

# Introduction à l'apprentissage profond en recherche d'Information

-

*Deep Learning and Information Retrieval*

- Retour sur la représentation des textes
  - représentation distribuée de mots, *word embedding*
- Brève introduction à l'apprentissage profond (Deep learning)
- Modèles neuronaux de RI (Neural IR Models)
- Modèles de transformateurs (*Transformers*)

Retour sur la représentation des textes :

Apprentissage de représentation de mots (textes)

- Bag of words, ngrams
  - Words as atomic units (noble, Brutus, ambiguous)
- Word representation
  - $(w_1, w_2, \dots, w_{|V|})$ , the whole vocabulary a million (or a billion) of words
  - Term vector :
    - term vector :  $(0, 0, 0, 0, 0, 0, \dots, 1, \dots, 0)$  (appelé **One hot vector**)
  - Text (document) is represented
    - $(0, 1, 0, 0, 1, 0, 0, \dots, 1, \dots, 0)$  (presence/absence of words)
    - $(0, 2, 0, 0, 3, 0, 0, \dots, 1, \dots, 5, \dots, 0)$  (term frequency)
    - $(0, 0.25, 0, 0, 0.29, 0, 0, \dots, 0.02, \dots, 0.6, \dots)$  : Weighted terms (tf\*idf)
  - Document collection is represented as a matrix

D : The noble Brutus  
was ambitious

	$d_o$	$d_1$	$d_2$	...	$d_j$	...	$d_{ D }$
$t_0$	1	0	0	0	2		1
$t_1$	0	0	0	0	1		1
$t_2$	1	1					
...							
...							
$t_{ V }$	1	0	0				1

- More sophisticated word weighting

- Information Theory

- [Amati et al 2004], [S. Clinchant & Gaussier, 2010]

- Probabilistic Model

- BM25 [Robertson et al. 1994]

- Statistical Language model [Ponte & Croft 1998], [ChengXiang et al 2002]

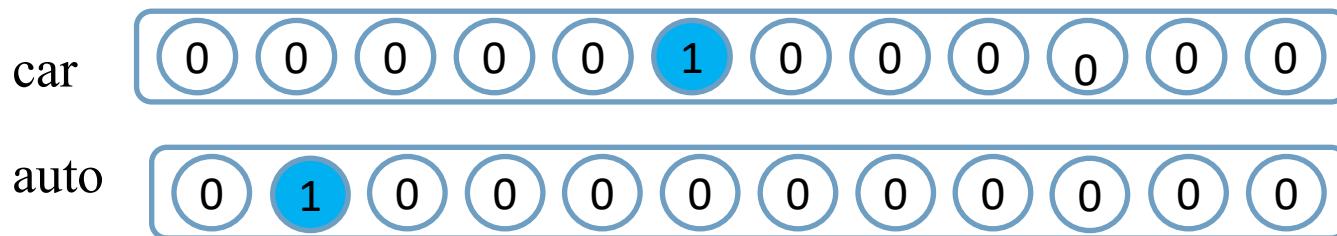
- Assigns a probability to a sequence of words

- $p(\text{"Today is Wednesday"}) \approx 0.001$

- $p(\text{"Today Wednesday is"}) \approx 0.000000000001$

- Limits of bag of words (One Hot vector)

- Synonymy, different word with the same meaning “Car” and “auto”,
  - Polysemy, same word but different meaning “salsa (dance)”, “salsa (sauce)”
  - Example :



- Beyond bag of words → Semantic representation (represent the meaning of a word)
  - Use thesaurus (WordNet) : containing a list of synonym sets and hypernyms
  - Distributional similarity based representations

## – Distributional similarity based representations

- "You shall know a word by the company it keeps"



Firth, J. R. (1957).

- Word is represented according to its context (its neighbors)

*...government debt problems turning into banking crises as happened in 2009...*

*...saying that Europe needs unified banking regulation to replace the hodgepodge...*

*...India has just given its banking system a shot in the arm...*

### • How to make neighbors represent words?

- Represent a word in a dense (low dimensional) vector (25 - 1000) built (learnt)
- LSA (LSI)based on Matrix Factorization (SVD) [S. Deerwester et al 1988]
- Neural Word Embedding: Word2Vec (CBOW, Skip-gram) [T. Mikolov et al 2013]

“car” : [0.015 , 0.001, -0.020, 0.015, 0.001]

Représentation continue (distribuée) des mots → Vers une représentation sémantique des mots

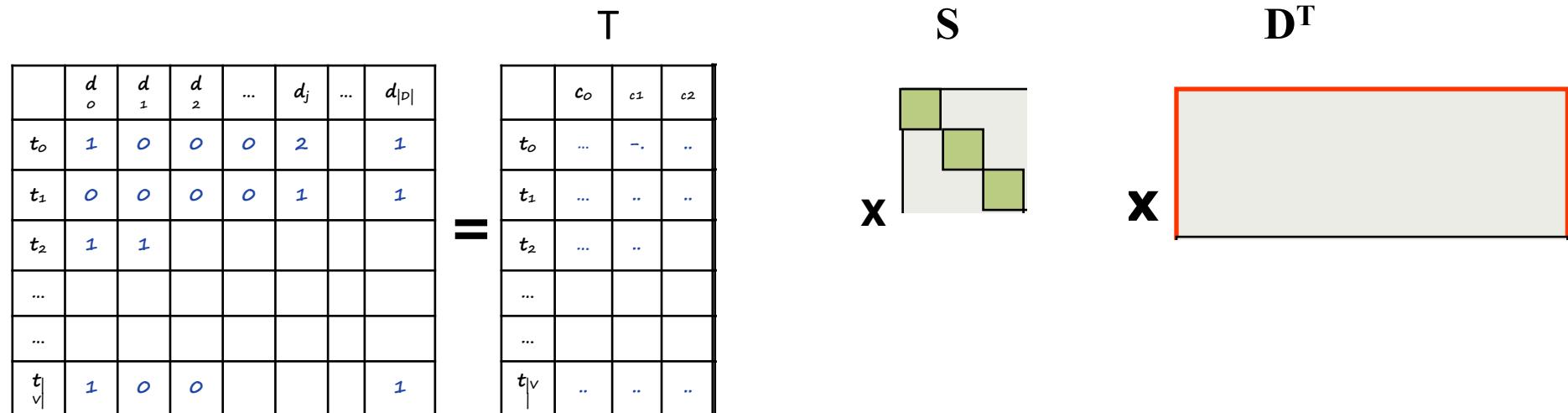
## • Illustration de la SVD

$$\underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_A = \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{V^T}$$

$$\underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_A = \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{V^T}$$

## • SVD (Singular Value Decomposition) $A=TS\mathbf{D}^T$

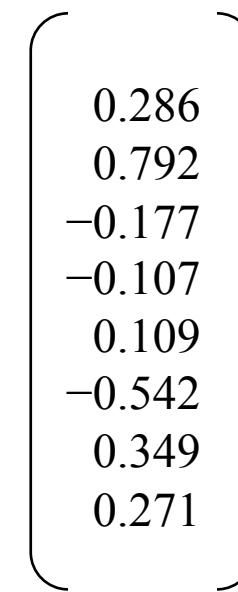
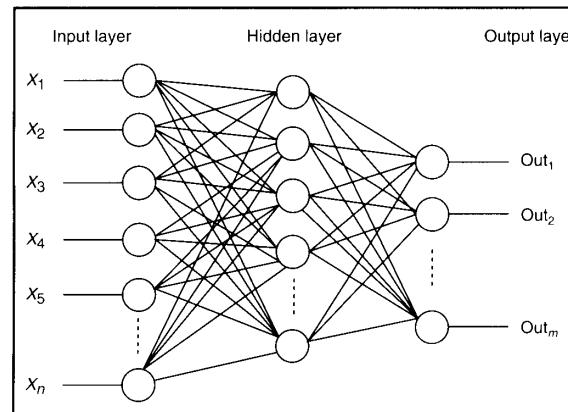
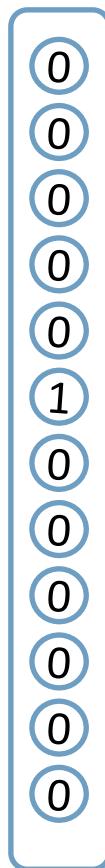
Select k (300) top singular values



- La matrice  $D$  représente les documents dans le nouvel espace vectoriel (espace des concepts)
- La matrice  $T$ : la représentation sémantique d'un terme
- La fonction qui permet le passage des termes vers les concepts est  $M = T[v,k].S^{-1}[k,k]$

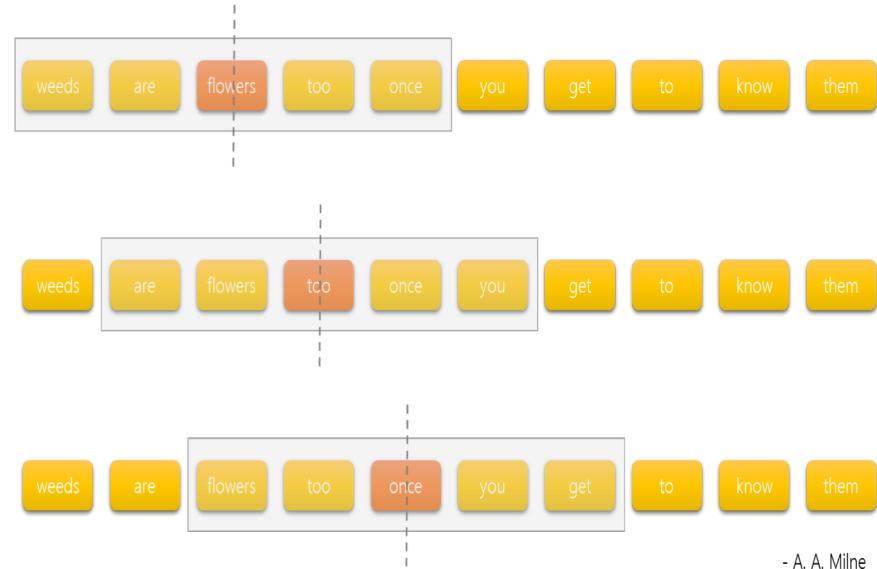
- Utiliser des réseaux de neurones pour apprendre des représentations continues (denses de mots)

car

 $|V| = 10^6$  $|V_{dense}| = 100$

## Word2vec

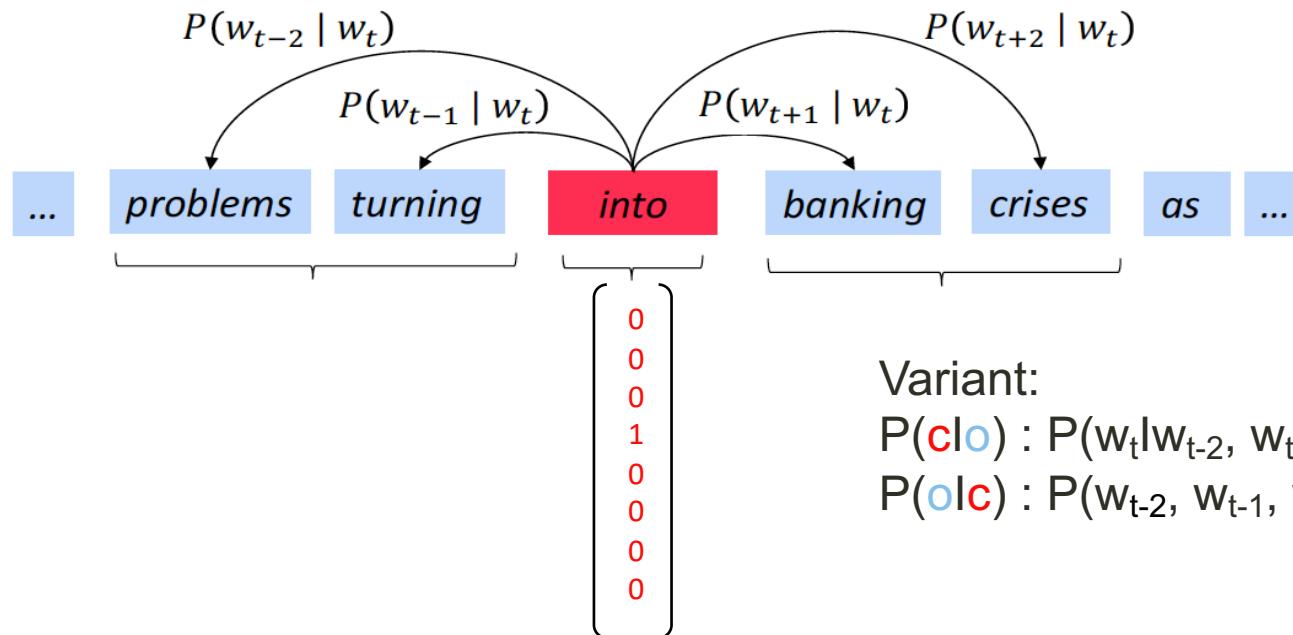
- Goal: simple (shallow) neural model learning from a large corpus of text (billion words)
- Every word in a fixed vocabulary is represented by a vector
- Go through each position  $t$  in the text, which has a center word  $c$  and context (“outside”) words  $o$
- Predict center ( $c$ ) word from its context( $o$ ) ( $P(o/c)$ ) or vice-versa ( $P(c/o)$ )



- A. A. Milne

- Word2vec (Mikolov et al. 2013), the main Idea:

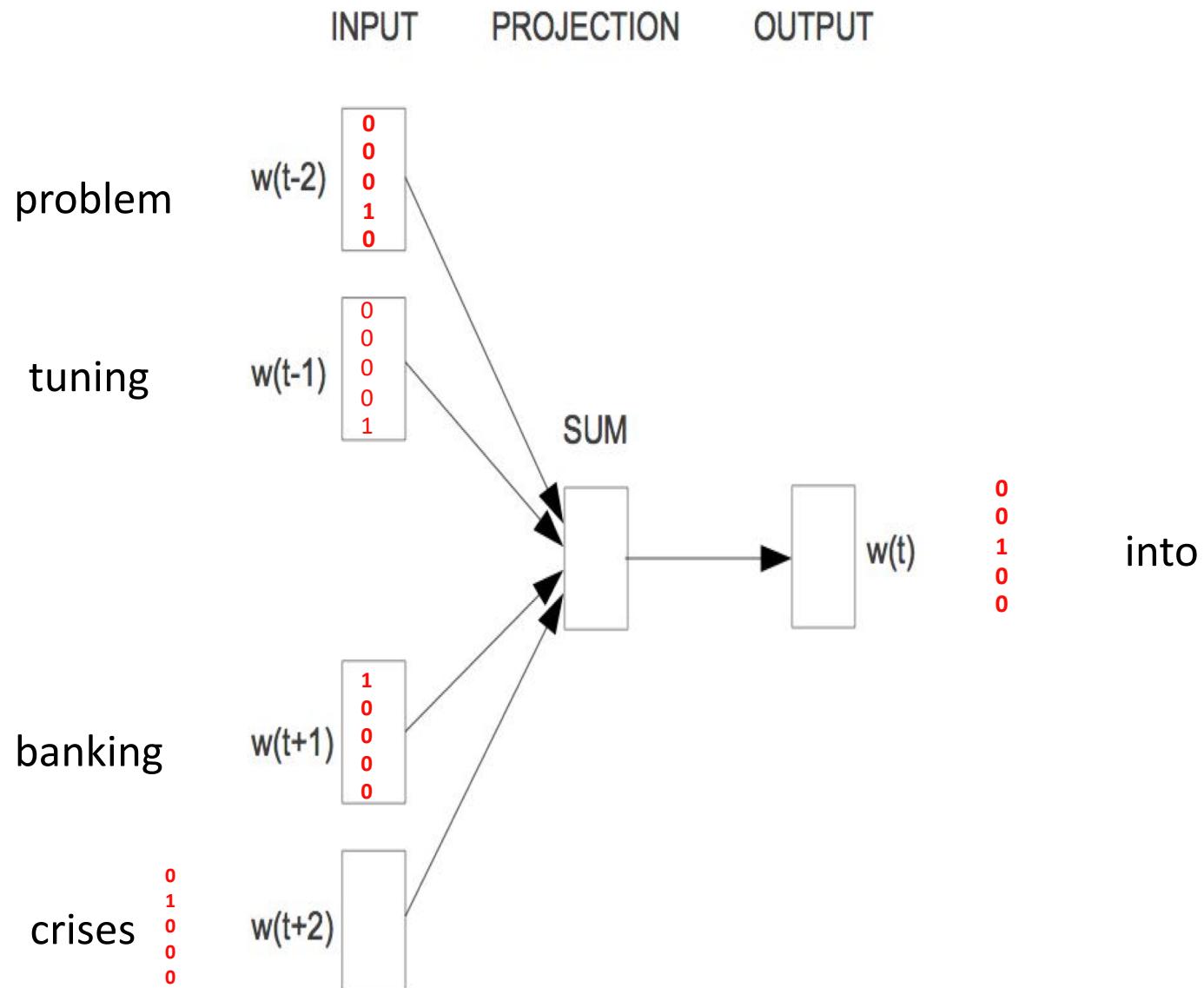
- Go through each position in the text, which has a **center word  $w_t$**  and **context ("outside") words  $w_{t-1}$**
- Use the similarity of the word vectors for **c** and **o** to calculate the probability  $P(\text{olc})$  or  $P(\text{clo})$



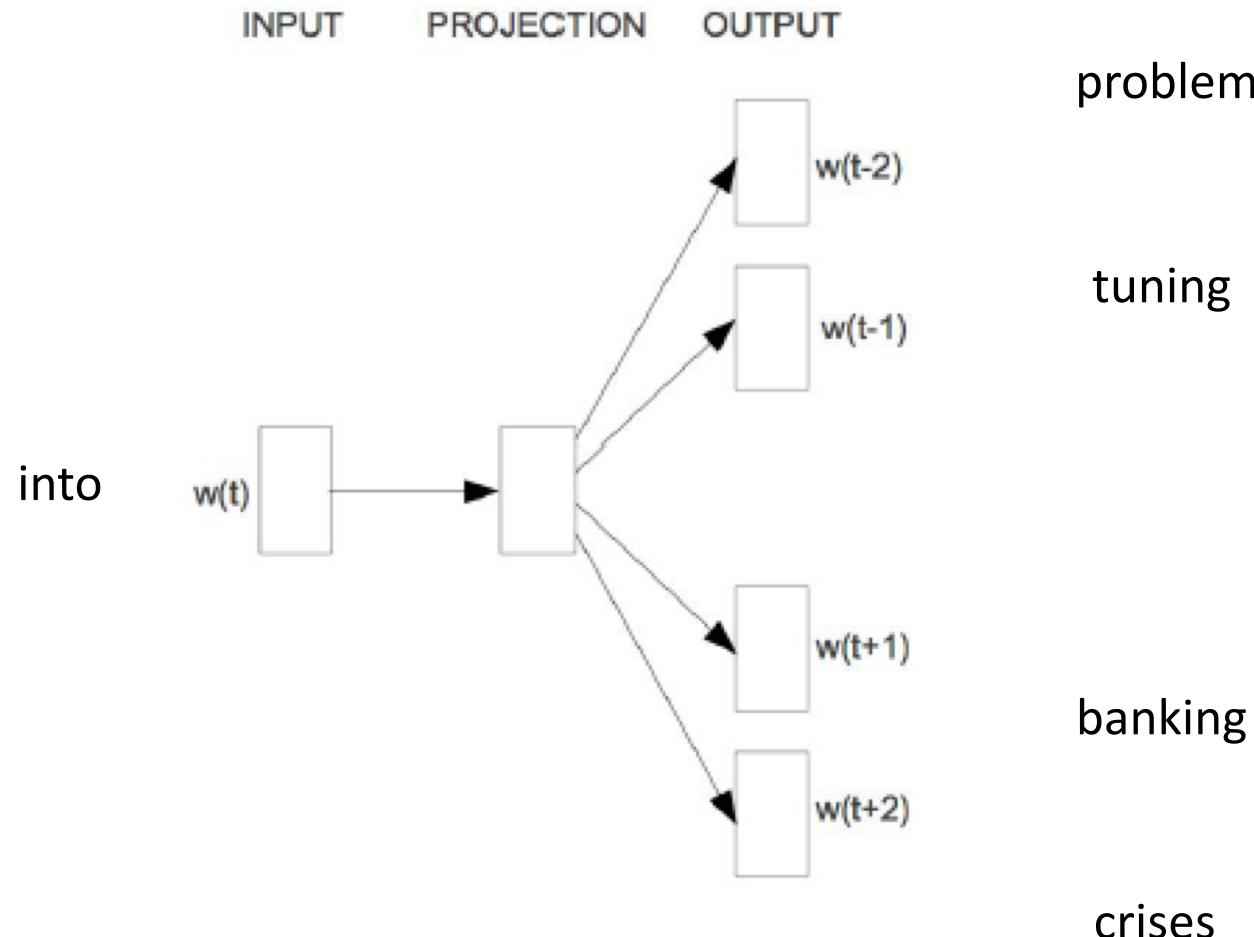
Variant:

$$\begin{aligned} P(\text{clo}) &: P(w_t | w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2}) \\ P(\text{olc}) &: P(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} | w_t) \end{aligned}$$

## • CBOW



## Skip-Gram



# Skipgram

$V \times 1$      $d \times V$      $d \times 1$

$w_t$      $W$      $= v_c$   
 $= Ww_t$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} \dots & 0.2 & \dots \\ \dots & -1.4 & \dots \\ \dots & 0.3 & \dots \\ \dots & -0.1 & \dots \\ \dots & 0.1 & \dots \\ \dots & 0.5 & \dots \end{bmatrix} \quad \begin{bmatrix} 0.2 \\ -1.4 \\ 0.3 \\ -0.1 \\ 0.1 \\ 0.5 \end{bmatrix}$$

↑  
one hot  
word  
symbol  
  
Looks up  
column of  
word embedding  
matrix as

$V \times d$   
 $W'$   
 $u_2$

$$[u_x^T v_c] \quad \text{softmax}(u_x^T v_c)$$

6.7	0.07
6.3	6.1
0.1	0.05
-6.7	0.01
-0.2	0.02
0.1	0.05
0.7	0.7

Truth  
 $w_{t-3}$

0
0
0
0
1
0

Softmax

$$p_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\text{softmax} \rightarrow \begin{bmatrix} 0.07 \\ 6.1 \\ 0.05 \\ 0.01 \\ 0.02 \\ 0.05 \\ 0.7 \end{bmatrix}$$

