

CS5242 : Neural Networks and Deep Learning

Lecture 8b: Recurrent Neural Networks Applications

Semester 1 2022/23

Ai Xin

aixin@comp.nus.edu.sg

Department of Computer Science
National University of Singapore (NUS)



Outline

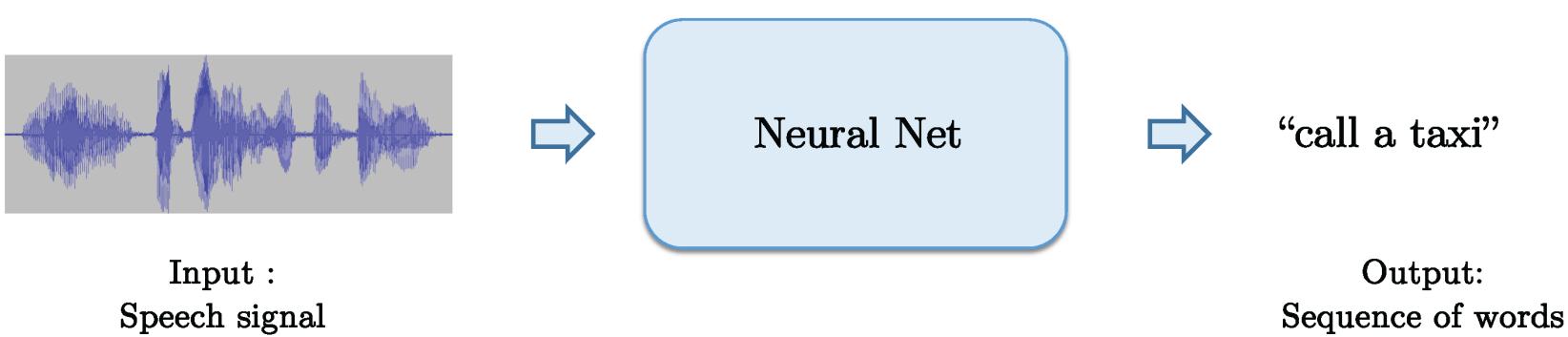
- Speech recognition
 - Introduction
 - Spectrogram
 - Inference
- Machine Translation
 - Introduction
 - Architecture
 - Inference
 - Training

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.

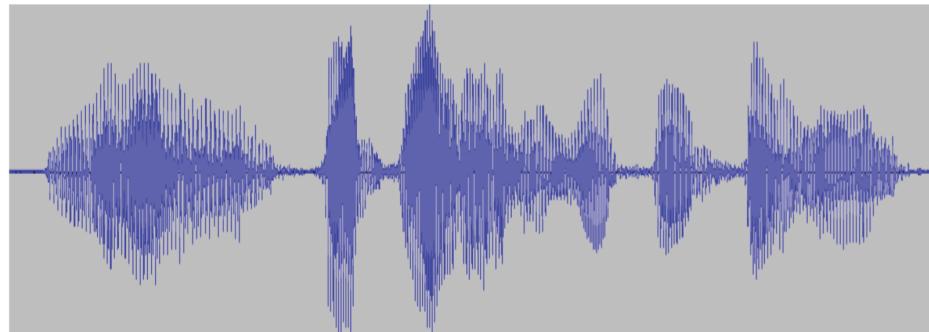


Outline

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

Spectrogram

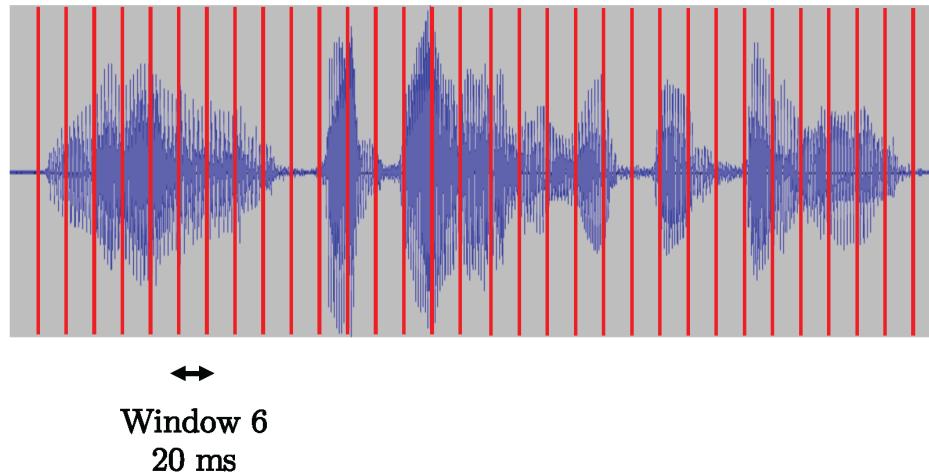
- We pre-process the raw audio signal to extract temporal and frequencial information.



