

Réseaux de neurones

IFT 780

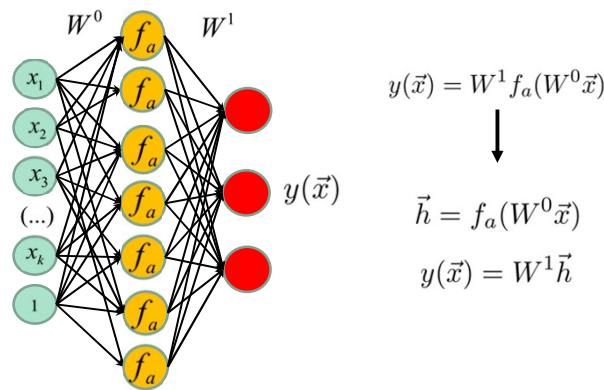
Réseaux récurrents

Par

Pierre-Marc Jodoin, Antoine Théberge

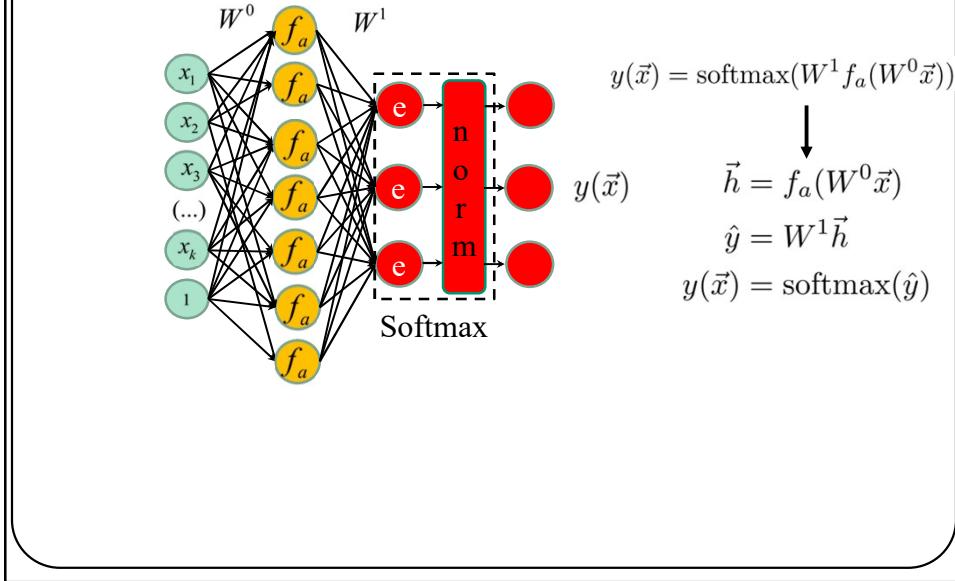
1

Réseau de neurones de base (régression)

 f_a : fonction d'activation

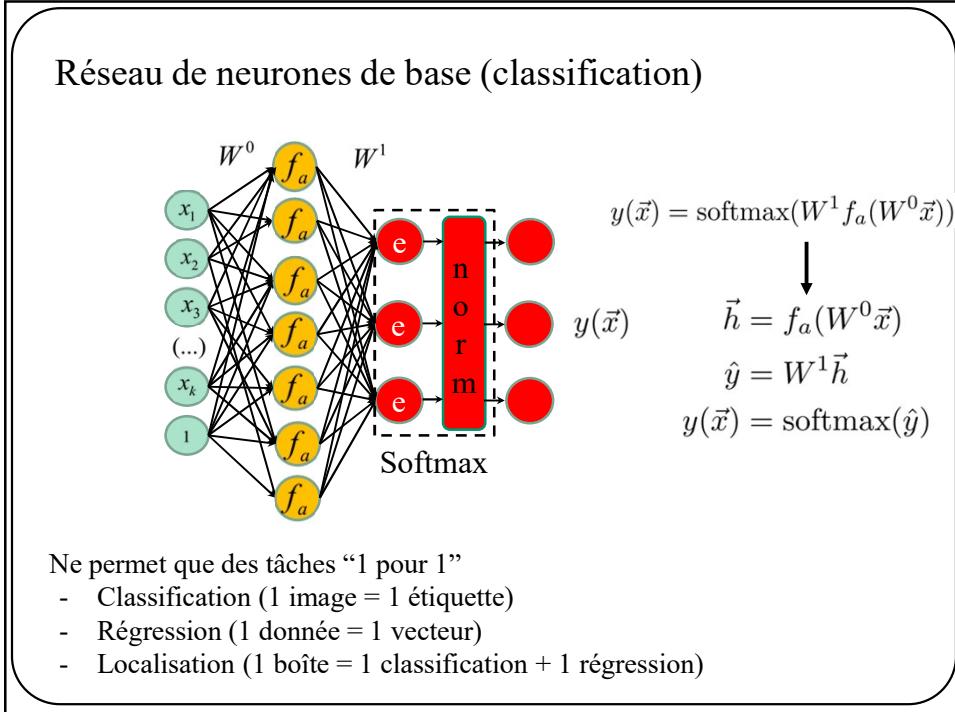
2

Réseau de neurones de base (classification)



3

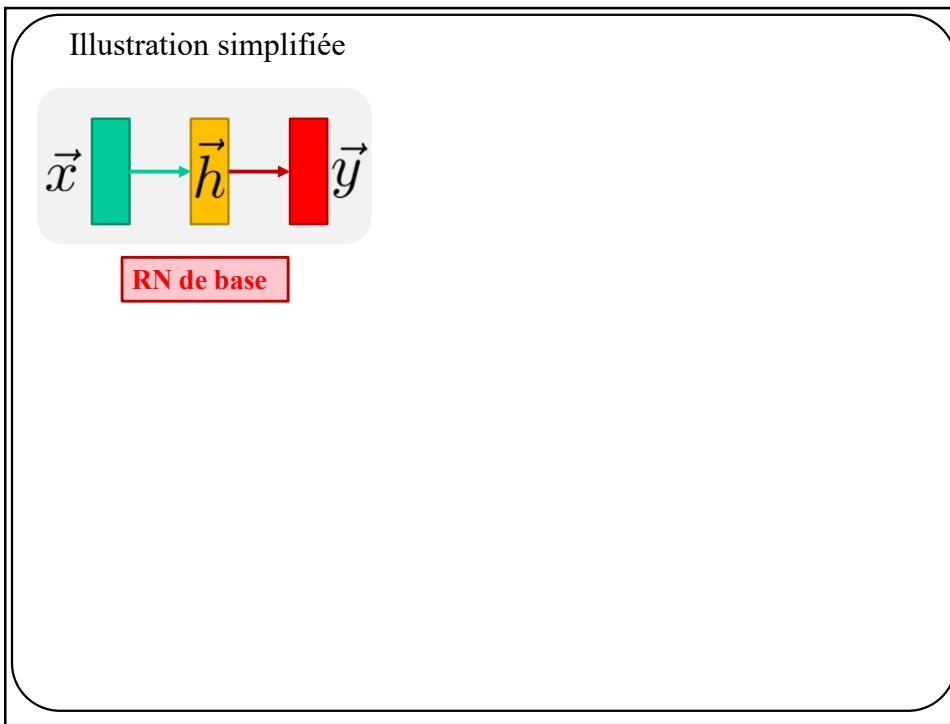
Réseau de neurones de base (classification)



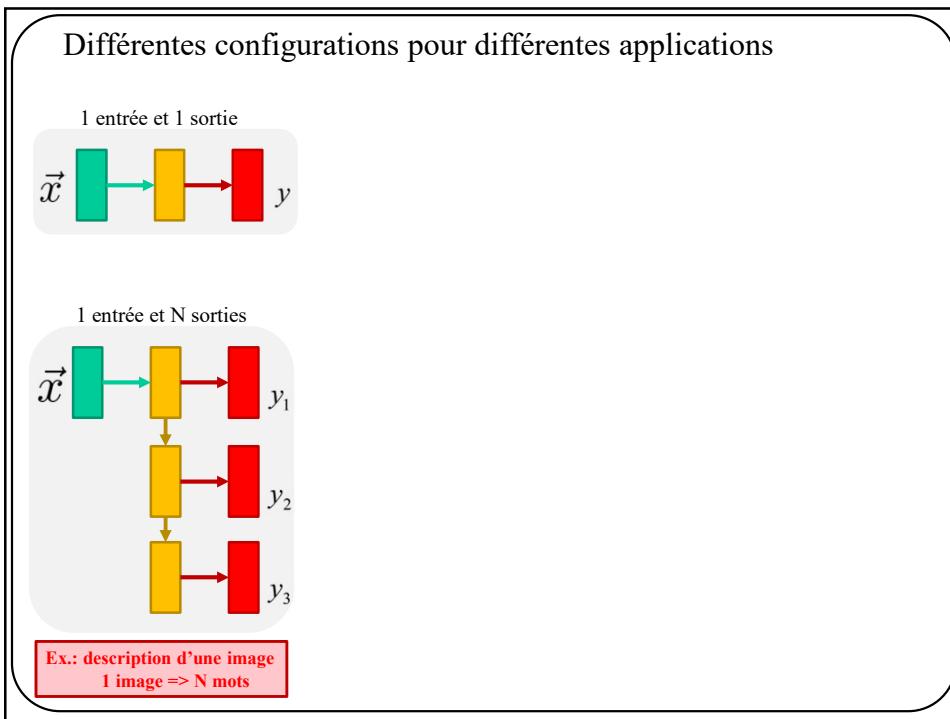
Ne permet que des tâches “1 pour 1”

- Classification (1 image = 1 étiquette)
- Régression (1 donnée = 1 vecteur)
- Localisation (1 boîte = 1 classification + 1 régression)

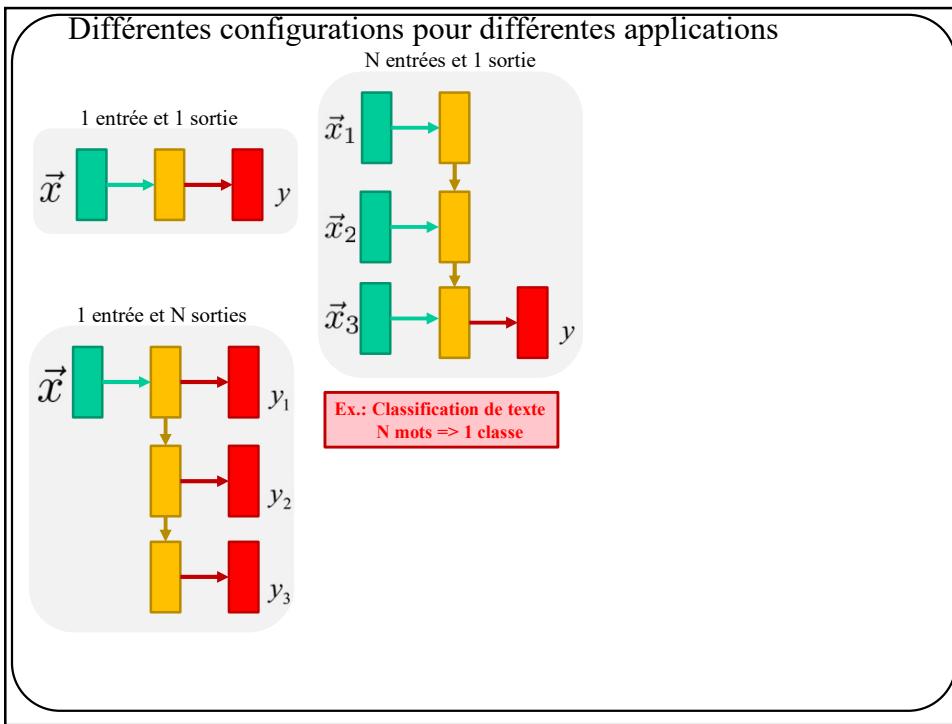
4



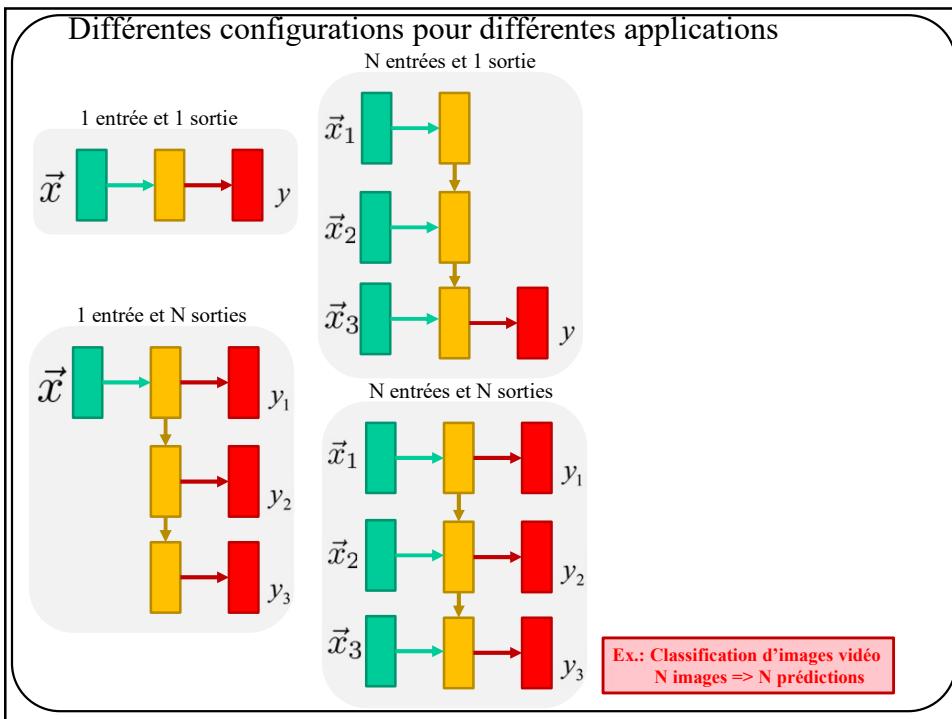
5



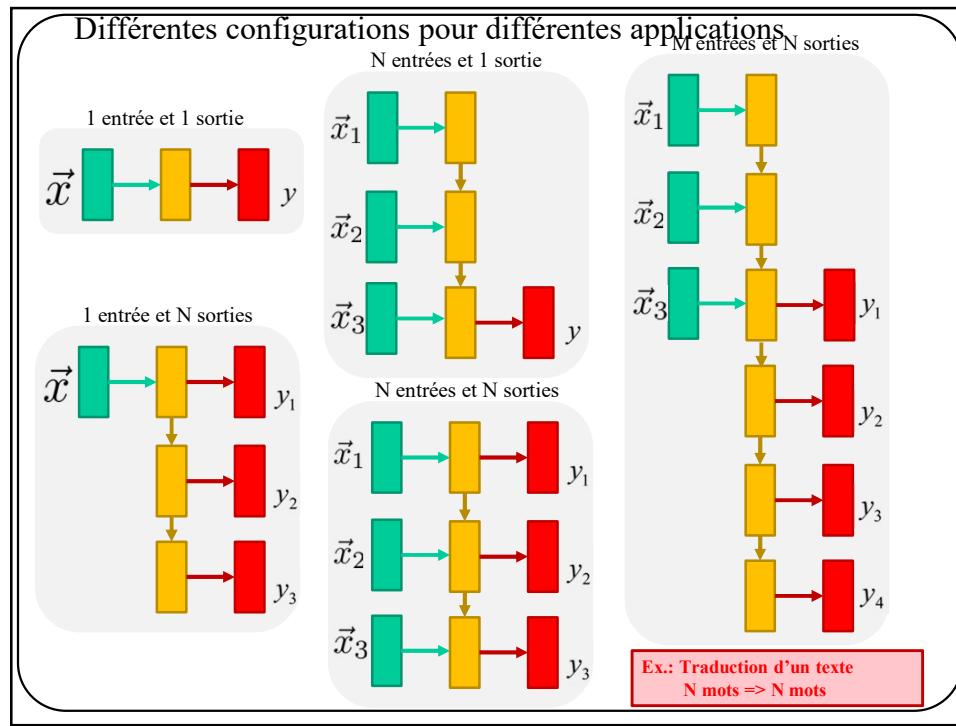
6



7

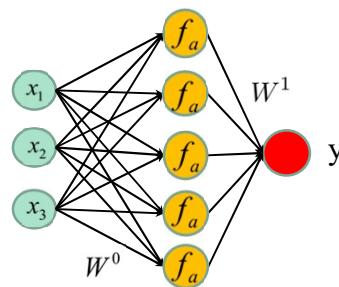


8



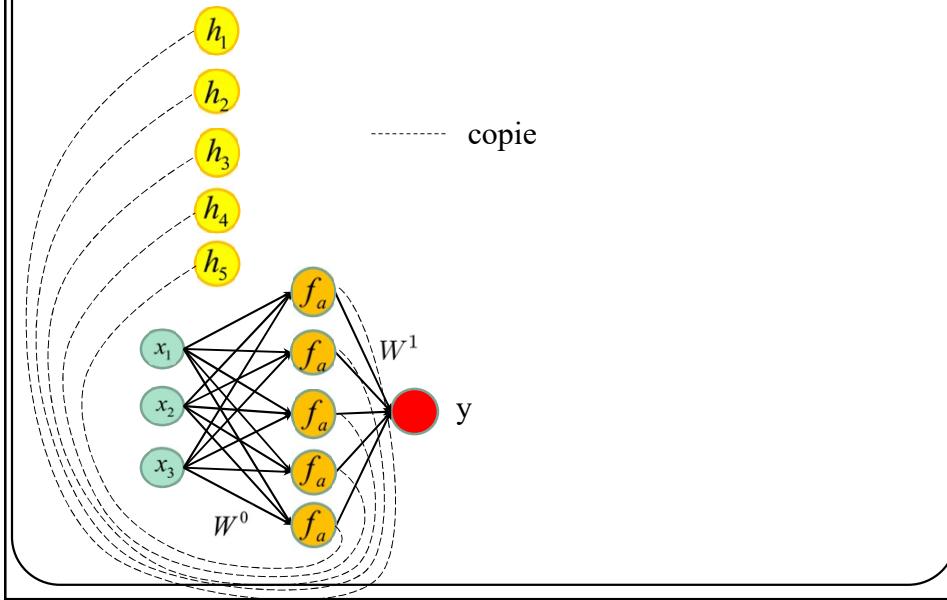
9

Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée



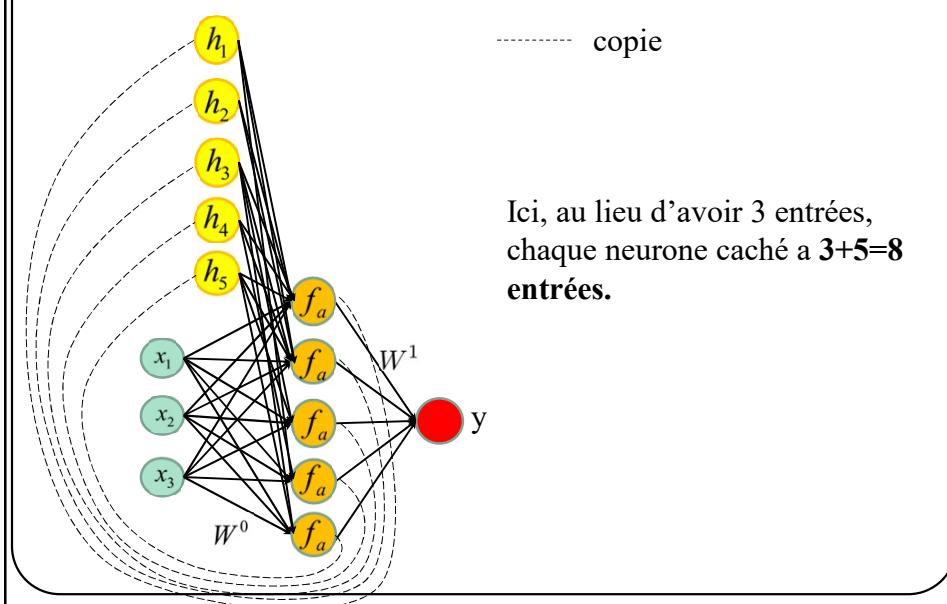
10

Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée



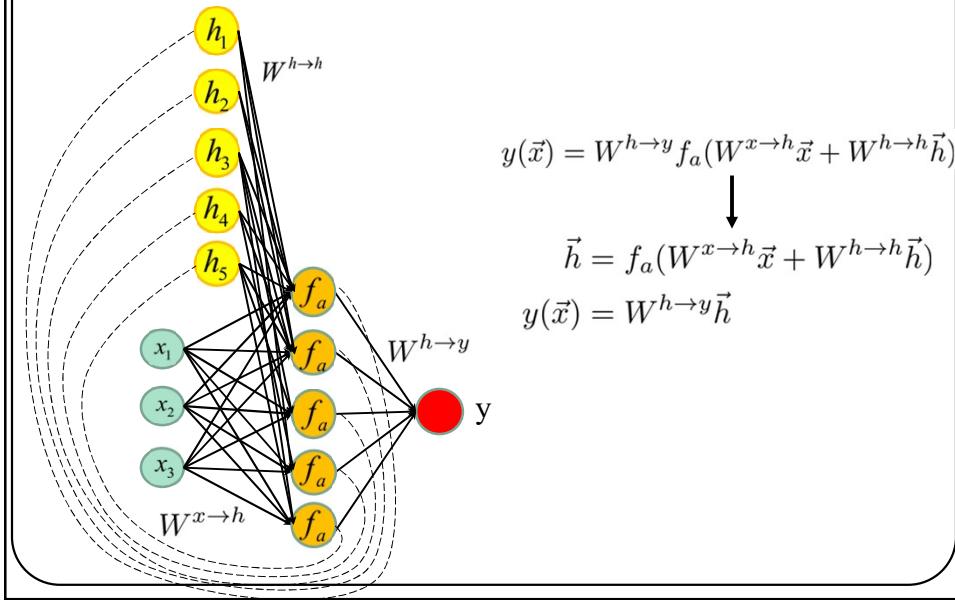
11

Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée



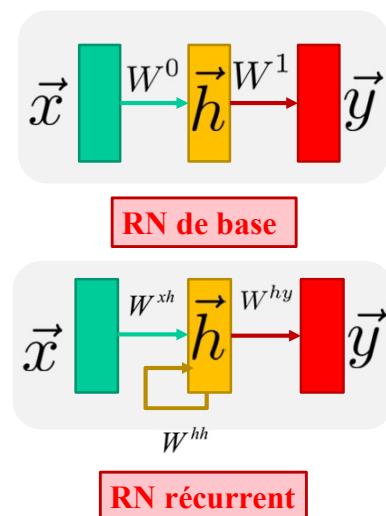
12

Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée



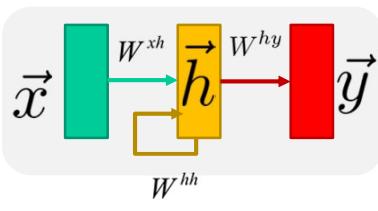
13

Illustration simplifiée



14

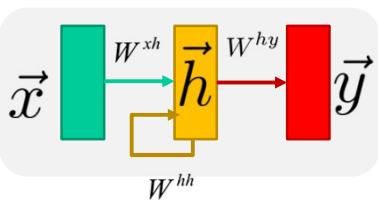
Dans le cas général avec K sorties (régression)



$$\begin{aligned} \vec{y}(\vec{x}) &= W^{hy} f_a(W^{xh} \vec{x} + W^{hh} \vec{h}) \\ \vec{h} &= f_a(W^{xh} \vec{x} + W^{hh} \vec{h}) \\ \vec{y}(\vec{x}) &= W^{hy} \vec{h} \end{aligned}$$

15

Dans le cas général avec K sorties (classification)

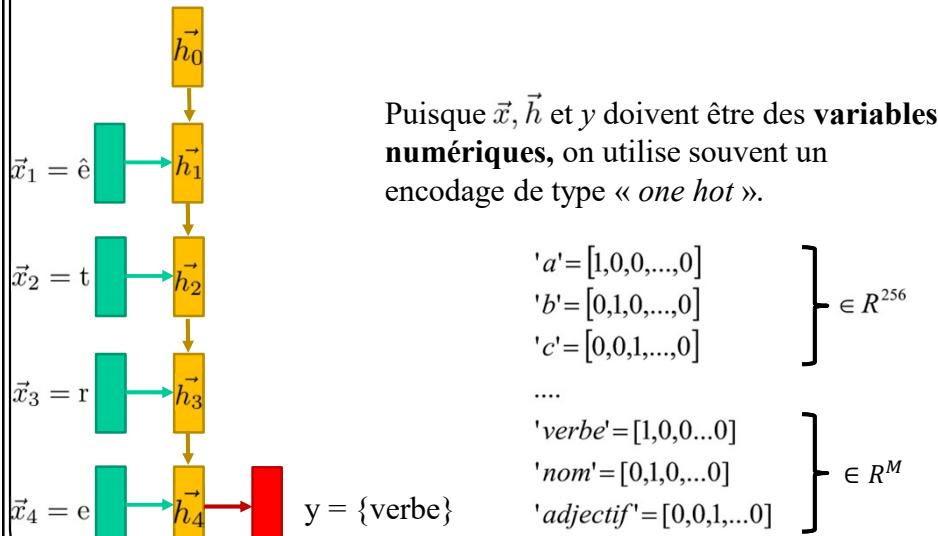


$$\begin{aligned} \vec{y}(\vec{x}) &= W^{hy} f_a(W^{xh} \vec{x} + W^{hh} \vec{h}) \\ \vec{h} &= f_a(W^{xh} \vec{x} + W^{hh} \vec{h}) \\ \hat{y} &= W^{hy} \vec{h} \\ \vec{y}(\vec{x}) &= \text{softmax}(\hat{y}) \end{aligned}$$

16

Exemple pour N entrées et 1 sortie:

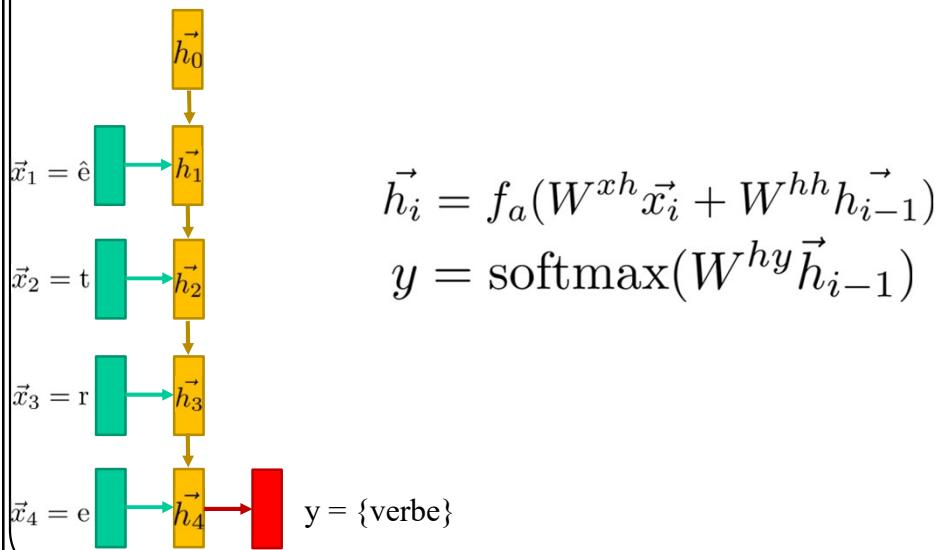
Analyse grammaticale (classification) : (\hat{e}, t, r, e) $\Rightarrow \{\text{verbe}\}$



17

Exemple pour N entrées et 1 sortie:

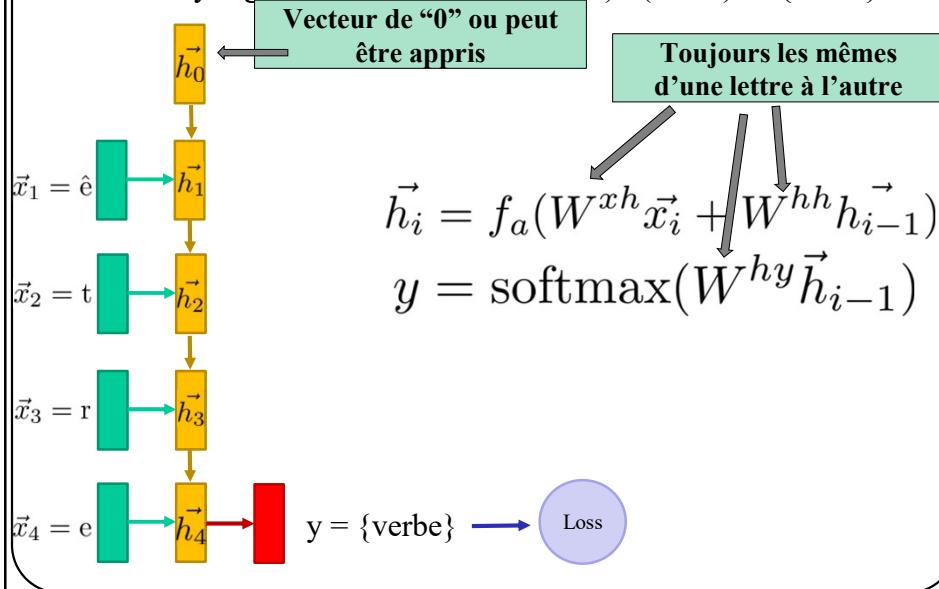
Analyse grammaticale (classification) : (\hat{e}, t, r, e) $\Rightarrow \{\text{verbe}\}$



18

Exemple pour N entrées et 1 sortie:

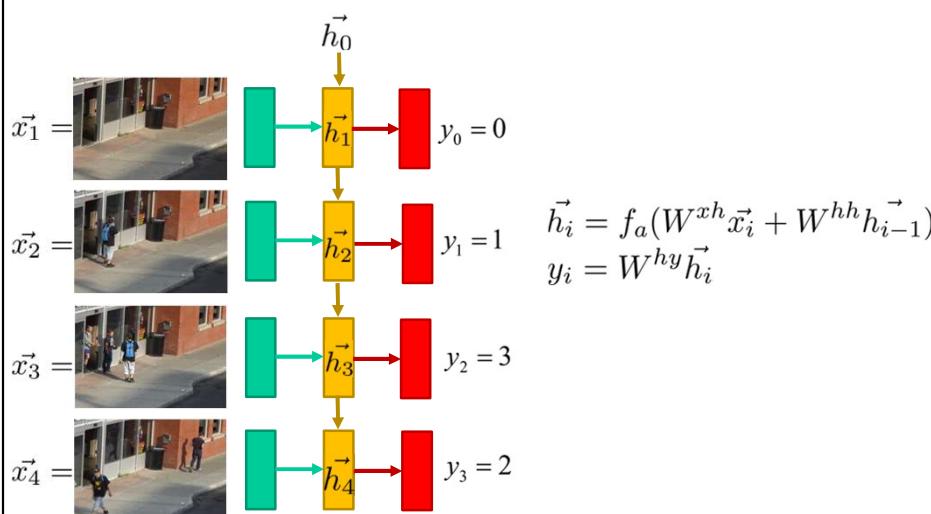
Analyse grammaticale (classification) : (é.t.r.e)=> {verbe}



19

Même idée pour N entrées et N sorties:

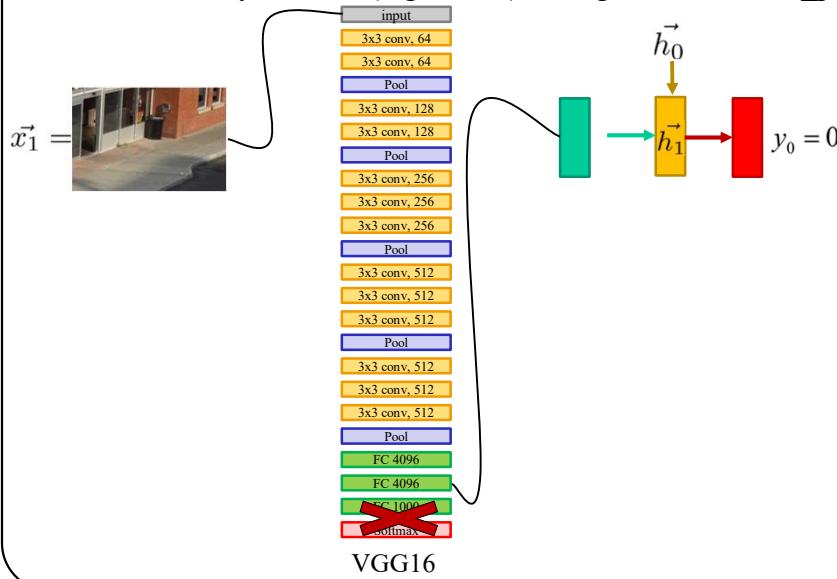
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



20

Même idée pour N entrées et N sorties:

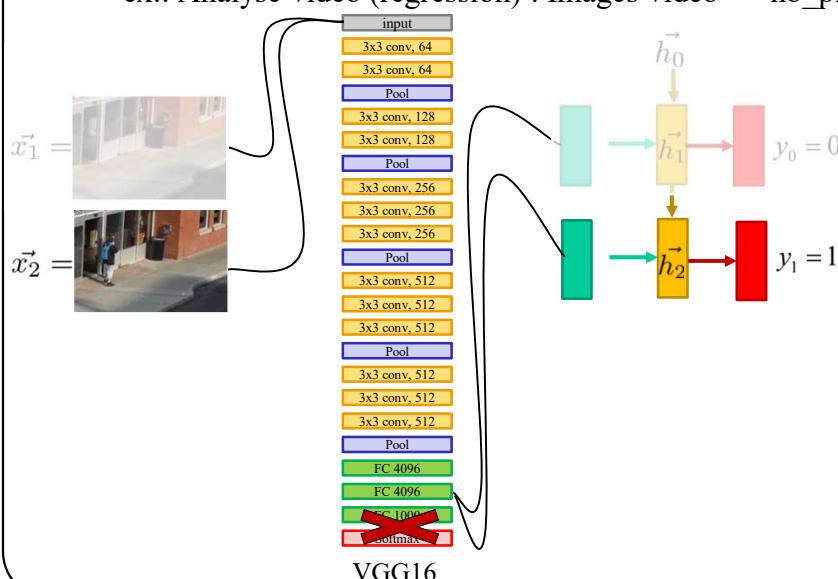
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



21

Même idée pour N entrées et N sorties:

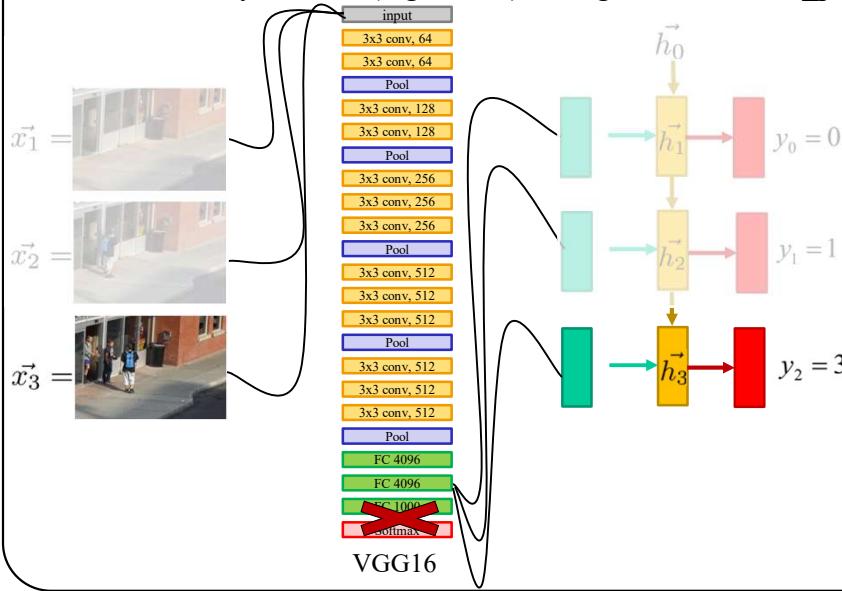
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



22

Même idée pour N entrées et N sorties:

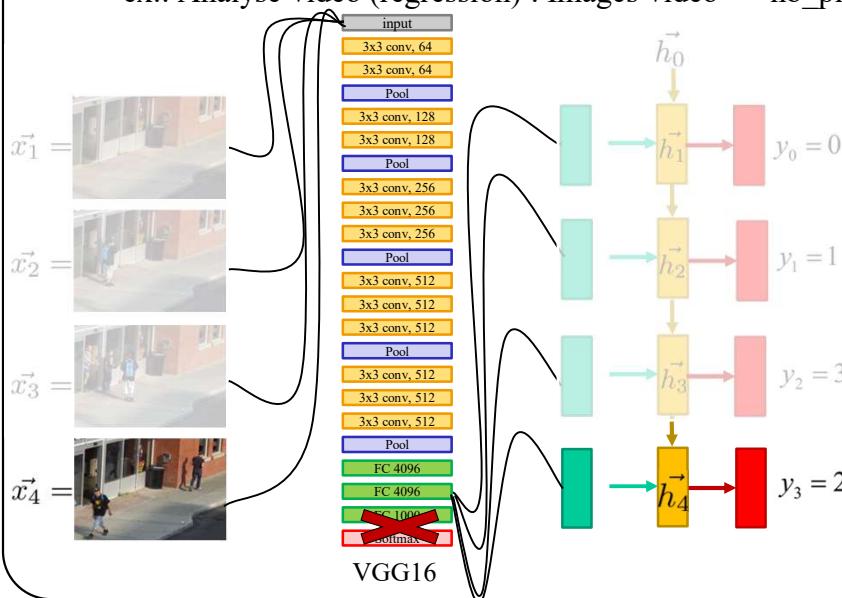
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



23

Même idée pour N entrées et N sorties:

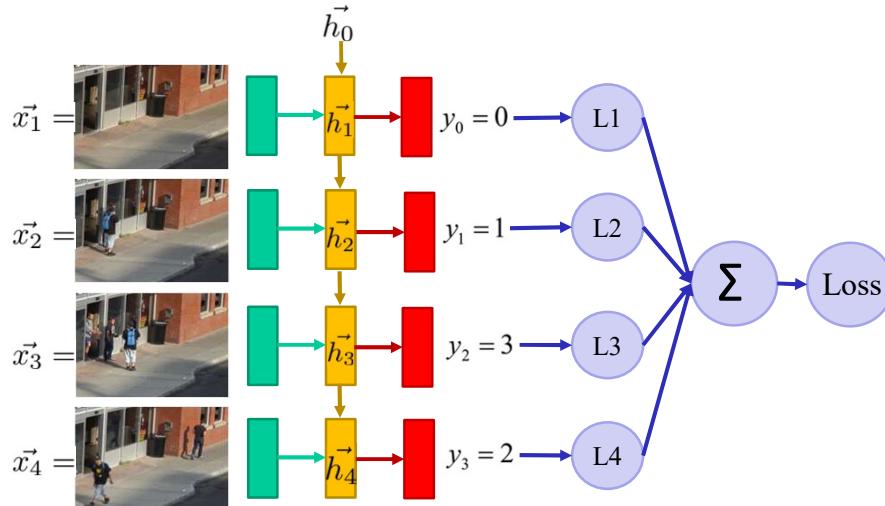
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



24

Même idée pour N entrées et N sorties:

ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



25

Autre exemple: **prédition de caractères** (modèle de langue)

Alphabet jouet : [a,e,m,s]

Représentation « one hot » jouet:

$$'a' = [1, 0, 0, 0]$$

$$'e' = [0, 1, 0, 0]$$

$$'m' = [0, 0, 1, 0]$$

$$'s' = [0, 0, 0, 1]$$

But : Entrainer un modèle à prédire les lettres du mot « **masse** ».

