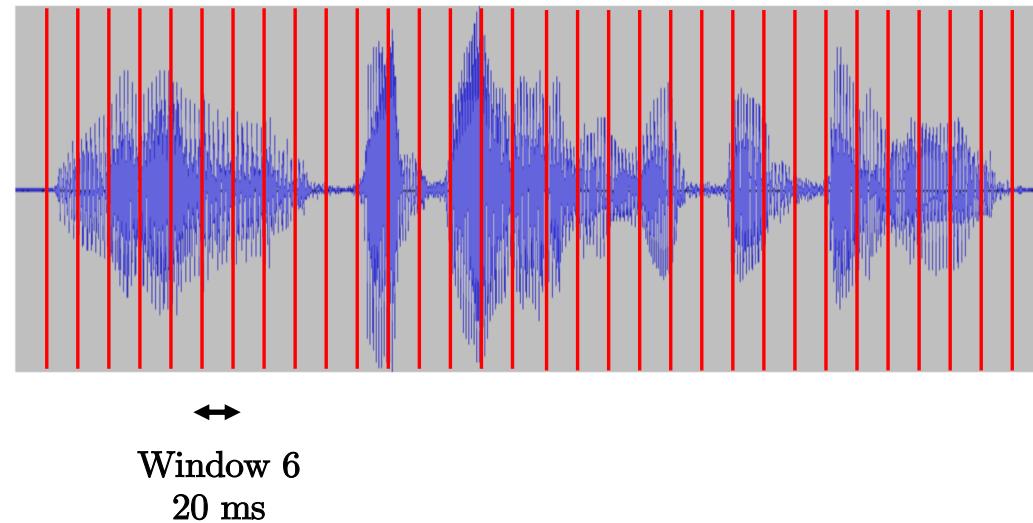
- We pre-process the raw audio signal to extract temporal and frequencial information.



Raw audio signal

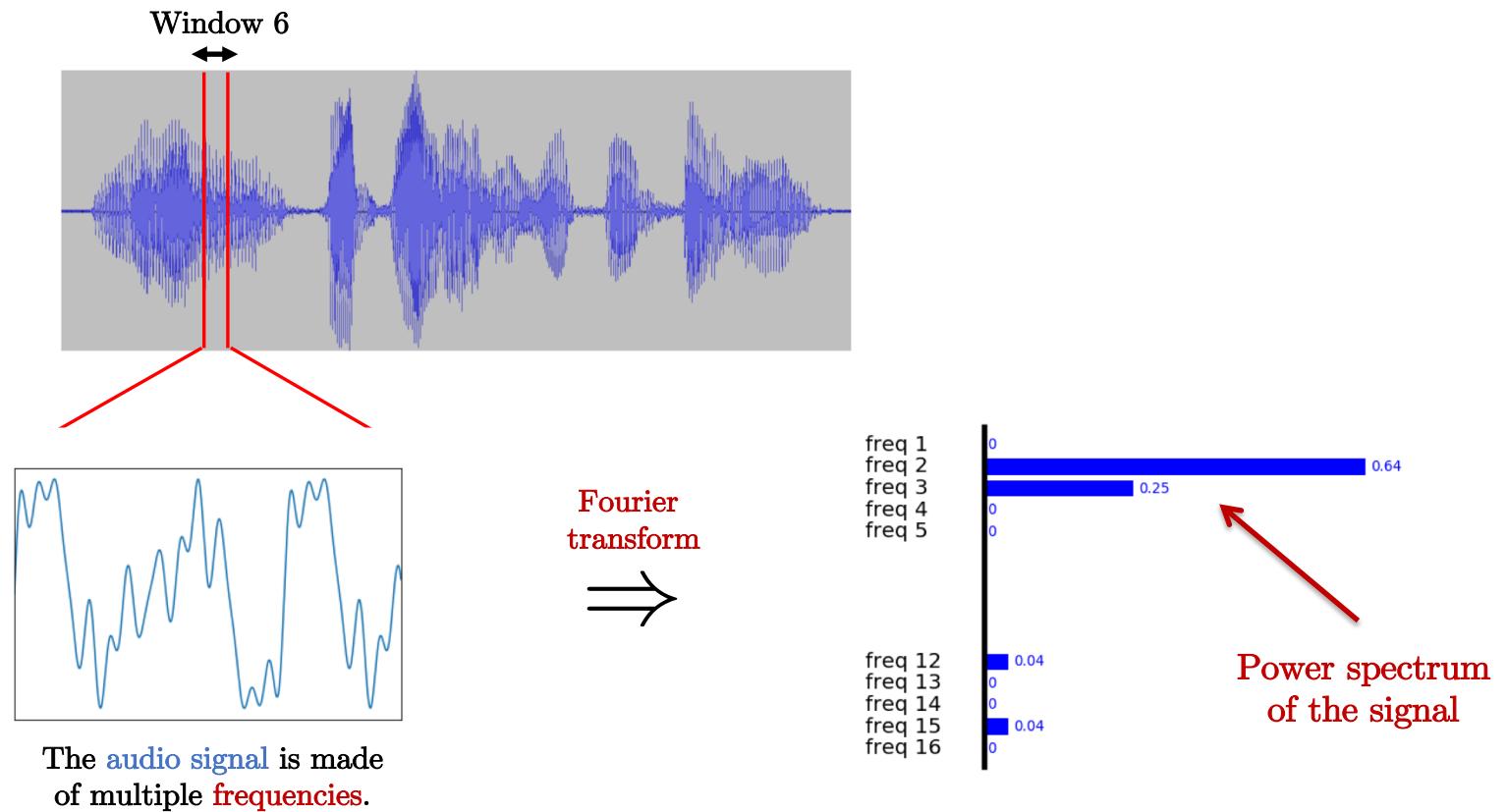
Spectrogram

- First, we cut the raw audio signal into small windows of 20 ms :



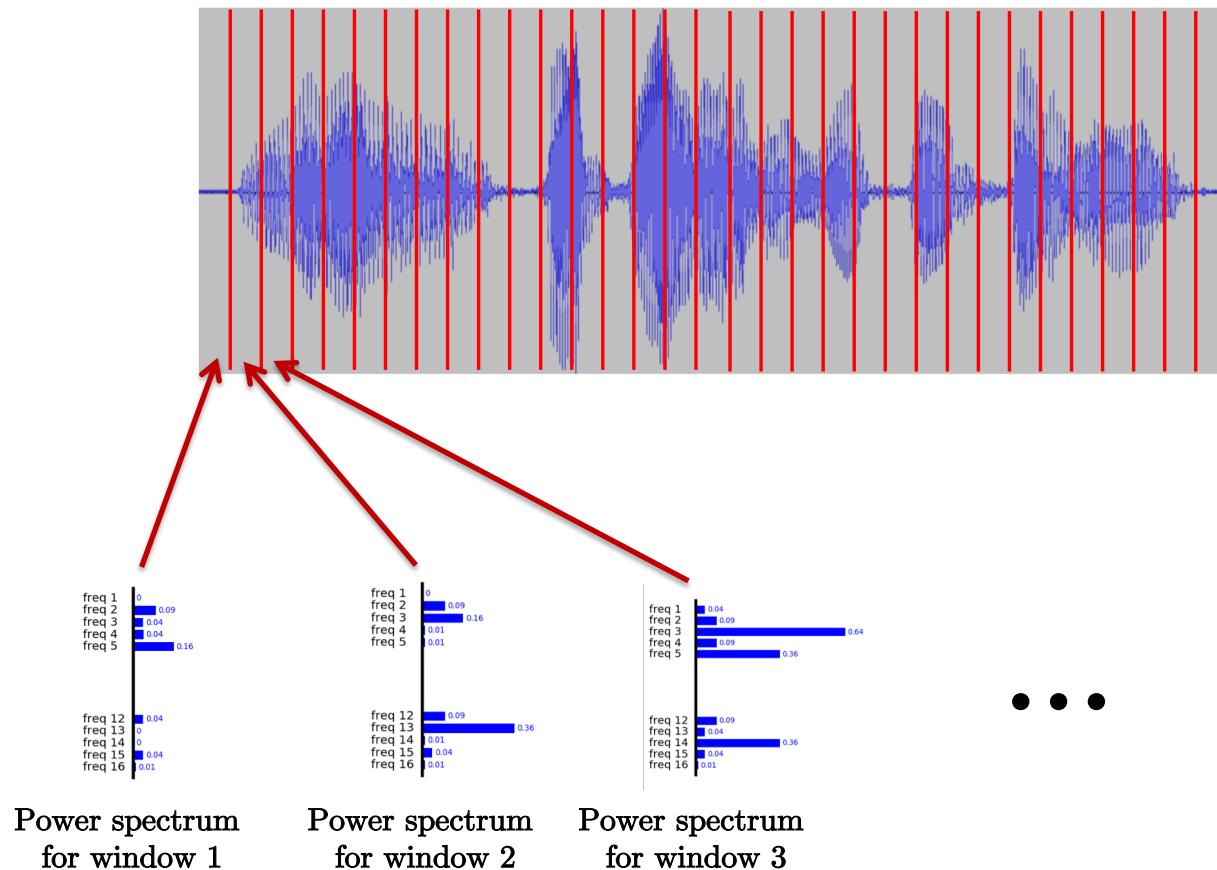
Spectrogram

- We compute the power spectrum of the signal (Fourier transform) for each window :



