

2014 NIPS paper:

An Autoencoder Approach to Learning Bilingual Word Representations

By

Stanislas Lauly

Supervisor : Hugo Larochelle

Work done with

2

- Sarath Chandar (*Indian Institute of Technology Madras, India*)
- Hugo Larochelle (*Université de Sherbrooke, Canada*)
- Mitesh M. Khapra (*IBM Research India*)
- Balaraman Ravindran (*Indian Institute of Technology Madras, India*)
- Vikas Raykar (*IBM Research India*)
- Amrita Saha (*IBM Research India*)

Outline

3

- Bilingual representation
 - ▣ We use it for cross-lingual classification task
- The models
 - ▣ Binary bag-of-words reconstruction training
 - ▣ Tree-based decoder training
- Cross-lingual classification task
- The data
- The results

Bilingual representation

4

- We want to learn text representation that is invariant to the language
 - ▣ Same text in two different languages
 - ▣ Same representation

- To do this we learn vectorial word representation

Bilingual representation: Cross-lingual classification task

- We want to classify documents for one language
- Labeled data available in an other language

INTERNET USAGE BY LANGUAGE 2007
& GROWTH, 2000-2007

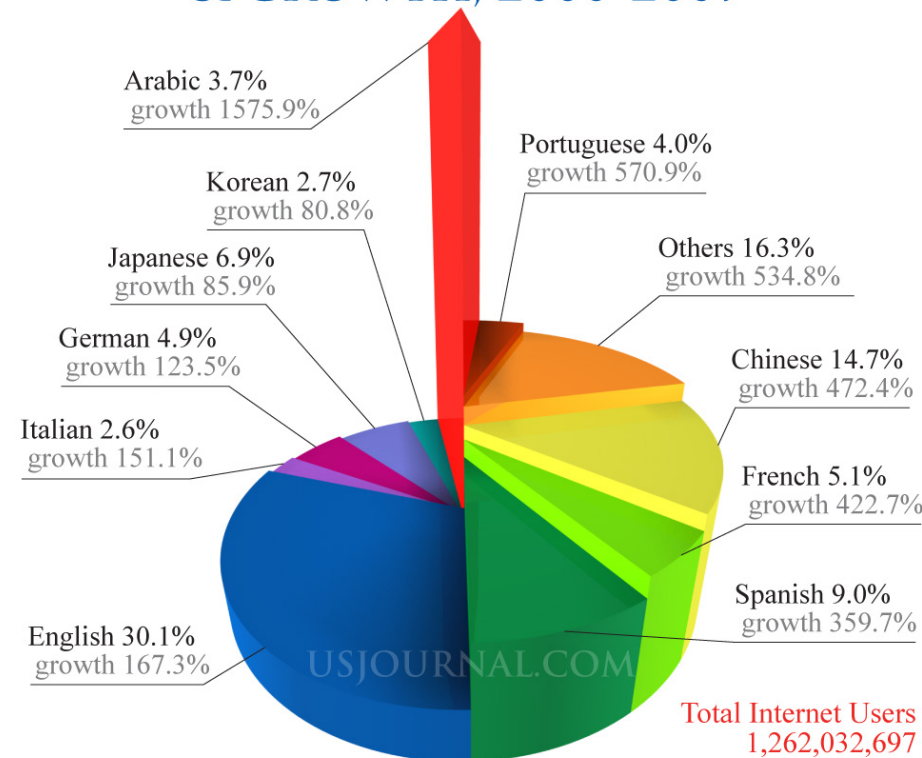


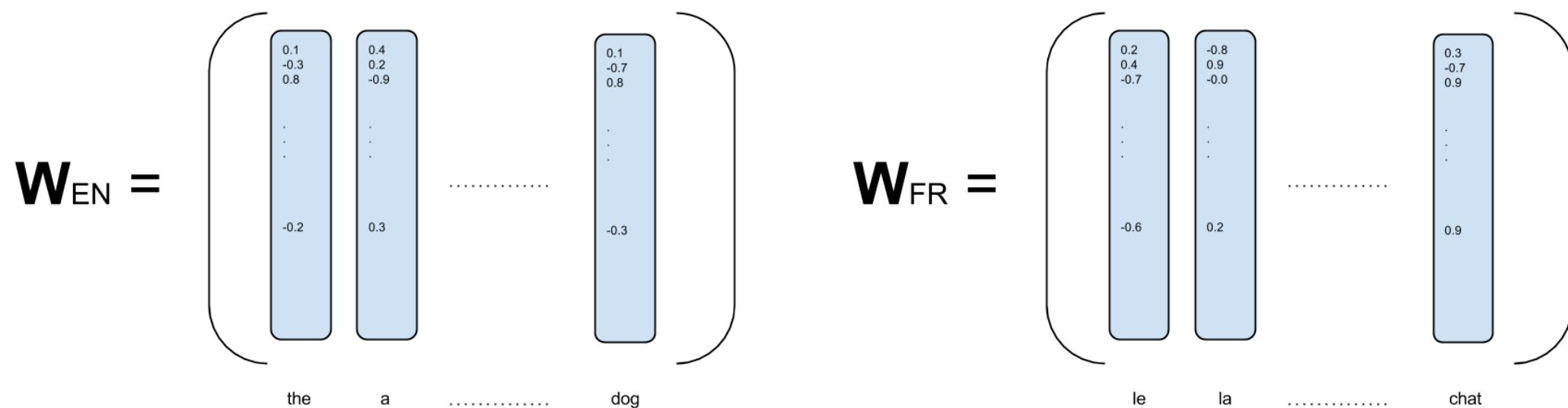
Image from Sarath Chandar

Source: www.internetworldstats.com
Graph by NingGeng Ong, USJournal.com

Bilingual representation

6

□ Vectorial representation



□ Text representation = sum of word vectors

Bilingual representation

7

- We use corpus of aligned sentences for languages **x** and **y**

x

y

It says that this should ... of relative stability.

Il précise que cela devrait ... de stabilité relative.

It will, I hope, be examined in a positive light.

Elle sera, je l'espère, examinée dans un ...

...

To this end ... as soon as possible.

En ce sens ... dans les plus brefs délais.

The models

8

- We explore autoencoders to learn the vectorial representation of the words
- Two types of autoencoders
 - ▣ Predict a binary bag-of-words
 - ▣ Predict a multinomial representation of the text
- Word alignment is not needed
- We use one language to predict the other

Binary bag-of-words reconstruction training

9

- Use a binary bag-of-words representation $\mathbf{v}(\mathbf{x})$

- Encoder

$$\mathbf{a}(\mathbf{x}) = \mathbf{c} + \mathbf{W}\mathbf{v}(\mathbf{x}), \quad \phi(\mathbf{x}) = \mathbf{h}(\mathbf{a}(\mathbf{x}))$$

