

Réseaux de neurones

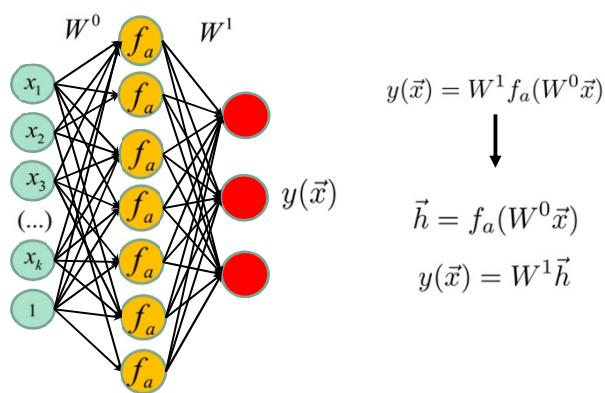
IFT 780

Réseaux récurrents

Par

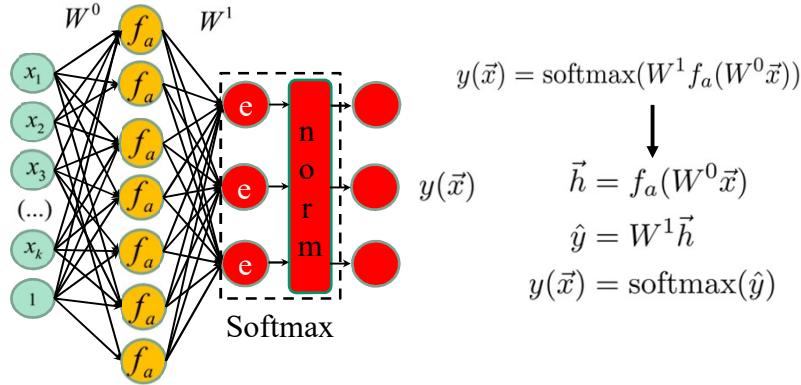
Pierre-Marc Jodoin, Antoine Théberge

Réseau de neurones de base (régression)

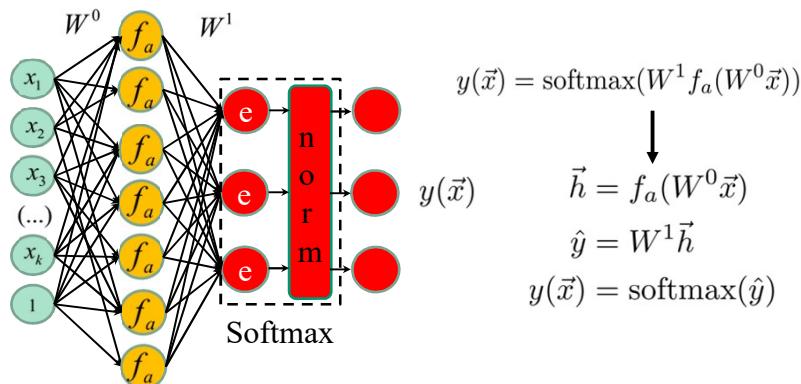


f_a : fonction d'activation

Réseau de neurones de base (classification)



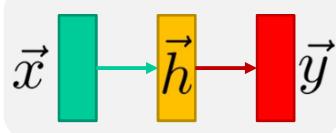
Réseau de neurones de base (classification)



Ne permettent que des tâches “1 pour 1”

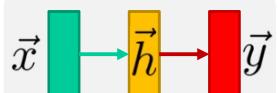
- Classification (1 image = 1 étiquette)
- Régression (1 donnée = 1 vecteur)
- Localisation (1 boîte = 1 classification + 1 régression)

Illustration simplifiée



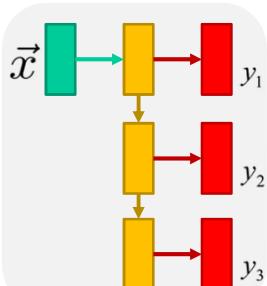
RN de base

Différentes configurations pour différentes applications

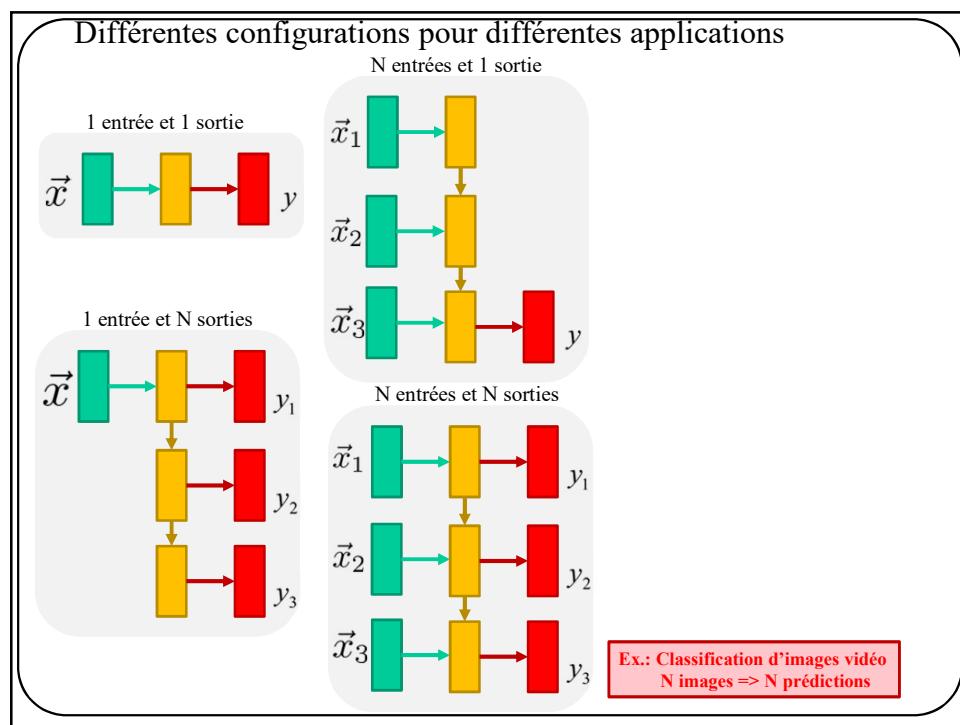
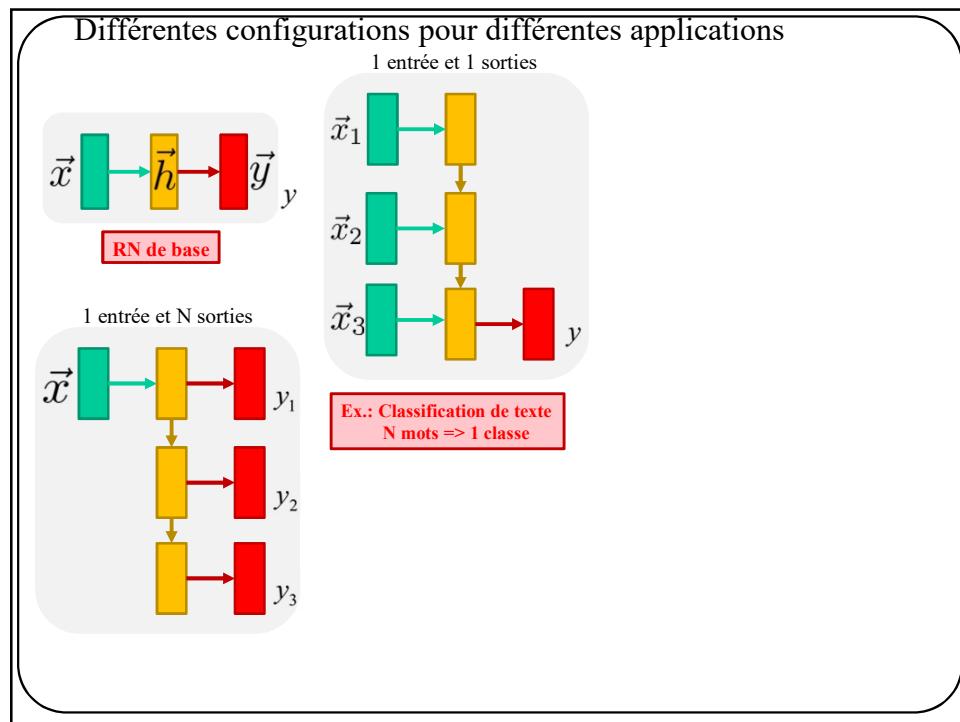


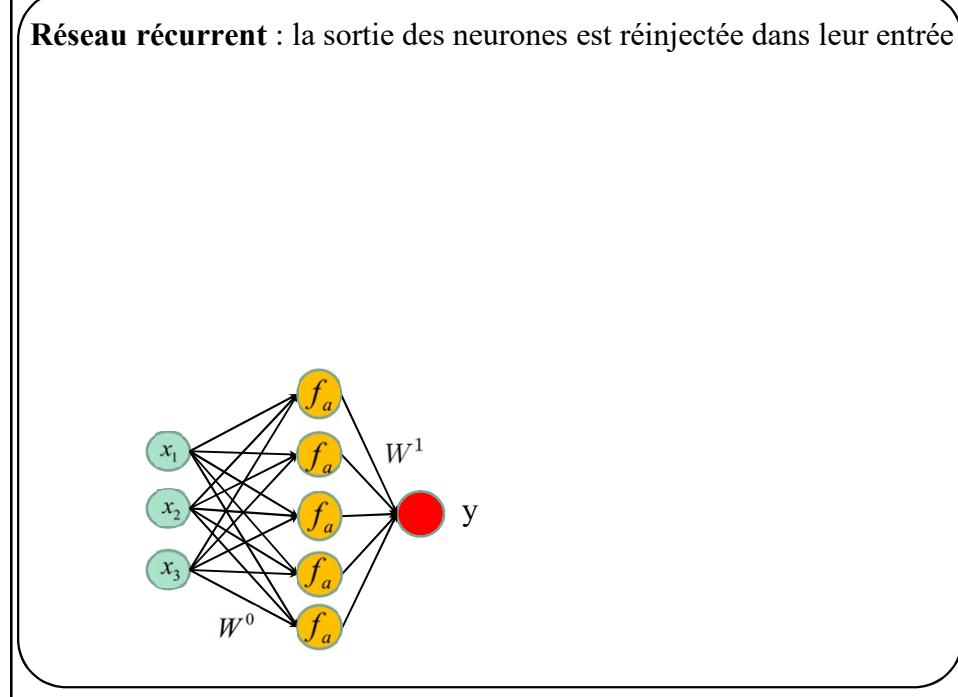
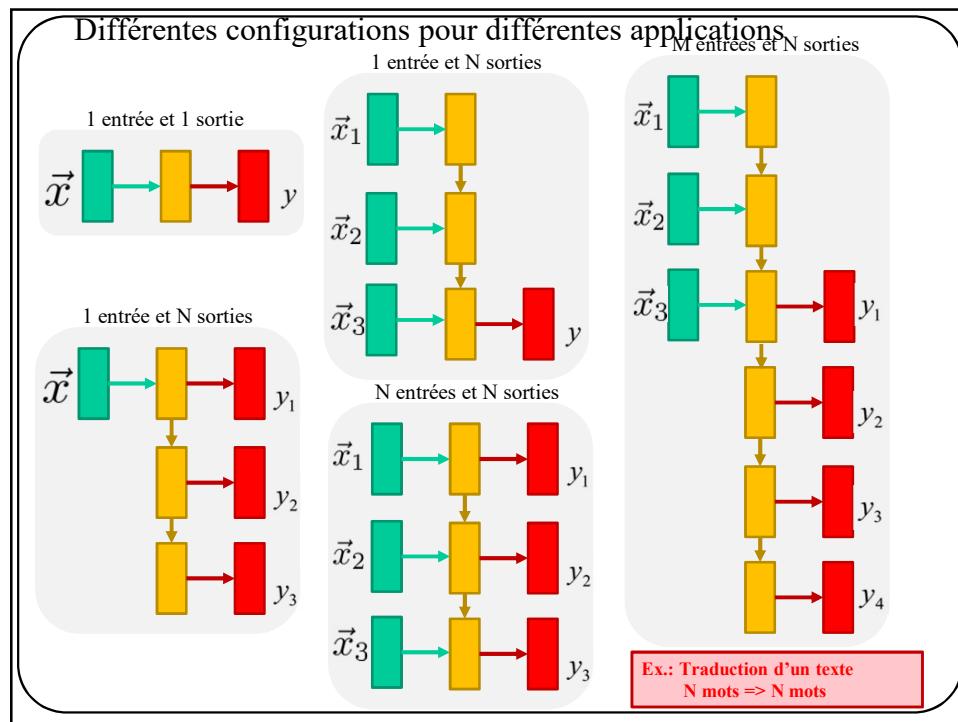
RN de base

1 entrée et N sorties

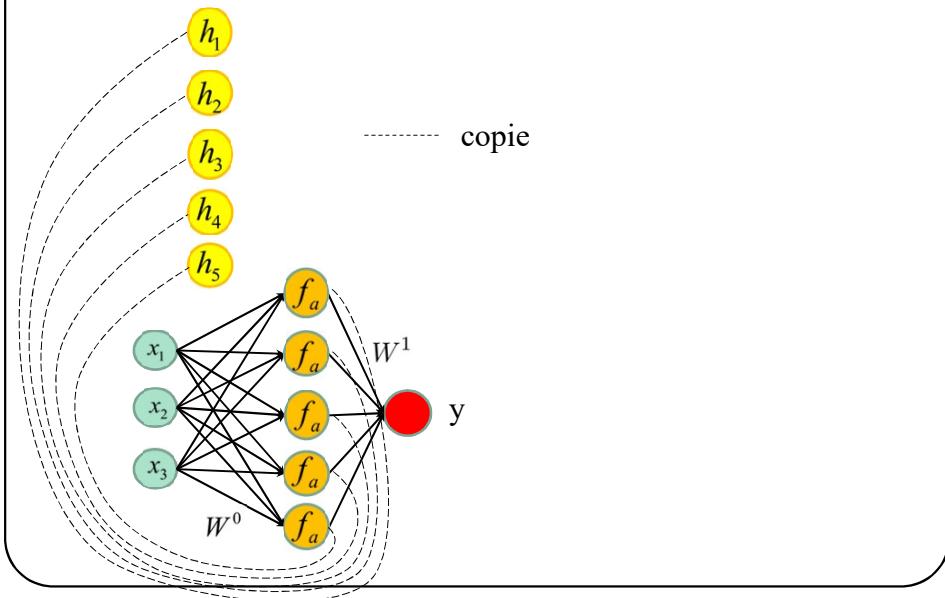


Ex.: description d'une image
1 image => N mots

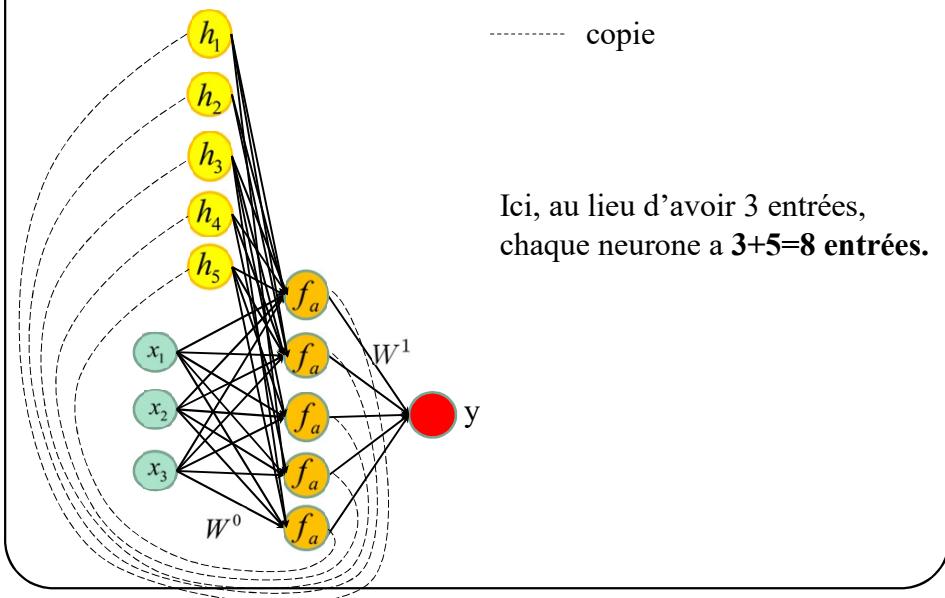




Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée



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



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

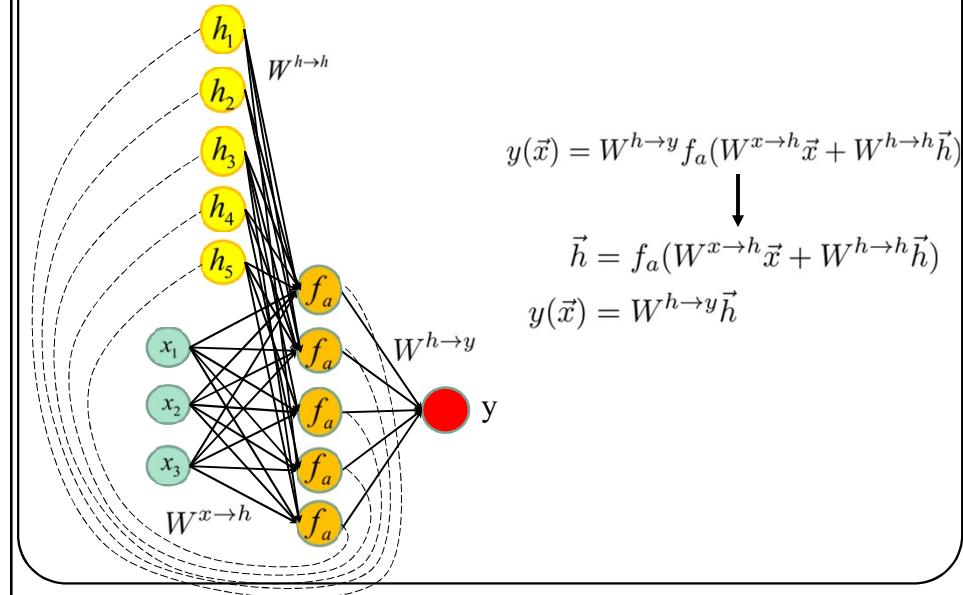
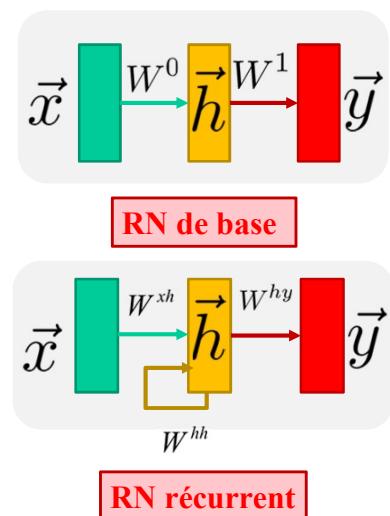
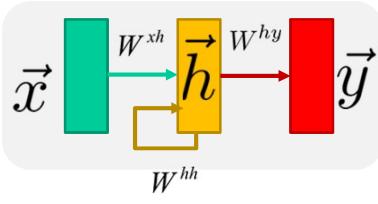


Illustration simplifiée

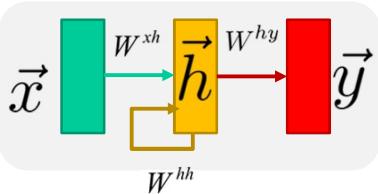


Dans le cas général avec K sorties (régression)



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

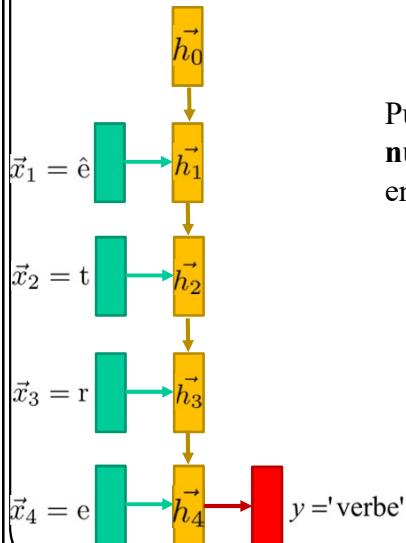
Dans le cas général avec K sorties (classification)



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

Exemple pour N entrées et 1 sortie:

Analyse grammaticale (classification) : (é.t.r.e)=>'verbe'



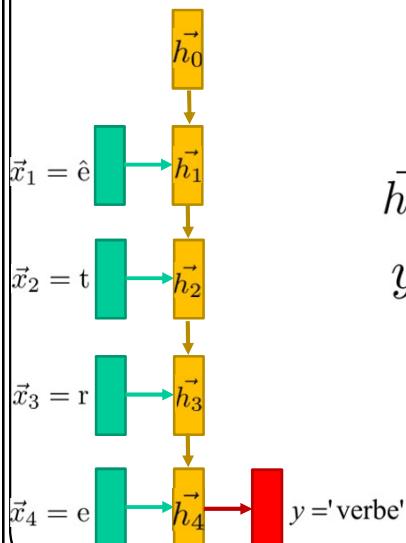
Puisque \vec{x} , \vec{h} et y doivent être des **variables numériques** on utilise souvent un encodage de type « *one hot* ».

$$\left. \begin{array}{l} 'a' = [1, 0, 0, \dots, 0] \\ 'b' = [0, 1, 0, \dots, 0] \\ 'c' = [0, 0, 1, \dots, 0] \\ \dots \\ 'verbe' = [1, 0, 0, \dots, 0] \\ 'nom' = [0, 1, 0, \dots, 0] \\ 'adjectif' = [0, 0, 1, \dots, 0] \end{array} \right\} \in R^{256}$$

$$\left. \begin{array}{l} 'a' = [1, 0, 0, \dots, 0] \\ 'b' = [0, 1, 0, \dots, 0] \\ 'c' = [0, 0, 1, \dots, 0] \\ \dots \\ 'verbe' = [1, 0, 0, \dots, 0] \\ 'nom' = [0, 1, 0, \dots, 0] \\ 'adjectif' = [0, 0, 1, \dots, 0] \end{array} \right\} \in R^M$$

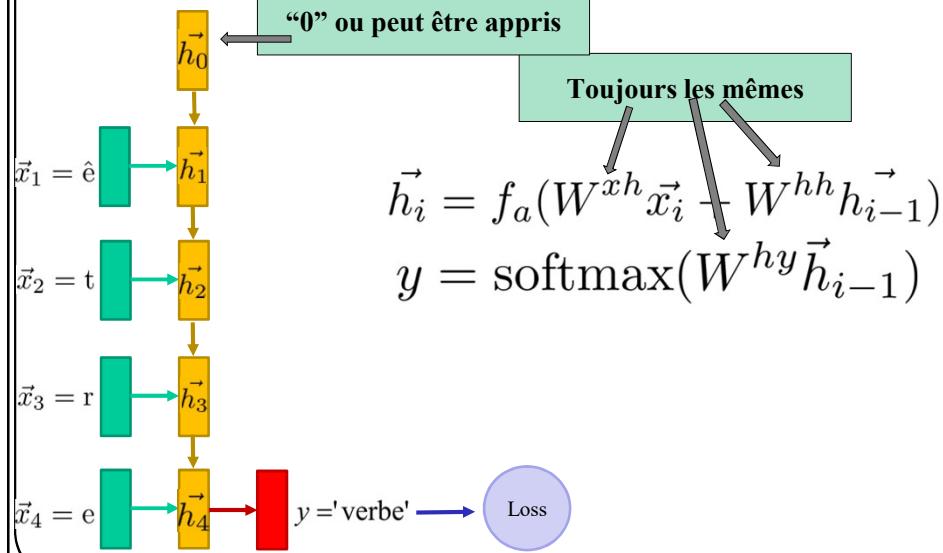
Exemple pour N entrées et 1 sortie:

Analyse grammaticale (classification) : (é.t.r.e)=>'verbe'



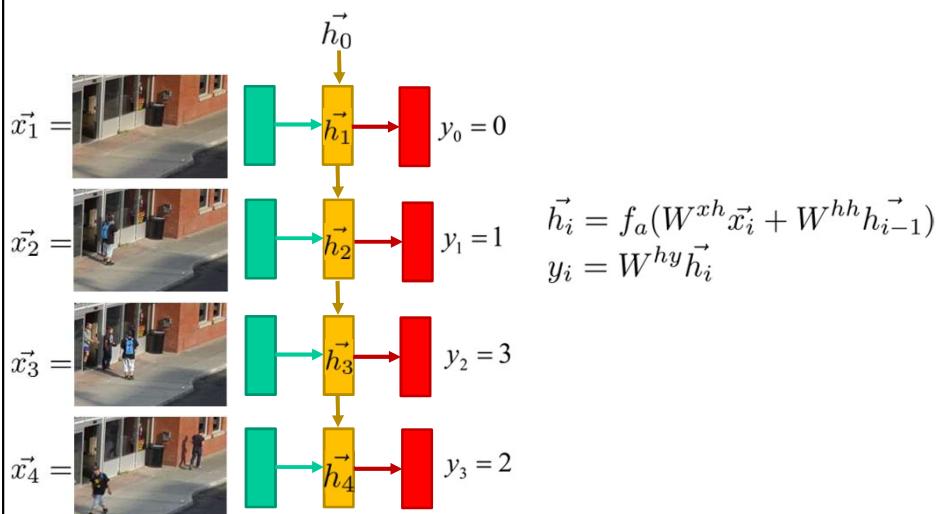
Exemple pour N entrées et 1 sortie:

Analyse grammaticale (classification) : (é.t.r.e)=>'verbe'

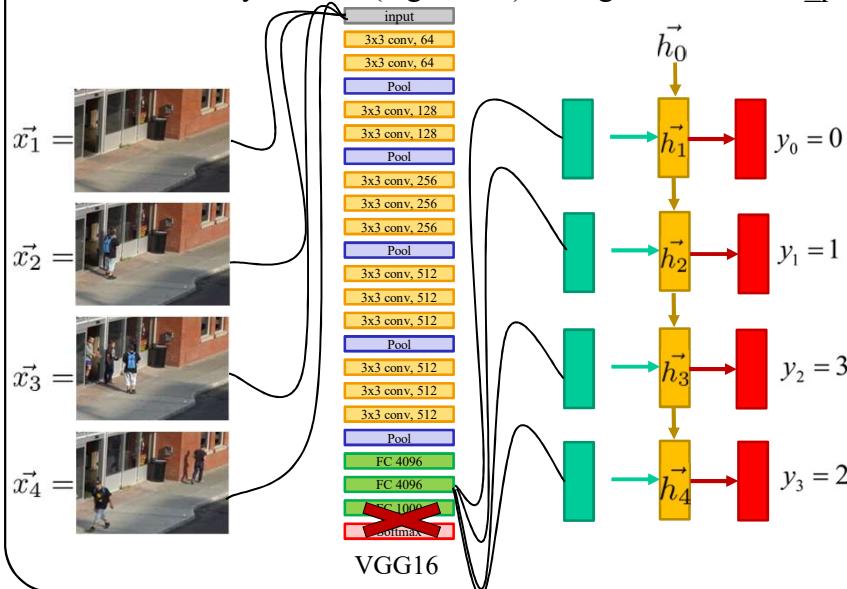


Même idée pour N entrées et N sorties:

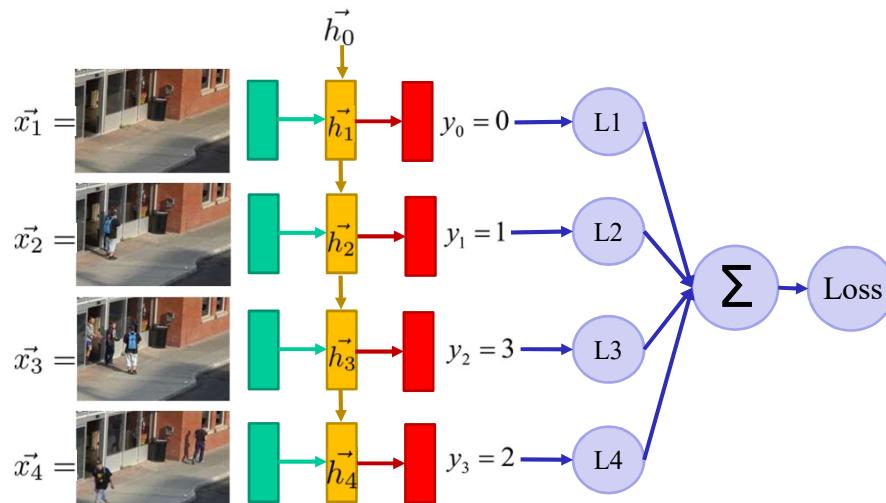
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



Même idée pour N entrées et N sorties:
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



Même idée pour N entrées et N sorties:
ex.: Analyse vidéo (régression) : Images vidéo => nb_piétons



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

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

Représentation « one hot » jouet:

$$\text{'a'} = [1, 0, 0, 0]$$

$$\text{'e'} = [0, 1, 0, 0]$$

$$\text{'m'} = [0, 0, 1, 0]$$

$$\text{'s'} = [0, 0, 0, 1]$$

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

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

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

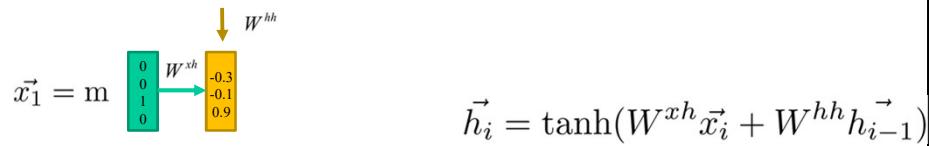
Entrainer un modèle à prédire les lettres du mot « **masse** ».