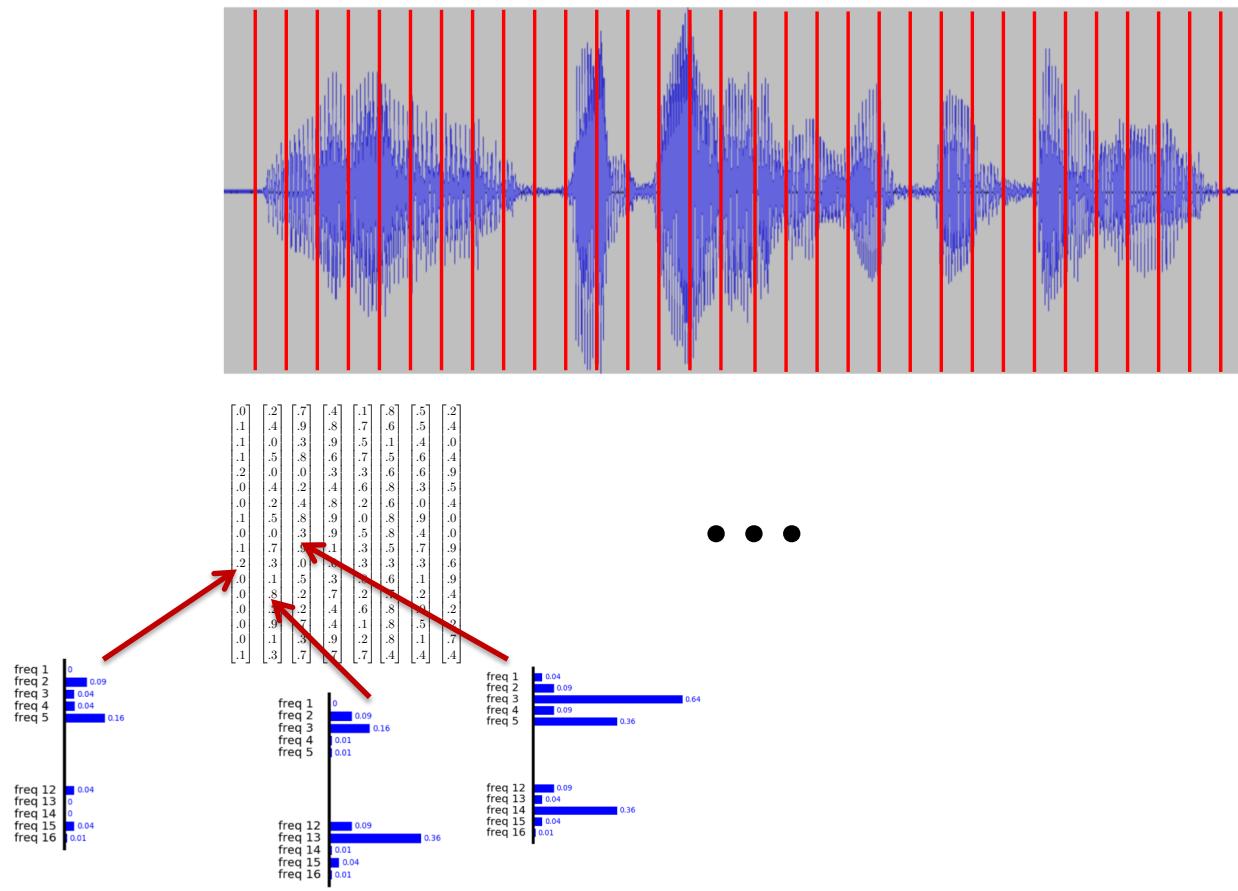
Spectrogram

- We compute the power spectrum for all 20ms windows :



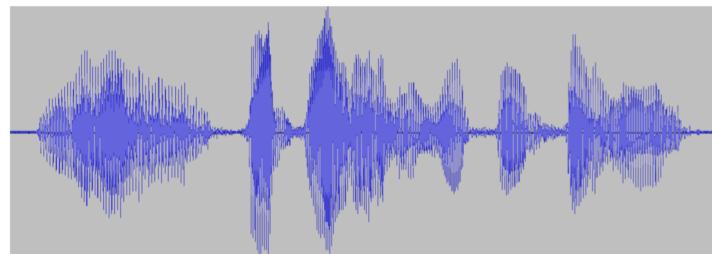
Spectrogram

- We compute the power spectrum of all 20ms windows :

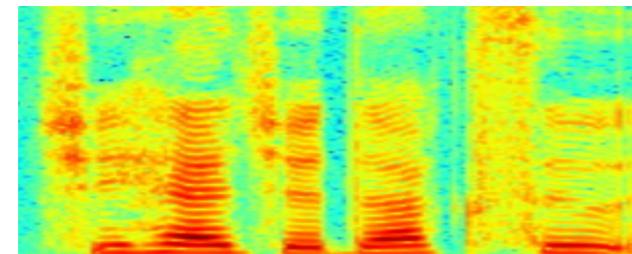


Spectrogram

- The **raw audio signal** has been converted into a **sequence of vectors** (size is between 100-200).
- This sequence of vectors is called the **spectrogram**.
- It can be computed very **quickly** with **FFT**.



Raw audio signal



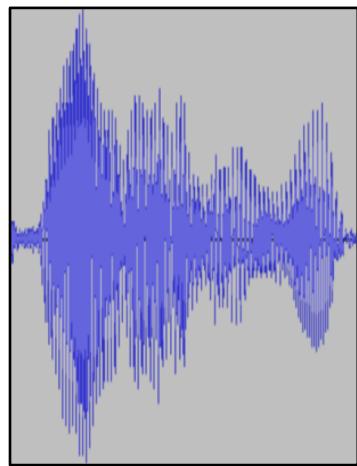
Spectrogram

Outline

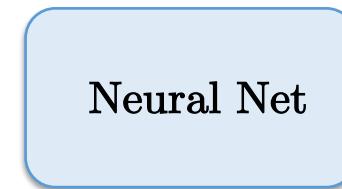
- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - Inference
 - Training

Inference

- Inference process :
 - Feed the sequence of vectors (spectrogram) to a trained RNN :



[.0]	[.2]	[.7]	[.4]	[.1]	[.8]	[.5]	[.2]
.1	.4	.9	.8	.7	.6	.5	.4
.1	.0	.3	.9	.5	.1	.4	.0
.1	.5	.8	.6	.7	.5	.6	.4
.2	.0	.0	.3	.3	.6	.6	.9
.0	.4	.2	.4	.6	.8	.3	.5
.0	.2	.4	.8	.2	.6	.0	.4
.1	.5	.8	.9	.0	.8	.9	.0
.0	.0	.3	.9	.5	.8	.4	.0
.1	.7	.9	.1	.3	.5	.7	.9
.2	.3	.0	.0	.3	.3	.3	.6
.0	.1	.5	.3	.8	.6	.1	.9
.0	.8	.2	.7	.2	.7	.2	.4
.0	.2	.2	.4	.6	.8	.0	.2
.0	.9	.7	.4	.1	.8	.5	.2
.0	.1	.3	.9	.2	.8	.1	.7
.1	.3	.7	.7	.7	.4	.4	.4



“Hello”

Input :
Audio signal

Spectrogram

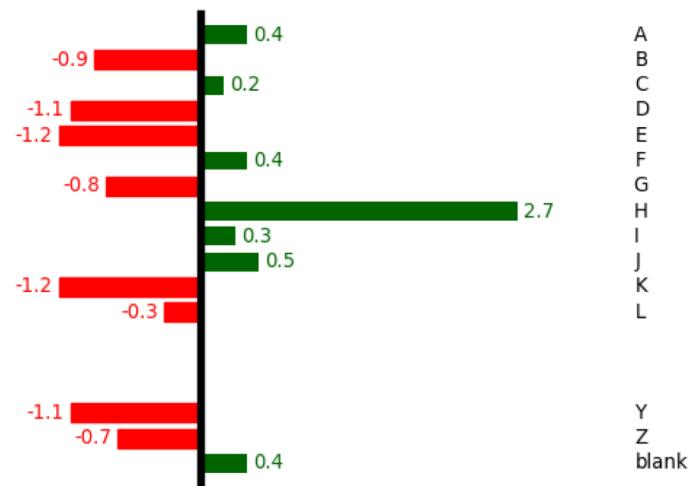
Output :
Sequence of words

Inference

- Dictionary/vocabulary is defined as follows :

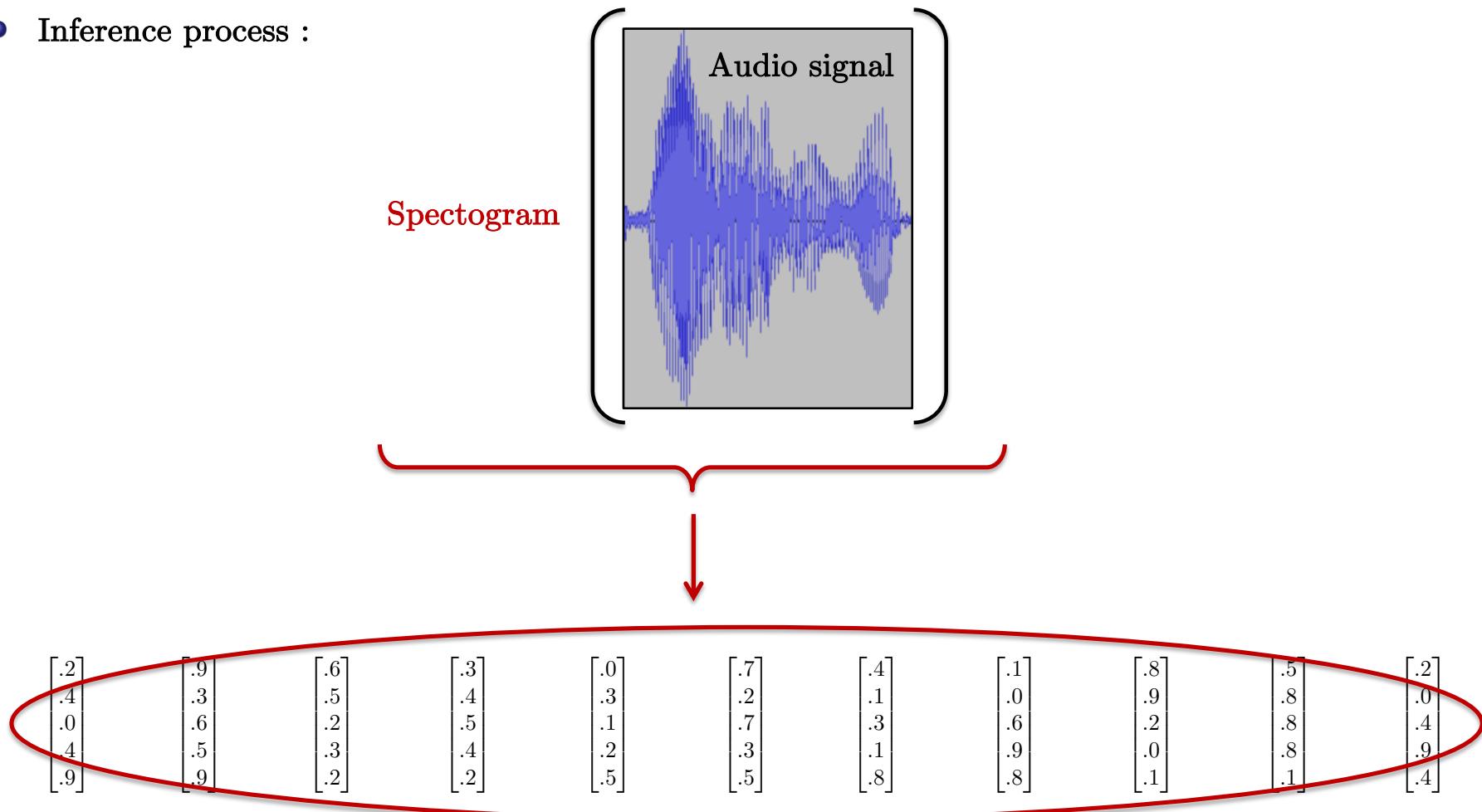
- 27 classes :

- The 26 letters of the alphabet
- Blank symbol
- Space Symbol



Inference

- Inference process :