Raw audio signal

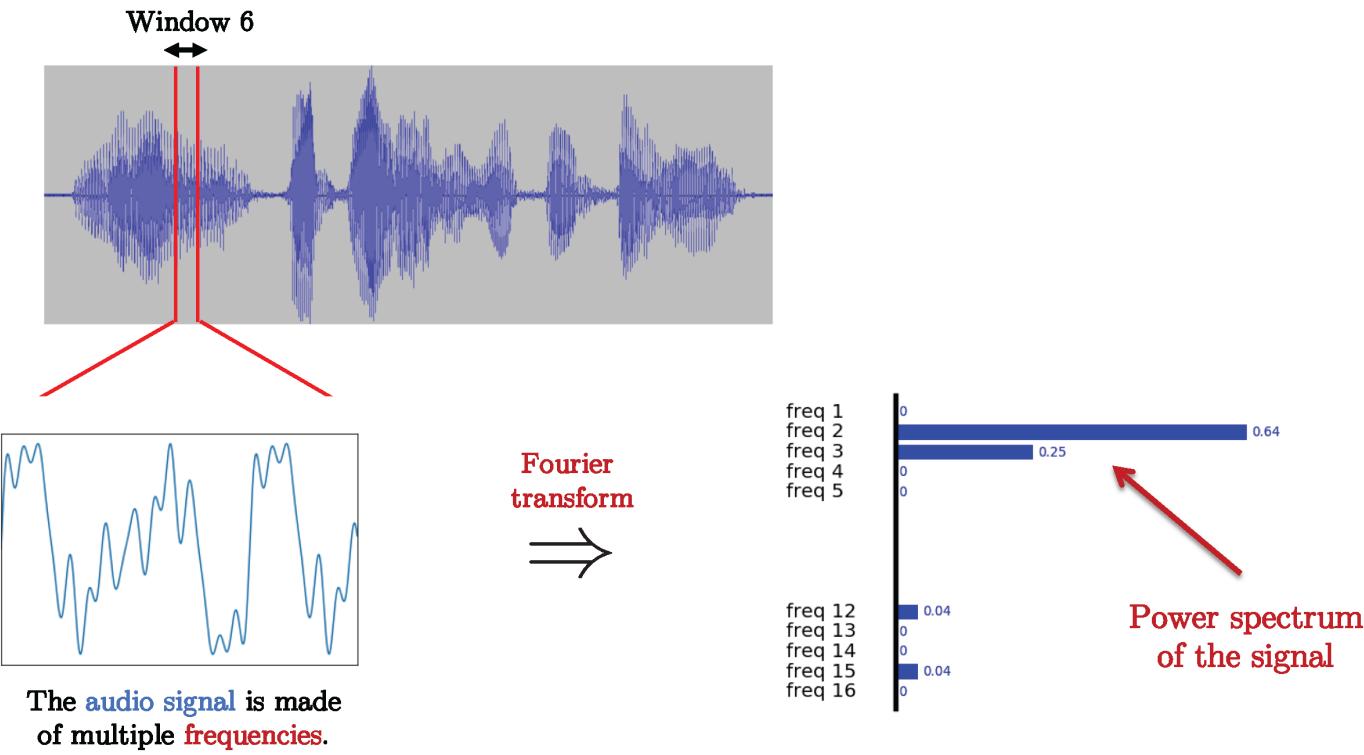
Spectrogram

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



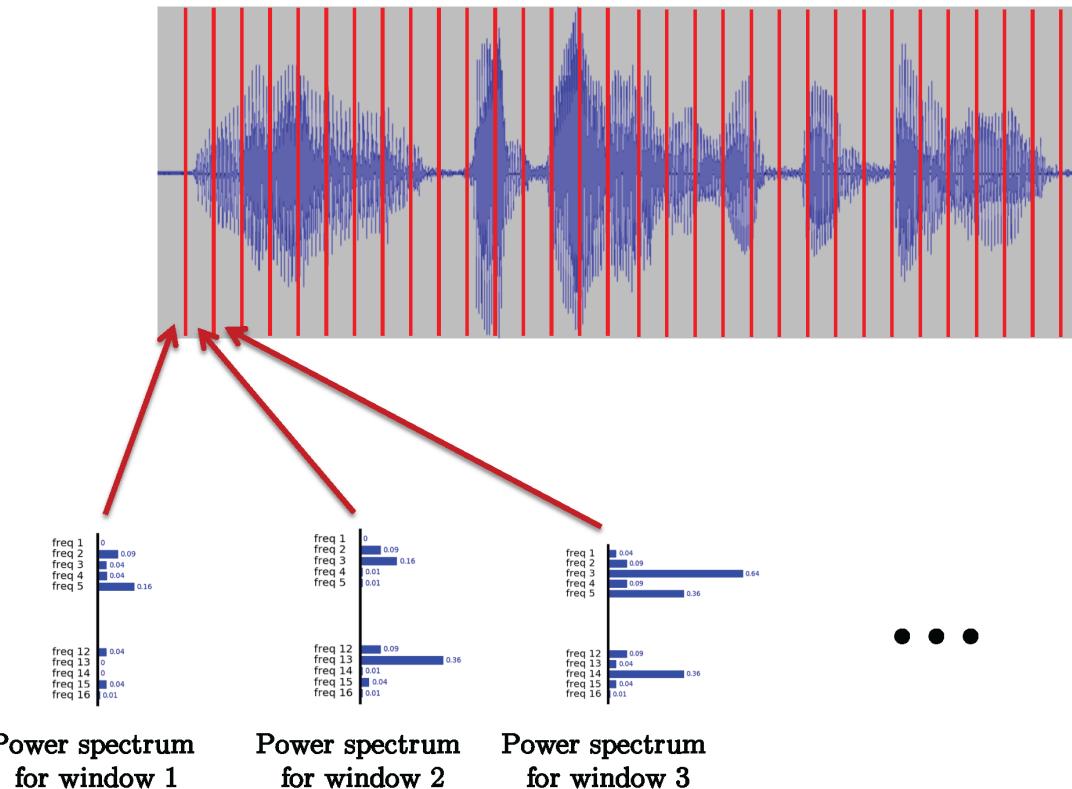
Spectrogram

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



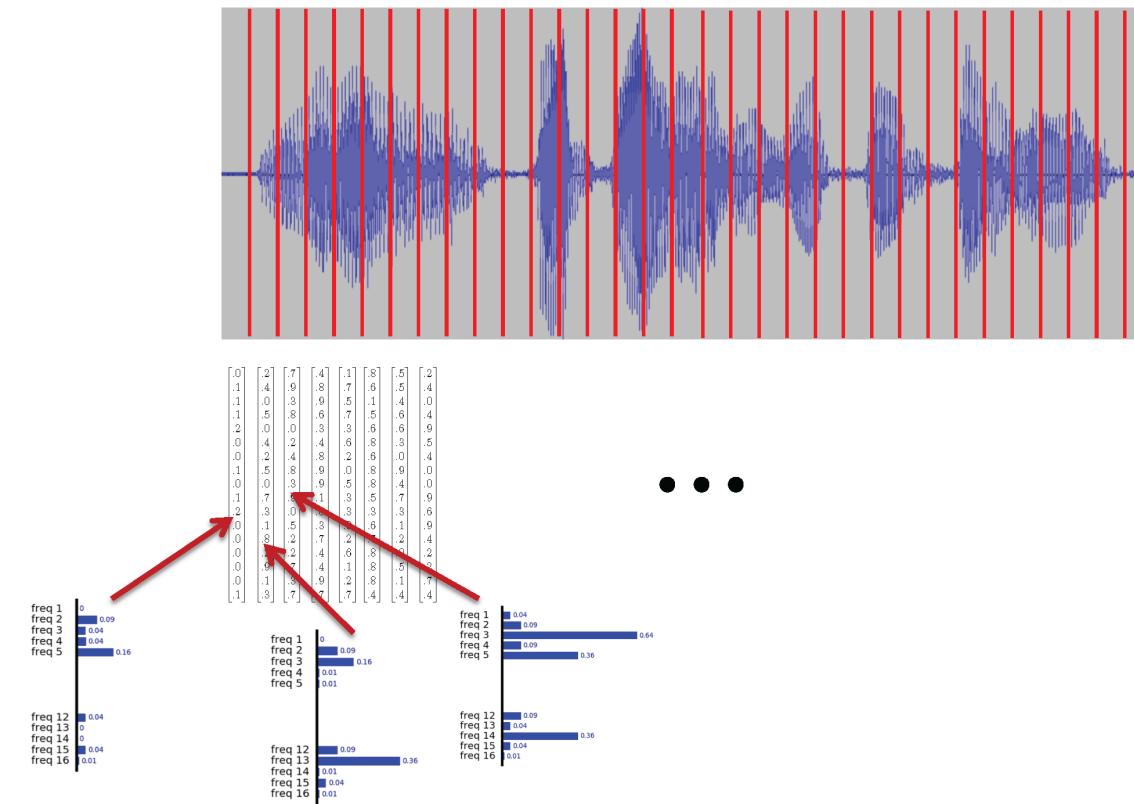
Spectrogram

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



Spectrogram

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

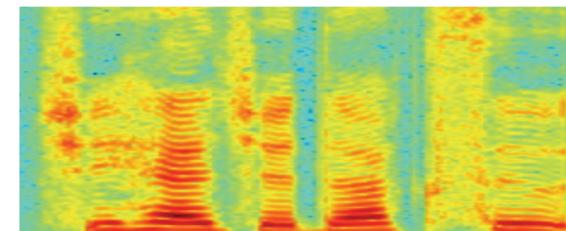


Spectrogram

- 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



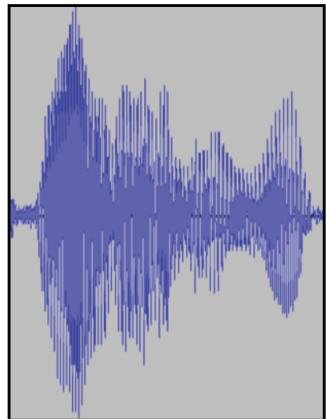
Spectrogram

Outline

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

Inference

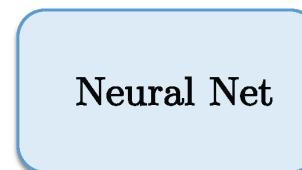
- Inference process :
 - Feed the sequence of vectors (spectrogram) to a trained RNN :



Input :
Audio signal

[.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

Spectrogram



“Hello”

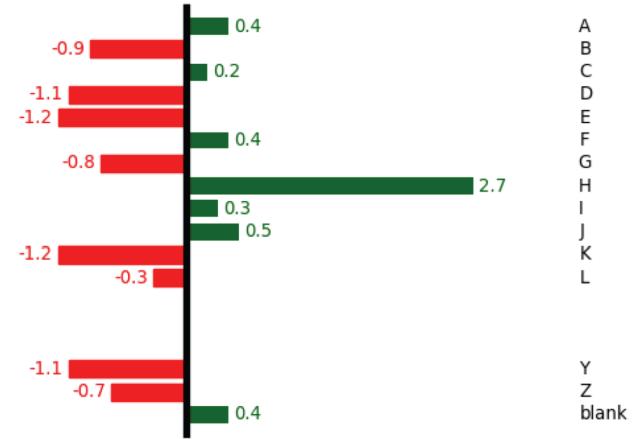
Output :
Sequence of words

Inference

- Dictionary/vocabulary is defined as follows :

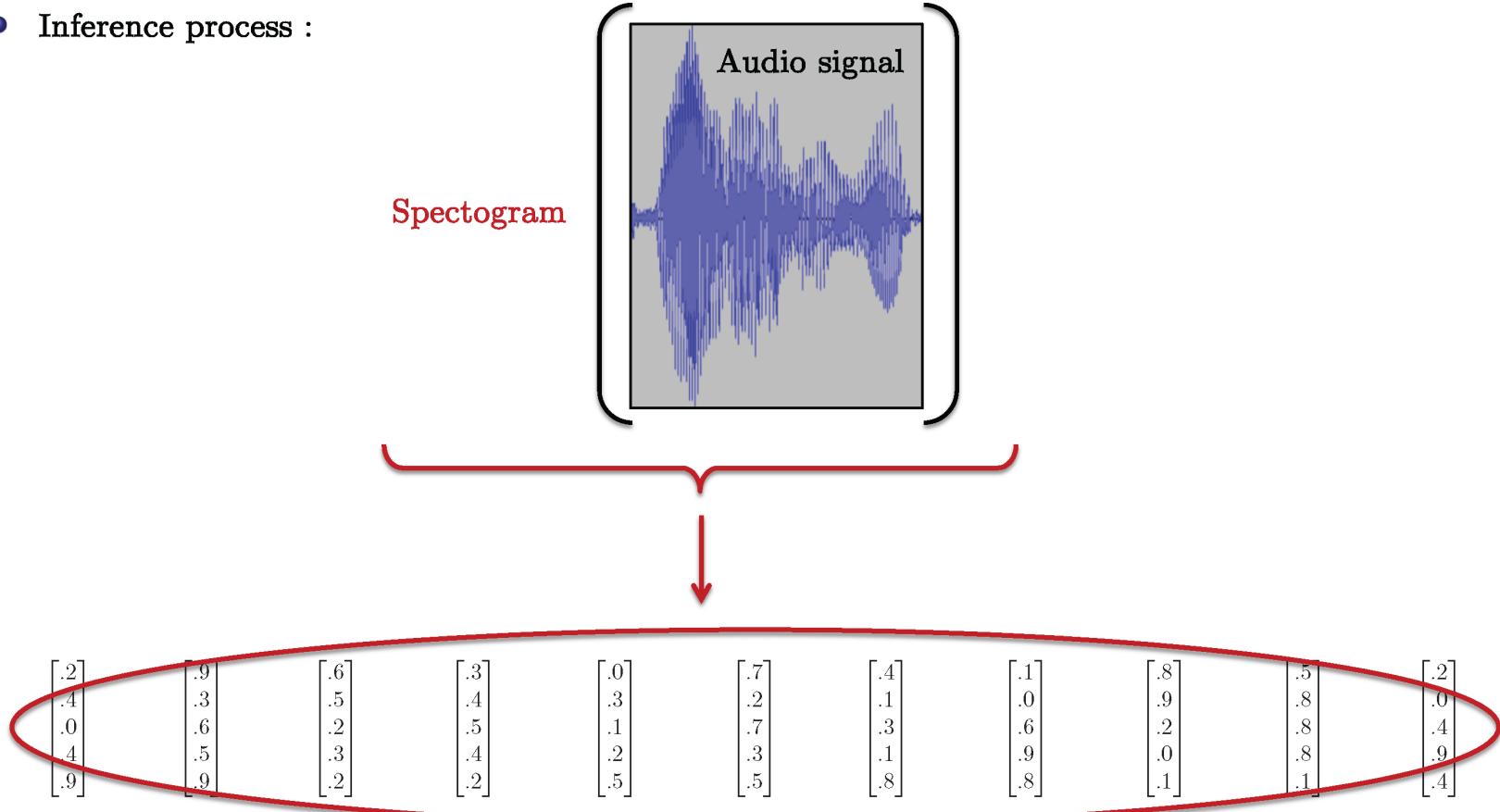
- 27 classes :

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



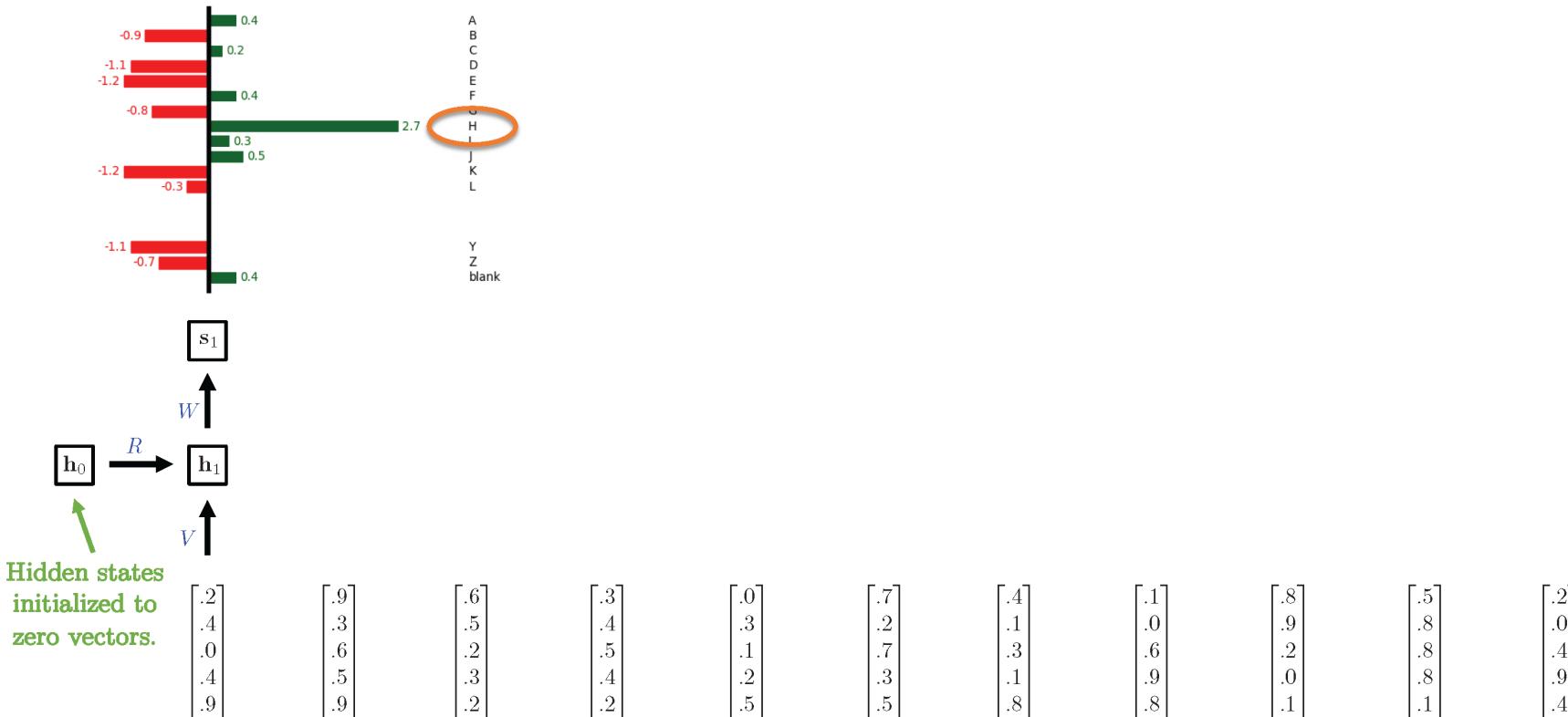
Inference

- Inference process :