$w_{t-2}$

0
1
0
0
0
0
0

Actual  
context  
words

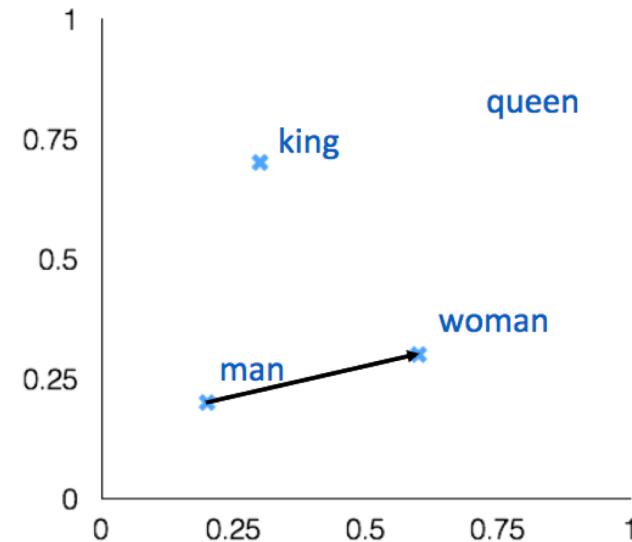
$$\text{softmax} \rightarrow \begin{bmatrix} 0.07 \\ 6.1 \\ 0.05 \\ 0.01 \\ 0.02 \\ 0.05 \\ 0.05 \end{bmatrix}$$

$w_{t-1}$

0
0
0
0
0
0
0

Output

man:woman :: king:?		
+	king	[ 0.30 0.70 ]
-	man	[ 0.20 0.20 ]
+	woman	[ 0.60 0.30 ]
<hr/>		
	queen	[ 0.70 0.80 ]



Test word analogies [Mikolov 2014]

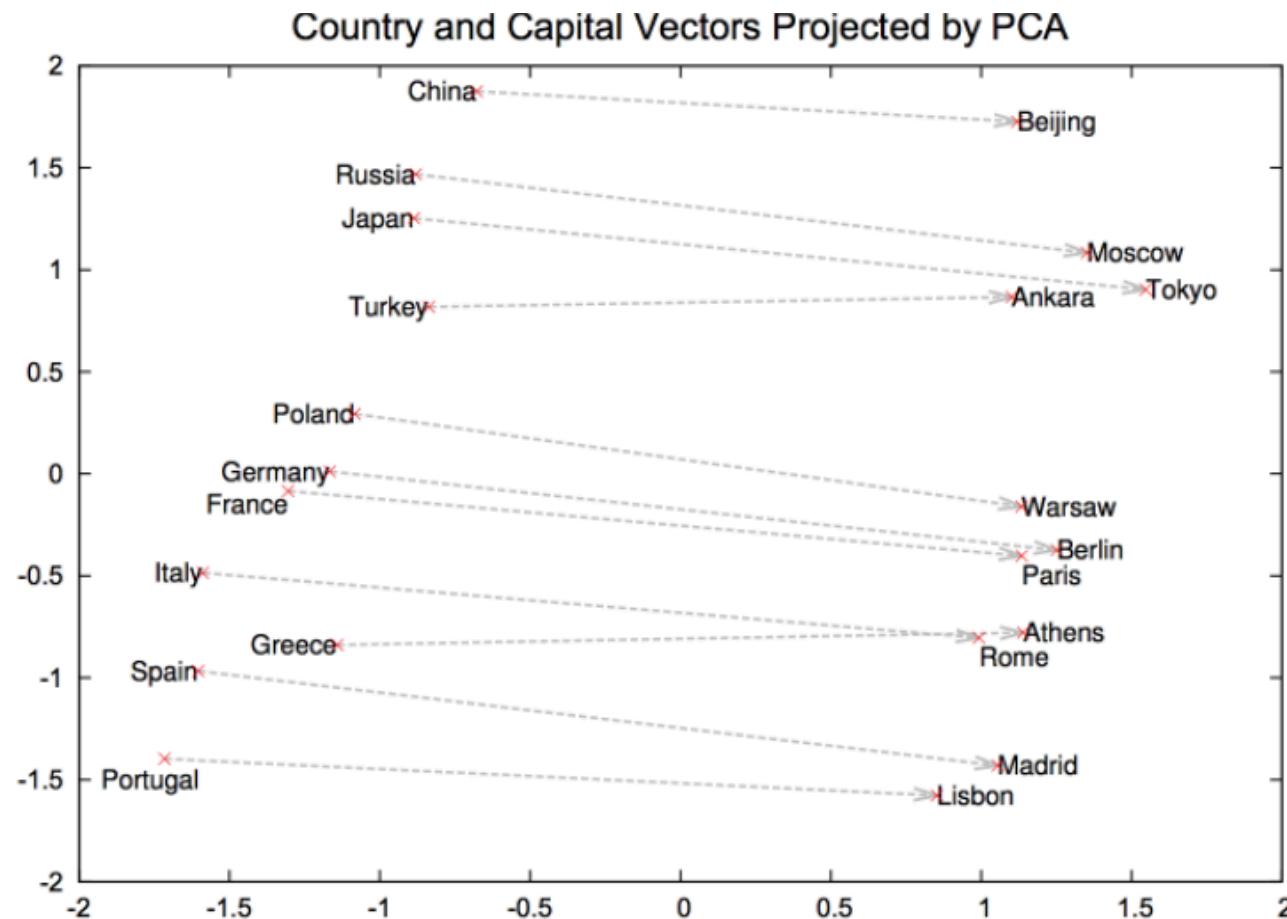


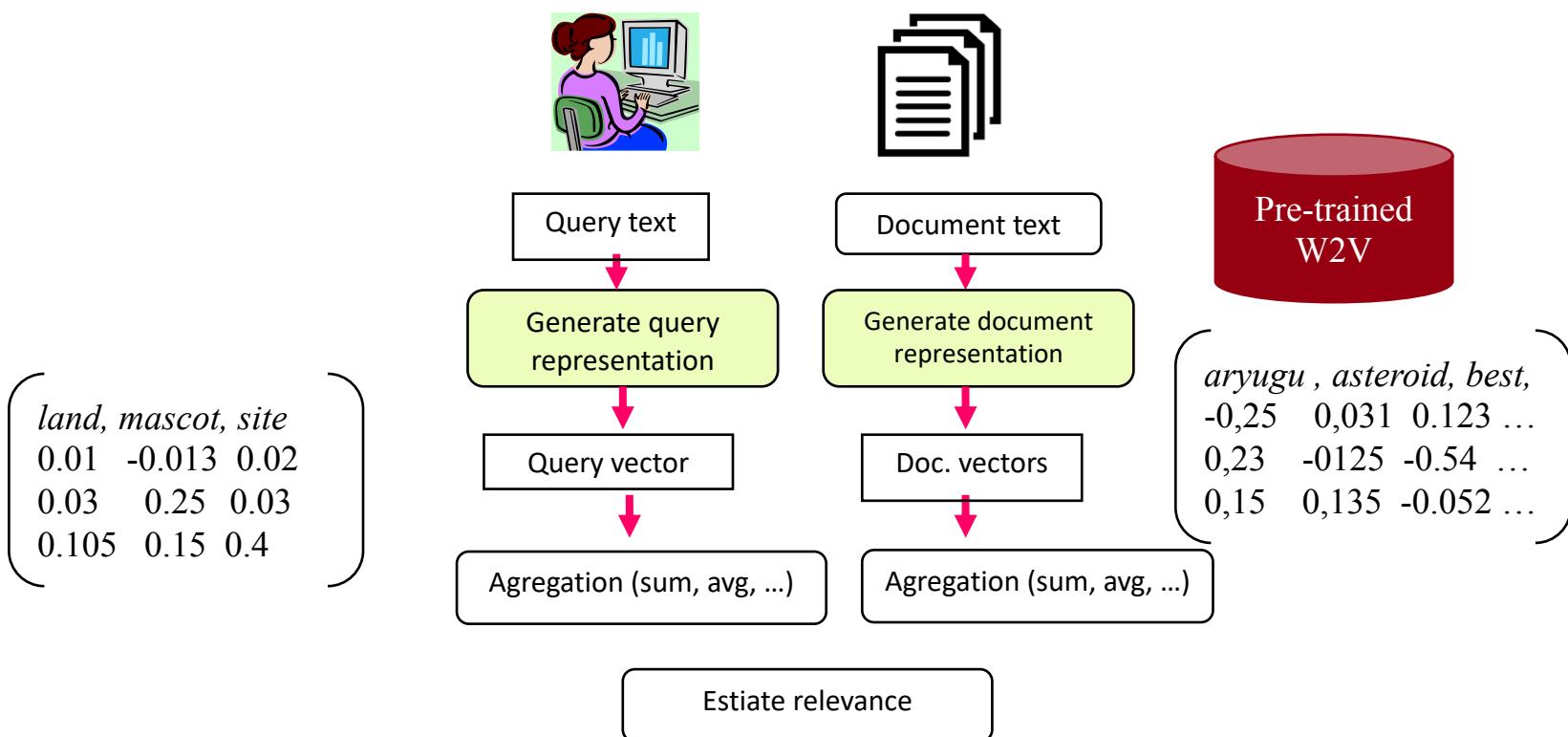
Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities. The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

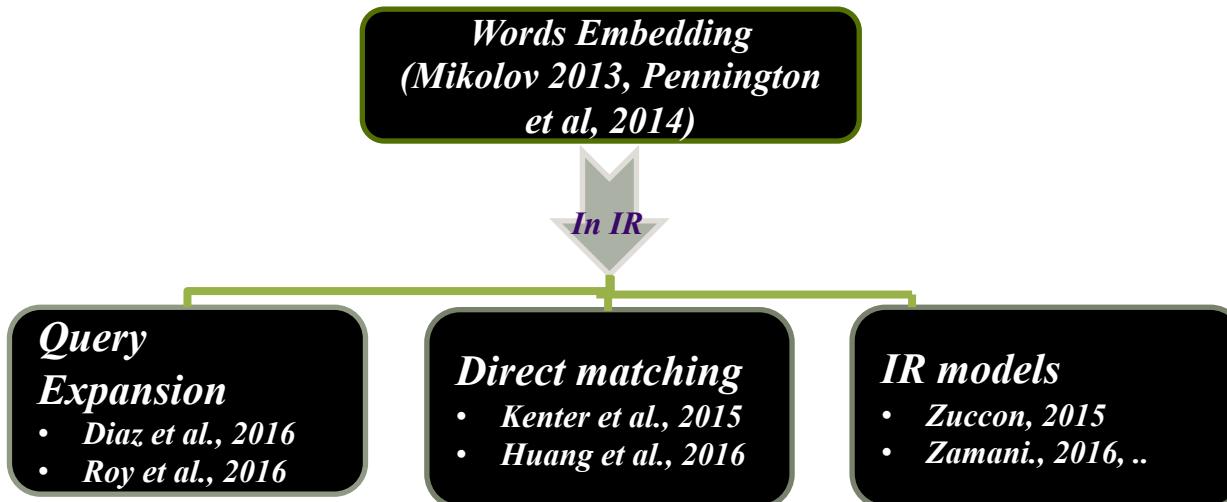
- Rappel de quelques definitions sur l'apprentissage profond
- Deep learning et RI
  - Retour sur la representation des documents
  - Utilisation des *Word Embedding* en RI
    - Un peu de pratique Python

- Un peu de pratique en Python
  - On peut construire son propre Word2Vec (opération fastidieuse) (il faut beaucoup, beaucoup de textes)
  - Utiliser des Word2Vec pré-entraînés (Google, Glove, ...)

- Une bonne majorité des modèles de représentation de textes utilise la représentation distribuée (continue) des termes.
- On peut apprendre différents types de représentations: termes, sous-termes (*n*-grammes), caractères, etc.
- Beaucoup de modèles de RI exploitent ce type de représentations

- Extend traditional IR models
  - Term weighting, language model smoothing, translation of vocab
- Expand query using embeddings (followed by non-neural IR)
  - Add words similar to the query
- IR models that work in the embedding space
  - Centroid distance, word mover's distance





$$rsv(q, d) = \sum_{t \in q} \sum_{t' \in V} P(t, t').P(t'/d)$$

Task	Studies
Ad-hoc Retrieval	ALMasri et al. (2016), Amer et al. (2016), BWESG (Vulic and Moens (2015)), Clinchant and Perronnin (2013), Diaz et al. (2016), GLM (Ganguly et al. (2015)), Mitra et al. (2016), Nalnisnick et al. (2016), NLTM (Zucccon et al. (2015)), Rekabsaz et al. (2016), Roy et al. (2016), Zamani and Croft (2016a), Zamani and Croft (2016b), Zheng and Callan (2015)

- Dual Embedding Space Model (Mitra et al 2016)
  - A document is represented by a centroid of its word vectors

$$\bar{\mathbf{D}} = \frac{1}{|D|} \sum_{\mathbf{d}_j \in D} \frac{\mathbf{d}_j}{\|\mathbf{d}_j\|}$$

- Query document similarity is average over query words cosine similarity

$$DESM(Q, D) = \frac{1}{|Q|} \sum_{q_i \in Q} \frac{\mathbf{q}_i^T \bar{\mathbf{D}}}{\|\mathbf{q}_i\| \|\bar{\mathbf{D}}\|}$$

- Extension des modèles de RI : Modèle translation de RI

- $rsv(q, d) = \sum_{t \in q} sim(t, t').P(t'|d)$

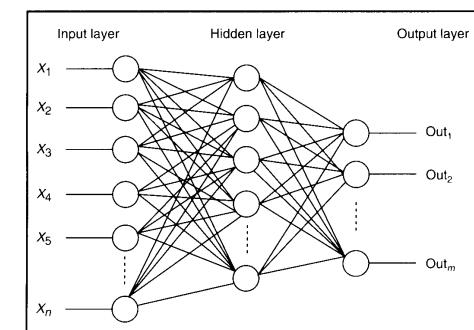
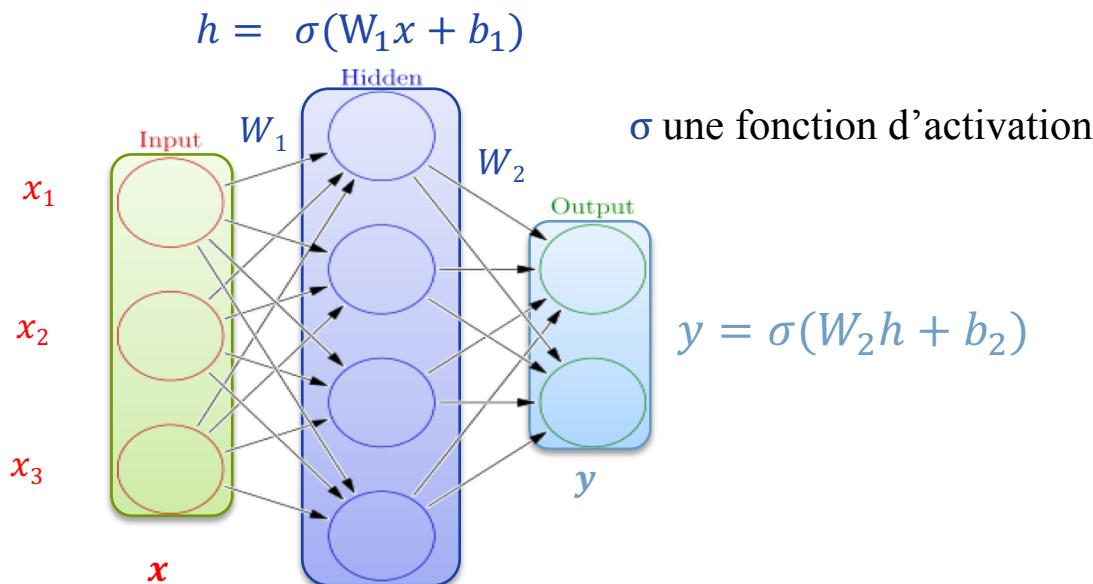
Compares query term with every document term

$$P(t|d) = \lambda P(t|d) + \alpha \sum_{t' \in d} P(t, t'|d)P(t') + \beta \sum_{t' \in N_t} P(t, t'|C)P(t') + (1 - \lambda - \alpha - \beta)P(t|C)$$

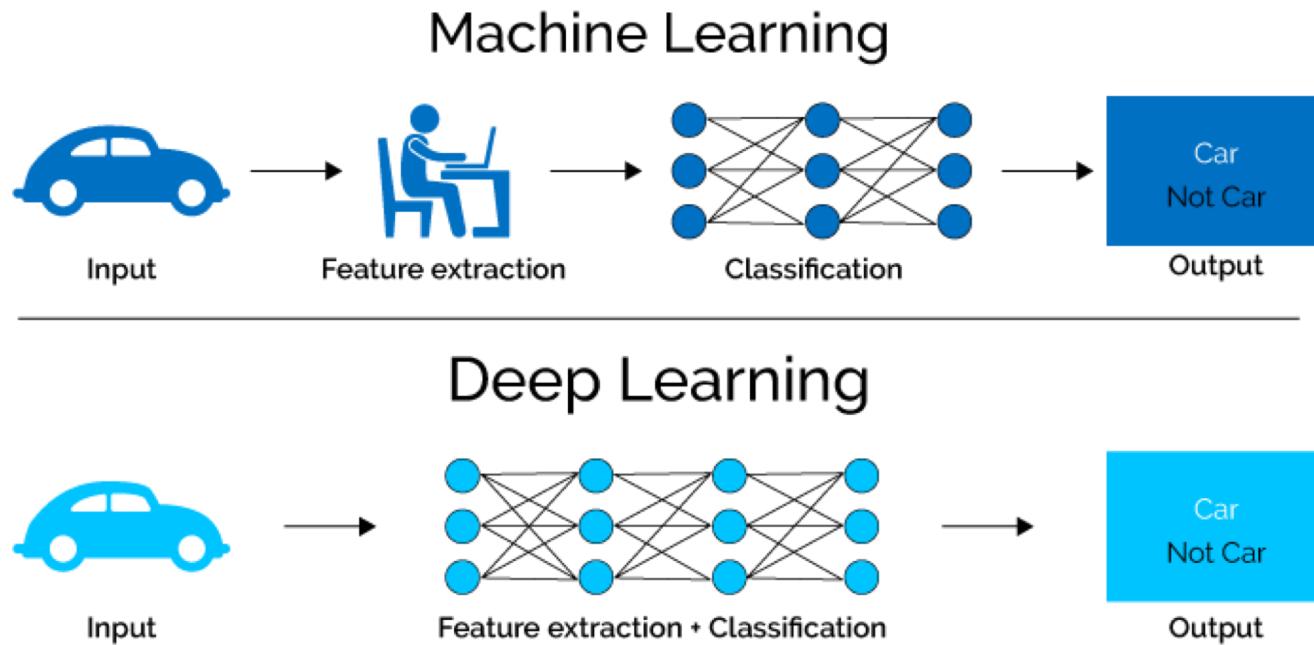
$$P(t, t'|d) = \frac{sim(t', t)}{\Sigma(d)} \frac{tf(t', d)}{|d|}$$

# Introduction à l'apprentissage profond ( Deep Learning)

- Apprentissage profond :
  - apprentissage basé sur des réseaux de neurones
  - objectif apprendre des représentations abstraites des données
- Au cœur de cet apprentissage : → Les réseaux de neurones
  - Un RN est une succession de couches de neurones, avec une entrée ( $X$ ), des couches cachées de neurones et une ou plusieurs sorties ( $Y$ ), le but est d'apprendre les «  $W_i$  » et «  $b_i$  »

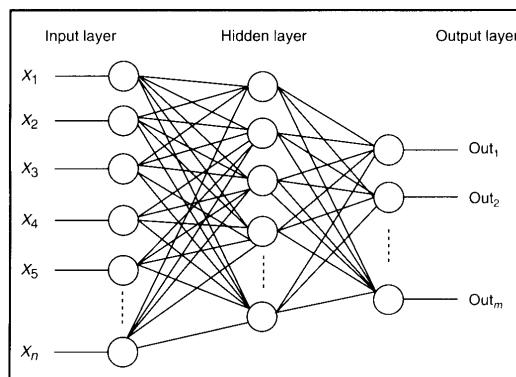


- Contrairement aux techniques d'apprentissage vues dans le cours précédent, les caractéristiques, données en entrée de l'algo d'apprentissage, sont construites automatiquement par apprentissage (on parle d'apprentissage de représentations)

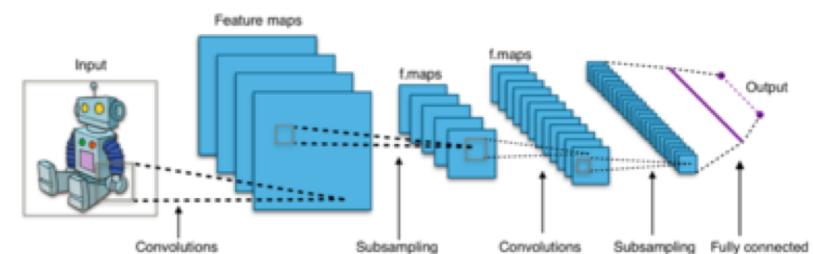


<https://www.xenonstack.com/blog/static/public/uploads/media/machine-learning-vs-deep-learning.png>

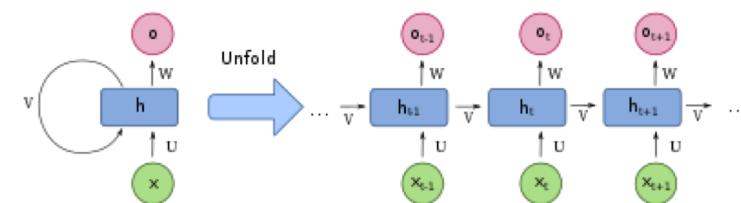
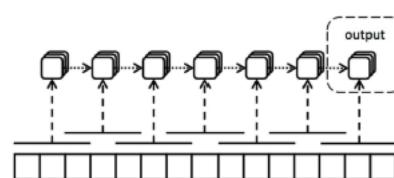
## Réseaux multicouches (multilayer perceptron)

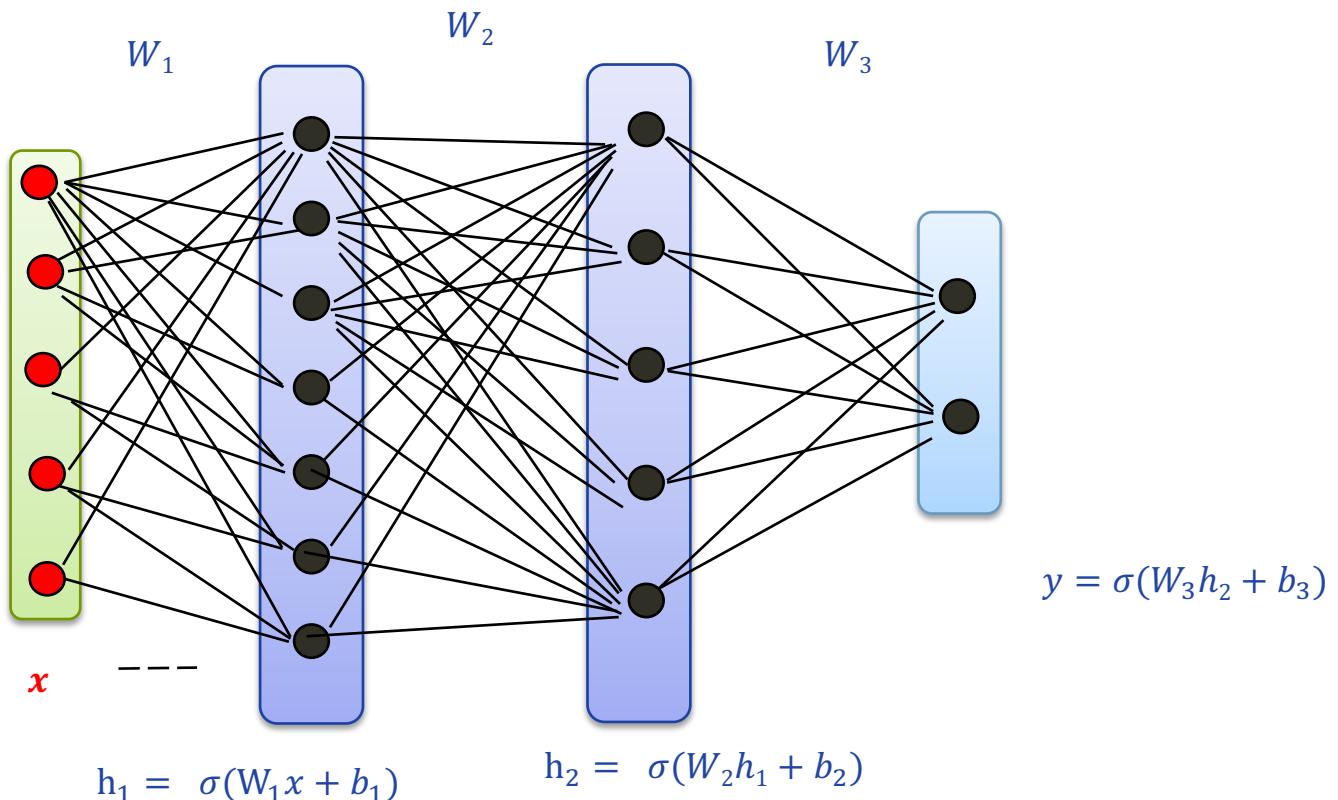


## Modèles à convolution ou convolutifs (CNN)



## Modèles récurrents (RNN)





$\sigma$  une fonction d’activation: *tanh, relu sigmoid, softmax, ...*

Voir un exemple sur github

These are the network parameters to be learned.