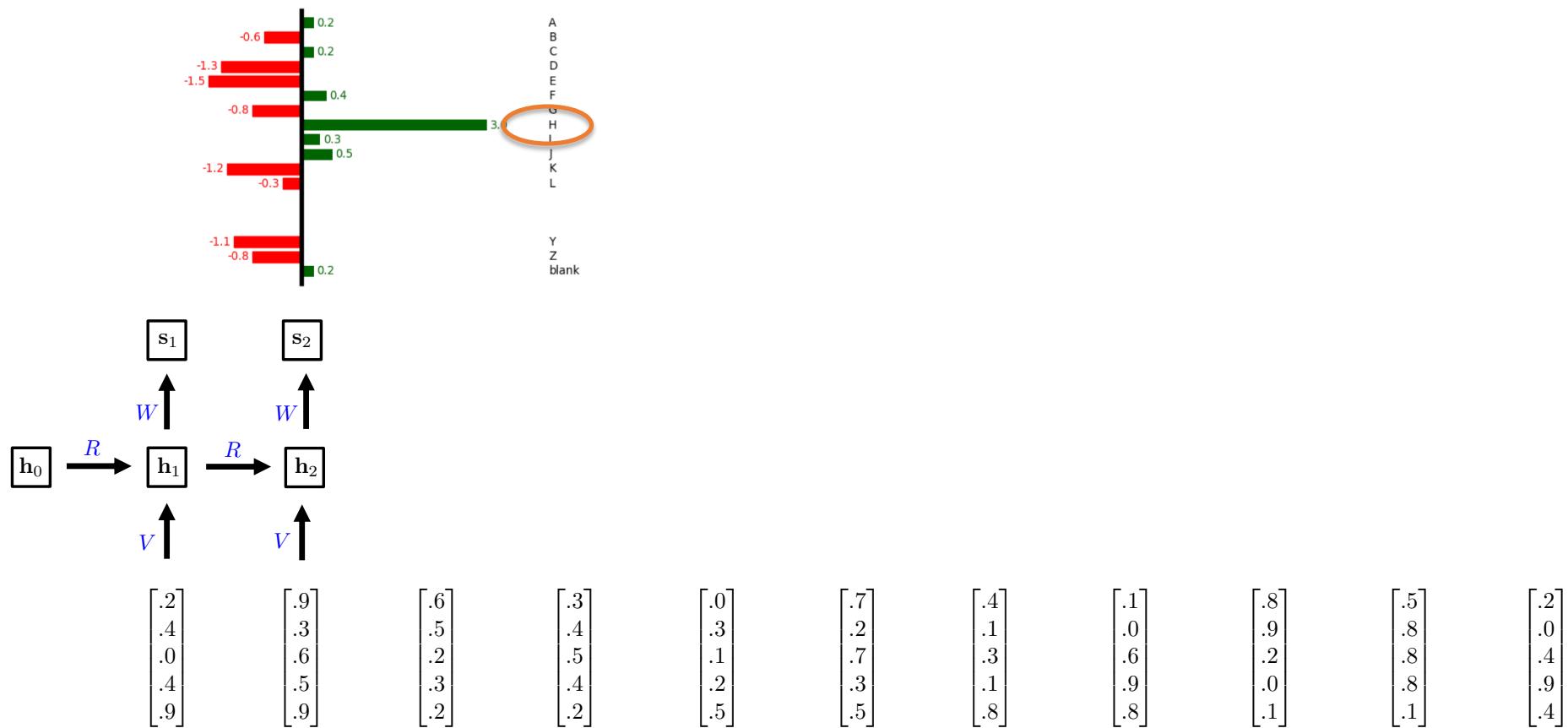
Inference

- Inference process : Softmax the letter corresponding to the input spectrogram vector.



Inference

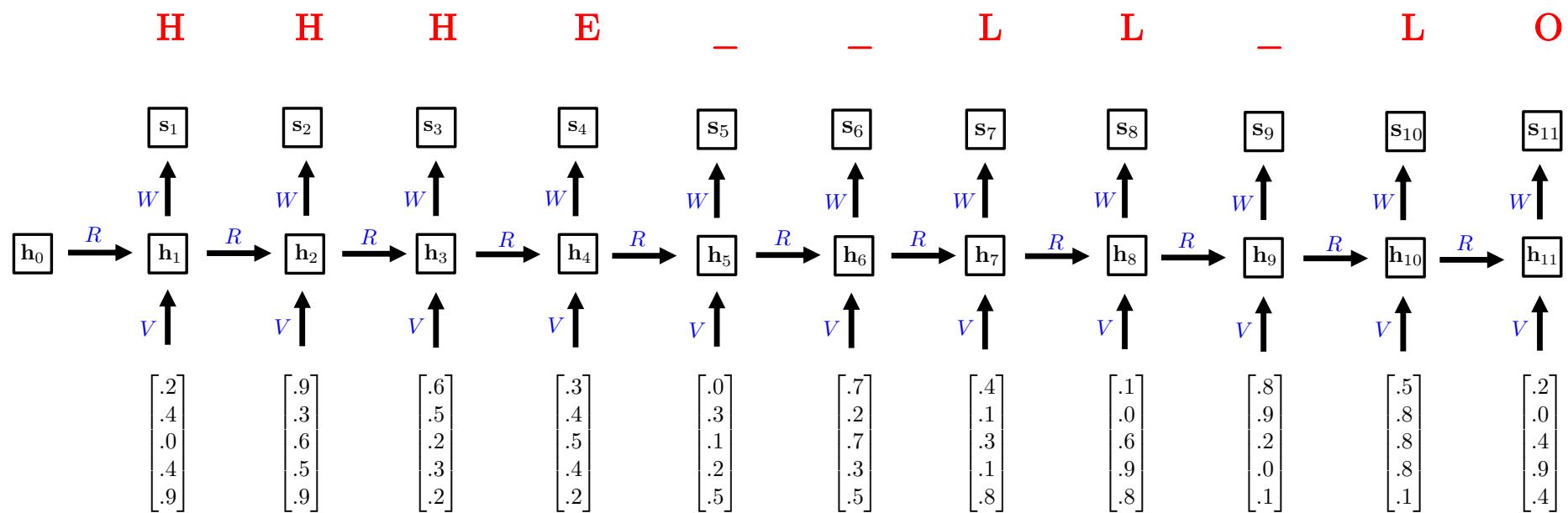
- Inference process : Softmax the letter corresponding to the input spectrogram vector.



Inference

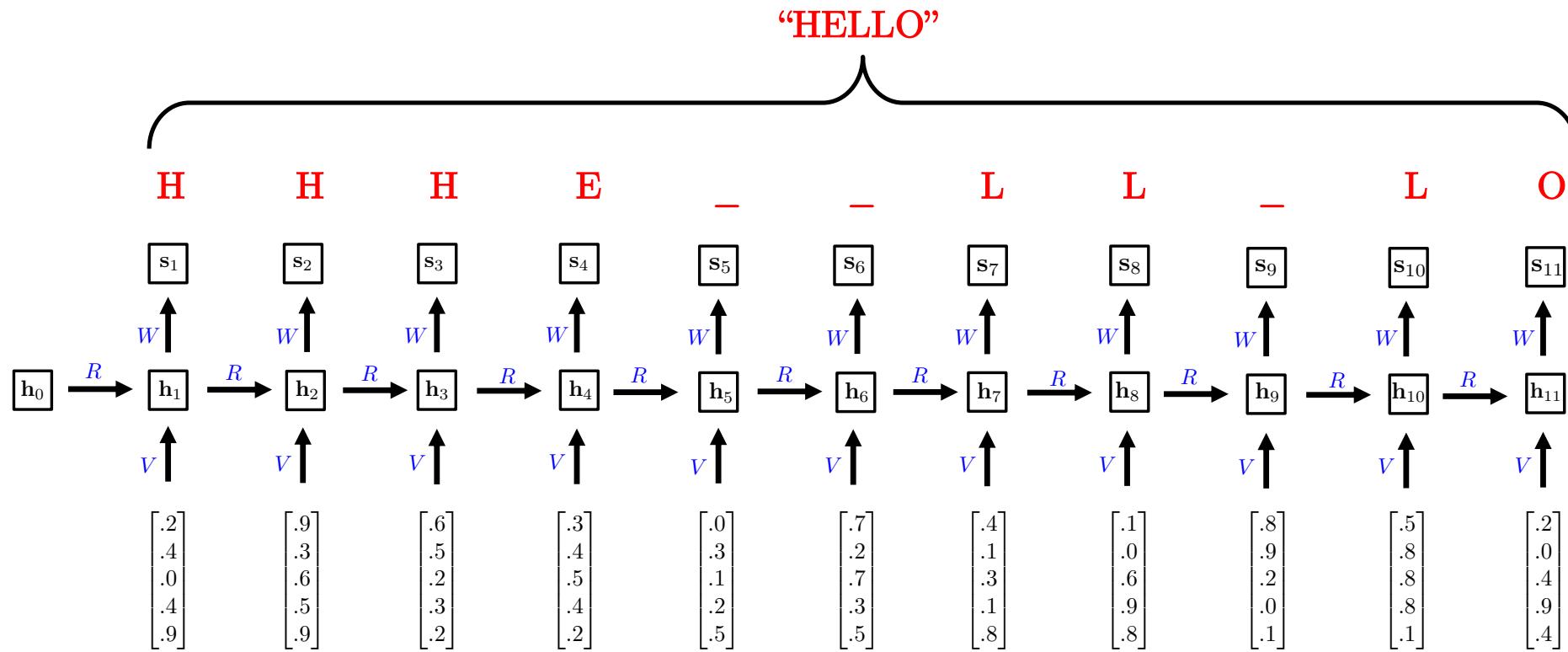
- Inference process :

- Forward is done.



Inference

- Inference process :
 - The last layer at the very top is **human hand-crafted** :
 - Connectionist Temporal Classification (CTC)



Outline

- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - **Introduction**
 - Architecture
 - Inference
 - Training

Neural machine translation

- Google Translate switched to Neural Machine Translation in Nov. 2016.
 - This led to a big improvement in quality !
- Harry Potter and The Chamber of Secrets : English ⇒ French ⇒ English

The screenshot shows the Google Translate interface with two main sections. The left section translates from English to French, and the right section translates from French back to English. A red circle highlights the word "squelched" in the English input, which is correctly translated to "s'éloignait dans". Another red circle highlights the word "walked away" in the final English output, which is also correct.

English: As Harry **squelched** long the deserted corridor he came across somebody who looked **just** as preoccupied as he was. Nearly Headless Nick, the ghost of Gryffindor Tower, was staring morosely out of a window, muttering under his breath, "... don't fulfill their requirements . . . half an inch, if that . . .".

French: Alors que Harry **s'éloignait dans** le couloir désert, il tomba sur quelqu'un qui semblait aussi préoccupé que lui. Presque sans tête Nick, le fantôme de la Tour de Gryffondor, regardait moralement par une fenêtre, murmurant dans sa barbe, "... ne remplit pas leurs exigences . . . un demi-pouce, si cela . . .".

English: "Hello, Nick," said Harry.