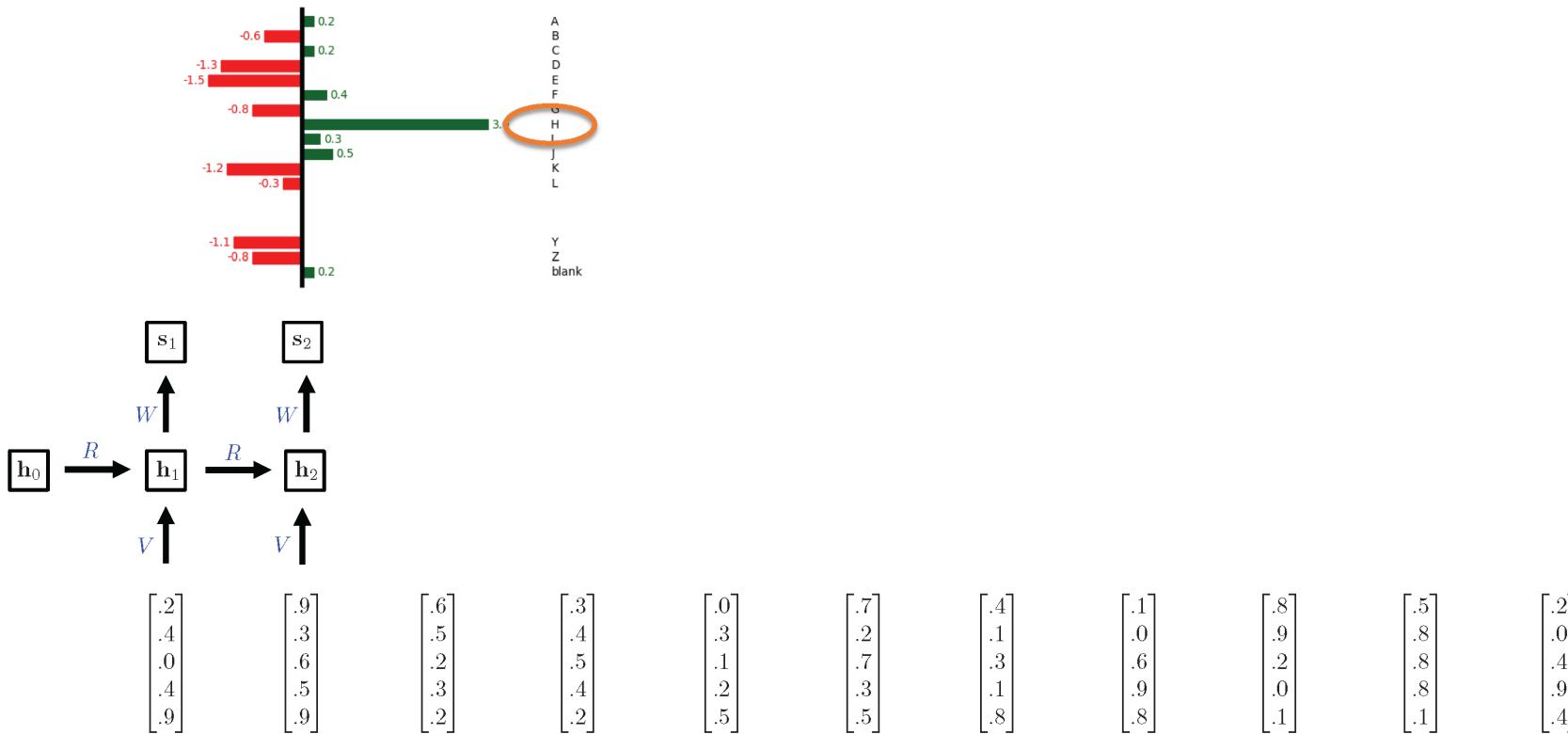
Inference

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



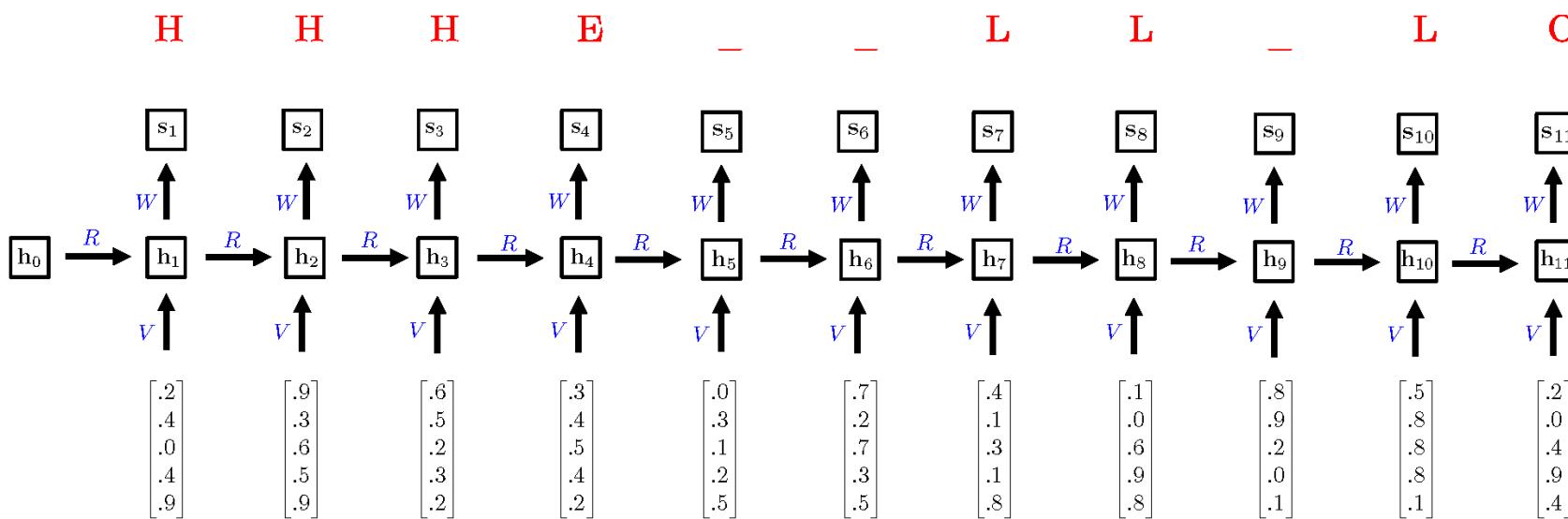
Inference

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



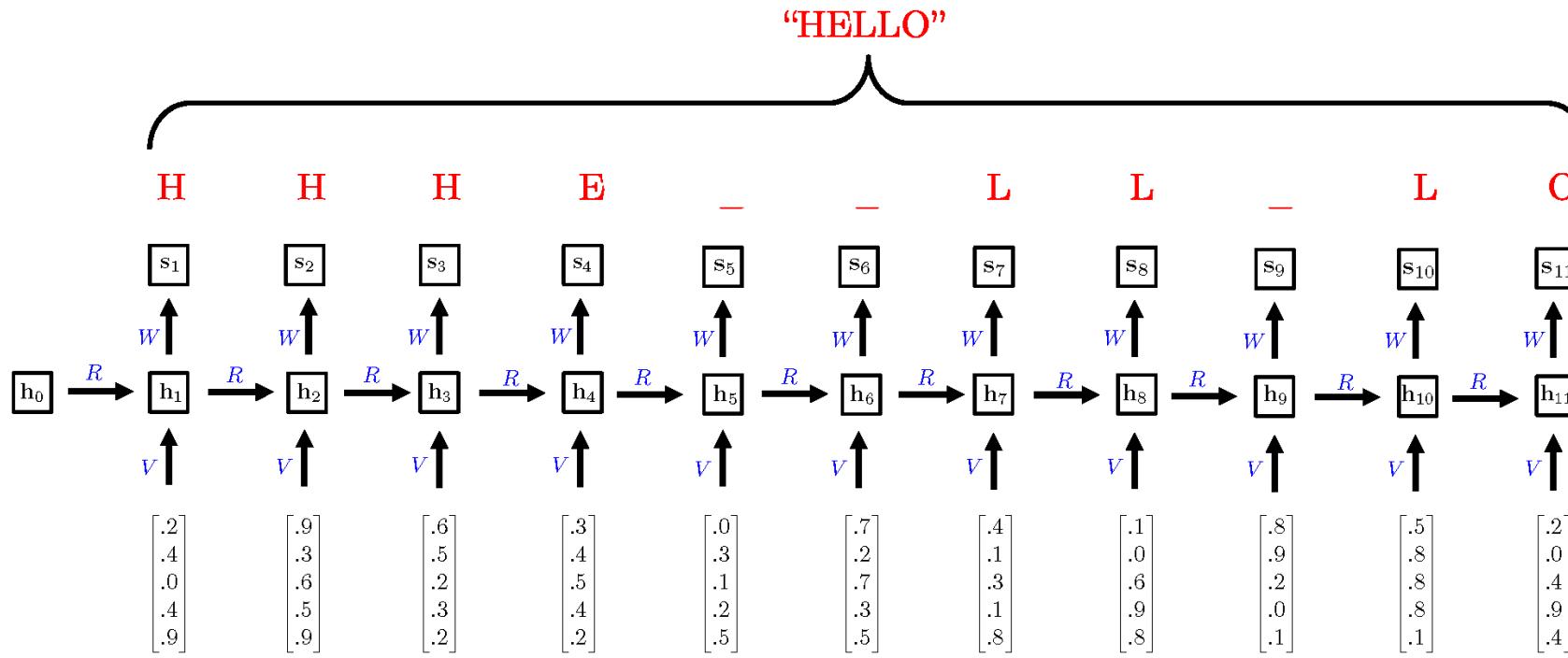
Inference

- 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 image displays three sequential screenshots of the Google Translate interface, illustrating the process of translating a passage from English to French and back to English.

Screenshot 1 (English to French): The input text is: "As Harry squelched along 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 ..."". The output translation is: "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 remplis pas leurs exigences ... un demi-pouce, si cela ...".

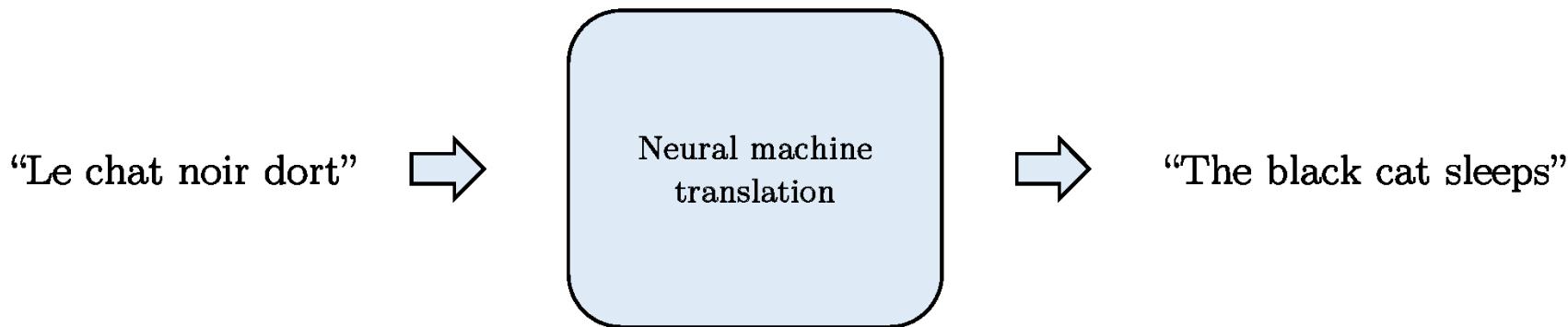
Screenshot 2 (French to English): The input text is: "Bonjour, Nick," dit Harry. The output translation is: "Hello, Nick," Harry said.