26

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

$$\vec{x}_1 = m \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \rightarrow$$

27

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

$$\vec{x}_1 = m \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{array}{c} W^{xh} \\ \downarrow \\ W^{hh} \end{array}$$

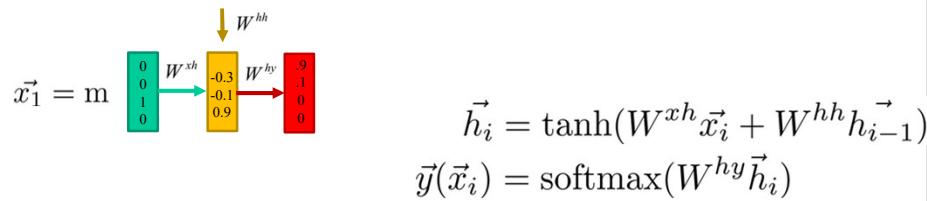
$$\vec{h}_i = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1})$$

28

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

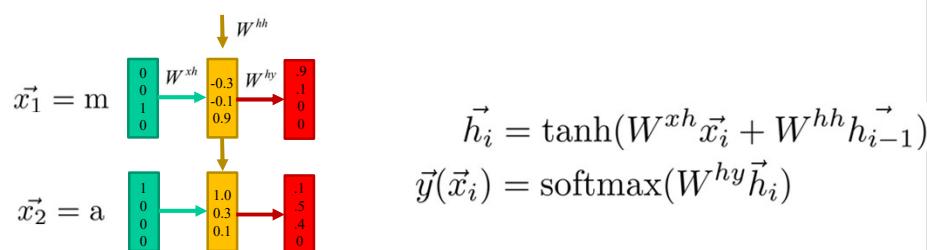


29

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

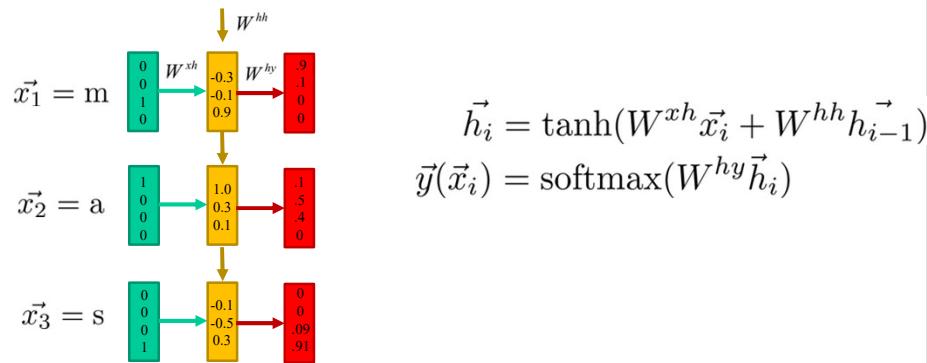


30

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

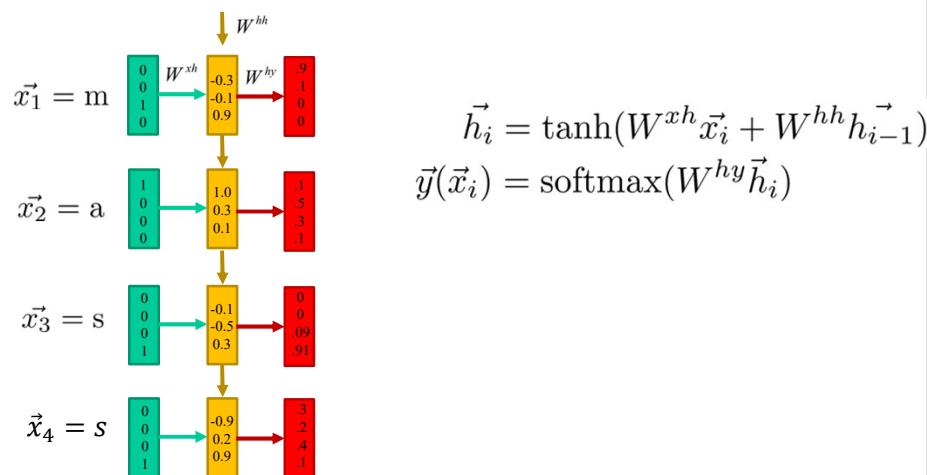


31

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

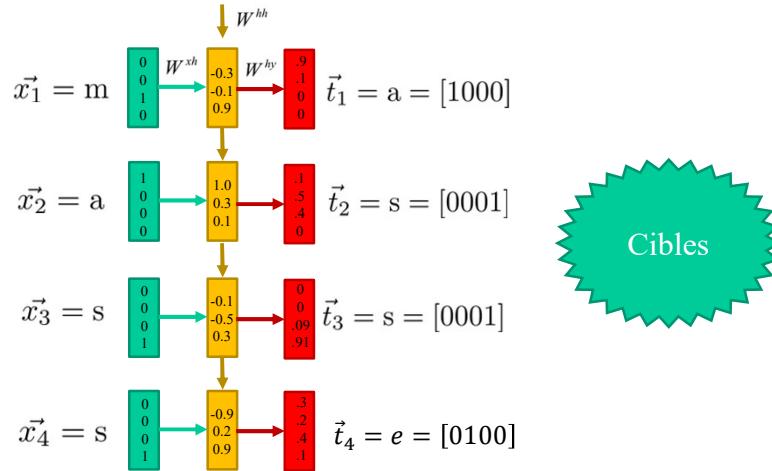


32

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

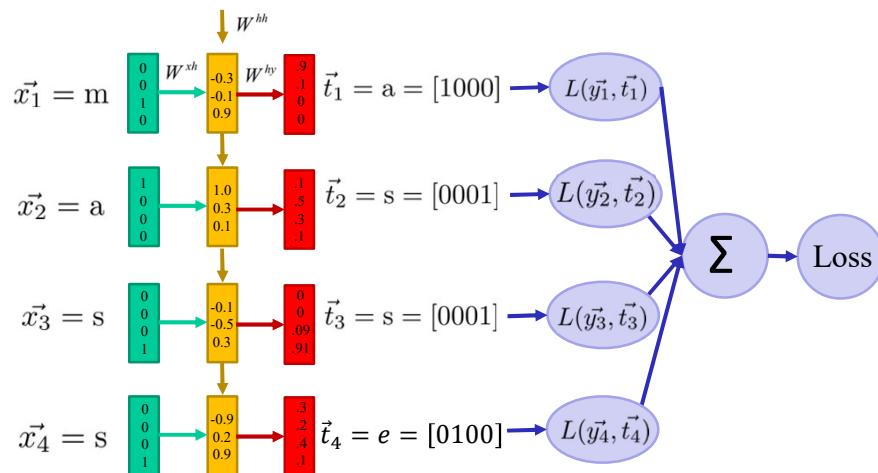


33

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

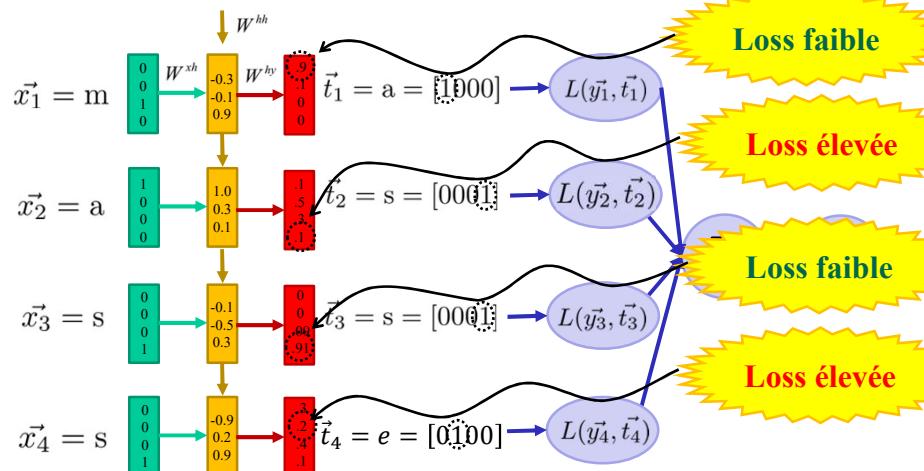


34

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

Entraîner un modèle à prédire les lettres du mot « **masse** ».

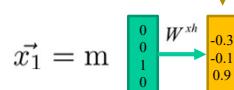


35

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



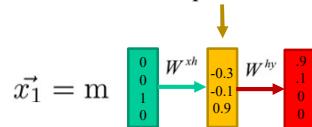
Étape 1 : Calcul de la couche cachée

36

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



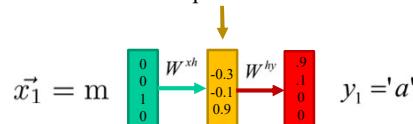
Étape 2 : Calcul de la sortie (softmax)

37

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



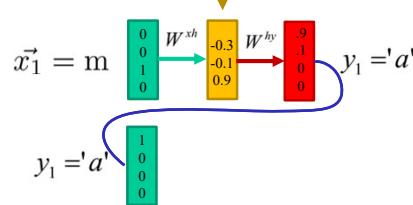
Étape 3 : Sélectionner le caractère le plus probable

38

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



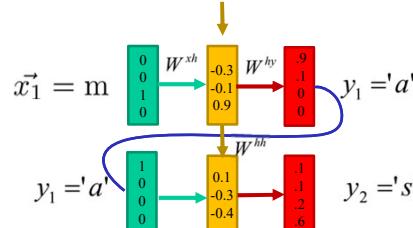
Étape 4 : Injecter le caractère prédit au début du réseau

39

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



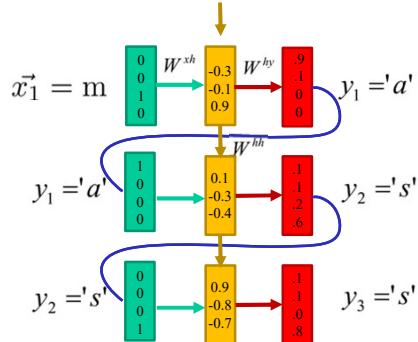
Et on recommence!

40

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

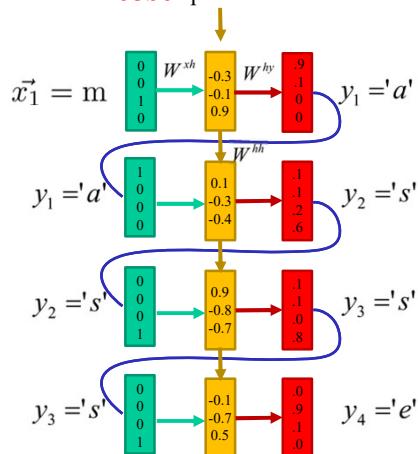


41

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



42

Autre exemple: prédition de caractères (modèle de langue)

Code python: “mini-char-RNN” de A. Karpathy

<https://gist.github.com/karpathy/d4dee566867f8291f086>

Un RNN en 112 lignes !

```
Minimal character-level language model with a Vanilla Recurrent Neural Network, in Python/numpy
git clone https://github.com/karpathy/min-char-rnn.git
cd min-char-rnn
python train.py

# This is a minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
# BSD License
#
# import numpy as np
#
# # data I/O
# data = open('input.txt', 'r').read() # should be simple plain text file
# chars = list(set(data))
# data_size, vocab_size = len(data), len(chars)
# print 'there are %s unique characters, %s unique.' % (data_size, vocab_size)
# char_to_ix = { ch:i for i,ch in enumerate(chars) }
# ix_to_char = { i:ch for i,ch in enumerate(chars) }
#
# # hyperparameters
# hidden_size = 100 # size of hidden layer of neurons
# seq_length = 25 # number of steps to unroll the RNN for
# learning_rate = 1e-1
#
# # model parameters
# Wih = np.random.rand(hidden_size, vocab_size)*0.01 # input to hidden
# Whh = np.random.rand(hidden_size, hidden_size)*0.01 # hidden to hidden
# Why = np.random.rand(vocab_size, hidden_size)*0.01 # hidden to output
# bh = np.zeros((hidden_size, 1)) # hidden bias
# by = np.zeros((vocab_size, 1)) # output bias
#
# def lossfun(inputs, targets, hprev):
#     """ Compute loss and gradient over a sequence """
#     if not inputs:
#         return np.zeros(0)
#     inputs, targets = list(inputs), list(targets)
#     hprev = np.array([0.]).repeat(hidden_size).reshape(1, -1)
#     return lossfun_step(inputs[0], targets[0], hprev)
#
# def lossfun_step(x, y, hprev):
#     """ Compute loss and gradient over a single step """
#     s = np.zeros((vocab_size, hidden_size)) # blank state
#     s[x] = 1
#     hprev = np.copy(hprev)
#     hprev[0] = 0
#     loss = 0
#     for i in range(len(inputs)):
#         s = np.tanh(s.dot(Wih) + hprev.dot(Whh) + Why + bh)
#         s = np.argmax(s)
#         hprev = s.reshape(1, -1)
#         loss += np.log(s[y])
#     return loss
#
# def forward_pass(inputs):
#     s = np.zeros((vocab_size, hidden_size)) # blank state
#     hprev = np.array([0.]).repeat(hidden_size).reshape(1, -1)
#     for i in range(len(inputs)):
#         s = np.tanh(s.dot(Wih) + hprev.dot(Whh) + Why + bh)
#         s = np.argmax(s)
#         hprev = s.reshape(1, -1)
#     return hprev
#
# def softmax(x):
#     e_x = np.exp(x - np.max(x))
#     return e_x / e_x.sum()
#
# def onehot(x, size):
#     if type(x) == int:
#         return np.zeros(size)
#     else:
#         return np.zeros((len(x), size))
#     return np.zeros((len(x), size))
#
# def sample(hprev, size=1):
#     s = np.tanh(hprev.dot(Wih) + Why + bh)
#     s = softmax(s)
#     r = np.random.uniform(0, 1)
#     for i in range(len(s)):
#         if r < s[i]:
#             return ix_to_char[i]
#             break
#     return ix_to_char[-1]
#
# def sample2(hprev, size=1):
#     s = np.tanh(hprev.dot(Wih) + Why + bh)
#     s = softmax(s)
#     r = np.random.uniform(0, 1)
#     for i in range(len(s)):
#         if r < s[i]:
#             return ix_to_char[i]
#             break
#     return ix_to_char[-1]
#
# def sample3(hprev, size=1):
#     s = np.tanh(hprev.dot(Wih) + Why + bh)
#     s = softmax(s)
#     r = np.random.uniform(0, 1)
#     for i in range(len(s)):
#         if r < s[i]:
#             return ix_to_char[i]
#             break
#     return ix_to_char[-1]
```

$$\begin{aligned} 'a' &= [1,0,0,\dots,0] \\ 'b' &= [0,1,0,\dots,0] \\ 'c' &= [0,0,1,\dots,0] \end{aligned} \quad \left. \right\} \in R^{256}$$

...

43

Autre exemple: prédition de caractères (modèle de langue)

Code python: “mini-char-RNN” de A. Karpathy

<https://gist.github.com/karpathy/d4dee566867f8291f086>

THE SONNETS

by William Shakespeare

From fairest creatures we desire increase,
But as the riper should by time decrease,
His tender heir might bear his memory:
But thou, contracted to thine own bright eyes,
Feed'st thy light's flame with self-sensual fuel,
Making a famine where abundance lies;
Thyself thy foe, to thy sweet self too cruel:
Thou art now the world's fresh ornament,
And only herald to the gaudy spring.
Within thine bosom thou dost hold
Abus'd that make'st waste in superfluous¹:
Pity the world, or else this glutton be,
To see the world's due, by the grave, and thee.

When forty winters shall besiege thy brow,
And dig deep trenches in thy beauty's field,
Thy youth's proud livery so garded on now,
Will be a satyr's mask, old age worth held:
Then being asked, where all thy beauty lies,
Where all the treasure of thy fury days,
To say, within thine own deep sunken eyes,
Were an all-eating shame, and thirties praise.
How much more pale, prouder, and more wan,
If these grizzled locks, This fair child of mine
Should sum my count, and make my old excuse,
Proving his beauty by succession thine!

¹ were to be new made when thou art old,
And see thy blood warm when thou feel it cold.

44

Autre exemple: prédition de caractères (modèle de langue)

tyntd-iafhatawiaoahrdemot lytdws e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tkldg t o idoe ns,smtt h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwyl fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.

↓ train more

"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.

Crédit: A. Karpathy, CS231

45

Texte généré une fois le modèle entraîné

PANDARUS:
Alas, I think he shall be come approached and the day
When little strain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:
Come, sir, I will make did behold your worship.

VIOLA:
I'll drink it.

VIOLA:
Why, Salisbury must find his flesh and thought
That which I am not ape, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell;
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:
O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.

Crédit: A. Karpathy, CS231

46

Entraînement sur le code source de Linux en C++

```

static void do_command(struct seq_file *m, void *v)
{
    int column = 32 << (cmd[2] & 0x80);
    if (state)
        cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (k & (1 << i))
            pipe = (in_use & UNXTHREAD_UNCCA) +
                ((count & 0x0000000fffff8) & 0x000000f) << 8;
        if (count == 0)
            subpid, ppc_md.kexec_handle, 0x20000000;
        pipe_set_bytes(i, 0);
    }
    /* Free our user pages pointer to place camera if all dash */
    subsystem_info = kof_changes[PAGE_SIZE];
    rek_controls(offset, idx, &offset);
    /* Now we want to deliberately put it to device */
    control_check_polarity(&context, val, 0);
    for (i = 0; i < COUNTER; i++)
        seq_puts(s, "policy ");
}

```

```

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e20.h>
#include <asm/system_info.h>
#include <asm/seewt.h>
#include <asm/pgproto.h>

#define REQ_PQ     vesa_slot_addr_pack
#define PPN_NOCOMP AFBR(0, load)
#define STACK_DBR(type)          (func)

#define SWAP_ALLOCATE(nr)      (e)
#define emulite_size() arch_get_unaligned_child()
#define access_rw(rw)  __asm volatile("movd %esp, %0" : "+r" (0)) \
if (_type & DQ_READ)

static void stat_PC_SPC __read_mostly offset(struct seq_argqueue,
pC[1]));

static void
os_prefix(unsigned long sys)
{
    #ifndef CONFIG_PREEMPT
    PUT_PARAM_RAID(z, sel) = get_state_state();
    set_pld_sum((unsigned long)state, current_state_str(),
                (unsigned long)->lr_full; low;
    }
}

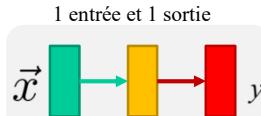
```

Crédit: A. Karpathy, CS231

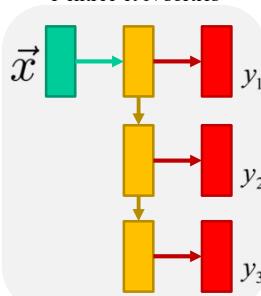
47

Différentes configurations pour différentes applications

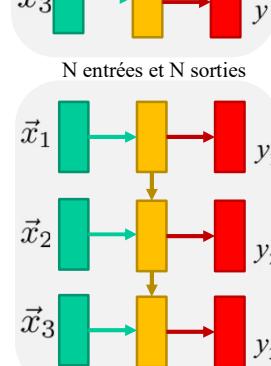
1 entrée et N sorties



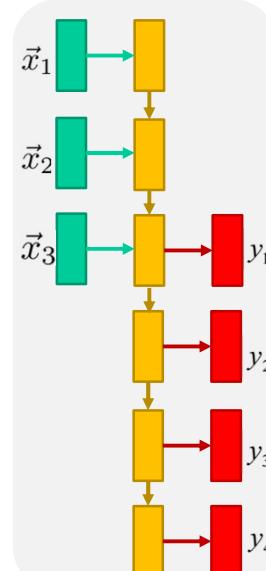
1 entrée et N sorties



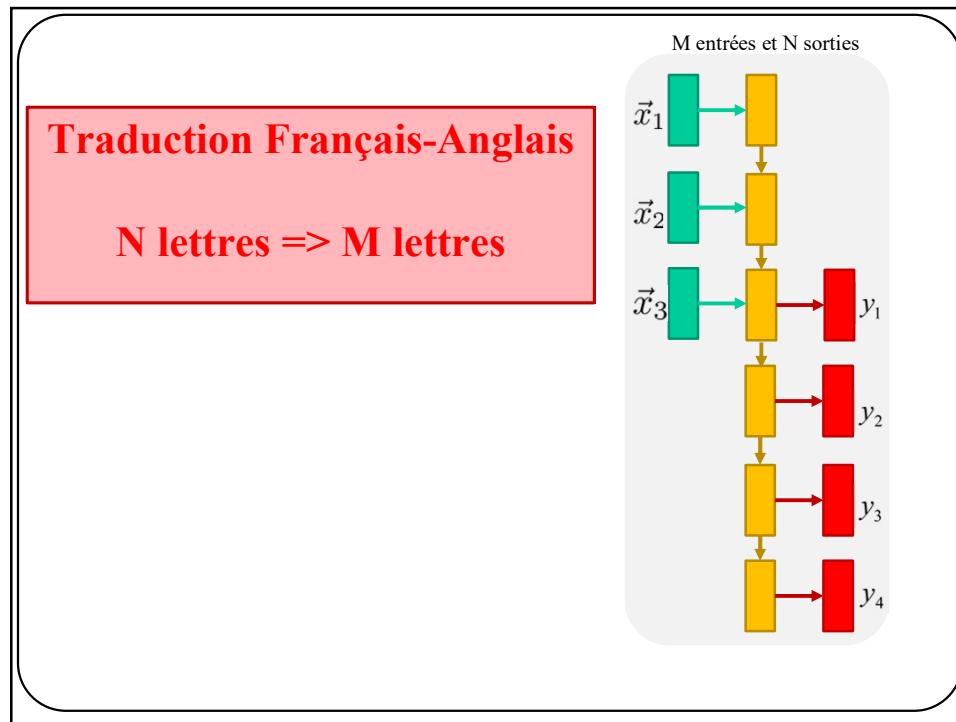
N entrées et N sorties



M entrées et N sorties



48



49

Autre exemple: traduction

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Traduire ‘assez’ -> ‘enough’

Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

Alphabet en: [<BoS>,e,g,h,n,o,u,<EoS>]

Pas le même nombre d’entrées que de sorties !
(BoS : Beginning of Sentence, EoS:End of Sentence).

50

Autre exemple: traduction

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Traduire ‘assez’ -> ‘enough’

Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

Alphabet en: [<BoS>,e,g,h,n,o,u,<EoS>]



51

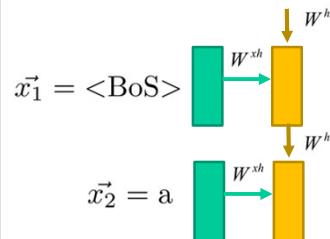
Autre exemple: traduction

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Traduire ‘assez’ -> ‘enough’

Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

Alphabet en: [<BoS>,e,g,h,n,o,u,<EoS>]



52

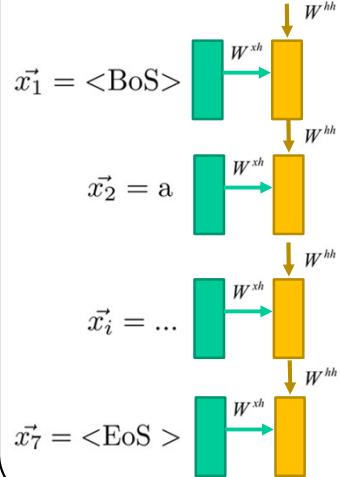
Autre exemple: traduction

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Traduire ‘assez’ -> ‘enough’

Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

Alphabet en: [<BoS>,e,g,h,n,o,u,<EoS>]



53

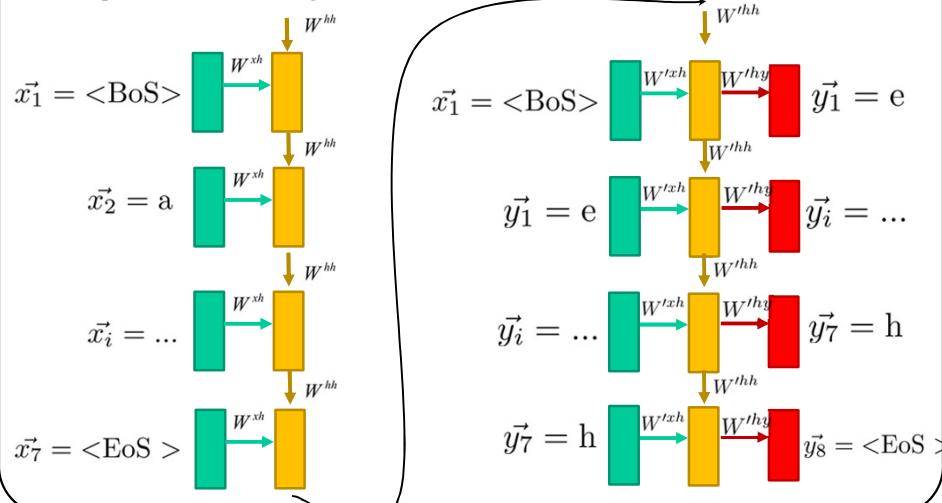
Autre exemple: traduction

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Traduire ‘assez’ -> ‘enough’

Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

Alphabet en: [<BoS>,e,g,h,n,o,u,<EoS>]



54

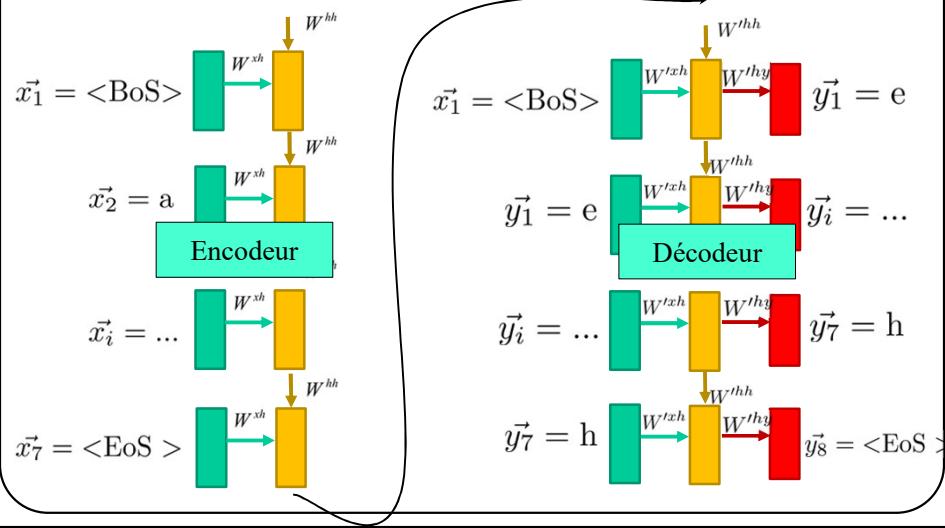
Autre exemple: traduction

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Traduire ‘assez’ -> ‘enough’

Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

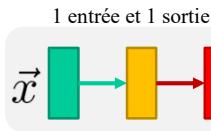
Alphabet en: [<BoS>,e,g,h,n,o,u,<EoS>]



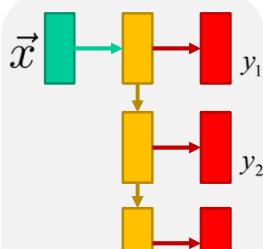
55

Différentes configurations pour différentes applications

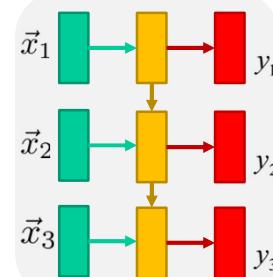
N entrées et 1 sortie



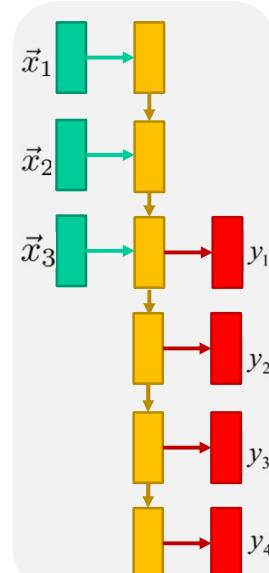
1 entrée et N sorties



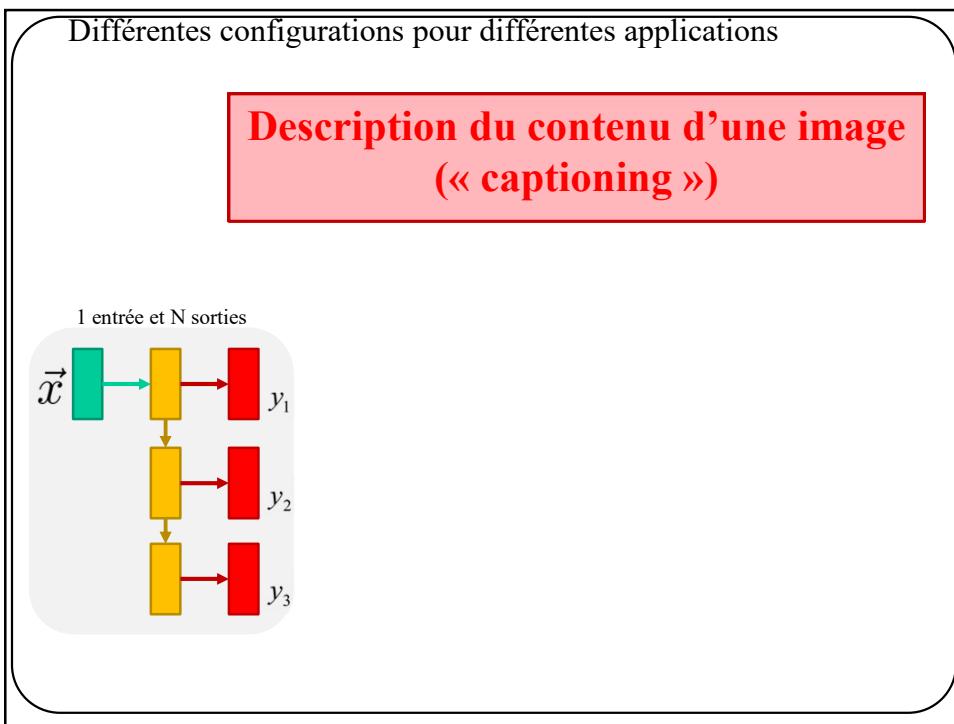
N entrées et N sorties



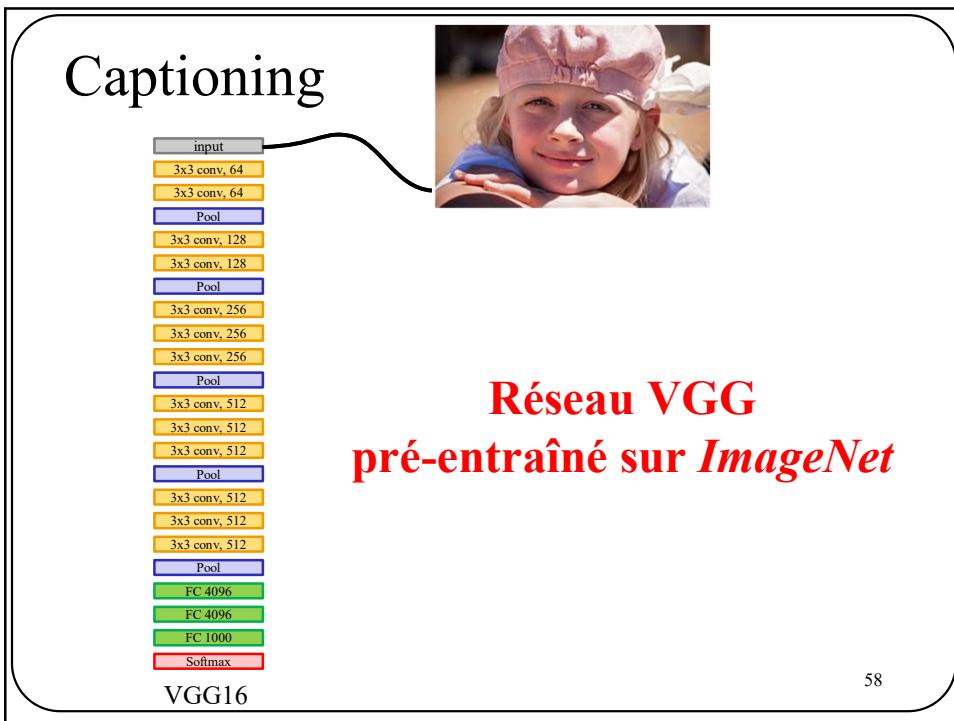
M entrées et N sorties



56

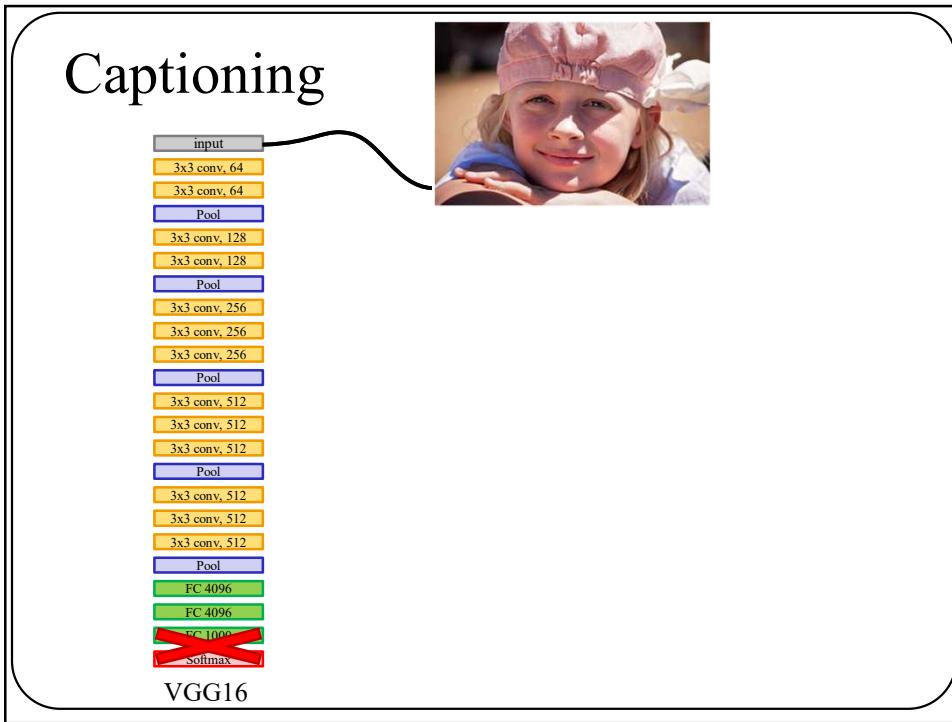


57

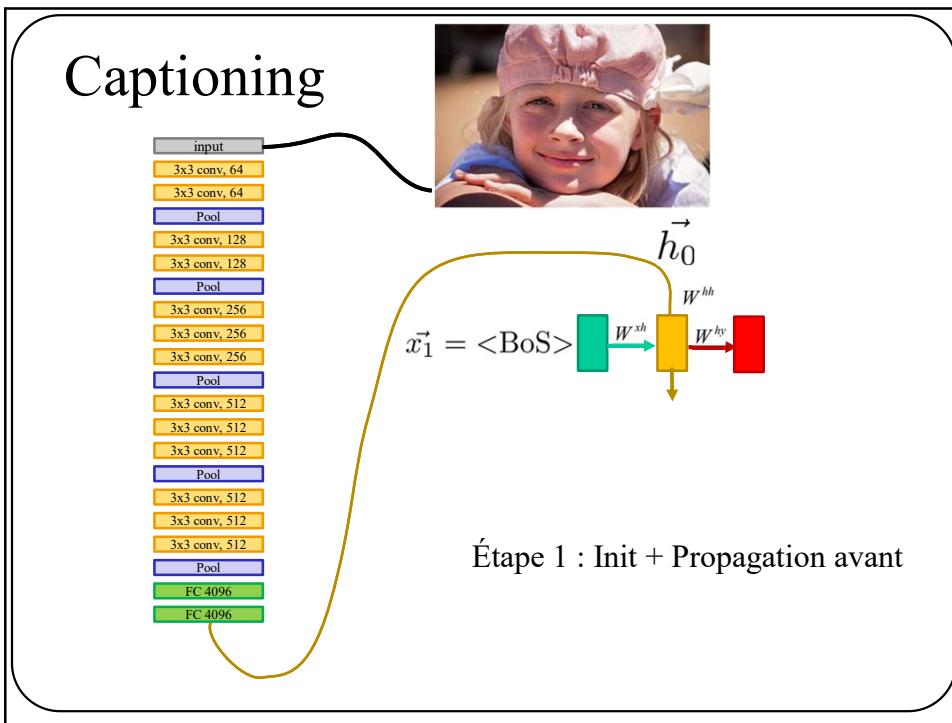


58

58

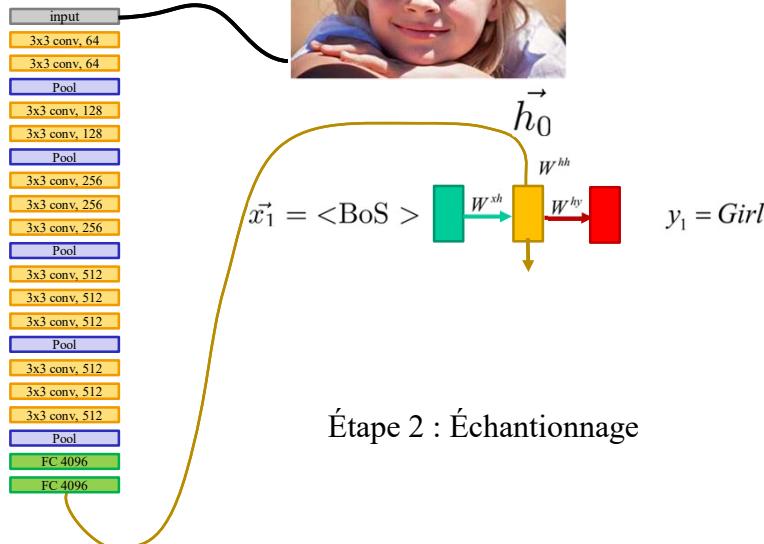


59



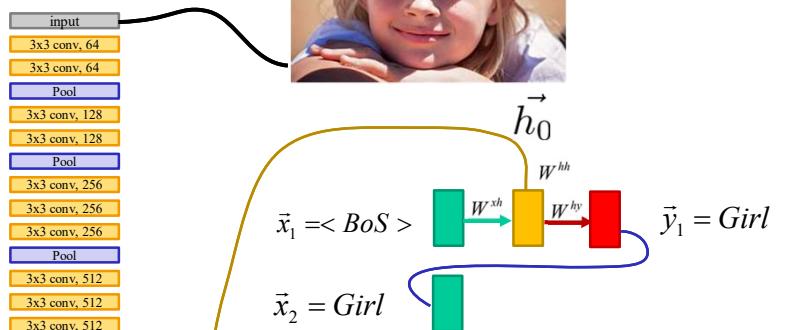
60

Captioning



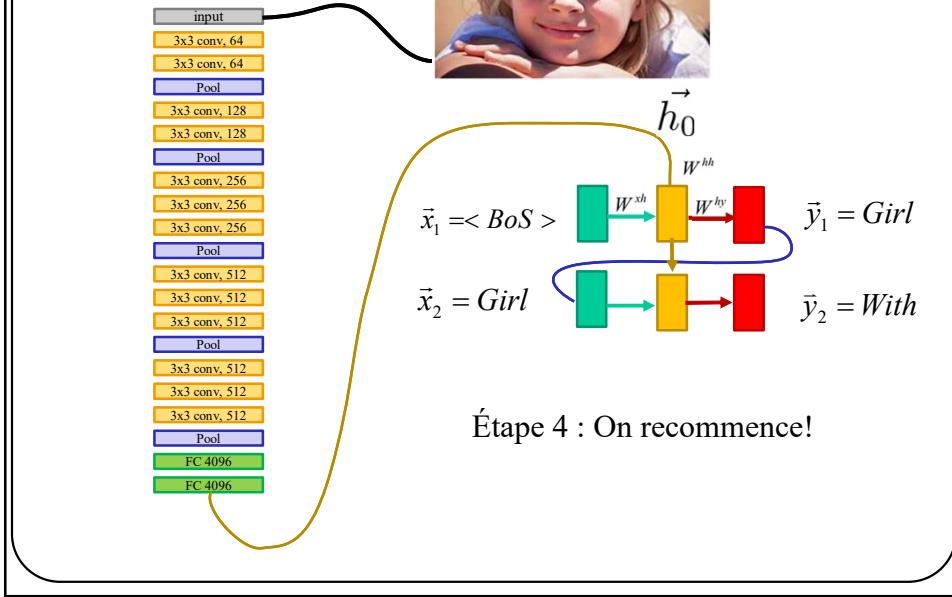
61

Captioning



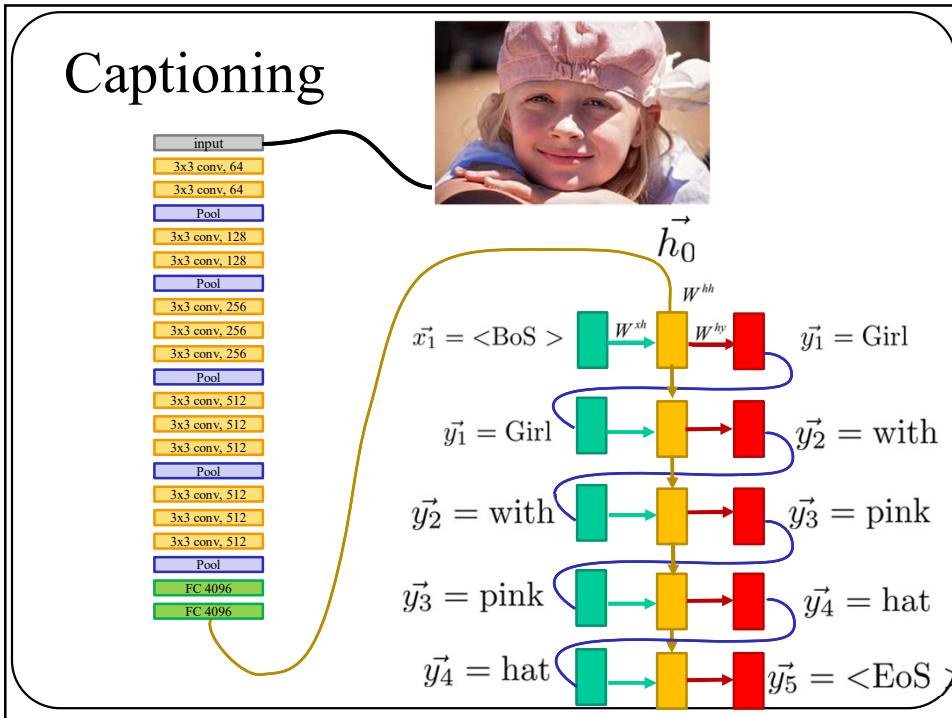
62

Captioning



63

Captioning



64

Exemples de résultats

<https://github.com/karpathy/neuraltalk2>



a elephant standing in a grassy field with trees in the background



a man riding a wave on top of a surfboard



a street sign on a pole in front of a building



a group of people playing a game with nintendo wil controllers



a couple of zebra standing on top of a dirt field

65

65

Exemples d'erreurs

<https://github.com/karpathy/neuraltalk2>



a man is throwing a frisbee in a park



a man riding a skateboard down a street



a laptop computer sitting on top of a wooden desk



A woman is holding a cat in her hand



A woman standing on a beach holding a surfboard

66

66

NeuralTalk and Walk <https://vimeo.com/146492001>



67

67

Analyse de texte

Souvent les modèles de langue utilisent l'encodage « one hot »