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

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

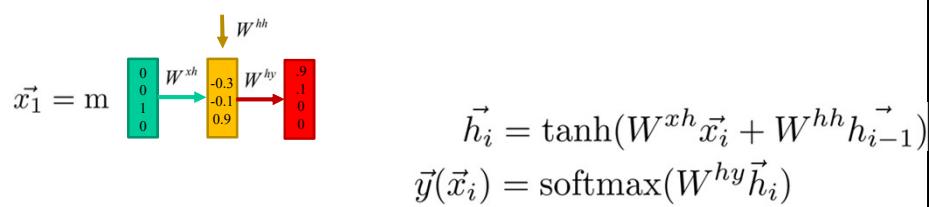
Entraîner un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

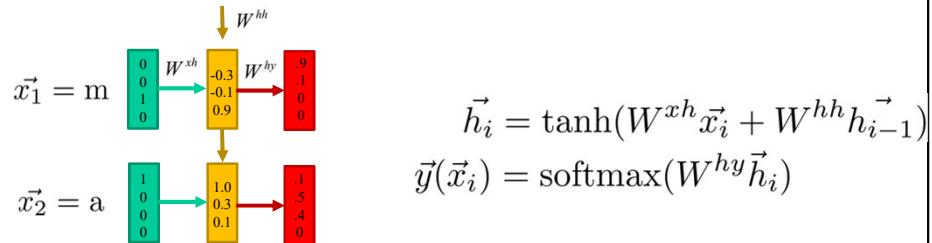
Entraîner un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

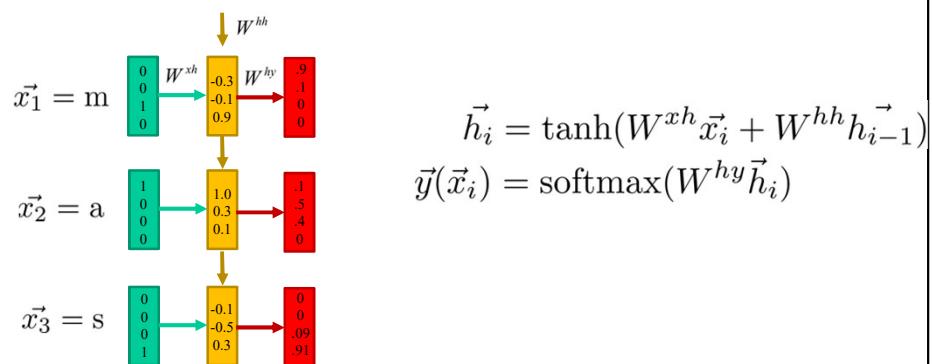
Entrainer un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

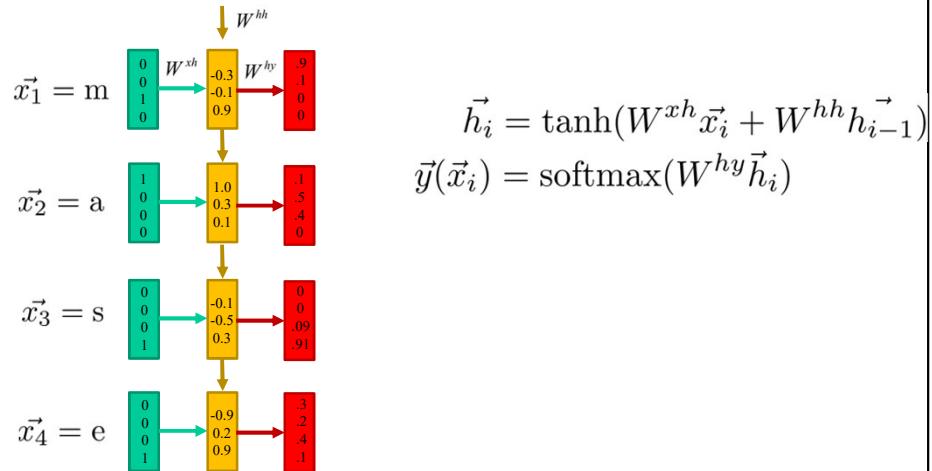
Entrainer un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

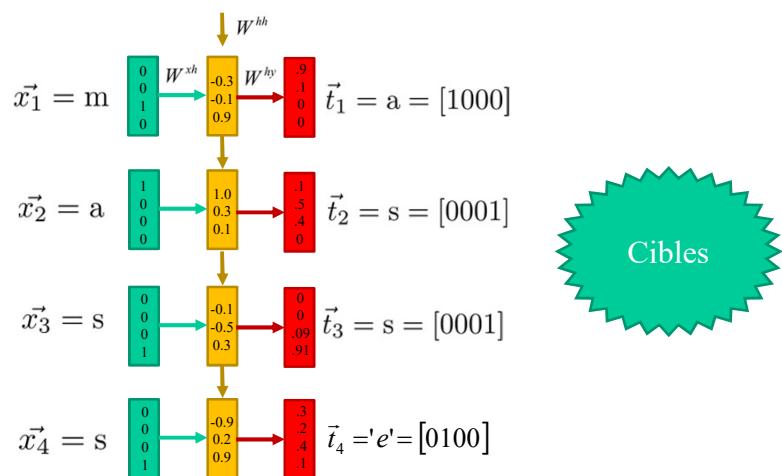
Entrainer un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

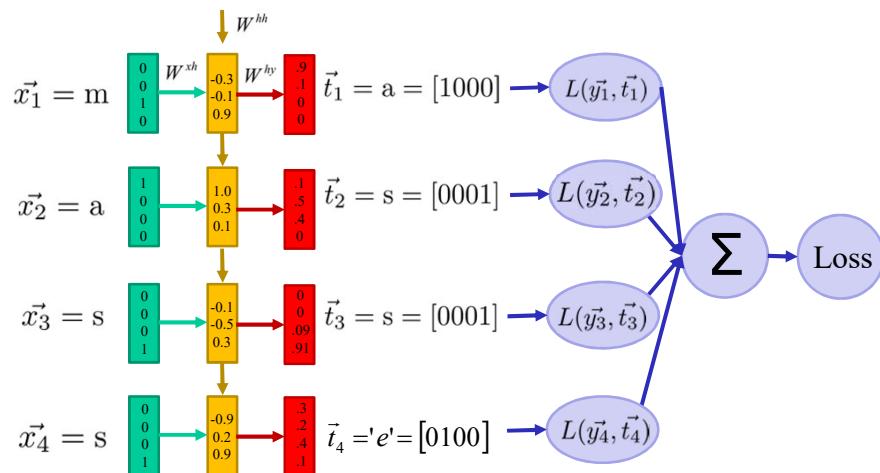
Entrainer un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédiction de caractères (modèle de langue)

Alphabet : [a,e,m,s]

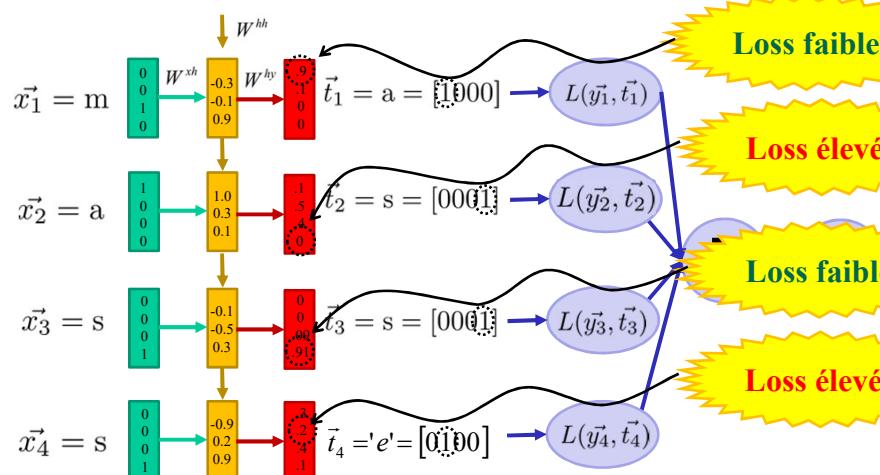
Entrainer un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédiction de caractères (modèle de langue)

Alphabet : [a,e,m,s]

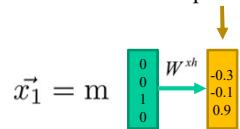
Entrainer un modèle à prédire les lettres du mot « **masse** ».



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

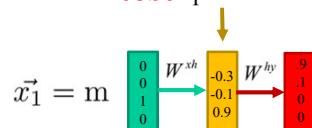


Étape 1 : Calcul de la couche cachée

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

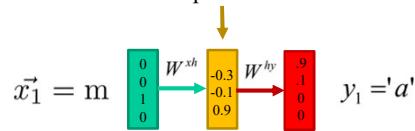


Étape 2 : Calcul de la sortie (softmax)

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

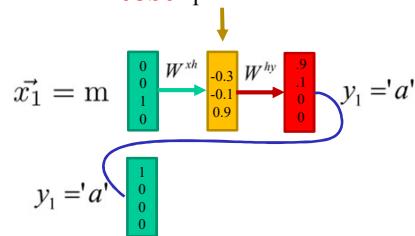


Étape 3 : Sélectionner le caractère le plus probable

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

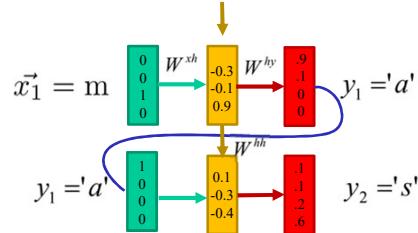


Étape 4 : Injecter le caractère prédit au début du réseau

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

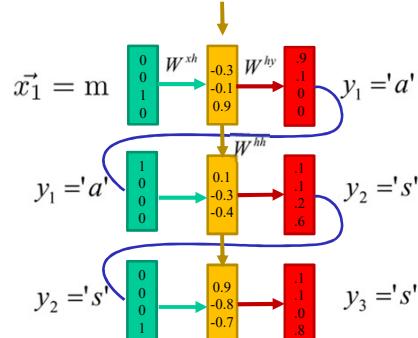


Et on recommence!

Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

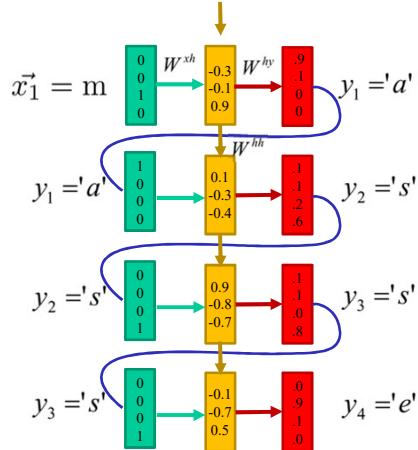
En test : prédire les lettres les unes après les autres



Autre exemple: prédition de caractères (modèle de langue)

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres



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
min-char-rnn.py
...
1 # Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
2 BSD License
3 ...
4 import numpy as np
5
6 # data I/O
7 data = open('input.txt', 'r').read() # should be simple plain text file
8 chars = list(data)
9 vocab_size = len(chars)
10 print('data has %d characters, %d unique.' % (data_size, vocab_size))
11 char_to_ix = { ch:i for i,ch in enumerate(chars) }
12 ix_to_char = { i:ch for i,ch in enumerate(chars) }
13
14 # hyperparameters
15 hidden_size = 200 # size of hidden layer of neurons
16 seq_length = 25 # number of steps to unroll the RNN for
17 learning_rate = 1e-1
18
19 # model parameters
20 Wih = np.random.rand(hidden_size, vocab_size)*0.01 # input to hidden
21 Whh = np.random.rand(hidden_size, hidden_size)*0.01 # hidden to hidden
22 Why = np.random.rand(vocab_size, hidden_size)*0.01 # hidden to output
23 bh = np.zeros(hidden_size, 1) # hidden bias
24 by = np.zeros((vocab_size, 1)) # output bias
25
26 def lossFun(inputs, targets, hprev):
27     """ Inputs, targets are both list of integers.
28     hprev is Hx array of initial hidden state
29     returns the loss, gradients on model parameters, and last hidden state
30     """
31     xs, hs, ys, ps = [], [], [], []
32     hprev = np.copy(hprev)
33     for t in range(len(inputs)):
34         x = inputs[t]
35         xenc = np.zeros((vocab_size, 1)) # encode in 1-of-k representation
36         xenc[[t]] = 1
37
38         hs.append(hprev)
39         hprev = np.tanh(np.dot(Wih, xenc) + np.dot(Whh, hprev) + bh)
40
41         ys.append(targets[t])
42         ps.append(ps)
43
44         loss += np.square(ys[t] - ps[t])
45
46     return loss, hs, ys, ps
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
```

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

...

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,
 That thereby beauty's rose might never die,
 But as the world is pass'd by every decease,
 His ende her might his last enemy,
 But thou, contracted to thine own bright eyes,
 Feed'st thy light's flame with self-substantial fuel,
 Making a famine where abundance lies,
 Thy self thy foe, to thy sweet self too cruel,
 Thou that art rose the world's fresh ornament,
 And only herald to the gaudy spring,
 Within thine own bad burkest thy content,
 And tender churl mak'st waste in niggarding:
 Pay the world, or else this glutton be,
 To eat the world's due, by the grave and thee.

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

...

When forty winters shall beside thy brow,
 And dig deep trenches in thy beauty's field,
 Thy youth's proud livery so garded now,
 Will be a tatter'd weed of small worth held:
 Then being gone, where thy beauty lies,
 Thou art the treacherous visage of thy days,
 To say, within thine own deep sunken eyes,
 Were an all-eating shame, and thriftless praise.
 How man's more praise deserve'd thy beauty's use,
 If thy more worthier sonnes That time ent'ren
 Shall sum my eame, and make my old excuse,
 Proving his beauty by succession thine!
 There were to be new made when thou art old,
 And see thy blood warm when thou feel'st it cold.

Autre exemple: prédition de caractères (modèle de langue)

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

↓ train more

"Tmont thithey" fomesscerliund
 Keushey. Thom here
 sheulke, ammerenith 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 belege, 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

Autre exemple: prédition de caractères (modèle de langue)

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 apes, 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

Autre exemple: prédition de caractères (modèle de langue)

Entraînement sur le code source de Linux en C++

Texte généré une fois le modèle entraîné

```

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 & 0x00000000fffffff0) & 0x0000000f) << 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 = &of_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 ");
}

```

```

#define <asm/io.h>
#define <asm/prom.h>
#define <asm/e20.h>
#define <asm/system_info.h>
#define <asm/estew.h>
#define <asm/ppproto.h>

#define REQ_FQ      vesa_slot_addr_pack
#define FPN_NOCOMP  AFSR(0, load)
#define STICK_DBR(type)  (func)

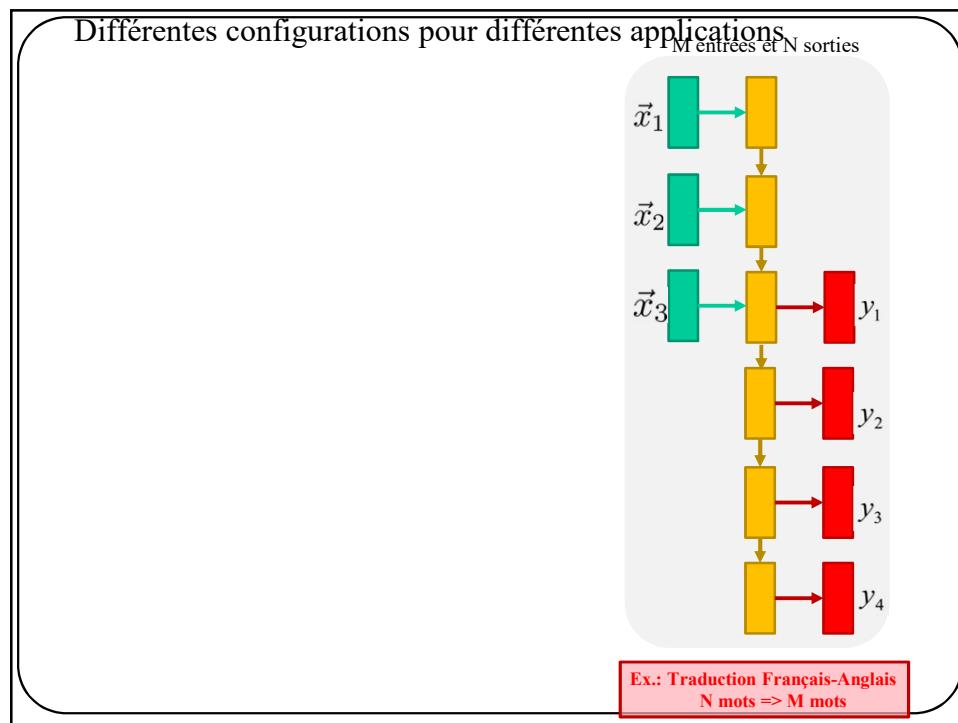
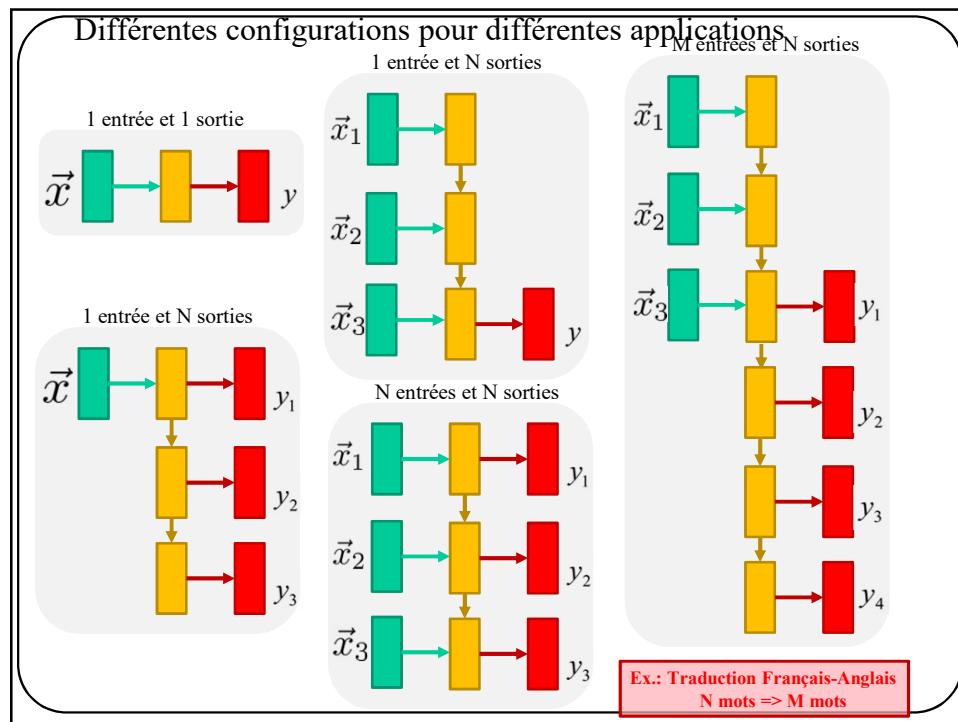
#define SMP_ALLOCATE(nr)      (a)
#define emulatc_size() arch_get_unaligned_child()
#define access_rw(TST)  __asm volatile("movd %esp, %0; t3" : "+r" (0)) \
    if (__TYPE & __D0_READ)

static void stat_pc_sec __read_mostly offsetof(struct seq_argqueue,
pC>[1]));