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

⋮ ⋮

Each filter detects a small pattern (3 x 3).

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

*Dot  
product*



1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

6 x 6 image

If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

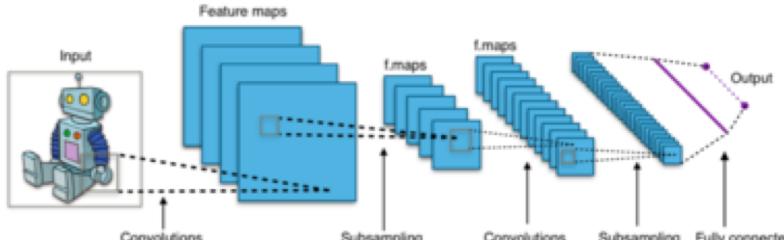
3

-3

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

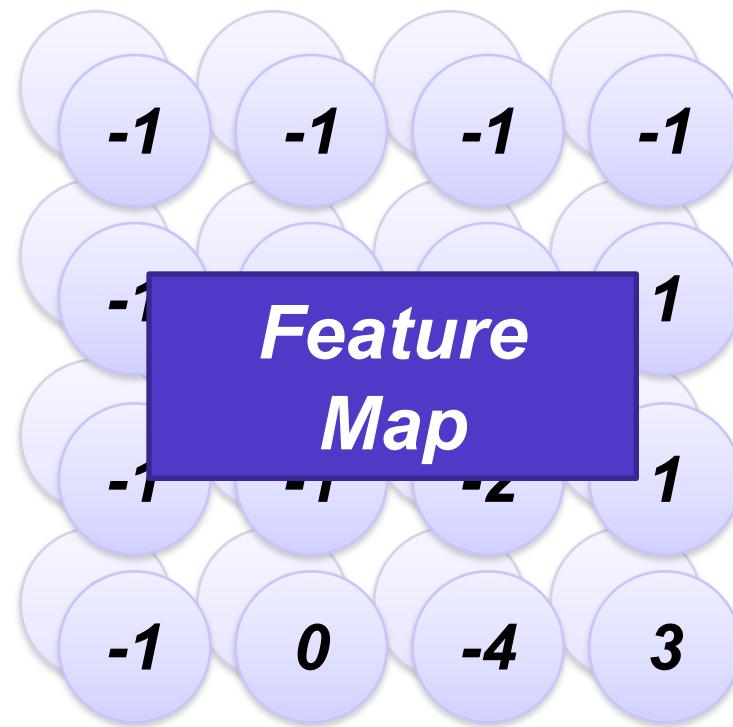
6 x 6 image



-1	1	-1
-1	1	-1
-1	1	-1

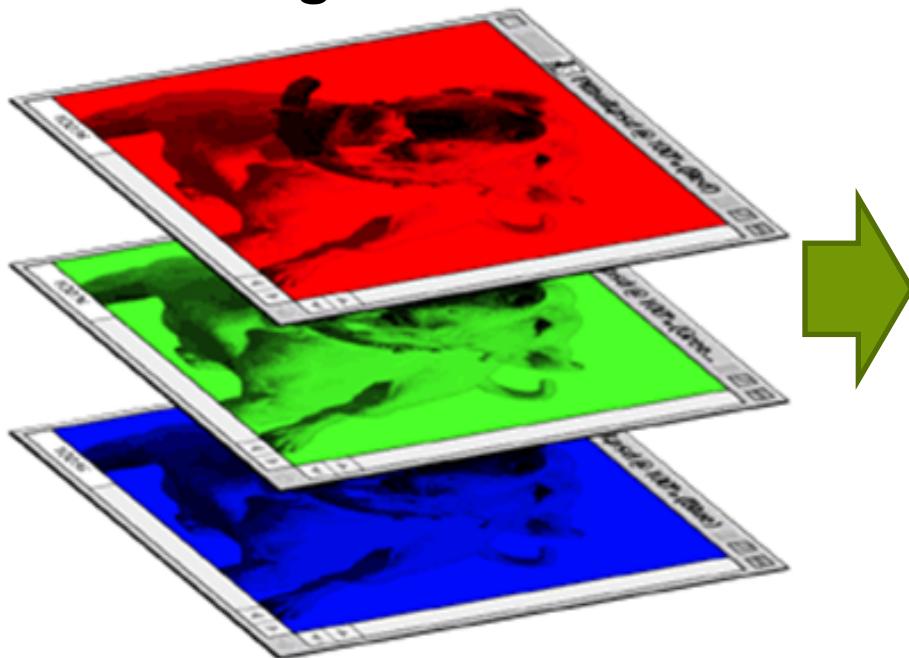
Filter 2

Repeat this for each filter



Two 4 x 4 images  
Forming 2 x 4 x 4 matrix

***Color image***



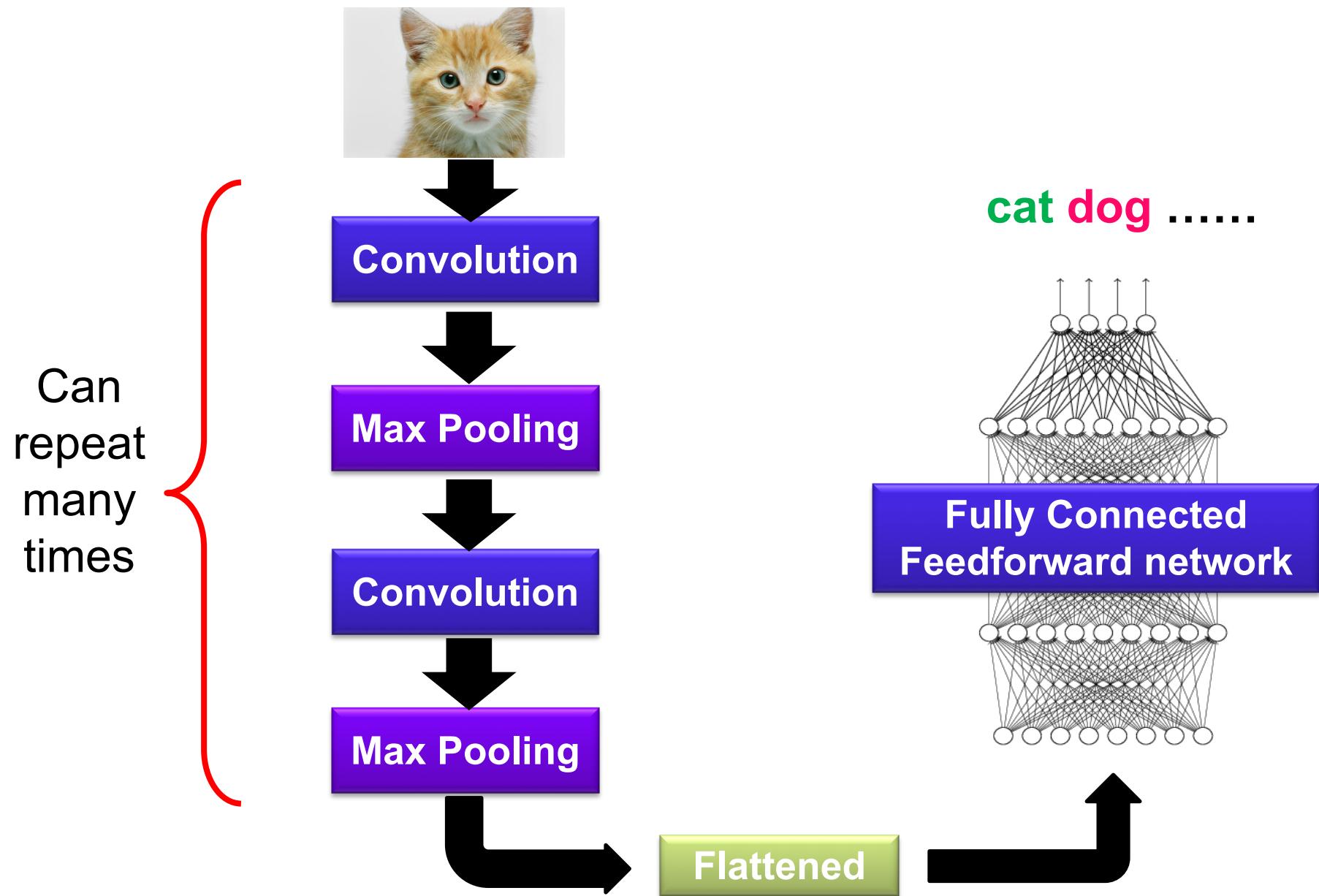
$$\begin{matrix} & & & & \\ & & & & \\ & & & & \\ \begin{matrix} 1 & -1 & -1 \\ -1 & 1 & -1 \\ -1 & -1 & 1 \end{matrix} & & & & \\ & & & & \\ & & & & \end{matrix}$$

Filter 1

$$\begin{matrix} & & & & \\ & & & & \\ & & & & \\ \begin{matrix} -1 & 1 & -1 \\ -1 & 1 & -1 \\ -1 & 1 & -1 \end{matrix} & & & & \\ & & & & \\ & & & & \end{matrix}$$

Filter 2

$$\begin{matrix} & & & & & & \\ & & & & & & \\ & & & & & & \\ \begin{matrix} 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \end{matrix} & & & & & & \\ & & & & & & \\ & & & & & & \end{matrix}$$

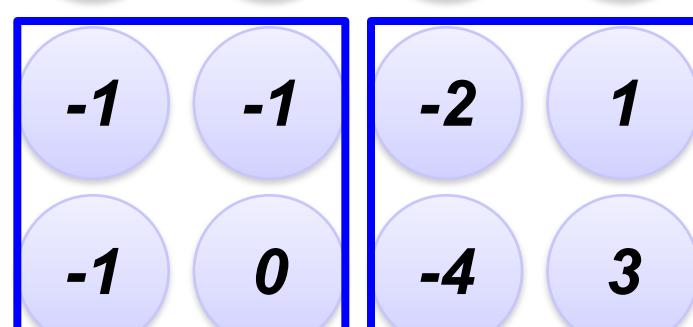
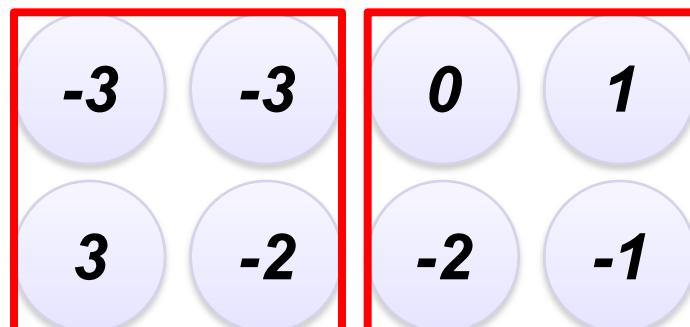
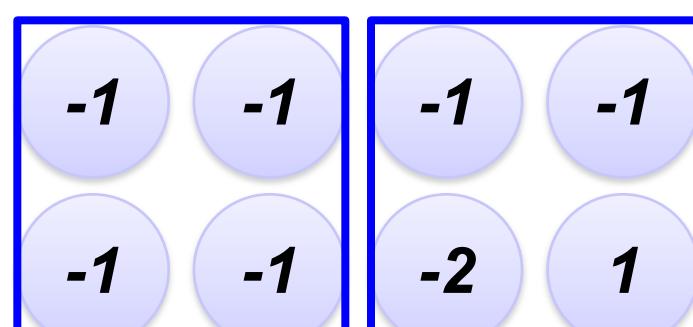
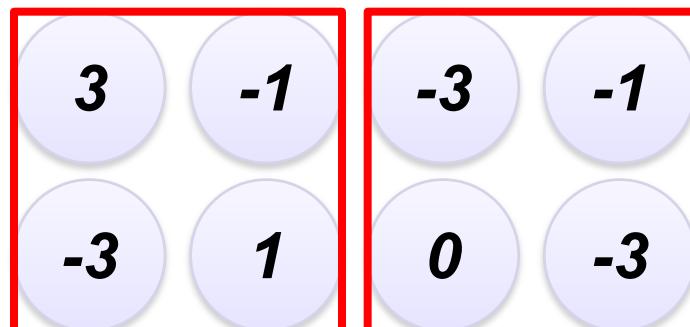


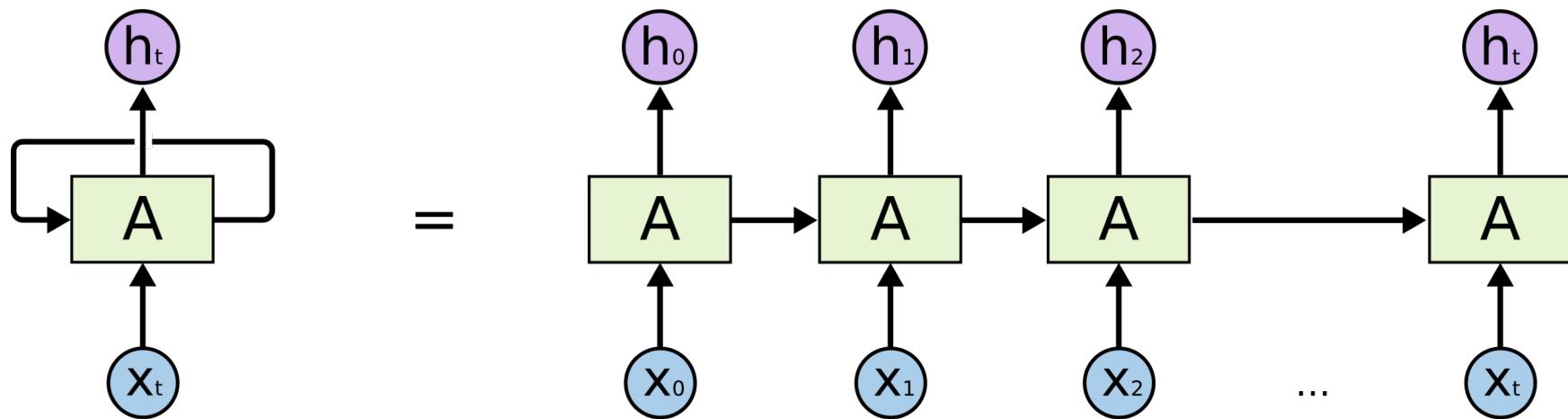
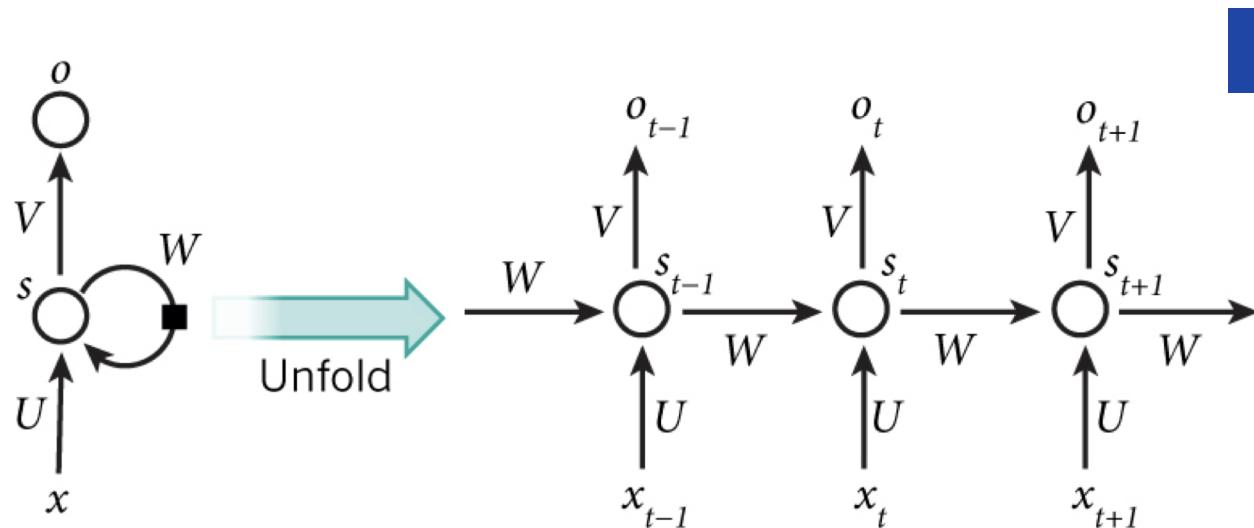
1	-1	-1
-1	1	-1
-1	-1	1

**Filter 1**

-1	1	-1
-1	1	-1
-1	1	-1

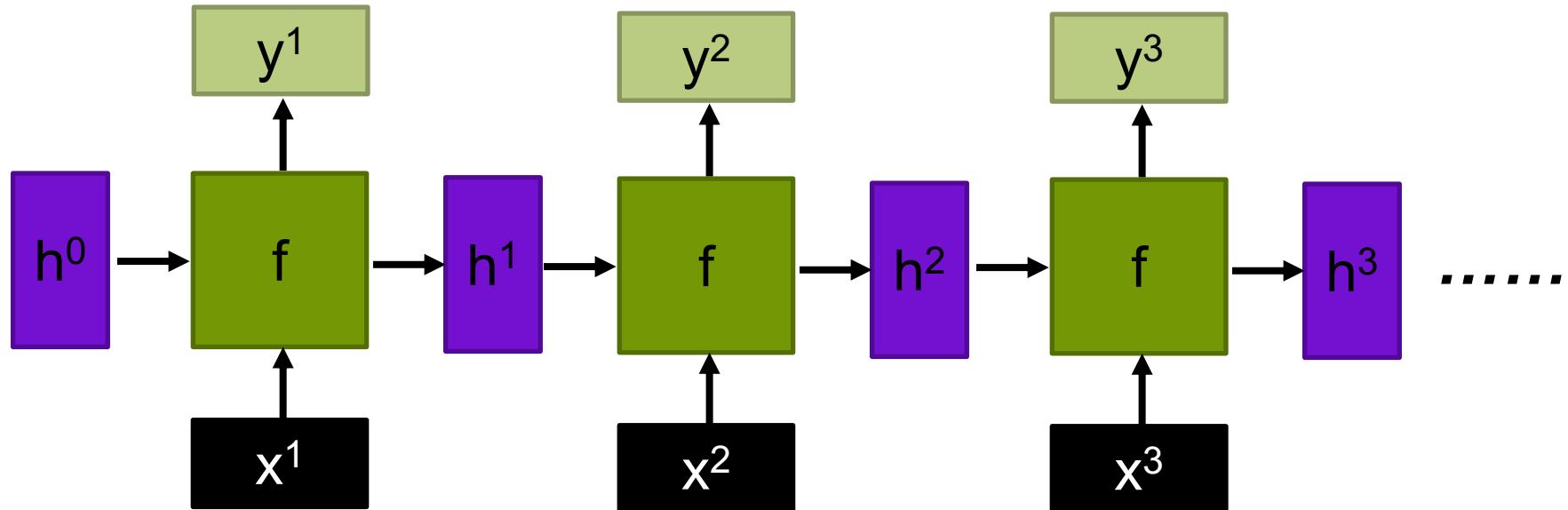
**Filter 2**





- Given function  $f: h^i, y = f(h, x)$

$h^i$  are vectors with the same dimension



No matter how long the input/output sequence is, we only need one function  $f$ .