Pour des **caractères**...

$$\begin{aligned} 'a' &= [1, 0, 0, \dots, 0] \\ 'b' &= [0, 1, 0, \dots, 0] \\ 'c' &= [0, 0, 1, \dots, 0] \end{aligned} \quad \left. \right\} \in R^{256}$$

...

68

68

Analyse de texte

Souvent les modèles de langue utilisent l'encodage « one hot »

Pour des **mots**...

$$\left. \begin{array}{l} \dots \\ 'grand' = [\dots, 1, 0, 0, \dots, 0] \\ 'grandement' = [\dots, 0, 1, 0, \dots, 0] \\ 'grandeur' = [\dots, 0, 0, 1, \dots, 0] \end{array} \right\} \in R^{10,000}$$

...

69

69

Prédiction sur des lettres vs. mots

$$\left. \begin{array}{l} 'a' = [1, 0, 0, \dots, 0] \\ 'b' = [0, 1, 0, \dots, 0] \\ 'c' = [0, 0, 1, \dots, 0] \end{array} \right\} \in R^{256} \quad \text{Prédiction sur des lettres}$$

...

...

$$\left. \begin{array}{l} 'grand' = [\dots, 1, 0, 0, \dots, 0] \\ 'grandement' = [\dots, 0, 1, 0, \dots, 0] \\ 'grandeur' = [\dots, 0, 0, 1, \dots, 0] \end{array} \right\} \in R^{10,000} \quad \text{Prédiction sur des mots}$$

...

70

Prédiction sur des lettres vs. mots

$$'a' = [1, 0, 0, \dots, 0]$$

$$'b' = [0, 1, 0, \dots, 0]$$

$$'c' = [0, 0, 1, \dots, 0]$$

...

...

$$'grand' = [0, 0, 0, \dots, 1, \dots, 0]$$

$$'granden' = [0, 0, 0, \dots, 1, \dots, 0]$$

$$'grandeu' = [0, 0, 0, \dots, 0, 1, \dots, 0]$$

...

En analyse des langues, un vecteur numérique associé à une séquence de caractères se nomme « JETON » (« token »)

lettres

mots

71

On peut aussi utiliser des fractions de mots

$$'e' = [0, 0, \dots, 1, \dots, 0]$$

...

$$'grand' = [0, 0, \dots, 1, \dots, 0]$$

...

$$'ment' = [0, 0, \dots, 1, \dots, 0]$$

...

$$\left. \begin{array}{l} \dots \\ \dots \\ \dots \\ \dots \end{array} \right\} \in \mathbb{R}^m$$

$'grand'$
 $'grand'+'e'$
 $'grand'+'e'+'ment'$

72

Limites des Jetons « one-hot »

Bien que simple, cet encodage a plusieurs **inconvénients**

1- Peu efficace en mémoire lorsque non compressés

ex.: 10,000 bits pour encoder le mot « **je** » dans une langue à 10,000 mots!

2- Pas de distance sémantique entre les Jetons:

Ex.

distance[one-hot('bon'), one-hot('bien')] = **distance**[one-hot('bon'), one-hot('trottoir')]

Or, on souhaiterait un **code** tel que

distance[code('bon'), code('bien')] << **distance**[code('bon'), code('trottoir')]
distance[code('Jean'), code('Chantal')] << **distance**[code('bon'), code('trottoir')]
distance[code('Inde'), code('Liban')] << **distance**[code('bon'), code('trottoir')]

74

Une solution est d'utiliser l'encodage **Word2Vec** de [Mikolov et al. '13]

Exemple jouet: on veut représenter ces 8 mots par des jetons à 4 éléments

	Jeton « one-hot »							
‘the’	1	0	0	0	0	0	0	0
‘quick’	0	1	0	0	0	0	0	0
‘brown’	0	0	1	0	0	0	0	0
‘fox’	0	0	0	1	0	0	0	0
‘jumps’	0	0	0	0	1	0	0	0
‘over’	0	0	0	0	0	1	0	0
‘lazy’	0	0	0	0	0	0	1	0
‘dog’	0	0	0	0	0	0	0	1

Dictionnaire
de Jetons

2	3	4	5
-1	-3	-2	2
11	6	4	-3
-4	8	-4	4
24	-6	42	17
91	13	14	-5
0	36	4	56
-1	0	1	35



75

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

Comment sélectionner le jeton d'un mot? En multipliant son vecteur One-hot par la matrice d'encodage (le dictionnaire!)

Ex: sélectionner le jeton de « brown »

$$\left(\begin{array}{ccccccccc} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right) \times \text{Dictionnaire Jeton} = \left(\begin{array}{cccc} 11 & 6 & 4 & -3 \end{array} \right)$$

(matrice d'encodage)

2	3	4	5
-1	-3	-2	2
11	6	4	-3
-4	8	-4	4
24	-6	42	17
91	13	14	-5
0	36	4	56
-1	0	1	35

76

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire jeton = matrice d'encodage

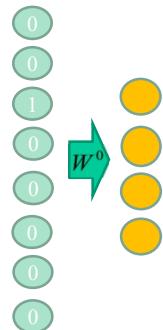
Première couche d'un réseau de neurones



=

matrice d'encodage

\vec{x} : brown



$\cdots W^0 \in R^{4 \times 8}$

77

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

Première couche d'un réseau de neurones

=

matrice d'encodage



$$\text{jeton}_{\vec{x}} = W^0 \vec{x}$$

78

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage



On pourra donc utiliser un réseau de neurones
pour calculer le contenu du dictionnaire

79

Word2Vec s'appuie sur 2 idées fondamentales

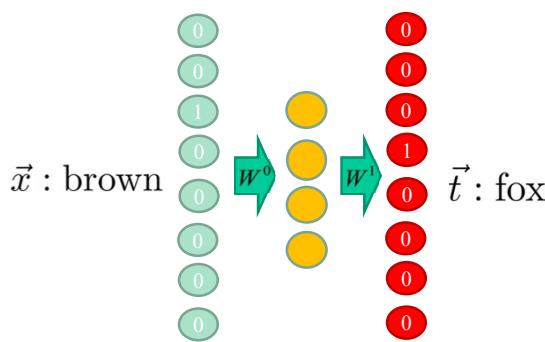
Idée 2: 2 mots proches dans un texte = 2 mots proches sémantiquement

Source Text	Training Samples
The quick brown fox jumps over the lazy dog.	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog.	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog.	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog.	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

Basé sur un corpus de texte, on va créer des **millions de paires de mots**

80

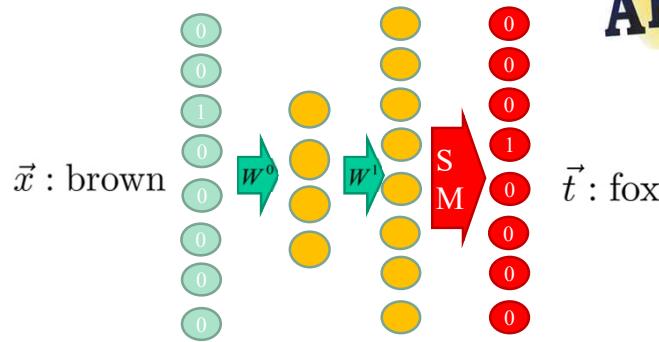
Word2Vec [Mikolov et al. '13]



Entraîner un réseau de neurones
à reproduire le 2^e mot partant du 1^{er}

81

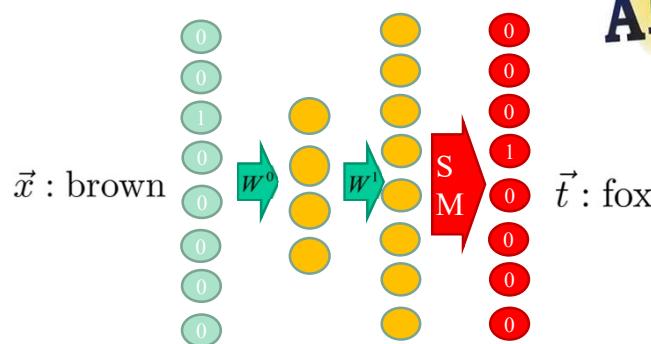
Word2Vec [Mikolov et al. '13]



Puisque la sortie est de type « *one-hot* »
on utilise un softmax

82

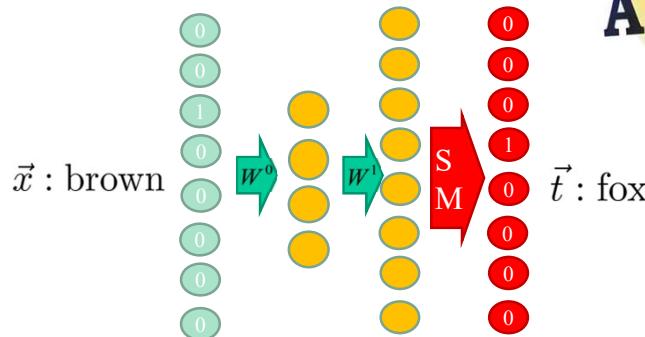
Word2Vec [Mikolov et al. '13]



$$y(\vec{x}) = \text{softmax}(W^1 f_a(W^0 x))$$

83

Word2Vec [Mikolov et al. '13]



Lorsqu'entraîné, utiliser W^0
comme dictionnaire

84

Word2Vec [Mikolov et al. '13]

Cet algorithme vient avec **d'autres détails**

- Réduire l'occurrence des mots fréquents et sémantiquement faibles (*the, of, for, this, or, and, ...*)
- Combiner des mots qui forment une entité (ex: *nations unies*)
- Divers trucs pour simplifier/accélérer l'entraînement

85

Distance sémantique entre deux mots = distance entre leur jeton

Word	First similar word	Second similar word	Third similar word
colosseum	rome (0.994)	roma (0.994)	coliseum (0.994)
colosseo	anfiteatro (0.995)	travel (0.994)	italia (0.994)
scala	aux (0.993)	camelias (0.992)	milano (0.992)
pompeii	retweeted (0.988)	nuovi (0.979)	settembre (0.978)
roma	rome (0.995)	metro (0.994)	colosseum (0.994)
italia	anfiteatro (0.995)	rome (0.995)	colosseo (0.994)
italy	travel (0.998)	davanti (0.997)	photography (0.997)

Word	Similar Words	Similarity	Word	Similar Words	Similarity
Linux	windows	0.85	Twitter	facebook	0.90
	redhat	0.83		instagram	0.86
	unix	0.83		netflix	0.84
	mac os	0.82		snapchat	0.82
	citrix	0.81		google	0.81
	serveurs	0.80		tweets	0.80
	microsoft	0.79		youtube	0.80
	ibm	0.79		linkedin	0.77
	windows server	0.79		maddyness	0.77
	env windows	0.79		tweet	0.77

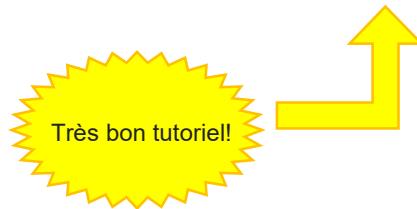
Ahmia, Oussama & Béchet, Nicolas & Marteau, Pierre-Francois. *Two Multilingual Corpora Extracted from the Tenders Electronic Daily for Machine Learning and Machine Translation Applications*.in LREC 2018

86

86

Word2Vec

<http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>



T.Mikolov et al. (2013). "Efficient Estimation of Word Representations in Vector Space", in ICLR 2013

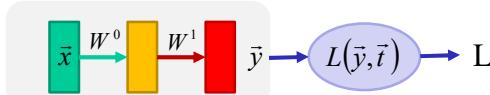
87

Comment entraîner un RNN?

88

Histoire de gradients

RN de classification avec entropie croisée



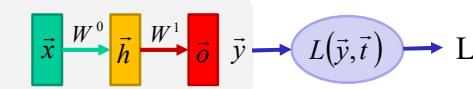
$$\vec{y}(\vec{x}) = S_M \left(W^1 \tanh \left(W^0 \vec{x} \right) \right)$$

$$L = L_{EC} (\vec{y}, \vec{t})$$

89

Histoire de gradients

Simple RN de classification avec entropie croisée



$$\vec{h} = \tanh(W^0 \vec{x})$$

$$\vec{o} = W^1 \vec{h}$$

$$\vec{y} = S_M(\vec{o})$$

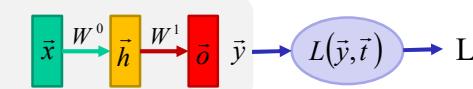
$$L = L_{CE}(\vec{y}, \vec{t})$$

↓
Propagation
avant

90

Histoire de gradients

Simple RN de classification avec entropie croisée



$$\vec{h} = \tanh(W^0 \vec{x})$$

$$\vec{o} = W^1 \vec{h}$$

$$\vec{y} = S_M(\vec{o})$$

$$L = L_{CE}(\vec{y}, \vec{t})$$

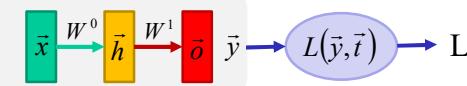
Pour entraîner le réseau
il faut calculer

$$\nabla_{W^0} L \text{ et } \nabla_{W^1} L$$

91

Histoire de gradients

Simple RN de classification avec entropie croisée



$$\vec{h} = \tanh(W^0 \vec{x})$$

$$\vec{o} = W^1 \vec{h}$$

$$\vec{y} = S_M(\vec{o})$$

$$L = L_{CE}(\vec{y}, \vec{t})$$

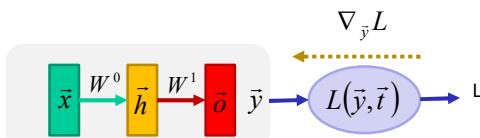
Dérivée en chaîne

$$\nabla_{W^1} L = \nabla_{\vec{y}} L \nabla_{\vec{o}} \vec{y} \nabla_{W^1} \vec{o}$$

$$\nabla_{W^0} L = \nabla_{\vec{y}} L \nabla_{\vec{o}} \vec{y} \nabla_{\vec{h}} \vec{o} \nabla_{W^0} \vec{h}$$

92

Histoire de gradients



$$\vec{h} = \tanh(W^0 \vec{x})$$

$$\vec{o} = W^1 \vec{h}$$

$$\vec{y} = S_M(\vec{o})$$

$$L = L_{CE}(\vec{y}, \vec{t})$$

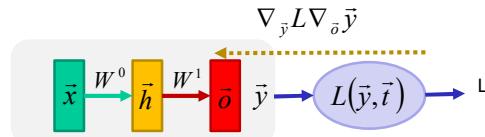
$$\nabla_{\vec{y}} L$$

$$\nabla_{\vec{y}} L = -\frac{\vec{t}}{\vec{y}}$$

Rétro-propagation

93

Histoire de gradients

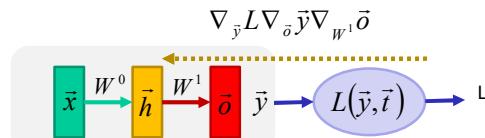


$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \uparrow$$

$$\begin{aligned}
 \nabla_{\vec{o}} \vec{y} &= I \vec{y}^T - \vec{y}^T \vec{y} \\
 \nabla_{\vec{y}} L &= -\frac{\vec{t}}{\vec{y}}
 \end{aligned}
 \quad \text{R\'etro-propagation}$$

94

Histoire de gradients

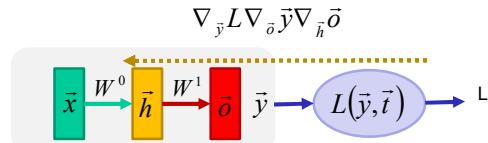


$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \uparrow$$

$$\begin{aligned}
 \nabla_{W^1} \vec{o} &= \vec{h} \\
 \nabla_{\vec{o}} \vec{y} &= I \vec{y}^T - \vec{y}^T \vec{y} \\
 \nabla_{\vec{y}} L &= -\frac{\vec{t}}{\vec{y}}
 \end{aligned}
 \quad \text{R\'etro-propagation}$$

95

Histoire de gradients



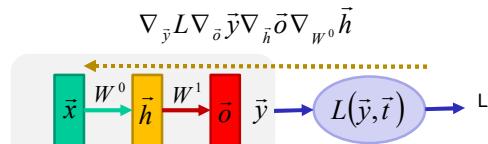
$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \uparrow$$

Rétro-propagation

$$\begin{aligned}
 \nabla_{\vec{h}} \vec{o} &= W^1 \\
 \nabla_{W^1} \vec{o} &= \vec{h} \\
 \nabla_{\vec{o}} \vec{y} &= I \vec{y}^T - \vec{y}^T \vec{y} \\
 \nabla_{\vec{y}} L &= -\frac{\vec{t}}{\vec{y}}
 \end{aligned}$$

96

Histoire de gradients



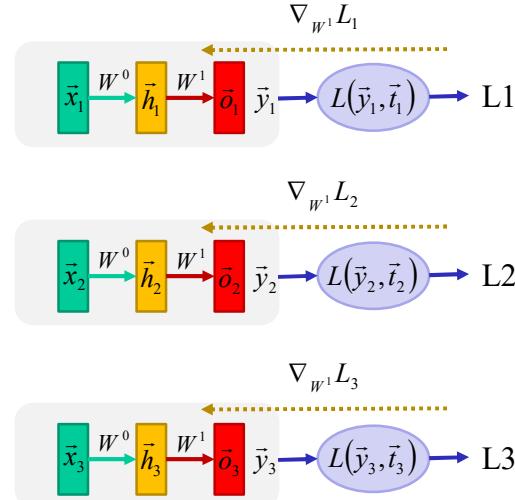
$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \uparrow$$

Rétro-propagation

$$\begin{aligned}
 \nabla_{W^0} \vec{h} &= 1 - \tanh^2(W^0 \vec{x}) \vec{x} \\
 \nabla_{\vec{h}} \vec{o} &= W^1 \\
 \nabla_{W^1} \vec{o} &= \vec{h} \\
 \nabla_{\vec{o}} \vec{y} &= I \vec{y}^T - \vec{y}^T \vec{y} \\
 \nabla_{\vec{y}} L &= -\frac{\vec{t}}{\vec{y}}
 \end{aligned}$$