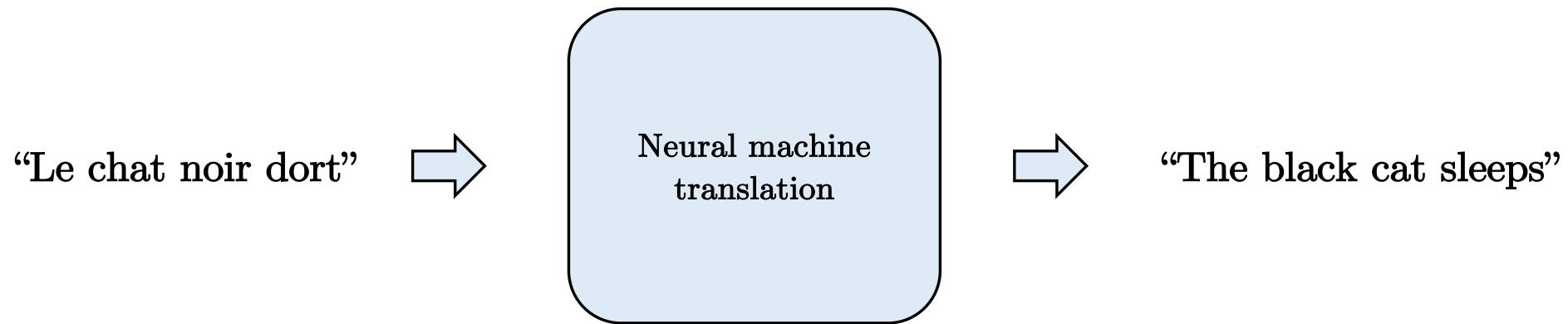
French: "Bonjour, Nick," dit Harry.

English: As Harry **walked away** into the deserted corridor, he came across someone who seemed **as** preoccupied as he was. Almost headless Nick, the ghost of the Gryffindor Tower, was looking morally through a window, whispering in his beard, "... does not meet their demands . . . half an inch, if that . . .".

Still far from professional translators, but some translators use it as first draft !

Neural machine translation

- The challenges of translation :
 - The words can appear in different order: **Alignment problem**.
 - Input and output sentences do not necessarily have the same number of words: **Length problem**.
 - Input and output sentences do not have the same number of words in their dictionary: **Vocabulary problem**.



Outline

- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - **Architecture**
 - Inference
 - Training

Architecture

- A **translation** system is composed of **two recurrent neural networks** :
 - **Encoder** network :
 - The encoder network will **summarize** the input sequence in language A with a **vector**.
 - **Decoder** network :
 - The decoder network will take the encoding vector representing the input sequence in language A and will decode it to an output **sequence** in language B.
- Full **end-to-end** machine learning systems :
 - No conceptual prior on languages – opposite to Chomsky's paradigm.
 - Big change in the NLP community !

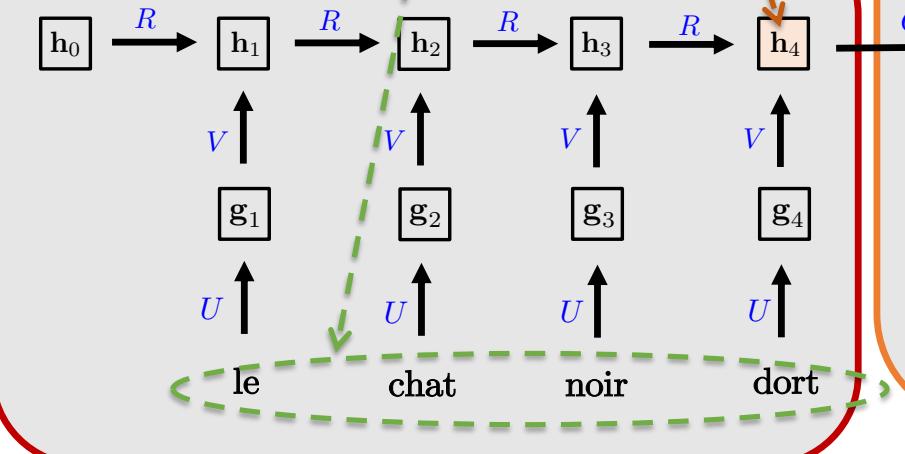
Architecture

Step 1: Take a sentence in French and encode it into a vector h of size 200.

Step 2: Take the vector created by the encoder and use it to generate a sentence in English.

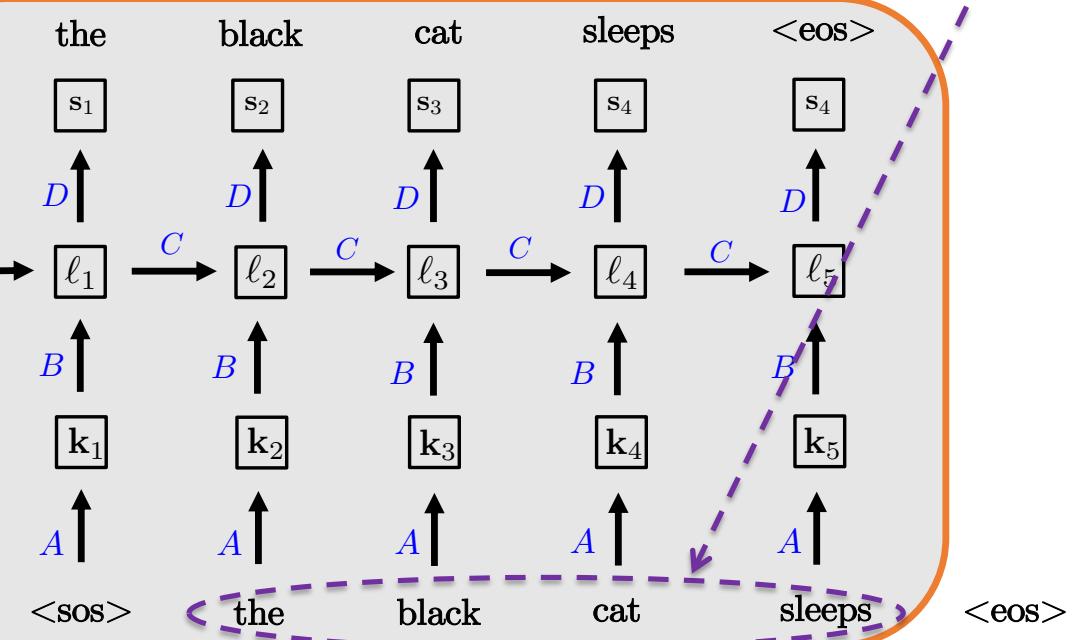
Encoder Network

Internal parameters are: U, V, R .



Decoder Network

Internal parameters are: A, B, C, D .



Architecture

- Dictionaries : We have **two vocabularies** :
 - **French** vocabulary has **35,000** words.
 - Each French word is a **one-hot-vector** with 35,000 entries.
 - **English** vocabulary has **32,000** words.
 - Each English word is a **one-hot-vector** with 32,000 entries.

Outline

- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - **Inference**
 - Training

Inference

- The network has been **trained**, let us use it for inference.
- First step** of inference:
 - Encode** the input sequence with an **history vector**.

History/encoding
vector h



Encoder
network



“le chat noir dort”

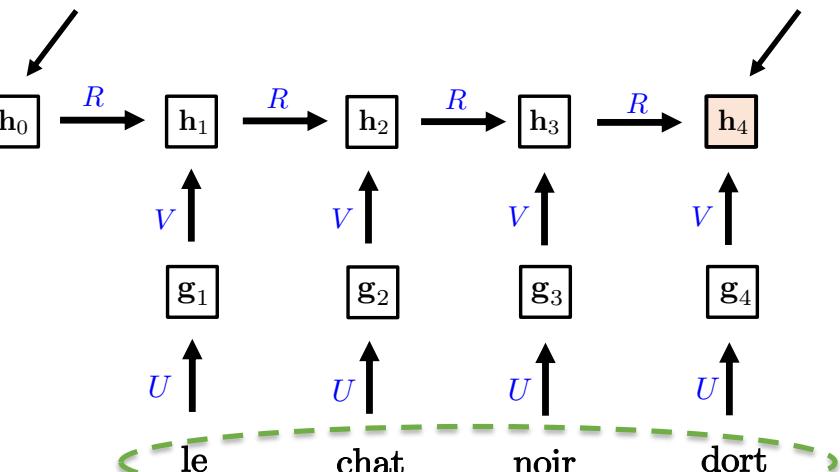
The hidden state is
initialized at zero.

The vector h_4 contains all
the **history of the sentence**.

Size = 200

Size = 200

One-hot-vector
Size = 35,000



French sentence that we want to translate.

Inference

- Second step of inference:
 - Decode the encoding/history vector into an output sequence.