static void
os_prefix(unsigned long sys)
{
    #ifndef CONFIG_PRESUME
    PUT_PARAM_RAID(2, sel) = get_state_state();
    set_pid_num((unsigned long)state, current_state_str(),
                (unsigned long)-1>lr_full; low;
    }

```

Crédit: A. Karpathy, CS231



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 : Begining of Sentence, EoS:End of Sentence).

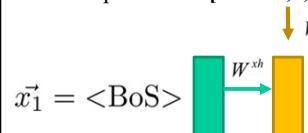
Autre exemple: traduction

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>]



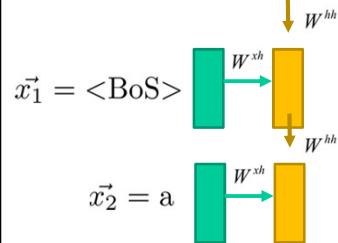
Autre exemple: traduction

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>]



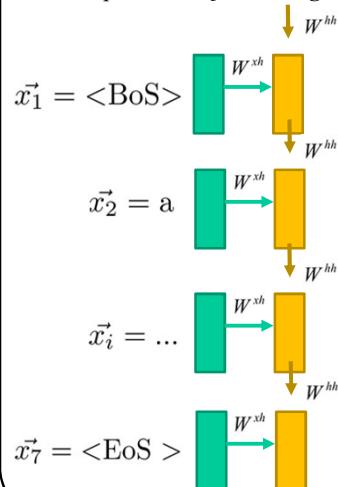
Autre exemple: traduction

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’

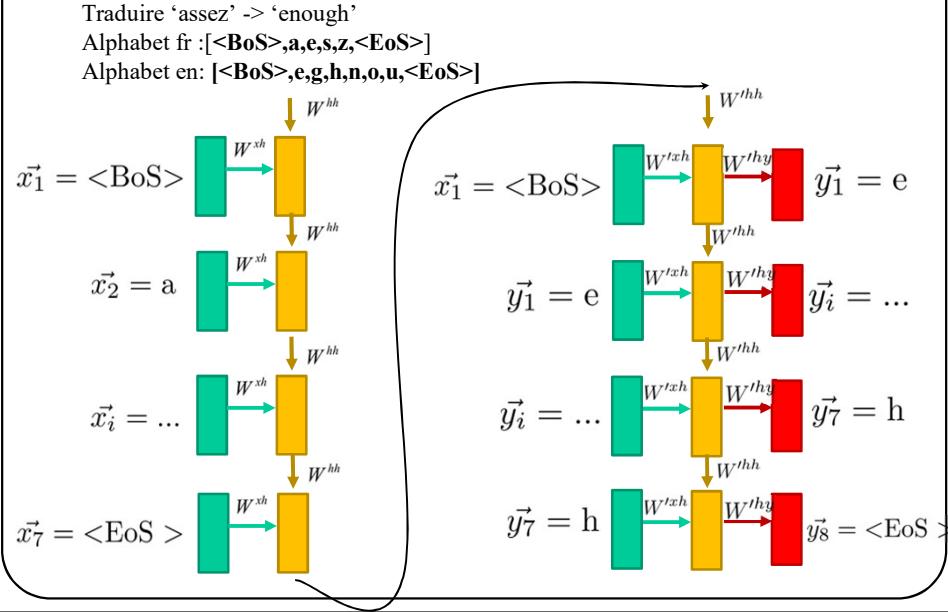
Alphabet fr :[<BoS>,a,e,s,z,<EoS>]

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



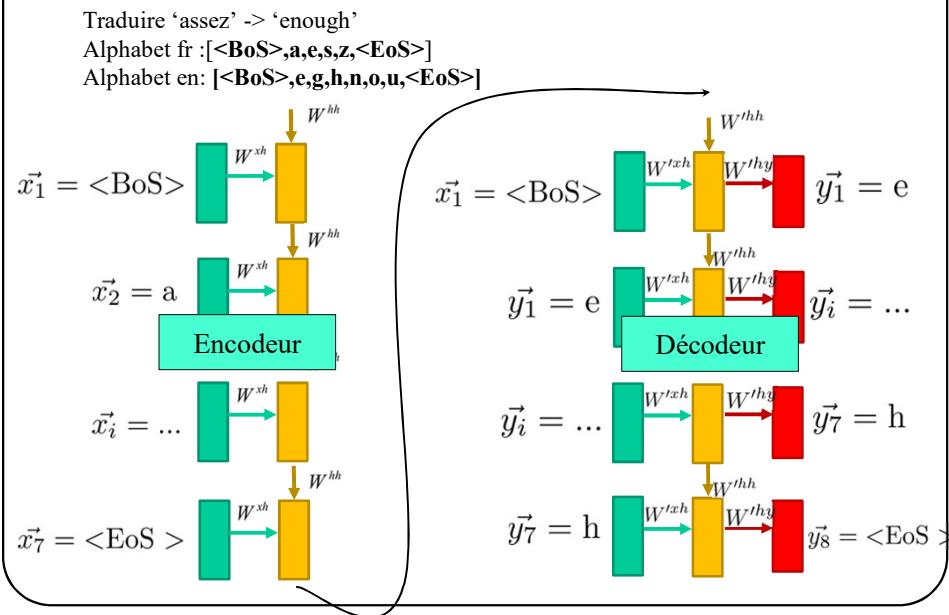
Autre exemple: traduction

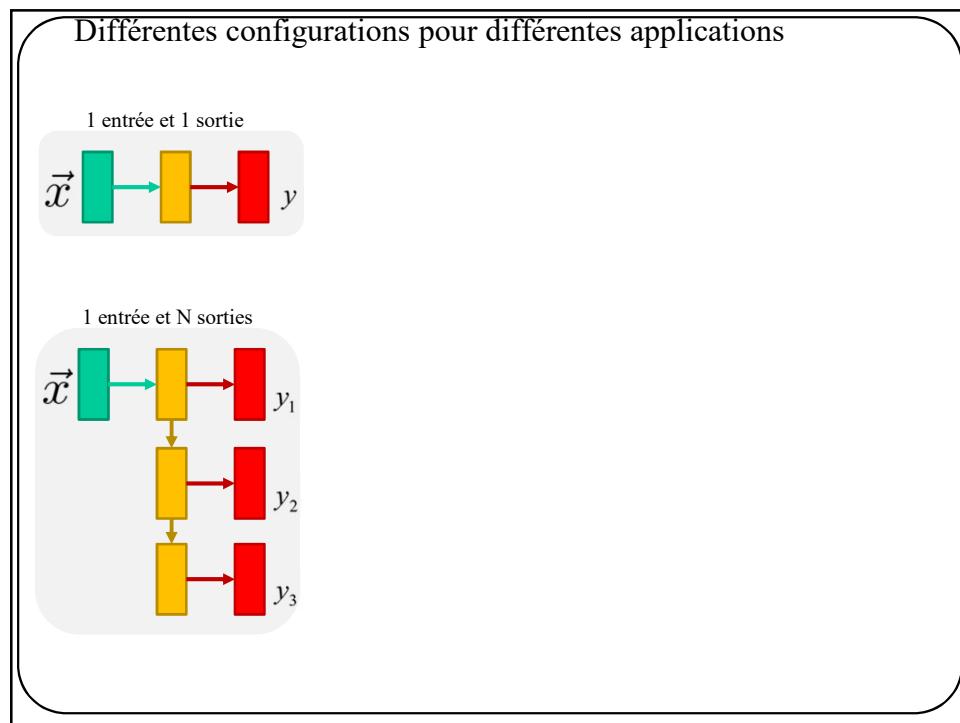
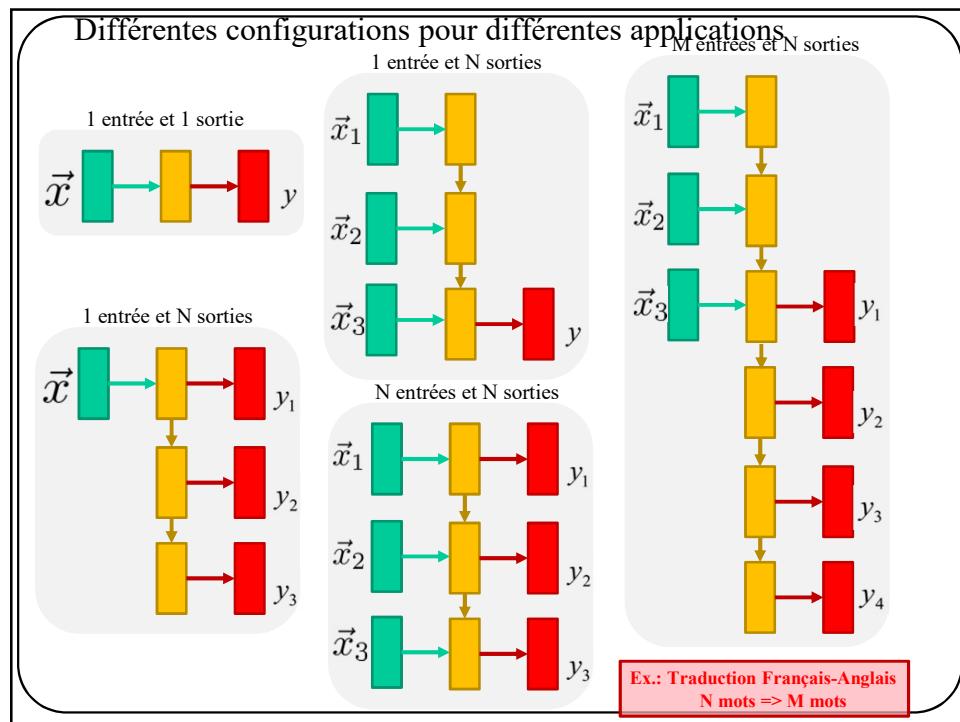
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

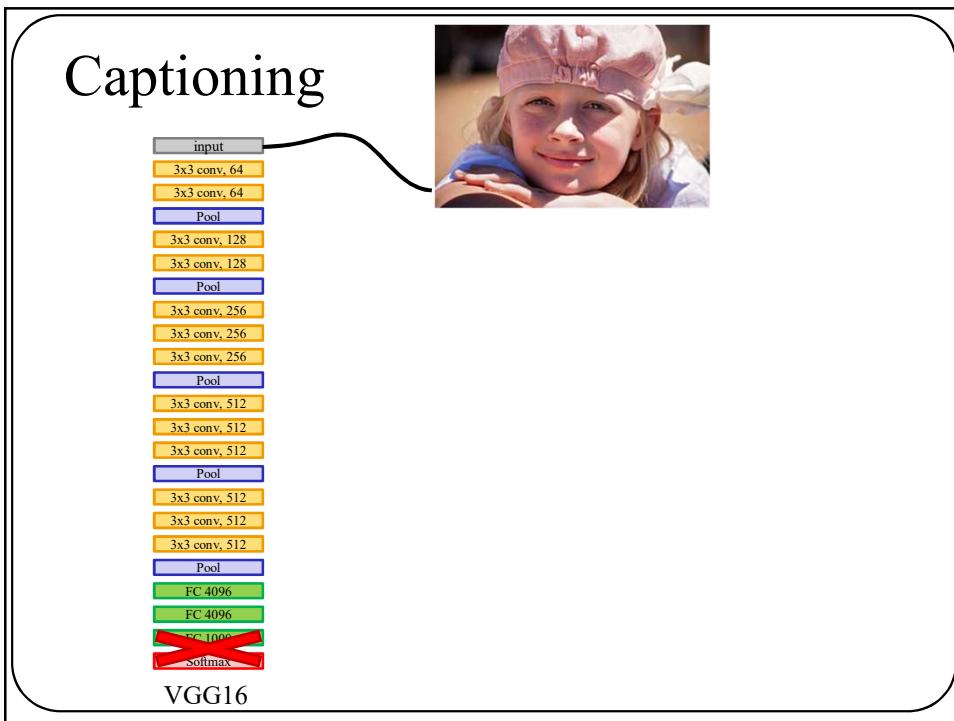
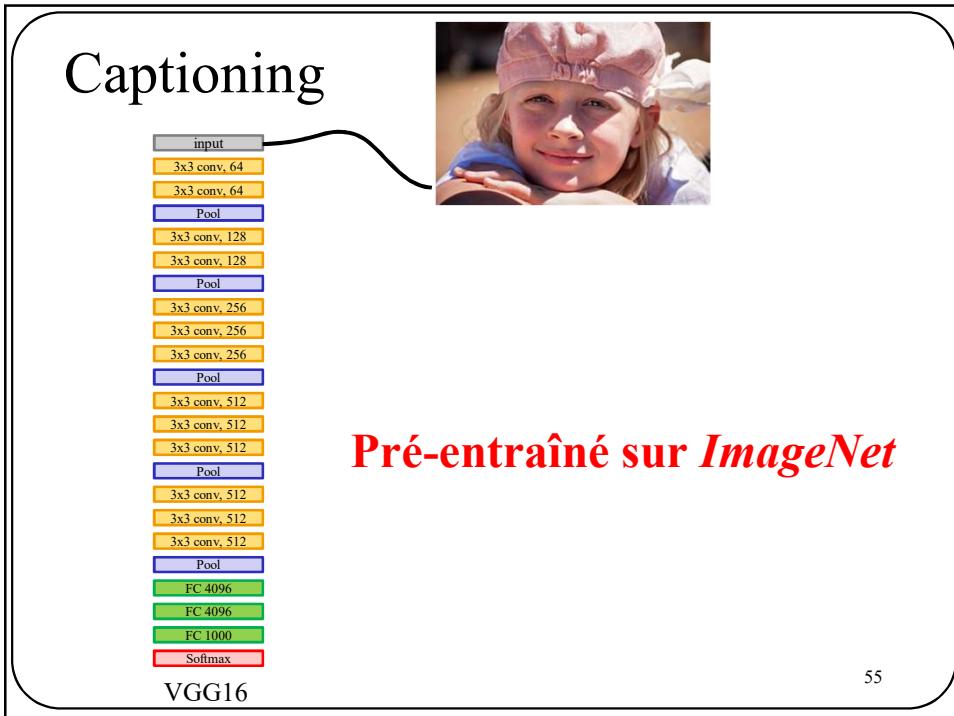


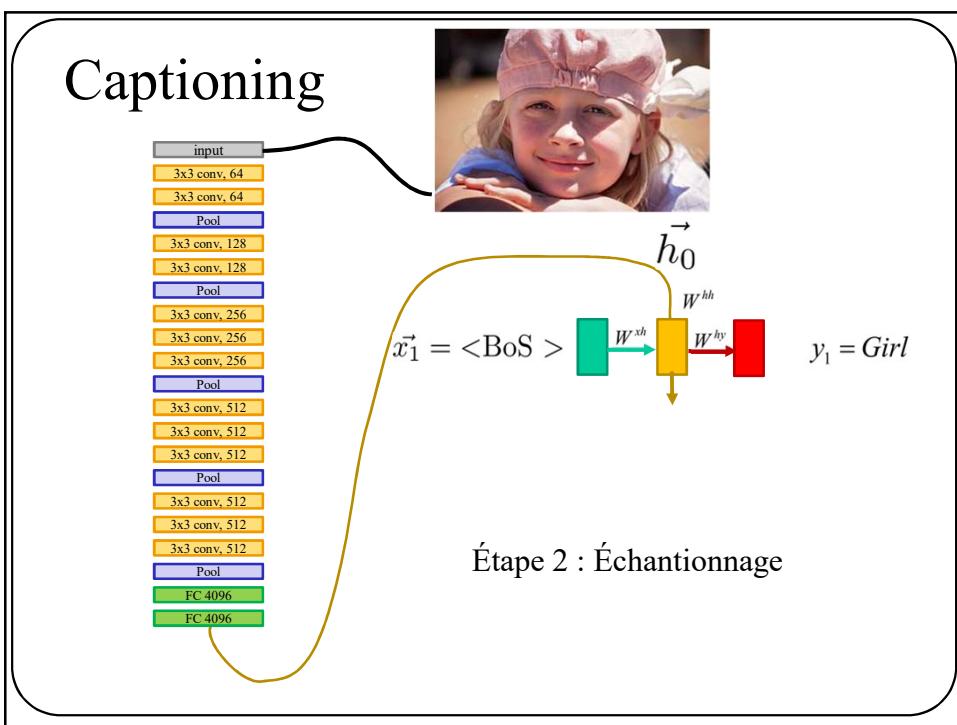
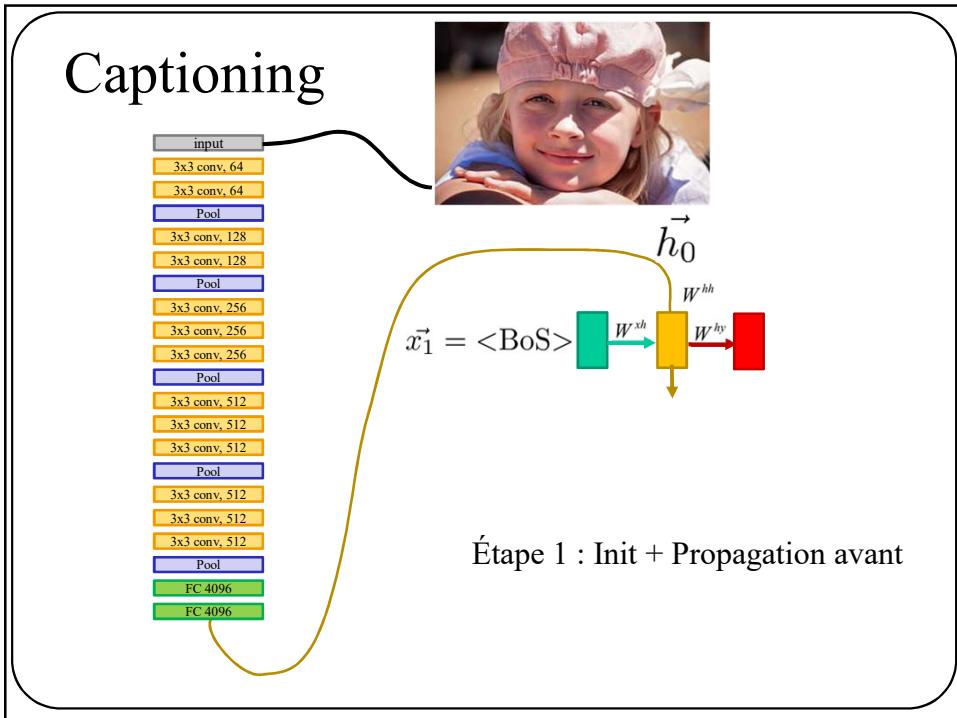
Autre exemple: traduction

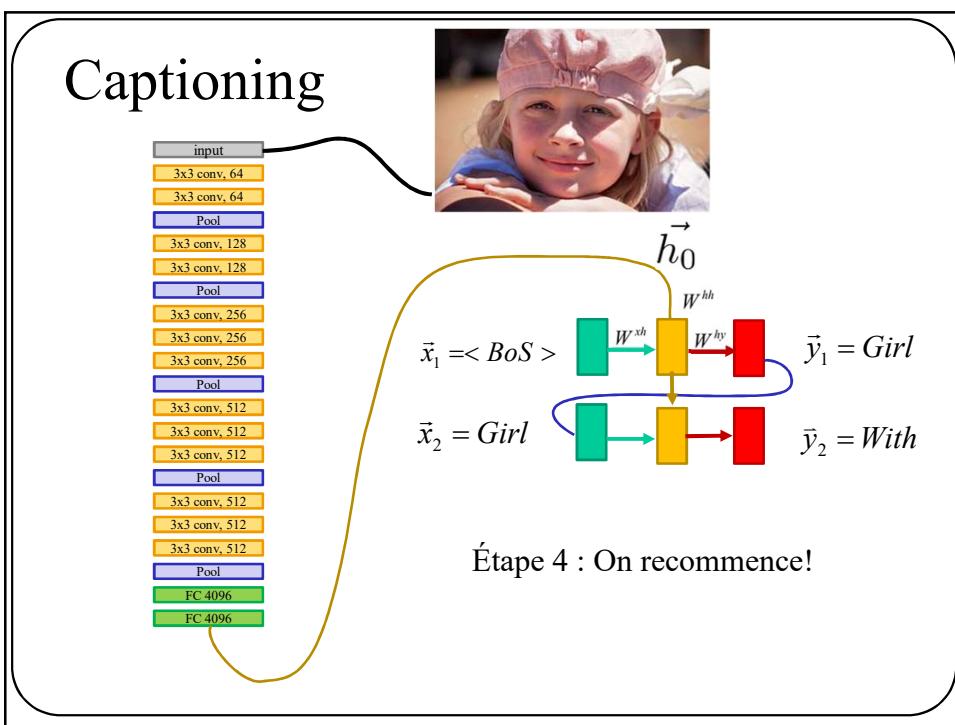
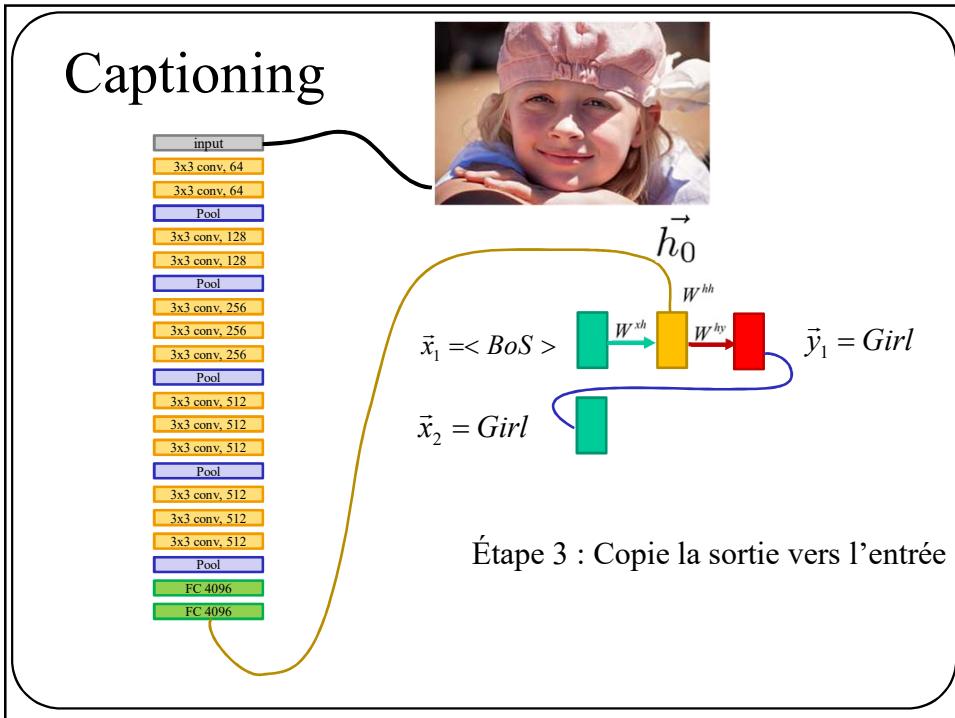
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

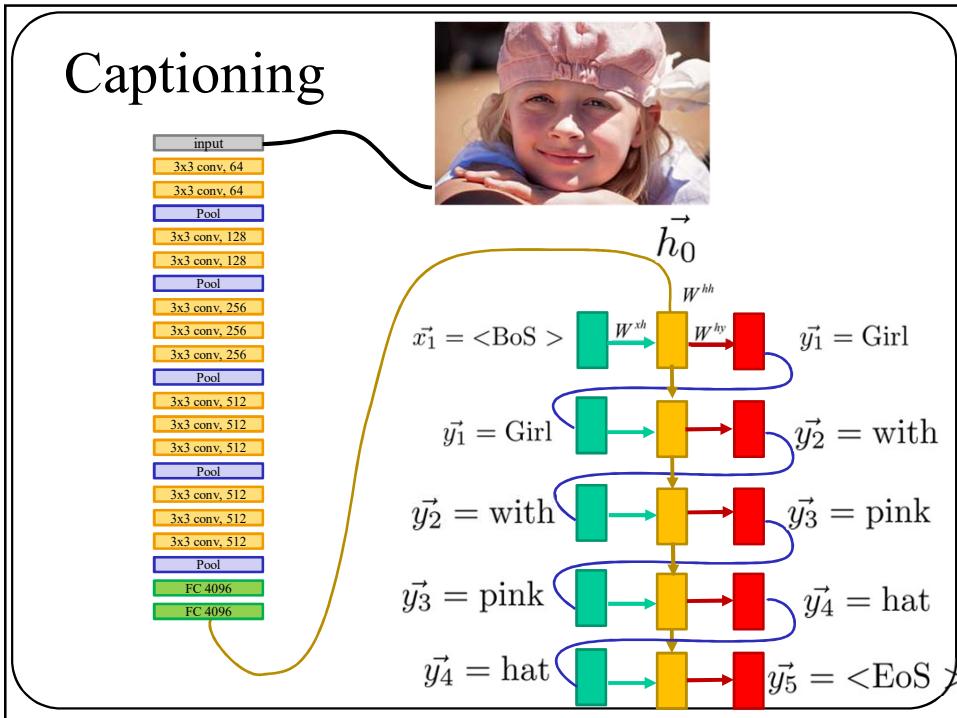












Exemples de résultats

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



an 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 wii controllers



a couple of zebra standing on top of a dirt field

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

63

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



a city street with a bridge and a bridge in the background



a display case filled with lots of different types of food

64

Analyse de texte

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

Pour des **caractères**...

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

65

Analyse de texte

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

Pour des **mots**...

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

66

Prédiction sur des lettres vs. mots

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

...

...

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

...

Prédiction sur des fractions de mots

$$\begin{aligned} 'e' &= [0, 0, \dots, 1, \dots, 0] \\ 'grand' &= [0, 0, \dots, 1, \dots, 0] \\ 'ment' &= [0, 0, \dots, 1, \dots, 0] \end{aligned} \quad \left. \right\} \in \mathbb{R}^m \quad \begin{array}{l} \dots \\ \dots \\ \dots \\ \dots \end{array} \quad \begin{array}{l} 'grand' \\ 'grand'+'e' \\ 'grand'+'e'+'ment' \end{array}$$

Tokenization (jeton-isation ?)

Idée: à partir d'un dictionnaire qui ne contient que des caractères, combiner les séquences fréquentes en jetons (*tokens*)

Les séquences fréquentes (comme les mots ou sous-mots fréquents) se voient attribuer un jeton. Les séquences peu fréquentes peuvent être bâties à partir de jetons.

Sennrich, R., Haddow, B., & Birch, A. (2015). Neural machine translation of rare words with subword units. arXiv preprint arXiv:1508.07909.

Algorithm 1 Learn BPE operations