Screenshot 3 (English to French again): The input text is: "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 morosely through a window, whispering in his beard, "... does not meet their demands ... half an inch, if that ...". The output translation is: "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 morosely through a window, whispering in his beard, "... does not meet their demands ... half an inch, if that ...".

Text at the bottom: 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

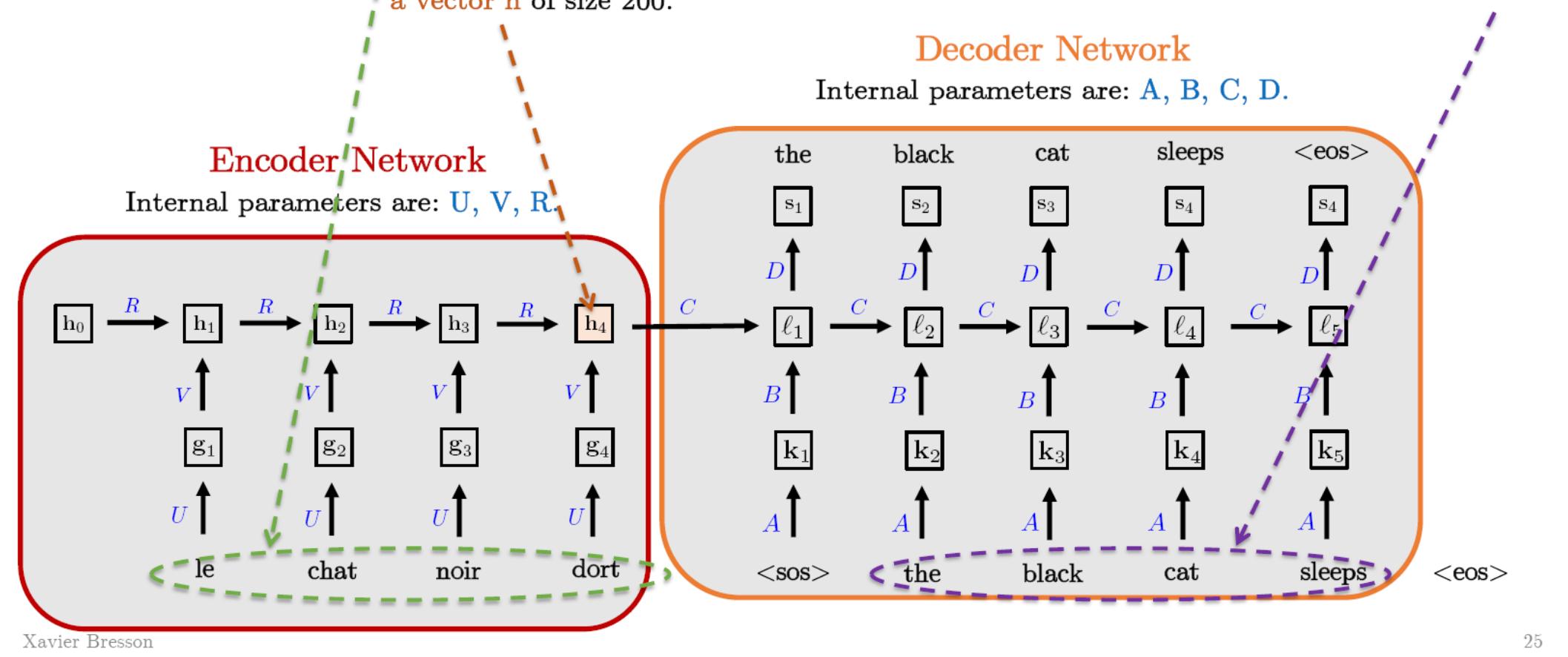
Architecture

- 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.



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



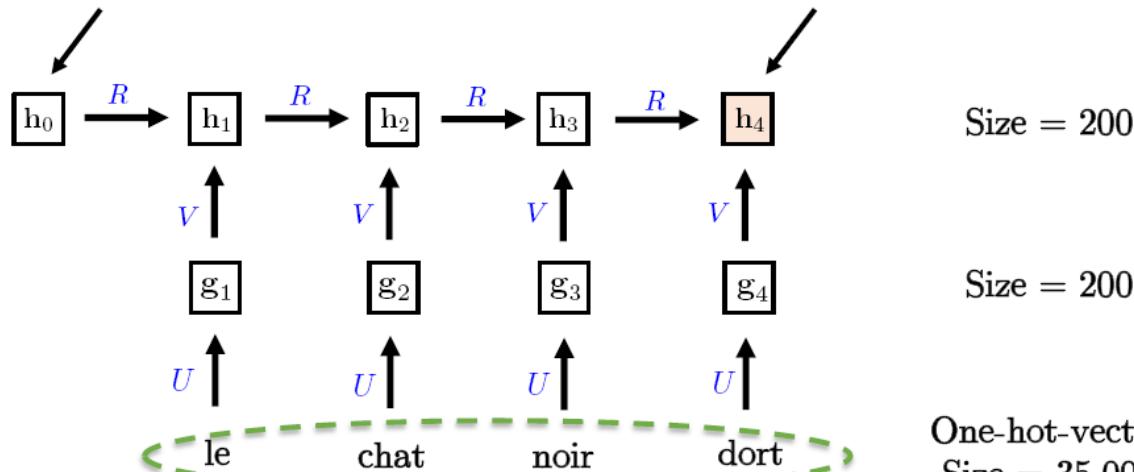
Encoder
network



“le chat noir dort”

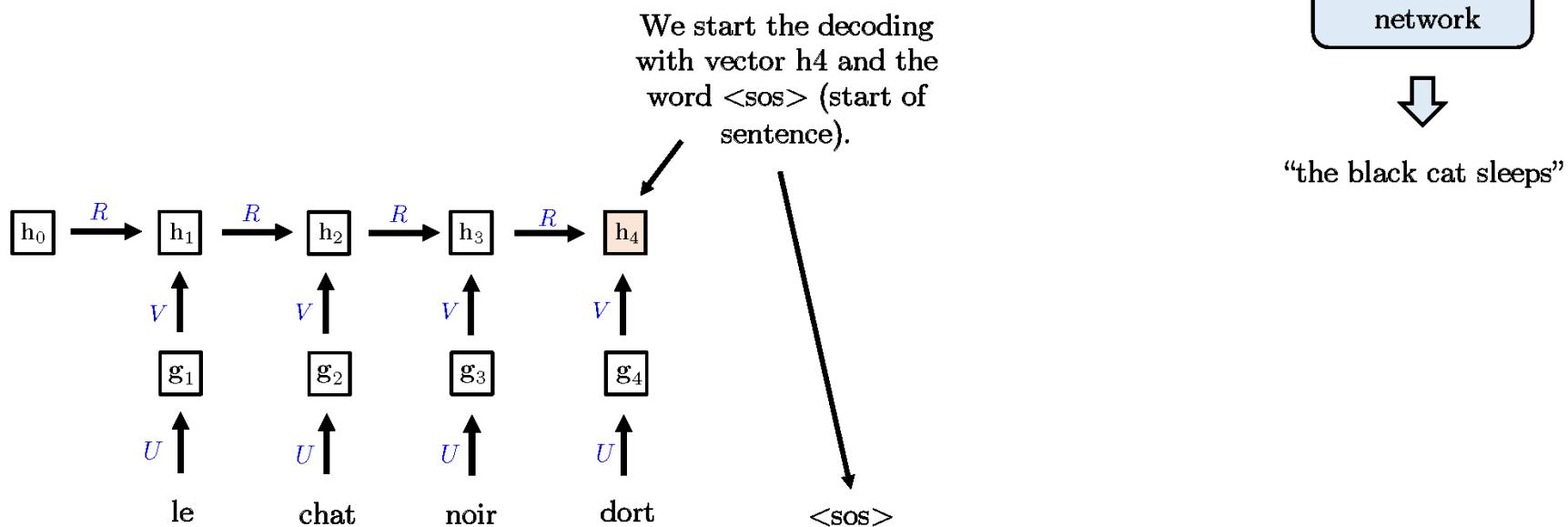
The hidden state is
initialized at zero.

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



Inference

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



History/encoding
vector h



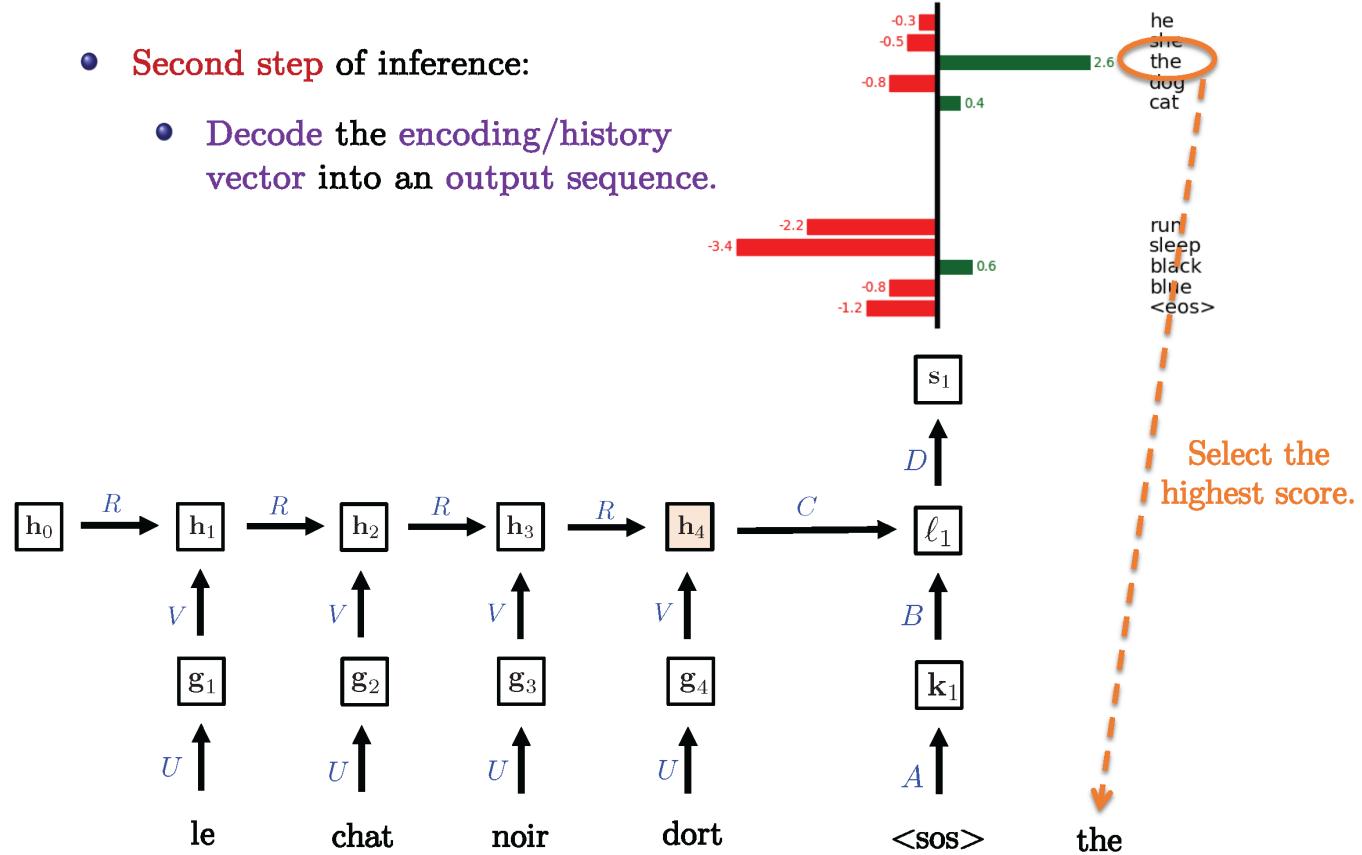
Decoder
network



"the black cat sleeps"