## RNN Diagram

... with shared (tied) weights

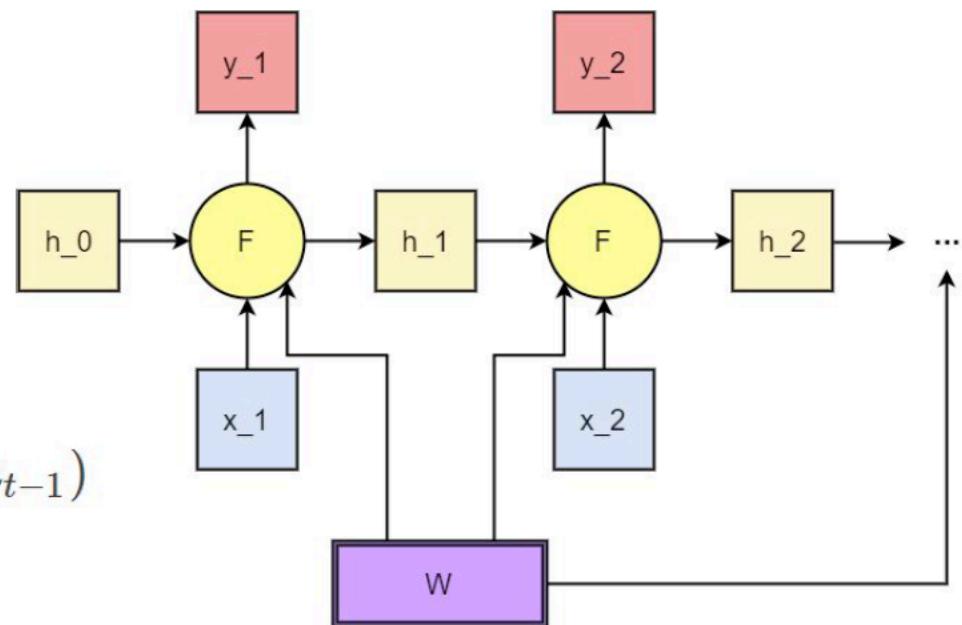
$$(h_1, y_1) = F(h_0, x_1, W)$$

$$(h_2, y_2) = F(h_1, x_2, W)$$

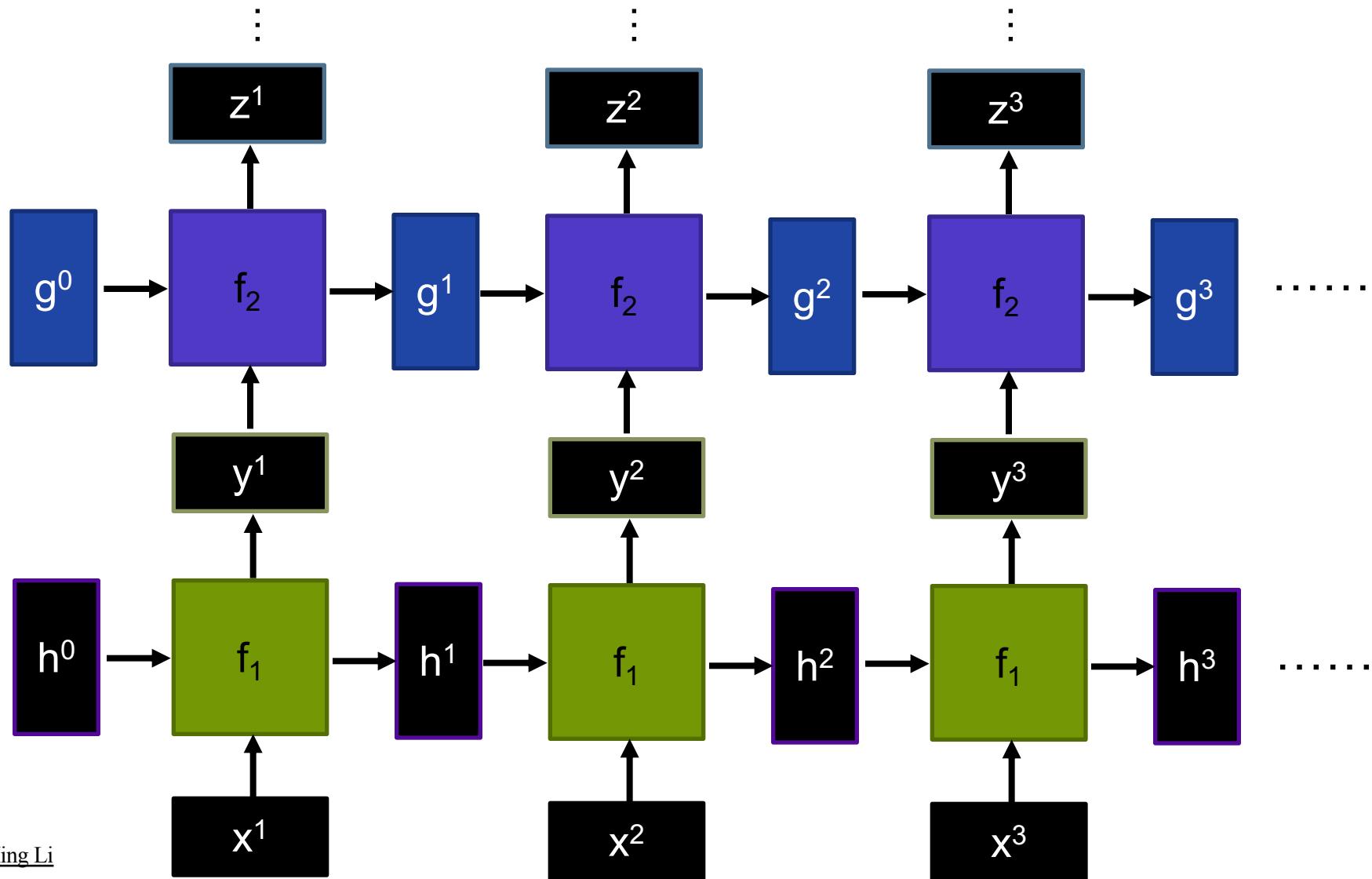
A simple implementation of F

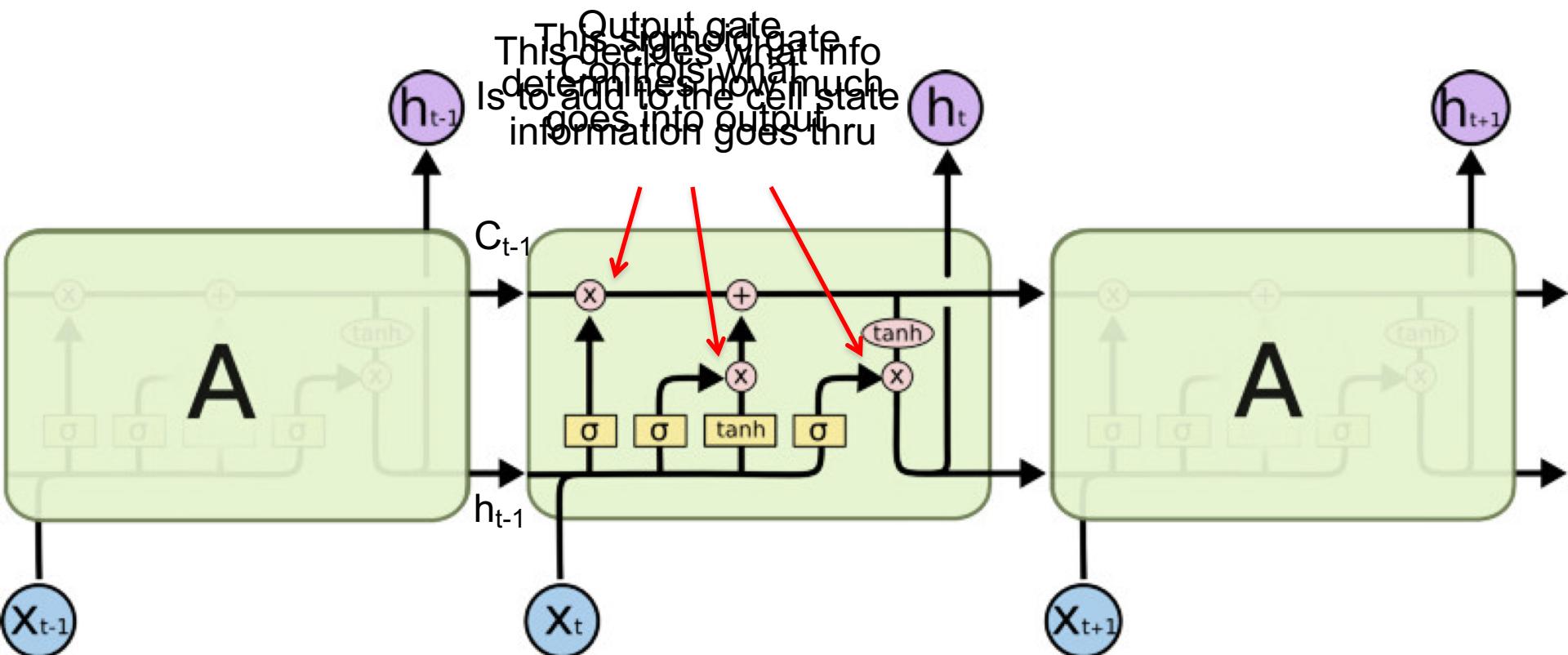
$$h_t = \tanh(W_{ih}x_t + W_{hh}h_{t-1})$$

$$y_t = W_{ho}h_t$$

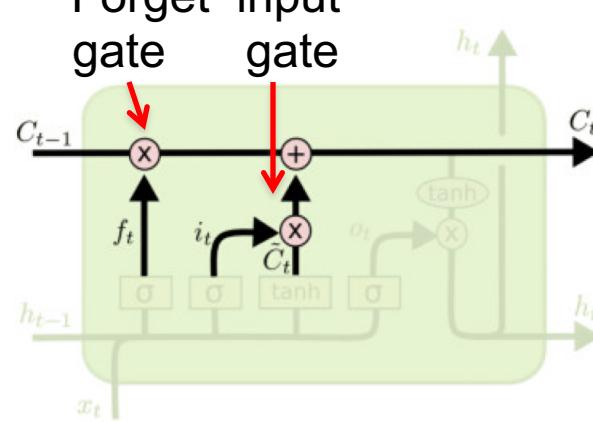


$$h^i, y = f_1(h, x), g^i, z = f_2(g, y)$$

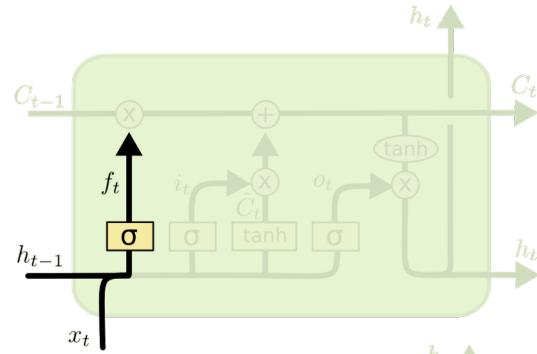




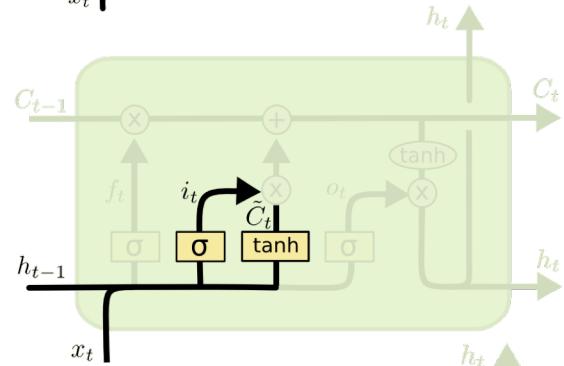
The core idea is this cell state  $C_t$ , it is changed slowly, with only minor linear interactions. It is very easy for information to flow along it unchanged.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

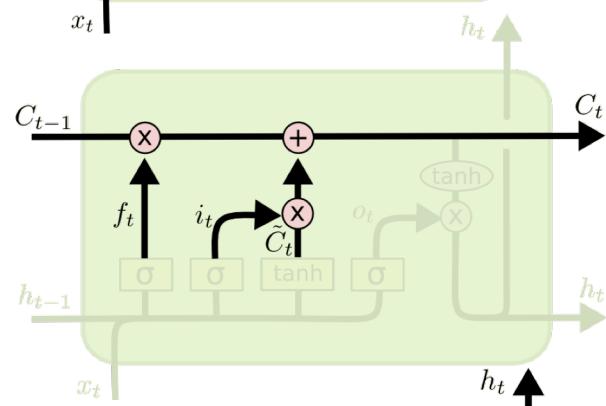


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

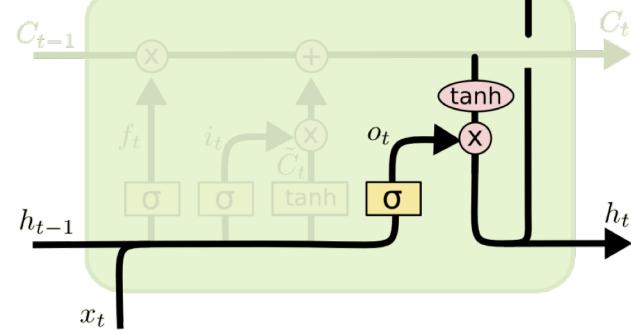


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

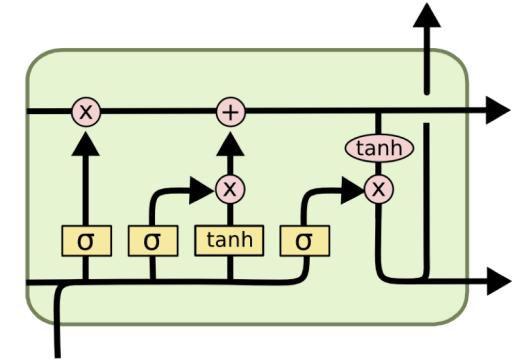


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

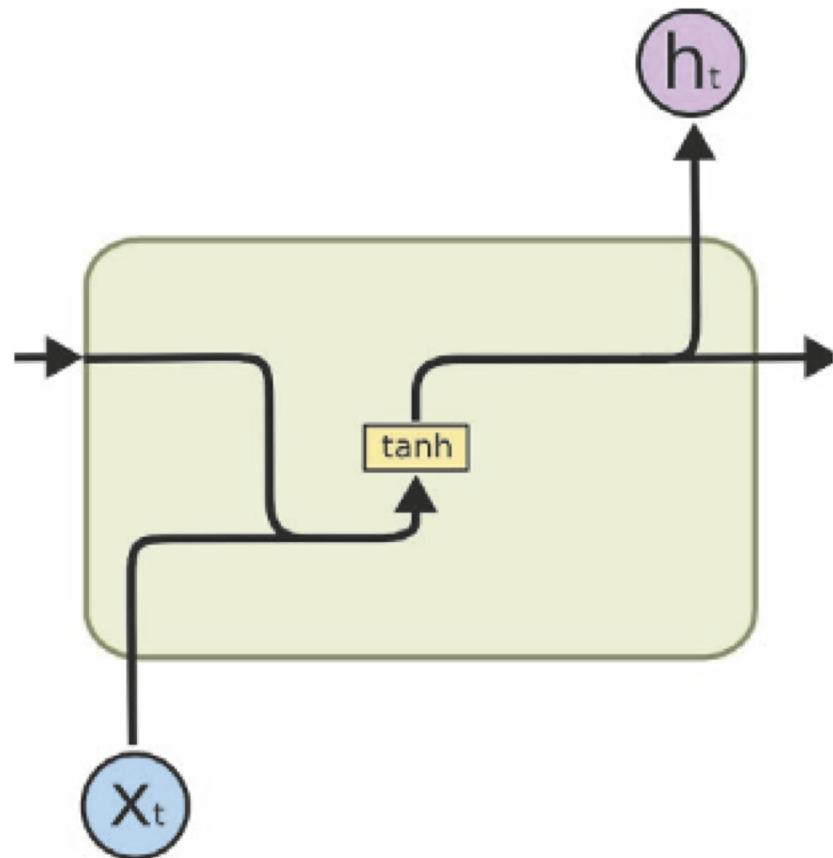
$$h_t = o_t * \tanh(C_t)$$



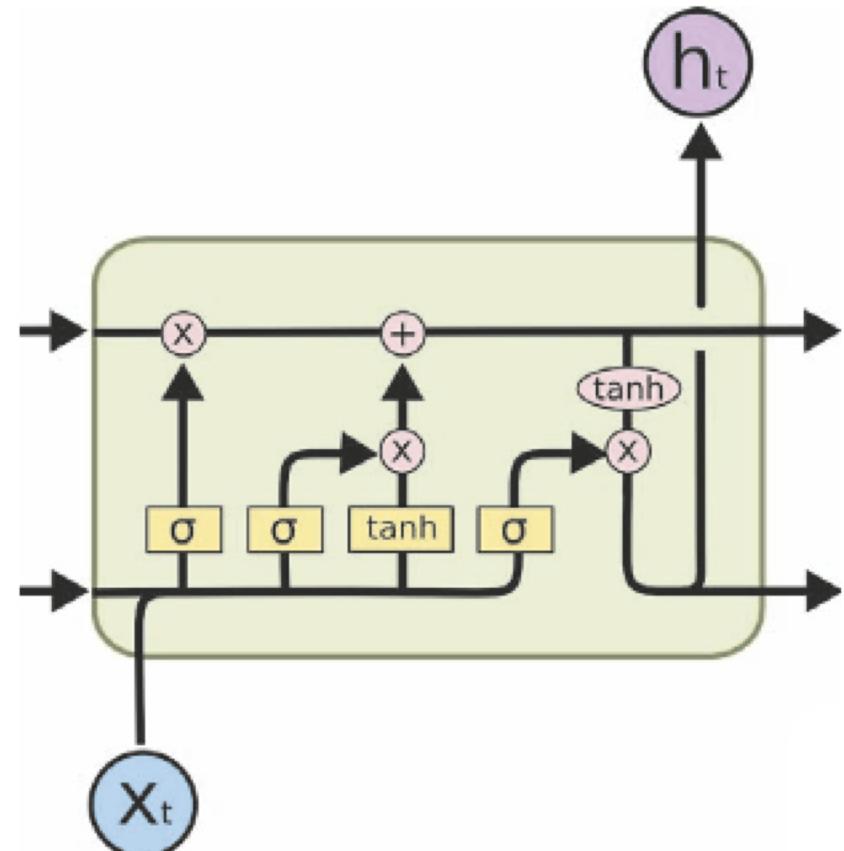
***i<sub>t</sub>* decides what component is to be updated.**  
***C'<sub>t</sub>* provides change contents**

***Updating the cell state***

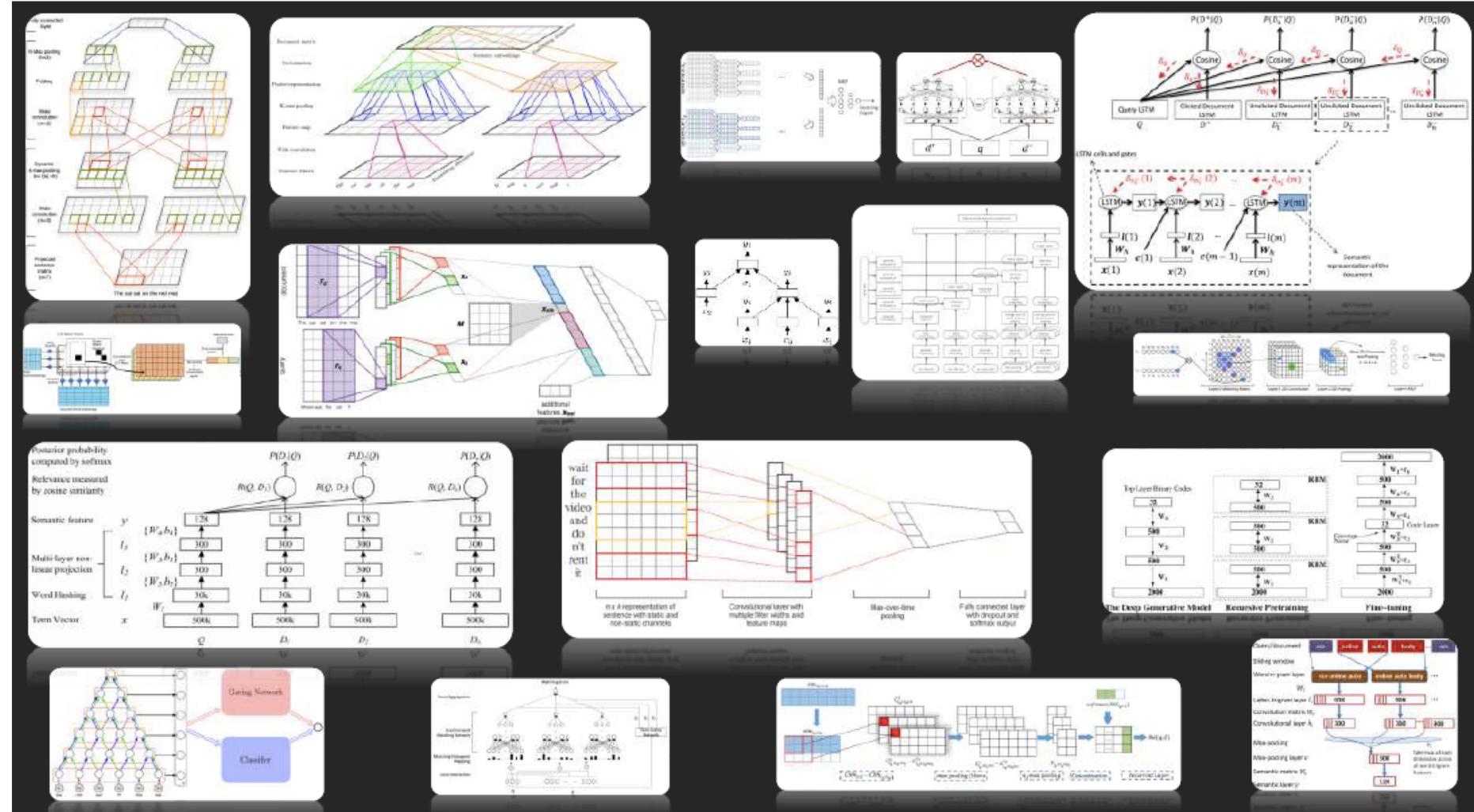
***Decide what part of the cell state to output***