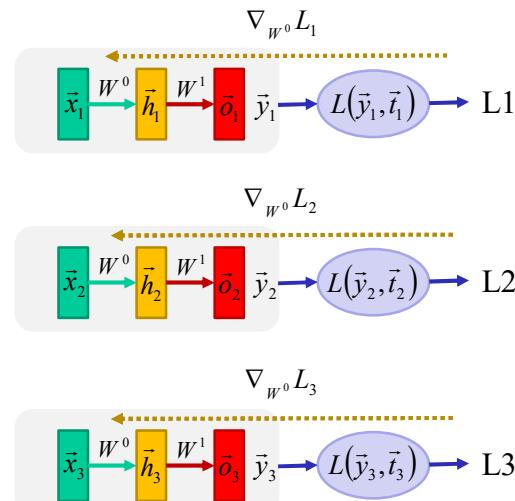
97

Ex.: 3 données, 3 rétro-propagations



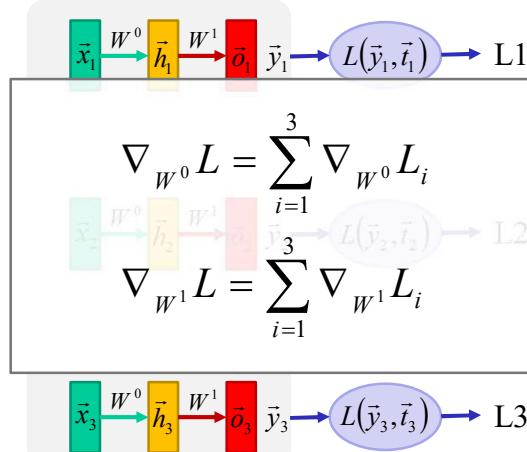
98

Ex.: 3 données, 3 rétro-propagations



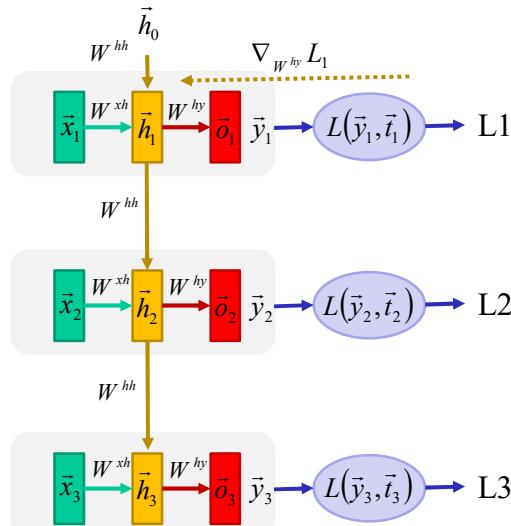
99

3 rétro-propagations



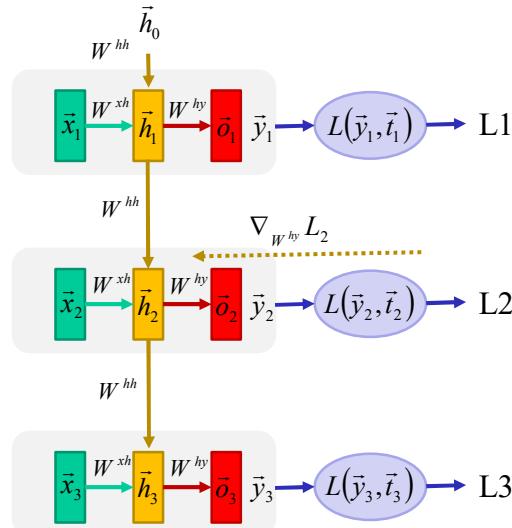
100

Réseau récurrent: gradient pour W^{hy}



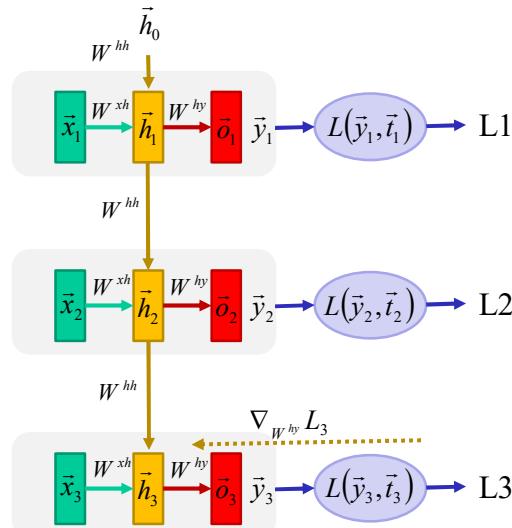
101

Réseau récurrent: gradient pour W^{hy}



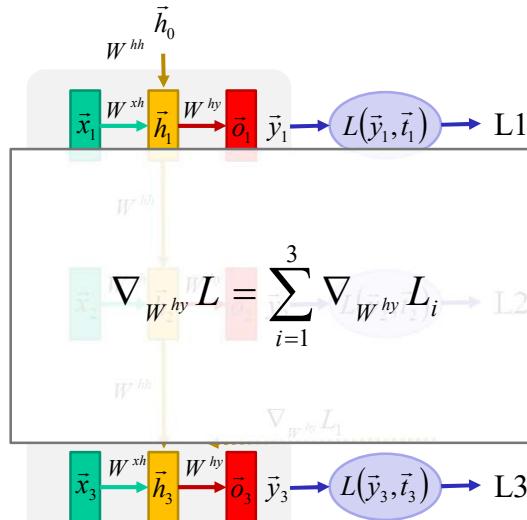
102

Réseau récurrent: gradient pour W^{hy}



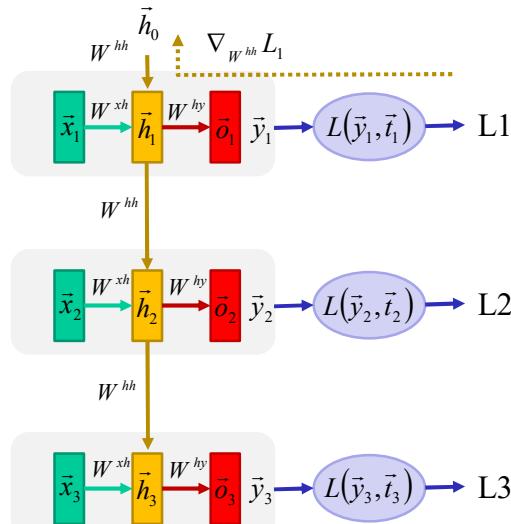
103

Réseau récurrent: gradient pour W^{hy}



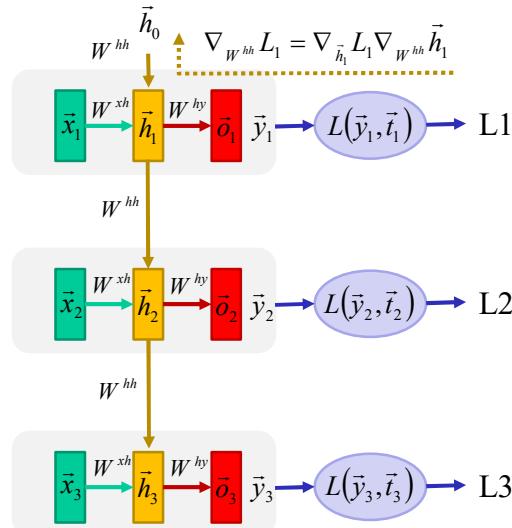
104

Réseau récurrent: gradient pour W^{hh}



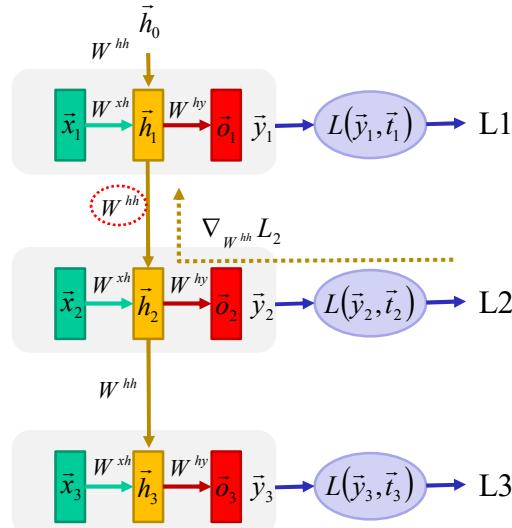
105

Réseau récurrent: gradient pour W^{hh}



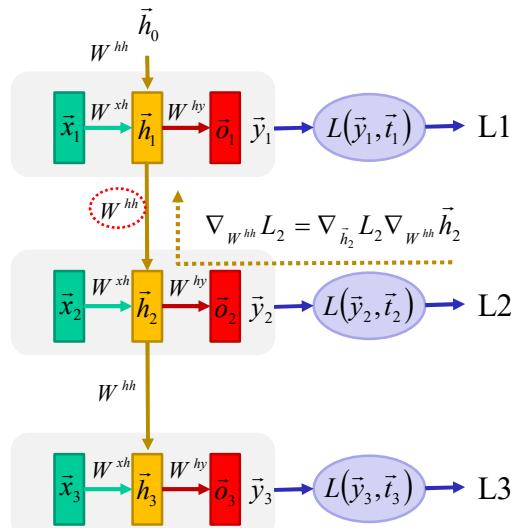
106

Réseau récurrent: gradient pour W^{hh}



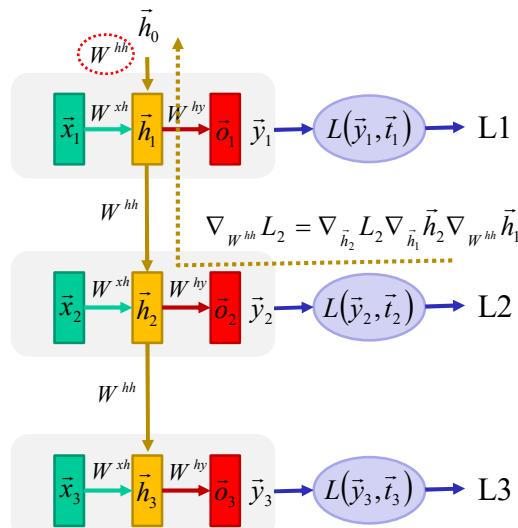
107

Réseau récurrent: gradient pour W^{hh}



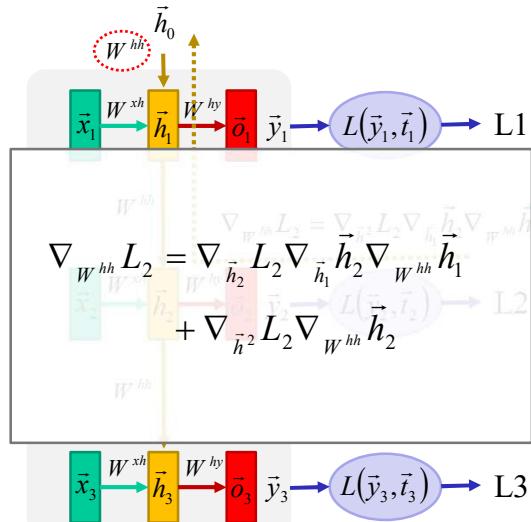
108

Réseau récurrent: gradient pour W^{hh}



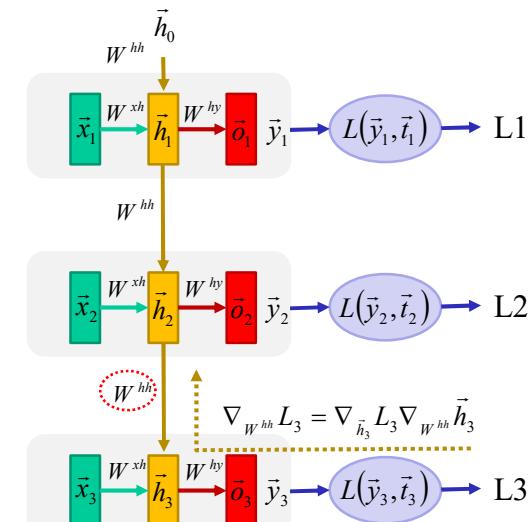
109

Réseau récurrent: gradient pour W^{hh}



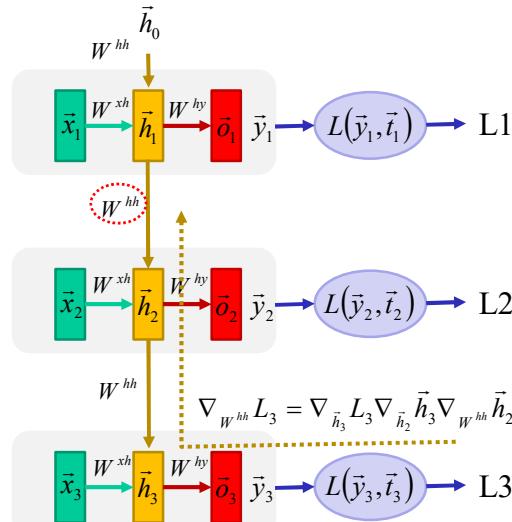
110

Réseau récurrent: gradient pour W^{hh}



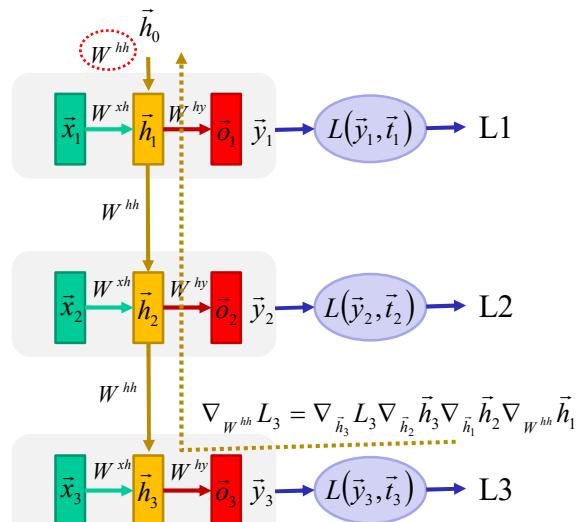
111

Réseau récurrent: gradient pour W^{hh}



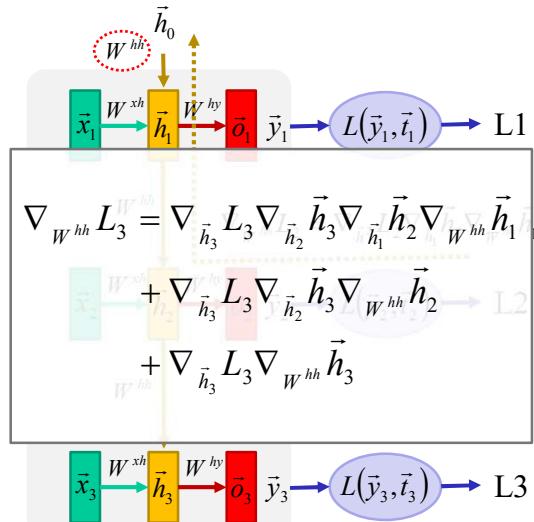
112

Réseau récurrent: gradient pour W^{hh}



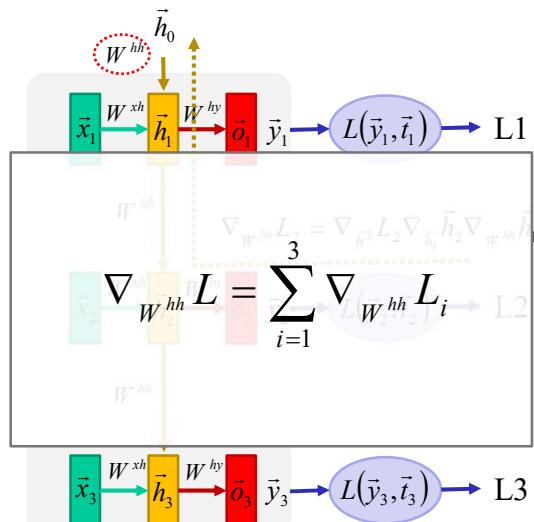
113

Réseau récurrent: gradient pour W^{hh}



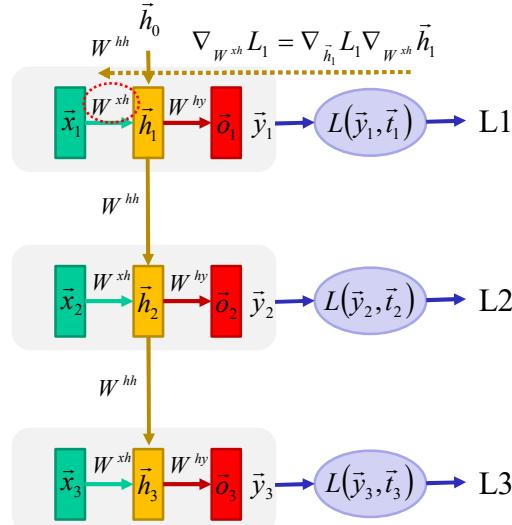
114

Réseau récurrent: gradient pour W^{hh}



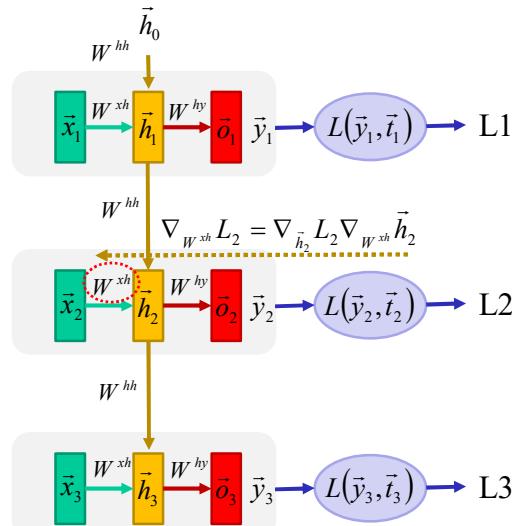
116

Réseau récurrent: gradient pour W^{xh}



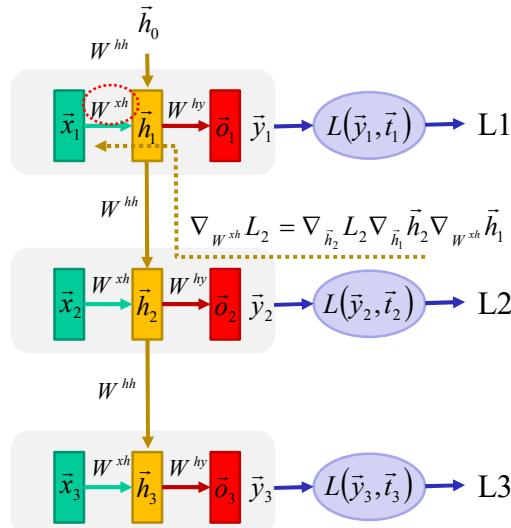
118

Réseau récurrent: gradient pour W^{xh}



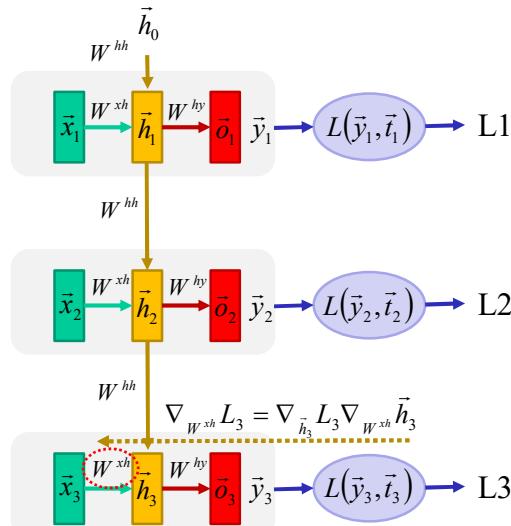
119

Réseau récurrent: gradient pour W^{xh}



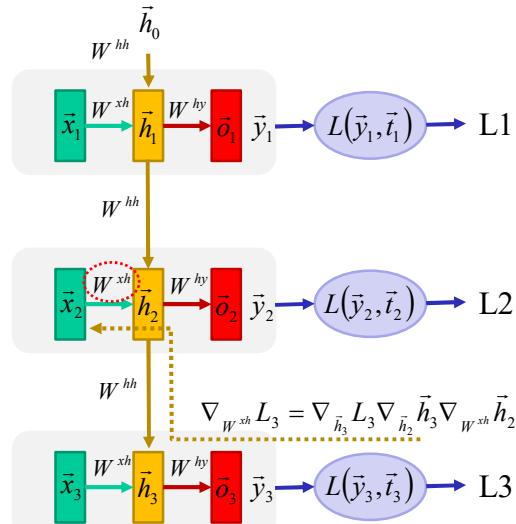
120

Réseau récurrent: gradient pour W^{xh}



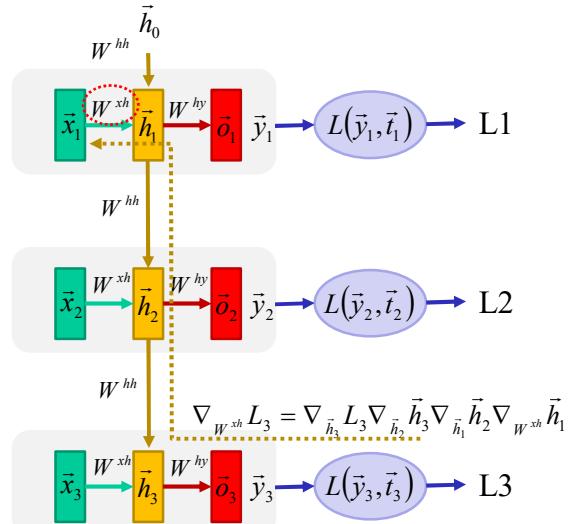
121

Réseau récurrent: gradient pour W^{xh}



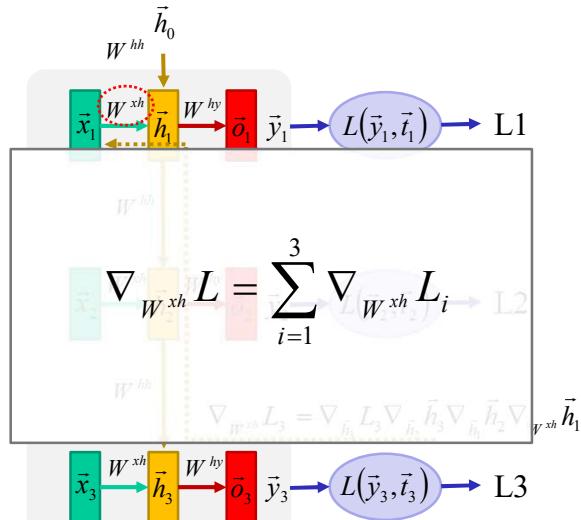
122

Réseau récurrent: gradient pour W^{xh}



123

Réseau récurrent: gradient pour W^{xh}



124

Réseau récurrent: calcul du gradient

Moins difficile qu'il n'y paraît.

```

44     # backward pass: compute gradients going backwards
45     dwxh, dwhh, dwhy = np.zeros_like(wxh), np.zeros_like(whh), np.zeros_like(why)
46     dbh, dy = np.zeros_like(bh), np.zeros_like(by)
47     dhnext = np.zeros_like(hs[0])
48     for t in reversed(range(len(inputs))): # backward pass
49         dy = np.copy(ps[t])
50         dy[targets[t]] -= 1 # backprop into y. see http://cs231n.github.io/neural-networks-case-study/#grad if confused here
51         dwhy += np.dot(dy, hs[t].T)
52         dy += dbh
53         dh = np.dot(why.T, dy) + dhnext # backprop into h
54         ddraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55         dbh += ddraw
56         dwhh += np.dot(draw, xs[t].T)
57         dwhh += np.dot(draw, hs[t-1].T)
58         dhnext = np.dot(whh.T, ddraw)
59     for dparam in [dwxh, dwhh, dwhy, dbh, dy]:
60         np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61     return loss, dwxh, dwhh, dwhy, dbh, dy, hs[len(inputs)-1]

```

Voir https://d2l.ai/chapter_recurrent-neural-networks/bptt.html pour plus d'informations

125

Les réseaux récurrents ont un inconvénient majeur:

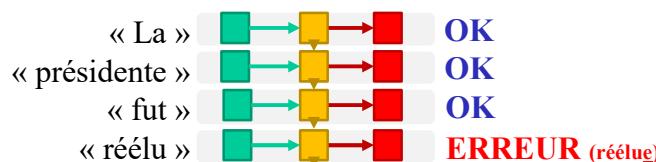
difficile à établir des
relations à longue distance

126

126

Exemples: analyse grammaticale

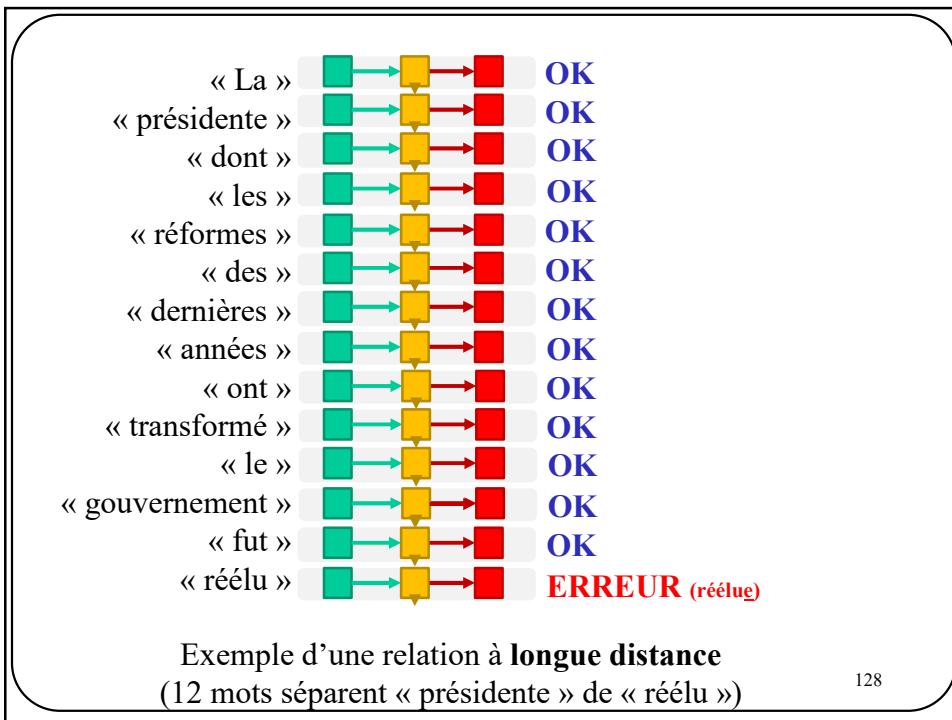
Entraîner un réseau à détecter des erreurs grammaticales



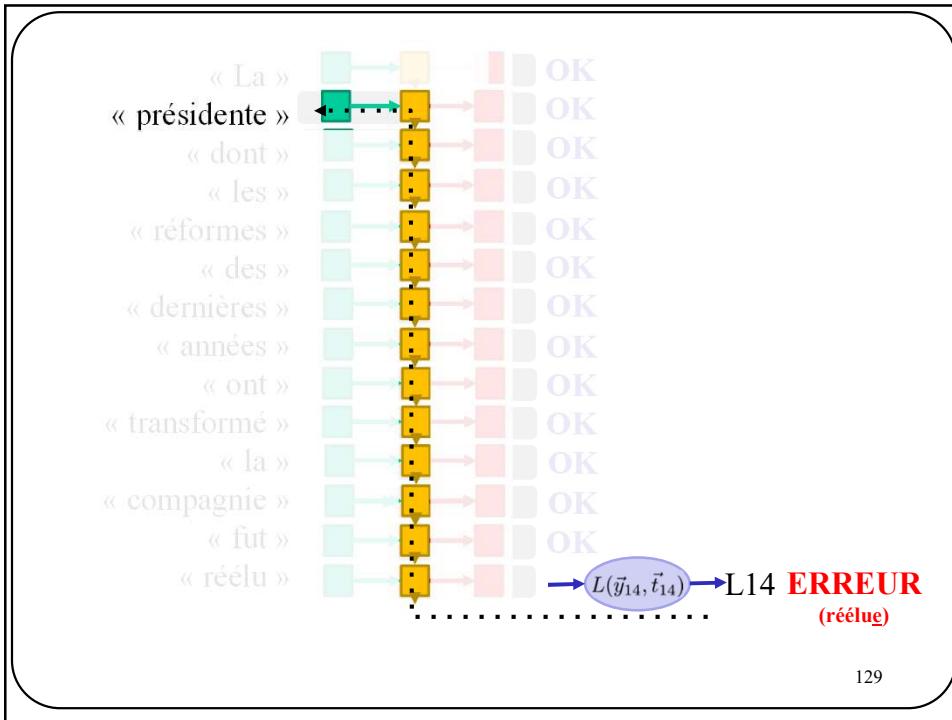
Exemple d'une relation à **courte distance**
(1 mot sépare « présidente » de « réélu »)

127

127

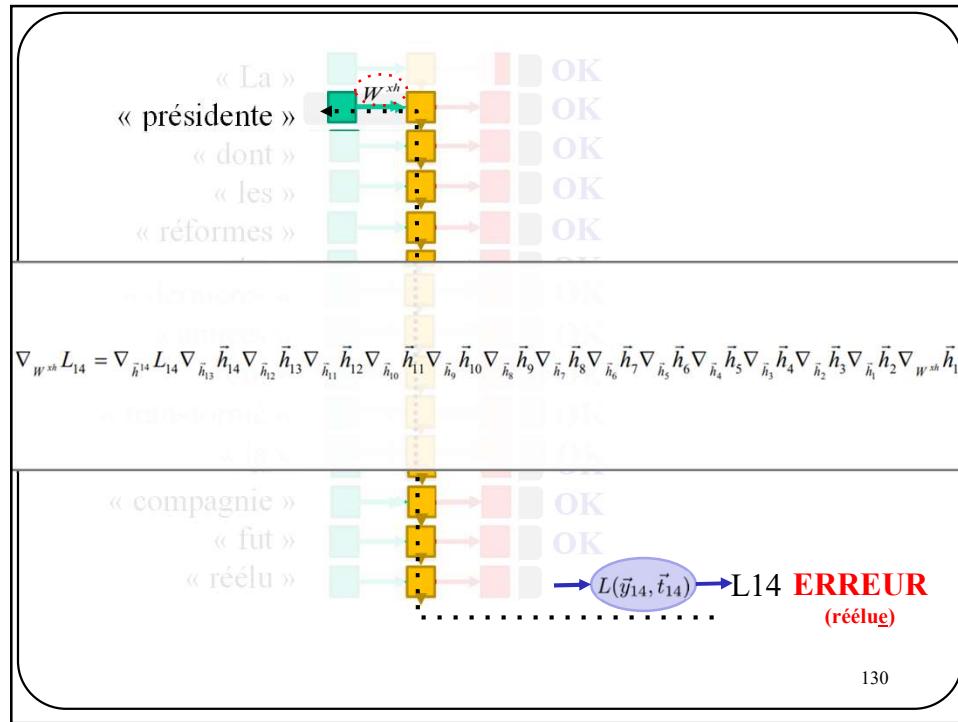


128

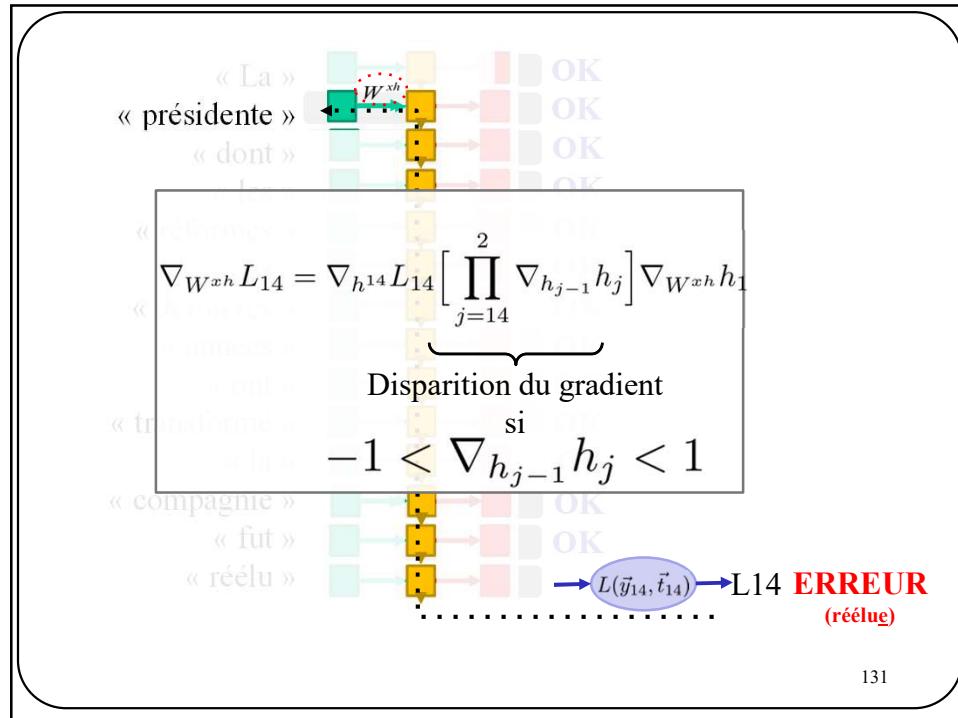


129

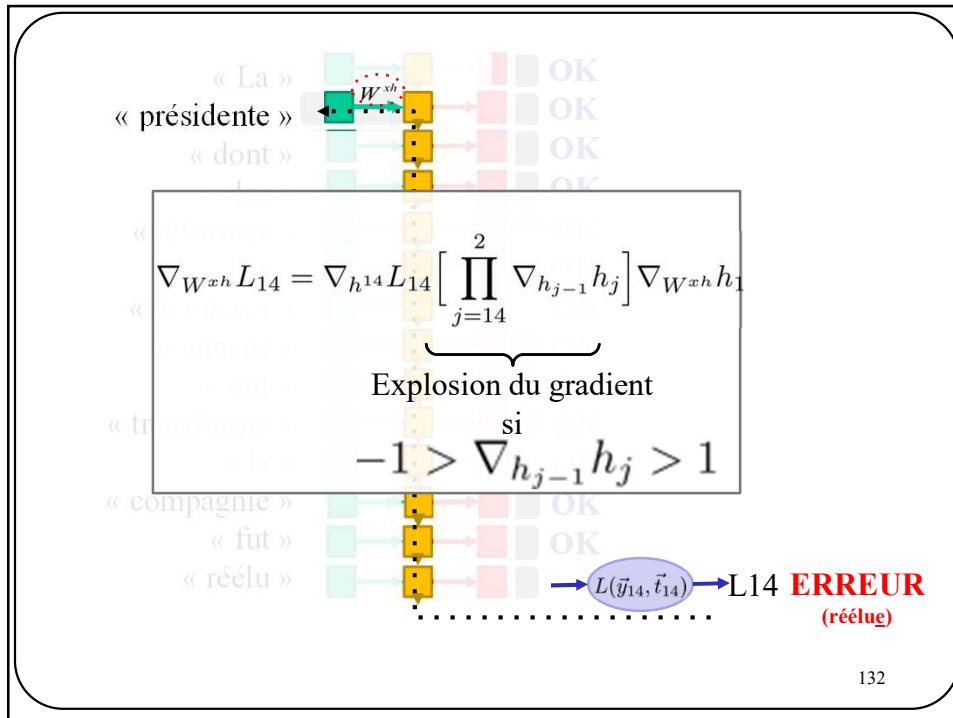
129



130



131



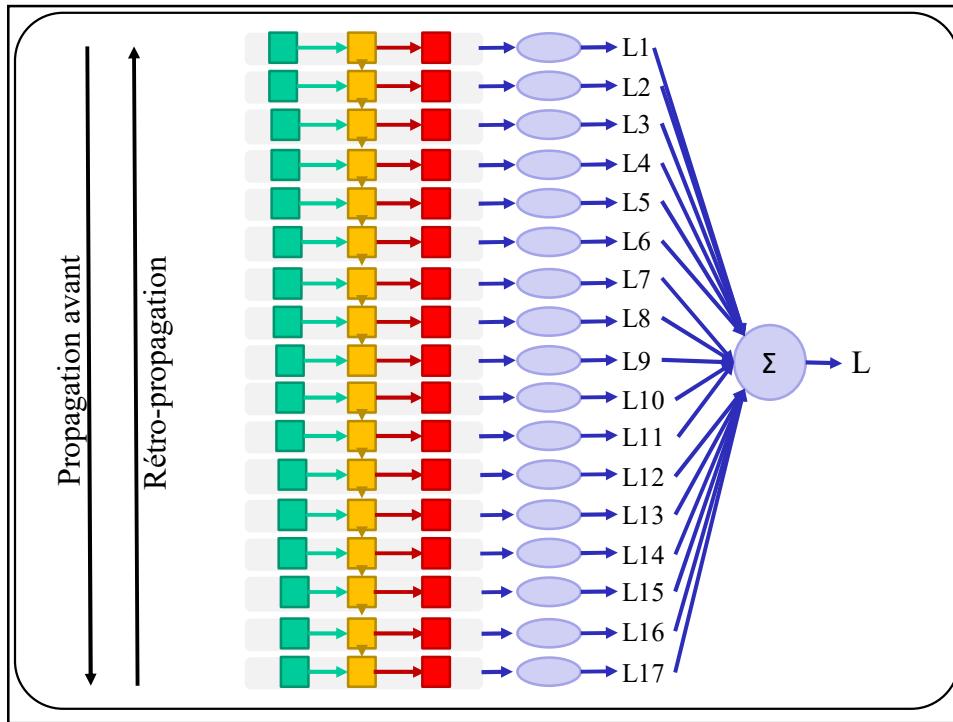
132

Problème connexe

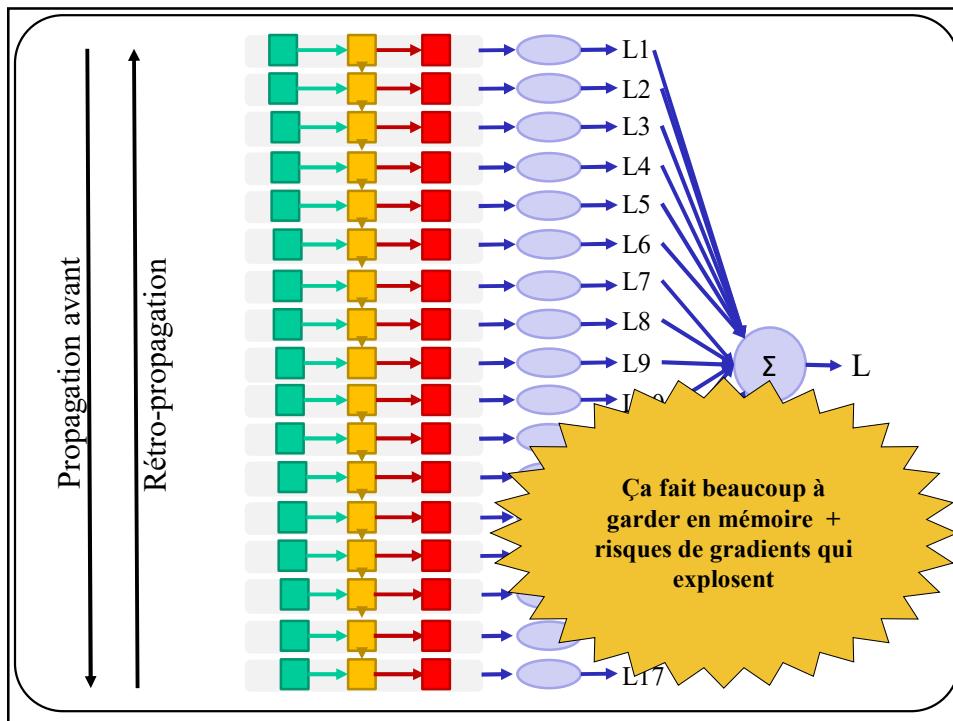
Gestion de la mémoire

133

133



134

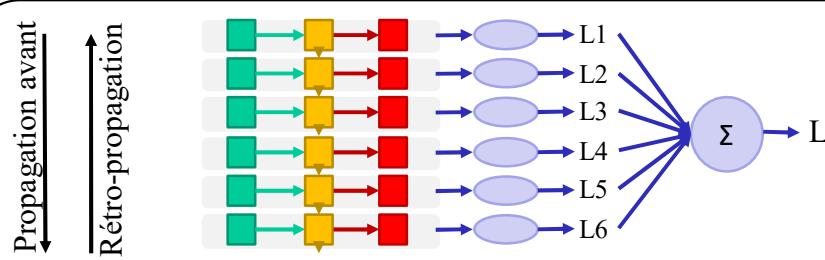


135

Solution pour la gestion de la mémoire **Fenêtres coulissantes**

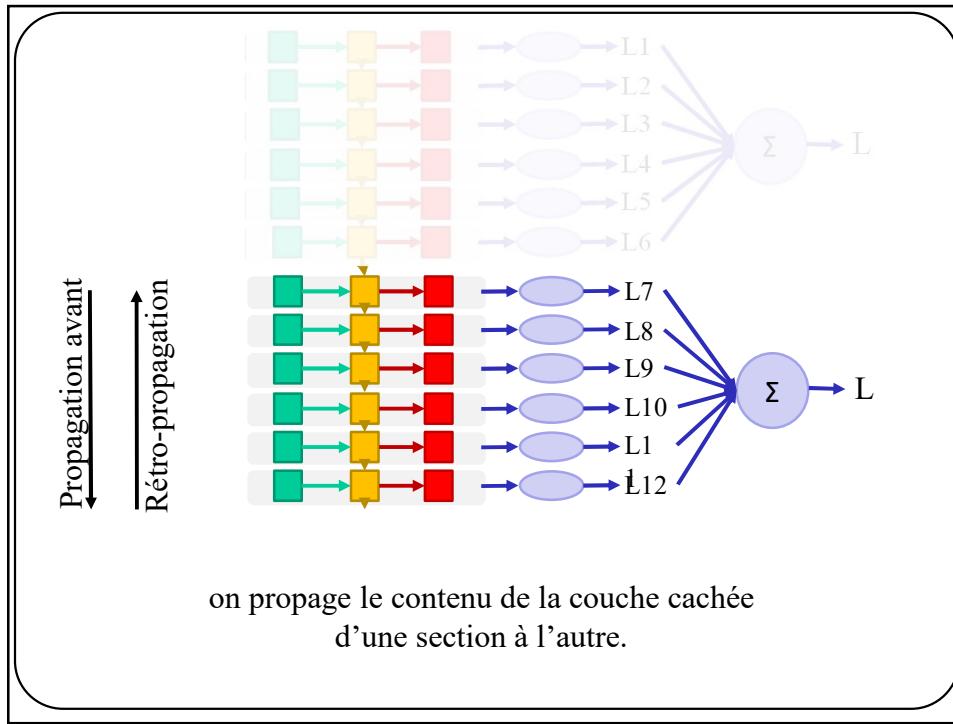
136

136

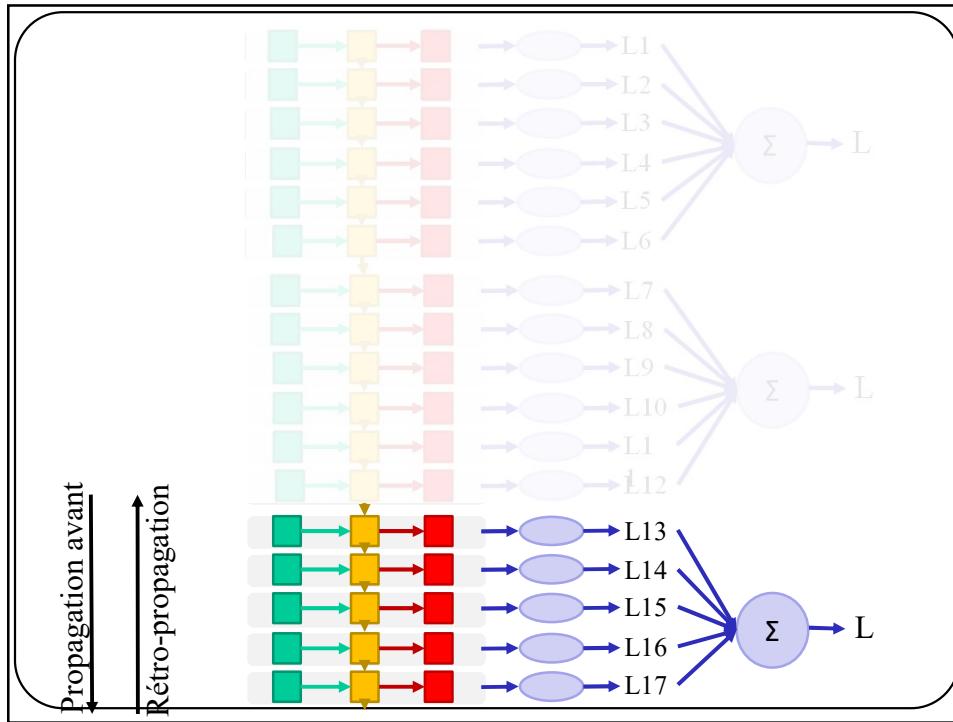


Lorsque les séquences sont trop longues
On entraîne le réseau par fenêtres coulissantes

137



138



139

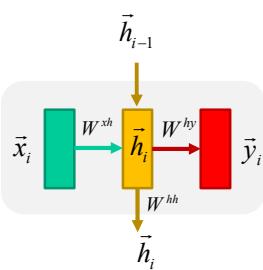
Solution à la disparition du gradient:

Gated Recurrent Unit : GRU
Long-Short Term Memory : LSTM

140

140

Illustration + formulation d'un RNN

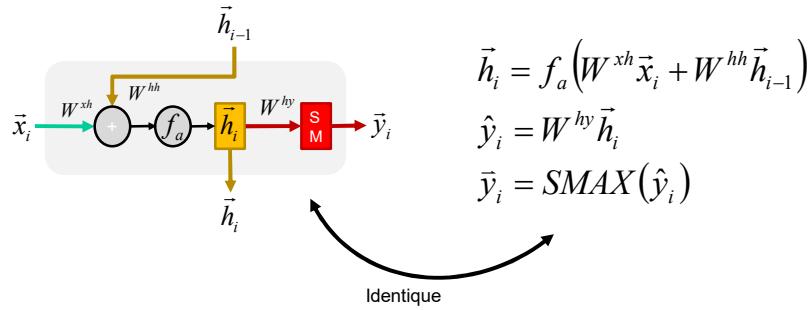


$$\begin{aligned}\vec{h}_i &= f_a(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \hat{y}_i &= W^{hy}\vec{h}_i \\ \bar{y}_i &= SMAX(\hat{y}_i)\end{aligned}$$

141

141

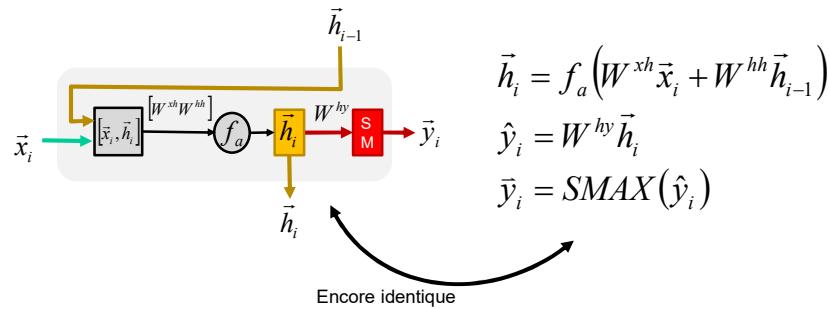
Autre illustration du même RNN



142

142

Autre illustration du même RNN



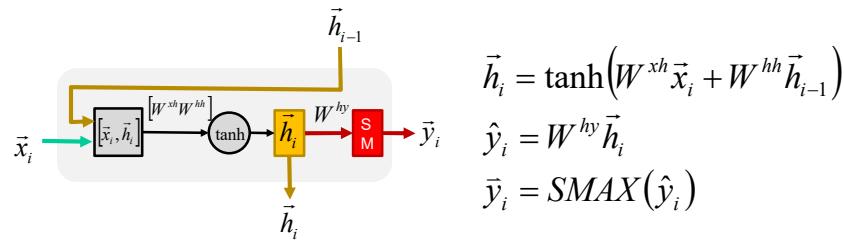
143

143

GRU (Gated Recurrent Unit)

Modif 1

$$f_a = \tanh$$



144

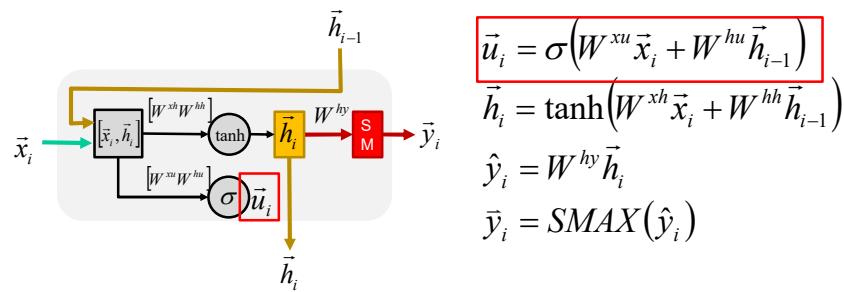
144

GRU (Gated Recurrent Unit)

Modif 2

Update gate

$$\sigma = \text{sigmoid}$$



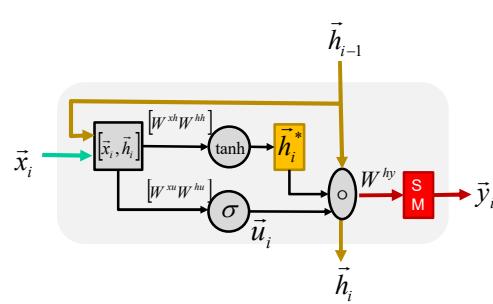
145

145

GRU (Gated Recurrent Unit)

Modif 2

Update gate



$$\begin{aligned}\vec{u}_i &= \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{h}_i^* &= \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \vec{h}_i &= \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \circ \vec{h}_{i-1} \\ \hat{y}_i &= W^{hy}\vec{h}_i \\ \bar{y}_i &= \text{SMAX}(\hat{y}_i)\end{aligned}$$

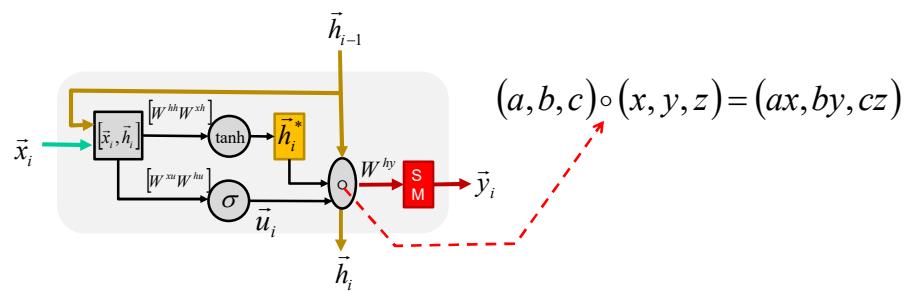
146

146

GRU (Gated Recurrent Unit)

Modif 2

« element-wise product »
ou
« Hadamard product »



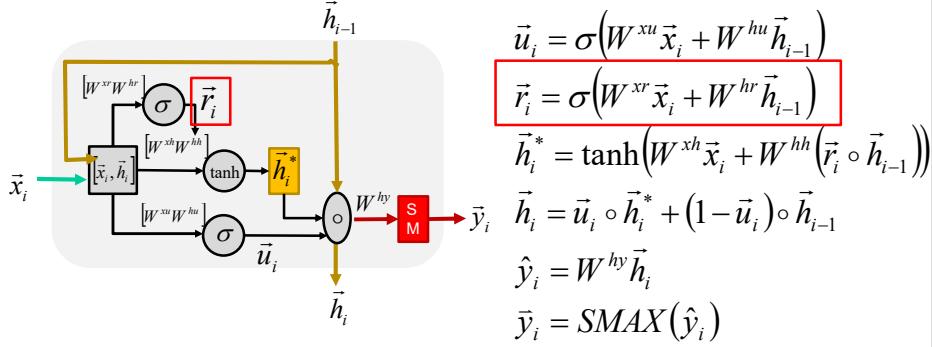
147

147

GRU (Gated Recurrent Unit)

Modif 3

Reset gate

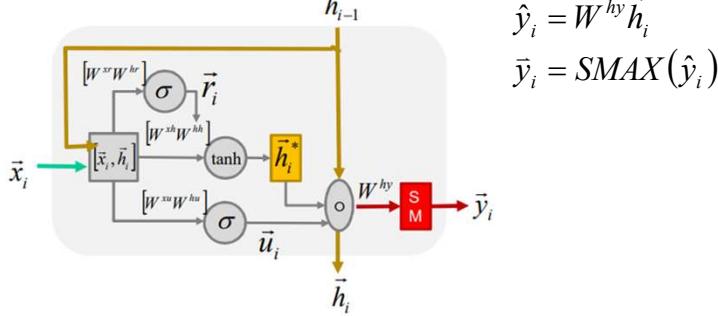


148

148

Comprendre les gates

$$SI \quad \begin{matrix} u_i \\ r_i \end{matrix} = 1 \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(W^{xu} \vec{x}_i + W^{hu} \vec{h}_{i-1}) \\ \vec{r}_i = \sigma(W^{xr} \vec{x}_i + W^{hr} \vec{h}_{i-1}) \\ \vec{h}_i^* = \tanh(W^{xh} \vec{x}_i + W^{hh} (\vec{r}_i \circ \vec{h}_{i-1})) \\ \vec{h}_i = \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \circ \vec{h}_{i-1} \\ \hat{y}_i = W^{hy} \vec{h}_i \\ \bar{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$