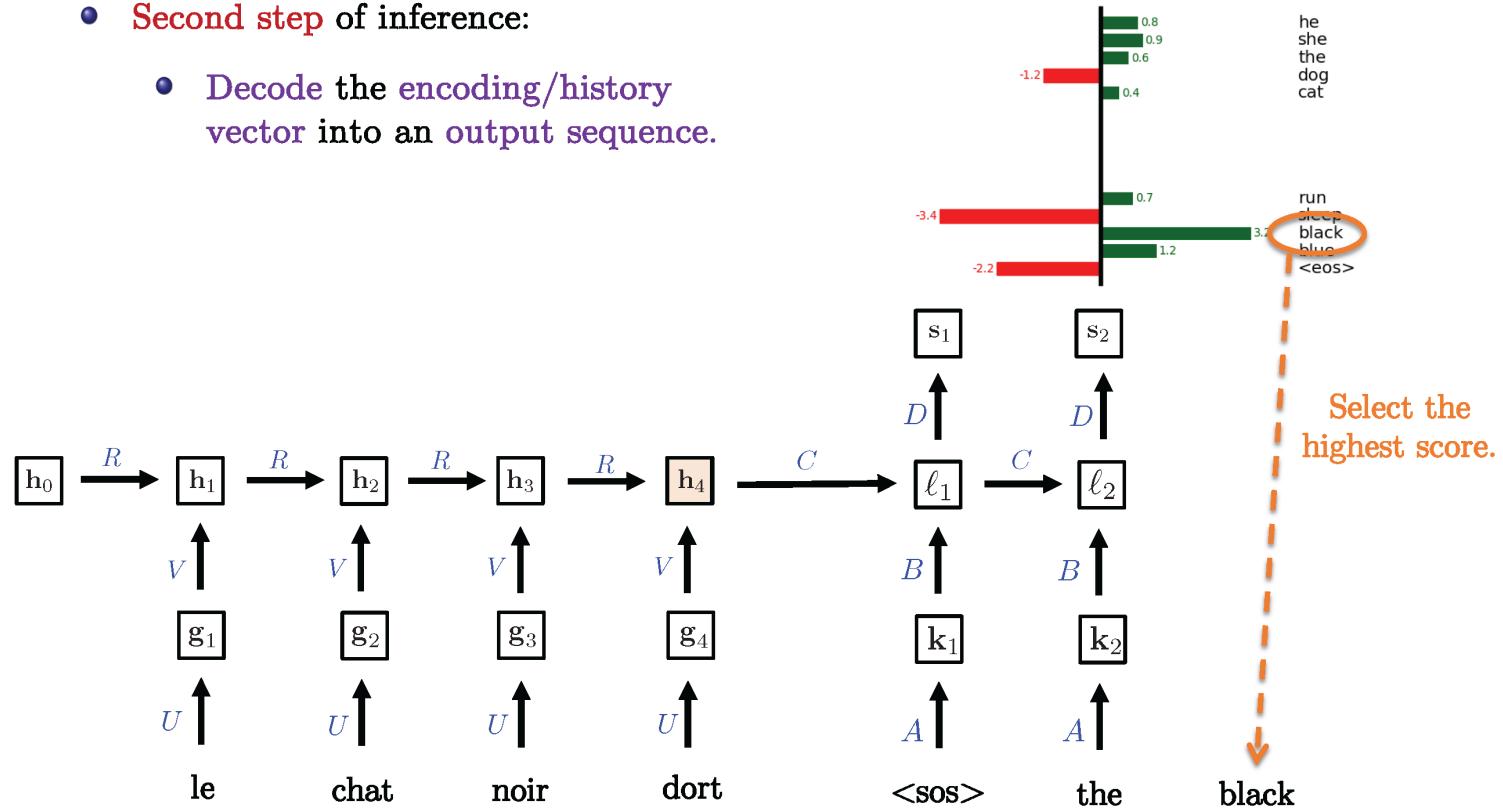
Inference

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



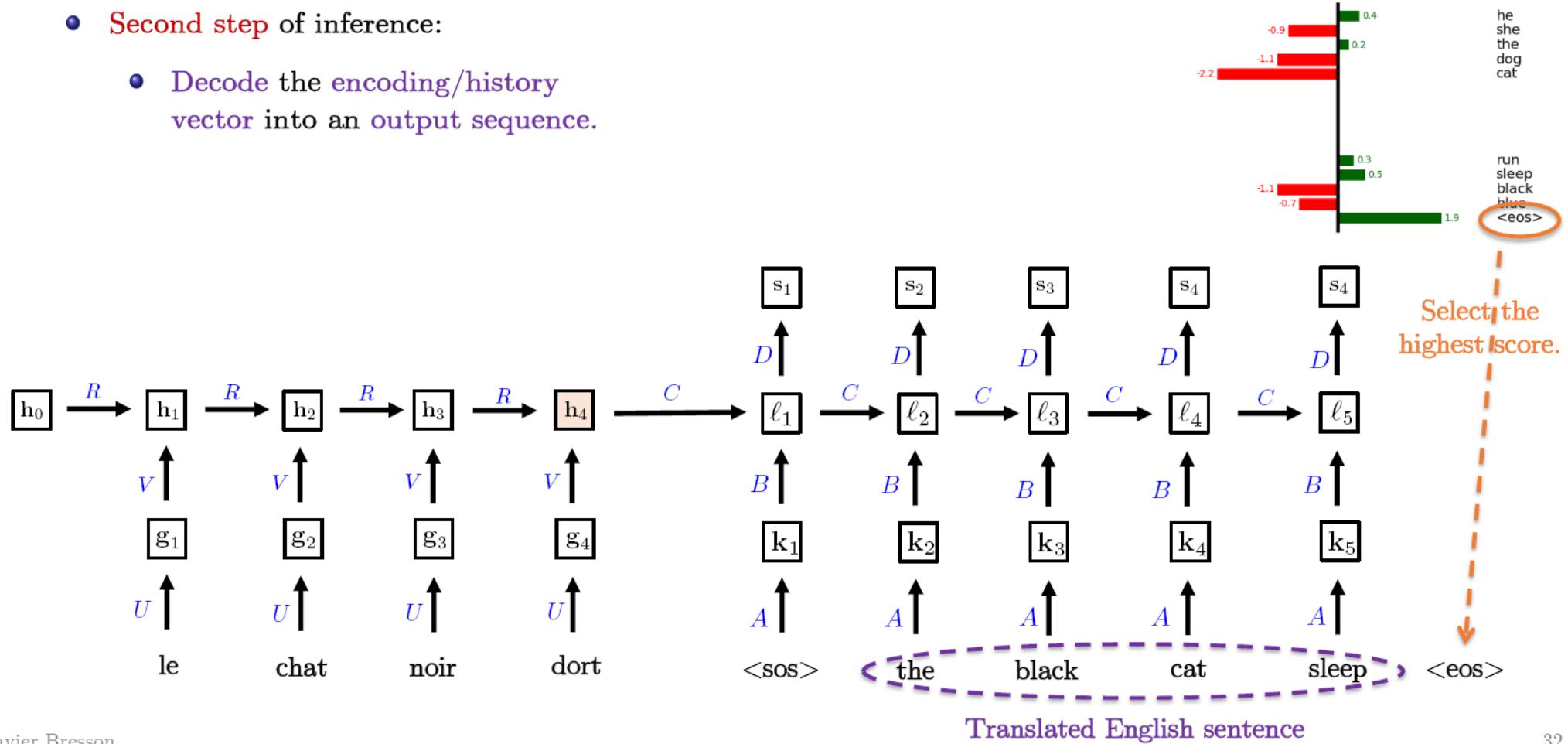
Inference

- 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** :

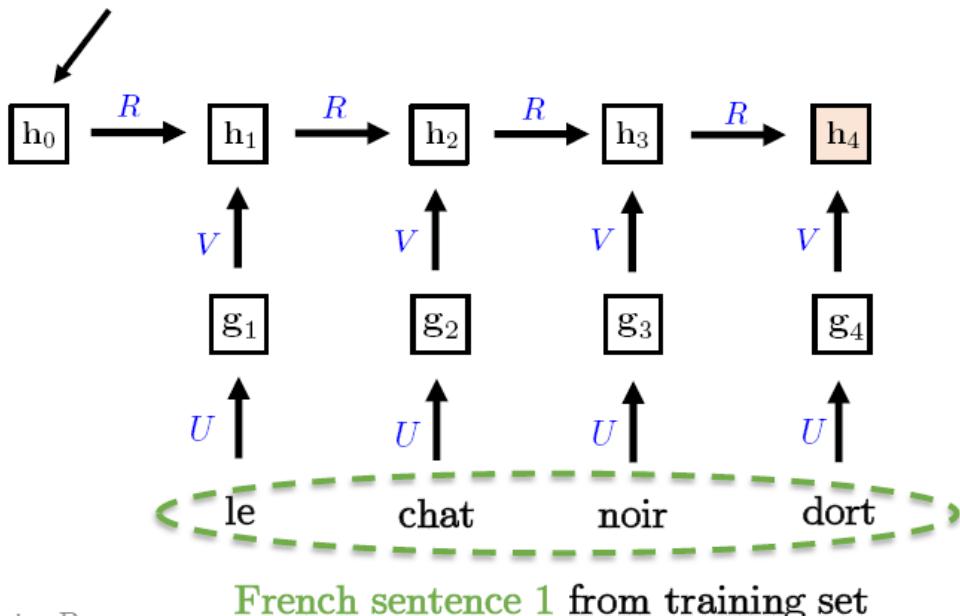
Sentence 1	le chat noir dort	<sos> the black cat sleep <eos>
Sentence 2	hier je suis aller a la plage	<sos> yesterday I went to the beach <eos>
Sentence 3	je m'appelle thomas	<sos> my name is thomas <eos>
⋮	⋮	⋮
Sentence 10,000,000	tu est un tres bon joueur de tennis	<sos> you are a very good tennis player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