(a) RNN



(b) LSTM

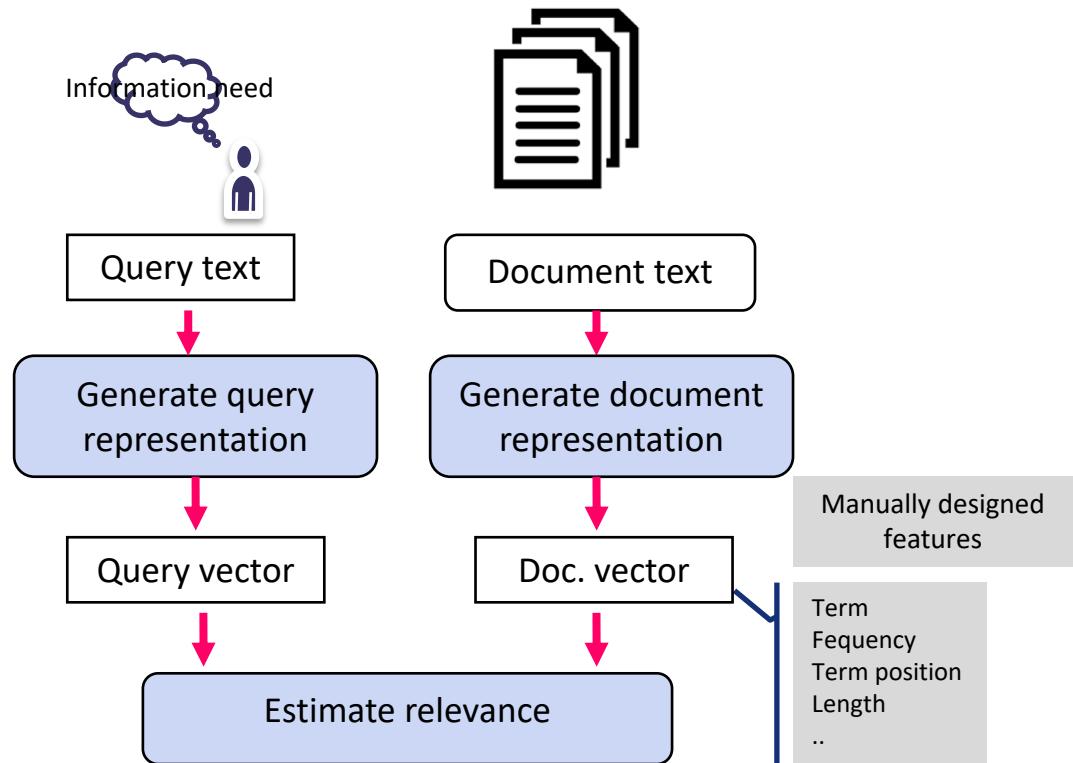


# Modèles neuronaux de RI

## Deep Neural IR

- Information retrieval (IR) in three simple steps:

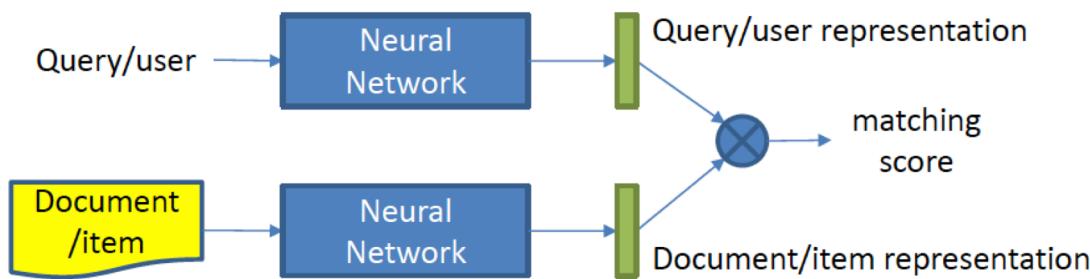
- Generate query representation
- Generate doc. representation
- Estimate relevance



$$\sum_{w \in q \cap d} \ln \frac{N - df(w) + 0.5}{df(w) + 0.5} \cdot \frac{(k_1 + 1) \times c(w, d)}{k_1((1 - b) + b \frac{|d|}{avdl}) + c(w, d)} \cdot \frac{(k_3 + 1) \times c(w, q)}{k_3 + c(w, q)}$$

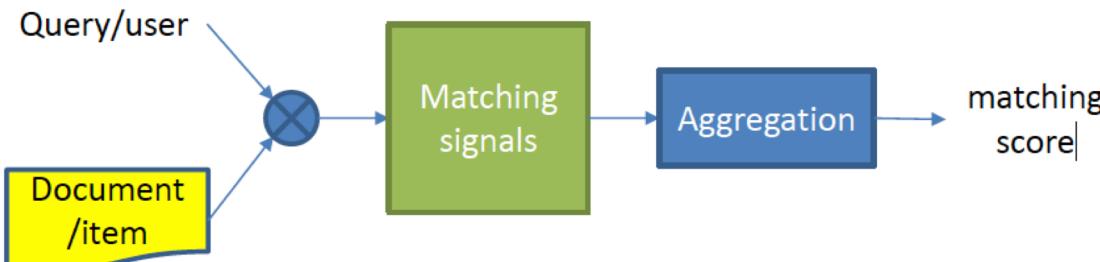
Learning 2 Rank (SVM, NN, ..)

- Apprentissage ds représentations Learning (Representation-based models)

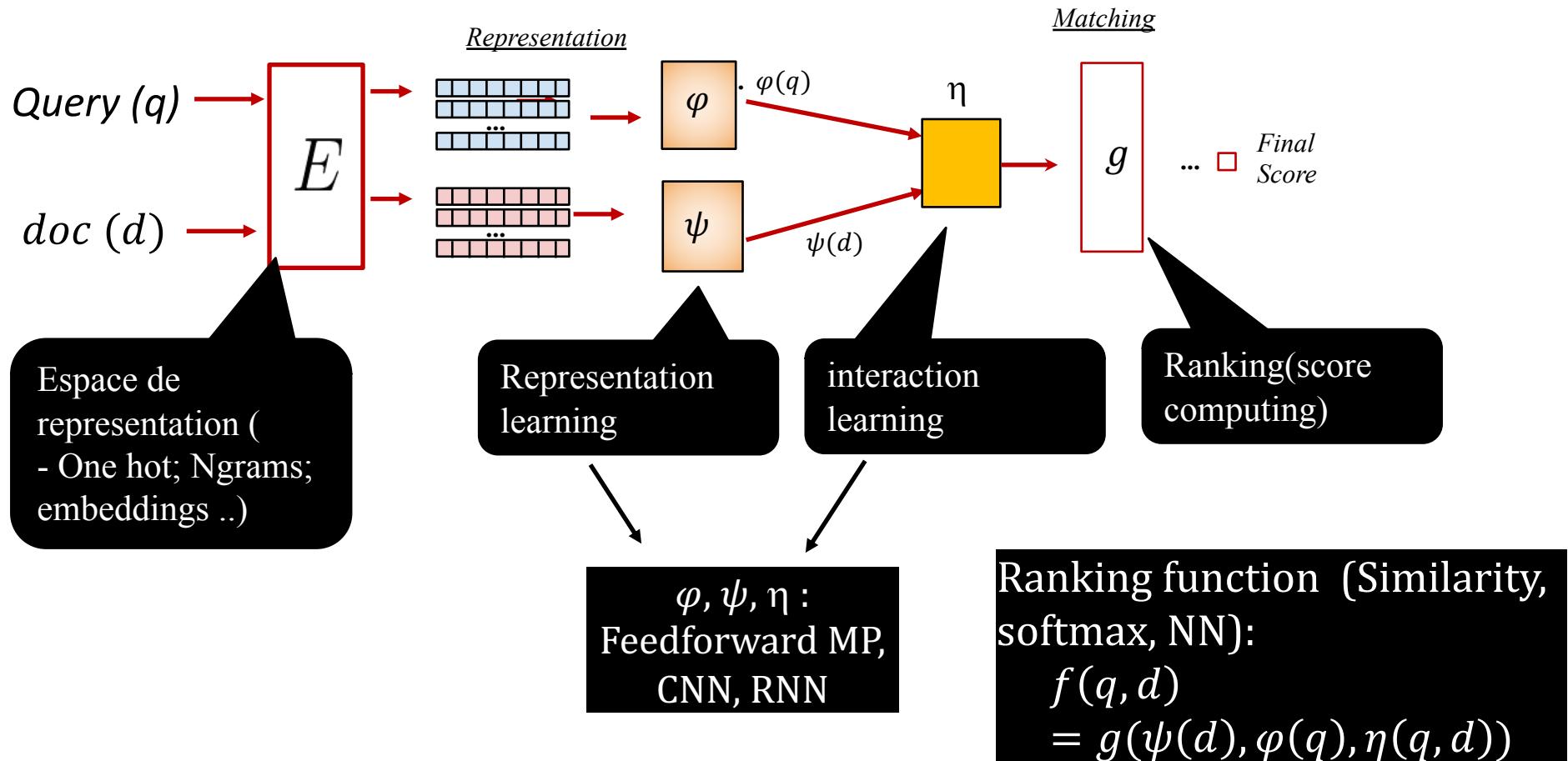


- Learn latent patterns that better represent a document and a query
- matching function is simple (cosine, MLP)

- Apprentissage de la fonction de «matching » (Interaction-based models)



Learn relevance from interactions between a document and a query rather than from individual representations



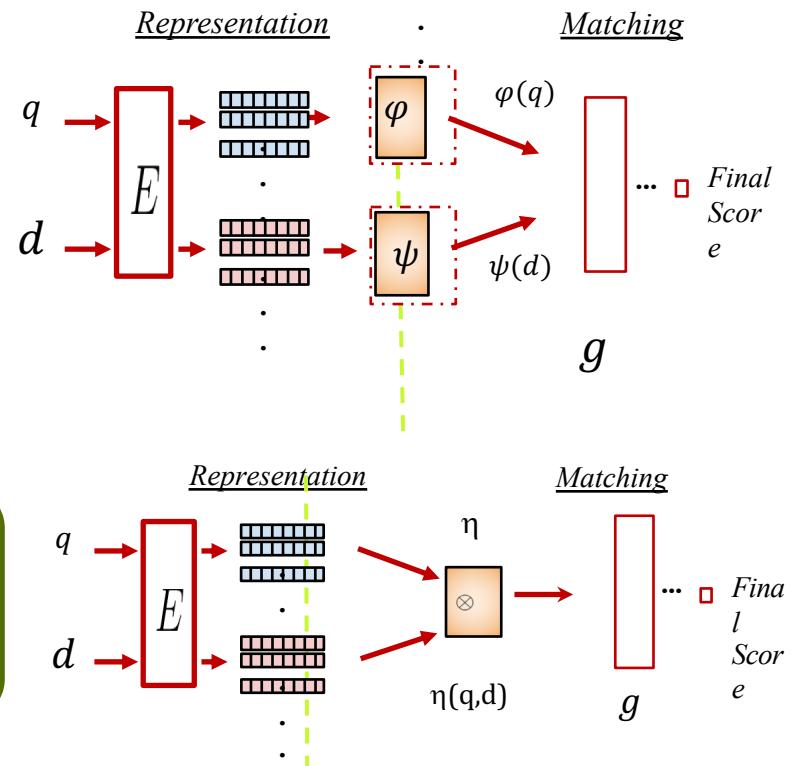
Neural text matching  
(Guo et al., 2016)

### Representation Focused

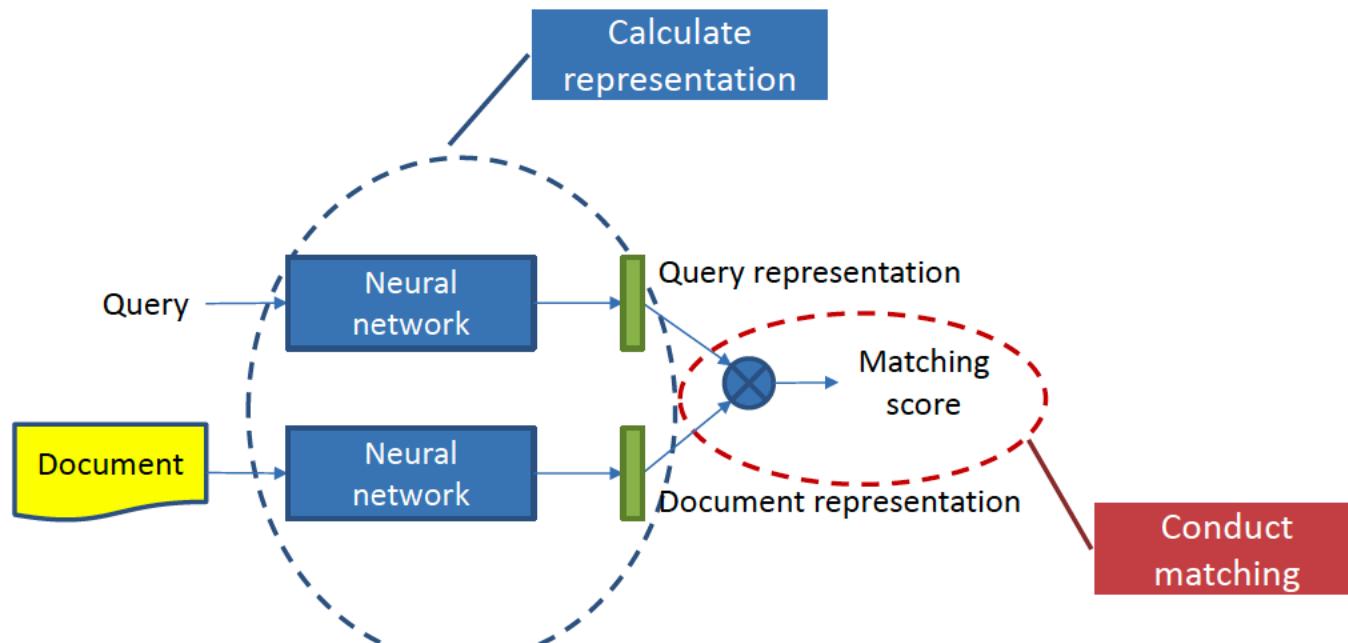
- DSSM (Huang et al., 2013)
- ARC-I (Hu et al. 2014)
- CLTR (Severyn et al., 2015)
- LSTM-RNN (Palangi rt al. 2016)
- DRCN (Kim et al., 2019)

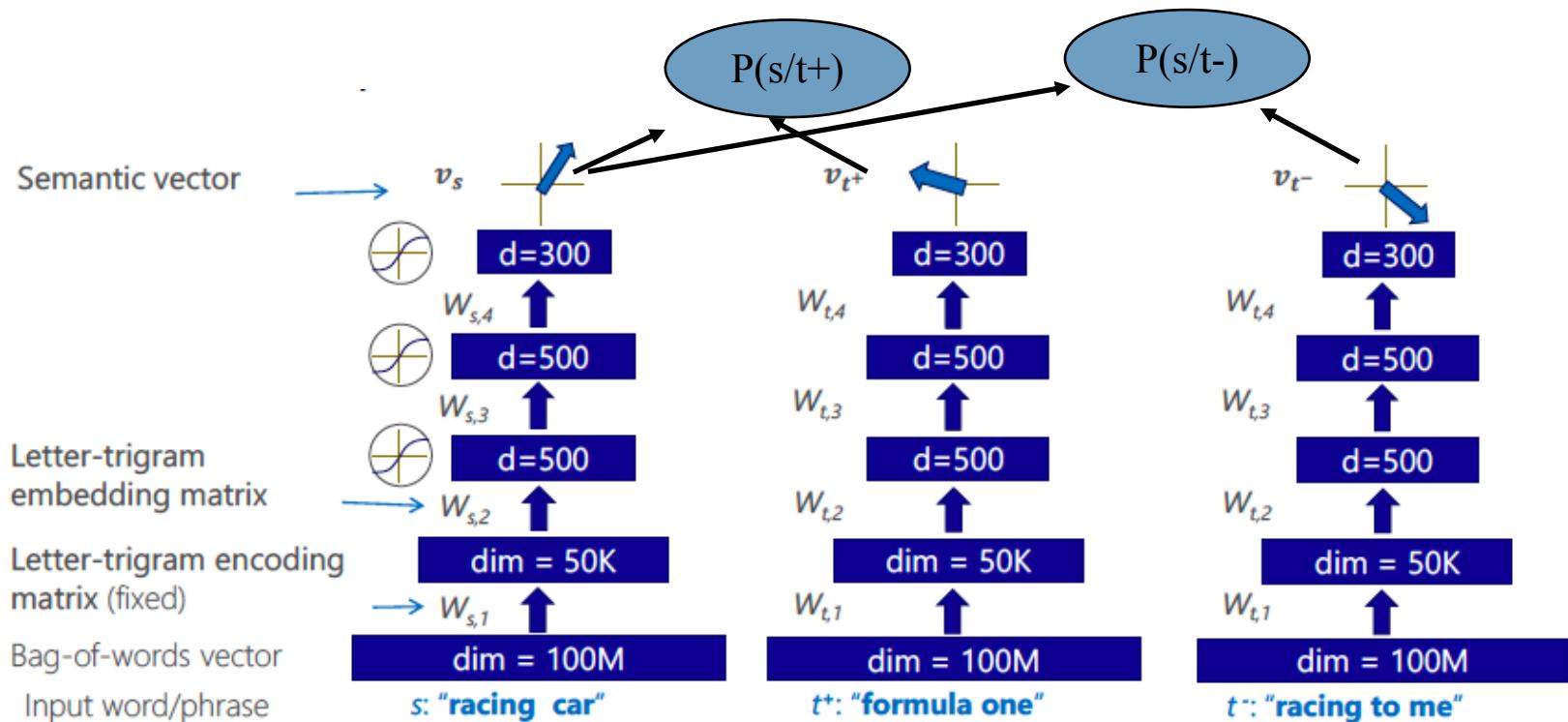
### Interaction Focused

- ARC-II (Hu et al., 2014)
- MatchPuramid (Pang et al. 2016)
- DRMM (Guo et al., 2016)
- DUET (Mitra et al., 2017)



- Étape 1 : Calculer (apprendre) les représentations des textes (requête et document)  $\phi(x)$
- Étape 2 : Calculer leur similarité (cosine, ...)  $F(\phi(q), \phi(d))$





DSSM [Huang et al., 2013, Shen et al., 2014)

© He et al., CIKM '14 tutorial

Learn a fixed-dimensional vector representation for each text separately and then perform matching within the latent space

- DSSM Model : (Autoencoder)

## DSSM Matching Function

- Cosine similarity between semantic vectors

$$S = \frac{x^T \cdot y}{|x| \cdot |y|}$$

- Training

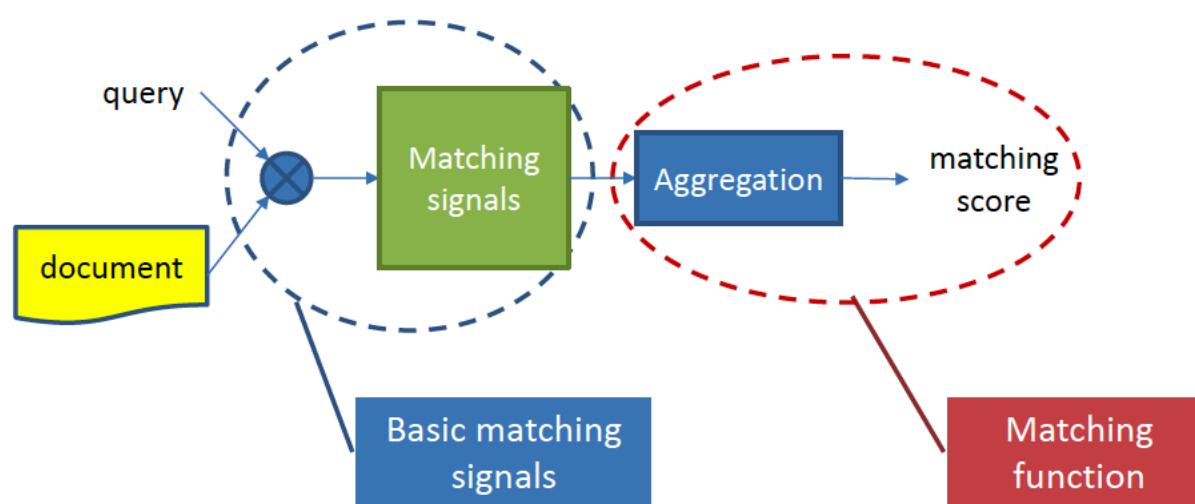
- A query  $q$  and a list of docs  $D = \{d^+, d_1^-, \dots, d_k^-\}$

- $d^+$  positive doc,  $d_1^-, \dots, d_k^-$  negative docs to query

- Objective:

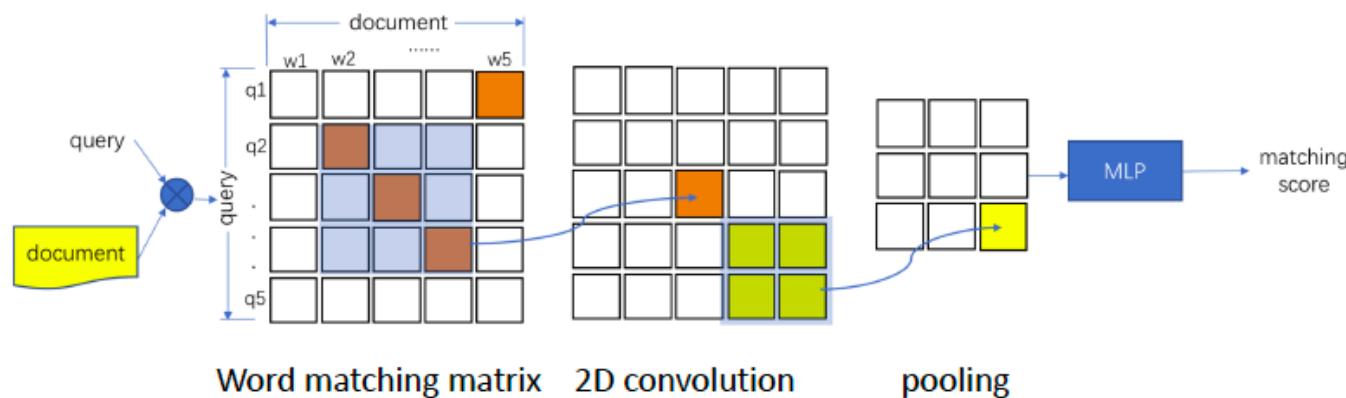
$$P(d^+|q) = \frac{\exp(\gamma \cos(q, d^+))}{\sum_{d \in D} \exp(\gamma \cos(q, d))}$$

- Étape 1 : construire des signaux de base de bas niveau
- Étape 2 : agréger les modèles de correspondance (matching)



Interaction-based models compute the interaction between each individual term of both query and documents.