```
import re, collections

def get_stats(vocab):
    pairs = collections.defaultdict(int)
    for word, freq in vocab.items():
        symbols = word.split()
        for i in range(len(symbols)-1):
            pairs[symbols[i],symbols[i+1]] += freq
    return pairs

def merge_vocab(pair, v_in):
    v_out = {}
    bigram = re.escape(' '.join(pair))
    p = re.compile(r'(%<!\S)' + bigram + r'(%!?\S)')
    for word in v_in:
        w_out = p.sub('', join(pair), word)
        v_out[w_out] = v_in[word]
    return v_out

vocab = {'l o w </w>': 5, 'l o w e r </w>': 2,
         'n e w e s t </w>': 6, 'w i d e s t </w>': 3}
num_merges = 10
for i in range(num_merges):
    pairs = get_stats(vocab)
    best = max(pairs, key=pairs.get)
    vocab = merge_vocab(best, vocab)
    print(best)
```

r·	→	r·
lo	→	lo
lo w	→	low
e r·	→	er·

Limites des « one-hot vectors »

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

1- Peu efficace en mémoire lorsque non compressé

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

2- Pas de distance sémantique entre les codes:

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')]

Word2Vec s'appuie sur 2 idées fondamentales

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

Idée 1: Dictionnaire = matrice d'encodage

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

	« one-hot »							Dictionnaire
‘the’	1	0	0	0	0	0	0	2 3 4 5
‘quick’	0	1	0	0	0	0	0	-1 -3 -2 2
‘brown’	0	0	1	0	0	0	0	11 6 4 -3
‘fox’	0	0	0	1	0	0	0	-4 8 -4 4
‘jumps’	0	0	0	0	1	0	0	24 -6 42 17
‘over’	0	0	0	0	0	1	0	91 13 14 -5
‘lazy’	0	0	0	0	0	0	1	0 36 4 56
‘dog’	0	0	0	0	0	0	1	-1 0 1 35



Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

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

Ex: sélectionner le code de « brown »

$$\begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} \text{Dictionnaire} \\ (\text{matrice d'encodage}) \end{pmatrix} = \begin{pmatrix} 11 & 6 & 4 & -3 \end{pmatrix}$$

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

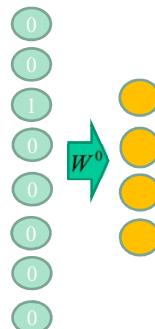
Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

Première couche d'un réseau de neurones

=
matrice d'encodage

\vec{x} : brown



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



Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

Première couche d'un réseau de neurones

=
matrice d'encodage

$$code_{\vec{x}} = W^0 \vec{x}$$



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

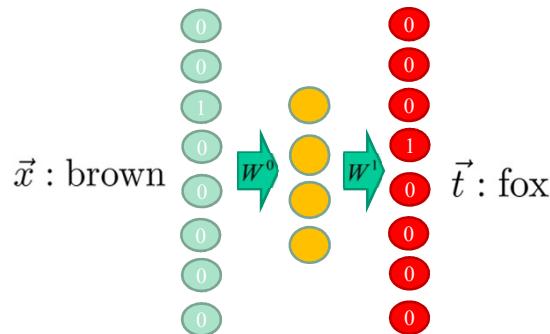
Word2Vec s'appuie sur 2 idées fondamentales

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)

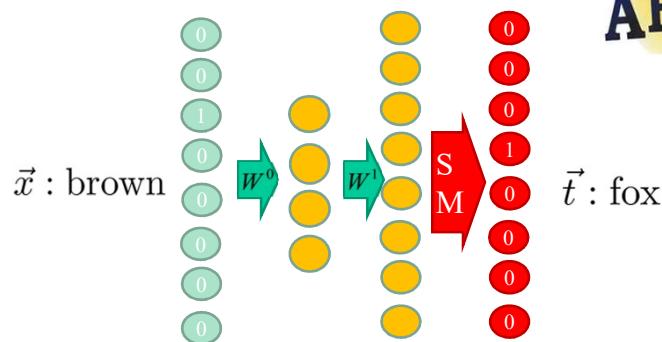
Basé sur un corpus de texte, on va créer des **millions de paires de mots**

Word2Vec [Mikolov et al. '13]



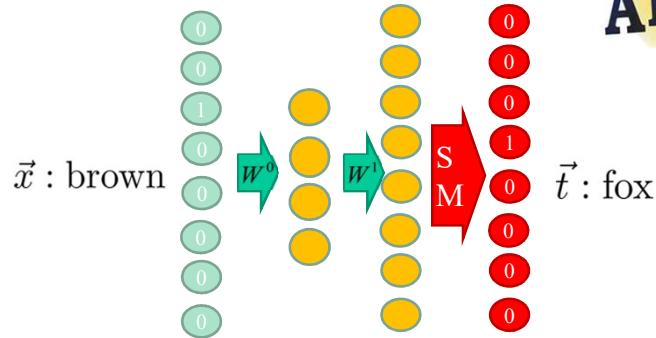
Entraîner un réseau de neurones
à reproduire le 2^e mot partant du 1^{er}

Word2Vec [Mikolov et al. '13]



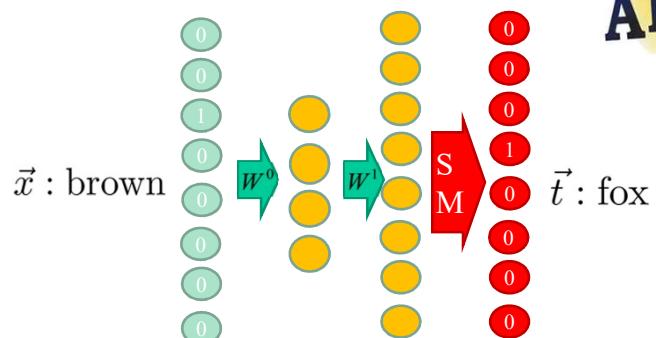
Puisque la sortie est de type « *one-hot* »
on utilise un softmax

Word2Vec [Mikolov et al. '13]



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

Word2Vec [Mikolov et al. '13]



Lorsqu'entraîné, utiliser W^0
comme dictionnaire

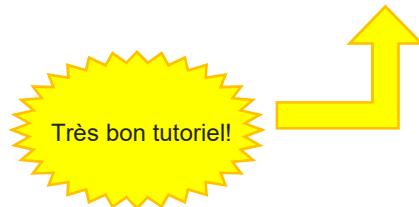
Word2Vec [Mikolov et al. '13]