History/encoding
vector h

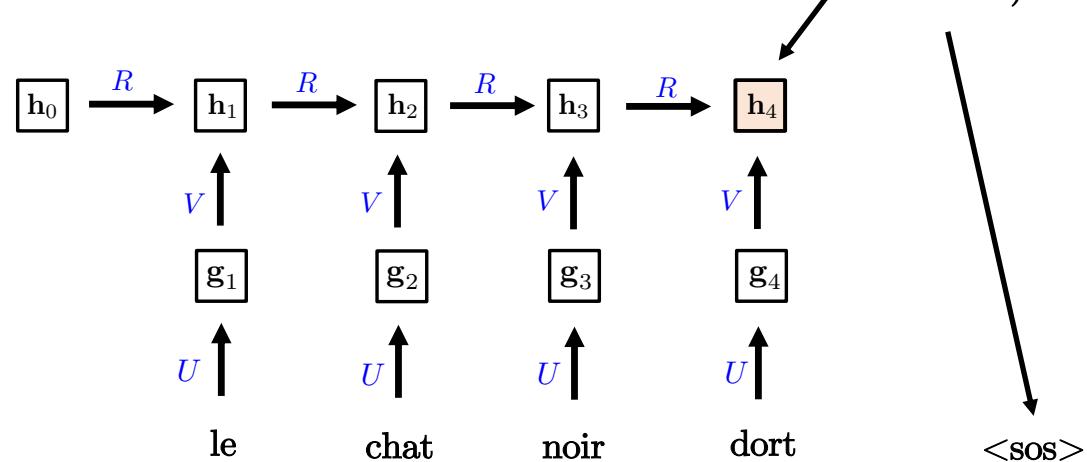


Decoder
network



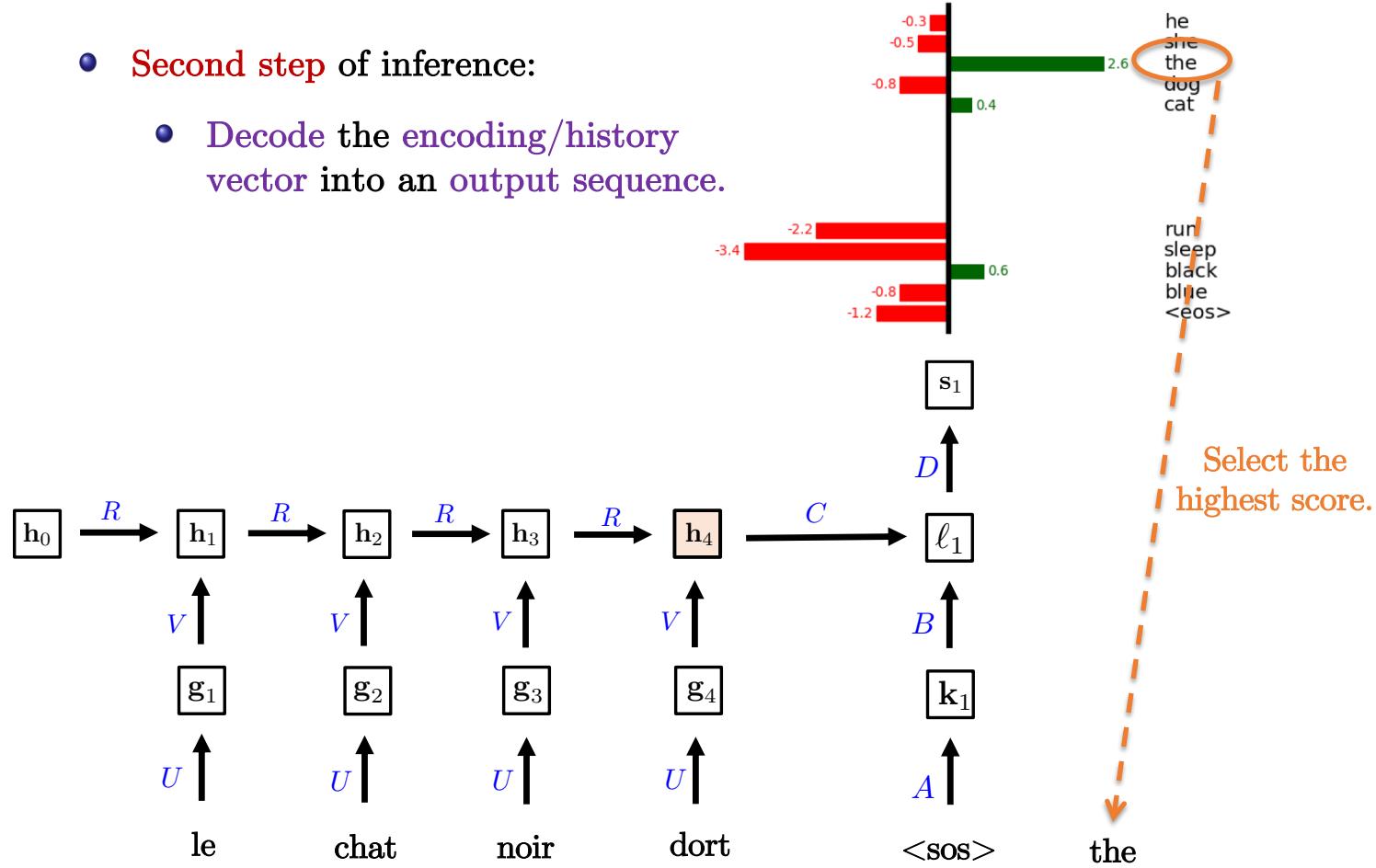
“the black cat sleeps”

We start the decoding
with vector h_4 and the
word <sos> (start of
sentence).



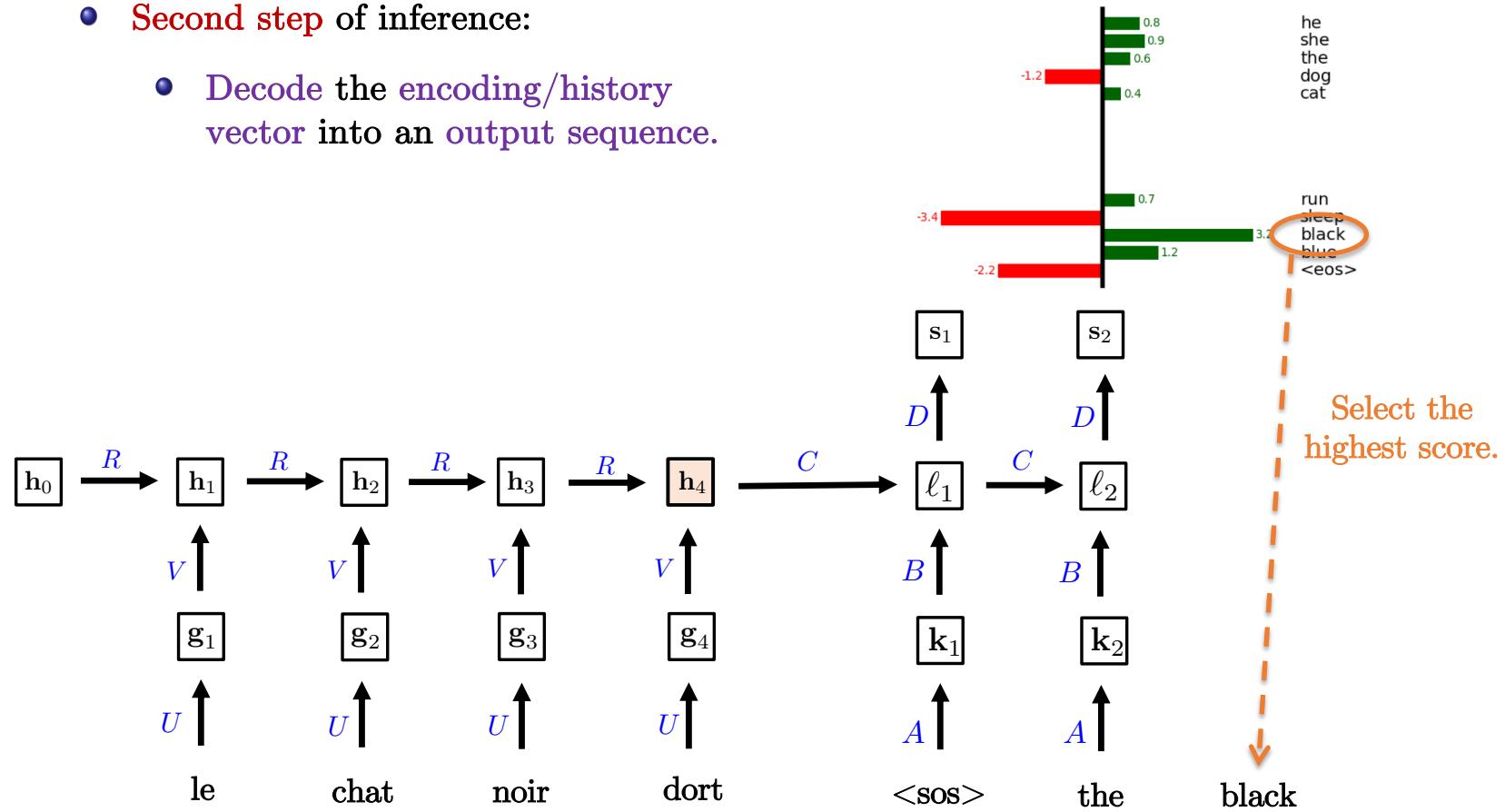
Inference

- Second step of inference:
 - Decode the encoding/history vector into an output sequence.



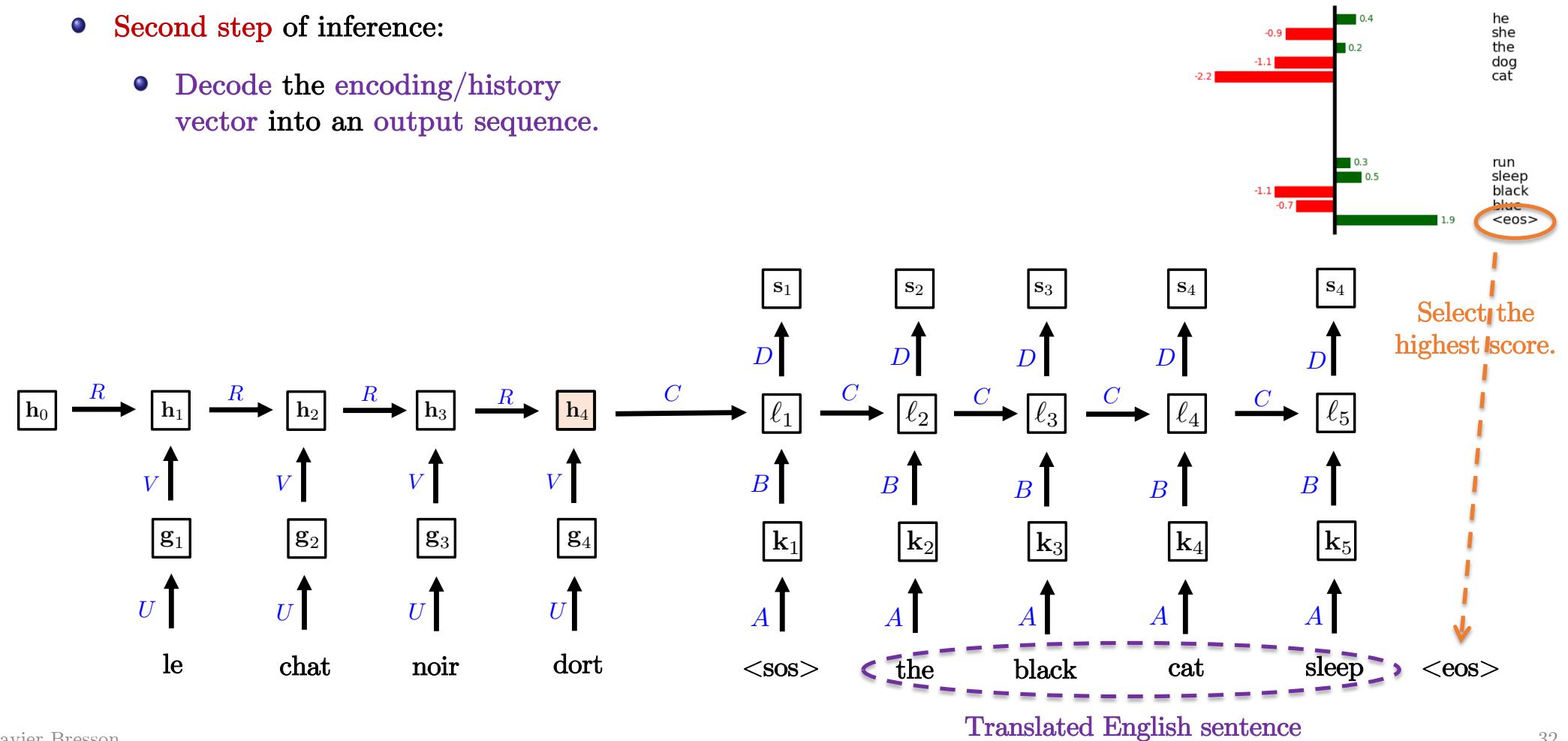
Inference

- Second step of inference:
 - Decode the encoding/history vector into an output sequence.



Inference

- Second step of inference:
 - Decode the encoding/history vector into an output sequence.