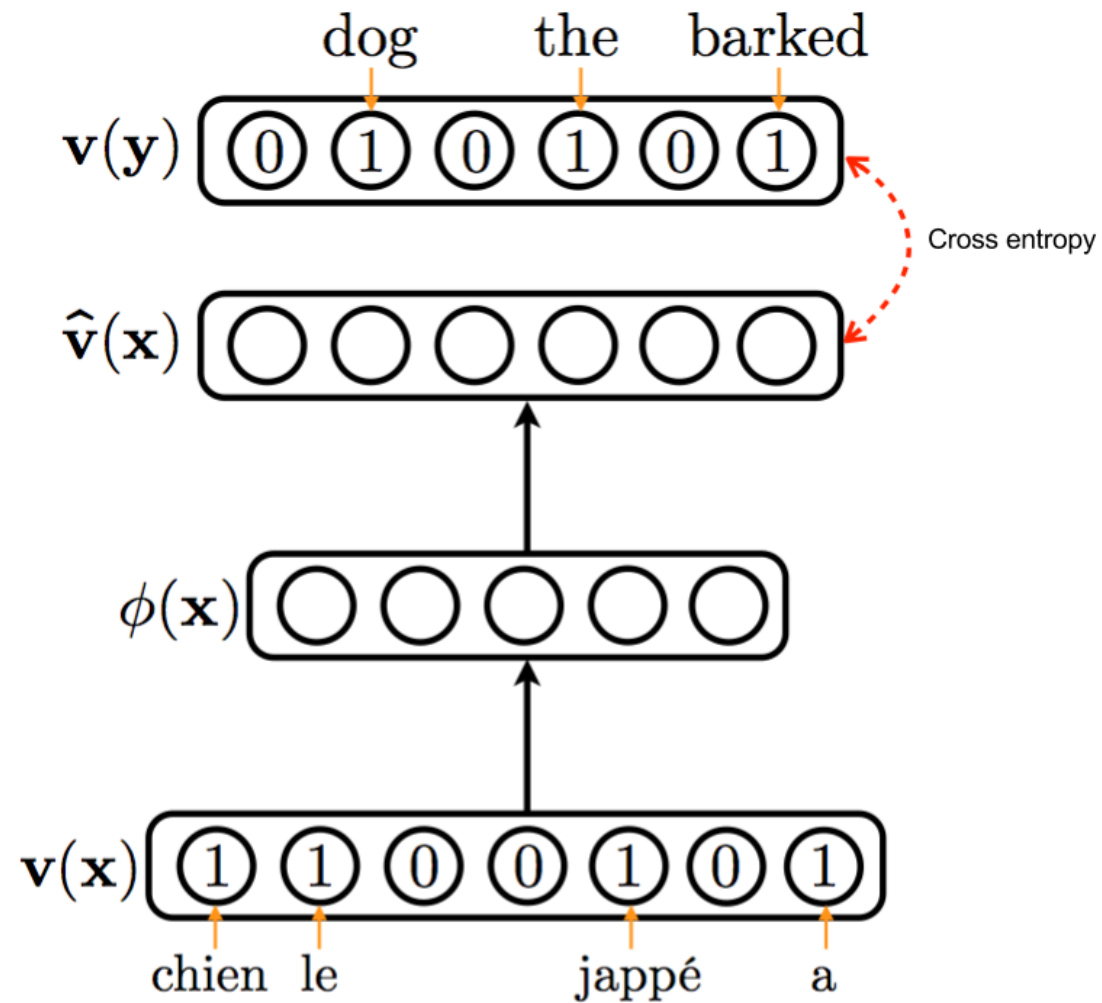
- Decoder

$$\hat{\mathbf{v}}(\mathbf{x}) = \text{sigm}(\mathbf{V}\phi(\mathbf{x}) + \mathbf{b})$$

- Cross entropy for the objective function

Binary bag-of-words reconstruction training

10



Tree-based decoder training

11

□ Encoder

$$\mathbf{a}(\mathbf{x}) = \mathbf{c} + \sum_{i=1}^{|\mathbf{x}|} \mathbf{W}_{\cdot, x_i}, \quad \phi(\mathbf{x}) = \mathbf{h}(\mathbf{a}(\mathbf{x}))$$