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

Limites du « one-hot vector »

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

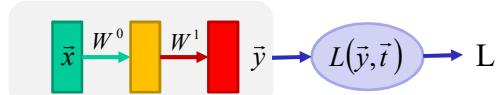


T.Mikolov et al. (2013). "Efficient Estimation of Word Representations in Vector Space", in ICLR 2013

Comment entraîner un RNN?

Histoire de gradients

RN de classification avec entropie croisée

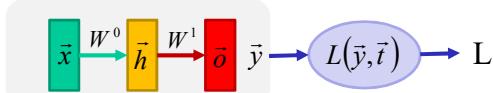


$$\vec{y}(\vec{x}) = S_M(W^1 \tanh(W^0 \vec{x}))$$

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

Histoire de gradients

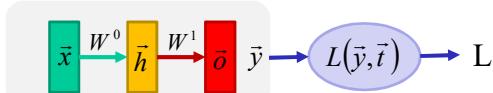
Simple RN de classification avec entropie croisée



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

Histoire de gradients

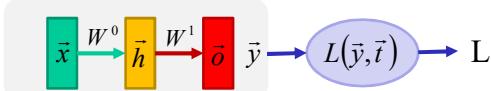
Simple RN de classification avec entropie croisée



$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \begin{array}{l} \text{Pour entraîner le réseau} \\ \text{il faut calculer} \\ \nabla_{W^0} L \text{ et } \nabla_{W^1} L \end{array}$$

Histoire de gradients

Simple RN de classification avec entropie croisée

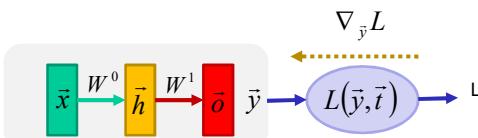


$$\begin{aligned}\vec{h} &= \tanh(W^0 \vec{x}) \\ \vec{o} &= W^1 \vec{h} \\ \vec{y} &= S_M(\vec{o}) \\ L &= L_{CE}(\vec{y}, \vec{t})\end{aligned}$$

Dérivée en chaîne

$$\begin{aligned}\nabla_{W^1} L &= \nabla_{\vec{y}} L \nabla_{\vec{o}} \vec{y} \nabla_{W^1} \vec{o} \\ \nabla_{W^0} L &= \nabla_{\vec{y}} L \nabla_{\vec{o}} \vec{y} \nabla_{\vec{h}} \vec{o} \nabla_{W^0} \vec{h}\end{aligned}$$

Histoire de gradients

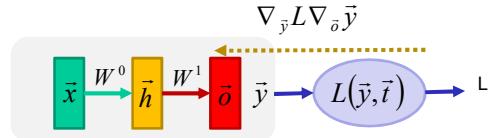


$$\begin{aligned}\vec{h} &= \tanh(W^0 \vec{x}) \\ \vec{o} &= W^1 \vec{h} \\ \vec{y} &= S_M(\vec{o}) \\ L &= L_{CE}(\vec{y}, \vec{t})\end{aligned}$$

Rétro-propagation

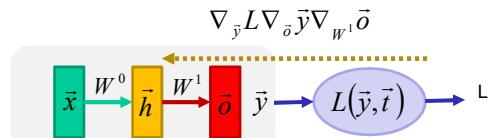
$$\nabla_{\vec{y}} L = -\frac{\vec{t}}{\vec{y}}$$

Histoire de gradients



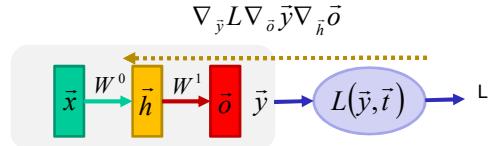
$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \begin{matrix} \uparrow \\ \downarrow \\ \text{R\'etro-propagation} \end{matrix} \quad \begin{aligned}
 \nabla_{\vec{o}} \vec{y} &= I \vec{y}^T - \vec{y}^T \vec{y} \\
 \nabla_{\vec{y}} L &= -\frac{\vec{t}}{\vec{y}}
 \end{aligned}$$

Histoire de gradients



$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \begin{matrix} \uparrow \\ \downarrow \\ \text{R\'etro-propagation} \end{matrix} \quad \begin{aligned}
 \nabla_{W^1} \vec{o} &= \vec{h} \\
 \nabla_{\vec{o}} \vec{y} &= I \vec{y}^T - \vec{y}^T \vec{y} \\
 \nabla_{\vec{y}} L &= -\frac{\vec{t}}{\vec{y}}
 \end{aligned}$$

Histoire de gradients

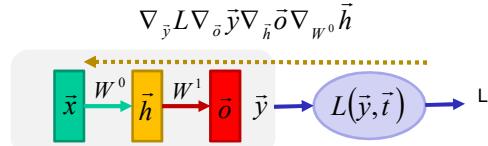


$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \uparrow$$

Rétro-propagation

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

Histoire de gradients

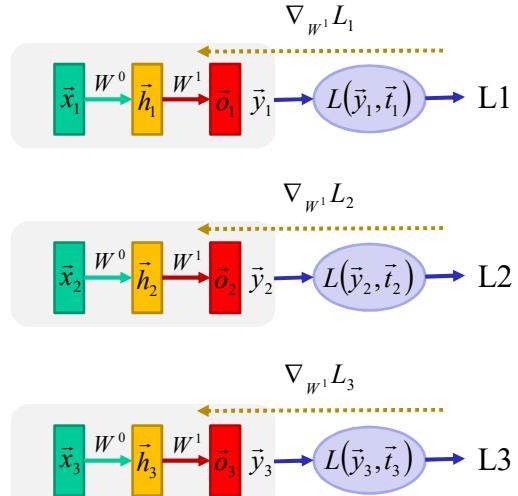


$$\begin{aligned}
 \vec{h} &= \tanh(W^0 \vec{x}) \\
 \vec{o} &= W^1 \vec{h} \\
 \vec{y} &= S_M(\vec{o}) \\
 L &= L_{CE}(\vec{y}, \vec{t})
 \end{aligned}
 \quad \downarrow \quad \uparrow$$

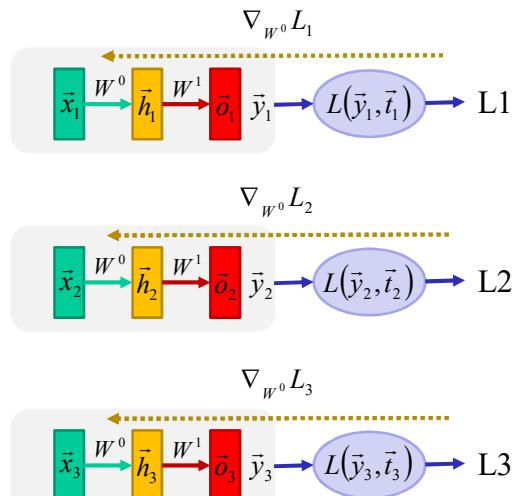
Rétro-propagation

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

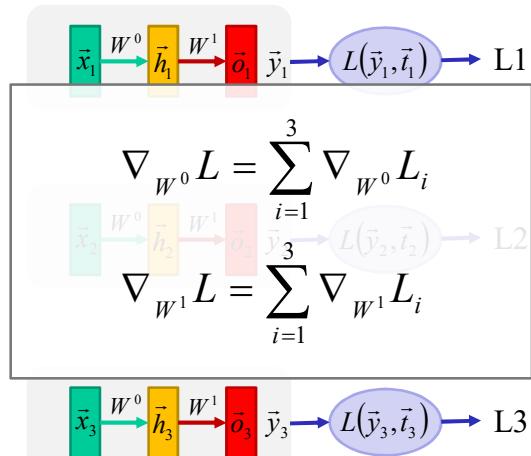
Ex.: 3 données, 3 rétro-propagations



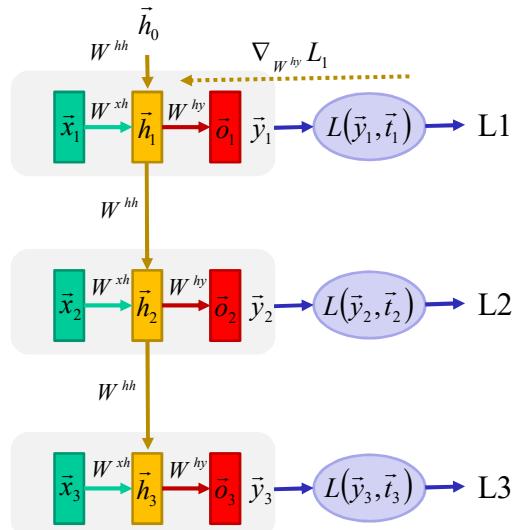