Outline

- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - Inference
 - Training

Training

- Let us **train** the network with **10 million French \Rightarrow English sentences** :

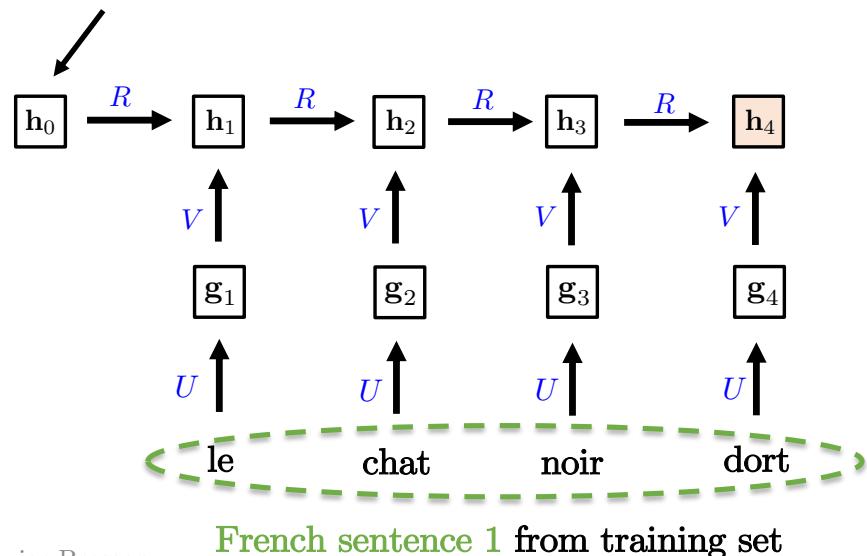
Sentence 1	le chat noir dort	<sos> the black cat sleep <eos>
Sentence 2	hier je suis aller a la plage	<sos> yesterday I went to the beach <eos>
Sentence 3	je m'appelle thomas	<sos> my name is thomas <eos>
.	.	.
Sentence 10,000,000	tu est un tres bon joueur de tennis	<sos> you are a very good tennis player <eos>

Training

- Let us feed the 1st French \Rightarrow English sentence :

< Labels: the, black, cat, sleeps, <eos> >
Translated English sentence 1

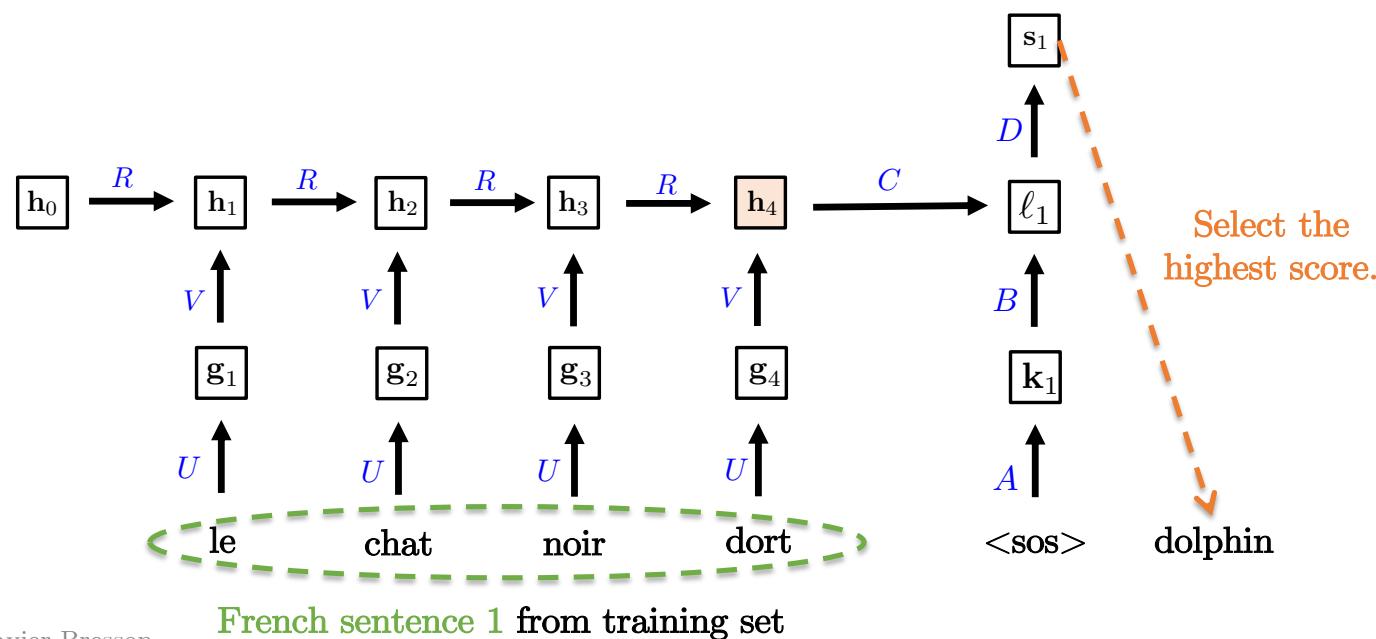
The hidden state is initialized at zero.



Training

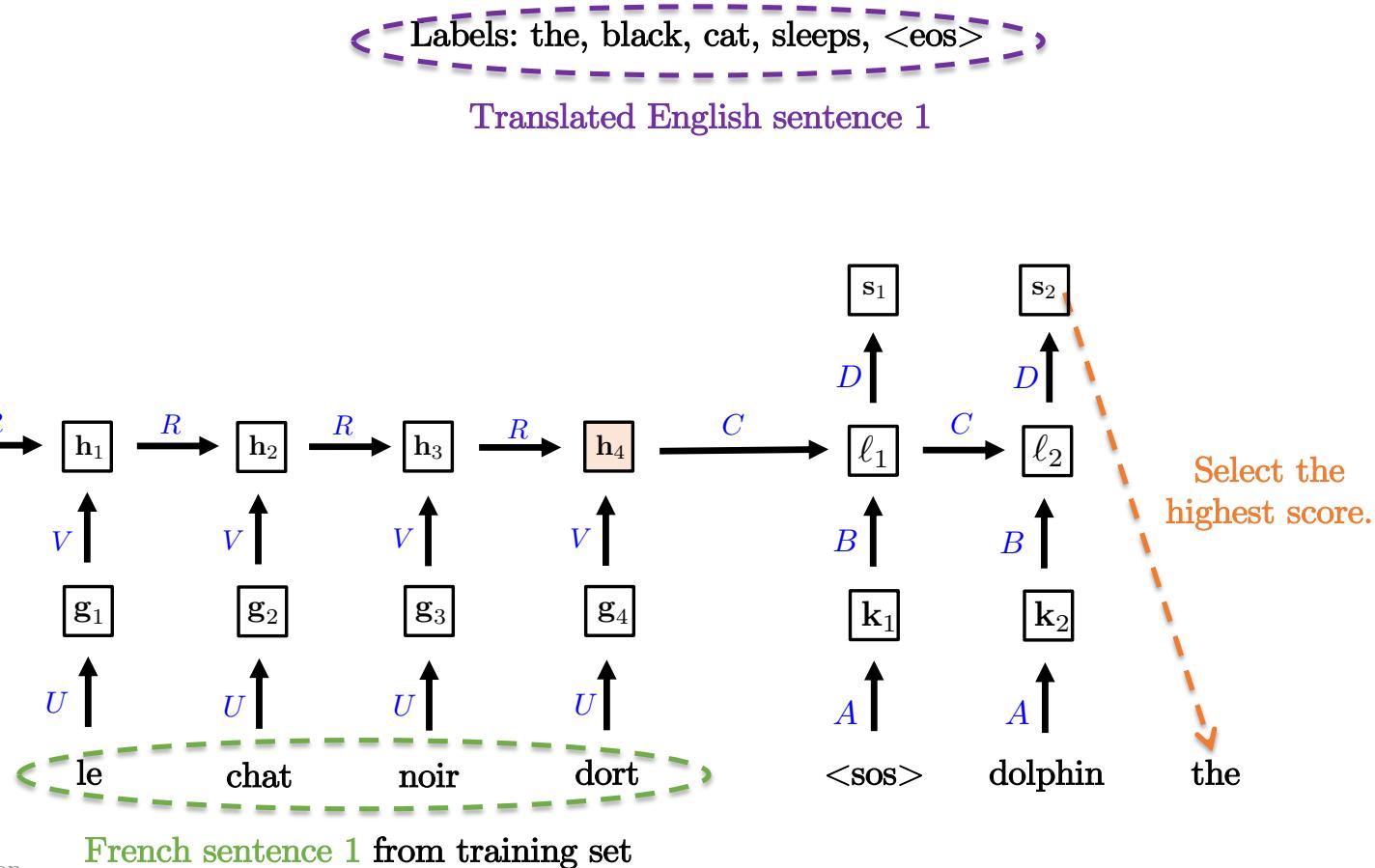
- Let us feed the 1st French \Rightarrow English sentence :

Labels: the, black, cat, sleeps, <eos>
 Translated English sentence 1



Training

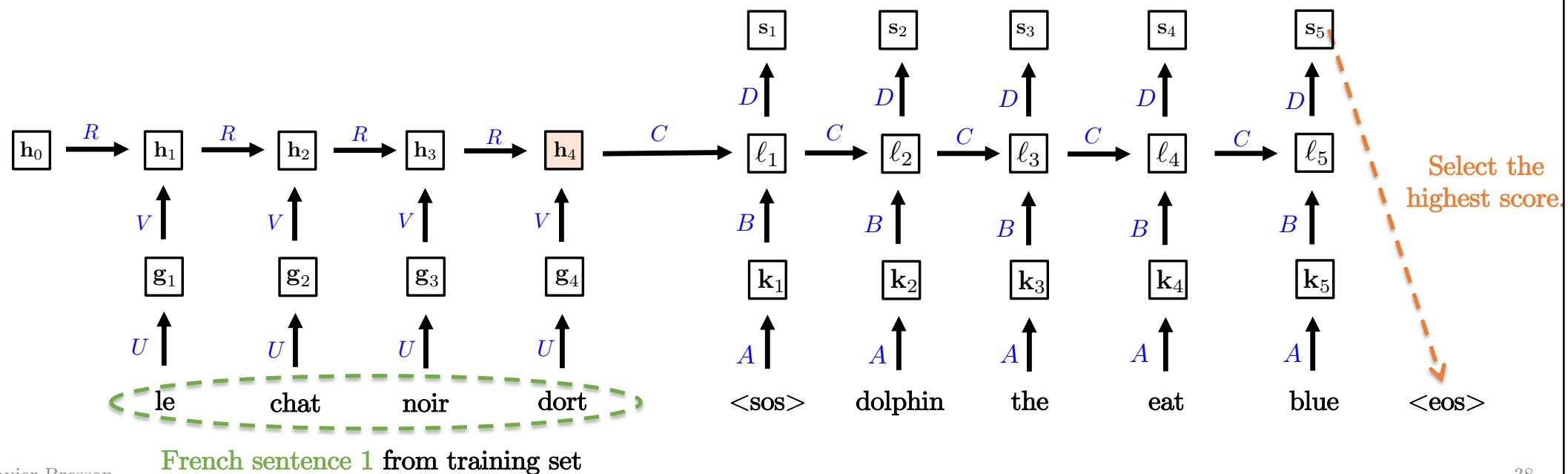
- Let us feed the 1st French \Rightarrow English sentence :