Training

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

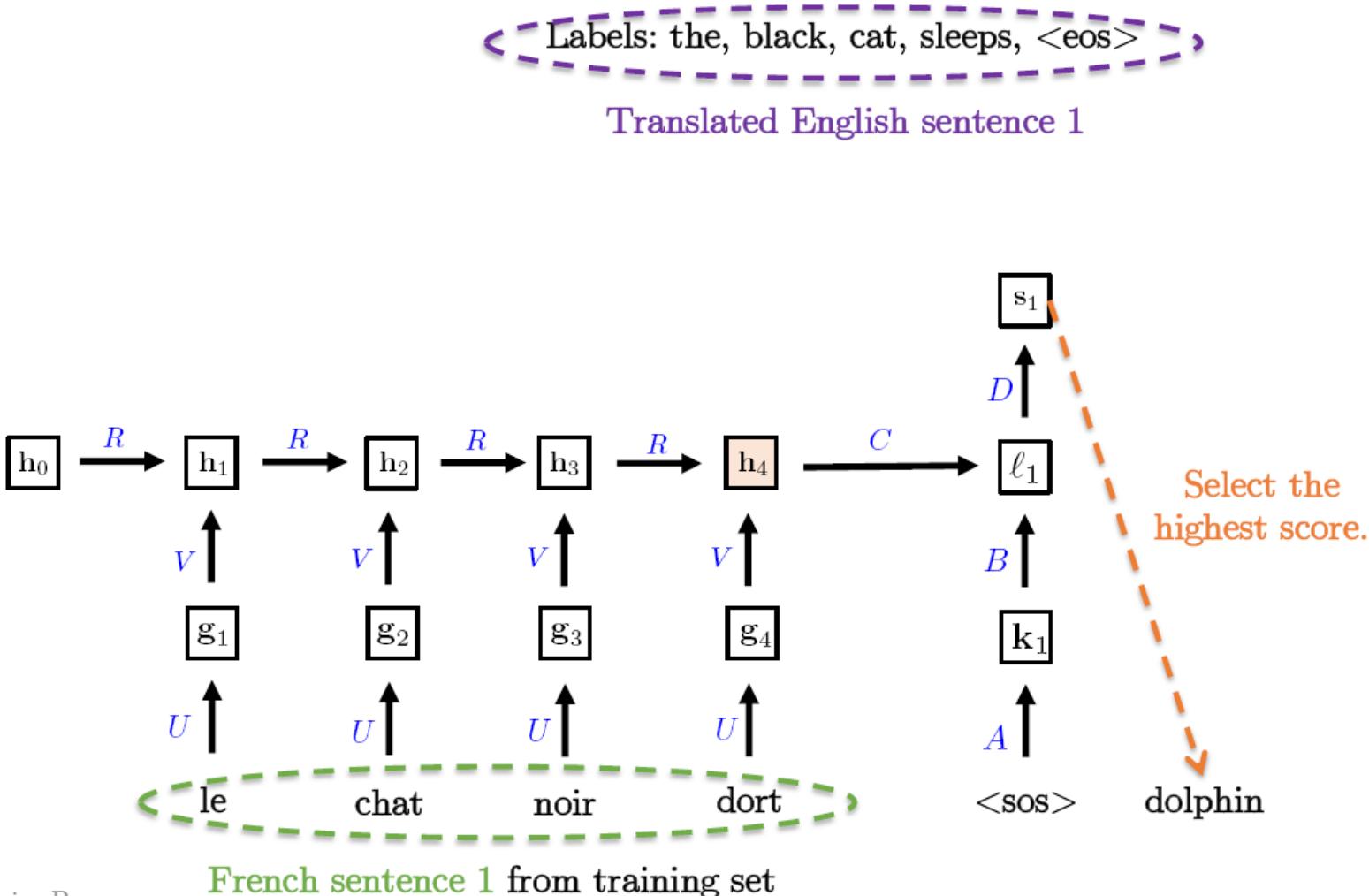
< Labels: the, black, cat, sleeps, <eos> >
Translated English sentence 1

The hidden state is initialized at zero.



Training

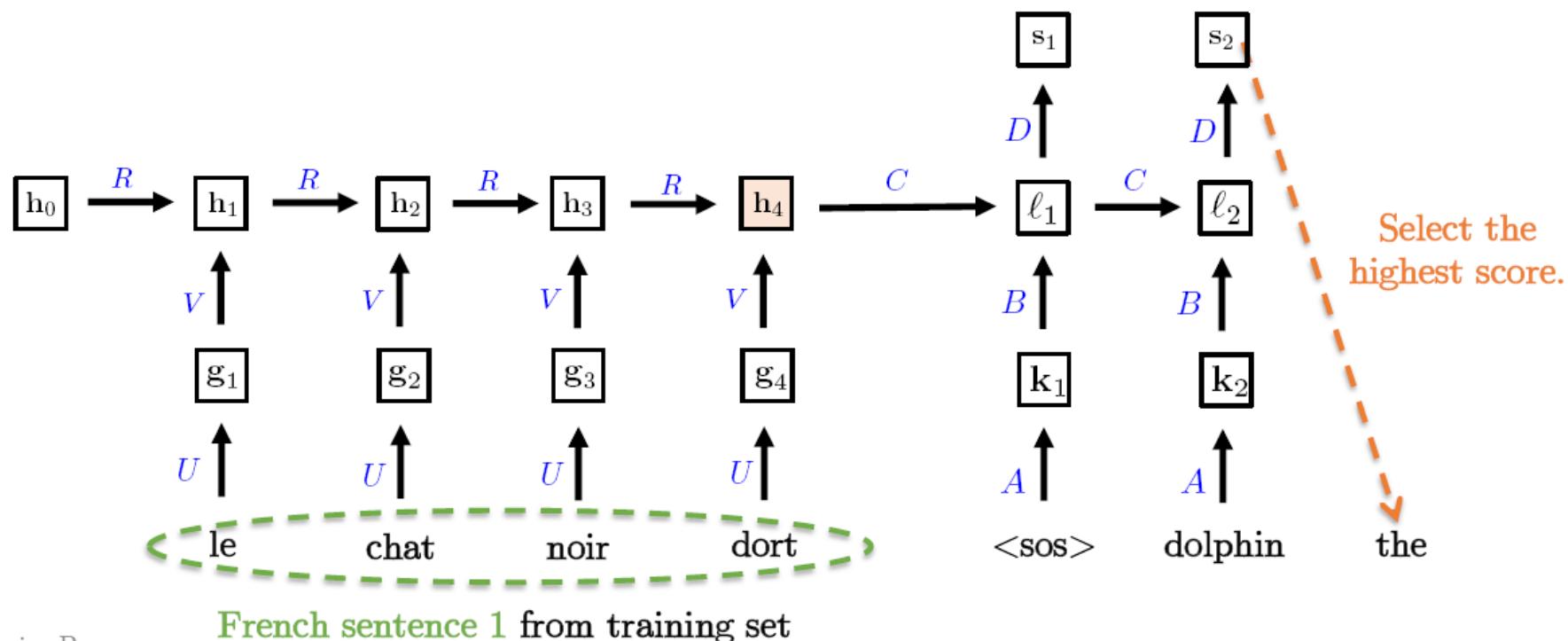
- Let us feed the 1st French \Rightarrow English sentence :



Training

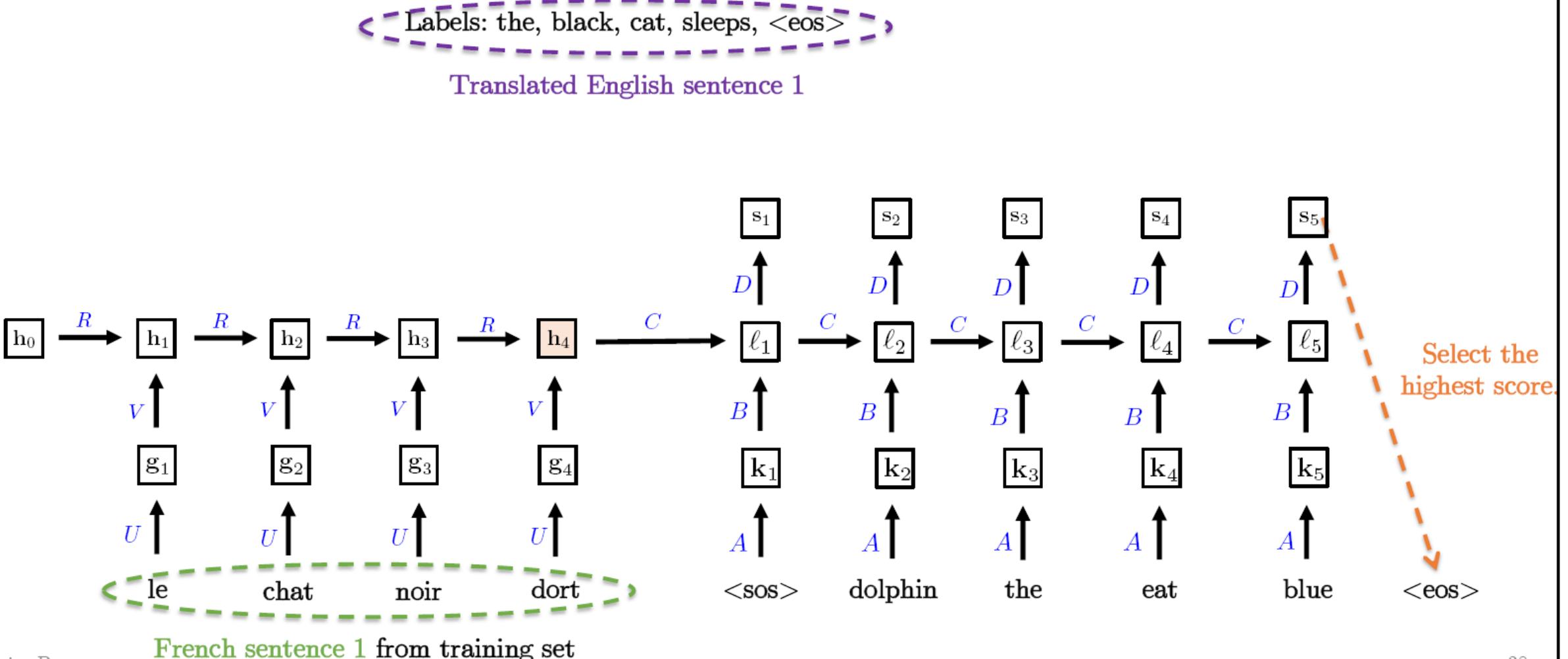
- Let us feed the 1st French \Rightarrow English sentence :