□ Decoder: Probabilistic binary tree for computing

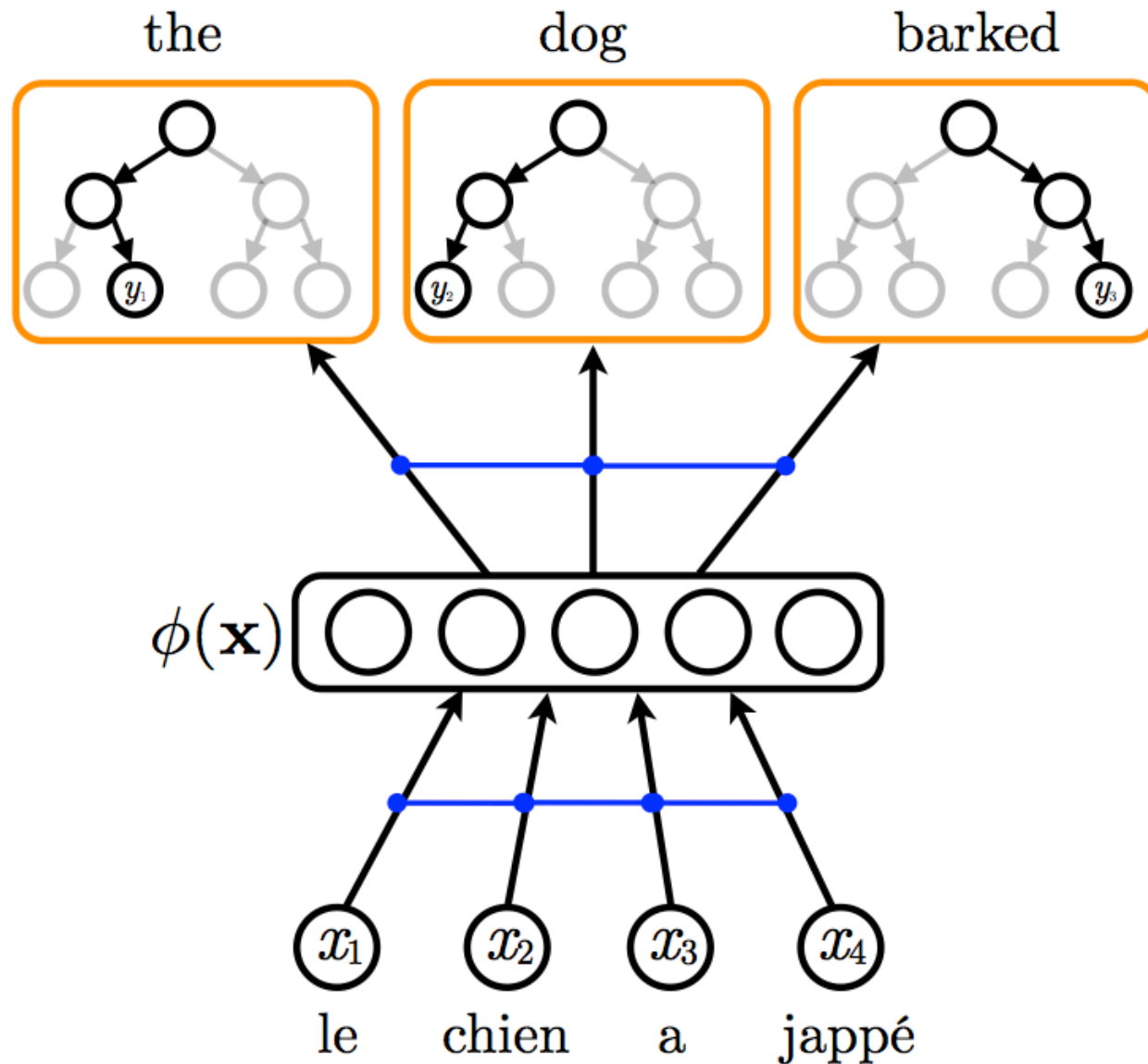
$$p(\hat{x} = x_i | \phi(\mathbf{x}))$$

▣ Each leaf of the tree is a word

□ Minimize the negative log probability

Tree-based decoder training

12



The models

13

- Reconstruction
 - x to x
 - y to y
 - x to y
 - y to x
- We can take advantage of unaligned data
- Maximize the correlation between the dimensions of the embedding's learned between languages

Cross-lingual classification task