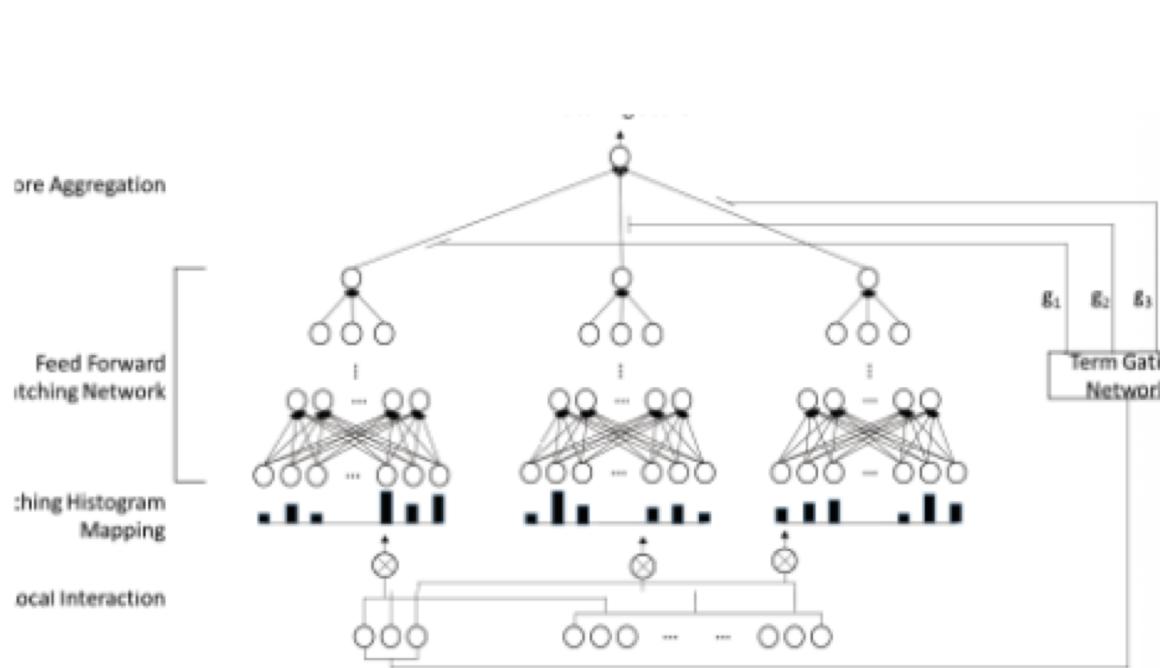
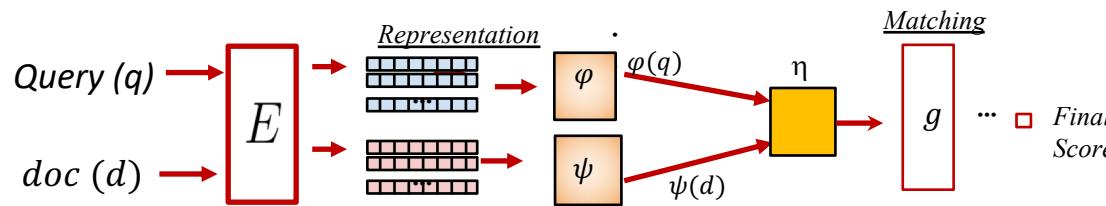
- Inspired by Image recognition
- Basic matching signals: word level matching matrix
- Matching function: 2D convolution+ MLP



84

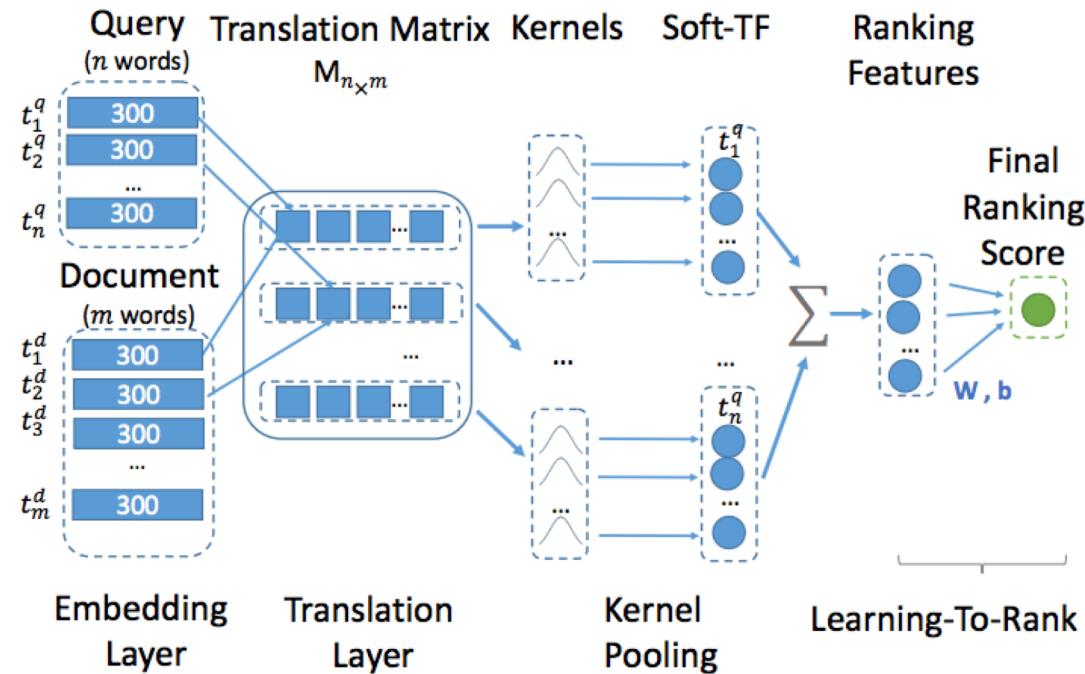
(Peng et al. AAA'16)

- Representation learning + matching function learning



(Guo et al., 2016)

- Compute word-interaction matrix, (s-cosine similarity of word embeddings)
- This results in k-dimensional vector.
- Aggregate the query term vectors into a fixed-dimensional query representation.

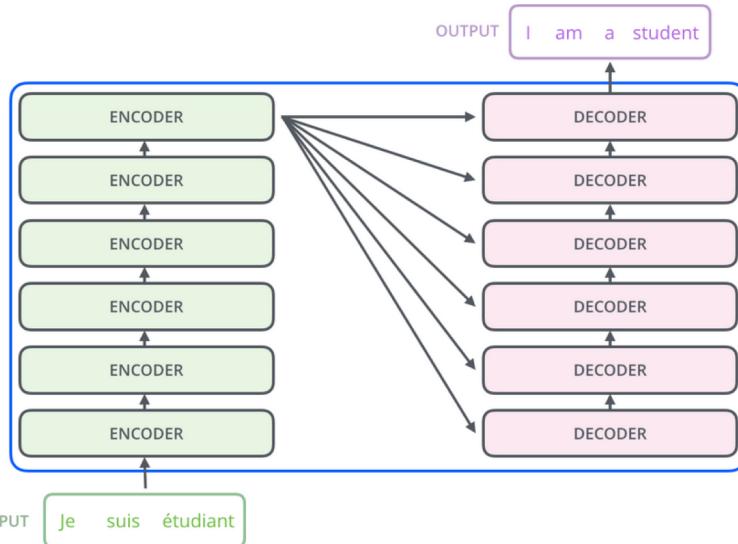


K-NRM [Xiong et al., 2017b]

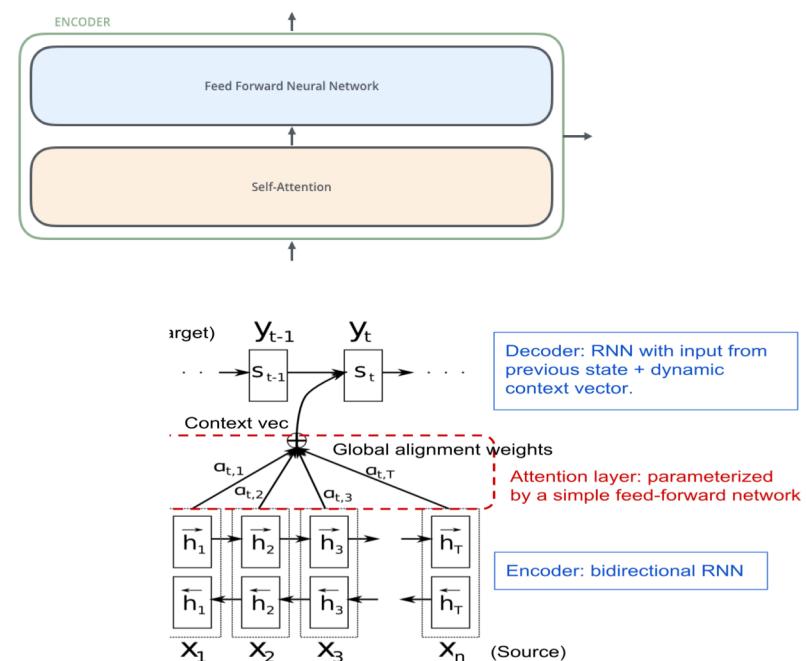
# Deep learning and IR

## La révolution des transformateurs à la BERT

- Transformateurs (*transformers*) ou (modèle auto attentif)(Self attention)
  - Catégorie de Deep ML utile pour le traitement de la langue (texte, séquences de mots) qui ne nécessite pas le traitement de la sequence dans l'ordre des mots (contrairement au RNN).
  - BERT, ELMO, ULMfit, OpenAI GPT
- Sont construits autour de trois notions:
  - Encodeur, Décodeur, Attention(self)



© J.Alammar



- Un panoplie de *transformateurs* avec des architectures différentes : (Codeur-self attention avec ou sans(décodeur))

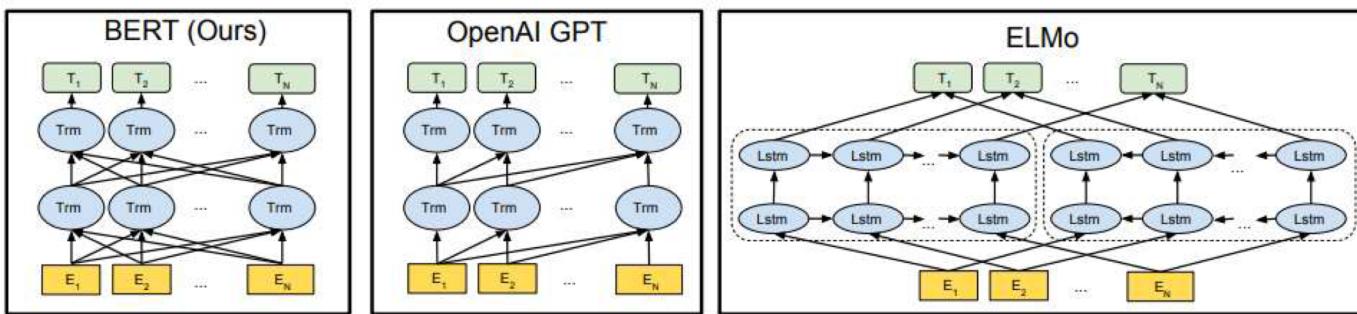
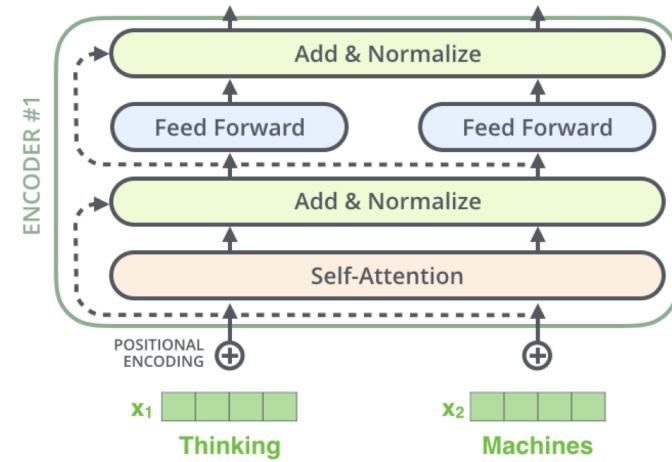
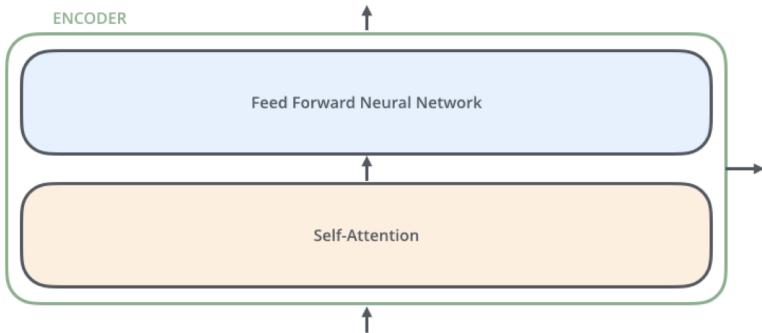
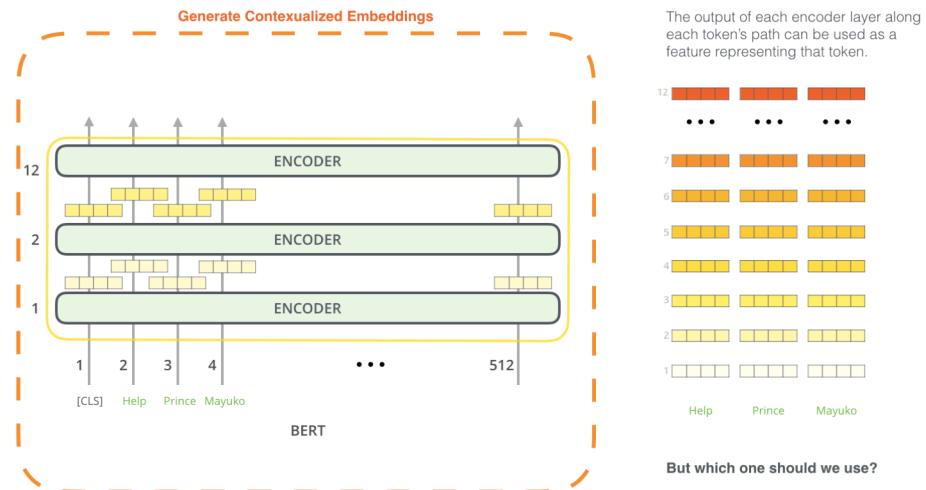
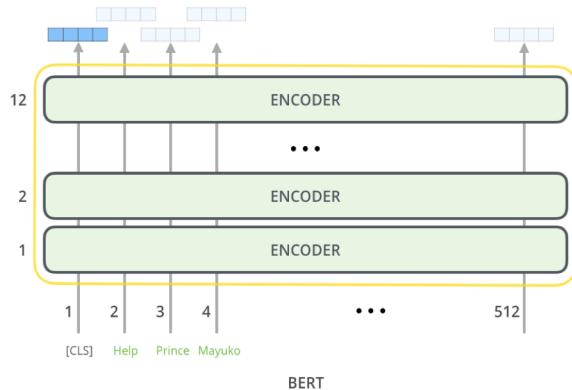


Figure 1: Differences in pre-training model architectures. BERT uses a bidirectional Transformer. OpenAI GPT uses a left-to-right Transformer. ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTM to generate features for downstream tasks. Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

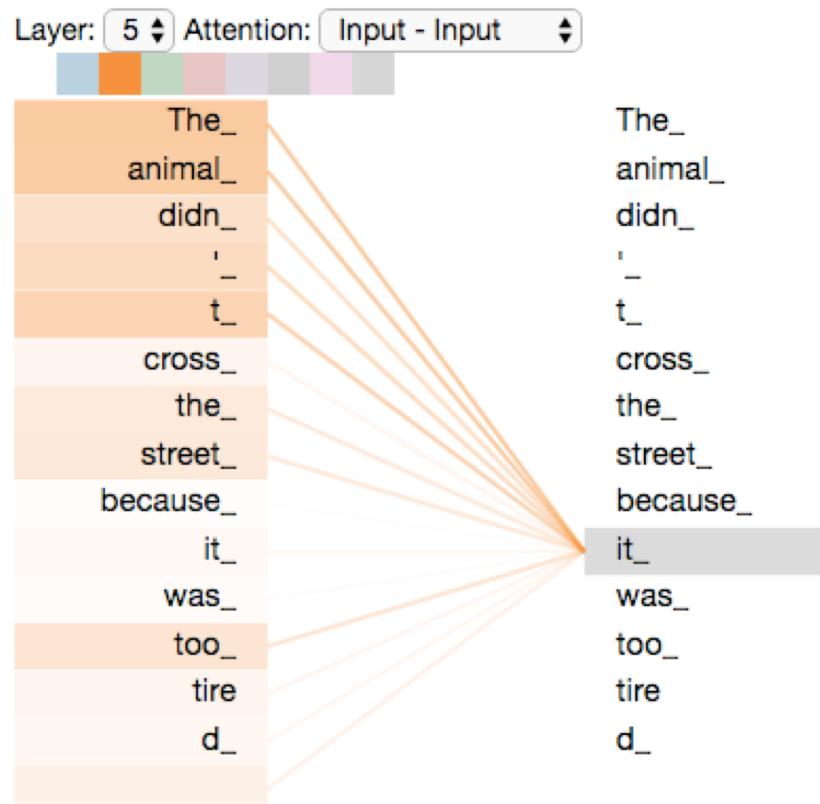
# Transformateurs (*Transformers*)

## Encodeur (*Encoder*)

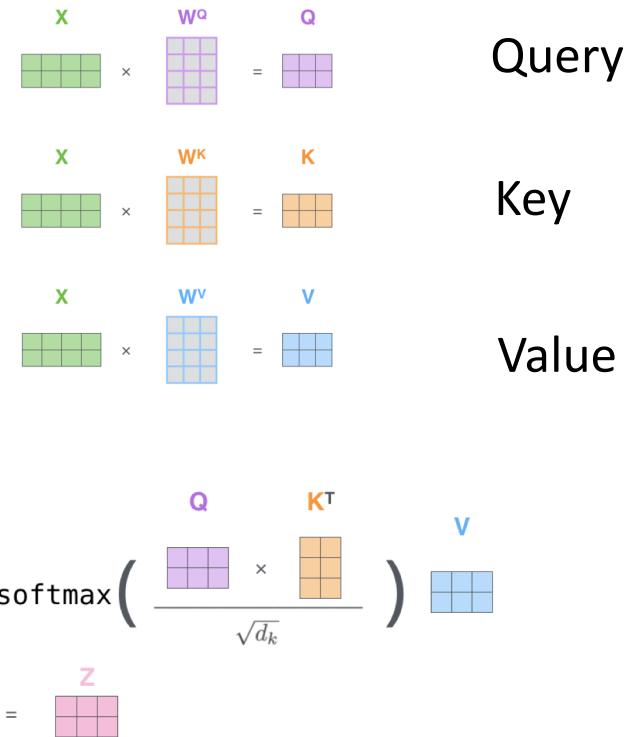
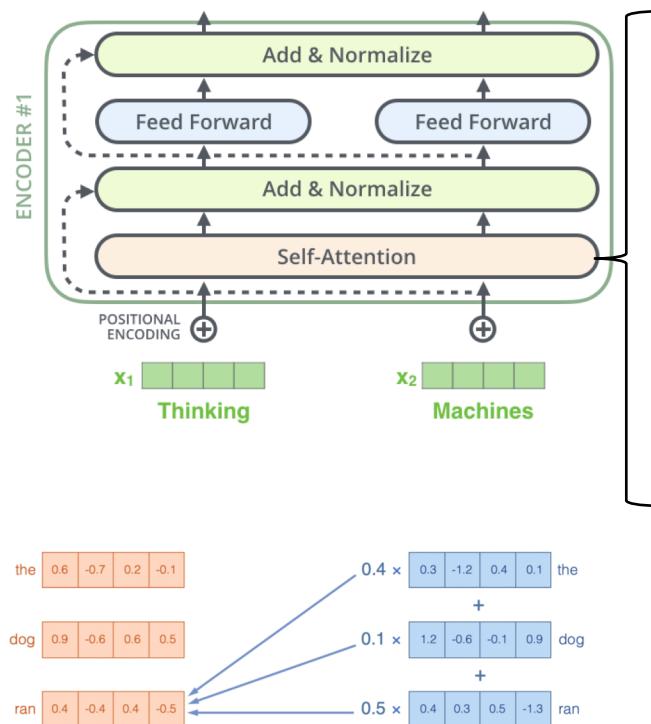


*"The animal didn't cross the **street** because it was too **wide**"*

*" The **animal** didn't cross the street because it was too **tired**"*



As the model processes each word (each position in the input sequence), self-attention allows it to look at other positions in the input sequence for clues that can help lead to a better encoding for this word