Training

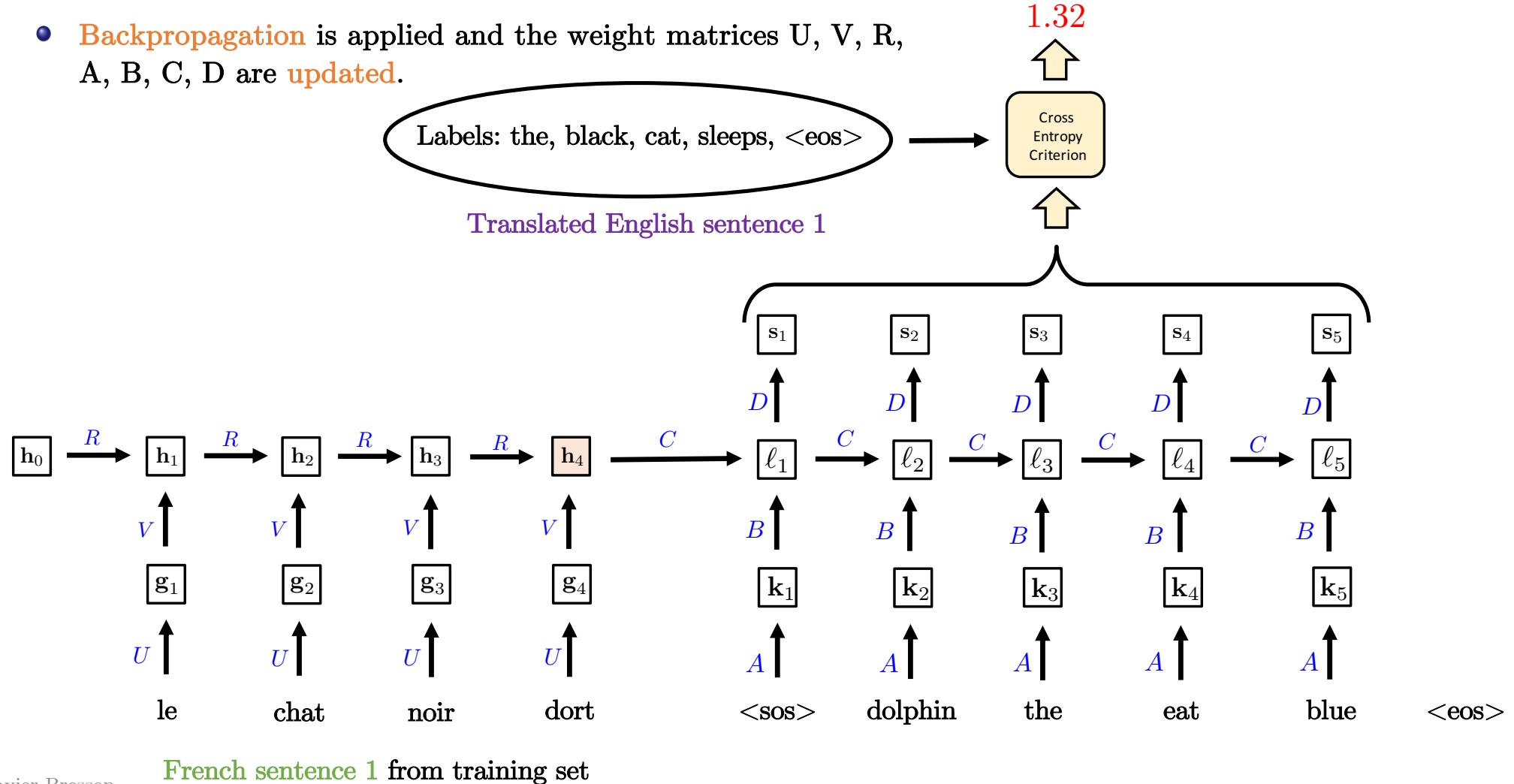
- Let us feed the 1st French \Rightarrow English sentence :

< Labels: the, black, cat, sleeps, <eos> >
 Translated English sentence 1



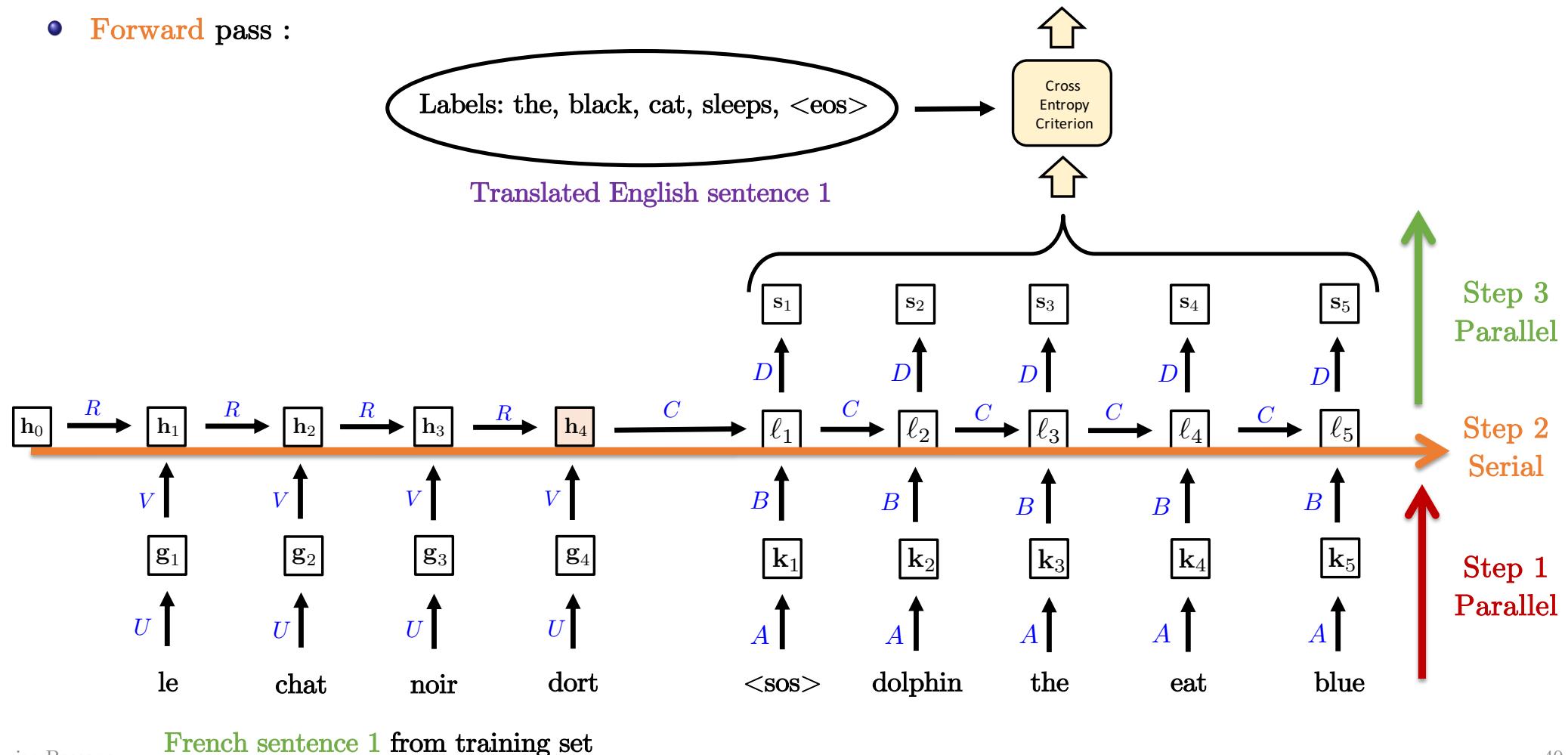
Training

- Backpropagation is applied and the weight matrices U , V , R , A , B , C , D are updated.



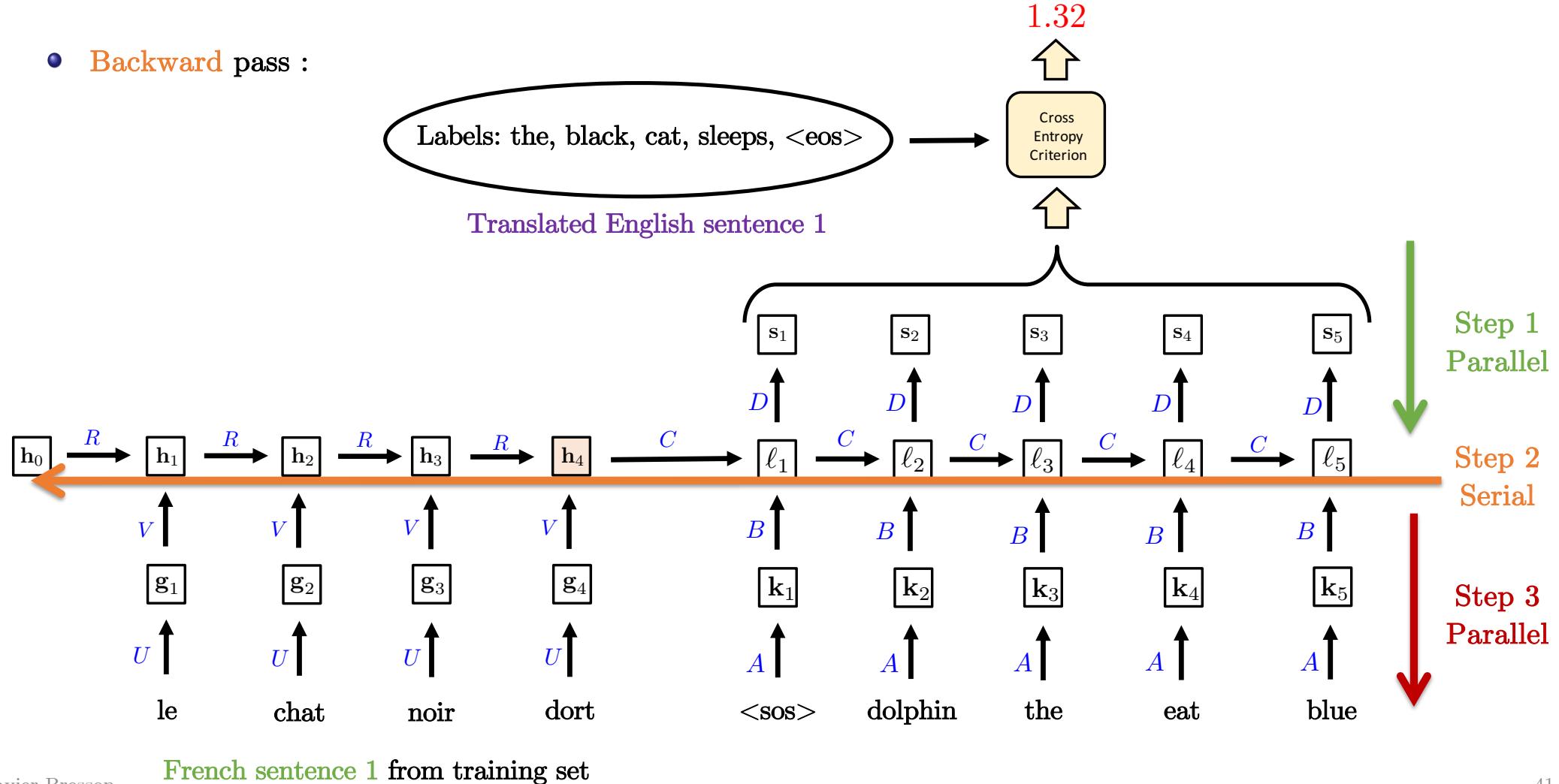
Training

- Forward pass :