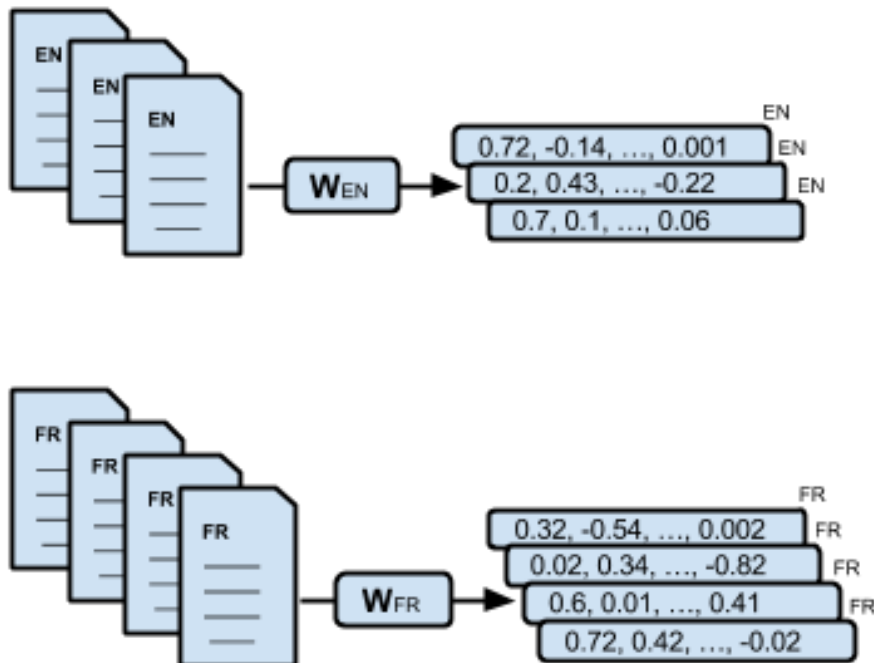
14

- We use the bilingual representation learned
- We want to classify documents for one language
- Labeled data available in an other language
 - ▣ Learn a classifier with one language
 - ▣ Apply that classifier on the other language

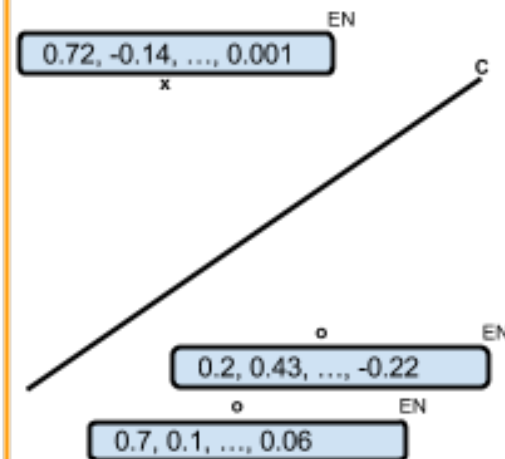
Classification task

15

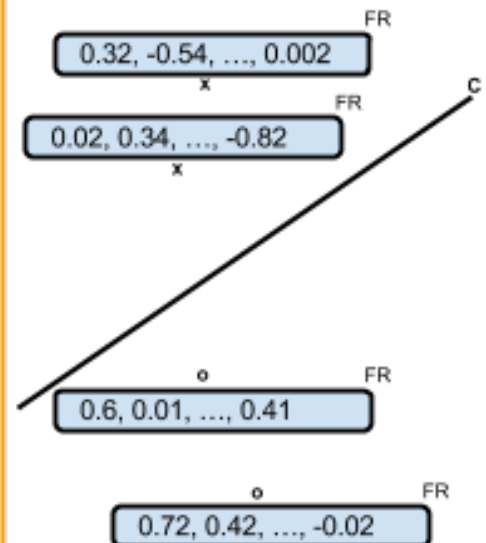
The model create a bilingual representation for documents in both languages