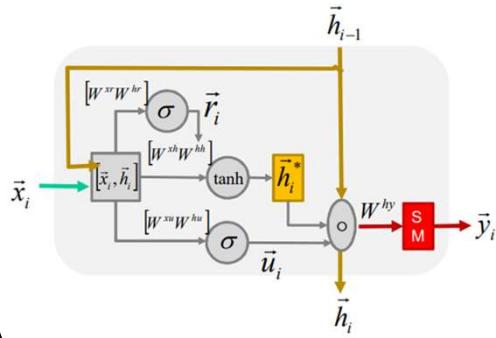


149

149

Comprendre les *gates*

$$\left. \begin{array}{l} SI \\ u_i = 1 \\ r_i = 1 \end{array} \right\} \quad \begin{aligned} \vec{u}_i &= \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{r}_i &= \sigma(W^{xr}\vec{x}_i + W^{hr}\vec{h}_{i-1}) \\ \vec{h}_i^* &= \tanh(W^{xh}\vec{x}_i + W^{hh}(\vec{r}_i \circ \vec{h}_{i-1})) \\ \vec{h}_i &= \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \times \vec{h}_{i-1} \\ \hat{y}_i &= W^{hy}\vec{h}_i \\ \bar{y}_i &= SMAX(\hat{y}_i) \end{aligned}$$

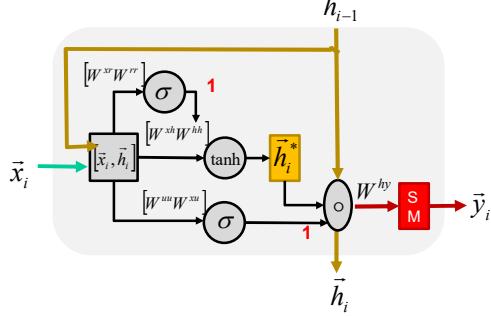


150

150

Comprendre les *gates*

$$\left. \begin{array}{l} SI \\ u_i = 1 \\ r_i = 1 \end{array} \right\} \quad \begin{aligned} \vec{u}_i &= \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{r}_i &= \sigma(W^{xr}\vec{x}_i + W^{hr}\vec{h}_{i-1}) \\ \vec{h}_i^* &= \tanh(W^{xh}\vec{x}_i + W^{hh}(\vec{r}_i \circ \vec{h}_{i-1})) \\ \vec{h}_i &= \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \times \vec{h}_{i-1} \\ \hat{y}_i &= W^{hy}\vec{h}_i \\ \bar{y}_i &= SMAX(\hat{y}_i) \end{aligned}$$

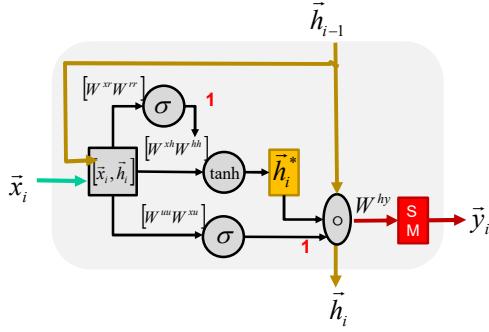


151

151

Comprendre les *gates*

$$SI \quad \begin{cases} \vec{u}_i = 1 \\ \vec{r}_i = 1 \end{cases} \quad \left\{ \begin{array}{l} \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \vec{h}_i = \vec{h}_i^* \\ \hat{y}_i = W^{hy}\vec{h}_i \\ \bar{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$

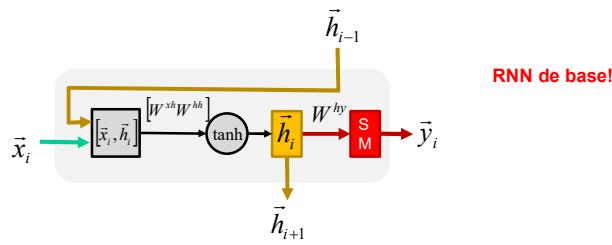


152

152

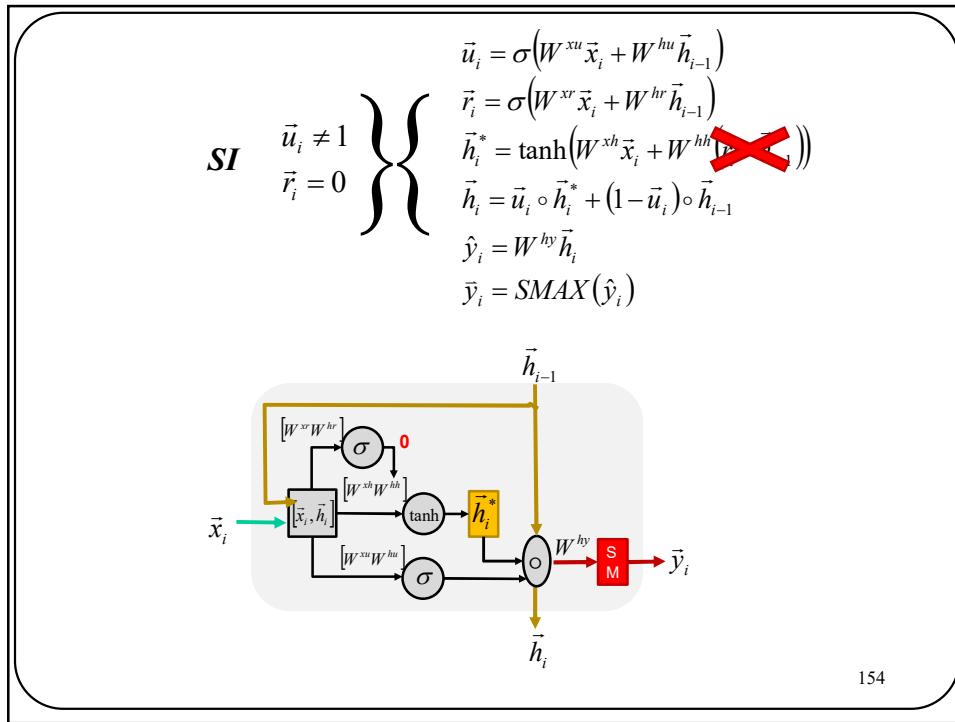
Comprendre les *gates*

$$SI \quad \begin{cases} \vec{u}_i = 1 \\ \vec{r}_i = 1 \end{cases} \quad \left\{ \begin{array}{l} \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \vec{h}_i = \vec{h}_i^* \\ \hat{y}_i = W^{hy}\vec{h}_i \\ \bar{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$

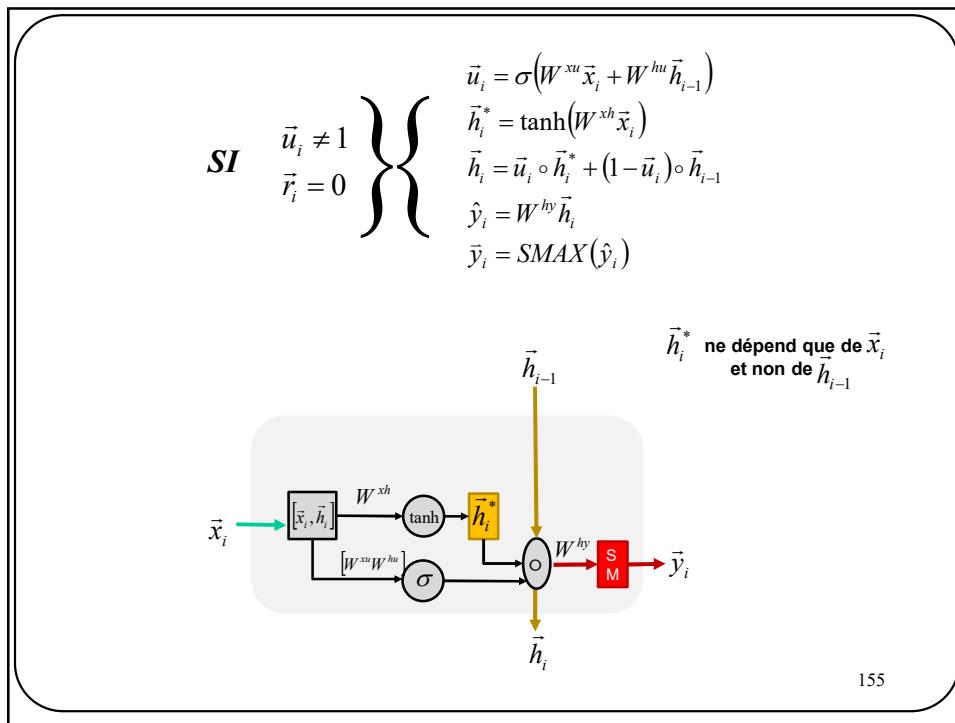


153

153



154



155

$$SI \quad \begin{cases} \vec{u}_i = 0 \\ \vec{r}_i = 0 \end{cases} \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(W^{ur}\vec{x}_i + W^{uh}\vec{h}_{i-1}) \\ \vec{r}_i = \sigma(W^{rr}\vec{x}_i + W^{rh}\vec{h}_{i-1}) \\ \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \vec{h}_i = \vec{r}_i \odot \vec{h}_i^* + (1 - \vec{u}_i) \odot \vec{h}_{i-1} \\ \vec{y}_i = W^{hy}\vec{h}_i \\ \vec{y}_i = SMAX(\vec{y}_i) \end{array} \right.$$

156

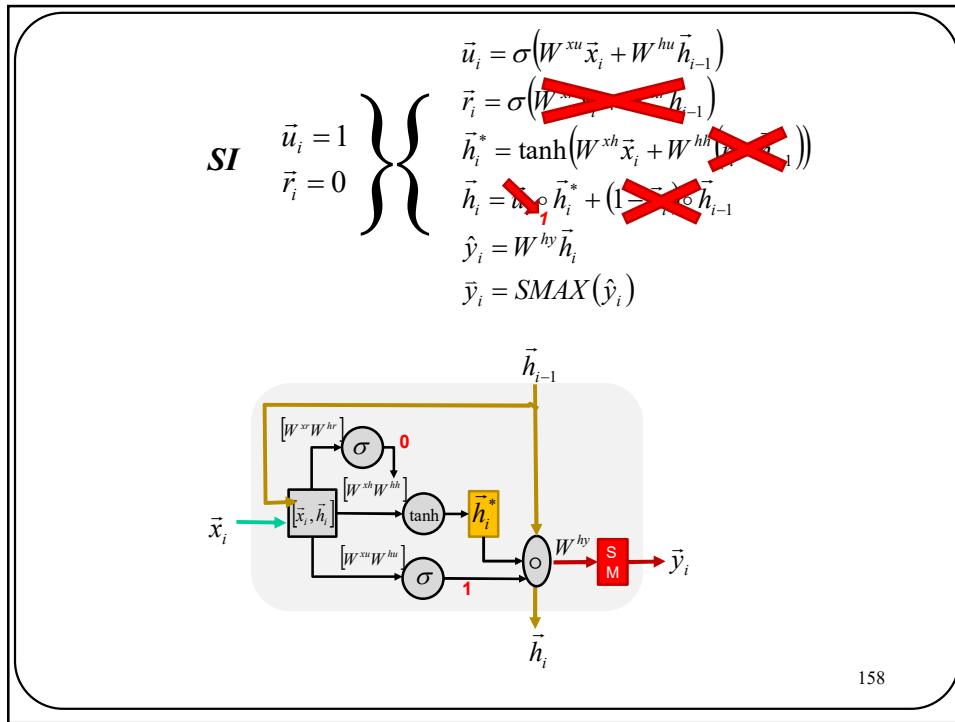
156

$$SI \quad \begin{cases} \vec{u}_i = 0 \\ \vec{r}_i = 0 \end{cases} \quad \left\{ \begin{array}{l} \vec{h}_i = \vec{h}_{i-1} \\ \hat{y}_i = W^{cy}\vec{h}_i \\ \vec{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$

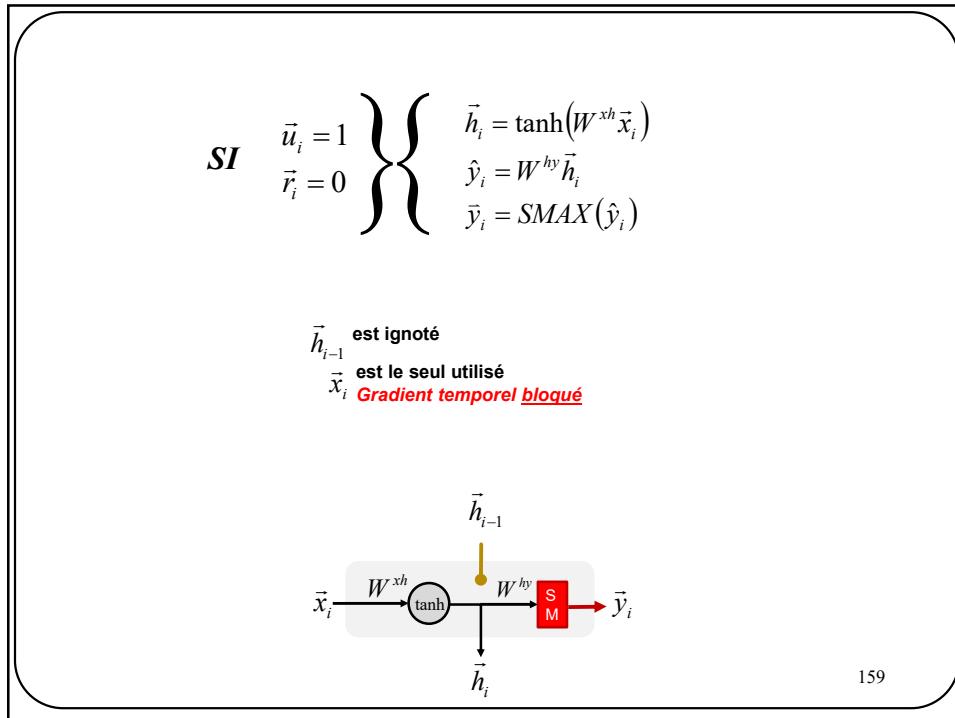
\vec{h}_{i-1} est recopié
 \vec{x}_i est ignoré
 $\text{Aucune disparition de gradient}$

157

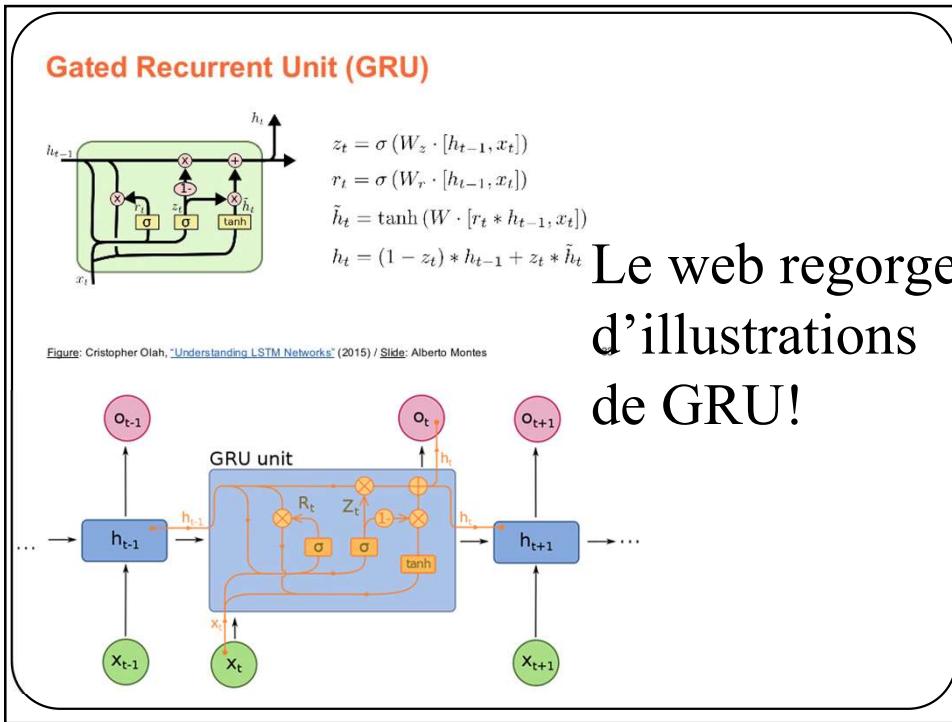
157



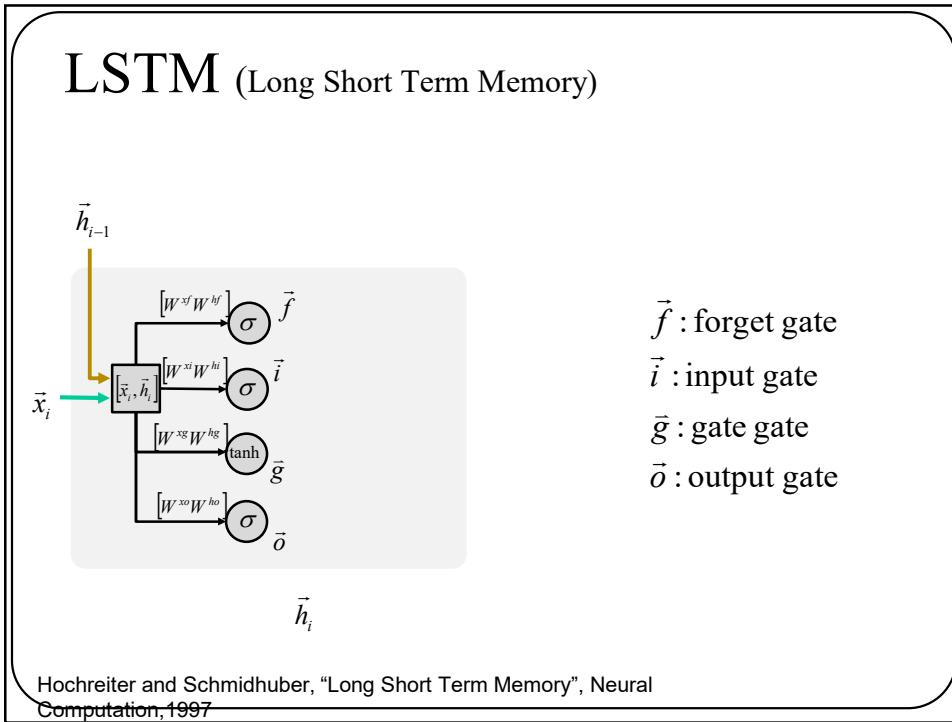
158



159

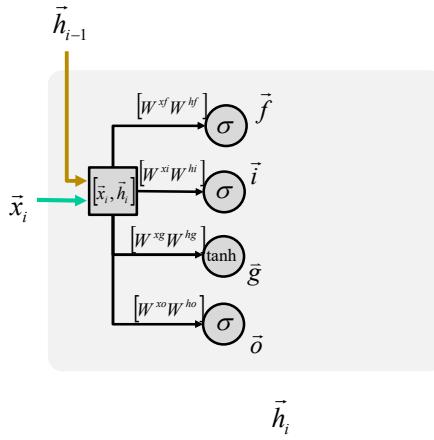


160



161

LSTM (Long Short Term Memory)



\vec{f} : forget gate

\vec{i} : input gate

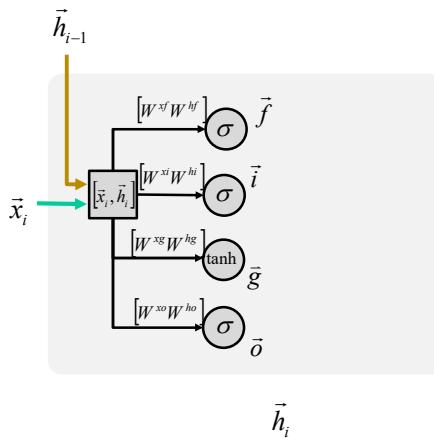
\vec{g} : gate gate

Modèle récurrent
le plus utilisé.
À bien comprendre !

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation, 1997

162

LSTM (Long Short Term Memory)



$$\vec{f} = \sigma(W^{xf} \vec{x}_i + W^{hf} \vec{h}_{i-1})$$

$$\vec{i} = \sigma(W^{xi} \vec{x}_i + W^{hi} \vec{h}_{i-1})$$

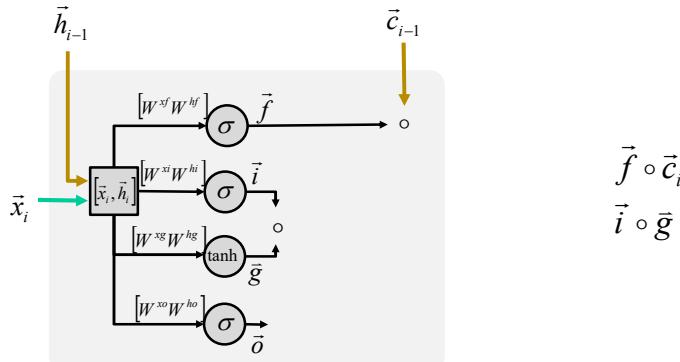
$$\vec{g} = \tanh(W^{xg} \vec{x}_i + W^{hg} \vec{h}_{i-1})$$

$$\vec{o} = \sigma(W^{xo} \vec{x}_i + W^{ho} \vec{h}_{i-1})$$

163

163

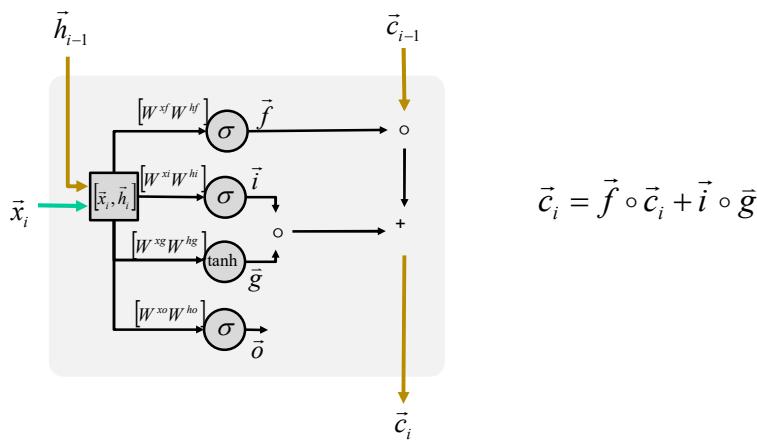
LSTM (Long Short Term Memory)



164

164

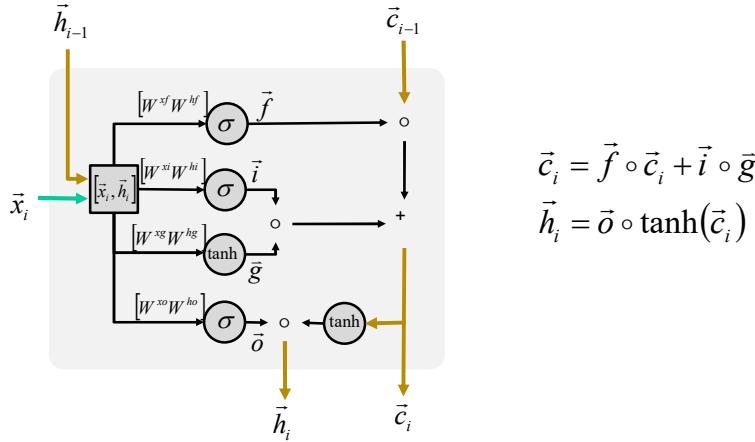
LSTM (Long Short Term Memory)



165

165

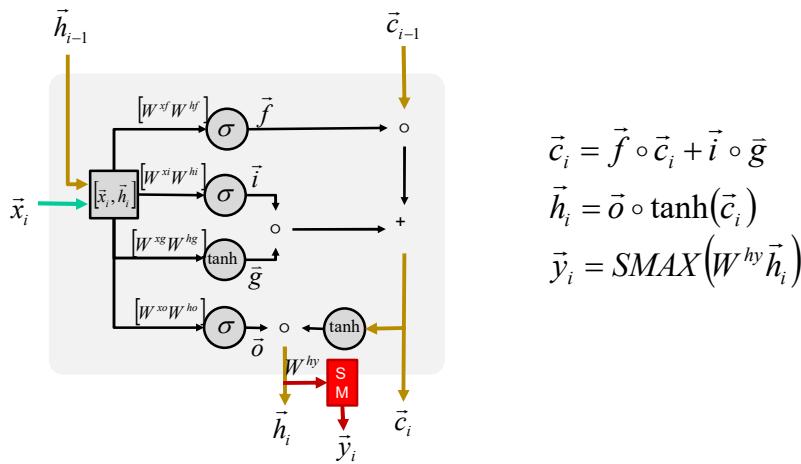
LSTM (Long Short Term Memory)



166

166

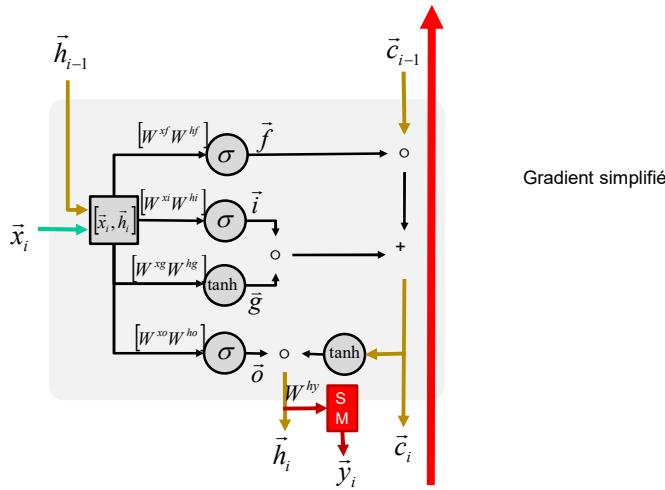
LSTM (Long Short Term Memory)



167

167

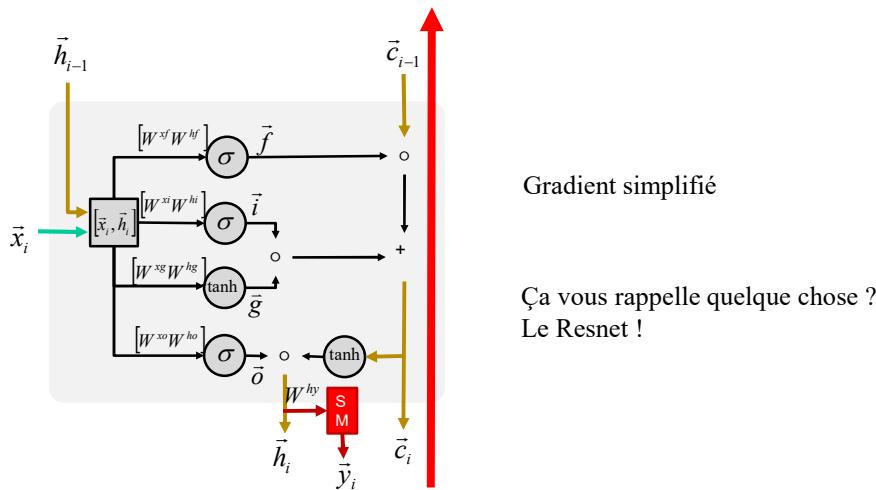
LSTM (Long Short Term Memory)



168

168

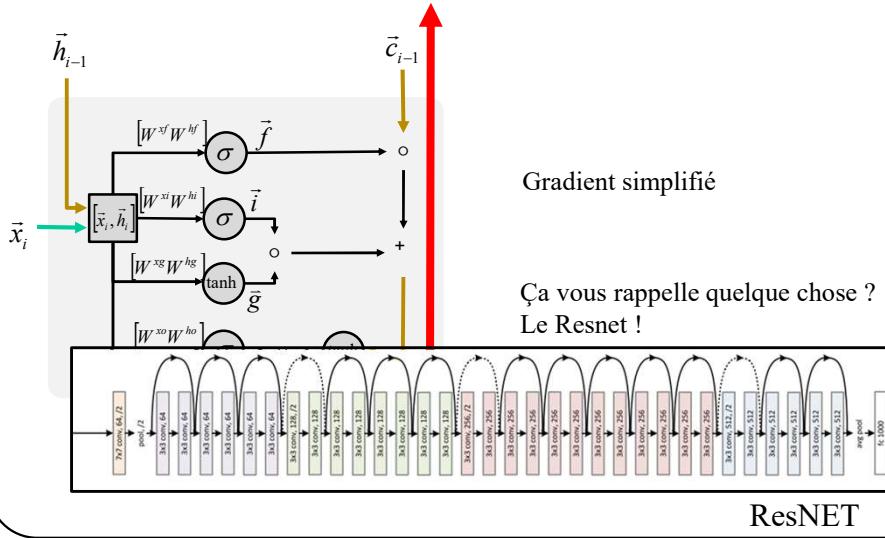
LSTM (Long Short Term Memory)



169

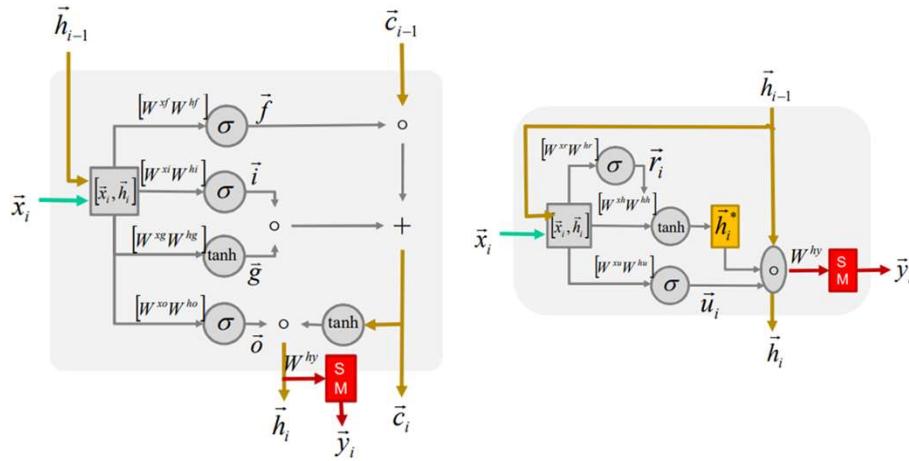
169

LSTM (Long Short Term Memory)



170

LSTM et GRU

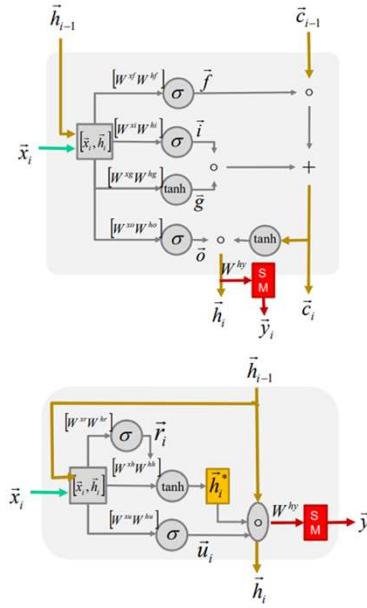


171

171

LSTM et GRU

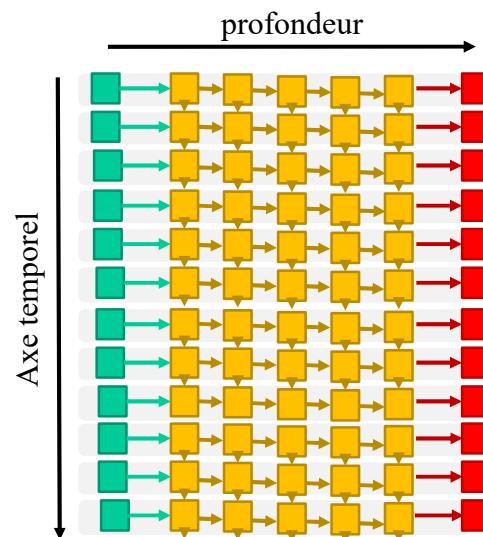
- Servent à protéger le gradient
- Conçus empiriquement
- GRU légèrement plus simple
- Les “gates” ne servent qu’à bloquer ou permettre à l’information (données ou temporelle) de passer



172

172

RNN multi-couches



173

Modèles d'attention

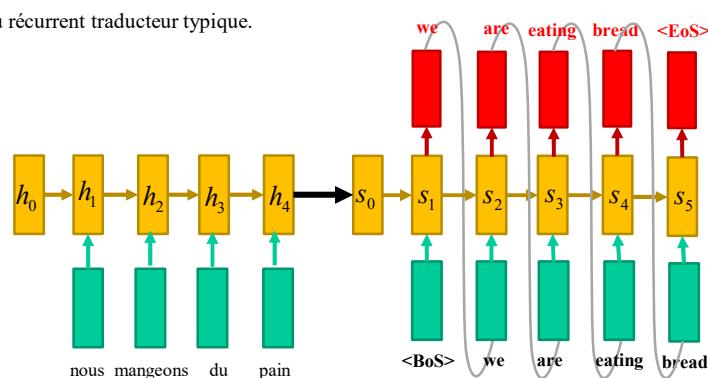
174

Seq2Seq:

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

'nous mangeons du pain' -> 'we are eating bread'

Réseau récurrent traducteur typique.



175

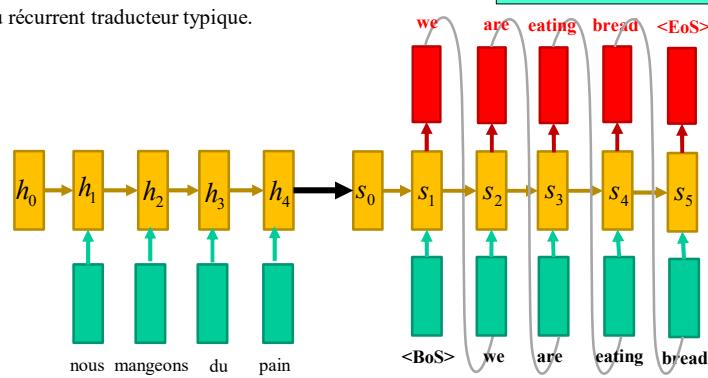
Seq2Seq:

'nous mangeons du pain' -> 'we are eating bread'

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Problème avec cette formulation ?

Réseau récurrent traducteur typique.



176

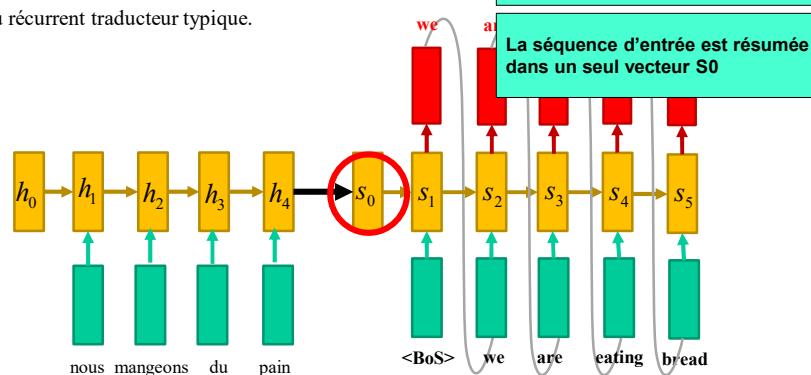
Seq2Seq:

'nous mangeons du pain' -> 'we are eating bread'

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Problème avec cette formulation ?

Réseau récurrent traducteur typique.

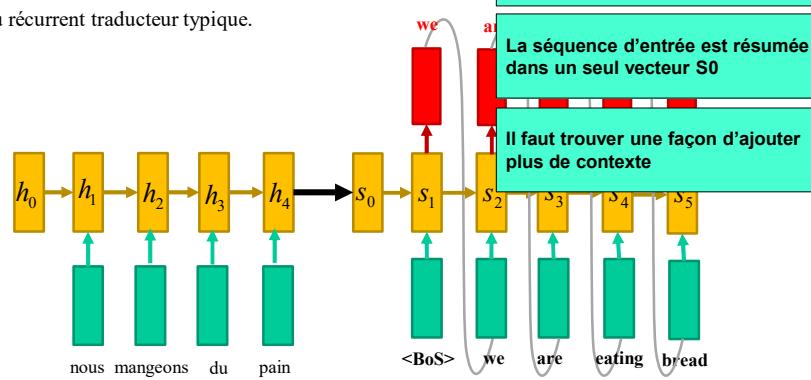


177

Seq2Seq:

'nous mangeons du pain' -> 'we are eating bread'

Réseau récurrent traducteur typique.