Ex.: 3 données, 3 rétro-propagations



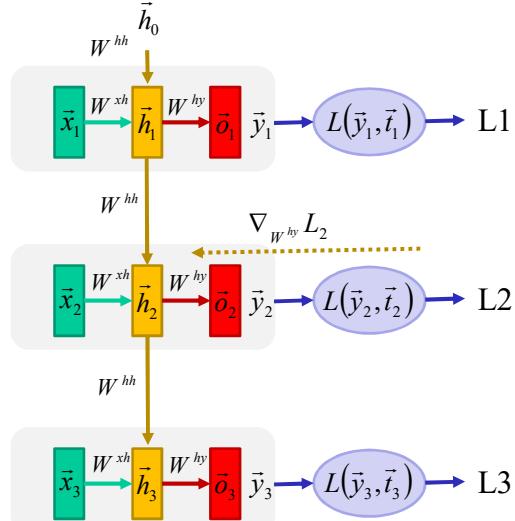
3 rétro-propagations



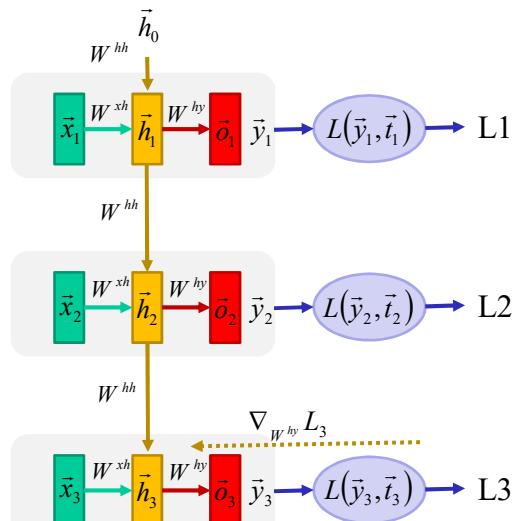
Réseau récurrent: gradient pour W^{hy}



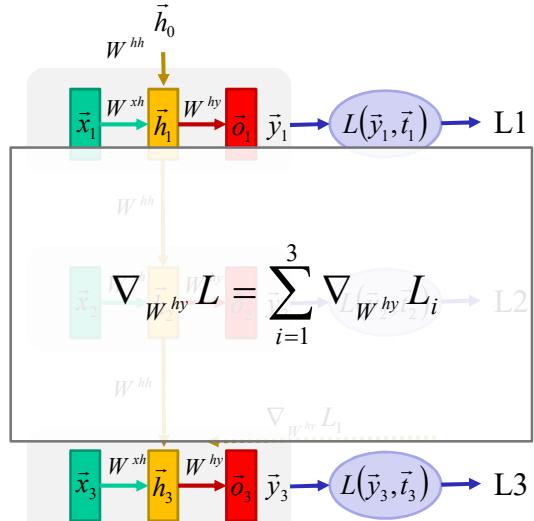
Réseau récurrent: gradient pour W^{hy}



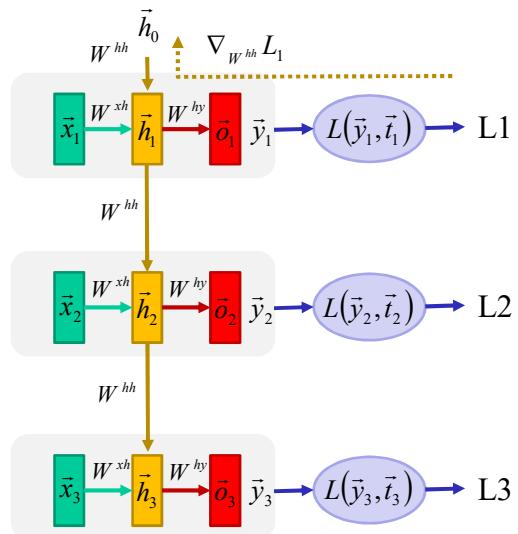
Réseau récurrent: gradient pour W^{hy}



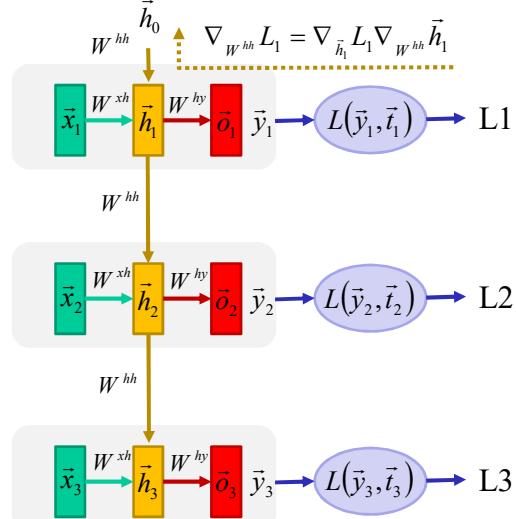
Réseau récurrent: gradient pour W^{hy}



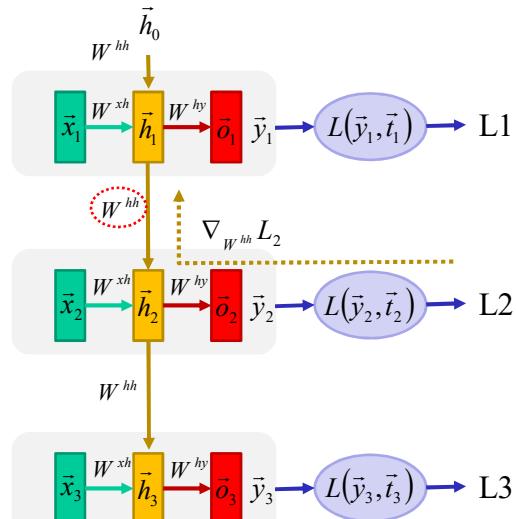
Réseau récurrent: gradient pour W^{hh}



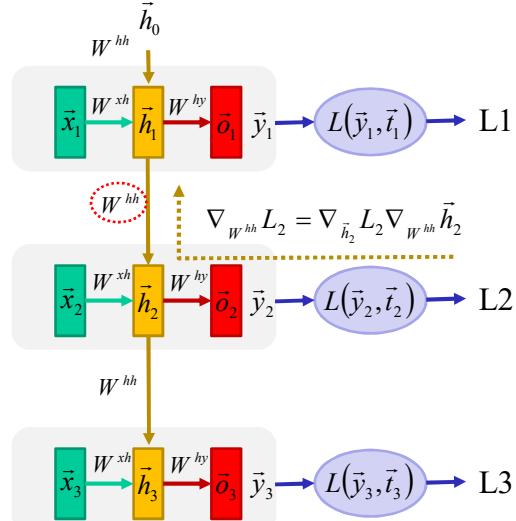
Réseau récurrent: gradient pour W^{hh}



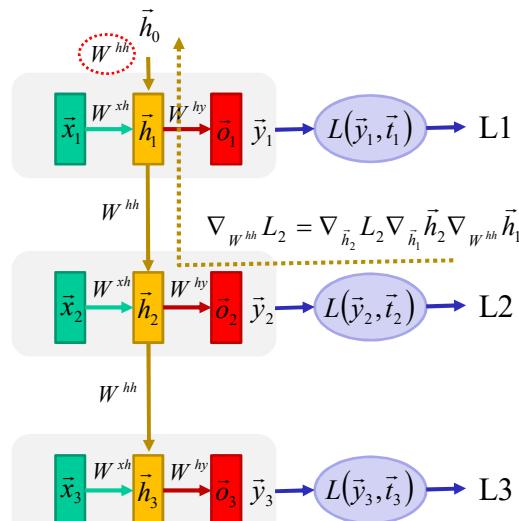
Réseau récurrent: gradient pour W^{hh}



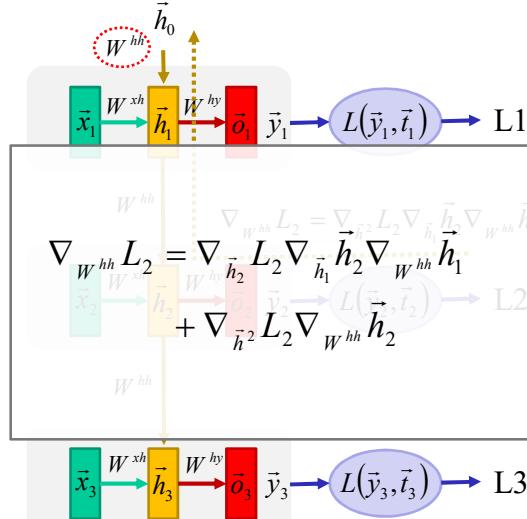
Réseau récurrent: gradient pour W^{hh}



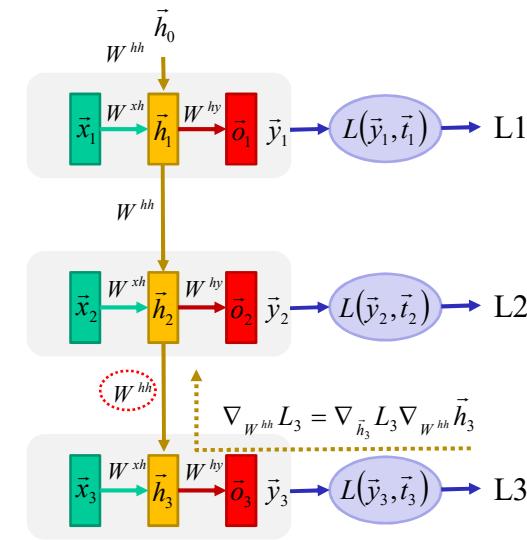
Réseau récurrent: gradient pour W^{hh}



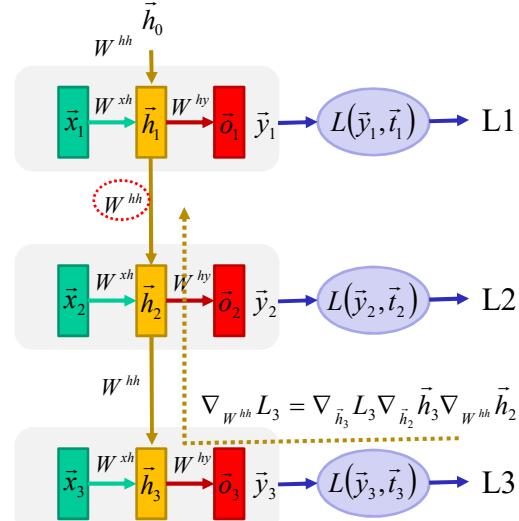
Réseau récurrent: gradient pour W^{hh}



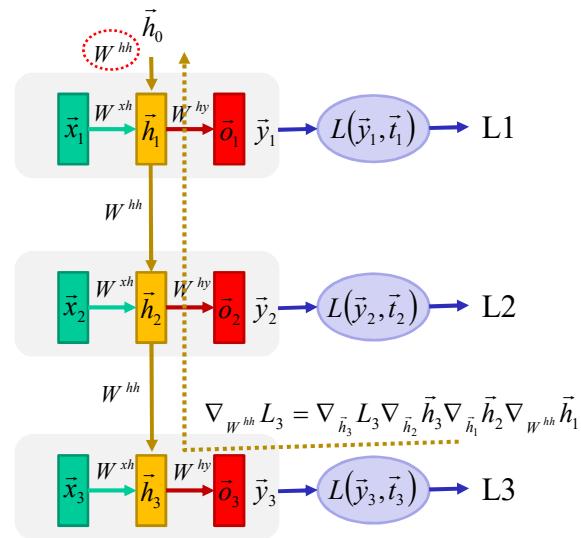
Réseau récurrent: gradient pour W^{hh}



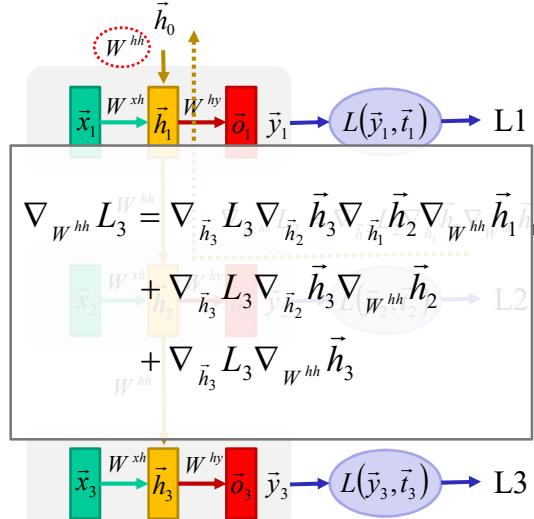
Réseau récurrent: gradient pour W^{hh}



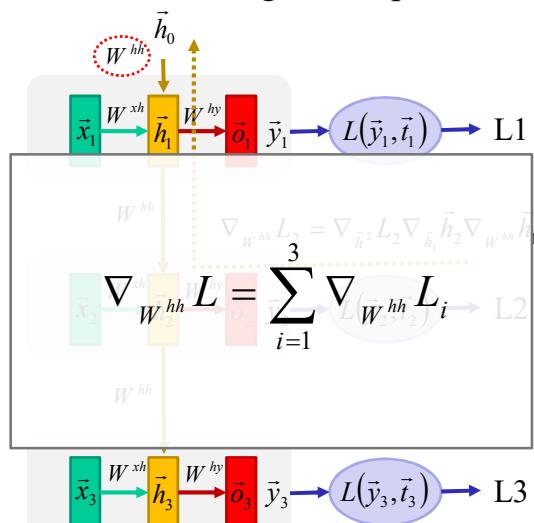
Réseau récurrent: gradient pour W^{hh}



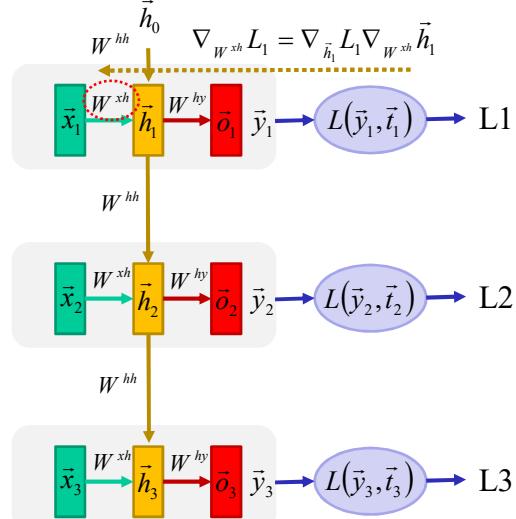
Réseau récurrent: gradient pour W^{hh}



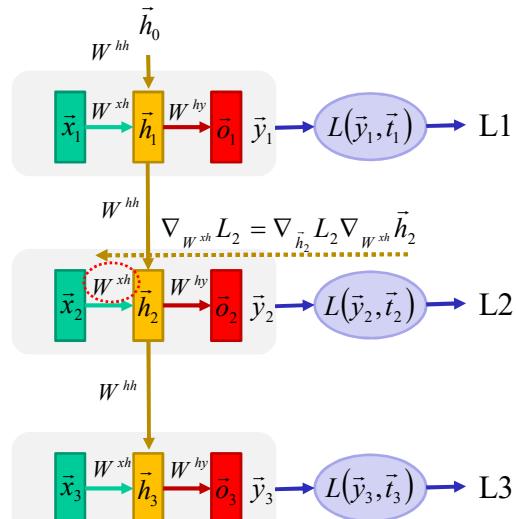
Réseau récurrent: gradient pour W^{hh}



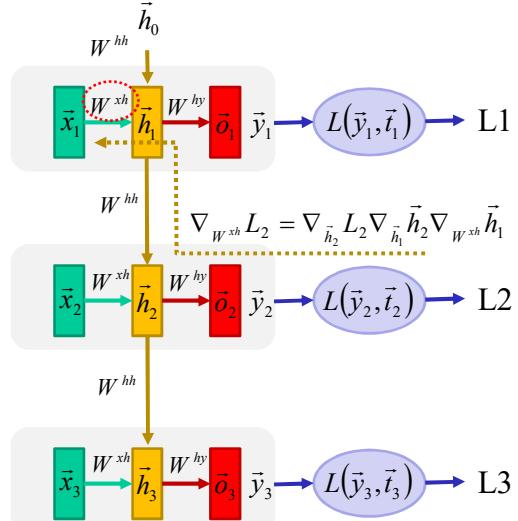
Réseau récurrent: gradient pour W^{xh}



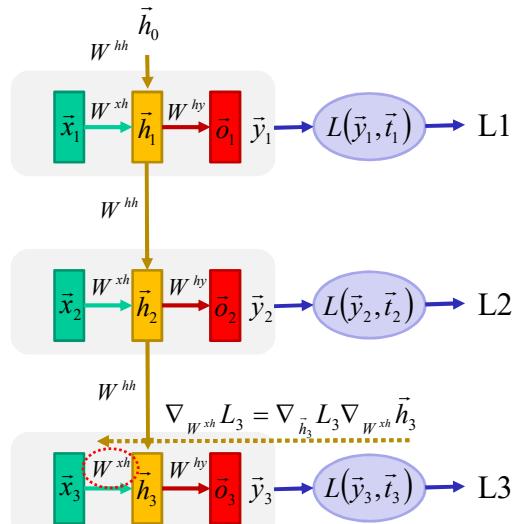
Réseau récurrent: gradient pour W^{xh}



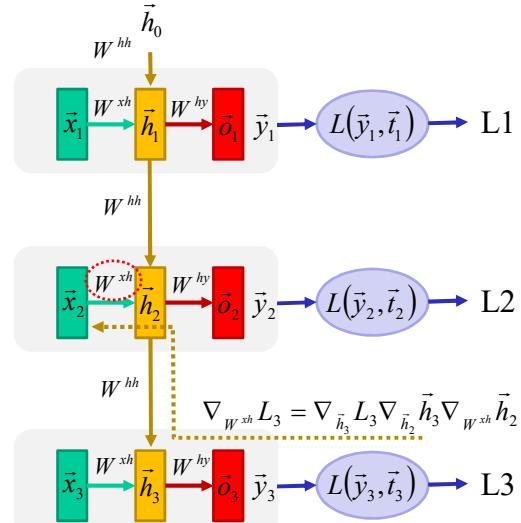
Réseau récurrent: gradient pour W^{xh}



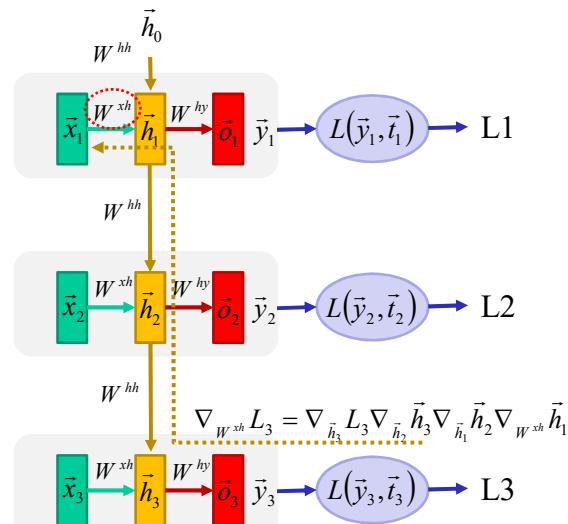
Réseau récurrent: gradient pour W^{xh}



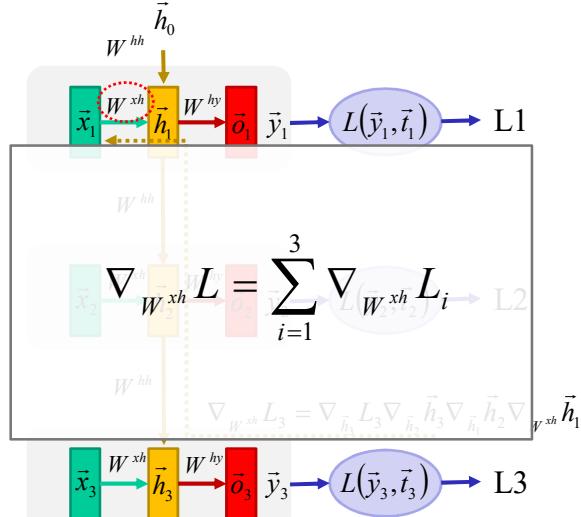
Réseau récurrent: gradient pour W^{xh}



Réseau récurrent: gradient pour W^{xh}



Réseau récurrent: gradient pour W^{xh}



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(xrange(len(inputs))): # loop over time
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 += db
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         dWxh += np.dot(ddraw, xs[t].T)
57         dWhh += np.dot(ddraw, 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

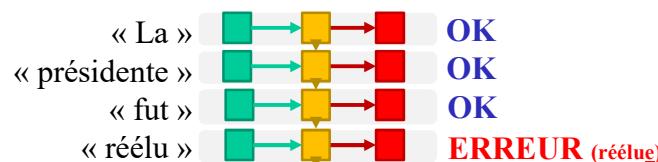
Les réseaux récurrents ont un inconvénient majeur:

difficile à établir des **relations à longue distance**

130

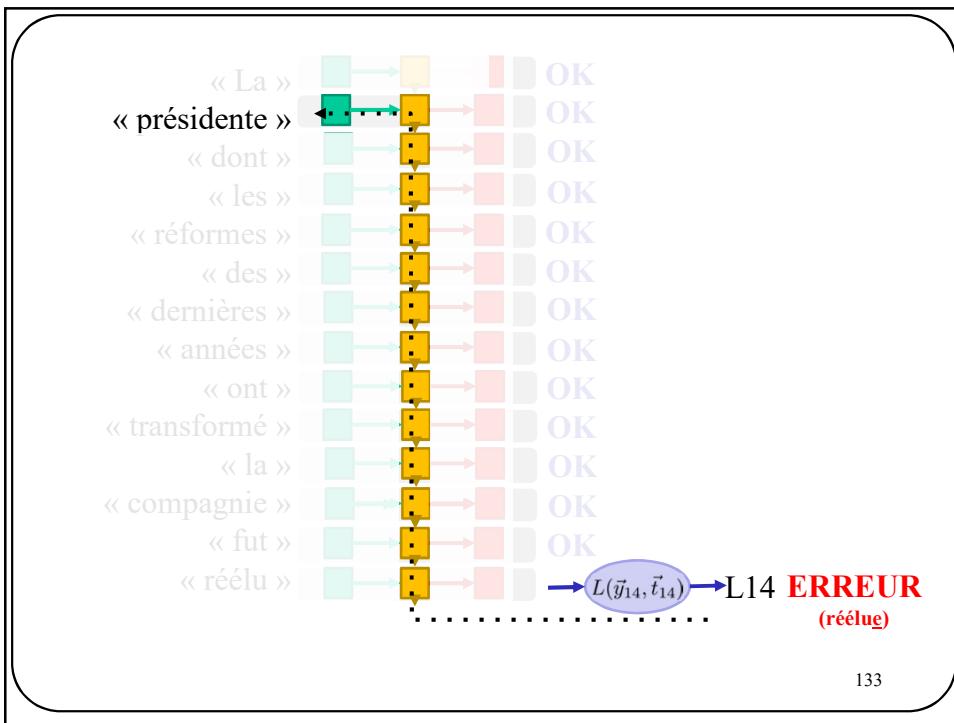
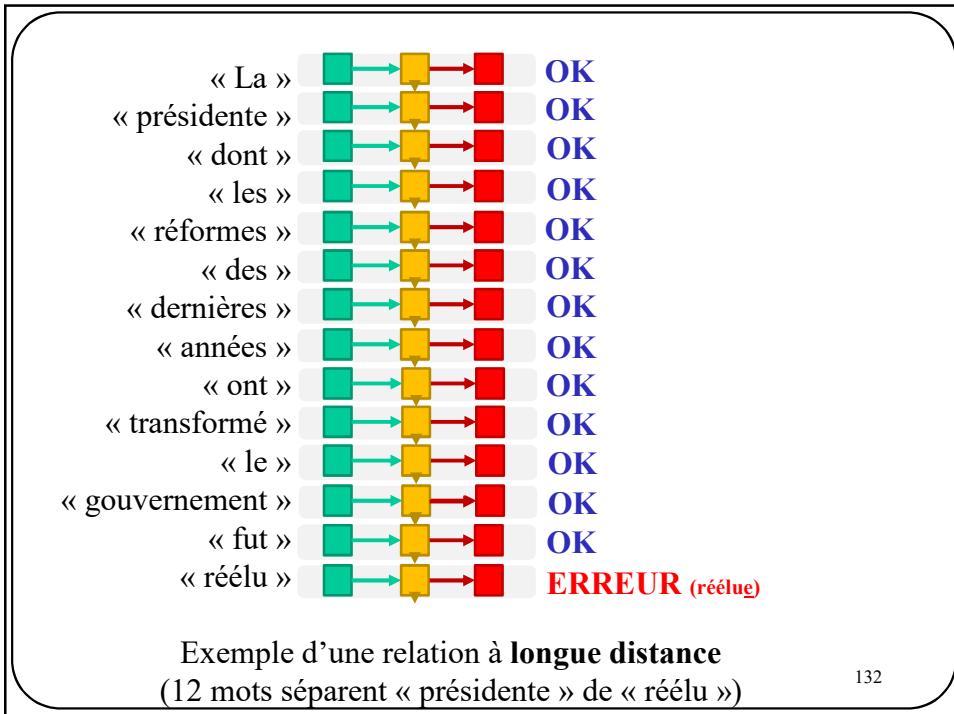
Exemples: analyse grammaticale

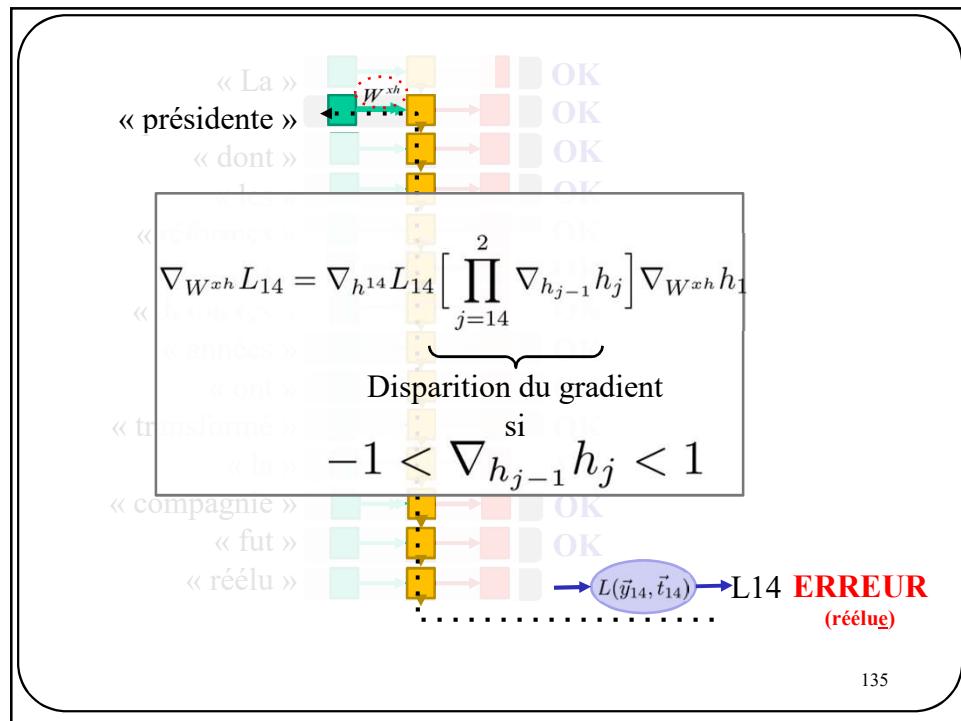
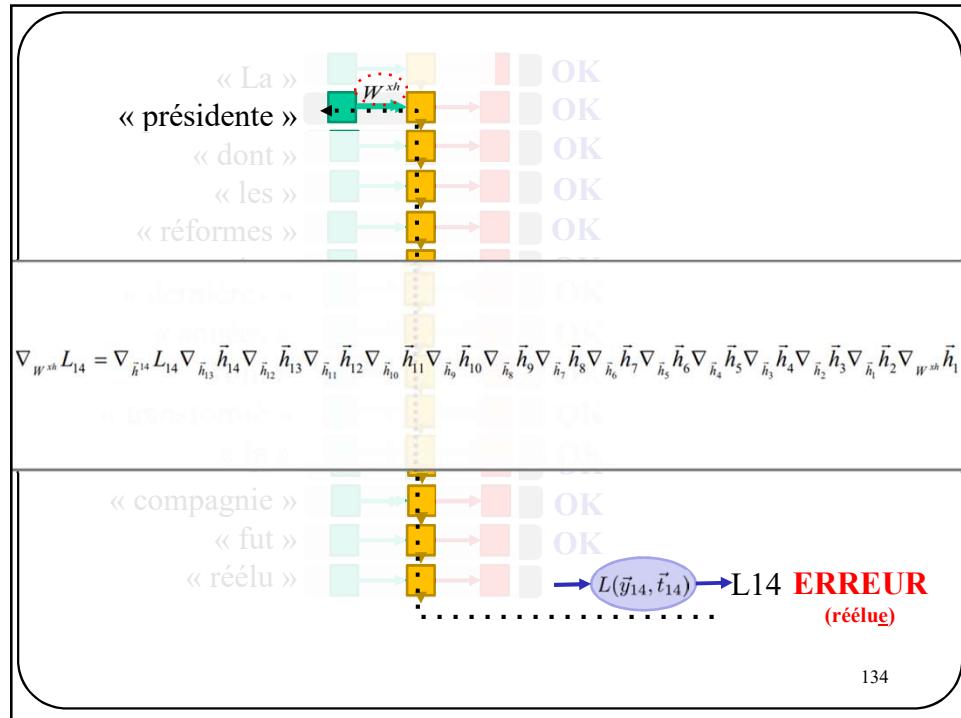
Entraîner un réseau à détecter des erreurs grammaticales

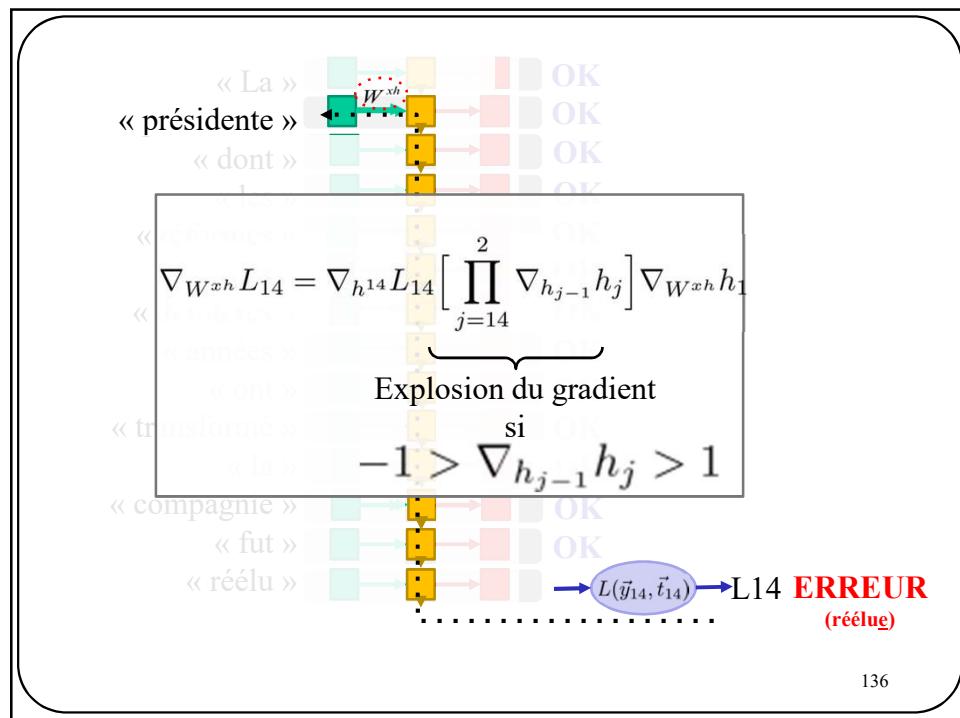


Exemple d'une relation à **courte distance**
(1 mot sépare « présidente » de « réélu »)

131



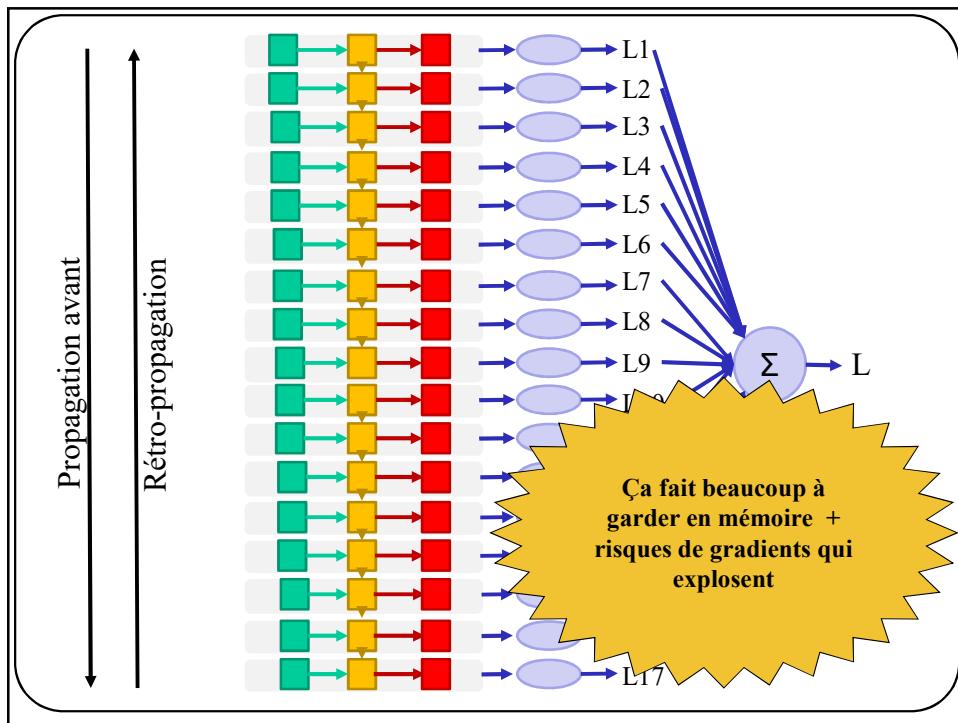
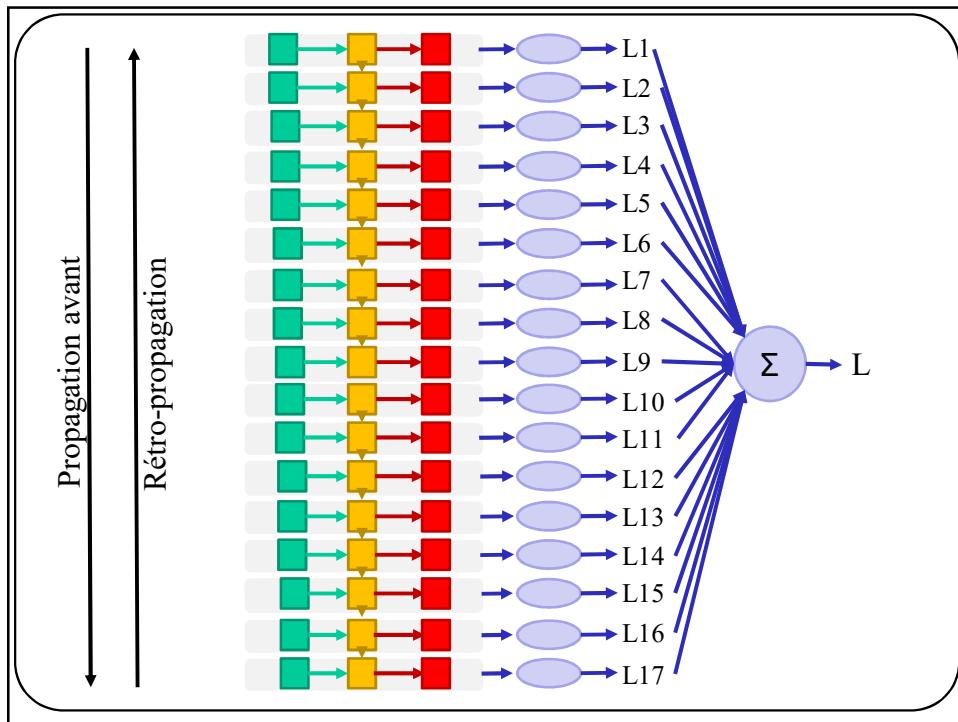




Problème connexe

Gestion de la mémoire

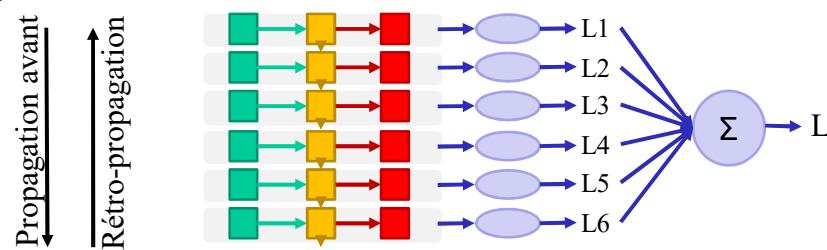
137



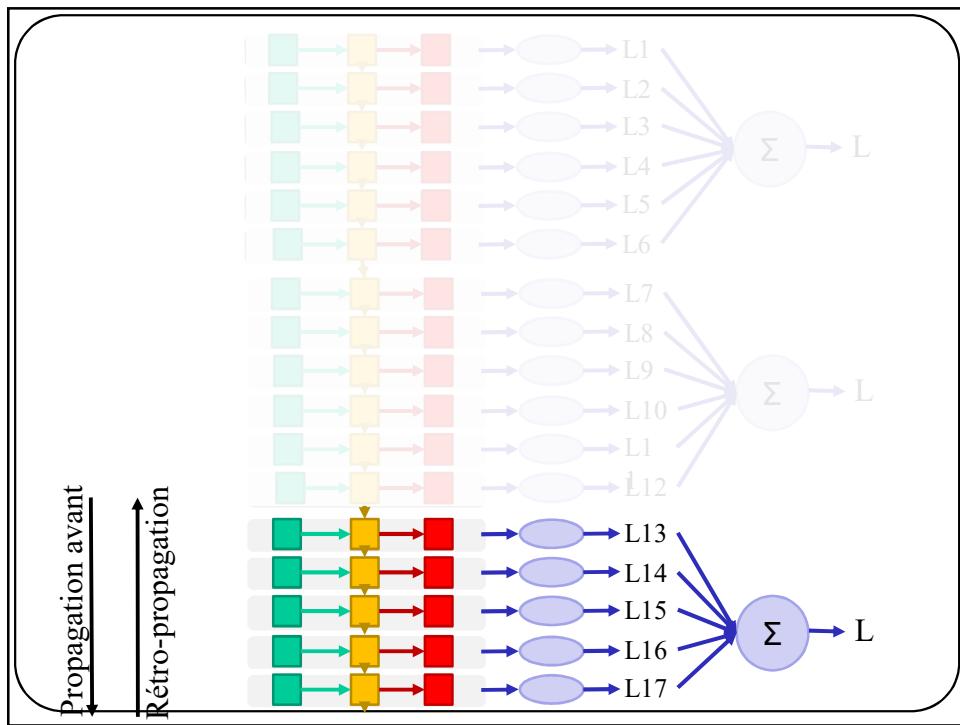
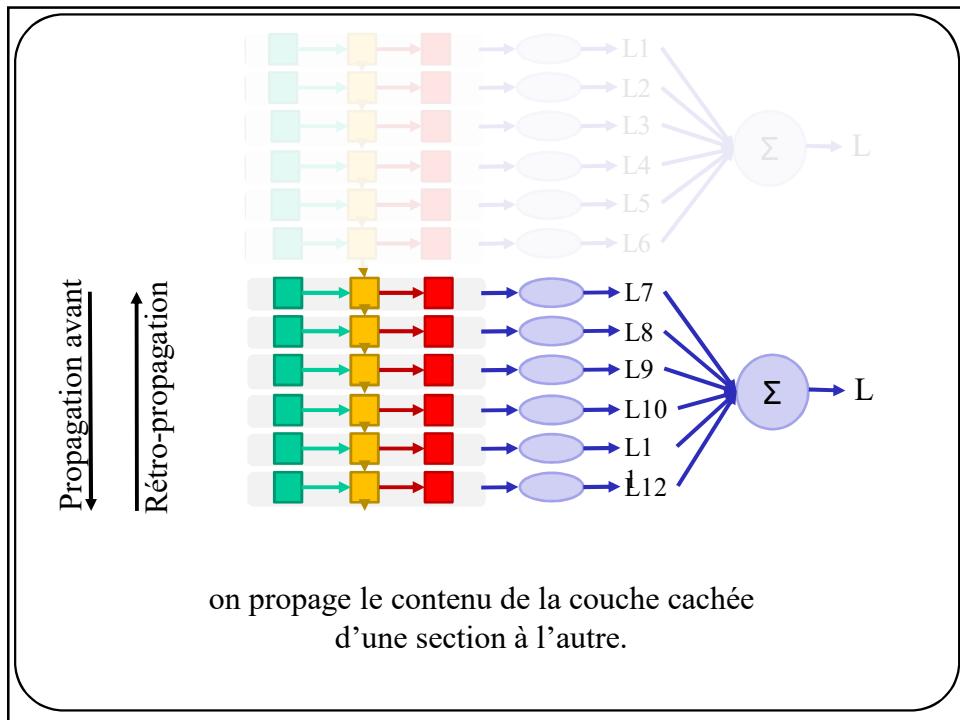
Solution pour la gestion de la mémoire

Fenêtres coulissantes

140



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

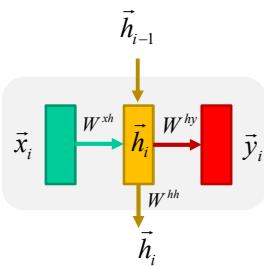


Solution à la disparition du gradient:

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