We learn a classifier on the representations of one language



We apply that classifier on the representations of the other language



Data

16

□ Europarl

- ▣ Transcript of the European parliament debate translated in multiple languages
- ▣ Dataset for learning the bilingual representation
- ▣ Each example is a pair (\mathbf{x}, \mathbf{y}) of the same sentence in two different language

□ Data: Reuters RCV1/RCV2

- ▣ Dataset for learning the classifier
- ▣ Different documents labeled for each language

Results

17

- Cross-lingual classification accuracy for 3 different pairs of languages, with 1 000 labeled examples.

| | EN → DE | DE → EN | EN → FR | FR → EN | EN → ES | ES → EN |
|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| BAE-tr | 81.8 | 60.1 | 70.4 | 61.8 | 59.4 | 60.4 |
| BAE-cr | 91.8 | 74.2 | 84.6 | 74.2 | 49.0 | 64.4 |
| Klementiev et al. | 77.6 | 71.1 | 74.5 | 61.9 | 31.3 | 63.0 |
| MT | 68.1 | 67.4 | 76.3 | 71.1 | 52.0 | 58.4 |
| Majority Class | 46.8 | 46.8 | 22.5 | 25.0 | 15.3 | 22.2 |

Results

18

- Cross-lingual classification accuracy for 1000 labeled examples (with 500K in training set).

| | EN -> DE | DE -> EN |
|---------------------|--------------|--------------|
| Hermann and Blunsom | 83.7% | 71.4% |
| BAE-cr | 87.9% | 76.7% |

Results

19

- Example English words along with the closest words both in English (EN) and German (DE), using the Euclidean distance between the embeddings learned by BAE-cr.

| Word | Lang | Nearest neighbors | Word | Lang | Nearest neighbors |
|-----------|------|------------------------------|-----------|------|-----------------------------------|
| january | EN | january, march, october | oil | EN | oil, supply, supplies, gas, fuel |
| | DE | januar, märz, oktober | | DE | öl, boden, befindet, gerät, erdöl |
| president | EN | president, i, mr, presidents | microsoft | EN | microsoft, cds, insider, ibm |
| | DE | präsident, präsidentin | | DE | microsoft, cds, warner |
| said | EN | said, told, say, believe | market | EN | market, markets, single |
| | DE | gesagt, sagte, sehr, heute | | DE | markt, marktes, märkte |

NADE Language Model

20

- Combine two approach
 - ▣ DocNADE (A Neural Autoregressive Topic Model)
 - ▣ Neural language model

DocNADE (A Neural Autoregressive Topic Model)