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



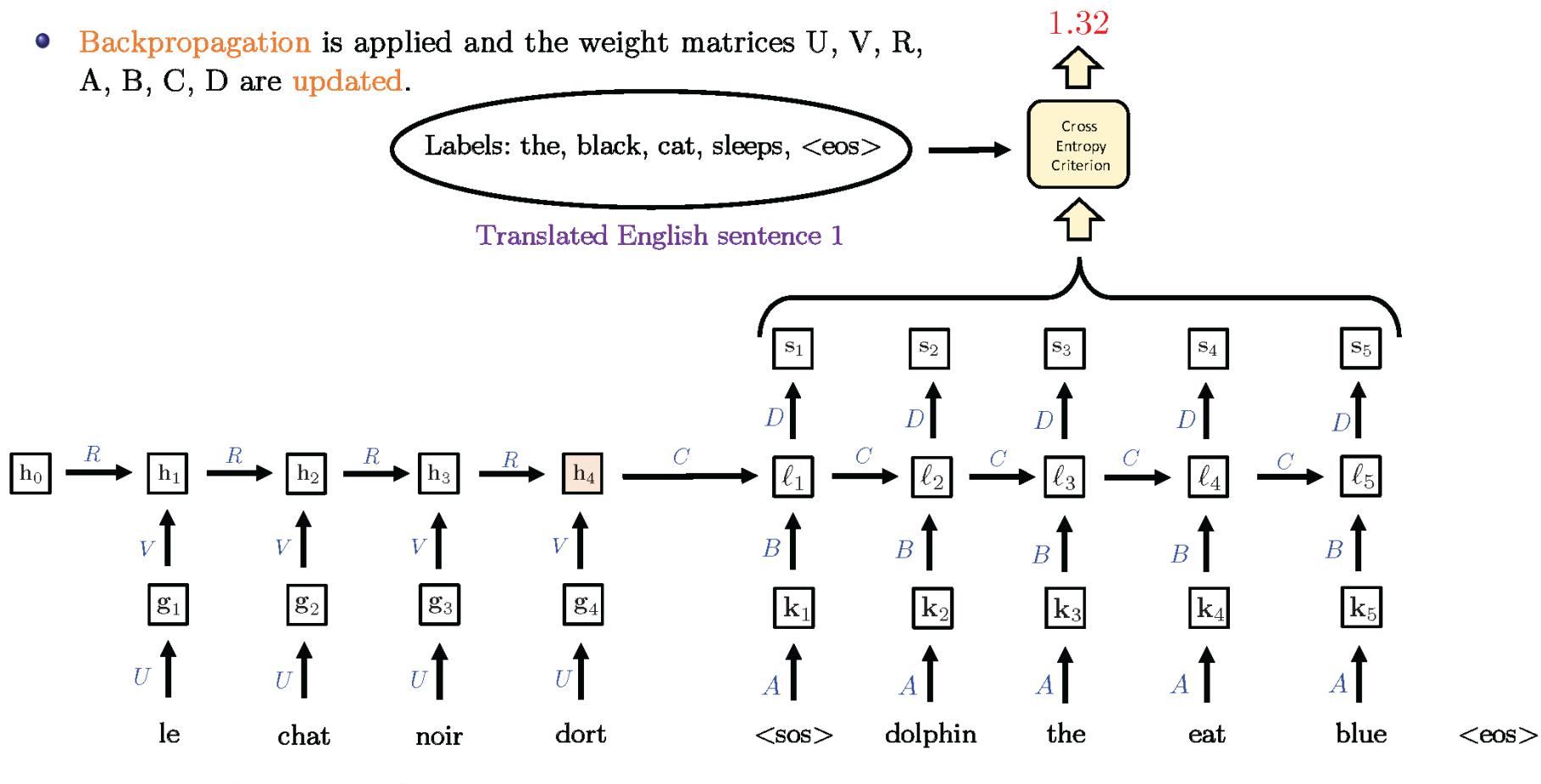
Training

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



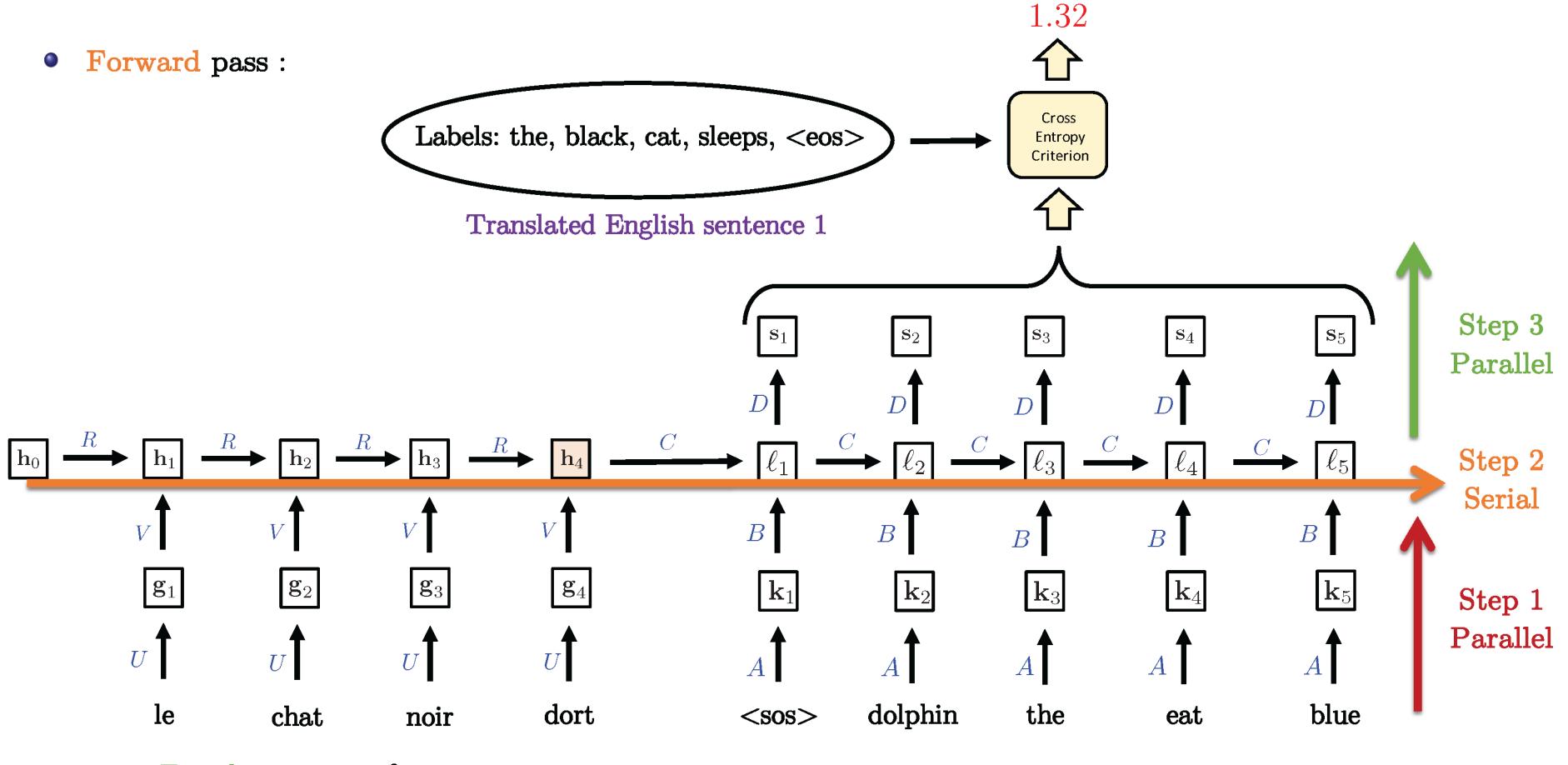
Training

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



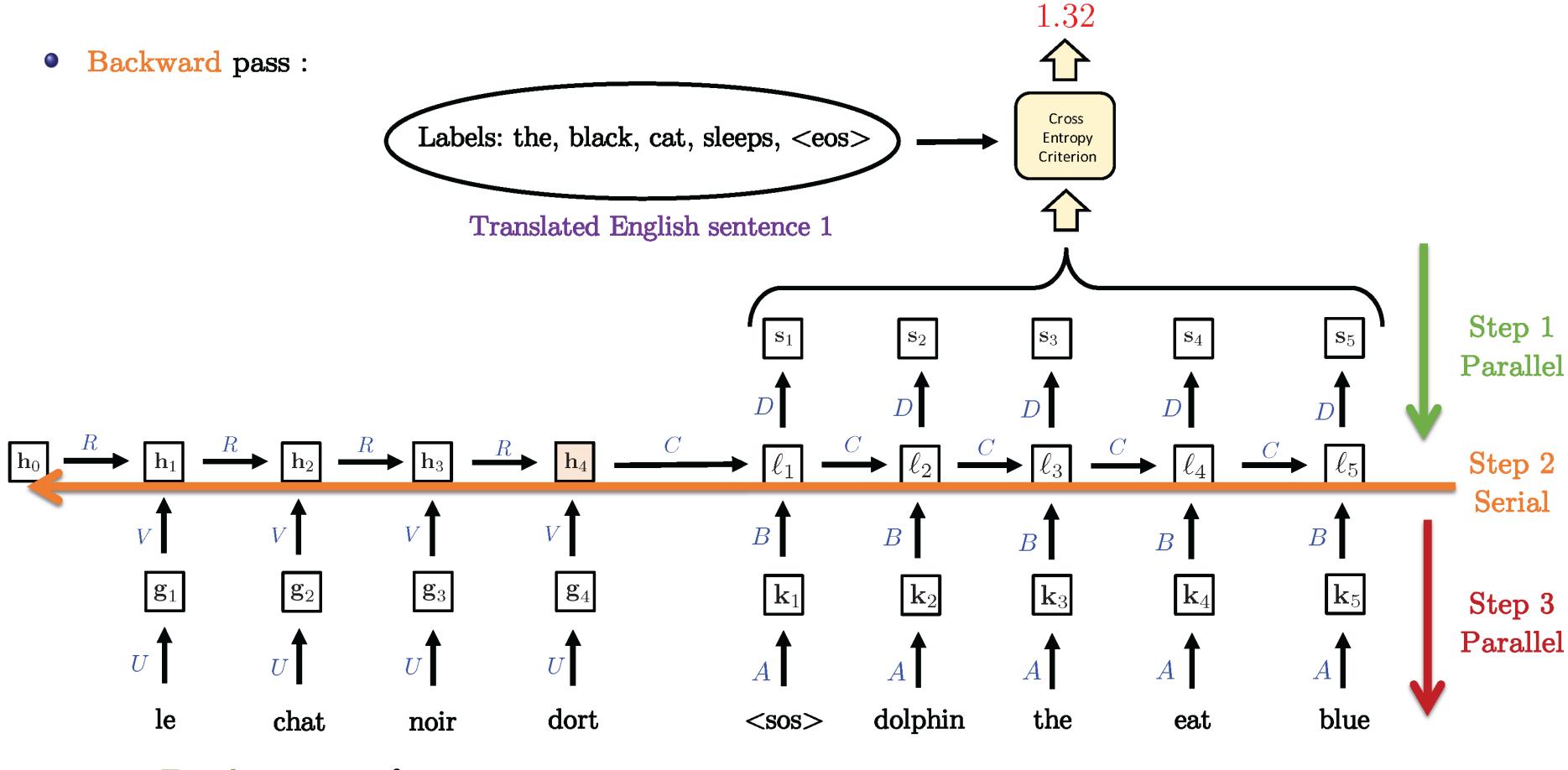
Training

- Forward pass :



Training

- Backward pass :



Training

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

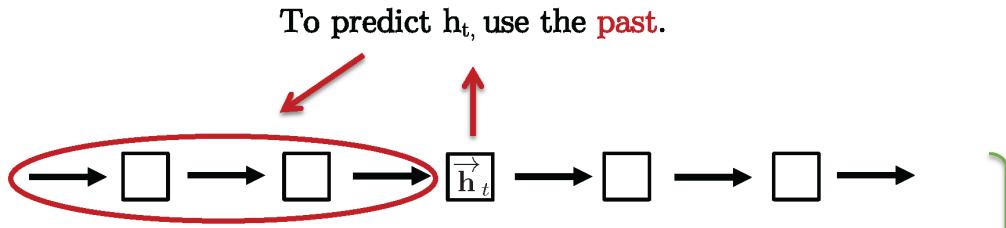
```
net = nn.LSTM( hidden_size , hidden_size , bidirectional=True )
```

- Use an attention mechanism.

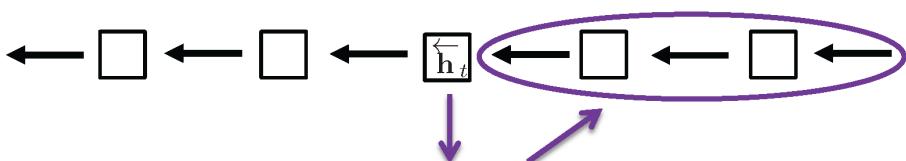
Training

- Bi-directional LSTM :

Causal signal w.r.t.
increasing time :

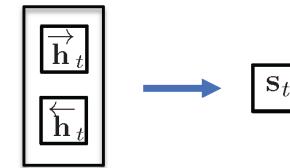


Causal signal w.r.t.
decreasing time :



Both uni-directional LSTMs are independent to compute their h .