# Self-attention

<https://arxiv.org/abs/1706.03762>

*q: query (to match others)*

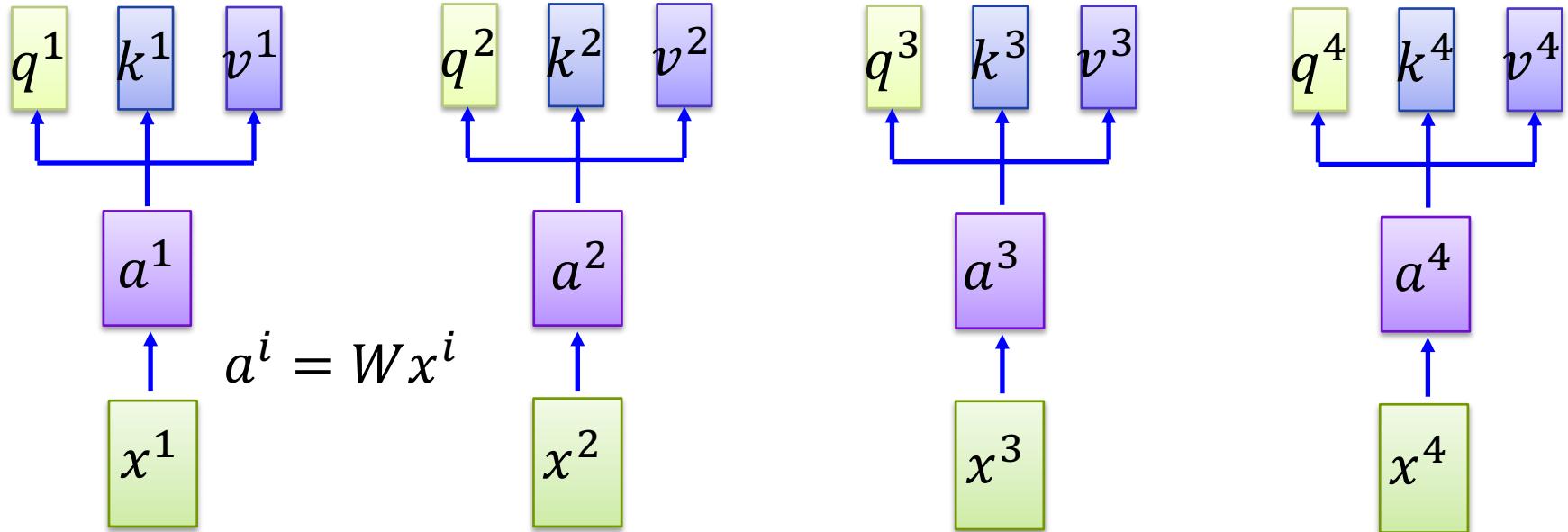
$$q^i = W^q a^i$$

*k: key (to be matched)*

$$k^i = W^k a^i$$

*v: information to be extracted*

$$v^i = W^v a^i$$



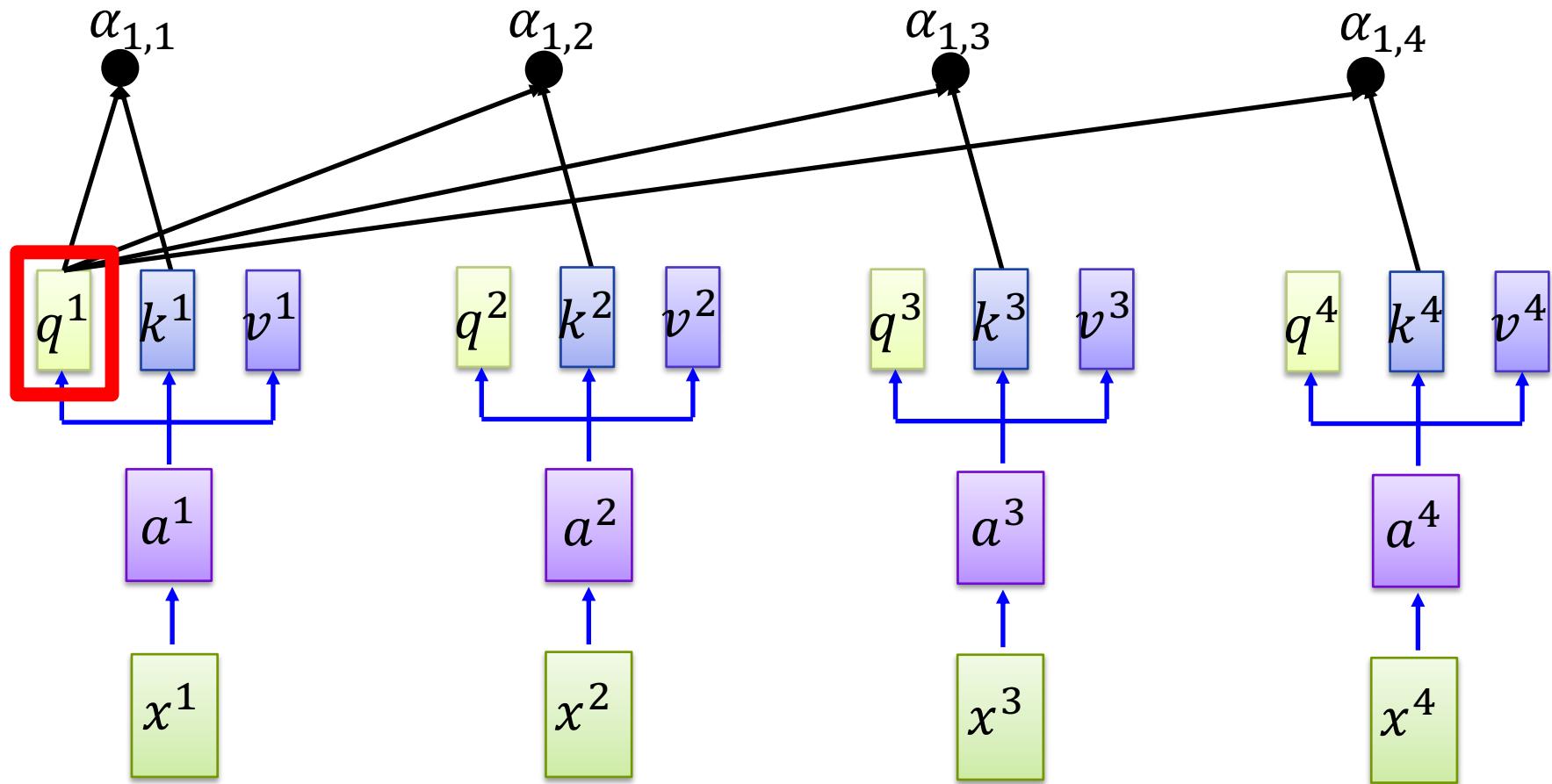
# Self-attention

*d is the dim of q and k*

*query q key k attention*

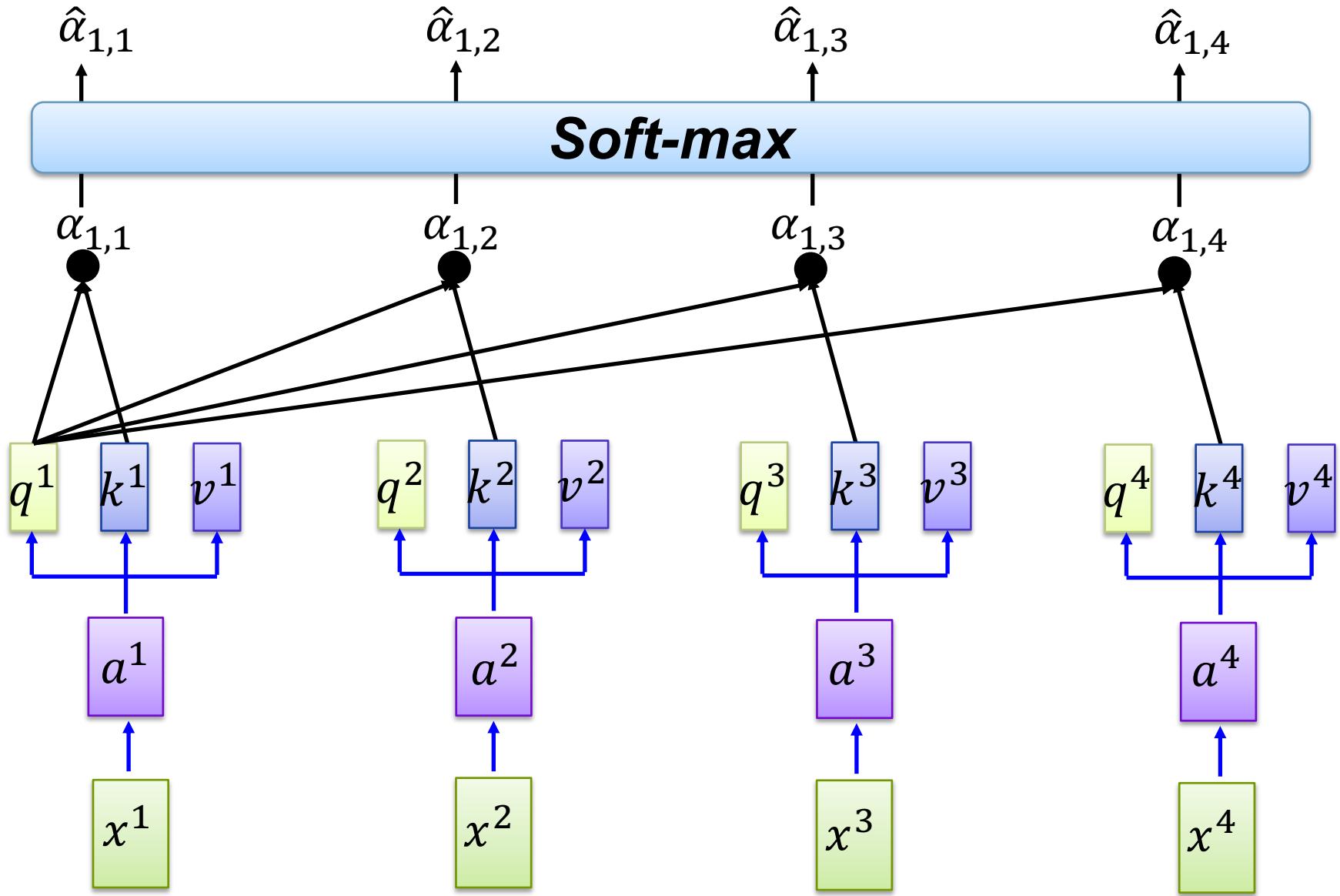
Scaled Dot-Product Attention:

$$\alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$$



# Self-attention

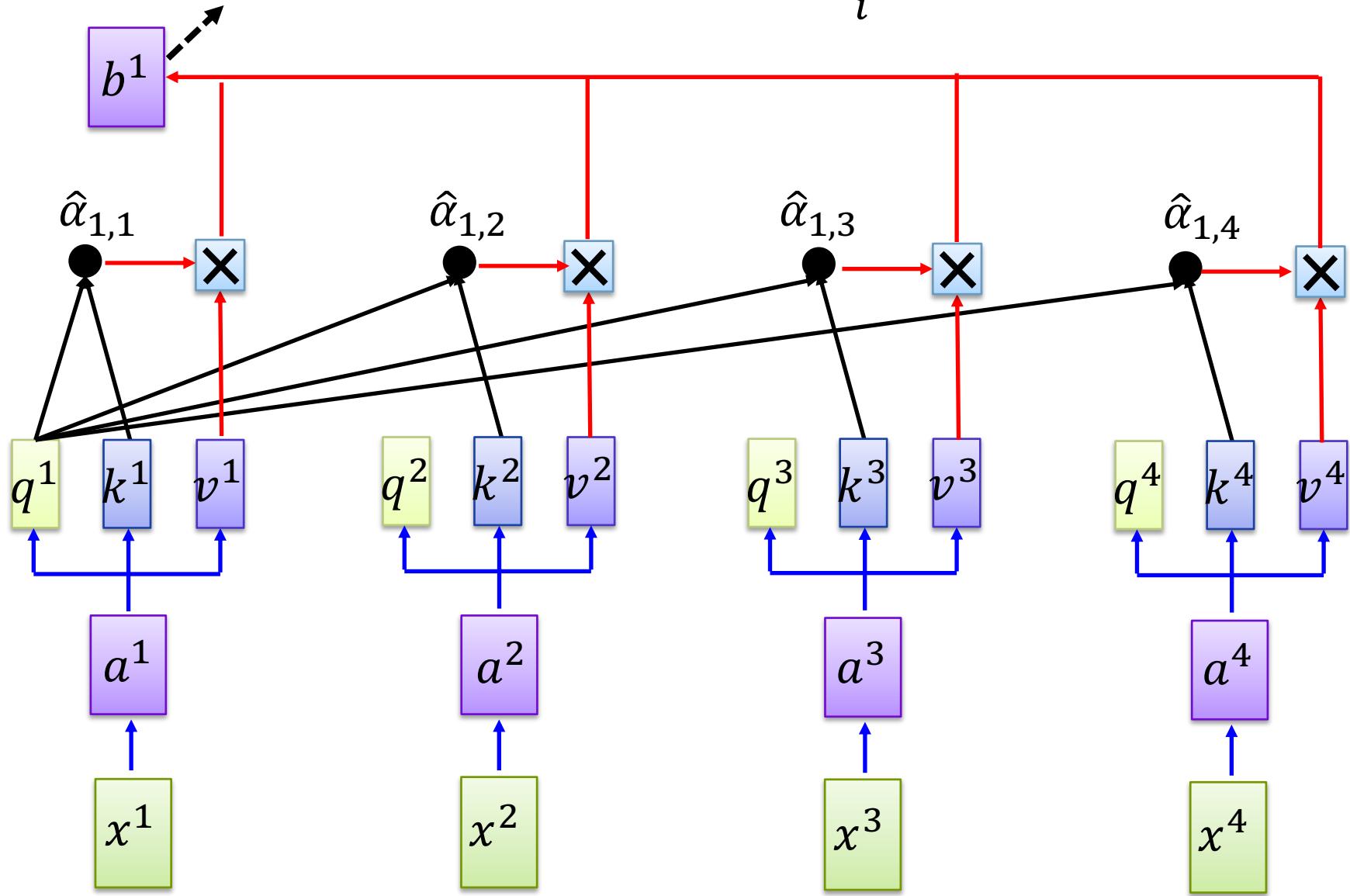
$$\hat{\alpha}_{1,i} = \exp(\alpha_{1,i}) / \sum_j \exp(\alpha_{1,j})$$



# Self-attention

$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$

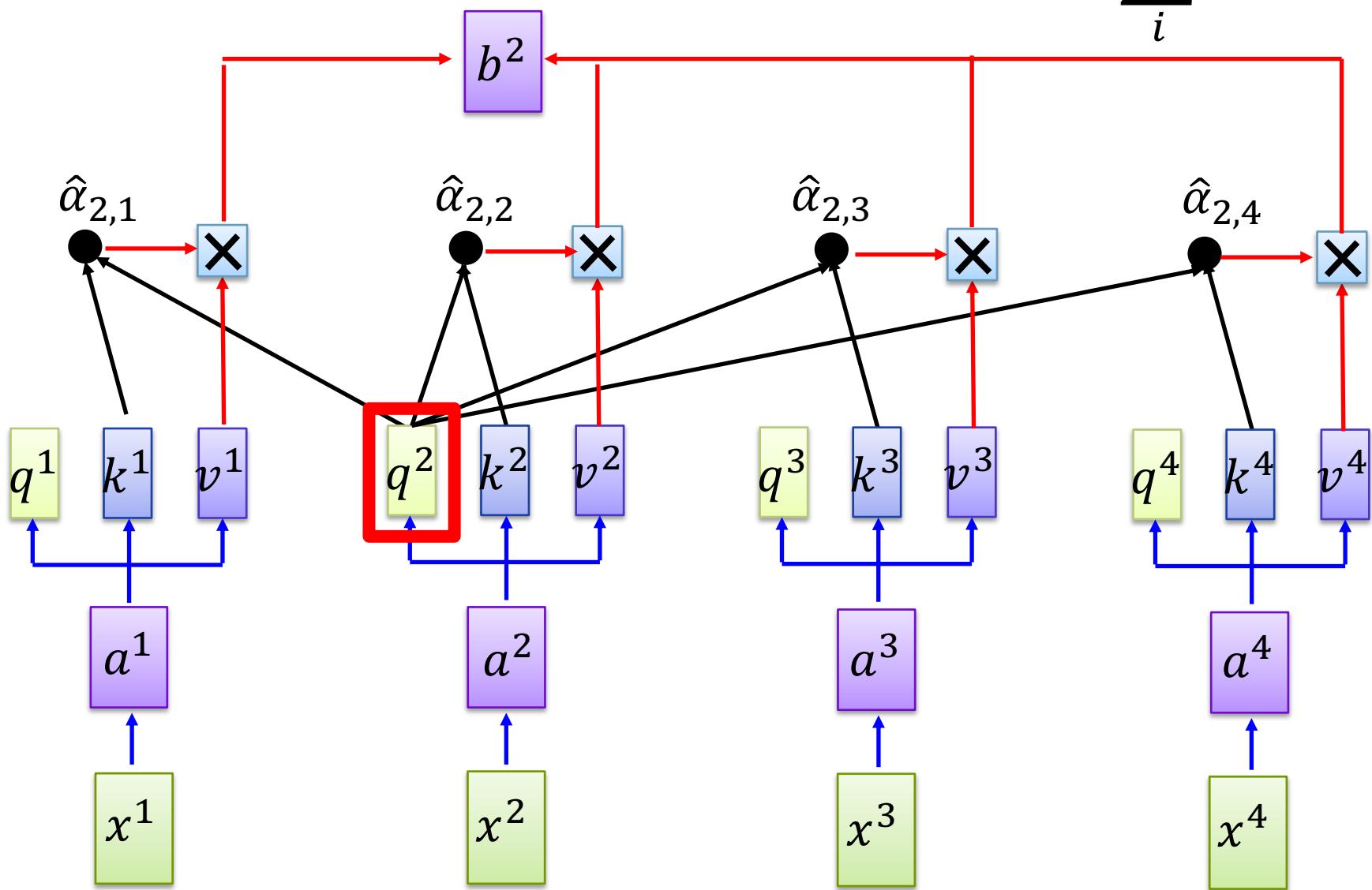
Considering the whole sequence



# Self-attention

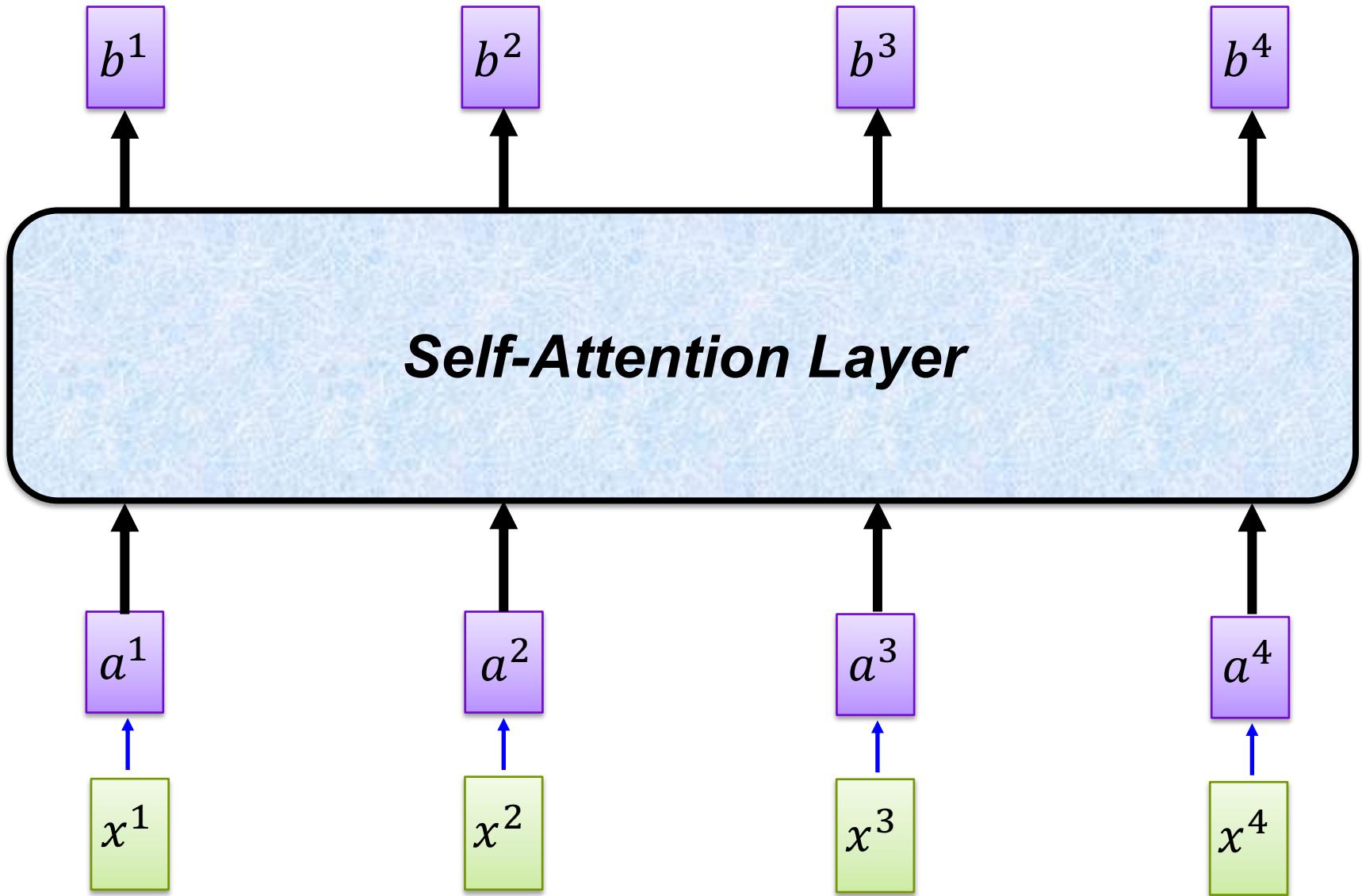
*query q key k attention*

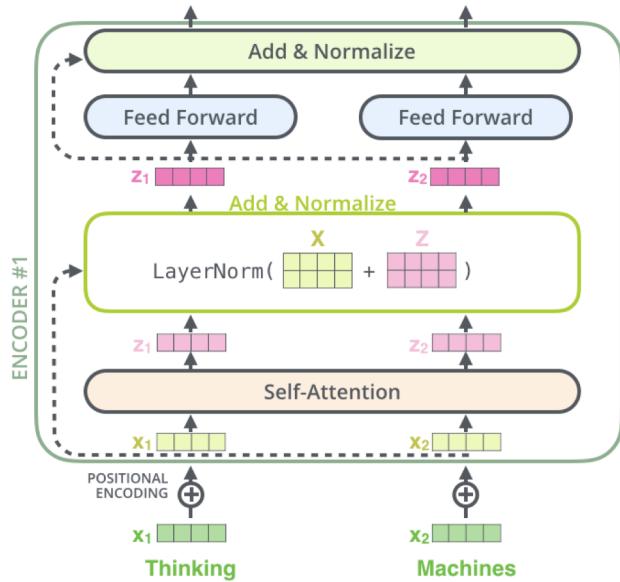
$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$



# Self-attention

$b^1, b^2, b^3, b^4$  can be parallelly computed.