144

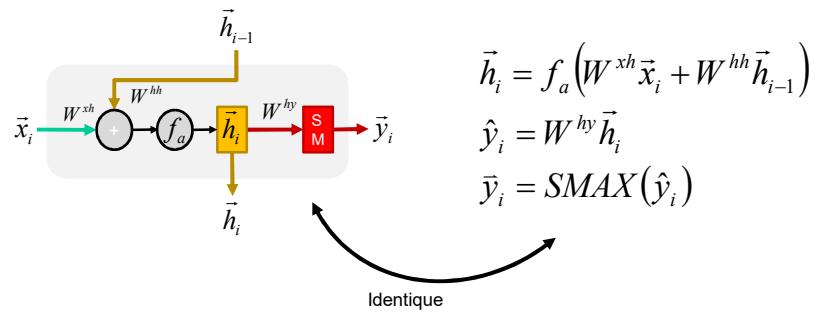
Illustration + formulation d'un RNN



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

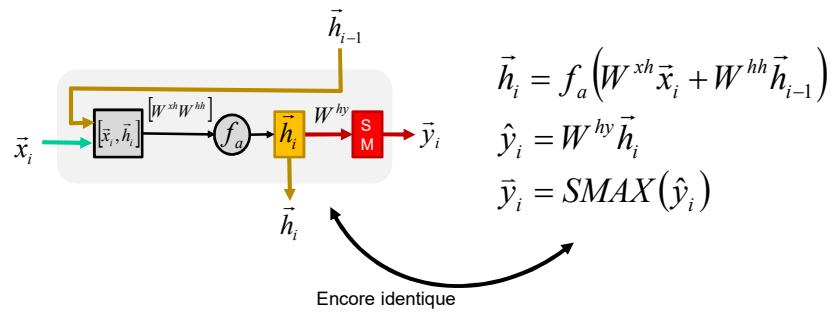
145

Autre illustration du même RNN



146

Autre illustration du même RNN

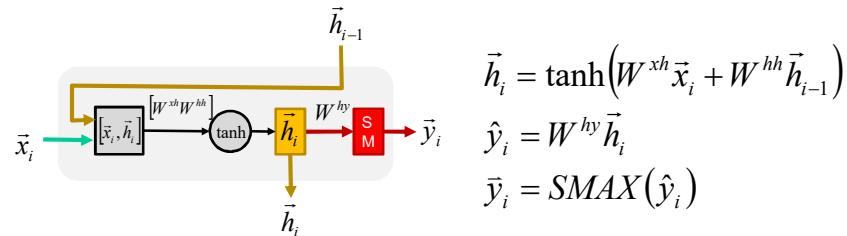


147

GRU (Gated Recurrent Unit)

Modif 1

$$f_a = \tanh$$



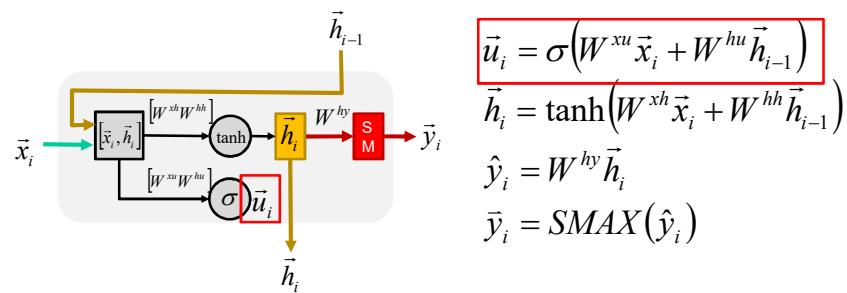
148

GRU (Gated Recurrent Unit)

Modif 2

Update gate

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

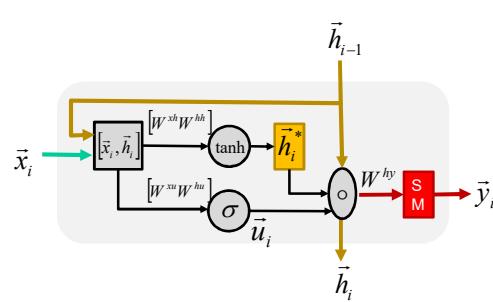


149

GRU (Gated Recurrent Unit)

Modif 2

Update gate



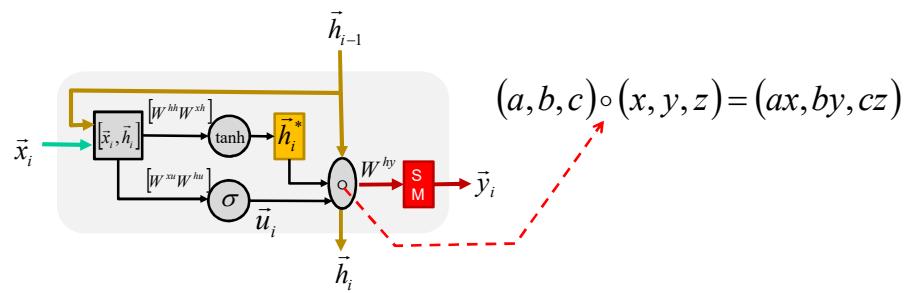
$$\begin{aligned}\vec{u}_i &= \sigma(W^{xu} \vec{x}_i + W^{hu} \vec{h}_{i-1}) \\ \vec{h}_i^* &= \tanh(W^{xh} \vec{x}_i + W^{hh} \vec{h}_{i-1}) \\ \vec{h}_i &= \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \circ \vec{h}_{i-1} \\ \hat{y}_i &= W^{hy} \vec{h}_i \\ \bar{y}_i &= \text{SMAX}(\hat{y}_i)\end{aligned}$$

150

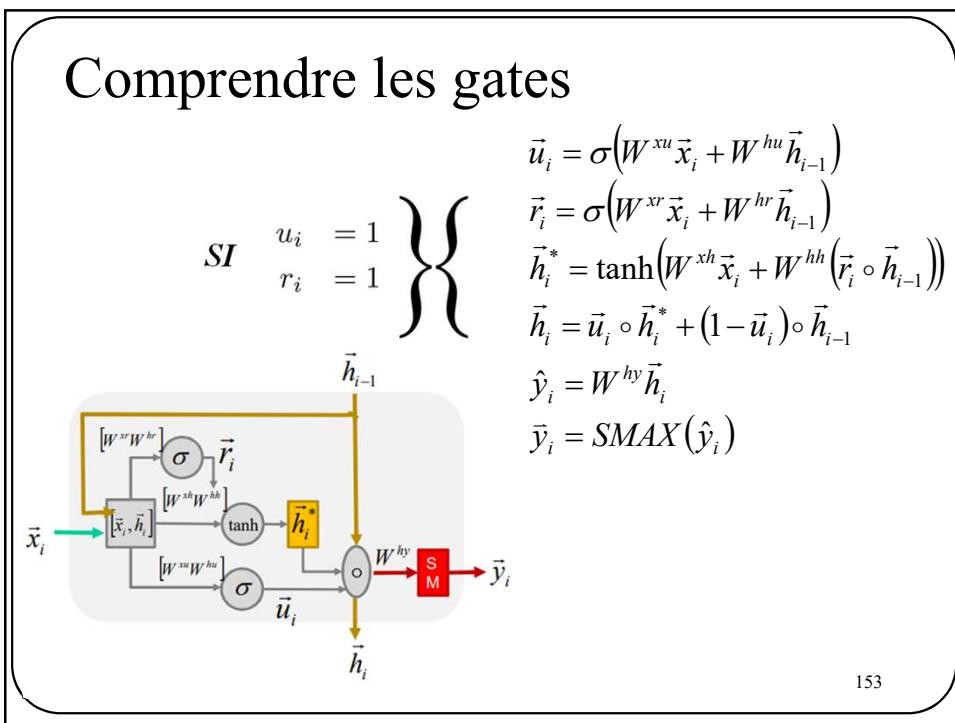
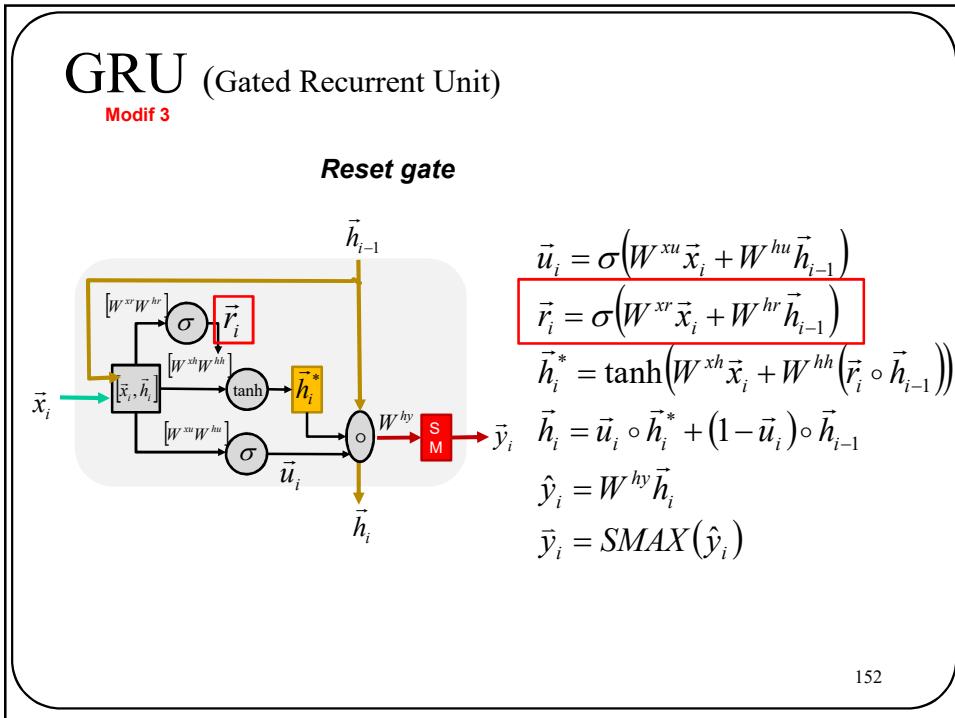
GRU (Gated Recurrent Unit)

Modif 2

« element-wise product »
ou
« Hadamard product »

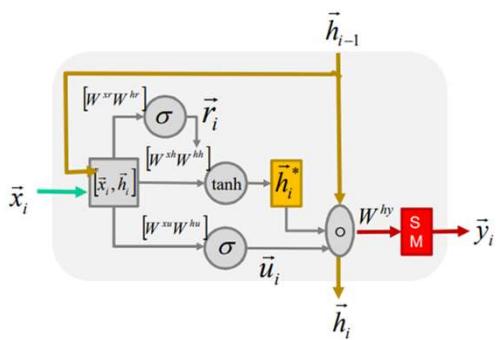


151



Comprendre les gates

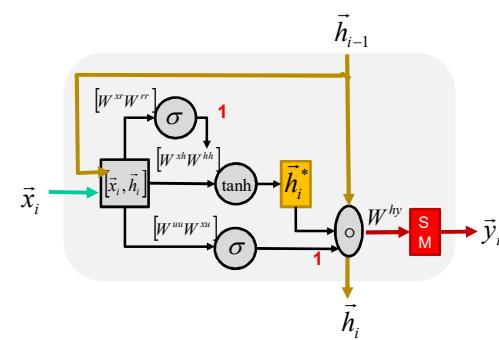
$$SI \quad \begin{cases} u_i = 1 \\ r_i = 1 \end{cases} \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{r}_i = \sigma(W^{xr}\vec{x}_i + W^{hr}\vec{h}_{i-1}) \\ \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}(\vec{r}_i \circ \vec{h}_{i-1})) \\ \vec{h}_i =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



154

Comprendre les gates

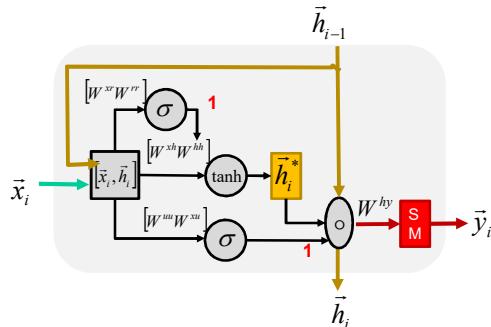
$$SI \quad \begin{cases} u_i = 1 \\ r_i = 1 \end{cases} \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{r}_i = \sigma(W^{xr}\vec{x}_i + W^{hr}\vec{h}_{i-1}) \color{red}{\rightarrow 1} \\ \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}(\vec{r}_i \circ \vec{h}_{i-1}))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



155

Comprendre les gates

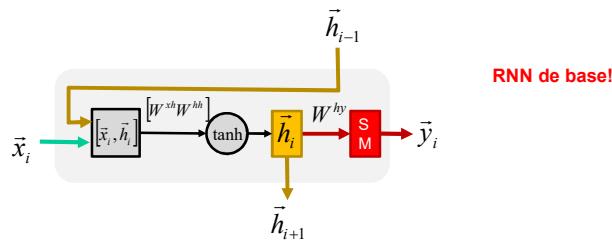
$$SI \quad \begin{cases} \vec{u}_i = 1 \\ \vec{r}_i = 1 \end{cases} \quad \left\{ \begin{array}{l} \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \vec{h}_i = \vec{h}_i^* \\ \hat{y}_i = W^{hy}\vec{h}_i \\ \vec{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$



156

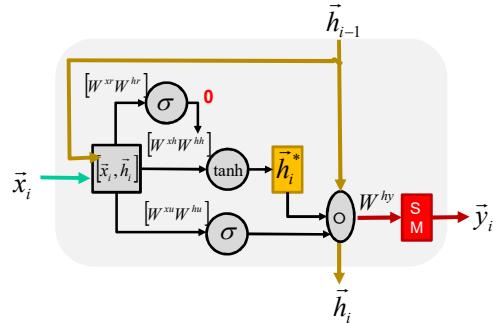
Comprendre les gates

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



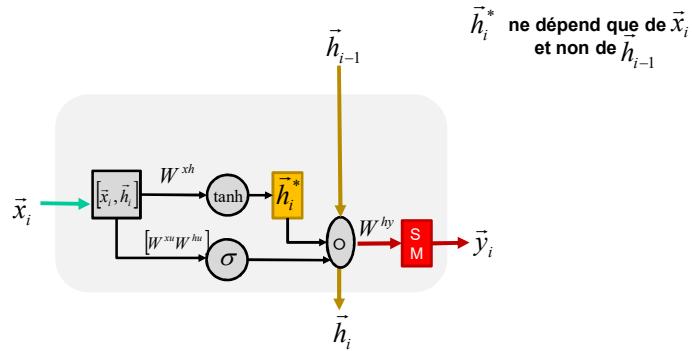
157

$$SI \quad \begin{cases} \vec{u}_i \neq 1 \\ \vec{r}_i = 0 \end{cases} \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{r}_i = \sigma(W^{xr}\vec{x}_i + W^{hr}\vec{h}_{i-1}) \\ \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \vec{h}_i = \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \circ \vec{h}_{i-1} \\ \hat{y}_i = W^{hy}\vec{h}_i \\ \vec{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$



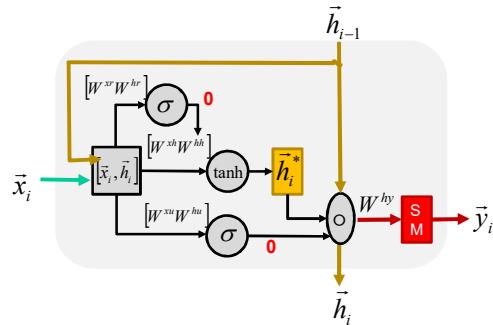
158

$$SI \quad \begin{cases} \vec{u}_i \neq 1 \\ \vec{r}_i = 0 \end{cases} \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(W^{xu}\vec{x}_i + W^{hu}\vec{h}_{i-1}) \\ \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i) \\ \vec{h}_i = \vec{u}_i \circ \vec{h}_i^* + (1 - \vec{u}_i) \circ \vec{h}_{i-1} \\ \hat{y}_i = W^{hy}\vec{h}_i \\ \vec{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$