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



Seq2Seq Model with Attention

Bahdanau Attention

$$\mathbf{c}_t = \sum_{i=1}^T \alpha_{t,i} \mathbf{h}_i = \sum_{i=1}^T \frac{\exp(e_{t,i})}{\sum_{j=1}^T \exp(e_{t,j})} \mathbf{h}_i$$

Decoder Hidden States

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$$

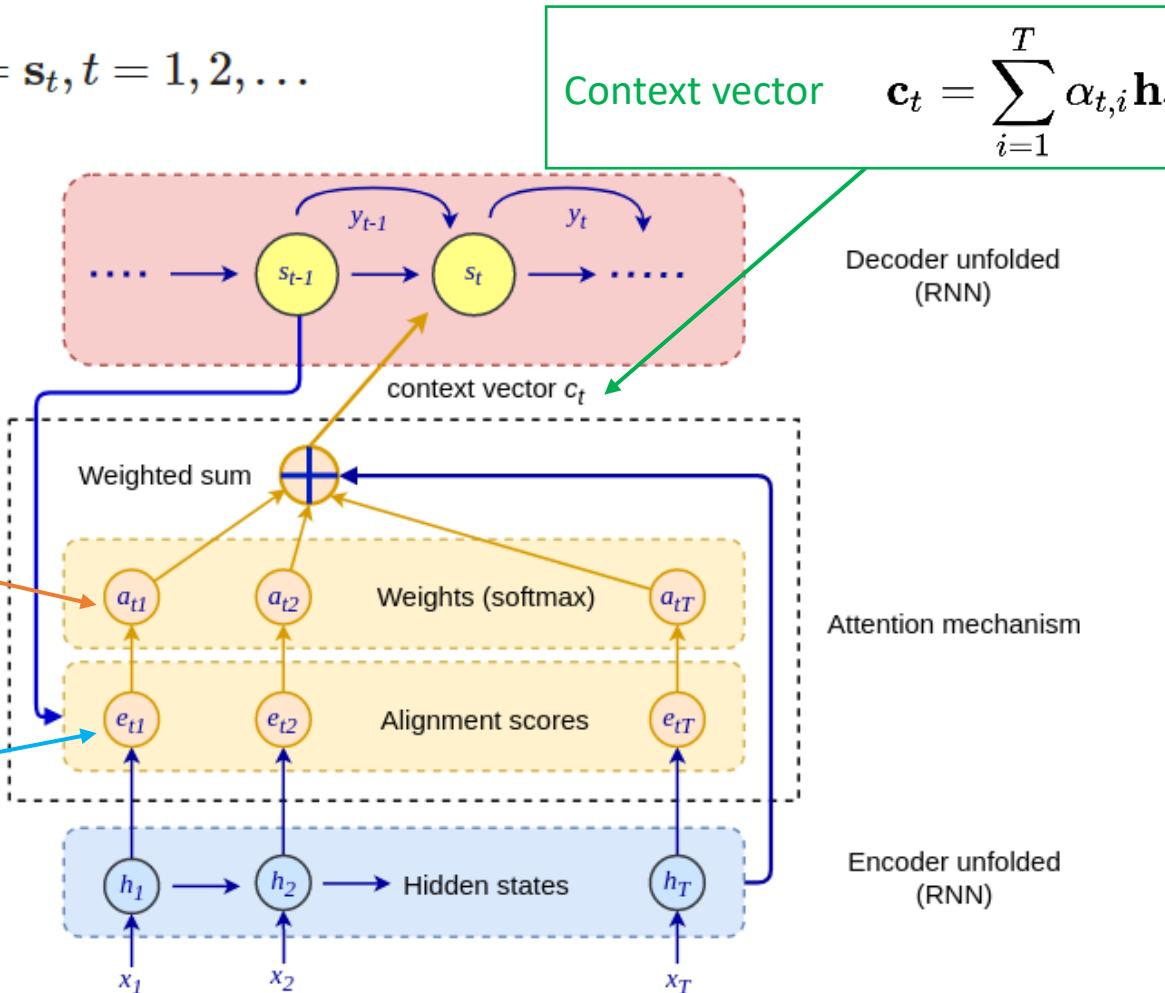
$$S = \mathbf{s}_t, t = 1, 2, \dots$$

Attention Weights

$$\alpha_{t,i} = \frac{\exp(e_{t,i})}{\sum_{j=1}^T \exp(e_{t,j})}$$

Alignment Scores

$$\begin{aligned} e_{t,i} &= a(\mathbf{s}_{t-1}, \mathbf{h}_i) \\ &= \mathbf{v}^T \tanh(\mathbf{W}[\mathbf{h}_i; \mathbf{s}_{t-1}]) \end{aligned}$$



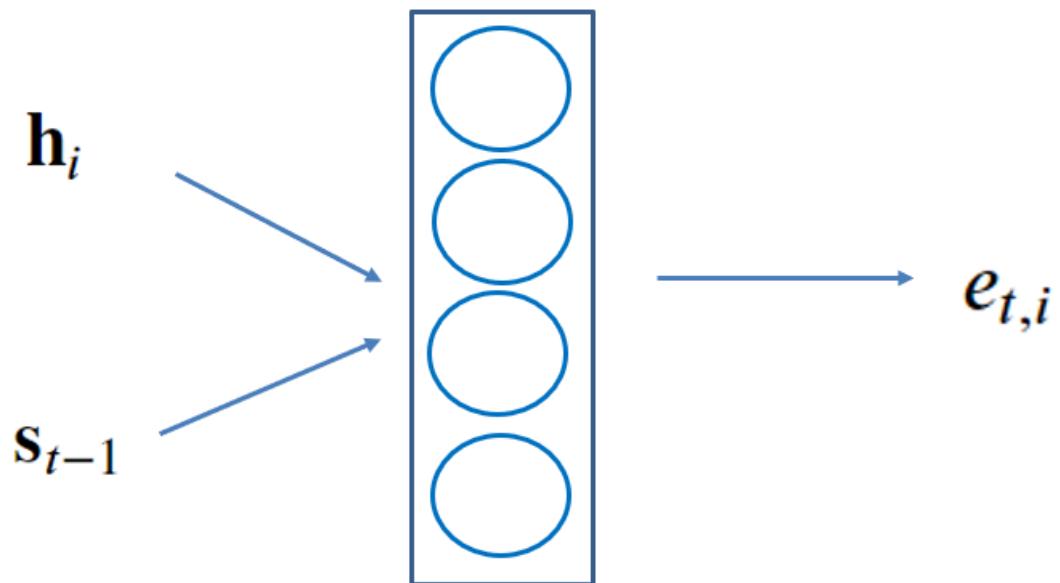
$$H = \mathbf{h}_i, i = 1, 2, \dots, T$$

Encoder Hidden States

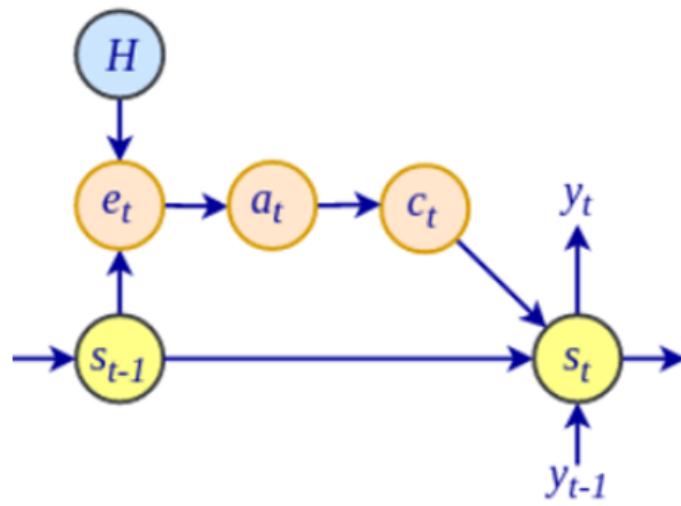
Alignment scores

$$e_{t,i} = a(\mathbf{s}_{t-1}, \mathbf{h}_i) = \begin{cases} \mathbf{s}_{t-1}^T \mathbf{h}_i & \text{dot} \\ \mathbf{s}_{t-1}^T \mathbf{W} \mathbf{h}_i & \text{general} \\ \boxed{\mathbf{v}^T \tanh(\mathbf{W}[\mathbf{h}_i; \mathbf{s}_{t-1}])} & \text{concat} \end{cases} \rightarrow \begin{array}{lll} \text{mutiply and sum} \\ \rightarrow \text{a linear model} \\ \rightarrow \text{a MLP} \end{array}$$

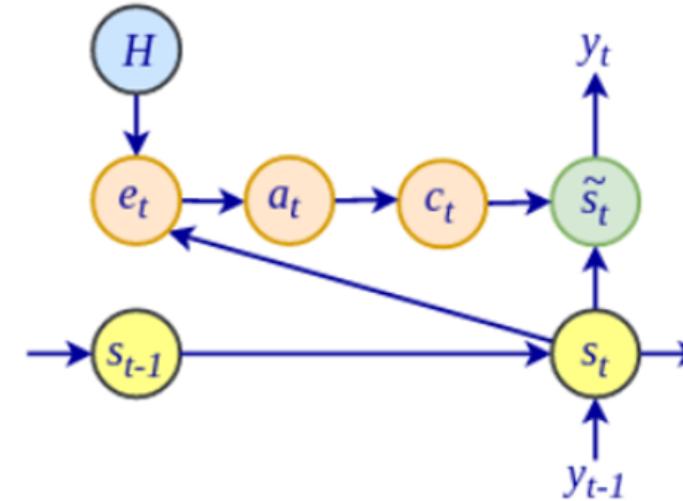
Bahdanau Attention



Bahdanau Attention



Luong Attention



	Bahdanau	Luong
Decoder Hidden State		$s_t = RNN_{decoder}(s_{t-1}, y_{t-1})$
Alignment Scores	$e_{t,i} = a(s_{t-1}, \mathbf{h}_i) = s_{t-1}^T \mathbf{h}_i$	$e_{t,i} = a(s_t, \mathbf{h}_i) = s_t^T \mathbf{h}_i$
Attention Weights	$\alpha_{t,i} = \text{softmax}(e_{t,i}/\mathbf{e}_t)$	$\alpha_{t,i} = \text{softmax}(e_{t,i}/\mathbf{e}_t)$
Context Vector	$\mathbf{c}_t = \sum_{i=1}^T \alpha_{t,i} \mathbf{h}_i$	$\mathbf{c}_t = \sum_{i=1}^T \alpha_{t,i} \mathbf{h}_i$
Decoder Hidden State (Attentional)	$s_t = RNN_{decoder}([s_{t-1}; \mathbf{c}_t], y_{t-1})$	$\tilde{s}_t = \tanh(W_c [\mathbf{c}_t; s_t])$
Final Output	$y_t = \text{softmax}(\mathbf{W}_y s_t)$	$y_t = \text{softmax}(\mathbf{W}_y \tilde{s}_t)$

Reference

- Bahdanau Attention Paper <https://arxiv.org/abs/1409.0473>
- Luong Attention Paper <https://arxiv.org/abs/1508.04025>
- Ivan Vasilev “Advanced Deep Learning with Python”