21

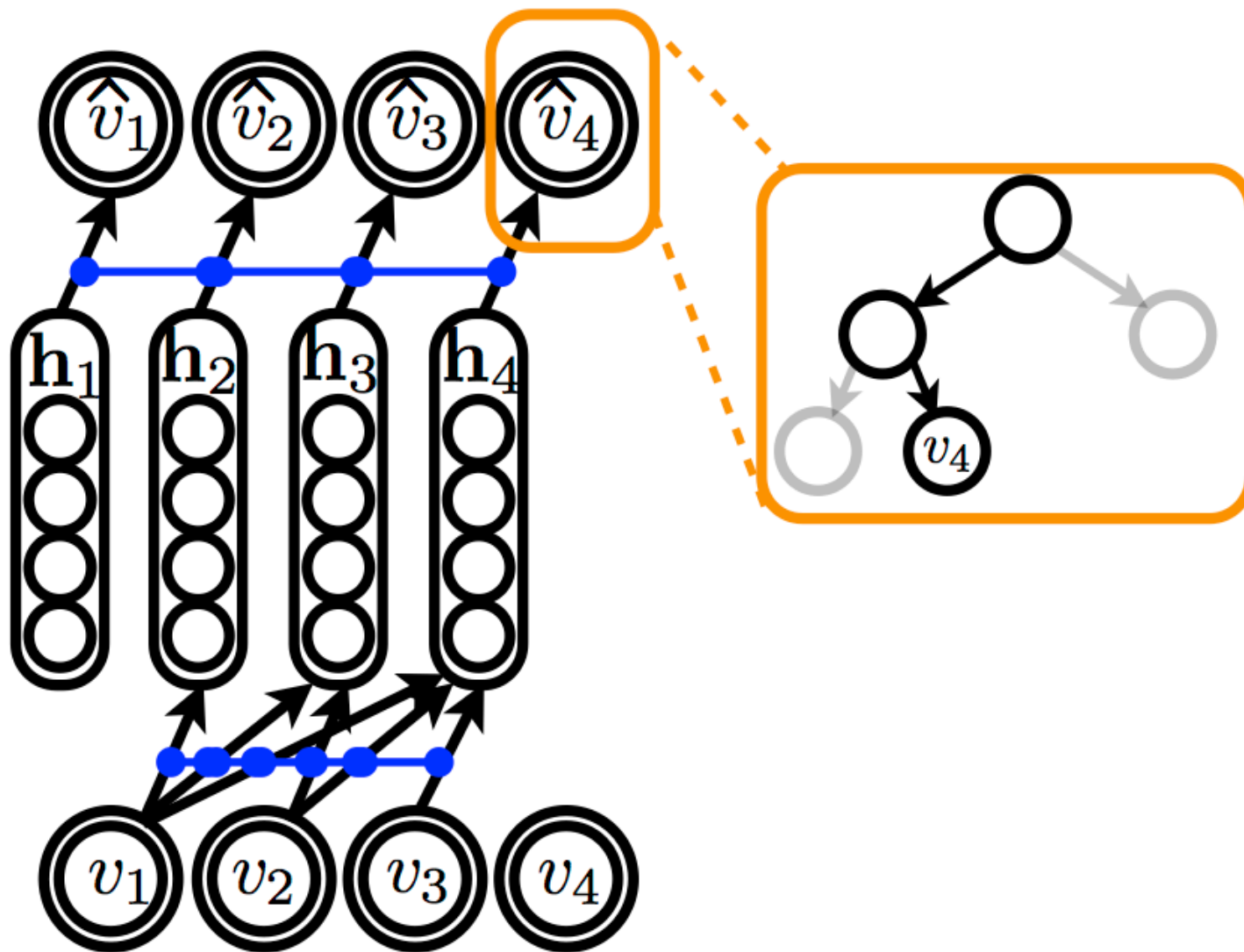
- Generative model.
- Multinomial observations $\mathbf{v} \in \{1, 2, \dots, N\}^D$.
- Extract a representation \mathbf{h}_i of all previous words $\mathbf{v}_{<i}$ for each v_i .

$$\mathbf{h}_i(\mathbf{v}_{<i}) = \text{sigm}(\mathbf{c} + \sum_{k < i} W_{:,v_k})$$

$$\mathbf{h}_{i+1}(\mathbf{v}_{<i+1}) = \text{sigm}(W_{:,v_i} + \mathbf{c} + \sum_{k < i} W_{:,v_k})$$