Training

- Backward pass :



Training

- State-of-the-art machine translation machine :
 - Stack 10 layers of bi-directional LSTM :

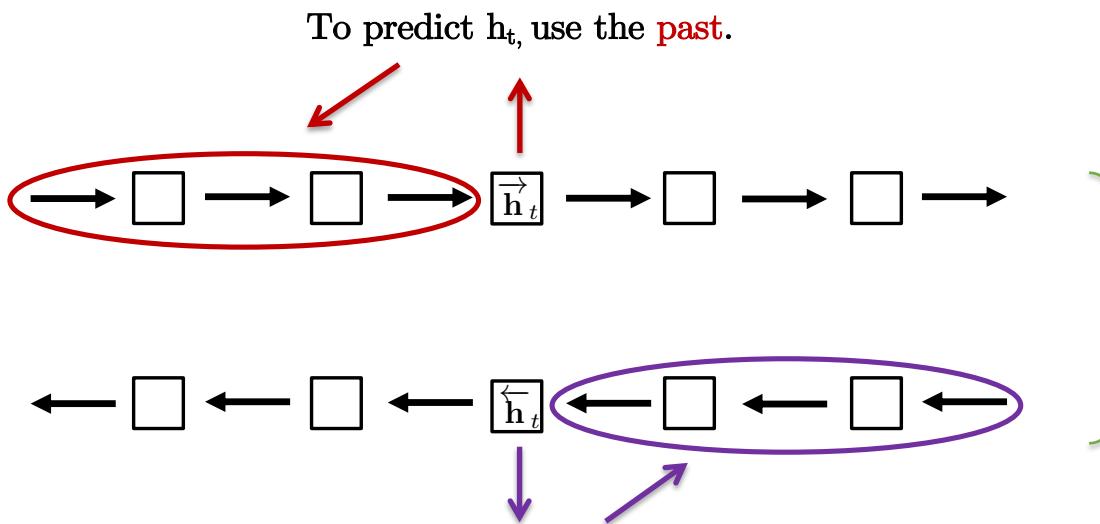
```
net = nn.LSTM( hidden_size , hidden_size , bidirectional=True )
```

- Use an attention mechanism.

Training

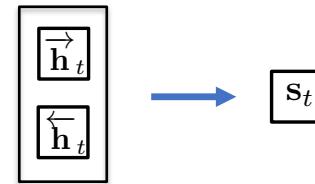
- Bi-directional LSTM :

Causal signal w.r.t.
increasing time :

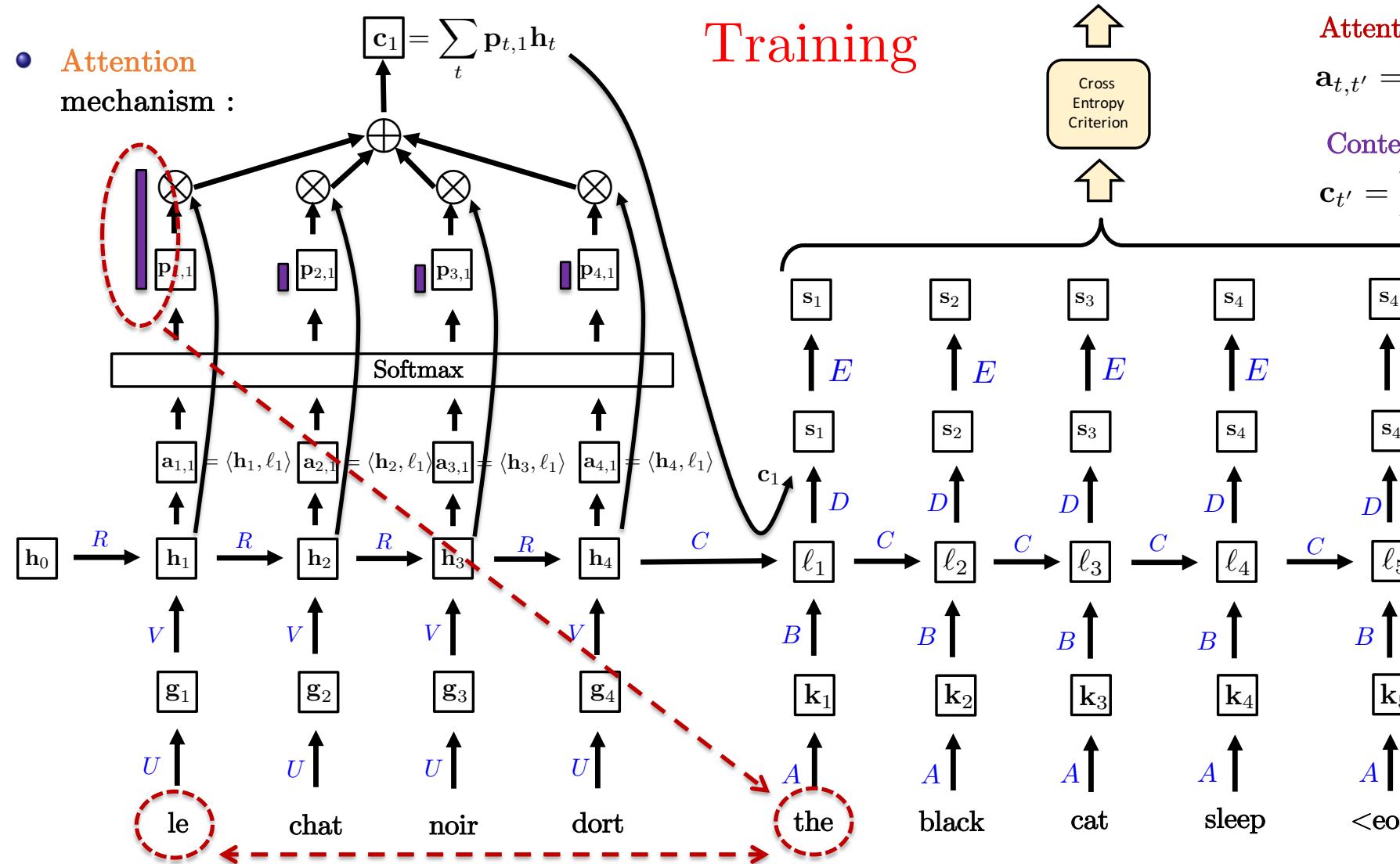


Causal signal w.r.t.
decreasing time :

Combine both
hidden vectors h to
better solve NLP
tasks (e.g. Q&A)



- Attention mechanism :

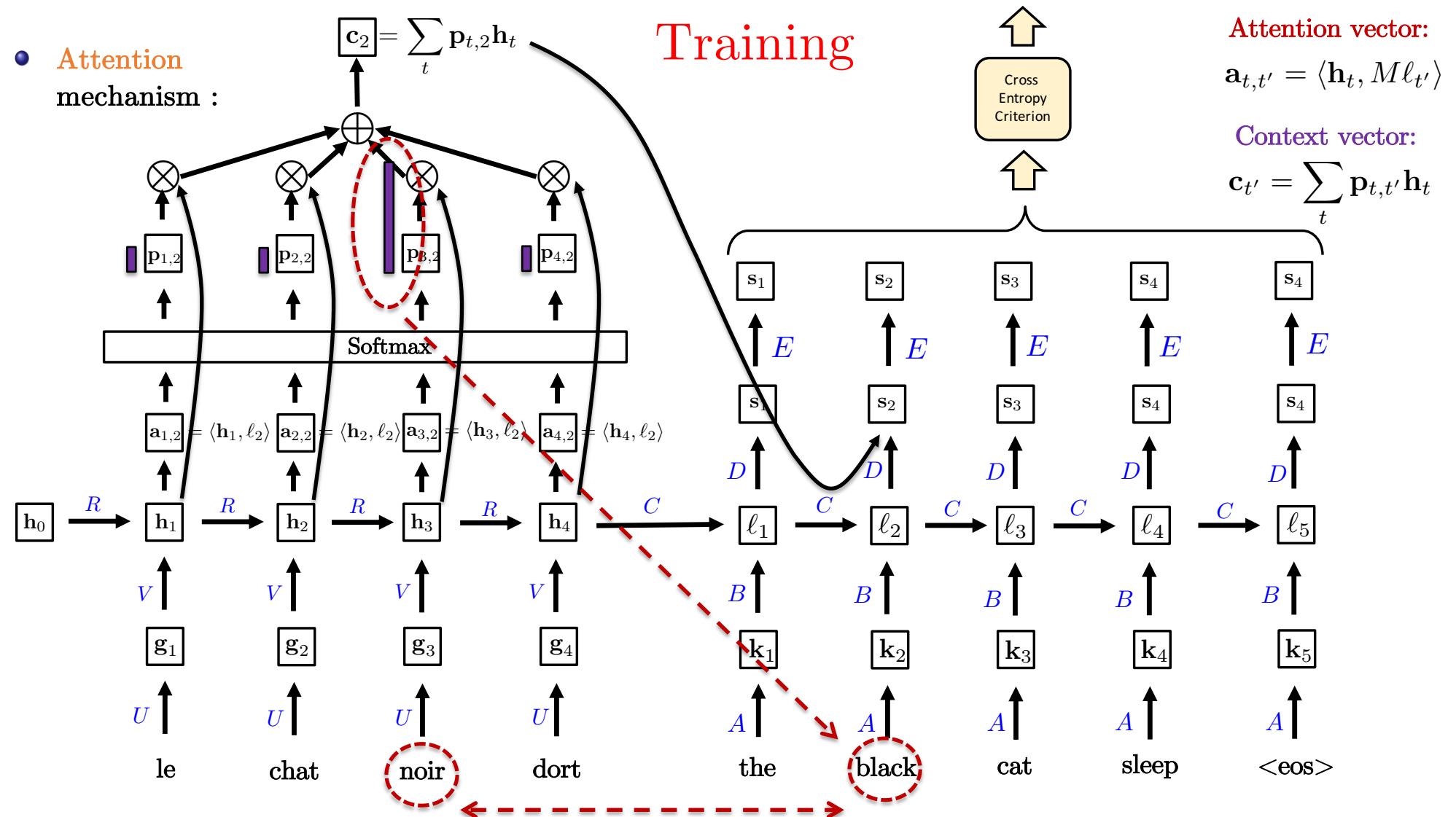


Attention vector:
 $\mathbf{a}_{t,t'} = \langle \mathbf{h}_t, M\ell_{t'} \rangle$

Context vector:

$$\mathbf{c}_{t'} = \sum_t \mathbf{p}_{t,t'} \mathbf{h}_t$$

- Attention mechanism :



End of RNNs ?

- Attention neural networks are taking over NLP !
 - 2019 breakthrough in AI

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Ilia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.



Questions?