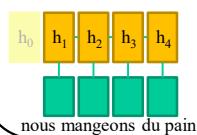
178

Seq2Seq avec attention

(version simplifiée)

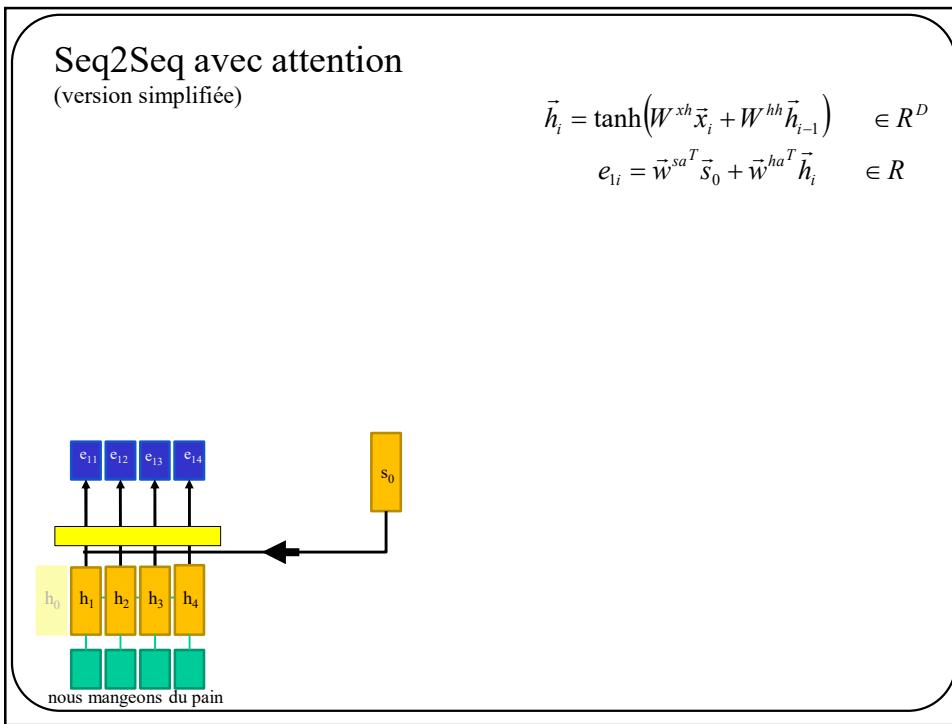
Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473.

$$\vec{h}_i = \tanh(W^{xh} \vec{x}_i + W^{hh} \vec{h}_{i-1}) \in R^D$$

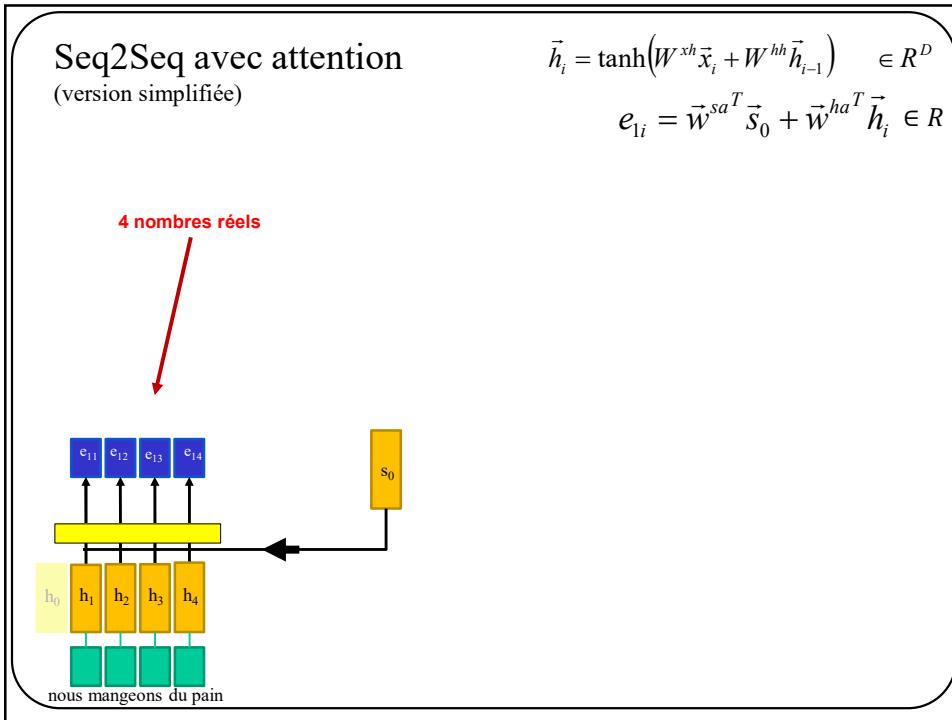


$s_0 = \text{init value}$

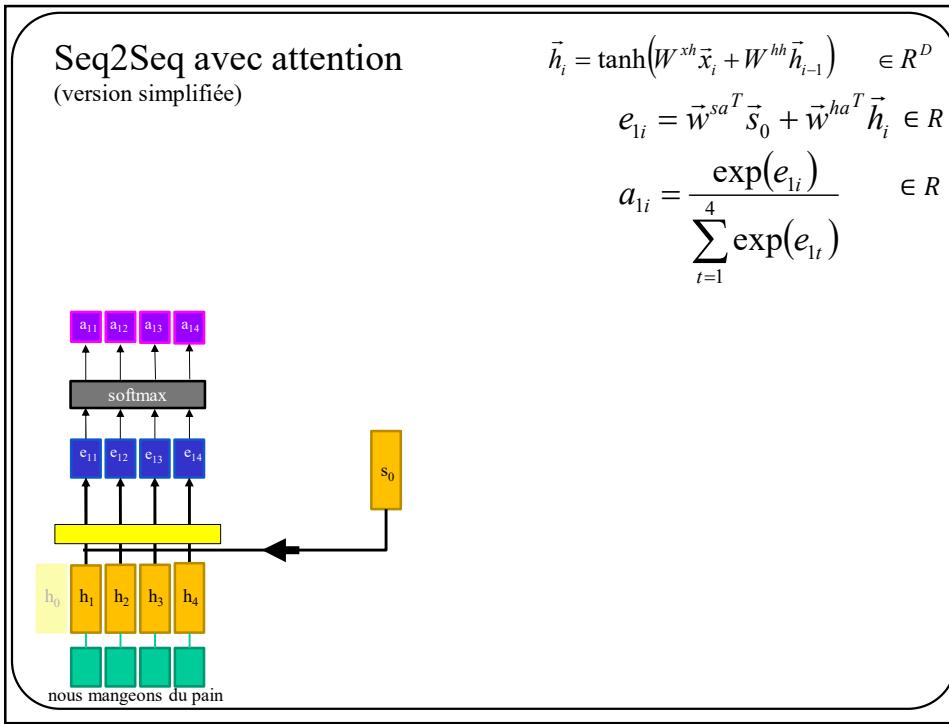
179



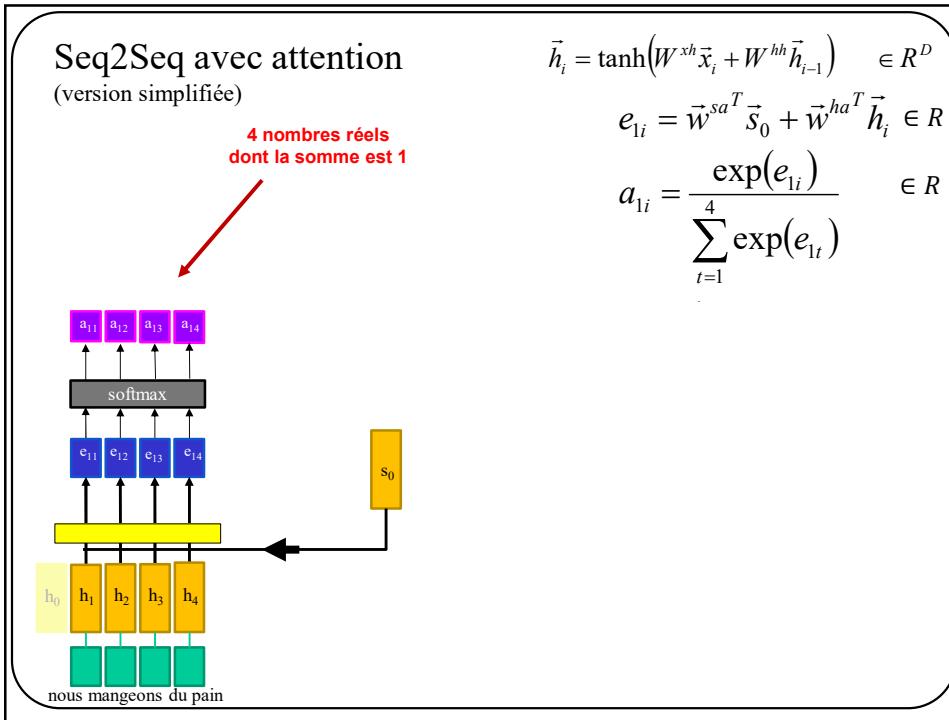
180



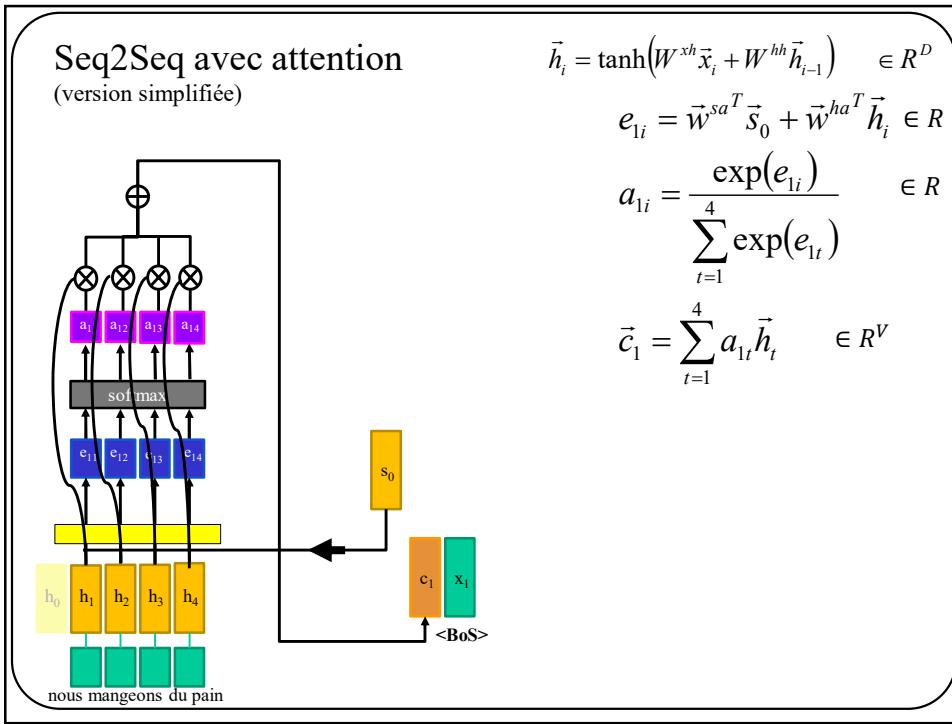
181



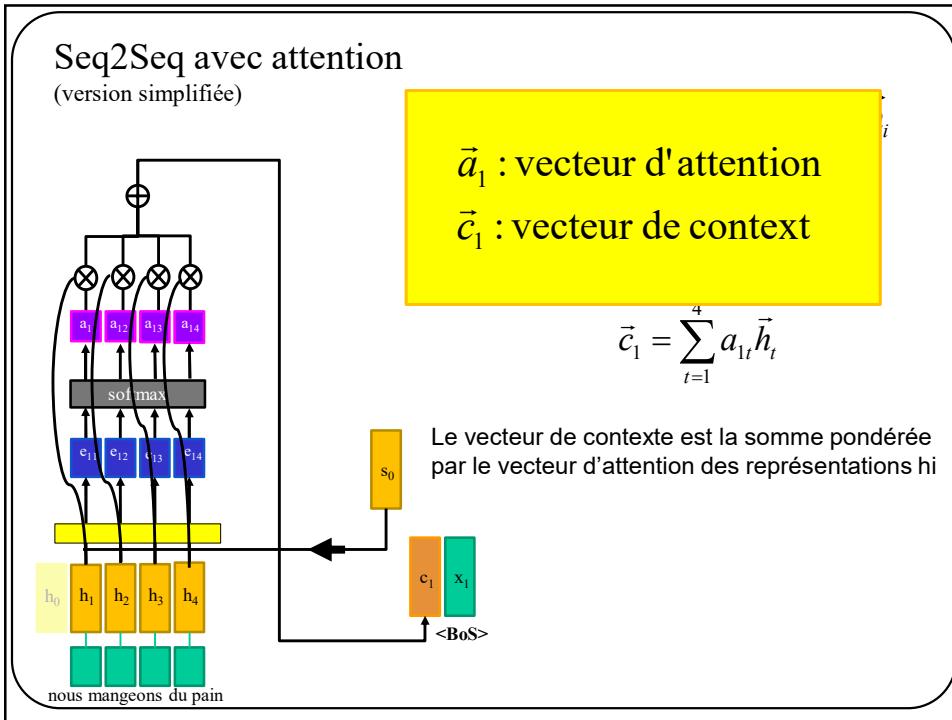
182



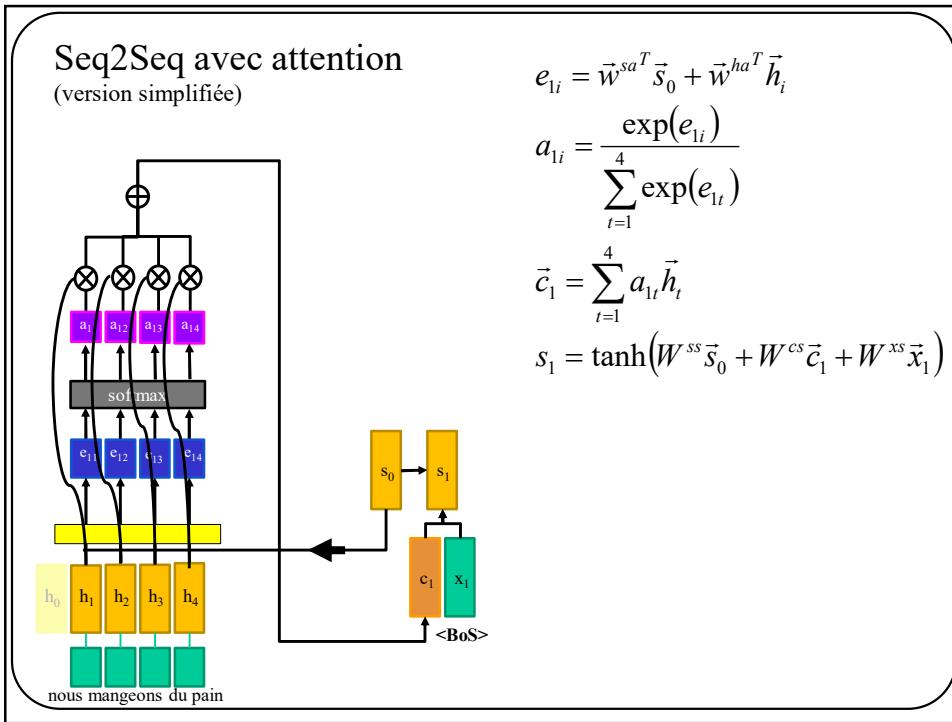
183



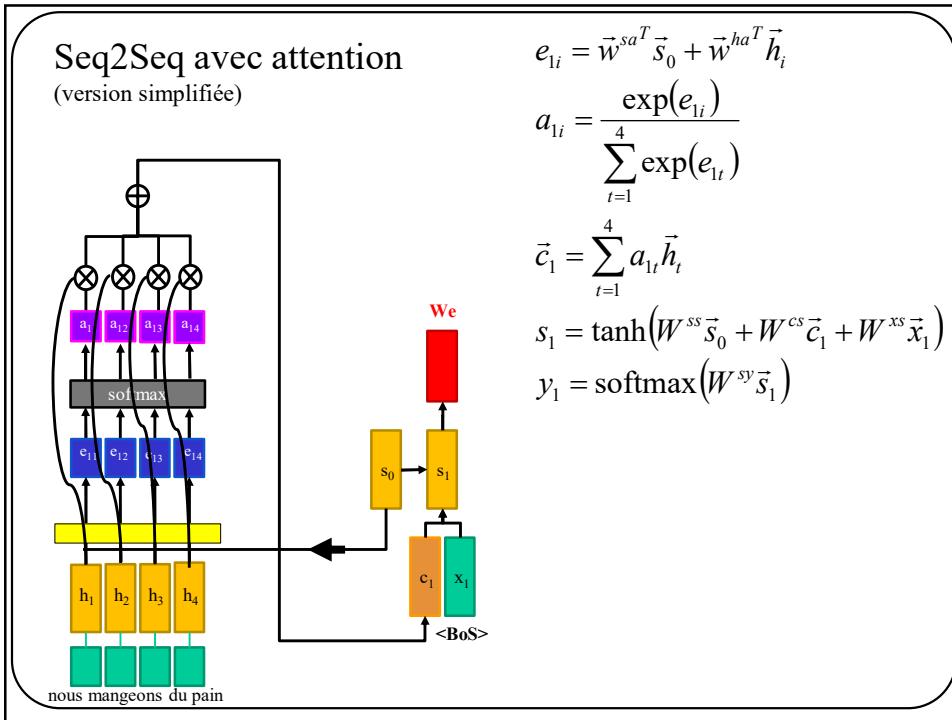
184



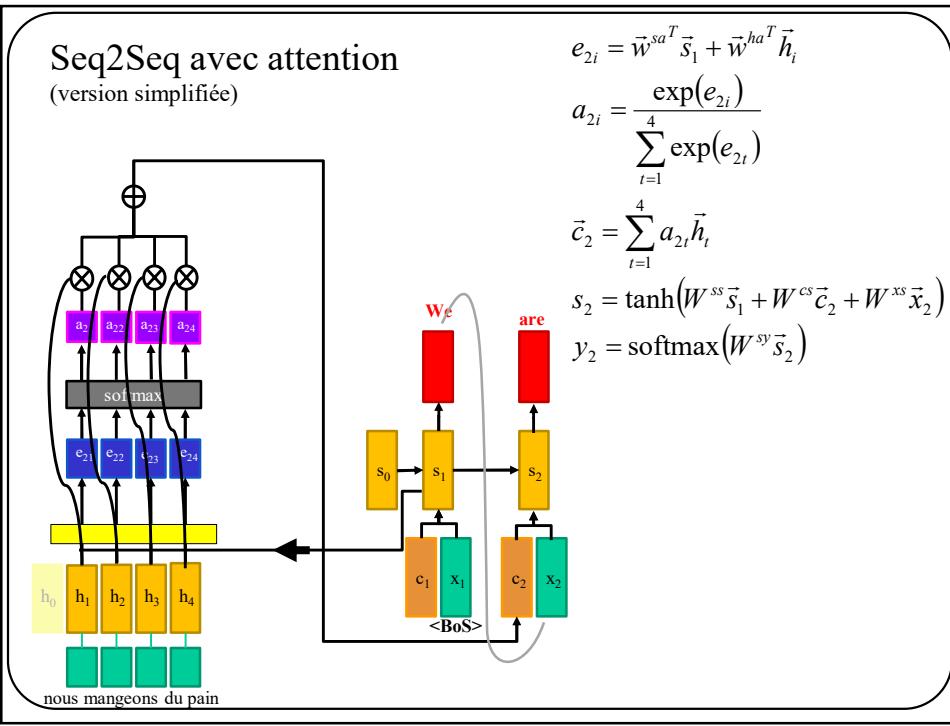
185



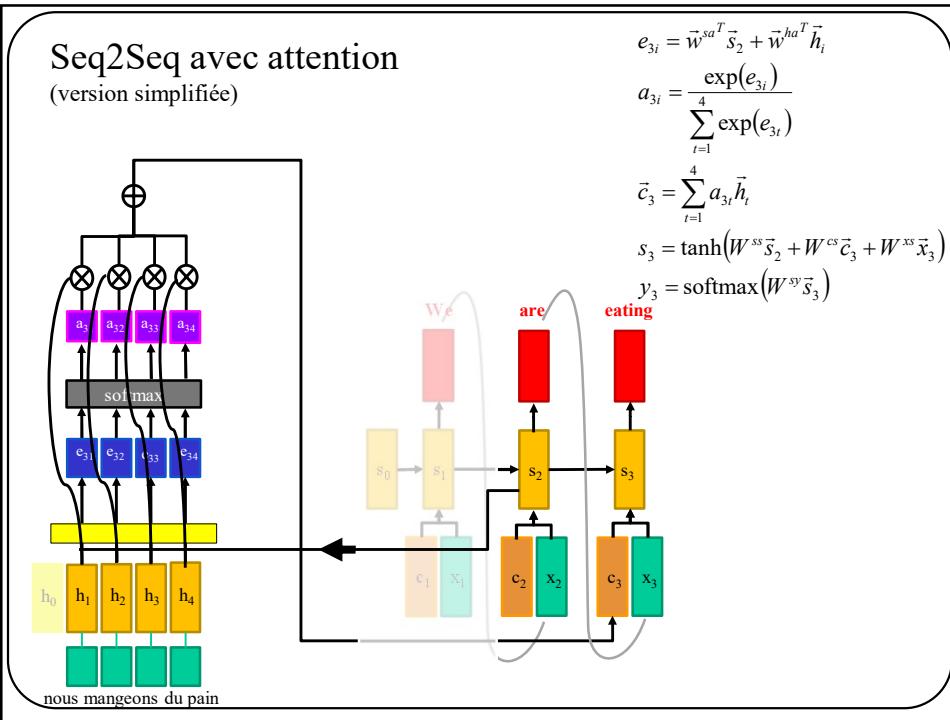
186



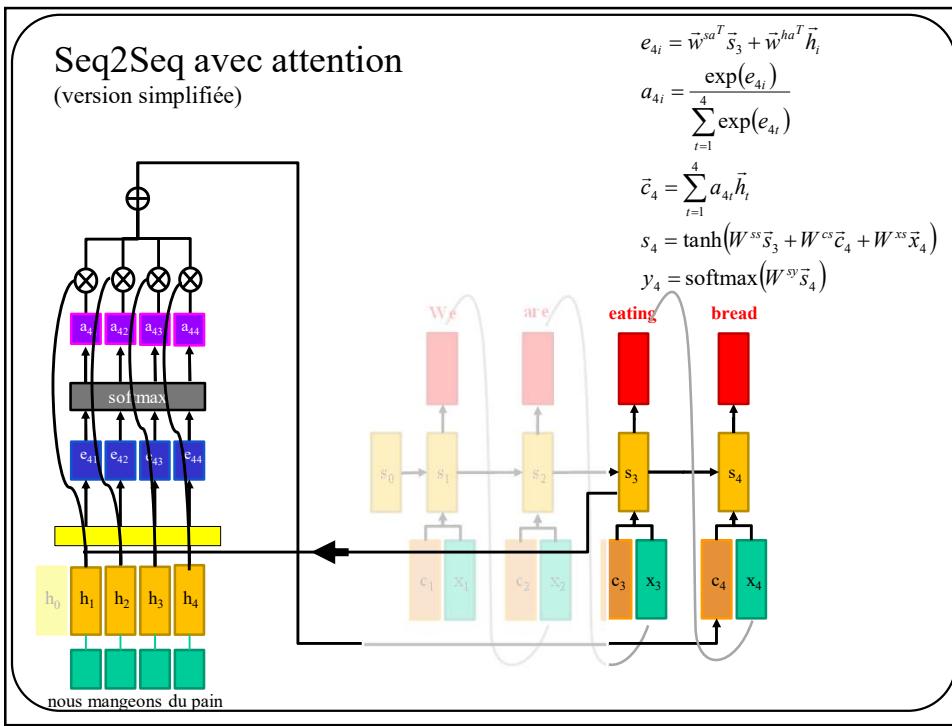
187



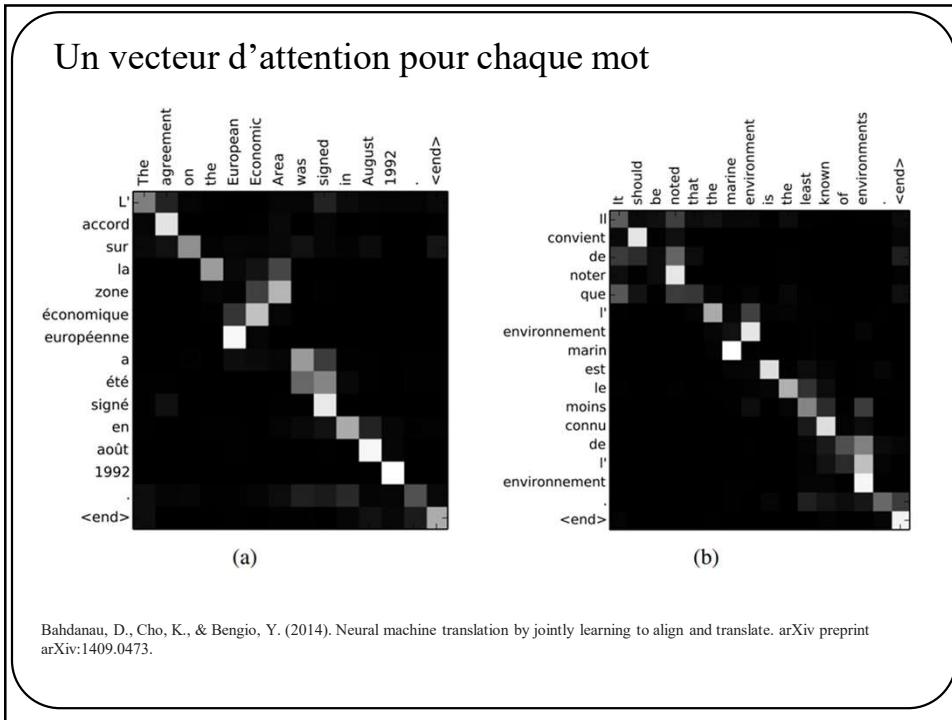
188



189



190



191

Seq2Seq avec attention

L'accord sur la zone économique européenne a été signé en août 1992.

The agreement on the European Economic Area was signed in August 1992.

(a)

Il convient de noter que l'environnement marin est le moins connu de l'environnement.

It should be noted that the marine environment is the least known of environments.

Ajoute de l'interprétabilité au modèle !

Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning. arXiv:1409.0473.

192

L'auto-attention (*self attention*)

204

204

Revenons à la base : multiplication matricielle

Considérons les 4 matrices suivantes

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3x4}$$

$$W^q = \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in R^{3x3}$$

$$W^k = \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in R^{3x3}$$

$$W^V = \begin{pmatrix} W^V_{11} & W^V_{12} & W^V_{13} \\ W^V_{21} & W^V_{22} & W^V_{23} \end{pmatrix} \in R^{2x3}$$

205

205

Revenons à la base : multiplication matricielle

Leur multiplication donne:

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3x4}$$

$$W^q X = Q = \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} & Q_{14} \\ Q_{21} & Q_{22} & Q_{23} & Q_{24} \\ Q_{31} & Q_{32} & Q_{33} & Q_{34} \end{pmatrix} \in R^{3x4}$$

$$W^q = \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in R^{3x3}$$

$$W^K X = K = \begin{pmatrix} K^x_{11} & K^x_{12} & K^x_{13} & K^x_{14} \\ K^x_{21} & K^x_{22} & K^x_{23} & K^x_{24} \\ K^x_{31} & K^x_{32} & K^x_{33} & K^x_{34} \end{pmatrix} \in R^{3x4}$$

$$W^k = \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in R^{3x3}$$

$$W^V X = V = \begin{pmatrix} V^x_{11} & V^x_{12} & V^x_{13} & V^x_{14} \\ V^x_{21} & V^x_{22} & V^x_{23} & V^x_{24} \end{pmatrix} \in R^{2x4}$$

206

206

Auto attention

X est une matrice de données pour laquelle chaque colonne i correspond au jeton d'une mot \vec{x}_i

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3 \times 4}$$

Dans cet exemple, 4 mots en entrée donc 4 colonnes dans X

Les jetons peuvent être obtenus par **Word2Vec**

207

207

Auto attention

$$W^q = \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in R^{3 \times 3}$$

$$W^k = \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in R^{3 \times 3}$$

$$W^v = \begin{pmatrix} W^v_{11} & W^v_{12} & W^v_{13} \\ W^v_{21} & W^v_{22} & W^v_{23} \end{pmatrix} \in R^{2 \times 3}$$

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3 \times 4}$$

W : Matrices de paramètres appris par **rétropropagation**

208

208

Auto attention

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3 \times 4}$$

$$W^q = \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in R^{3 \times 3}$$

$$W^k = \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in R^{3 \times 3}$$

$$W^v = \begin{pmatrix} W^v_{11} & W^v_{12} & W^v_{13} \\ W^v_{21} & W^v_{22} & W^v_{23} \end{pmatrix} \in R^{2 \times 3}$$

Matrices de paramètres appris par rétropropagation

Pour ces 3 matrices, le nombre de colonnes (3) doit être égale au nombre de lignes dans X (3)

209

209

Auto attention

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3 \times 4}$$

$$W^q = \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in R^{3 \times 3}$$

$$W^k = \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in R^{3 \times 3}$$

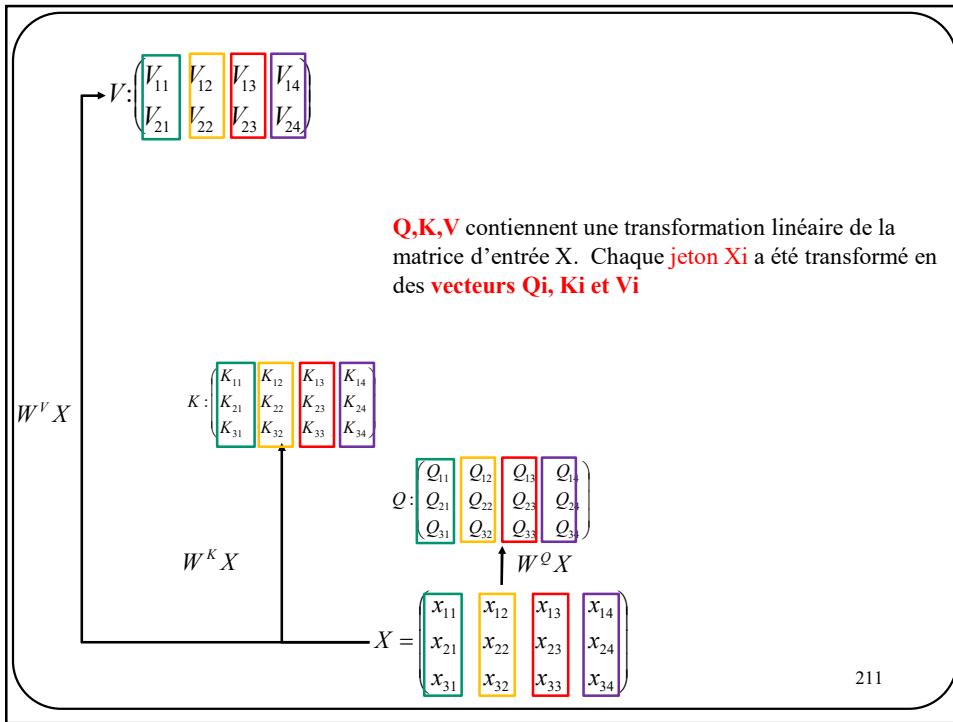
$$W^v = \begin{pmatrix} W^v_{11} & W^v_{12} & W^v_{13} \\ W^v_{21} & W^v_{22} & W^v_{23} \end{pmatrix} \in R^{2 \times 3}$$

Matrices de paramètres appris par rétropropagation

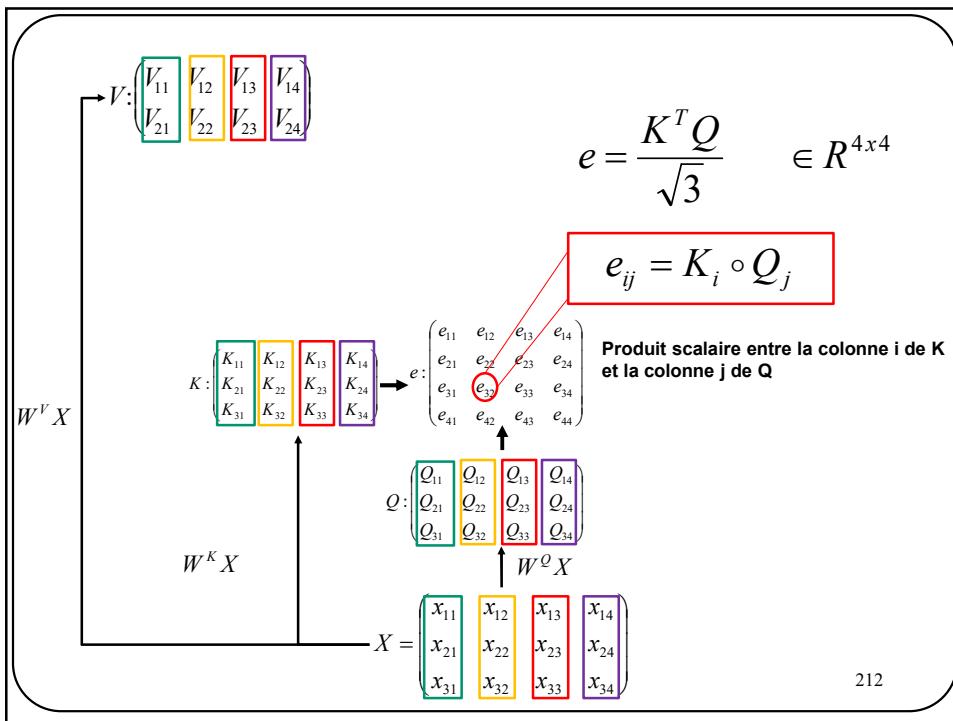
Pour ces 3 matrices, le nombre de ligne (3,3,2) est arbitraire

210

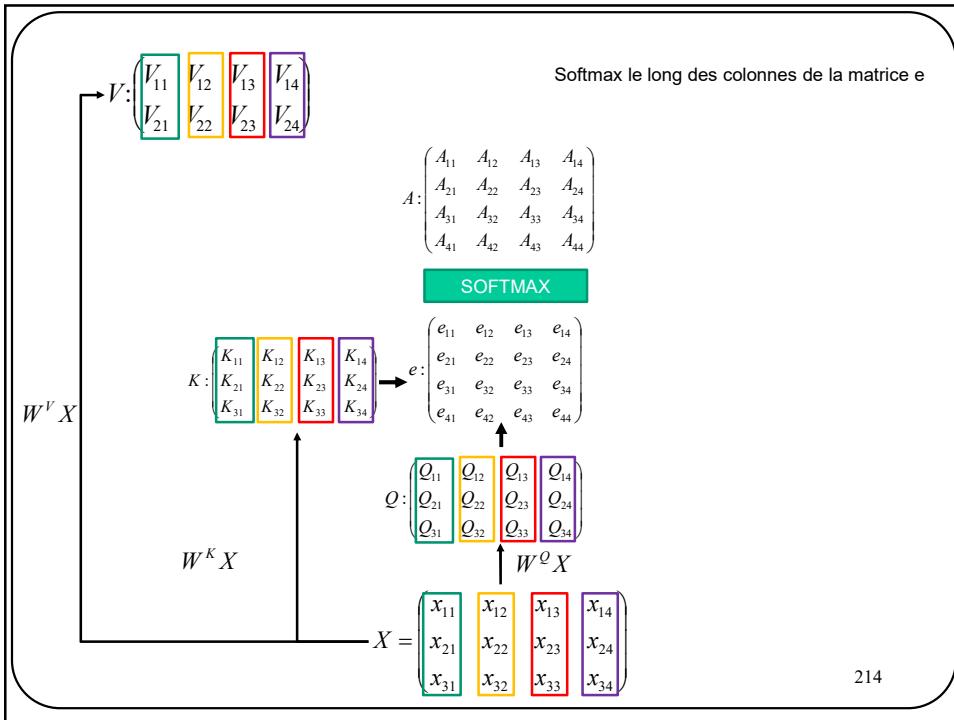
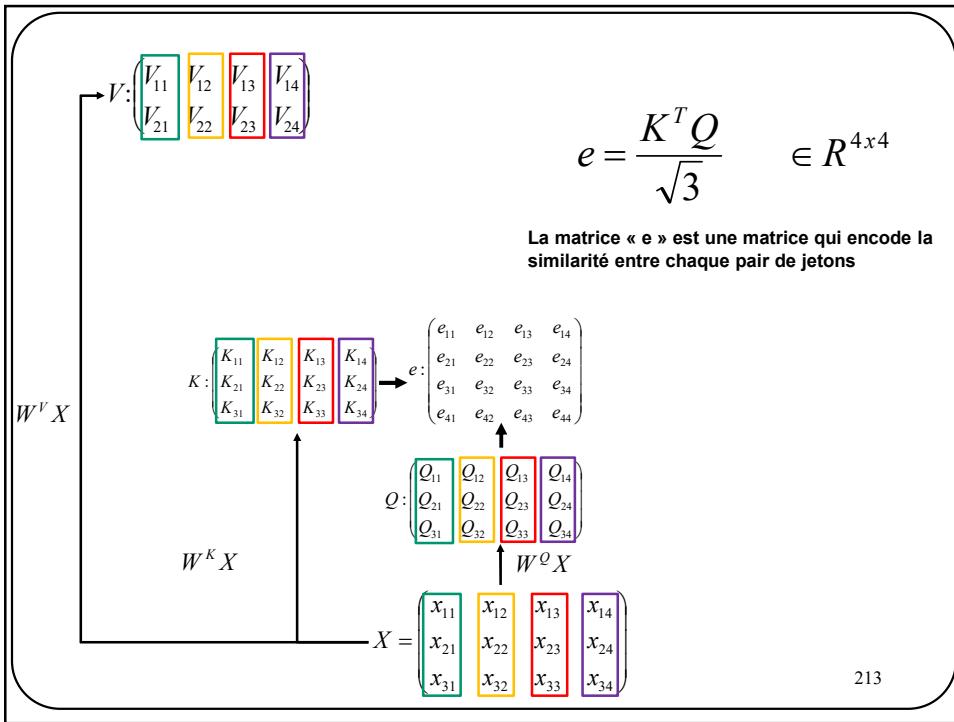
210