1) This is our input sentence\*

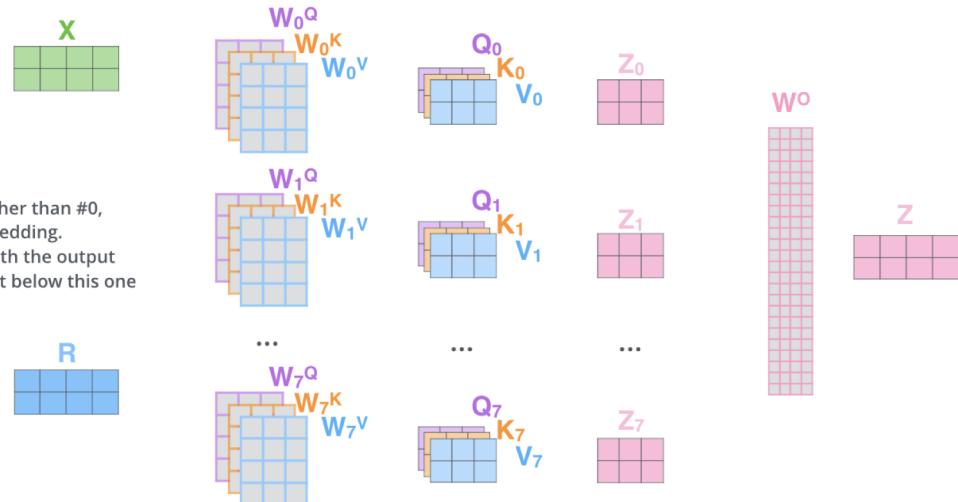
2) We embed each word\*

3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices

4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer

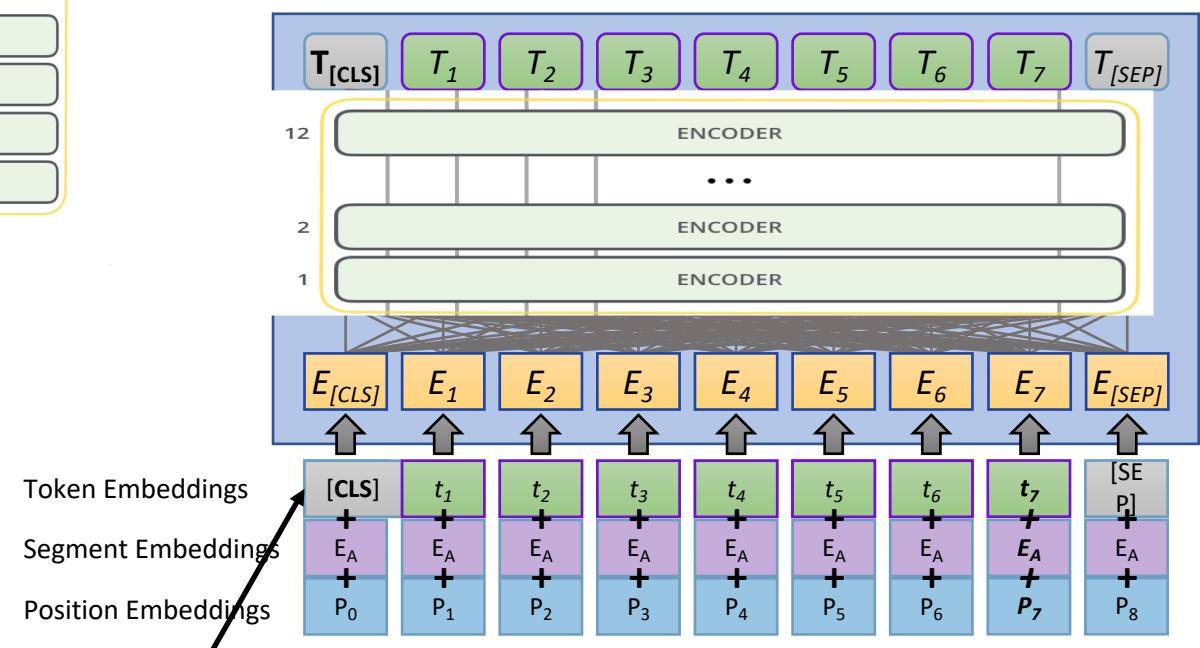
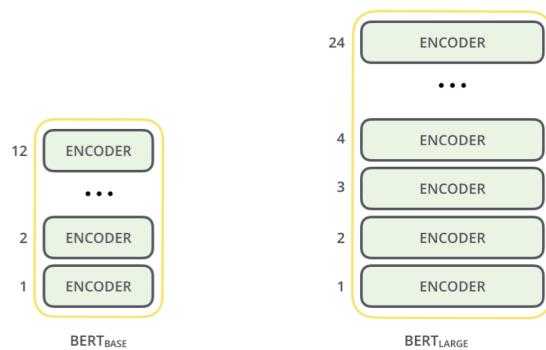
\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



- BERT is an attention mechanism that learns contextual relations between words (or sub-words) in a text (apprend **des modèles de langues**)
- BERT = un Encodeur (pas de décodeur), développé par Google
- BERT (and BERT like models RoBERT, ALBERT, ...) est la nouvelle “star” du TAL (traitement automatique des langues) et plus généralement des tâches exploitant le texte (recherche d’information, traduction automatique, système questions réponses, génération automatique de textes, ...)

En pratique, Deux BERT ont été proposés

- BASE : 12 encoders, , 12 Attention, Dim E/S (768)
- LARGE : 24 encoders, 16 attention heads, dimension vecteur d'E/S 1024



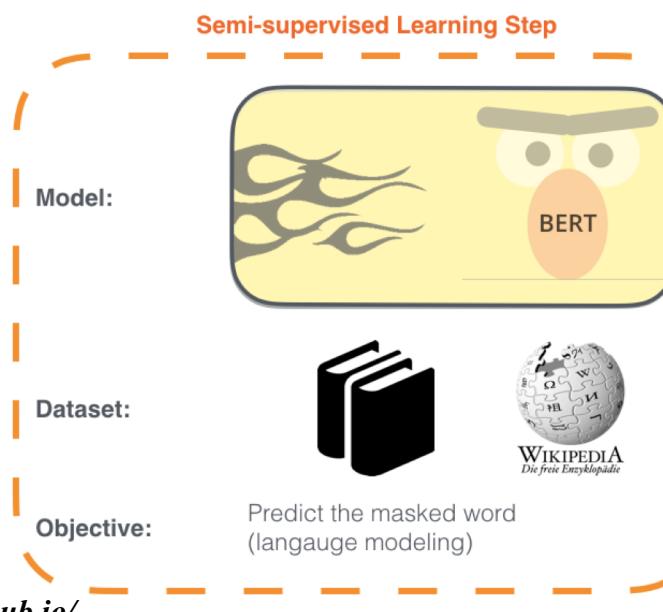
© alammar.github.io/ Le premier **token** d'entrée est un jeton spécial [CLS].

## En pratique

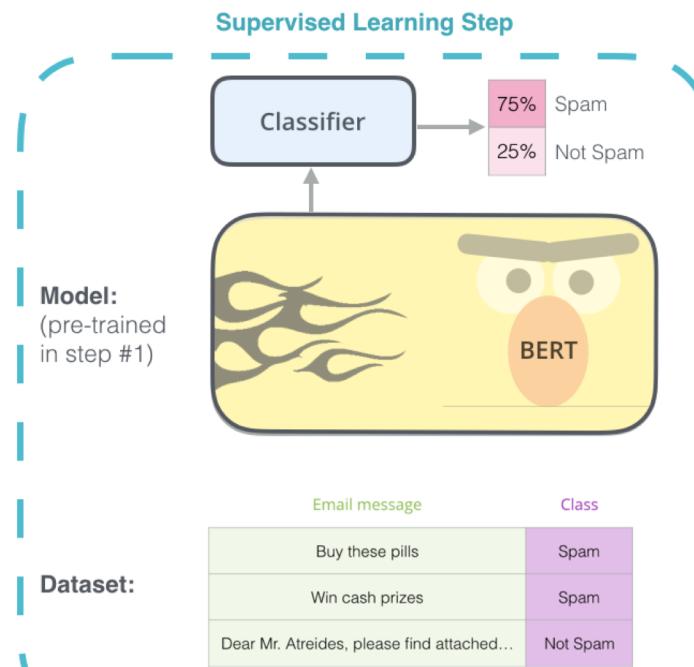
- Etape 1: on apprend un modèle de langue
- Etape 2 : on exploite le modèle sur une tâche supervisée (finetuning)

### 1 - Semi-supervised training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



### 2 - Supervised training on a specific task with a labeled dataset.



Apprentissage d'un modèle de langue (dans le cadre de BERT),  
2 tâches : masked word prediction (Masked Language Model) et Next Sentence Prediction

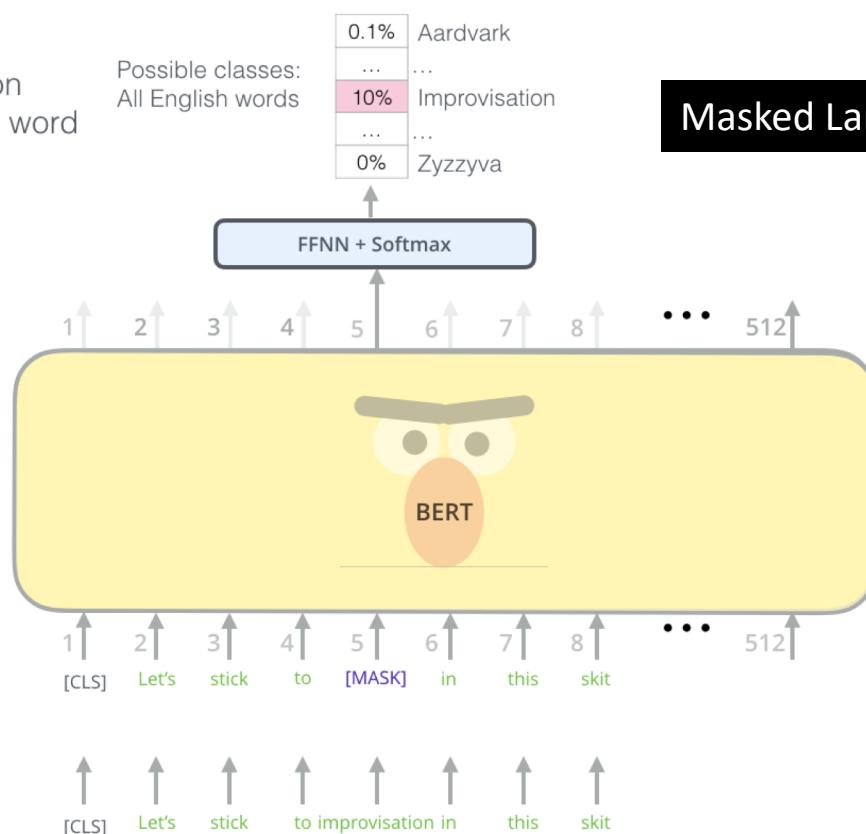
Use the output of the  
masked word's position  
to predict the masked word

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zzyzyva

Masked Language Model

Randomly mask  
15% of tokens

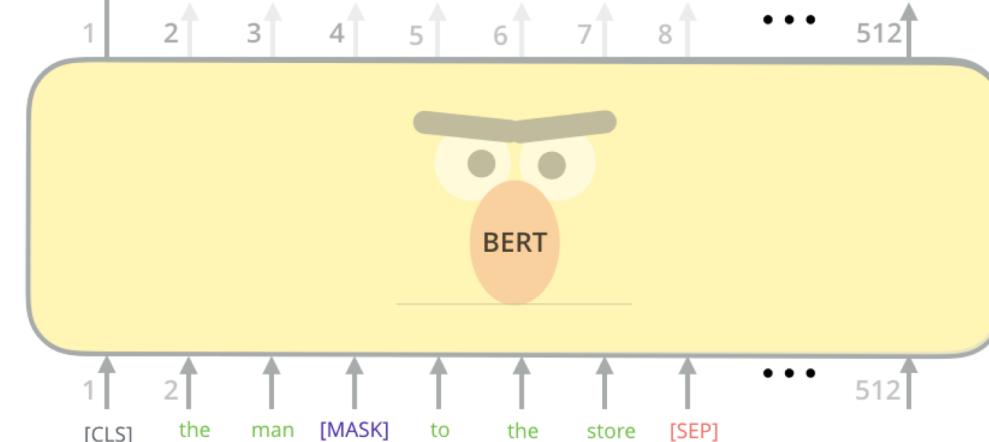


Input

Predict likelihood  
that sentence B  
belongs after  
sentence A



Next Sentence prediction  
(NSP)



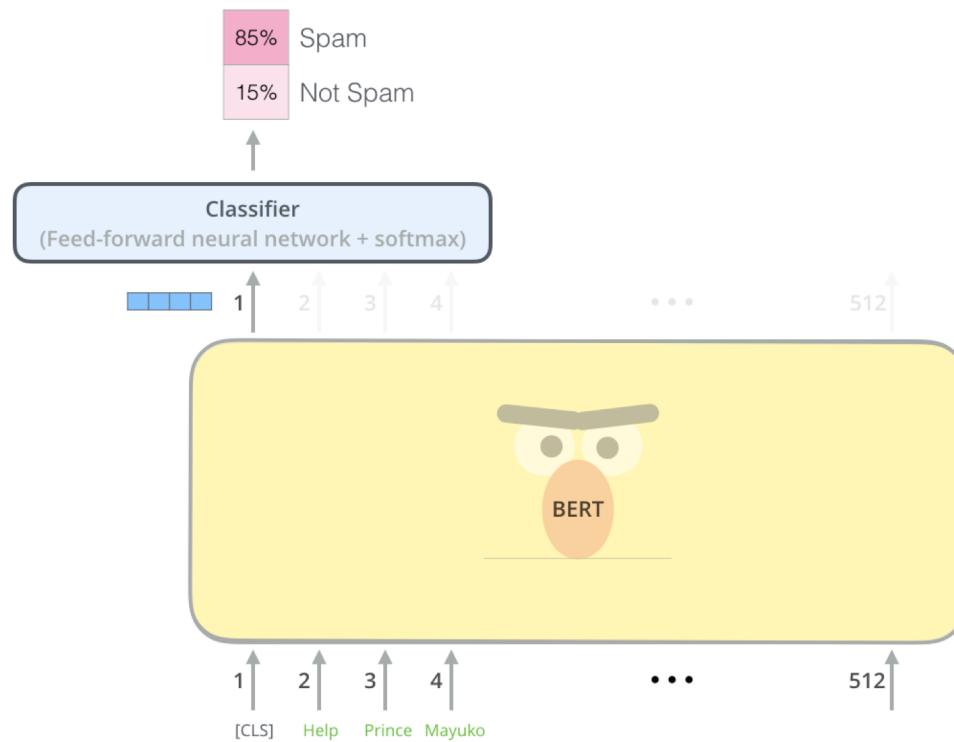
Tokenized  
Input

1 [CLS] 2 the 3 man 4 [MASK] 5 to 6 the 7 store 8 [SEP] ... 512

Input

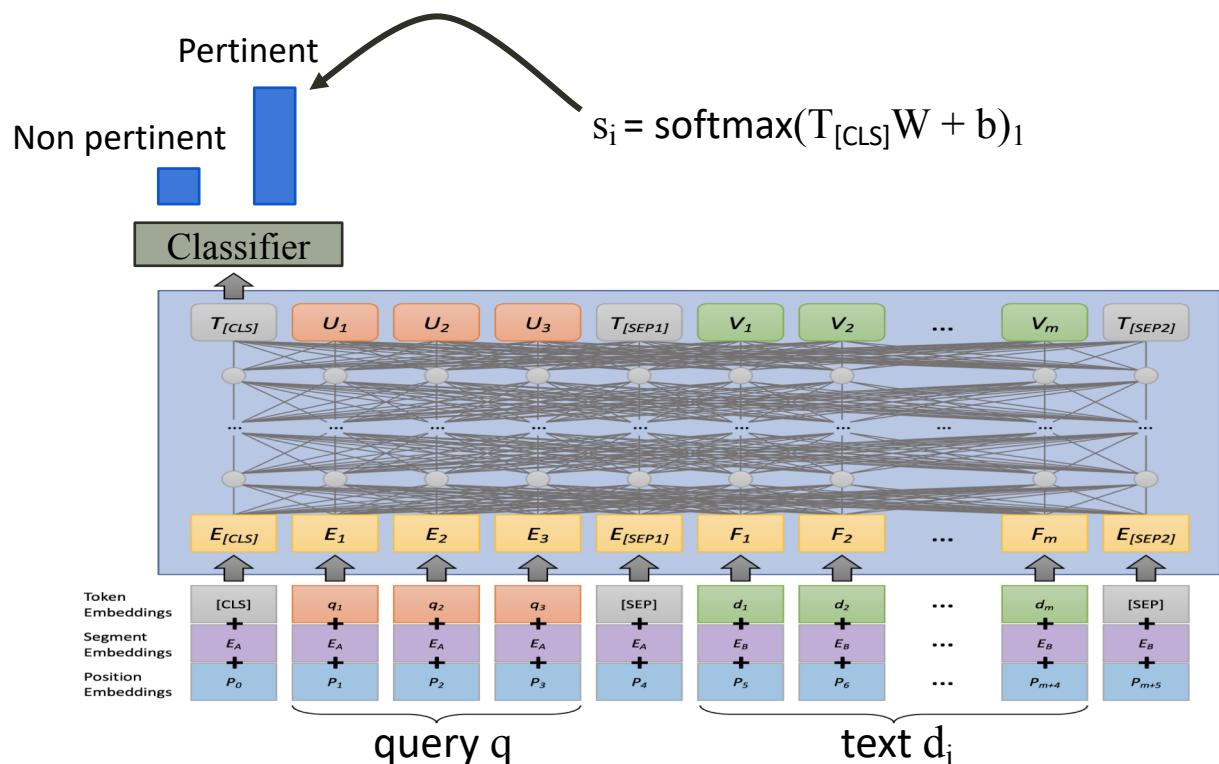
[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]  
Sentence A Sentence B

- Affinage (utilisation de BERT pour entraîner sa propre tâche)



## Affinage (Fine-tuning)

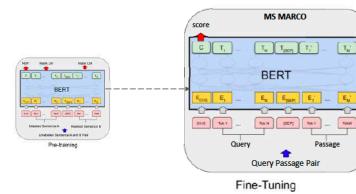
- Utiliser un BERT pré-entraîné (prendre un checkpoint de chez Google ou de huggingface)
- Prendre une collection d'apprentissage : Requêtes-documents Pertinents et non pertinents.
- Fournir en entrée de BERT requête+“partie du doc” (Pb pas plus de 512 tokens)



- Sélectionner une liste de documents (100, 1000) en utilisant un modèle classique (BM25)
- Faire du fine-tuning
  - Utiliser un BERT affine “*finetuné*” (*c'est mieux*)
  - Découper le document en petits morceaux pour faire au total (requête+document+SEP+CLS pas plus 512 mots)
  - Fournir en entrée de BERT requête+“partie du doc”
  - Bert renvoie un score de pertinence

### BERT FOR RANKING

BERT (and other large-scale unsupervised language models) are demonstrating dramatic performance improvements on many IR tasks



Rodrigo Nogueira, and Kyunghyun Cho. *Passage Re-ranking with BERT*. In arXiv. 2019.

## BERT FOR RANKING

Impact across both academia and industry

### Bing says it has been applying BERT since April

The natural language processing capabilities are now applied to all Bing queries globally.

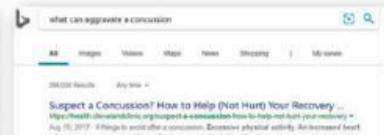
George Nguyen on November 19, 2019 at 1:38 pm

Bing has been using BERT to improve the quality of search results since April. Microsoft has stated. The transformer models are now applied to every Bing query globally.

BEFORE



AFTER



GeorgeCoutinho - Getty Images/Corbis via Getty Images

Google [recently announced](#) one of the biggest updates to its search algorithm in recent years. By using new neural networking techniques to better understand the intentions behind queries, Google says it can now offer more relevant results for searches in 10 languages in the U.S. In English, jobs support for other languages and features coming later. For featured snippets, the update is already live globally.

In the world of search updates, when algorithm changes are often far more subtle, an update that affects 13% of search queries is pretty big. And it will likely have the world's SEO experts up at night.

Google notes that this update will work best for longer, more conversational queries — and in many ways, that's how Google would really like you to search. These days, because it's easier to understand a full sentence than a sequence of keywords.

Can you get medicine for someone pharmacy?

BEFORE



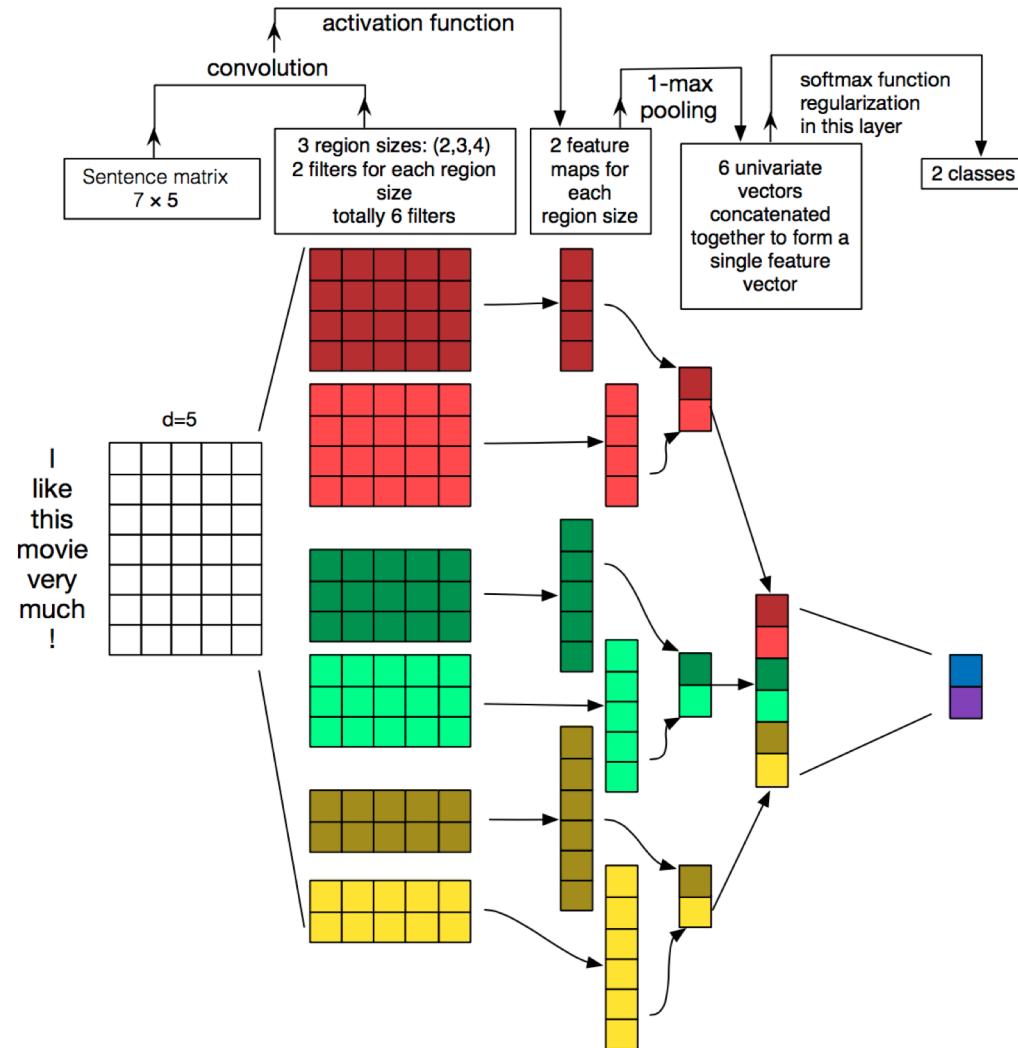
AFTER



- Utiliser le notebook sur Google Colab : [BERT FineTuning with Cloud TPUs.](#)
- Panoplie de modèles préentraînés :  
<https://huggingface.co/>

# Outils pour faire du Deep learning





Zhang and Wallace, 2015; Sentence classification