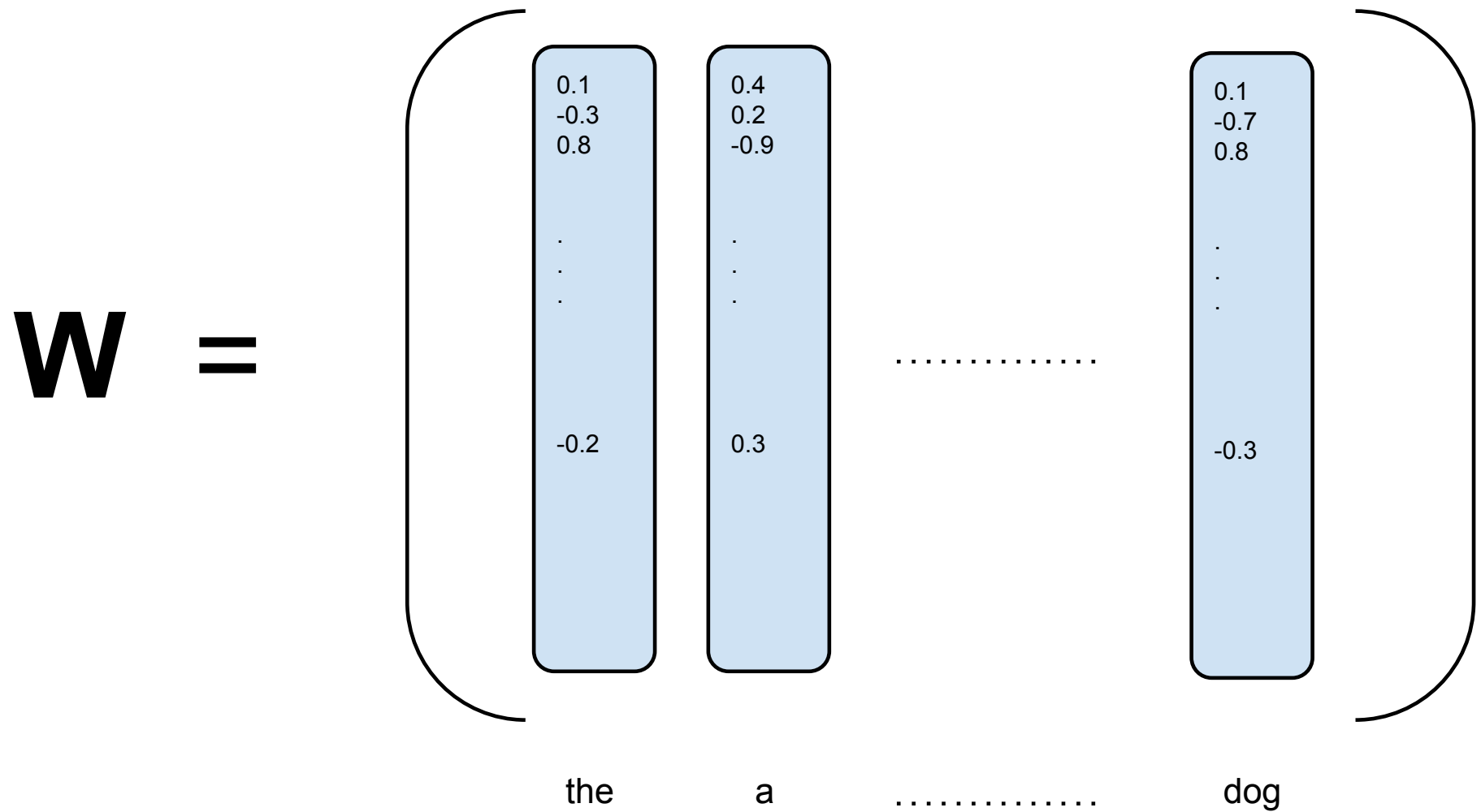
DocNADE

22



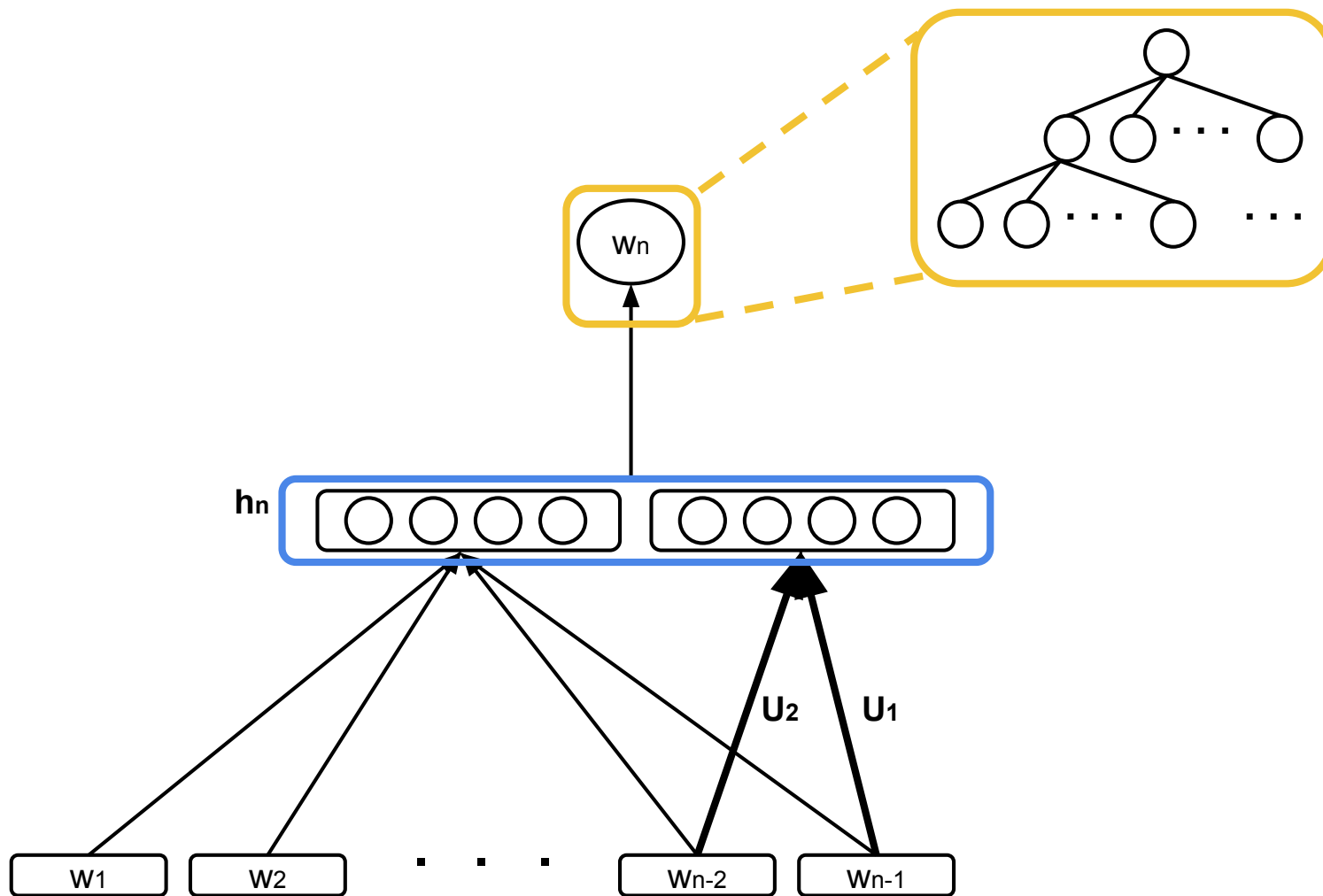
Vectorial representation

23



NADE Language Model

24



Preliminary results

25

□ Perplexity on AP news dataset

| | Perplexity |
|---------------------|------------|
| Mnih & Hinton | 112.1 |
| NADE language model | 113.8 |

Thank you !

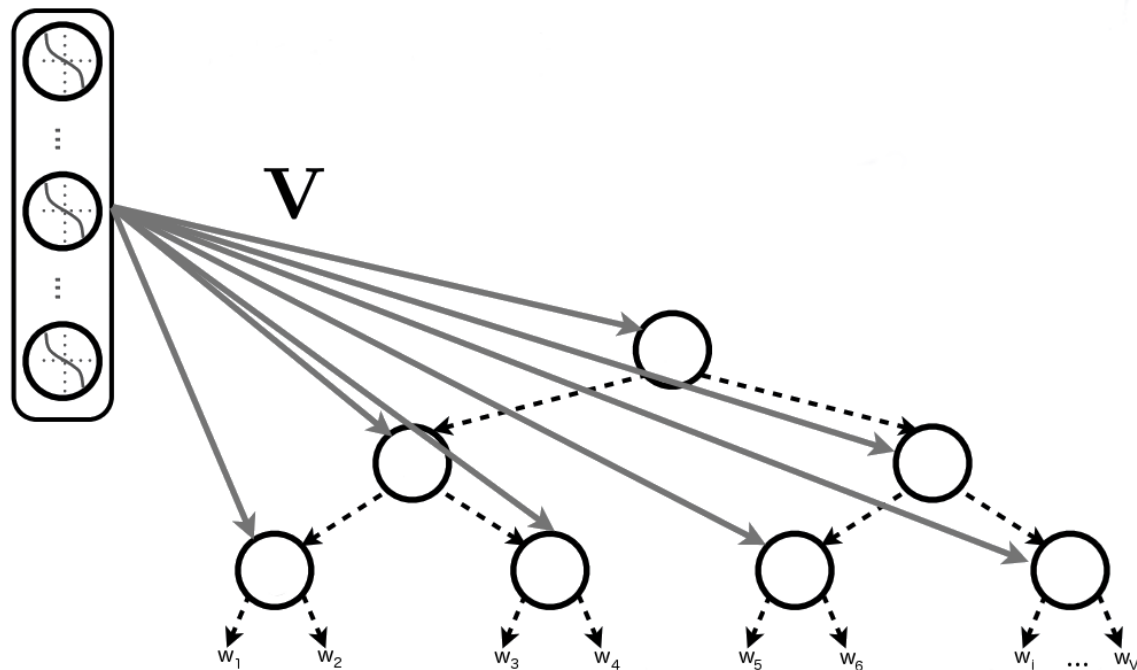
26

Question time !

Tree-based decoder training

27

- Each leaf of the tree is a word
- $\mathbf{l}(w)$ is the sequence of nodes for the word w
- $\pi(w)$ is the sequence of left/right choice
- $\mathbf{l}(w)_1$ is the root $\pi(w)_1$ is its left/right choice



Tree-based decoder training

28

□ $V_{l(w)_m,:}$ is the binary logistic regression for the node $l(w)_m$

$$p(w | \mathbf{h}) = \prod_{m=1}^{|\pi(w)|} p(\pi(w)_m | \mathbf{h})$$

$$p(\pi(w)_m = 1 | \mathbf{h}) = \text{sigm}(b_{l(w)_m} + V_{l(w)_m,:} \mathbf{h})$$