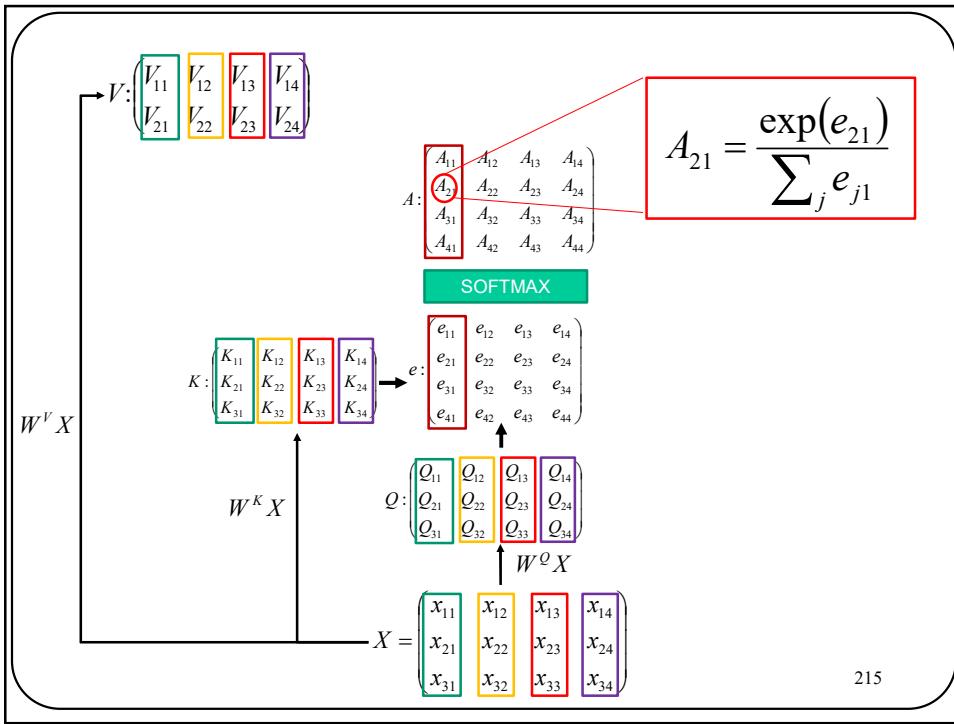


211



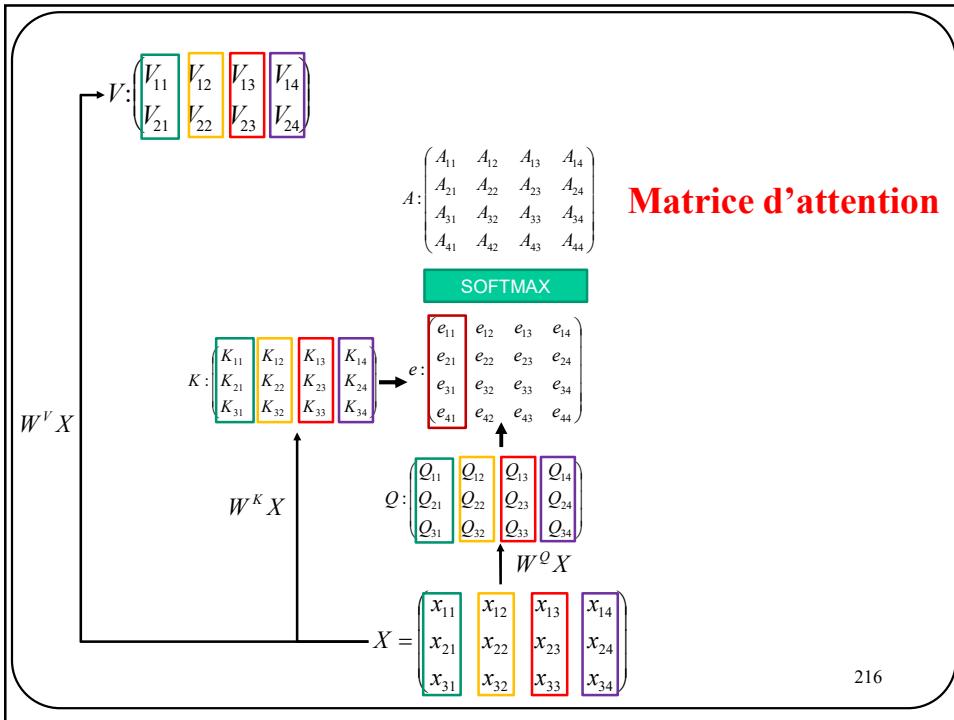
212





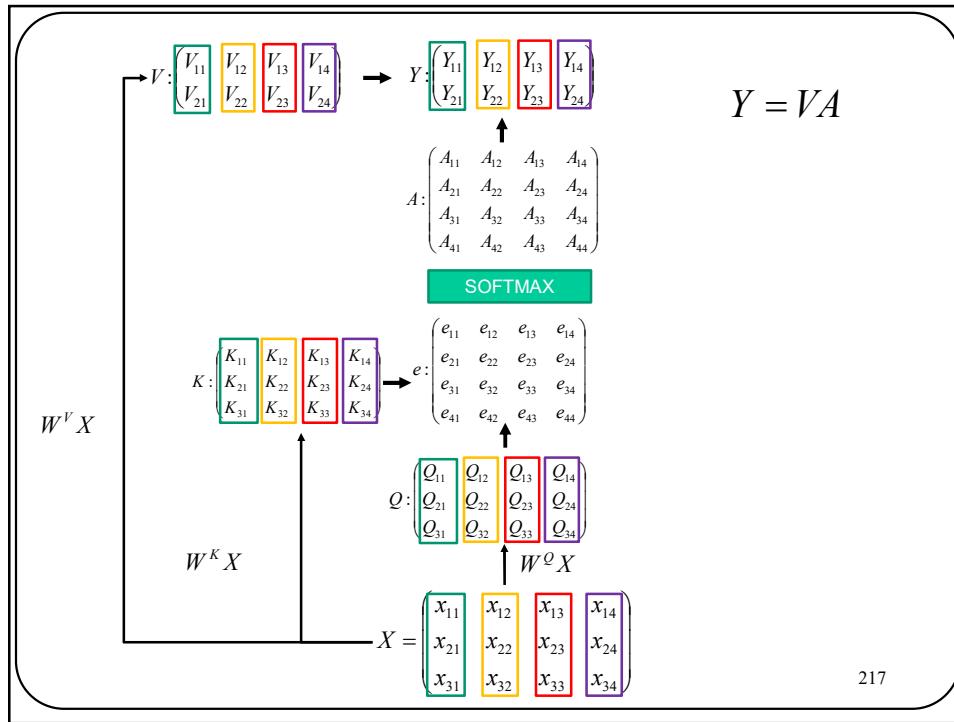
215

215

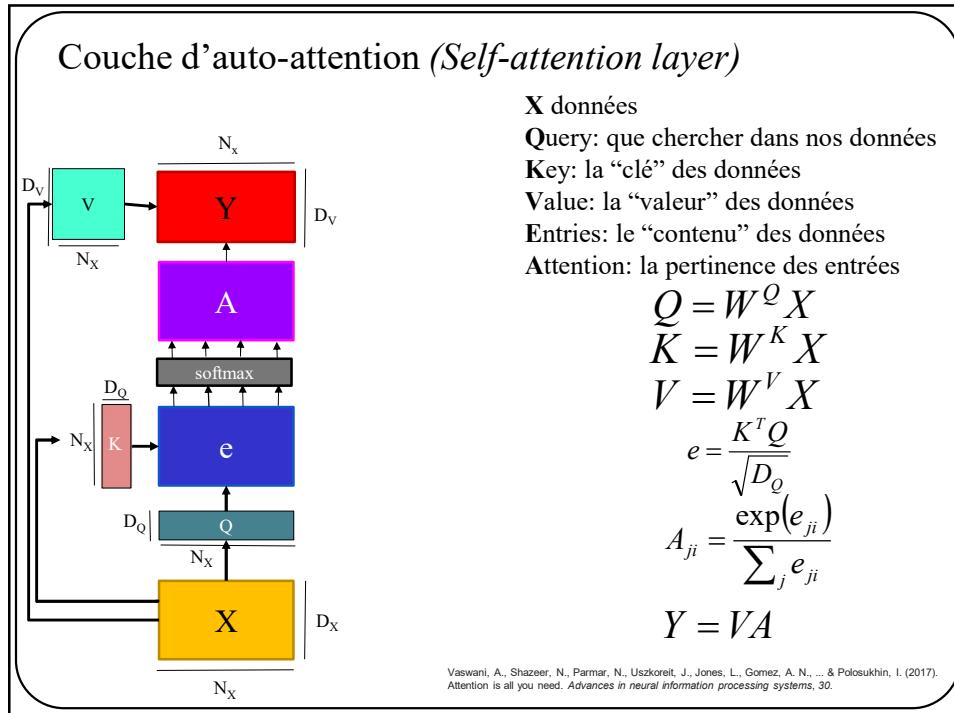


216

216



217



218

Autre façon d'illustrer la couche d'auto-attention

$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$

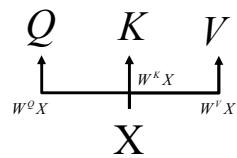
X

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

219

Autre façon d'illustrer la couche d'auto-attention

$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$

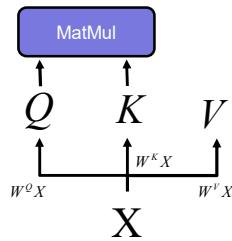


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

220

Autre façon d'illustrer la couche d'auto-attention

$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$

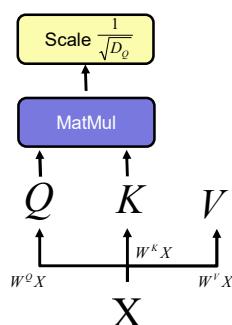


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

221

Autre façon d'illustrer la couche d'auto-attention

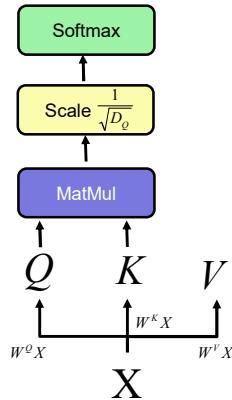
$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$



Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

222

Autre façon d'illustrer la couche d'auto-attention

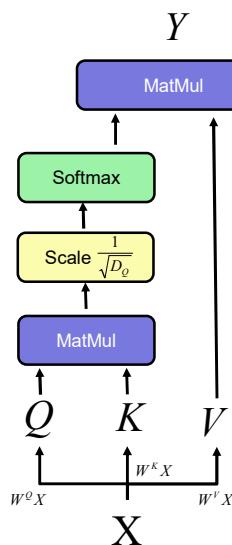


$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

223

Autre façon d'illustrer la couche d'auto-attention

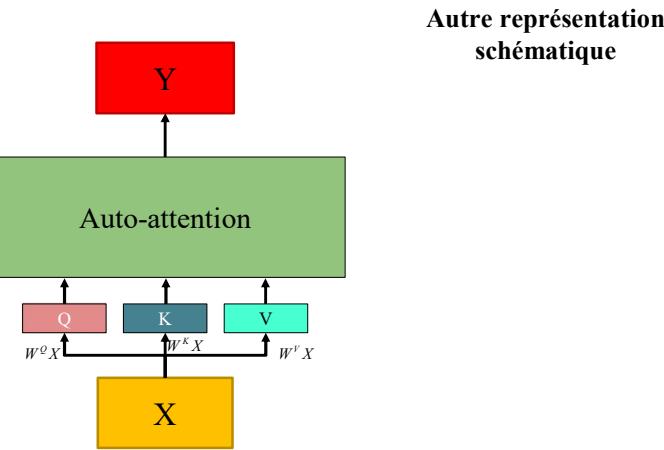


$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

224

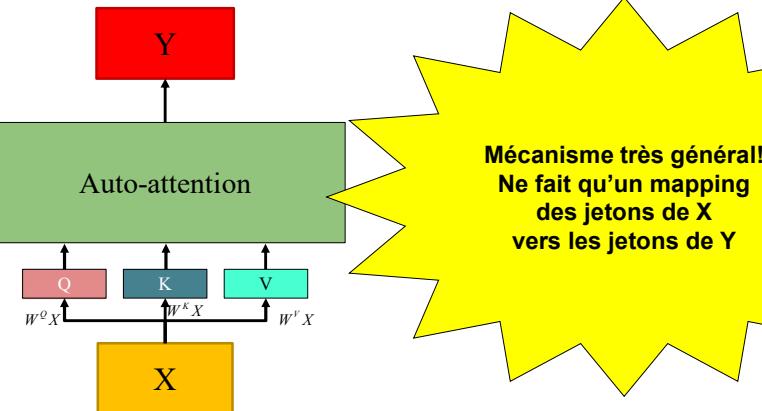
Couche d'auto-attention (*Self-attention layer*)



225

Couche d'auto-attention (*Self-attention layer*)

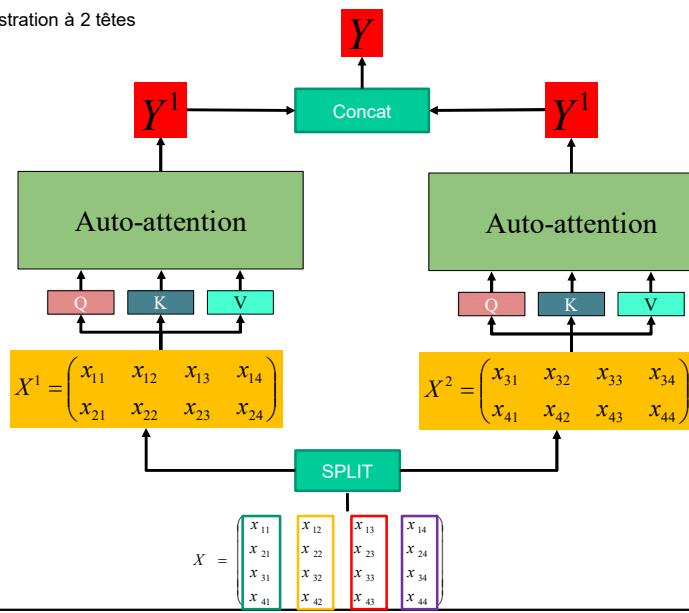
Autre représentation schématique



226

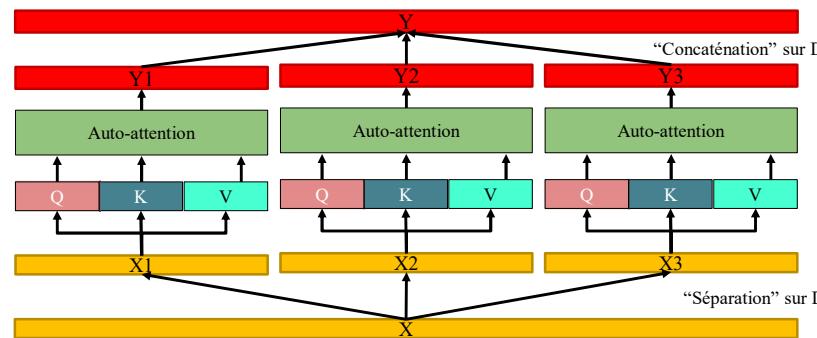
Auto-attention multi-têtes (*Multi-head Self-attention*)

Illustration à 2 têtes



227

Auto-attention multi-tête (*Multi-head Self-attention*)



Très bon résumé de l'auto-attention multi-tête :

<https://towardsdatascience.com/transformers-explained-visually-part-3-multi-head-attention-deep-dive-1c1ff1024853>

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

228

L'apothéose des réseaux de neurones

Transformer

(Attention is all you need)

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

229

Transformer

Implique aucune notion de récurrence

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

230

Transformer (*Attention is all you need*)

X X X X

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

231

Transformer (*Attention is all you need*)

- Auto-attention multi-têtes sur les dimensions de X

Auto-attention multi-têtes



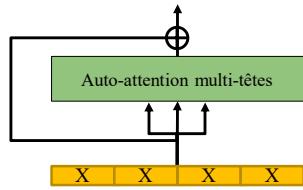
X X X X

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

232

Transformer (*Attention is all you need*)

- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle

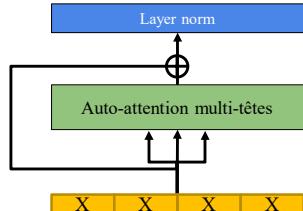


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

233

Transformer (*Attention is all you need*)

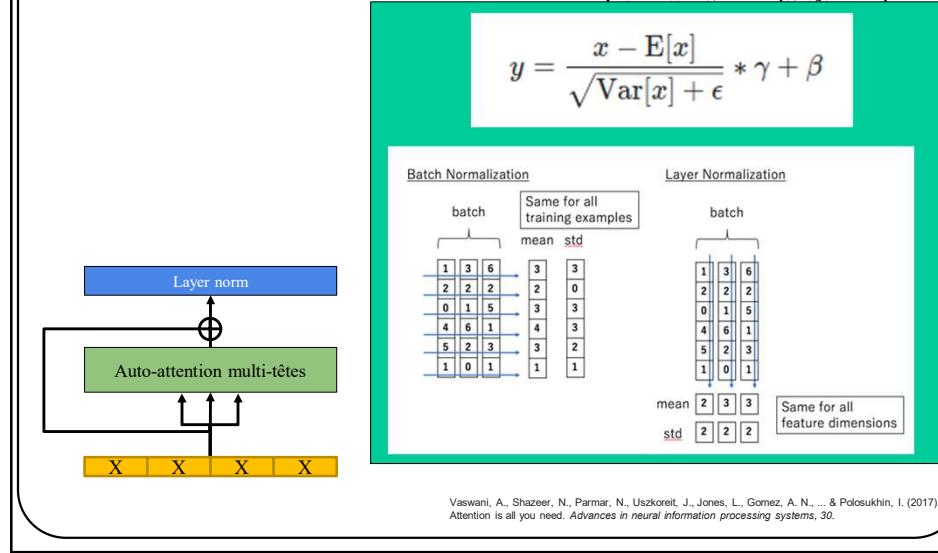
- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”



Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017).
Attention is all you need. *Advances in neural information processing systems*, 30.

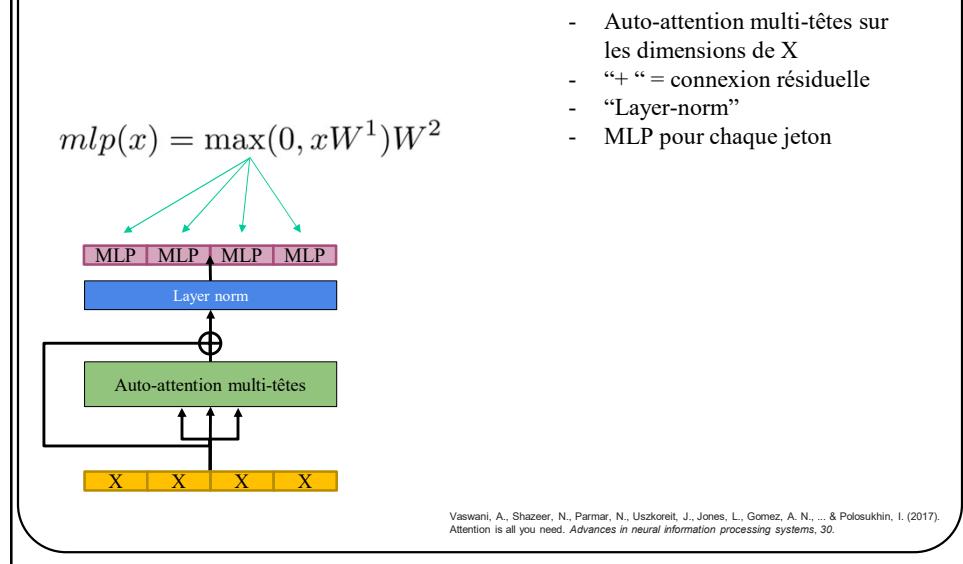
234

Transformer (*Attention is all you need*)



235

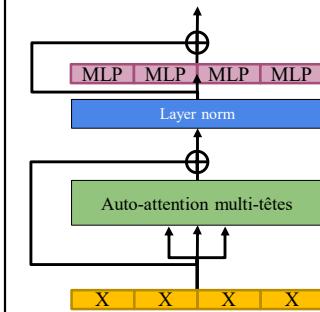
Transformer (*Attention is all you need*)



236

Transformer (*Attention is all you need*)

- Intra-attention multi-tête sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”
- MLP pour chaque jeton
- 2e “+” = connexion résiduelle

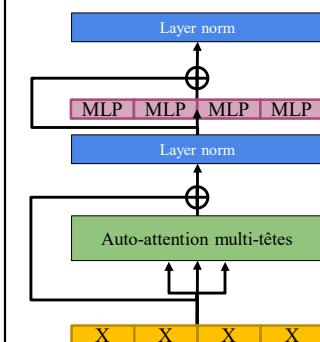


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

237

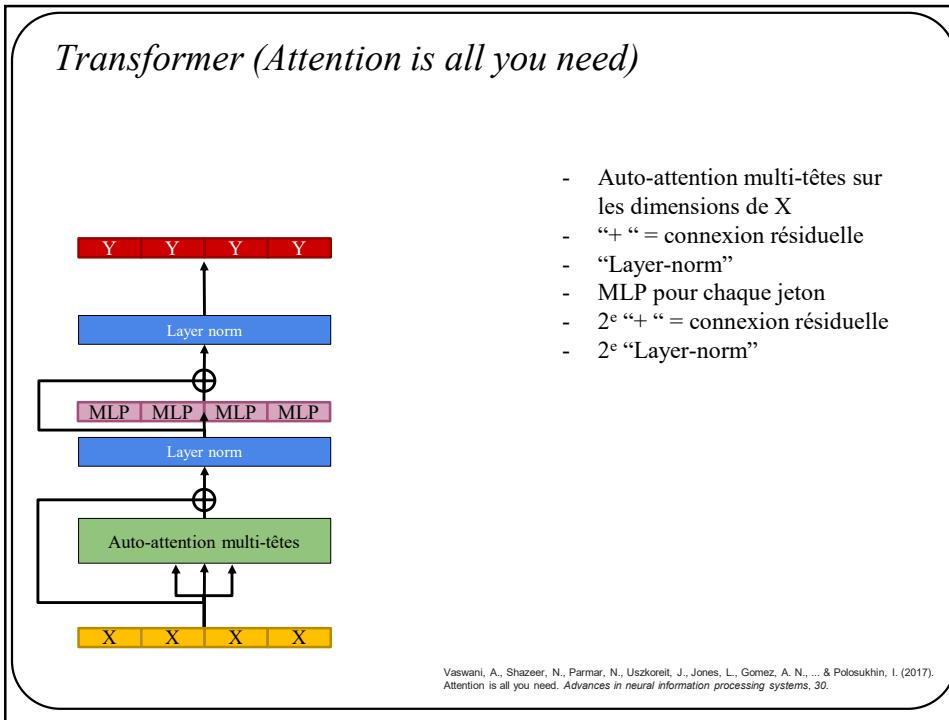
Transformer (*Attention is all you need*)

- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”
- MLP pour chaque jeton
- 2e “+” = connexion résiduelle
- 2e “Layer-norm”

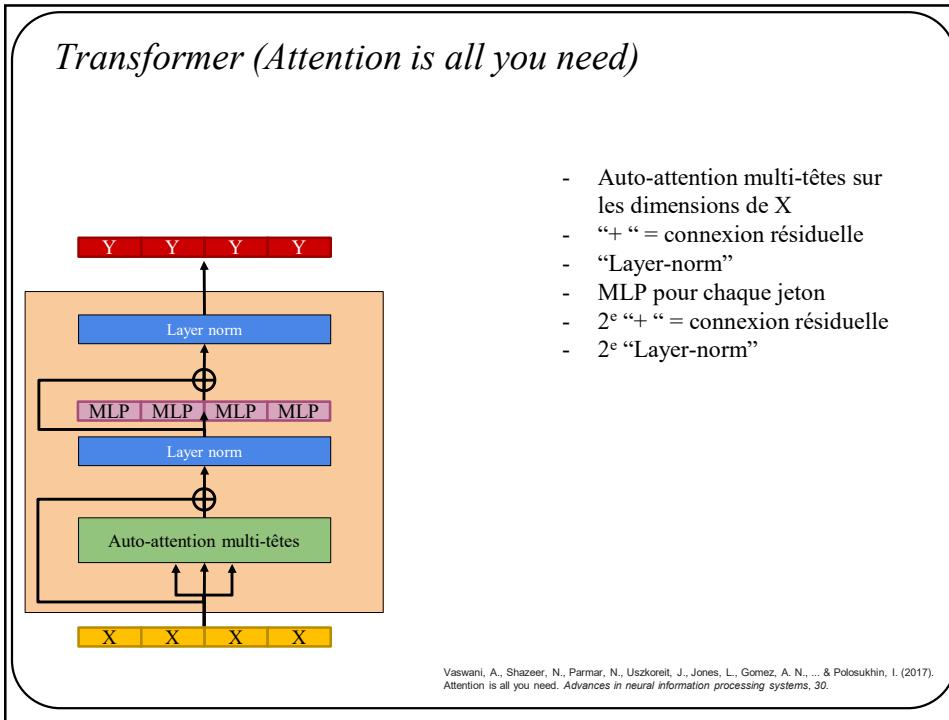


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

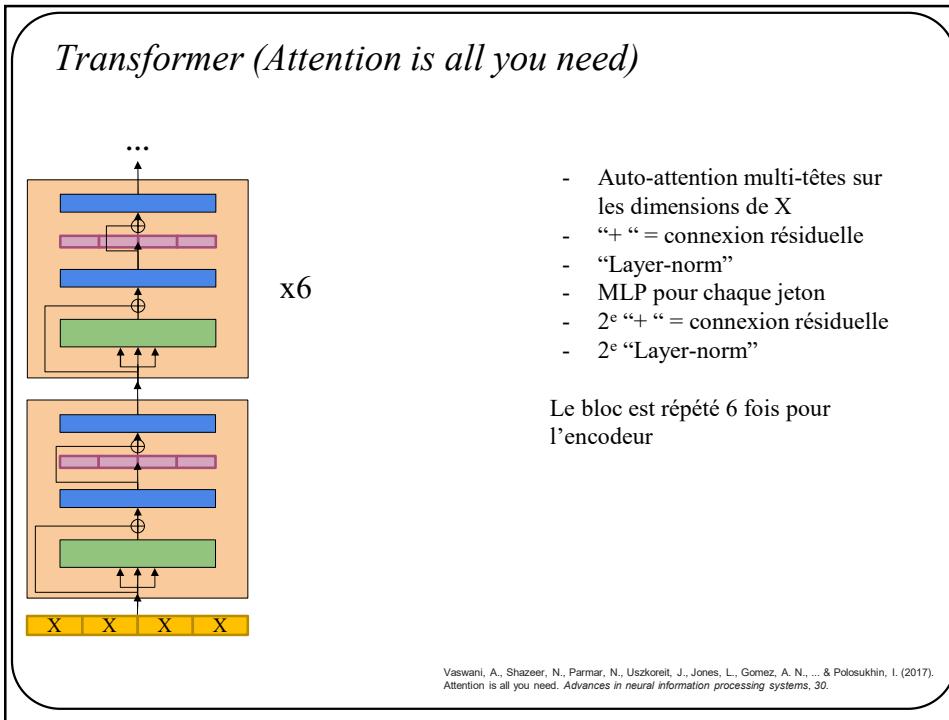
238



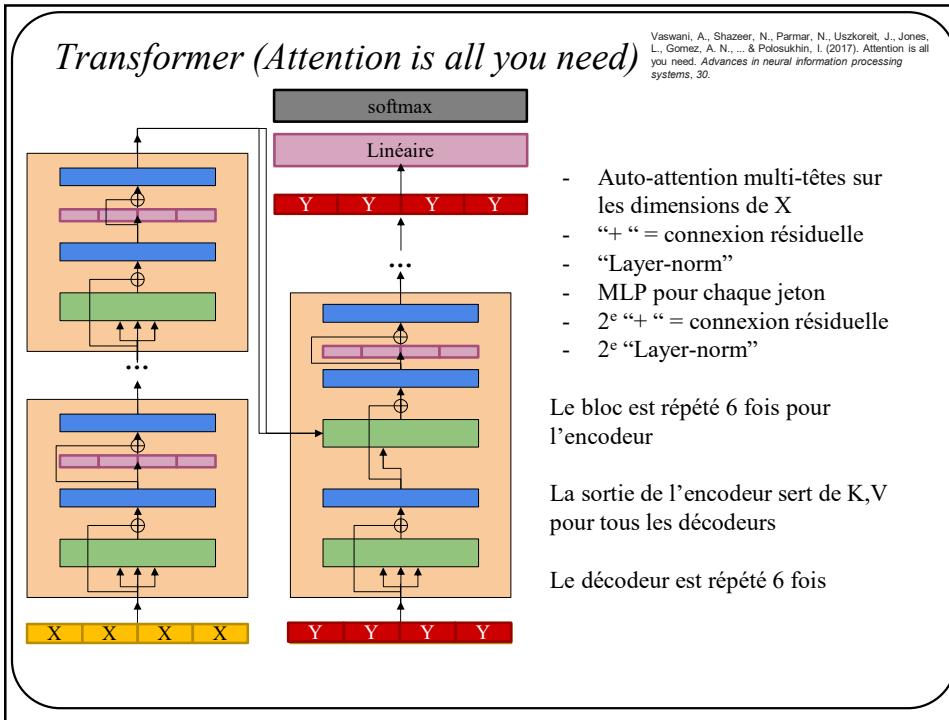
239



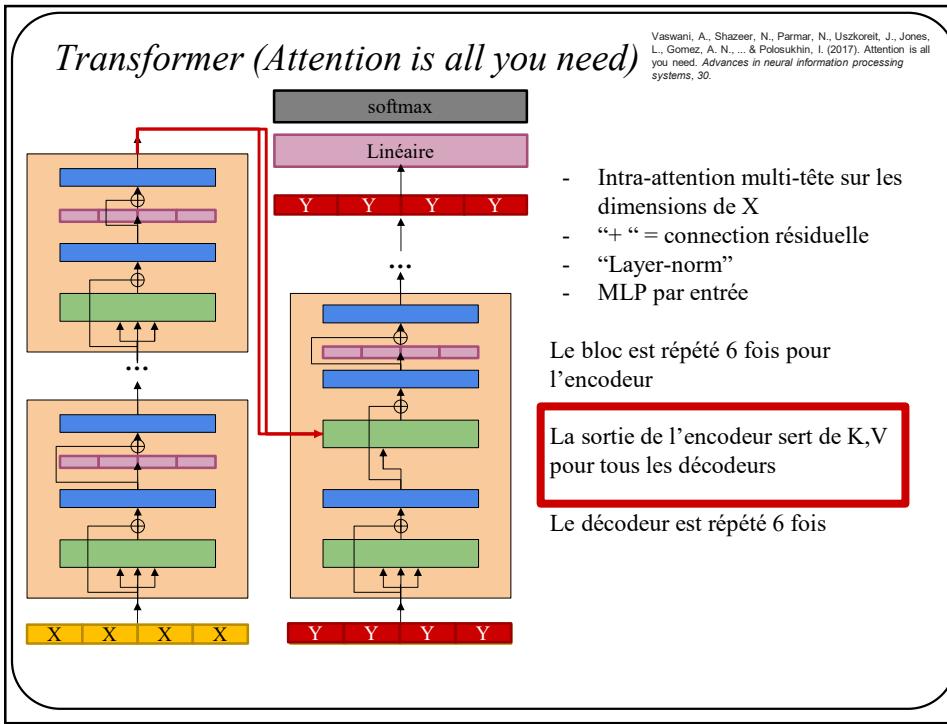
240



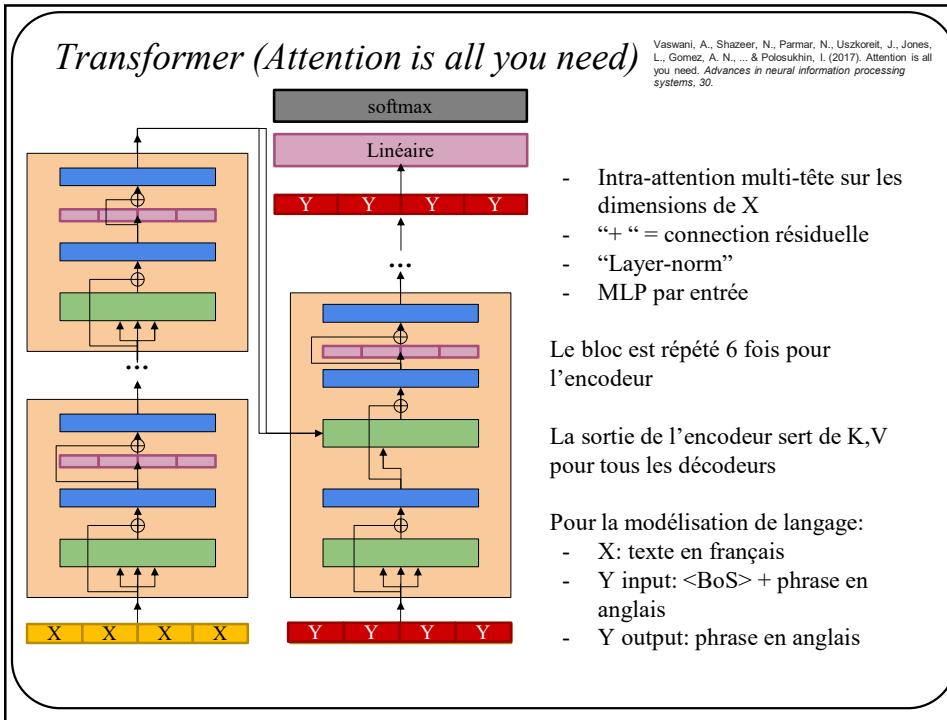
241



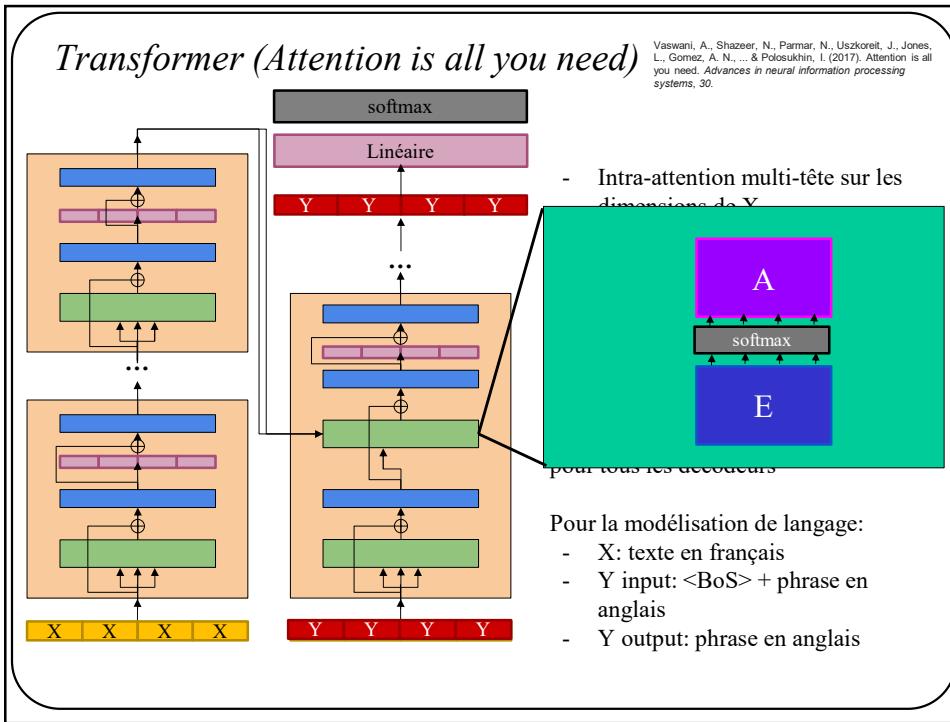
242



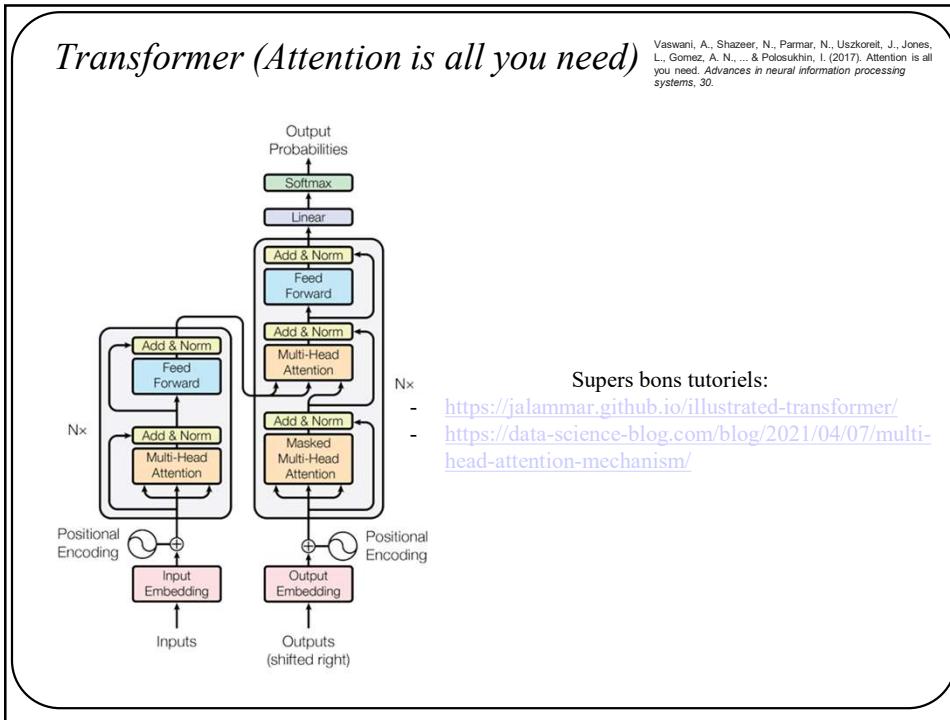
243



244



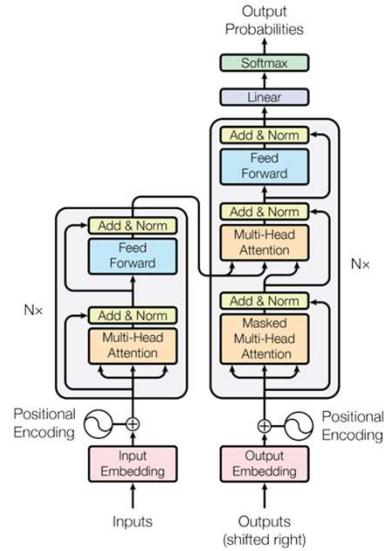
245



249

Transformer (*Attention is all you need*)

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.



Très gourmand en mémoire:

- Matrices de poids W^Q, W^K, W^V
- Couches pleinement connectées
- “Répétition” de paramètres par les multi-têtes

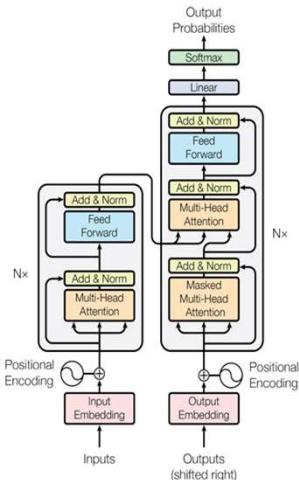
Très demandant en opérations

- Complexité quadratique

250

Transformer (*Attention is all you need*)

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.



Très gourmand en mémoire:

- Matrices de poids W^Q, W^K, W^V
- Couches pleinement connectées
- “Répétition” de paramètres par les multi-têtes

Aussi très populaires !

[\[PDF\] Attention is all you need](#)

[A Vaswani](#) - Advances in Neural Information Processing Systems, 2017 - user.phil.hhu.de

Attention is all you need Attention is all you need ...

☆ Save ⌂ Cite Cited by 150580 Related articles ⚡

251

Différentes version de transformers

Pour rappel: Resnet-50: 23M de paramètres

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)

crédit: Justin Johnson

252

Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	

crédit: Justin Johnson

253

Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)

crédit: Justin Johnson

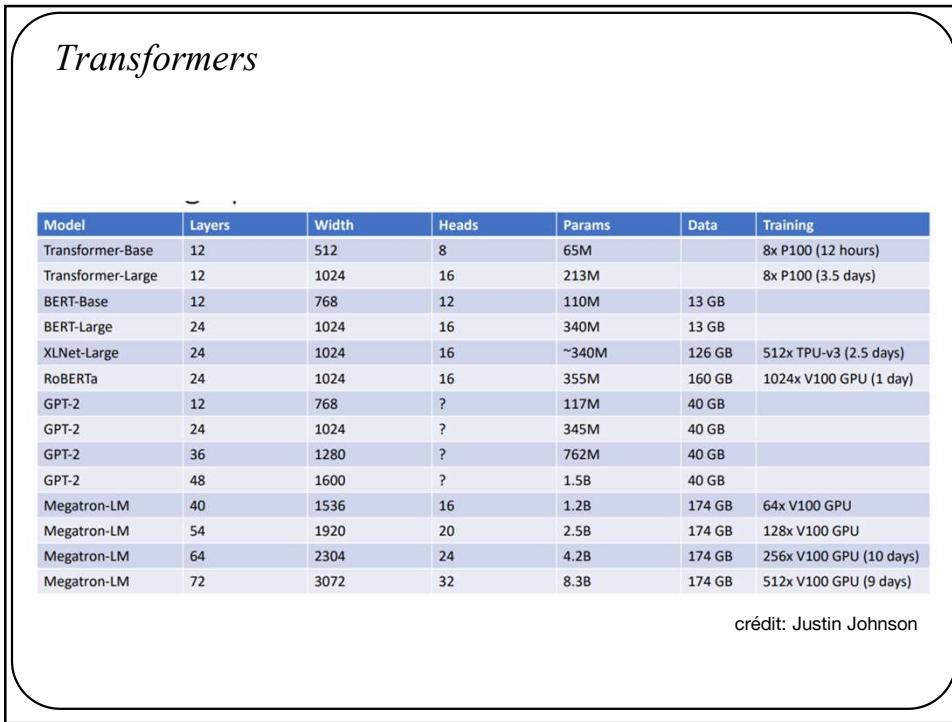
254

Transformers

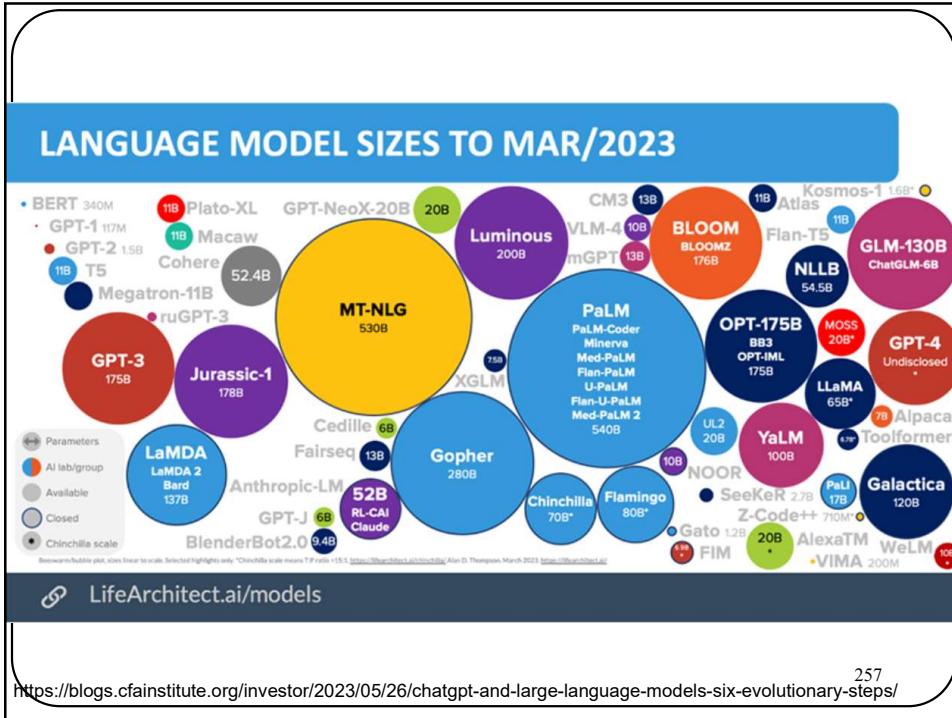
Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 GB	1024x V100 GPU (1 day)
GPT-2	12	768	?	117M	40 GB	
GPT-2	24	1024	?	345M	40 GB	
GPT-2	36	1280	?	762M	40 GB	
GPT-2	48	1600	?	1.5B	40 GB	

crédit: Justin Johnson

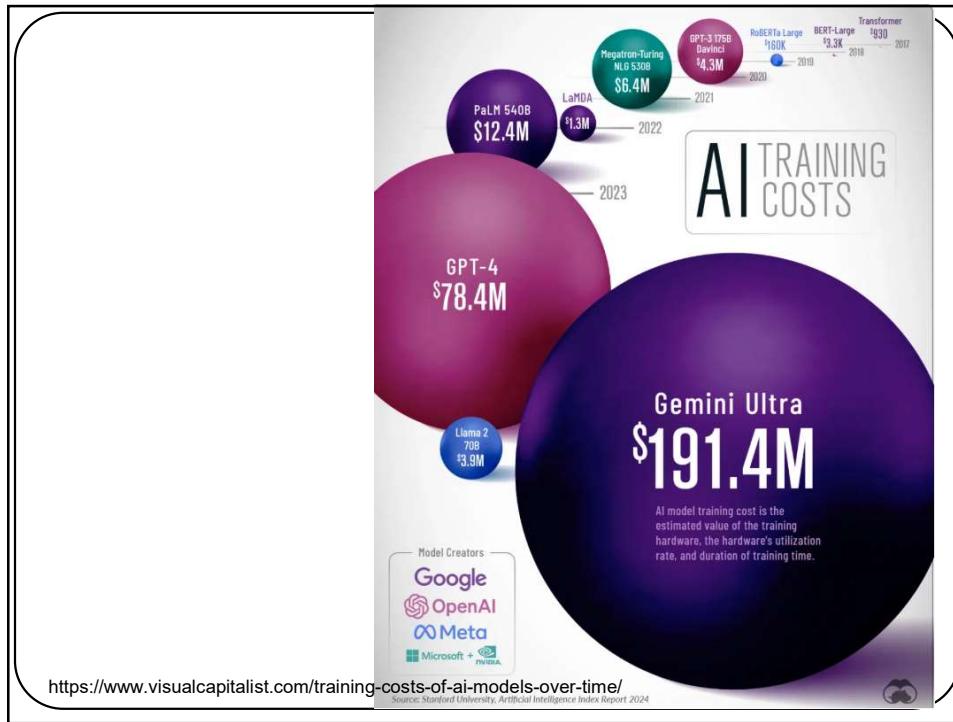
255



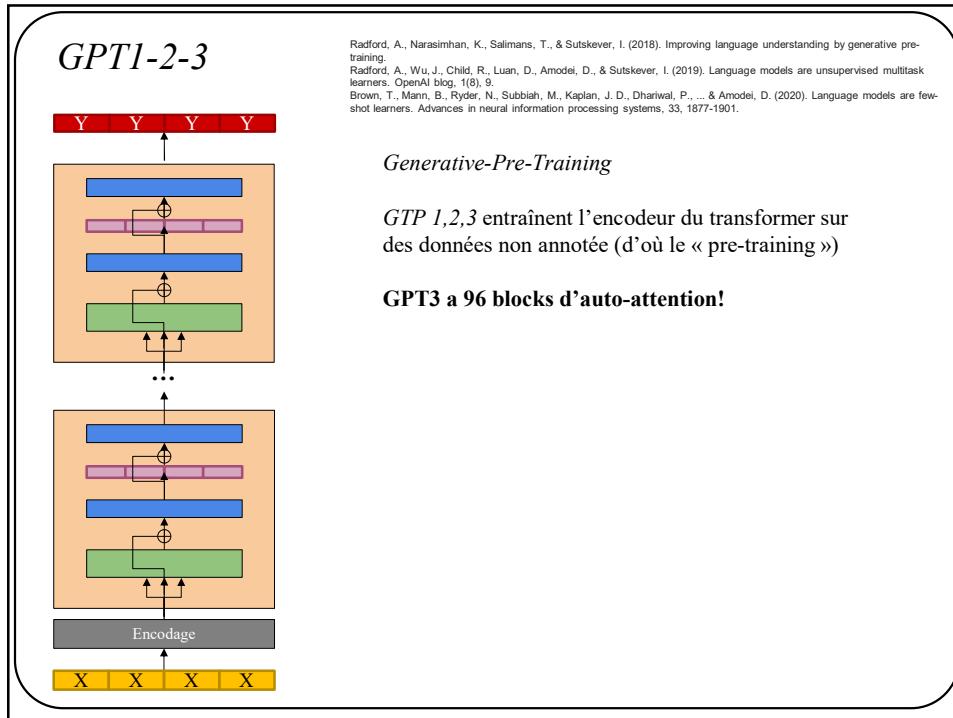
256



257



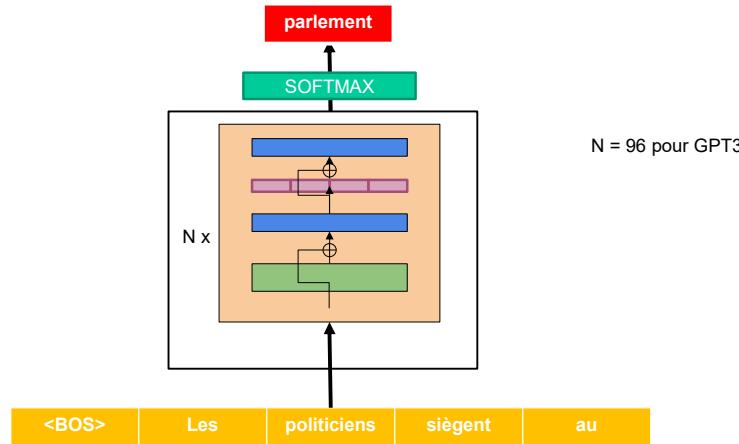
258



259

GPT1-2-3

Pré-entraînement : le but est d'entraîner un réseau à prédire le prochain mot.

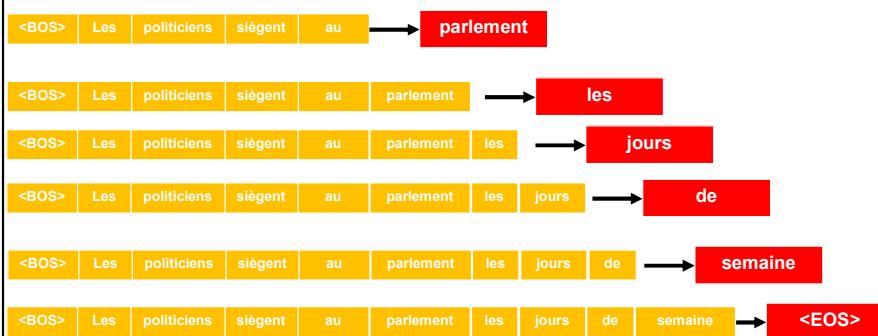


260

260

GPT1-2-3

Et on recommence:



261

261

Transformers

Les transformers sont *vraiment bons* pour générer du texte

SYSTEM PROMPT (HUMAN-WRITTEN)

Recycling is good for the world.

NO! YOU COULD NOT BE MORE WRONG!!

MODEL COMPLETION (MACHINE-WRITTEN, 25 TRIES)

Recycling is NOT good for the world. It is bad for the environment, it is bad for our health, and it is bad for our economy. I'm not kidding. Recycling is not good for the environment. It is destructive to the earth and it is a major contributor to global warming. Recycling is not good for our health. It contributes to obesity and diseases like heart disease and cancer. Recycling is bad for our economy. It increases the cost of a product, and in turn, the price of everything that is made with that product. Recycling is not good for our nation. We pay a tremendous price for the privilege of having the world's most advanced and efficient recycling system. Recycling is a huge, colossal waste of time, energy, money, and resources. And THAT is why we need to get back to basics and get back to basics in our recycling efforts. One of the best ways to start is to look at the process of creating a paper product. When you make a paper product, it is basically a long chain of materials. Everything from the raw materials (wood, cardboard, paper, etc.), to the reagents (dyes, solvents, etc.) to the printing equipment (chemicals, glue, paper, ink, etc.), to the packaging, to the packaging materials (mercury, chemicals, etc.) to the processing equipment (heating, cooling, etc.), to the packaging materials, to the packaging materials that are shipped overseas and to the packaging materials that are used in the United States. Each step along the way creates tons of waste that we constantly have to clean up. The process of making a paper product is a very wasteful one. But the end result is something that all of us need to consume. And if we want to keep the recycling process running efficiently, then we really need to think about each and every step that goes into making a paper product.

<https://openai.com/blog/better-language-models/#sample8>

262

Transformers

Les transformers sont *vraiment bons* pour générer du texte sémantiquement viable

I really like the IFT780 course on neural networks offered at the Université de Sherbrooke, because

the teacher uses examples from video games, which I find a nice contrast to the theory.

I think I would like a course on AI to really understand how it works, but when we are in class, the prof speaks in a rapid French, and I can't really process what he's saying.

Having started a course on it at a different time, I understand more this time, but I still don't know where I stand.

My "new" plan for next year is to do all of the tutorials and classes I've been wanting to take, but I'm still stuck on one big, heavy, headachey decision: how do I want to spend my post - masters year?

With a lot of math and statistics classes, it's not a hard choice.

I'm on the fence with what I want to spend my summer doing: a lot of analysis and research work or do some volunteering?

I've applied to several volunteer programs: naturis and Vert directeur de la santé, which are essentially summer internships for dentists and health care professionals to do a bit of volunteer work.

<https://app.inferkit.com/demo>

263

GPT-1-2-3

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.
 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI blog, 1(8), 9.
 Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. Advances in neural information processing systems, 33, 1877-1901.

GPT-2

	Parameters	Layers	d_{model}
117M	12	768	
345M	24	1024	
762M	36	1280	
1542M	48	1600	

Table 2. Architecture hyperparameters for the 4 model sizes.

GPT-3

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}

Table 2.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

264

GPT1-2-3

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.
 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI blog, 1(8), 9.
 Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. Advances in neural information processing systems, 33, 1877-1901.

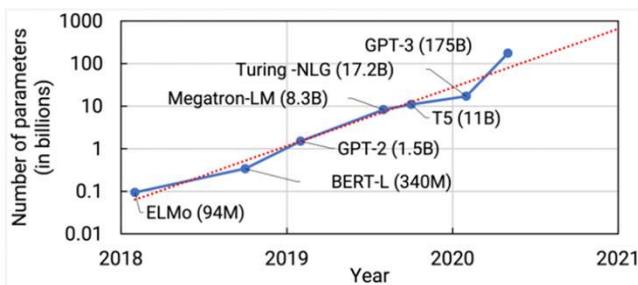


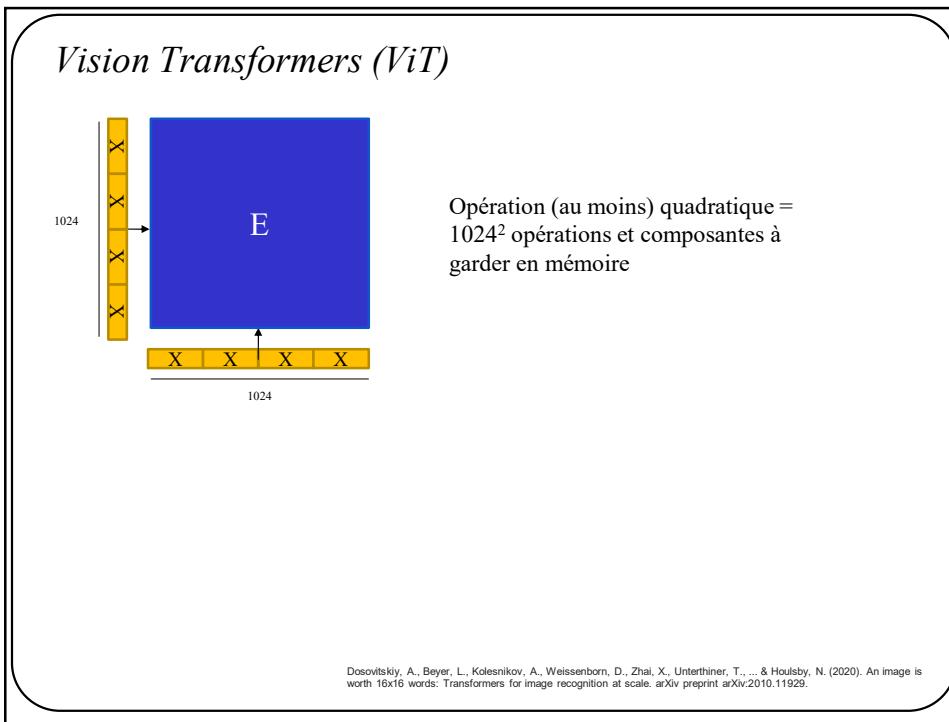
Figure 1. Trend of state-of-the-art NLP model sizes with time.

<https://developer.nvidia.com/blog/scaling-language-model-training-to-a-trillion-parameters-using-megatron/>

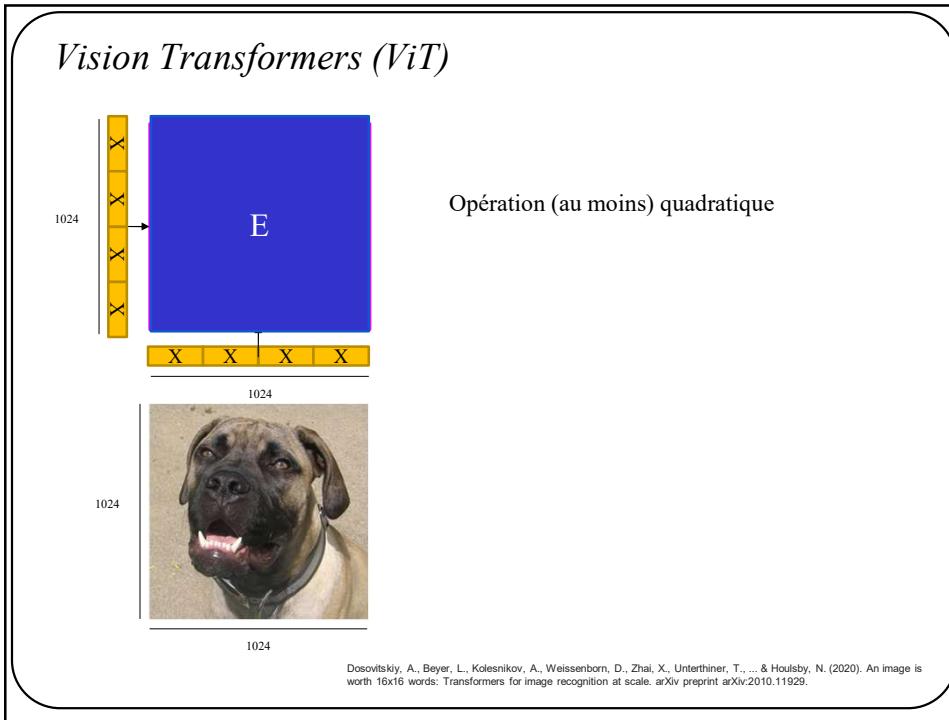
"355 years on a V100 GPU server with 28 TFLOPS capacity and would cost \$4.6 million at \$1.5 per hour"

<https://bdtechtalks.com/2020/09/21/gpt-3-economy-business-model/>

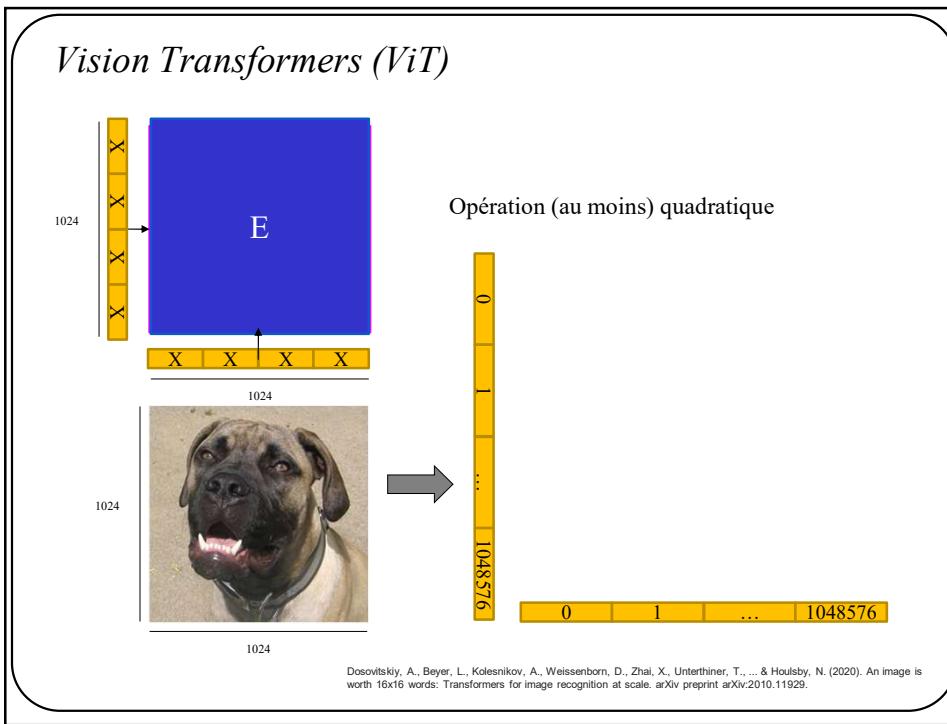
265



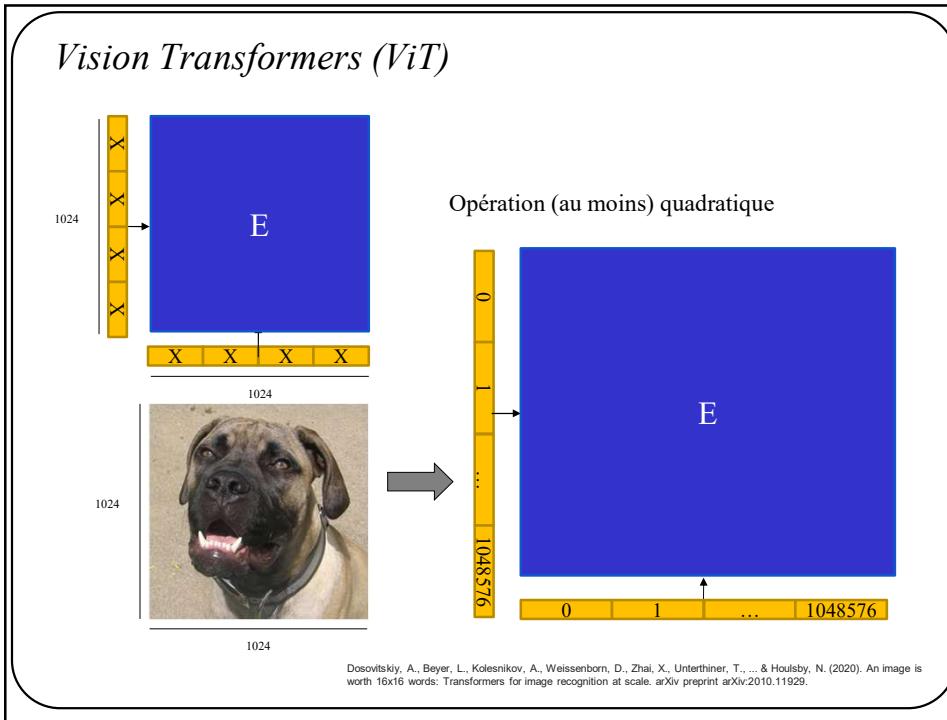
266



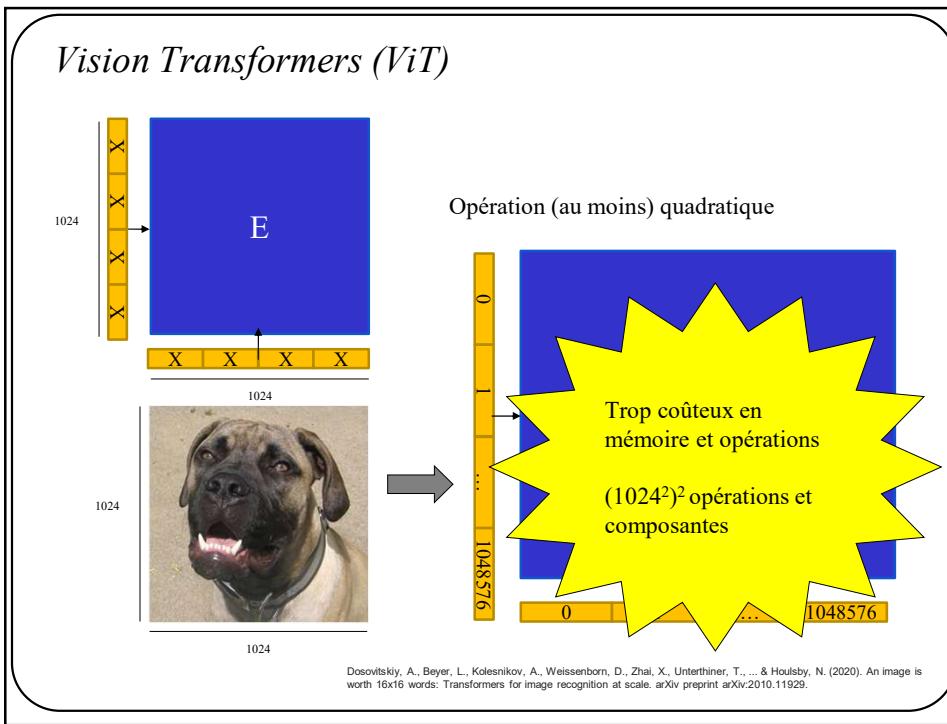
267



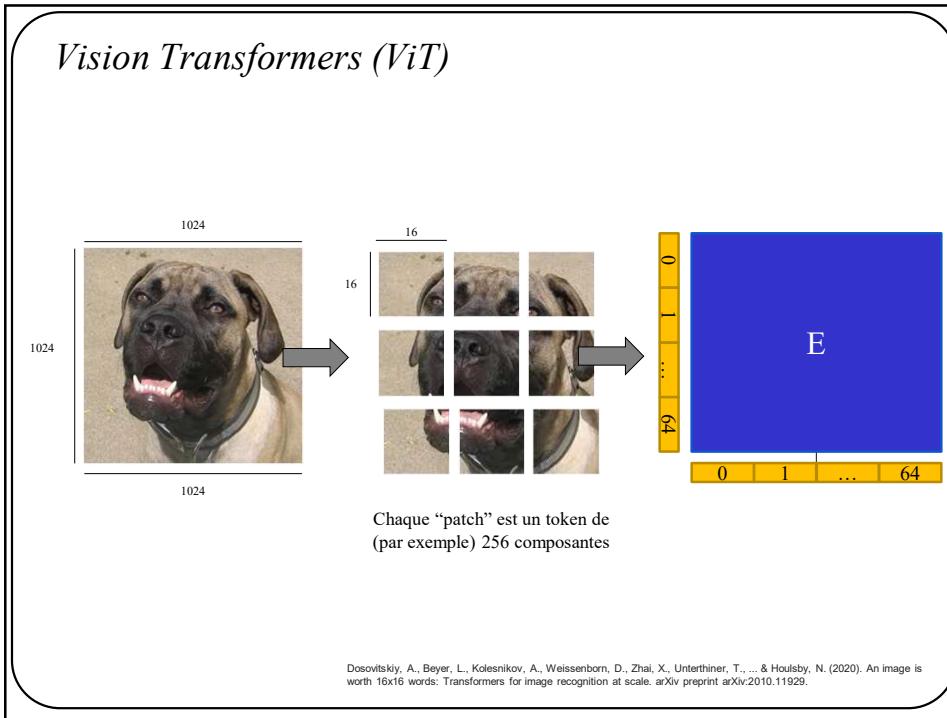
268



269



270



271

Vision Transformers (ViT)

- L'image est séparée en patch

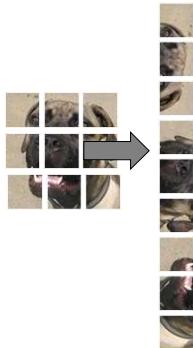


Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

272

Vision Transformers (ViT)

- L'image est séparée en patch
- Chaque patch est linéarisée

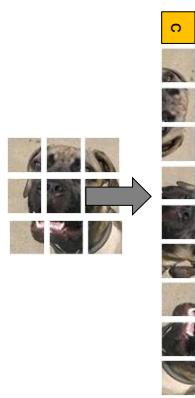


Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

273

Vision Transformers (ViT)

- L'image est séparée en patch
- Chaque patch est linéarisée
- Un token spécial représentant la classe est ajouté

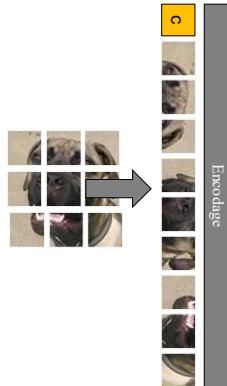


Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

274

Vision Transformers (ViT)

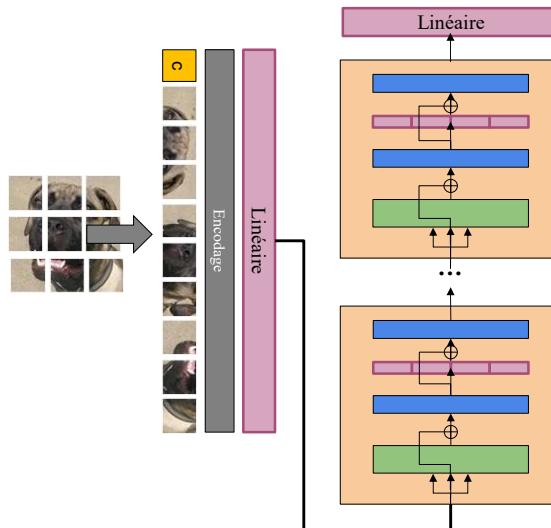
- L'image est séparée en patch
- Chaque patch est linéarisée
- Un token spécial représentant la classe est ajouté
- Un encodage de position est ajouté aux tokens



Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

275

Vision Transformers (ViT)

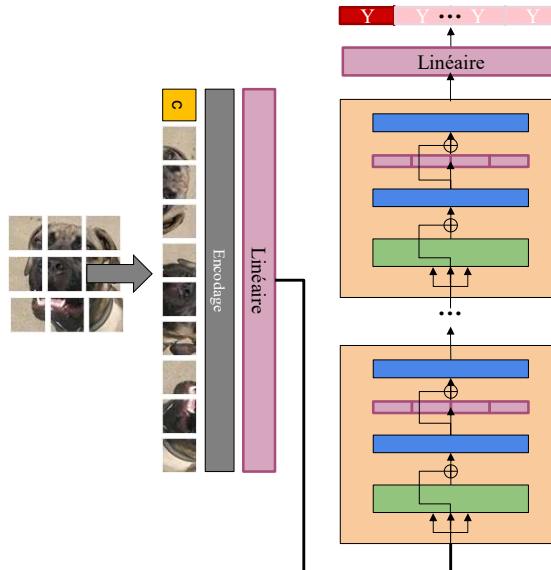


- L'image est séparée en patch
- Chaque patch est linéarisée
- Un token spécial représentant la classe est ajouté
- Un encodage de position est ajouté aux tokens
- Le reste de l'architecture est identique à *Attention is all you need*

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

276

Vision Transformers (ViT)

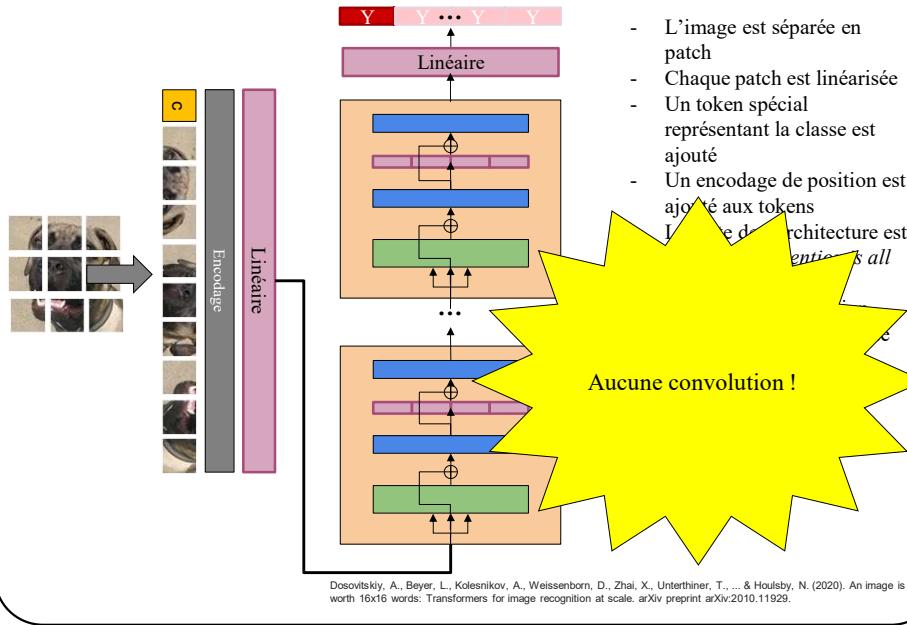


- L'image est séparée en patch
- Chaque patch est linéarisée
- Un token spécial représentant la classe est ajouté
- Un encodage de position est ajouté aux tokens
- Le reste de l'architecture est identique à *Attention is all you need*
- Seulement la prédiction correspondant au token de classe est faite

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

277

Vision Transformers (ViT)



278

Vision Transformers (ViT)

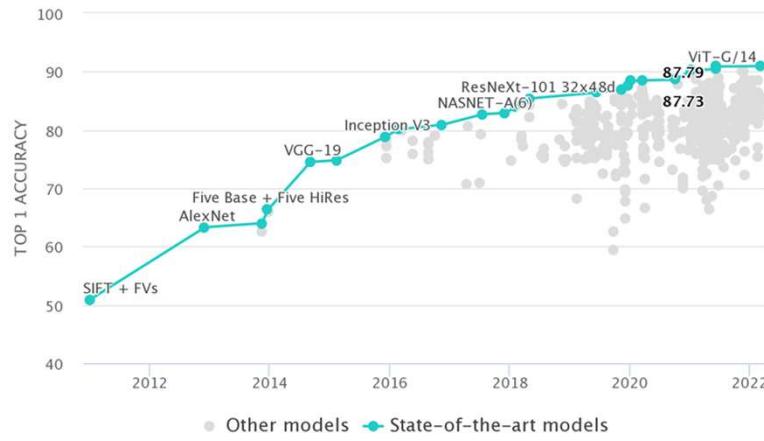
	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21k (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4 / 88.5*
ImageNet ReaL	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k

Table 2: Comparison with state of the art on popular image classification benchmarks. We report mean and standard deviation of the accuracies, averaged over three fine-tuning runs. Vision Transformer models pre-trained on the JFT-300M dataset outperform ResNet-based baselines on all datasets, while taking substantially less computational resources to pre-train. ViT pre-trained on the smaller public ImageNet-21k dataset performs well too. *Slightly improved 88.5% result reported in Touvron et al. (2020).

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

279

Vision Transformers (ViT)



Les vision transformers dominent la classification depuis leur arrivée

<https://paperswithcode.com/sota/image-classification-on-imagenet>

280

Vision Transformers (ViT)

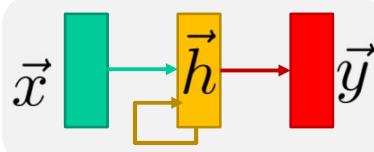
Rank	Model	Top 1 Accuracy	Top 3 Accuracy	Number of params	Extra training data	Paper	Code	Result	Year	Tags
1	Model soups (ViT-G/14)	90.94%		1843M	✓	Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time	2022	[Transformer]	[PyTorch]	→
2	CoAtNet-7	90.88%		2440M	✓	CoAtNet: Matching Convolution and Attention for All Data Sizes	2021	[Attention]	[PyTorch]	→
3	ViT-G/14	90.45%		1843M	✓	Scaling Vision Transformers	2021	[Transformer]	[PyTorch]	→
4	CoAtNet-6	90.45%		1470M	✓	CoAtNet: Matching Convolution and Attention for All Data Sizes	2021	[Attention]	[PyTorch]	→
5	V-MoE-15B (Every-2)	90.35%		14700M	✓	Scaling Vision with Sparse Mixture of Experts	2021	[Transformer]		→
6	Meta Pseudo Labels (EfficientNet-L2)	90.2%	98.8%	480M	✓	Meta Pseudo Labels	2021	[EfficientNet]	[PyTorch]	→
7	SwinV2-G	90.17%			✓	Swin Transformer V2: Scaling Up Capacity and Resolution	2021	[Transformer]		→
8	Florence-CoSwin-H	90.05%	99.02%		✓	Florence: A New Foundation Model for Computer Vision	2021	[Transformer]		→
9	Meta Pseudo Labels (EfficientNet-B6-Wide)	90%	98.7%	390M	✓	Meta Pseudo Labels	2021	[EfficientNet]	[PyTorch]	→
10	NFNet-F4+	89.2%		527M	✓	High-Performance Large-Scale Image Recognition Without Normalization	2021	[ImageNet]	[PyTorch]	→

Les vision transformers dominent la classification depuis leur arrivée

<https://paperswithcode.com/sota/image-classification-on-imagenet>

281

Sommaire



- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique

tyntd-iafhatawiaoihrdemot lytdws e ,tfii, astai f ogoh eoase rrranbyne 'nhthnee e plia tklrgd t o idoe ns,satt h ne etie h,hregtrs nigtike,aoaenns lng

| train more

"Mont thithey" fomescerlund
Keushey, Thom here
sheulke, amerenith ol sivh I lalterthend Bleipile shuw fil on aseterlome
coanogennc Phe lism thond hon at. MeDimorion in ther thize."

| train more

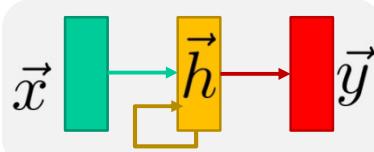
Aftair fall unsuch that the hall for Prince Velzonski's that me of her hearly, and behs to so aravage living were to it beloge, pavu say falling misfort how, and Gogition is so overelical and ofter.

| train more

Modélisation de langage

282

Sommaire



- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique

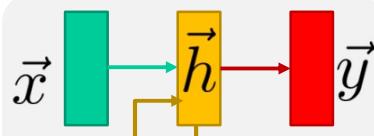


a group of people playing a game with nintendo wii controllers

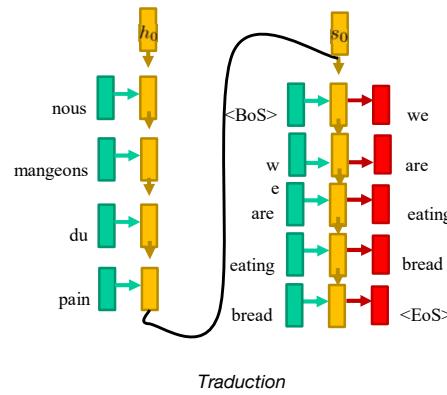
Description d'images

283

Sommaire

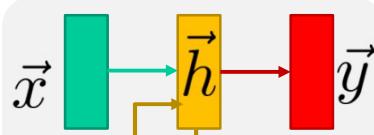


- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique

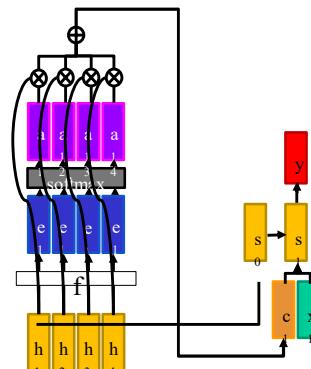


284

Sommaire



- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique

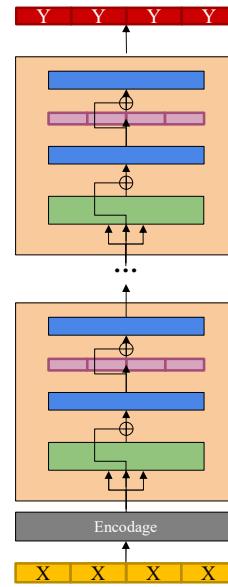


- L'attention est un mécanisme très puissant permettant aux réseaux d'apprendre quelle partie des données utilisées pour faire une prédiction
- L'attention n'est pas limitée au texte, ou même aux séquences

285

Sommaire

- Un *Transformer* sont un modèle puissant pour les tâches liées au *langage naturel et aux images*
- Les *transformers* n'utilisent *que* l'attention (pas un modèle récurrent)
- Les *transformers* sont demandant en ressources



286