159

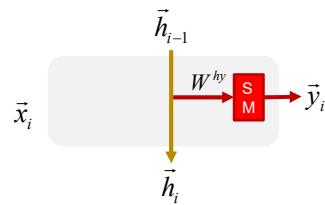
$$SI \quad \begin{cases} \vec{u}_i = 0 \\ \vec{r}_i = 0 \end{cases} \quad \left\{ \begin{array}{l} \vec{u}_i = \sigma(\cancel{W^{uu}\vec{x}_i + W^{uu}\vec{h}_{i-1}}) \\ \vec{r}_i = \sigma(\cancel{W^{rr}\vec{x}_i + W^{rr}\vec{h}_{i-1}}) \\ \vec{h}_i^* = \tanh(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\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



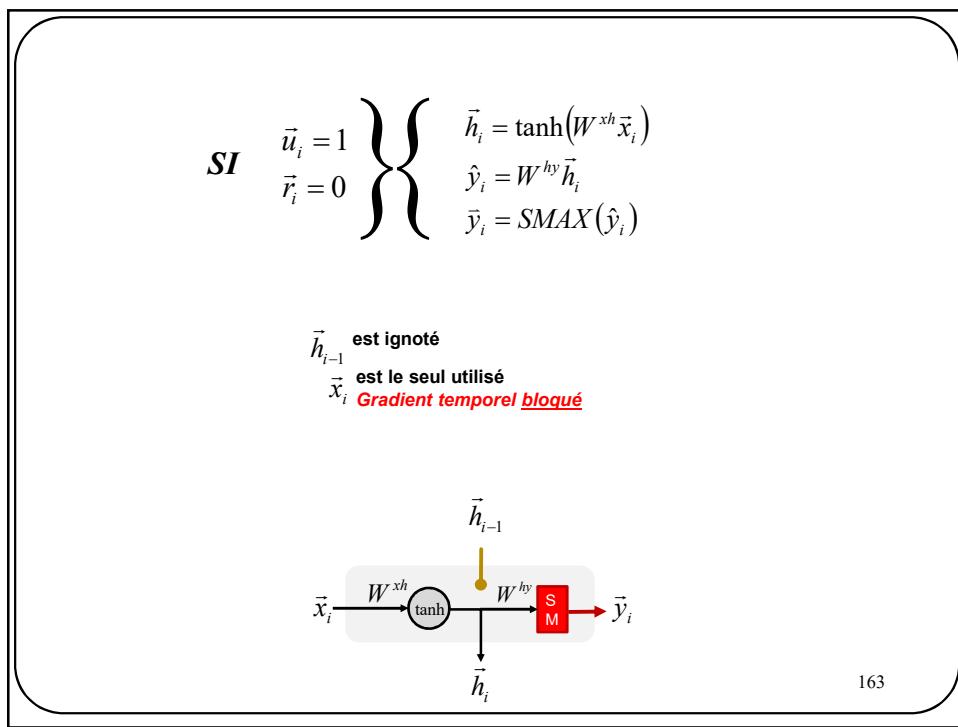
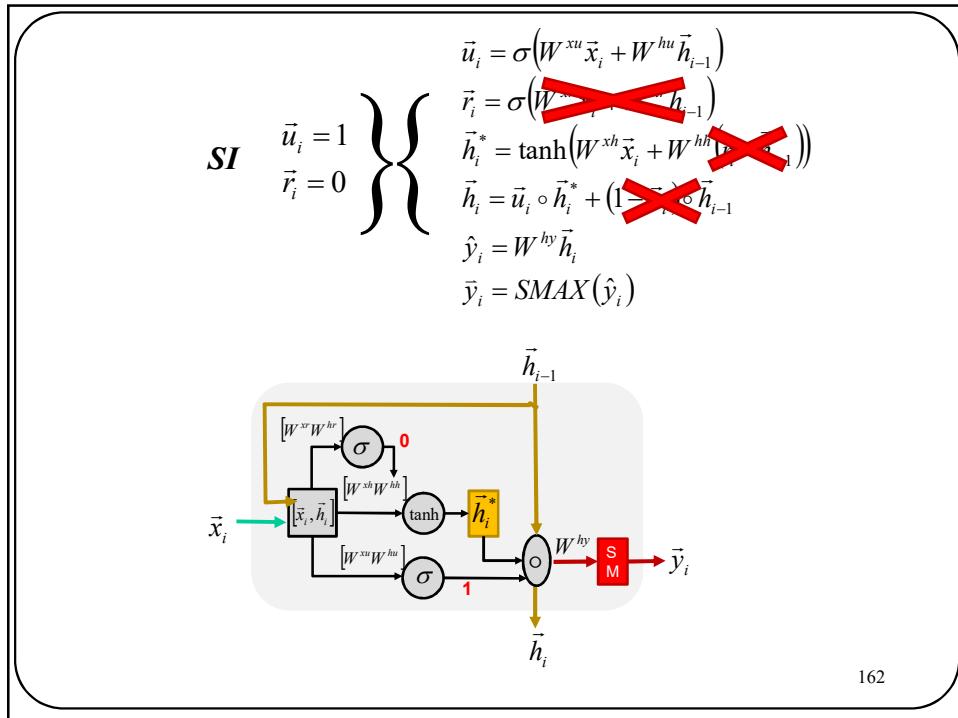
160

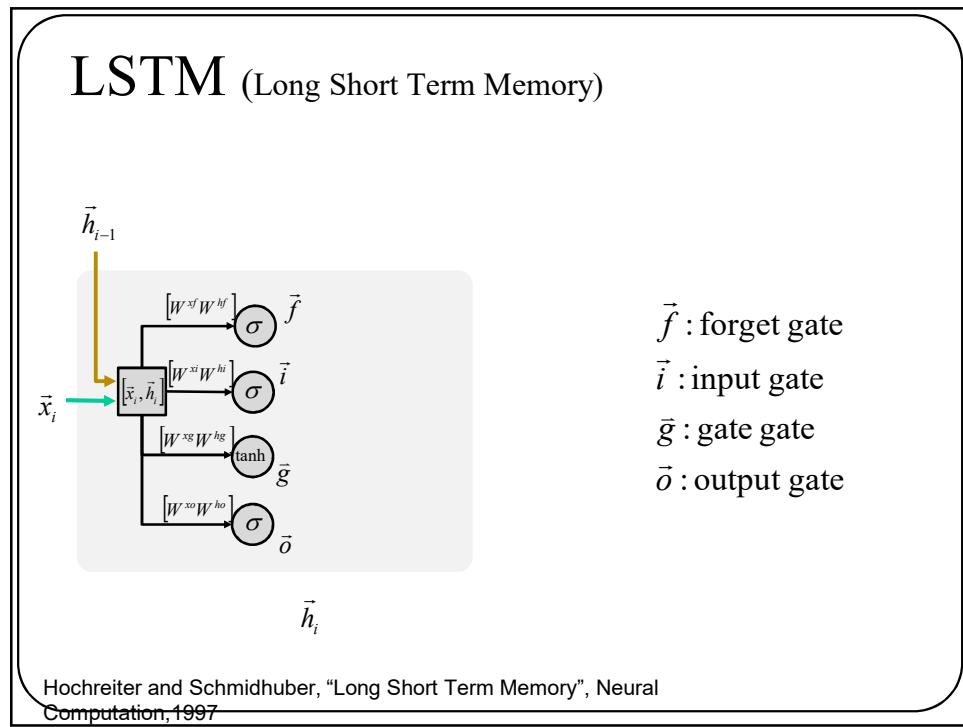
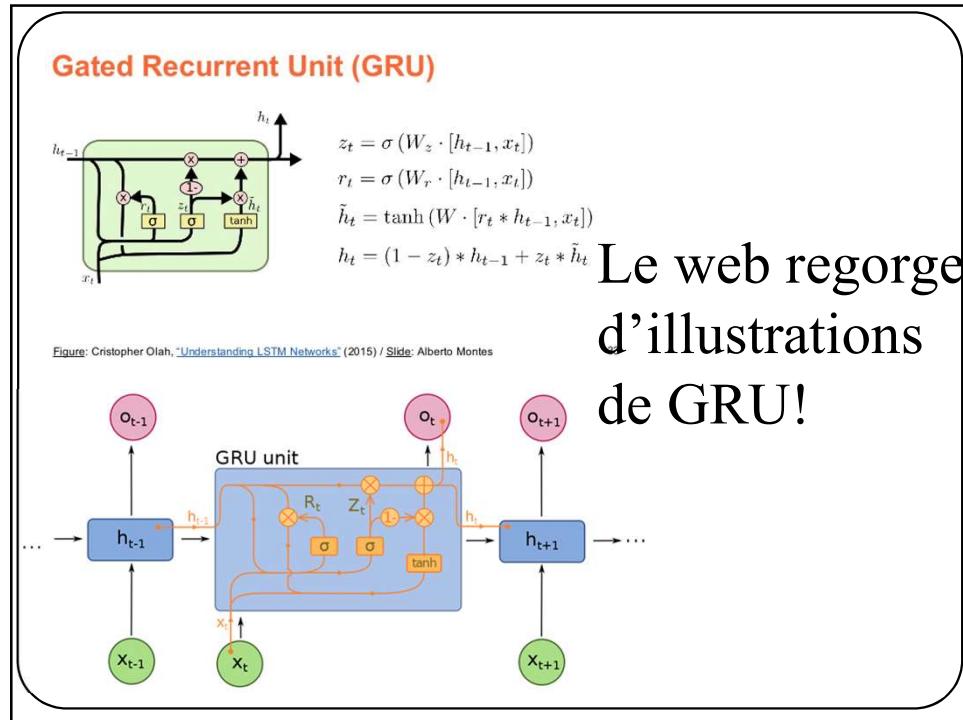
$$SI \quad \begin{cases} \vec{u}_i = 0 \\ \vec{r}_i = 0 \end{cases} \quad \left\{ \begin{array}{l} \vec{h}_i = \vec{h}_{i-1} \\ \hat{y}_i = W^{cy}\vec{h}_i \\ \vec{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$

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

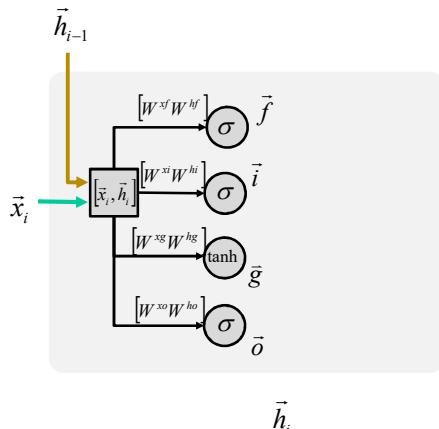


161





LSTM (Long Short Term Memory)



\vec{f} : forget gate

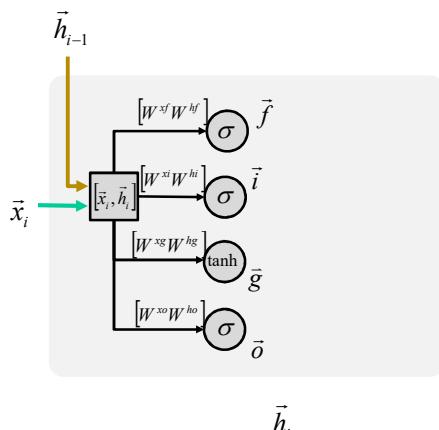
\vec{i} : input gate

\bar{g} : candidate state

Modèle récurrent
le plus utilisé.
À bien
comprendre !

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

LSTM (Long Short Term Memory)



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

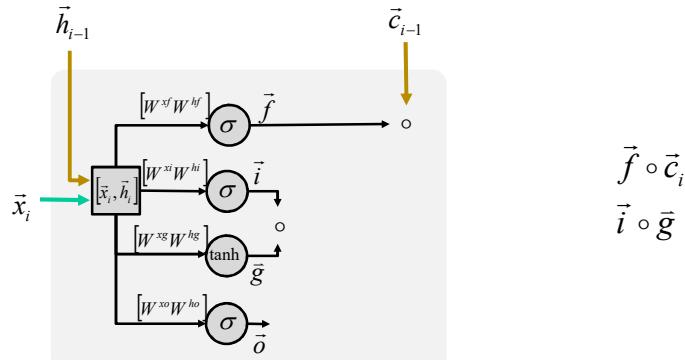
$$\vec{i} = \sigma(W^{xi}\vec{x}_i + W^{hi}\vec{h}_{i-1})$$

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

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

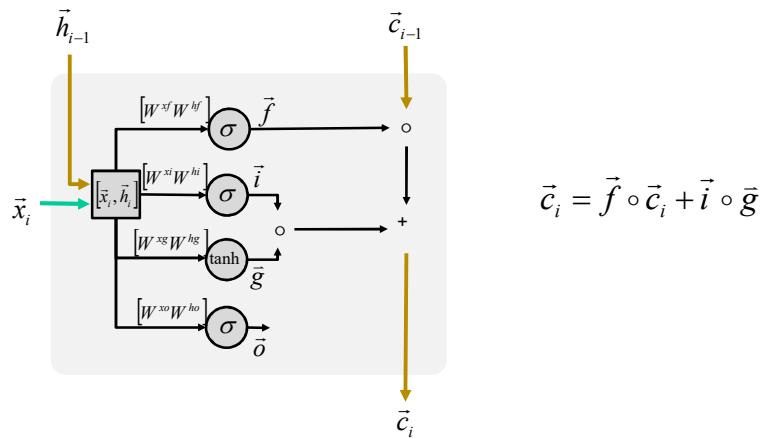
167

LSTM (Long Short Term Memory)



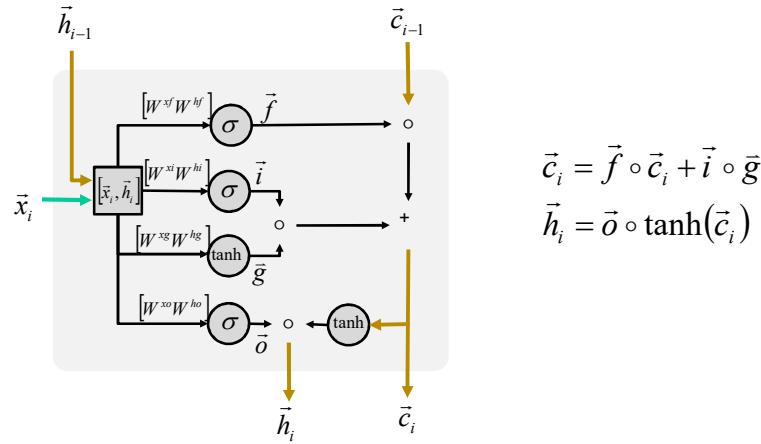
168

LSTM (Long Short Term Memory)



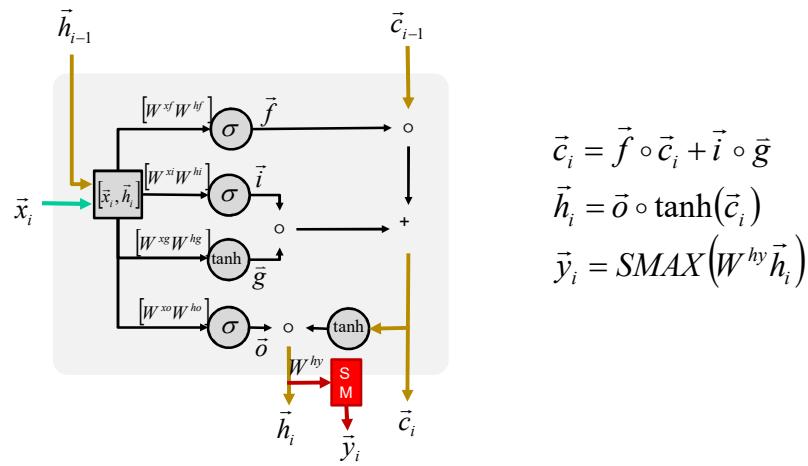
169

LSTM (Long Short Term Memory)



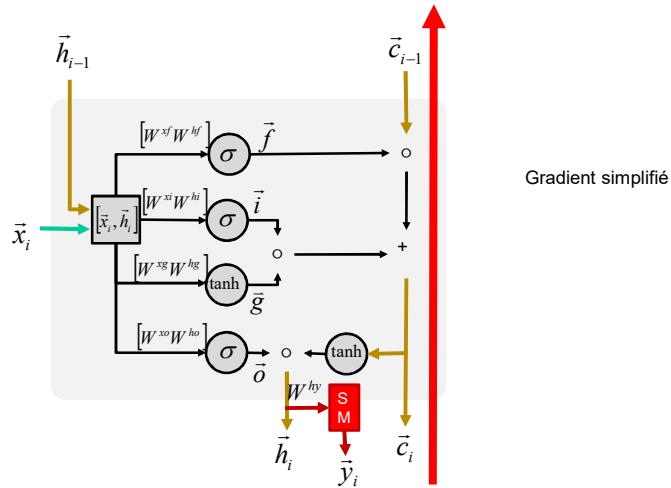
170

LSTM (Long Short Term Memory)



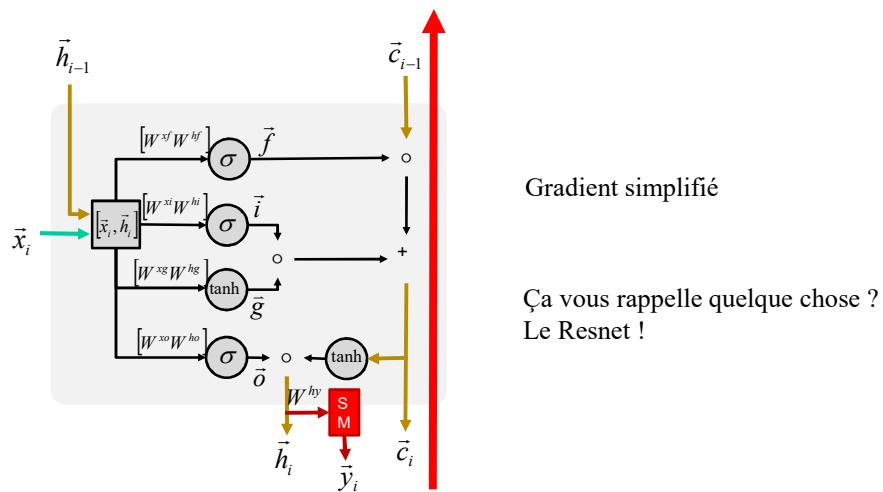
171

LSTM (Long Short Term Memory)

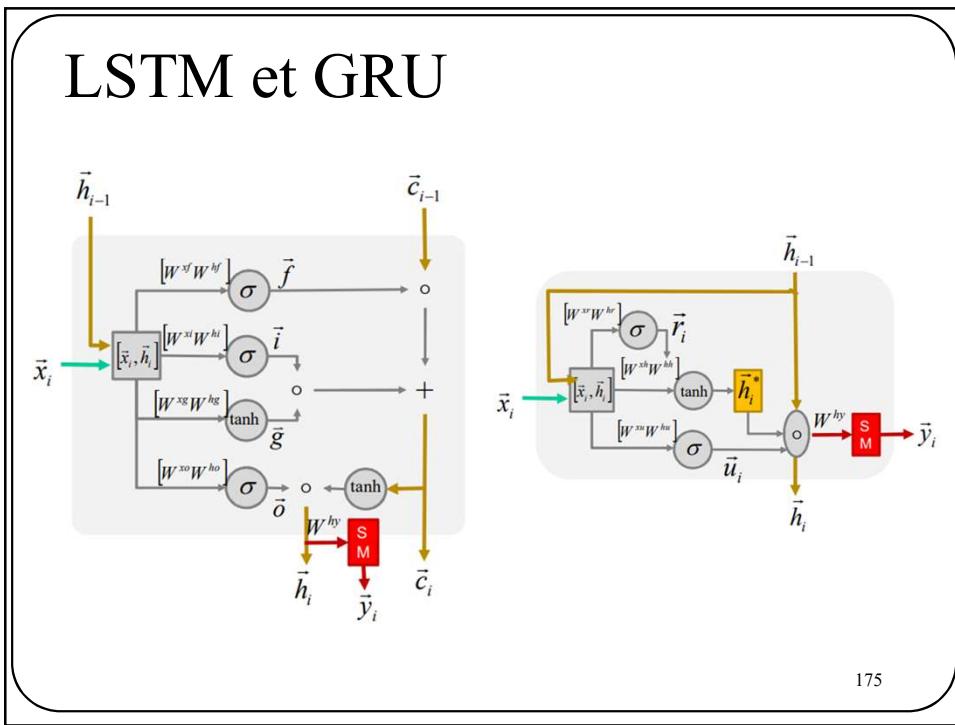
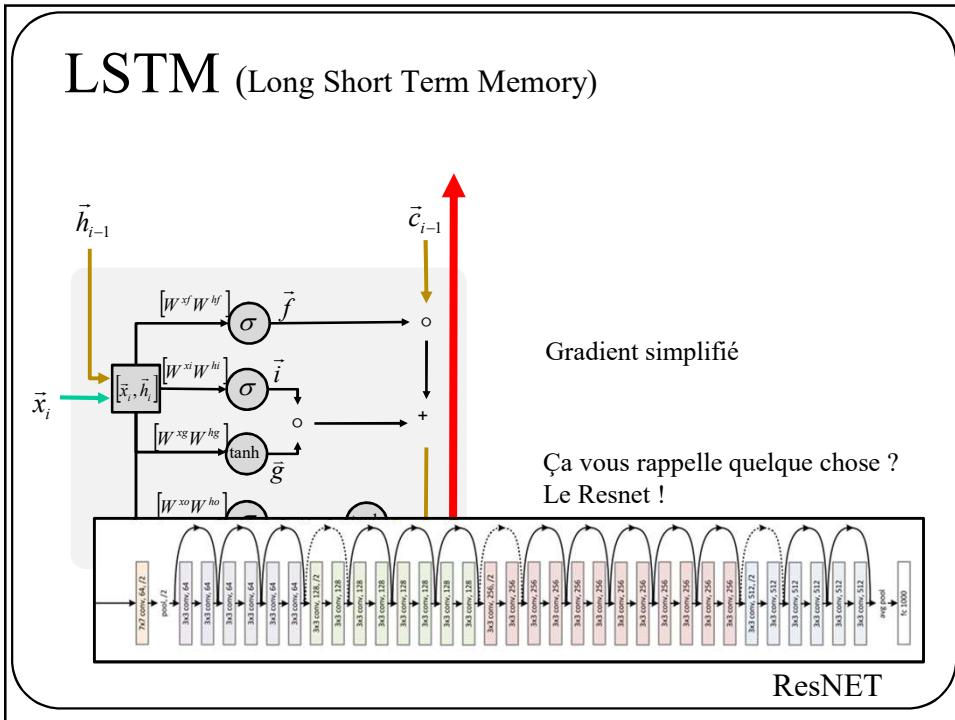


172

LSTM (Long Short Term Memory)

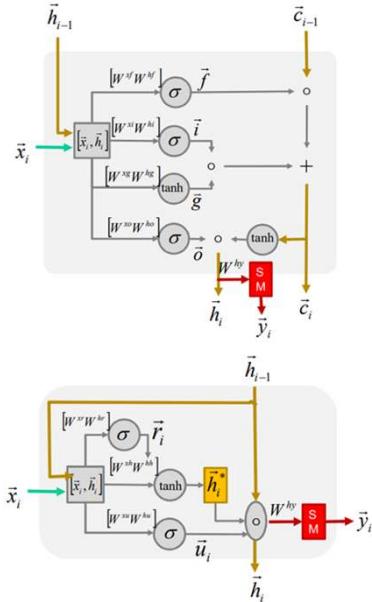


173



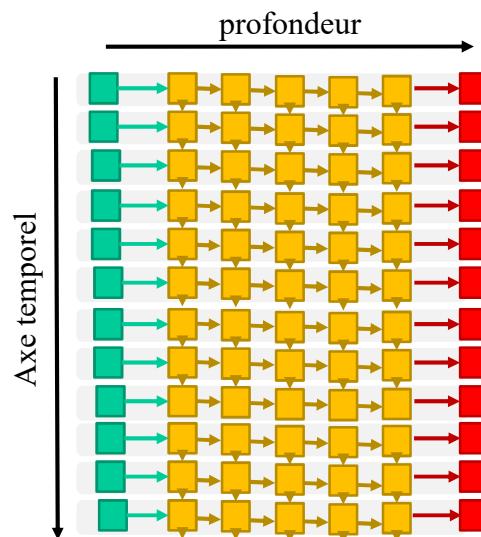
LSTM et GRU

- 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



176

RNN multi-couches



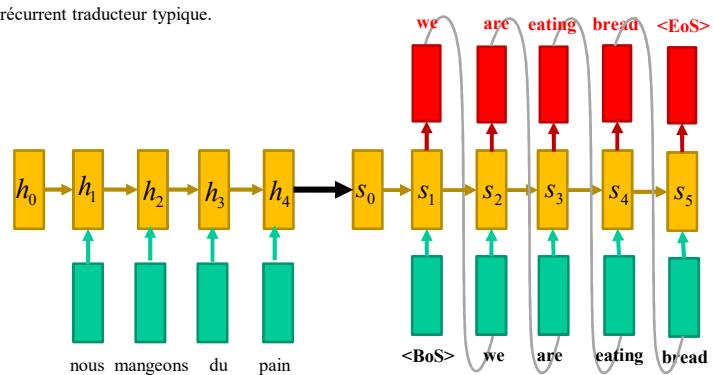
Modèles d'attention

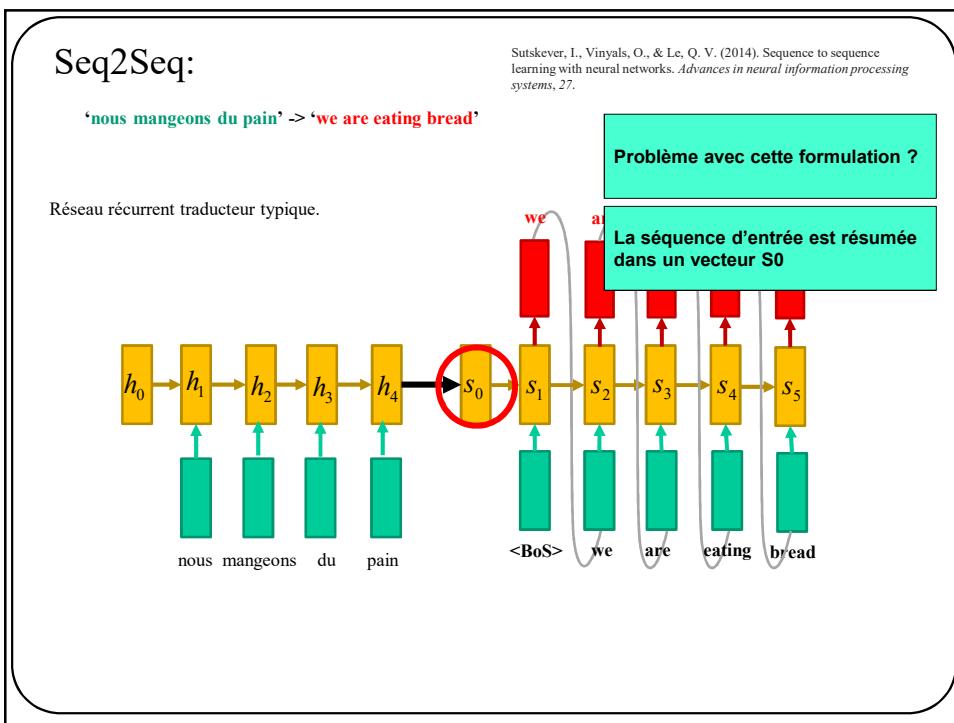
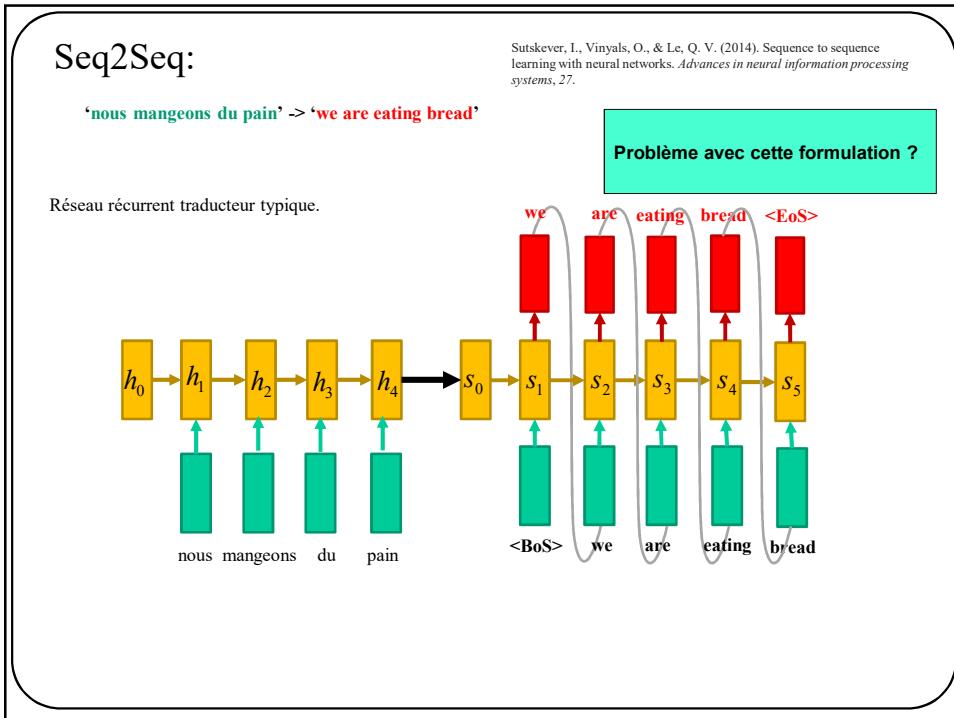
Seq2Seq:

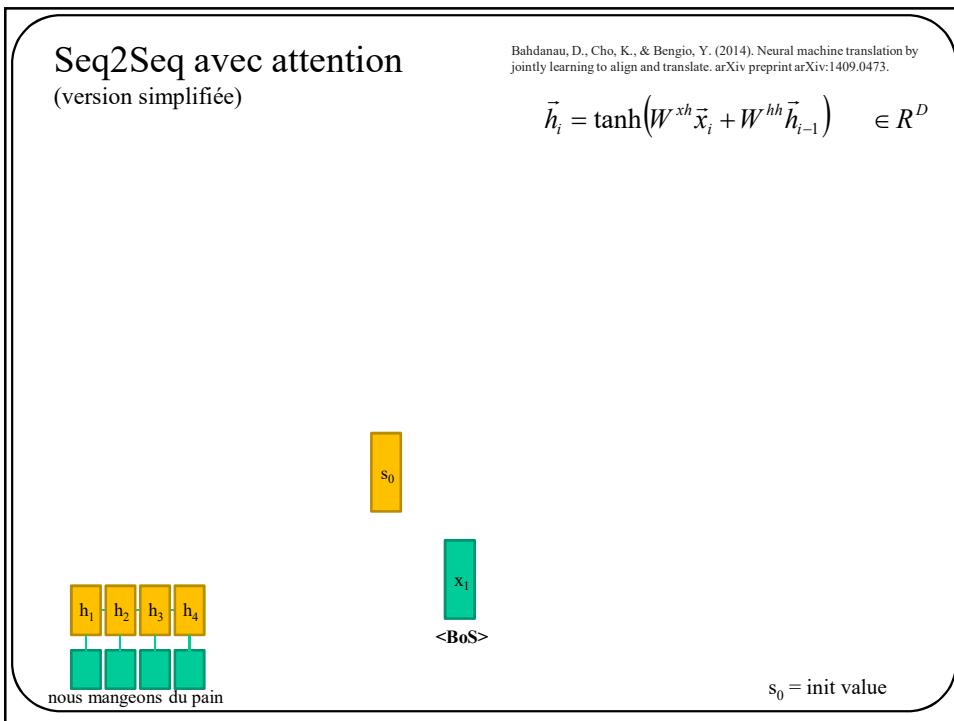
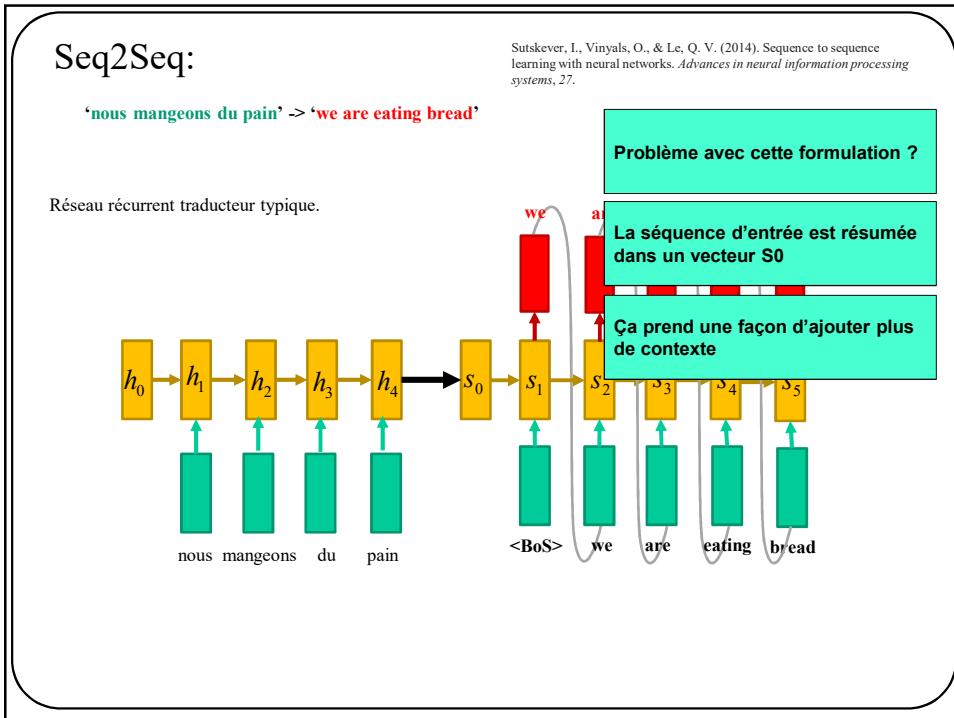
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

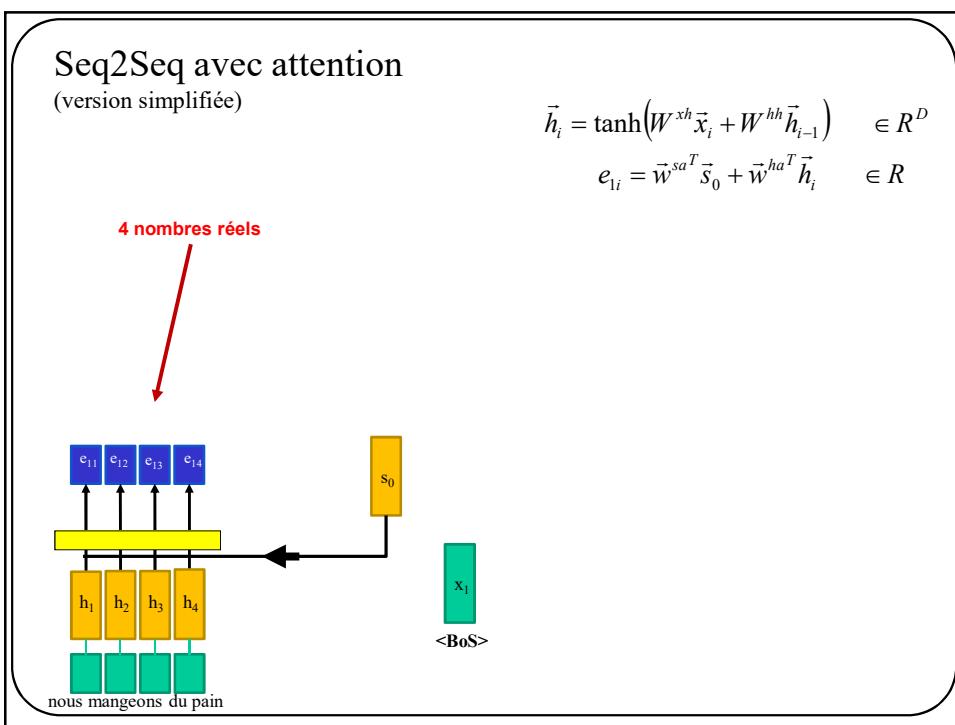
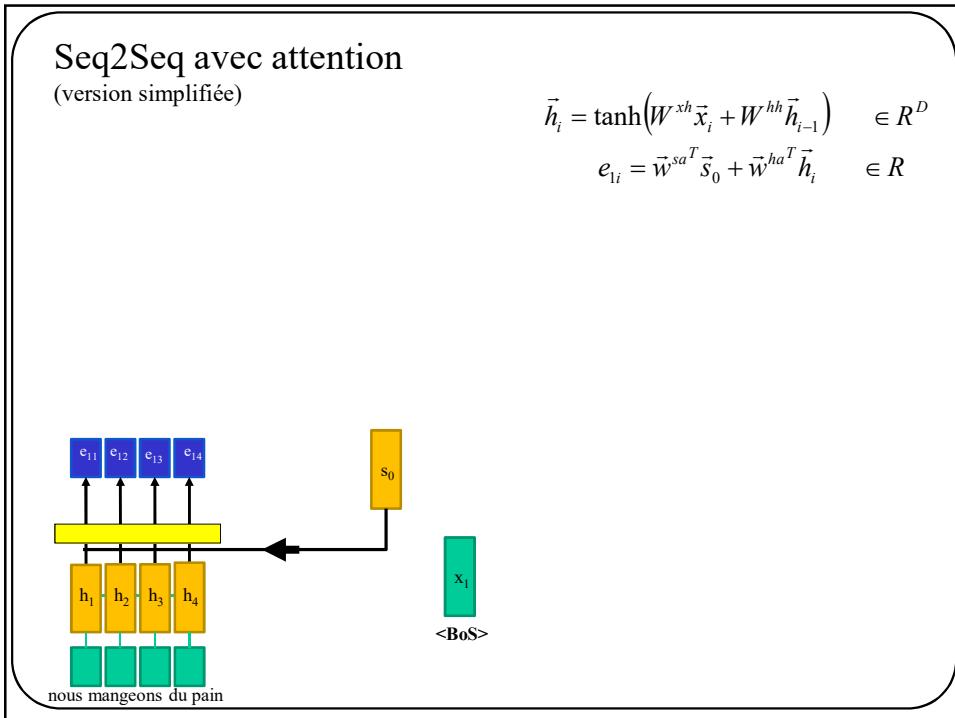
‘nous mangeons du pain’ > ‘we are eating bread’

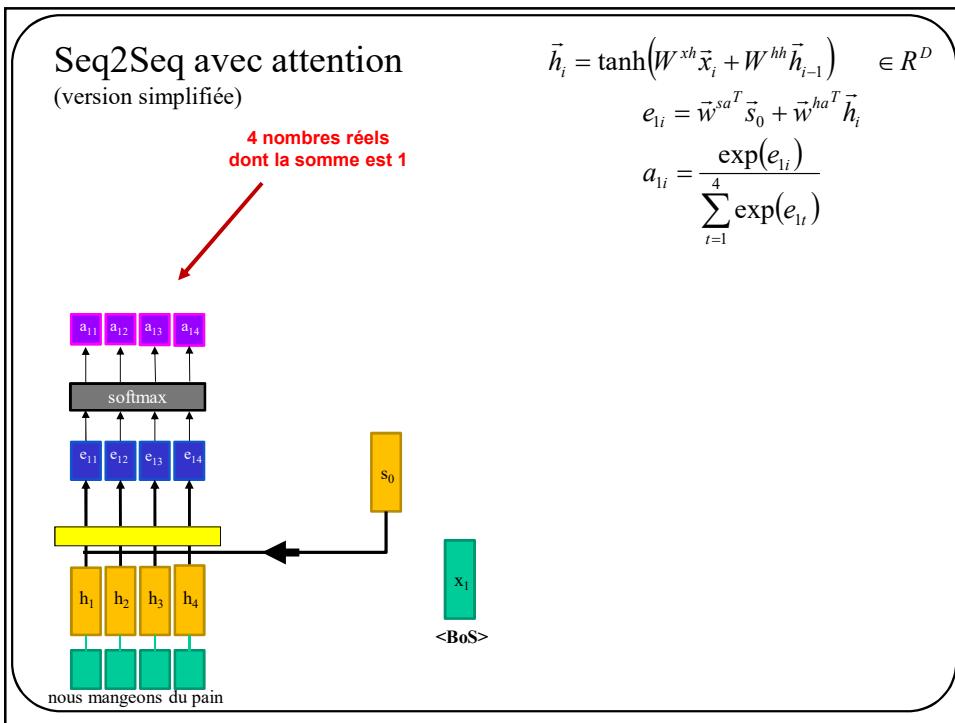
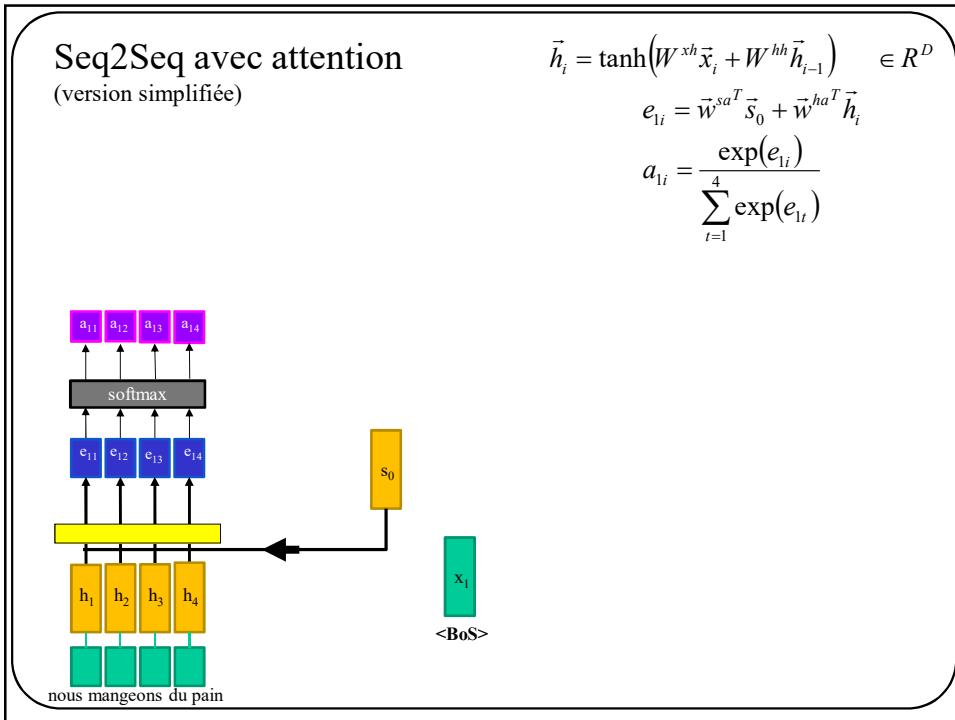
Réseau récurrent traducteur typique.

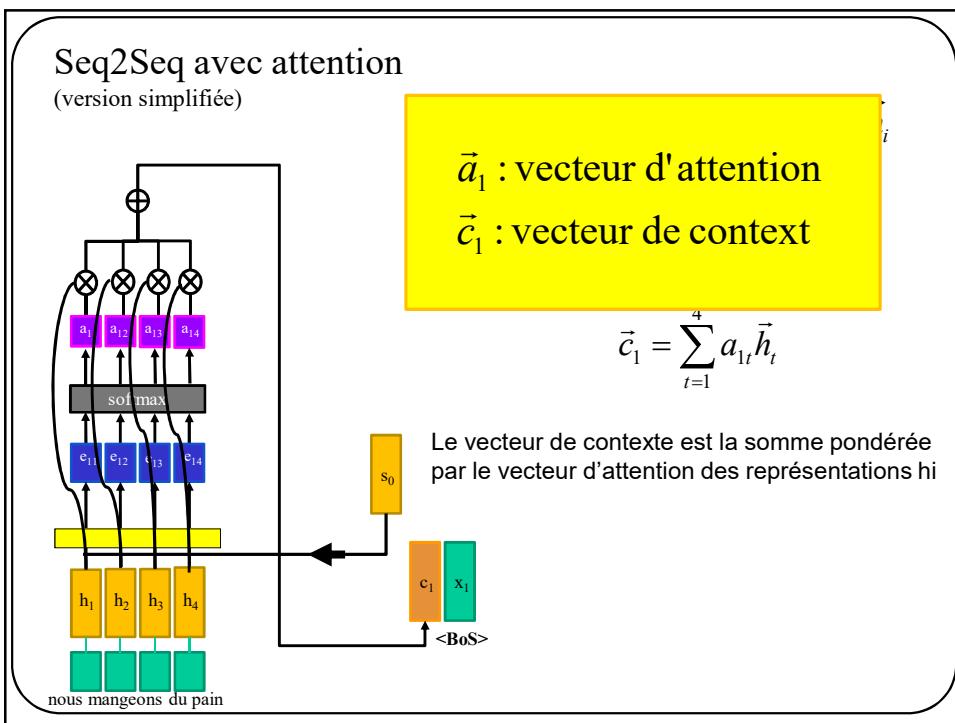
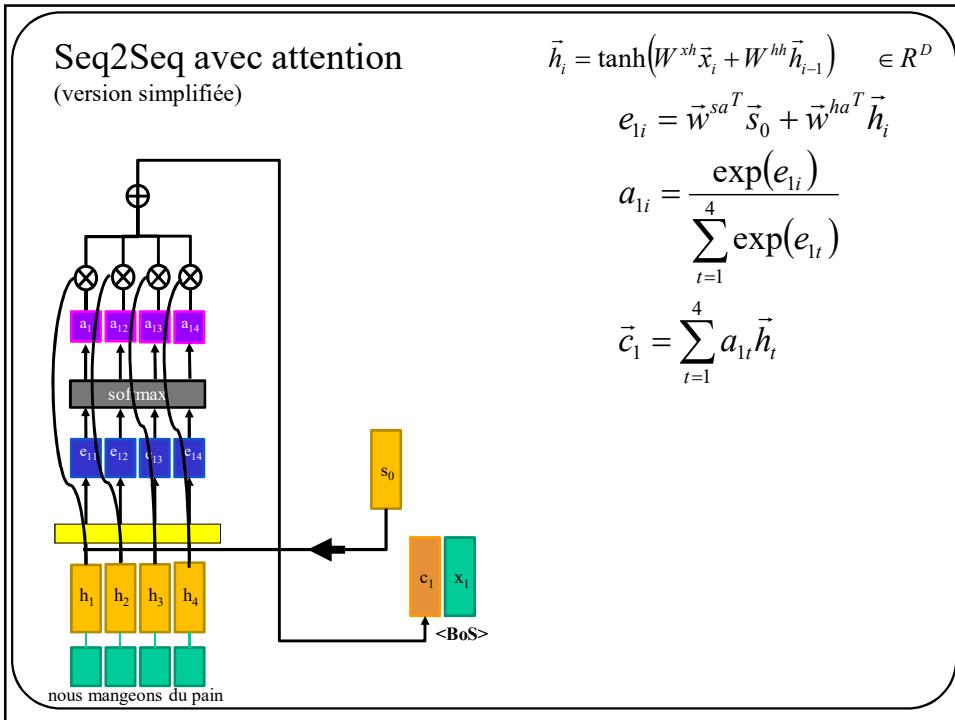


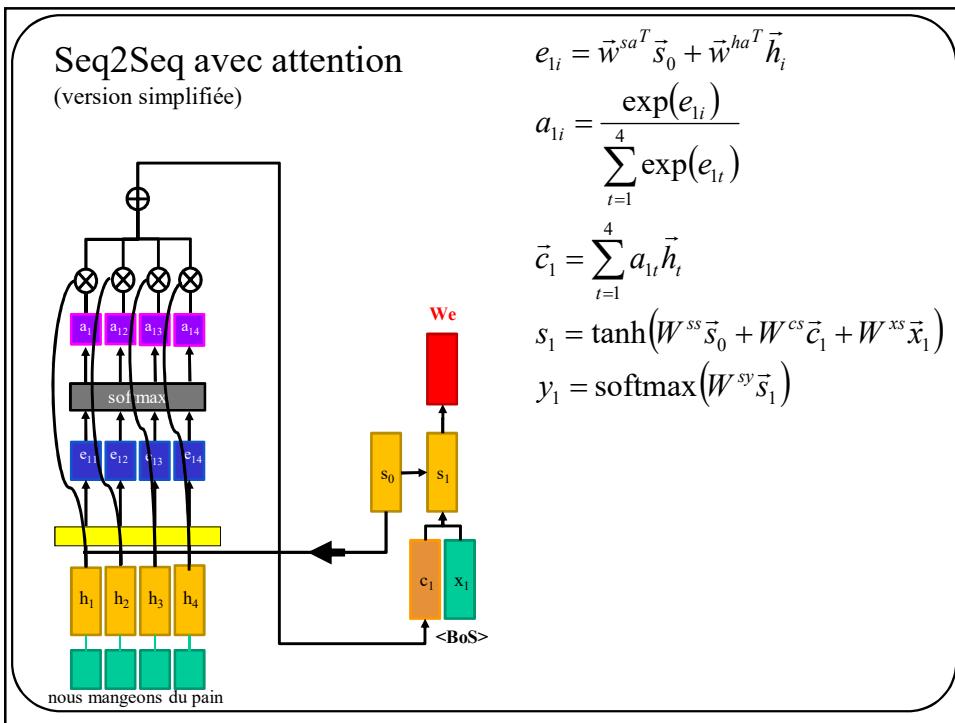
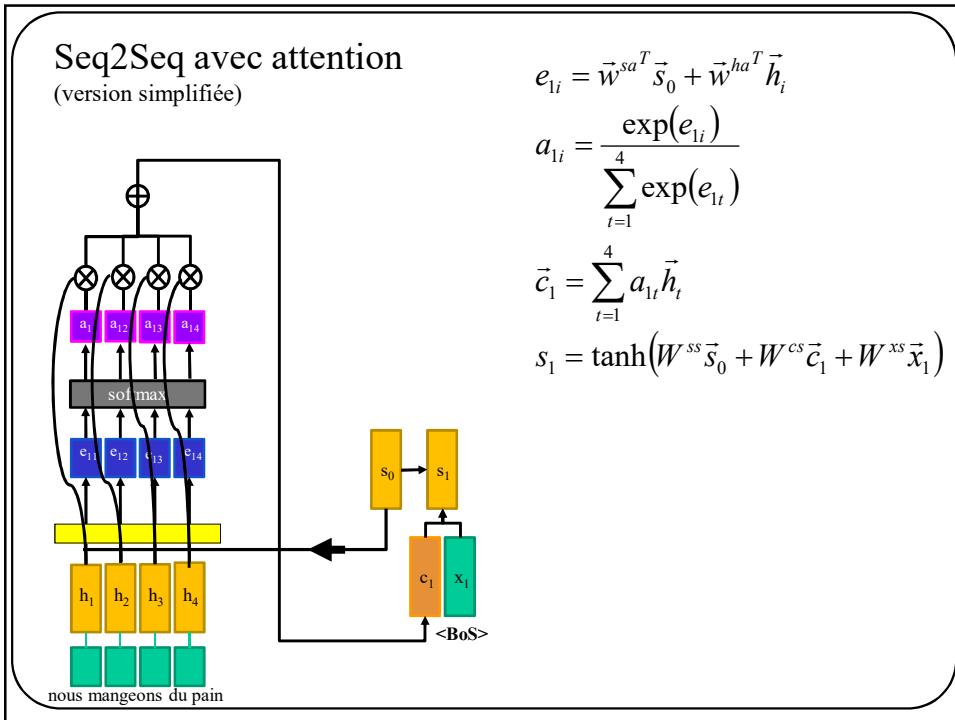


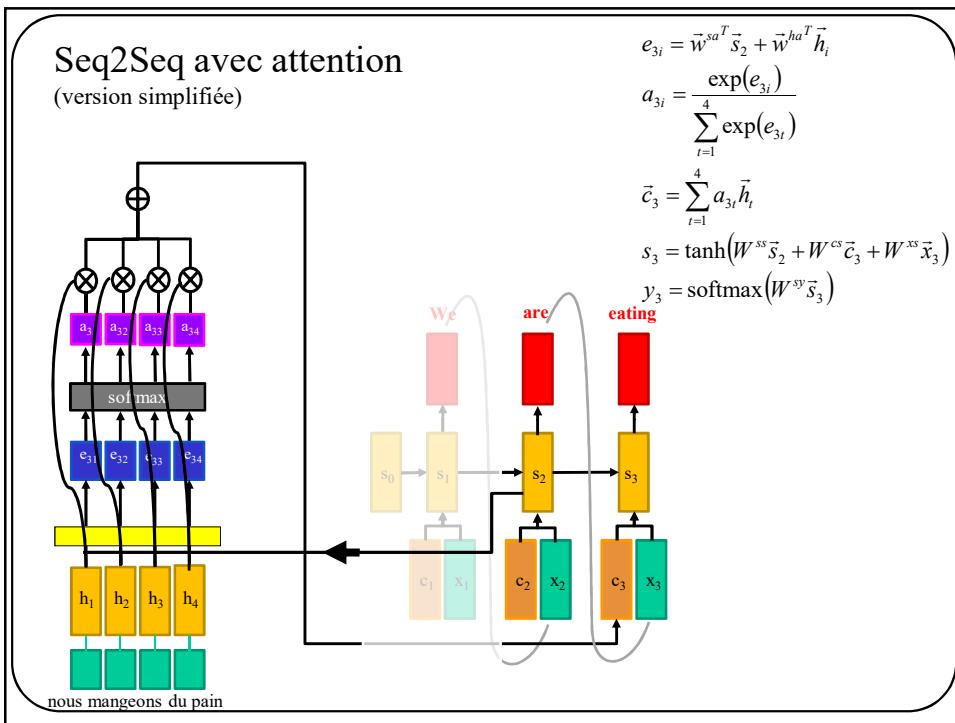
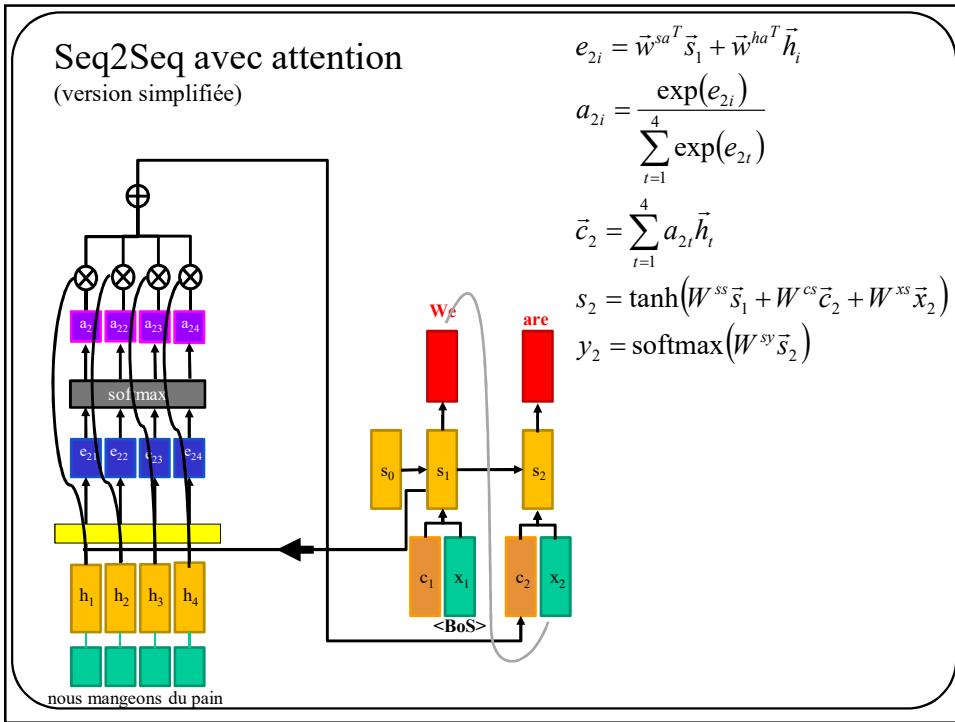


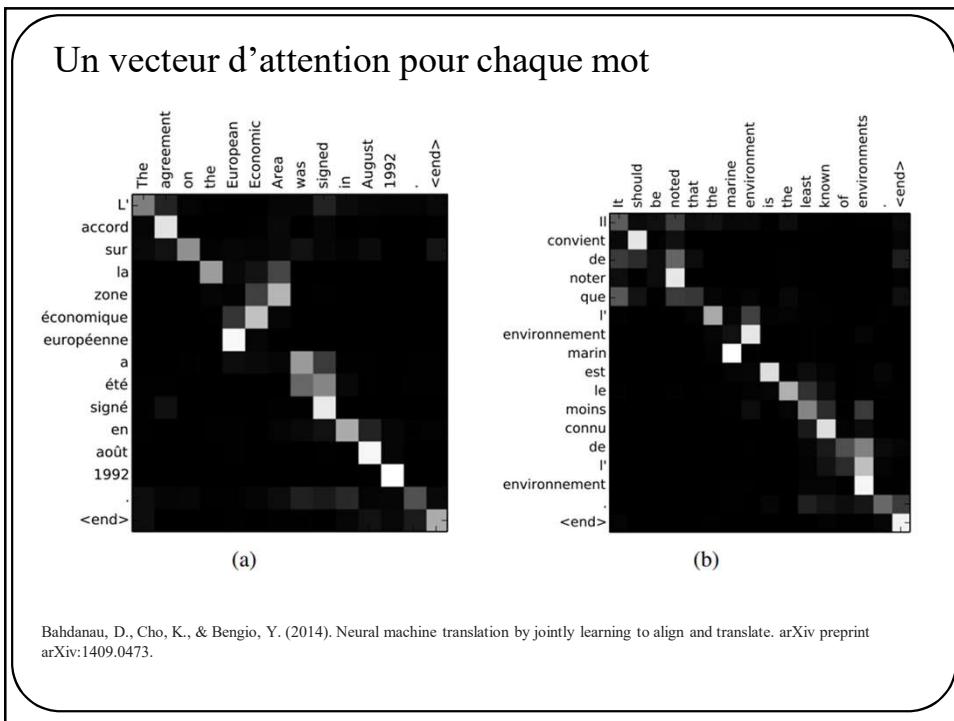
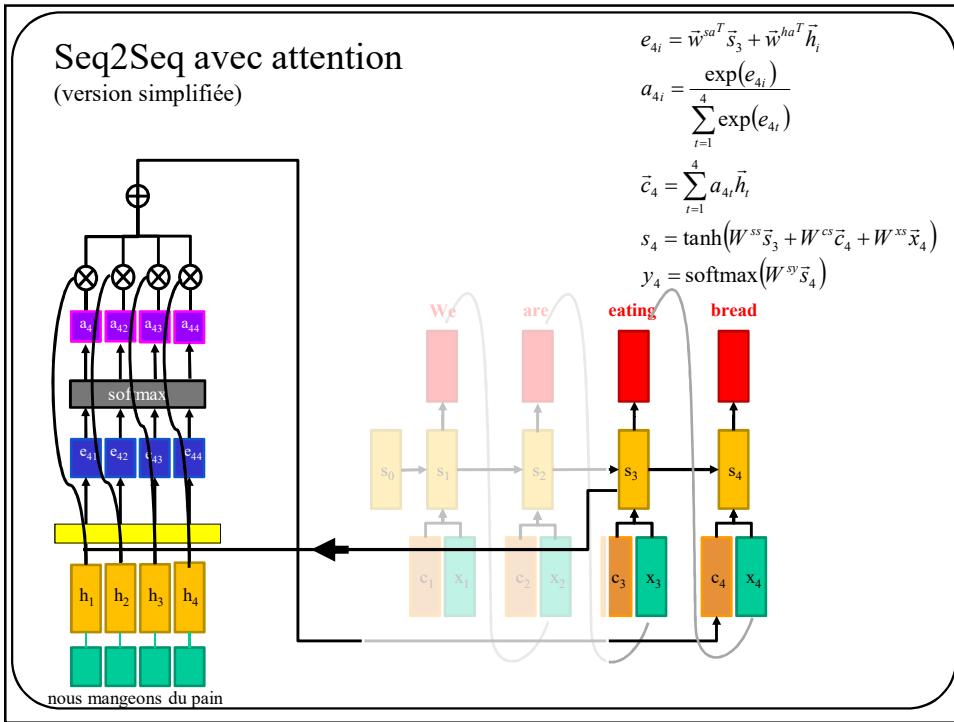




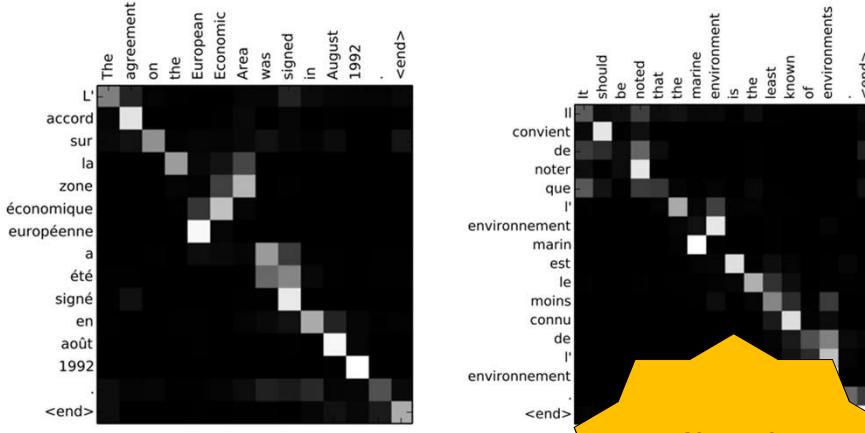








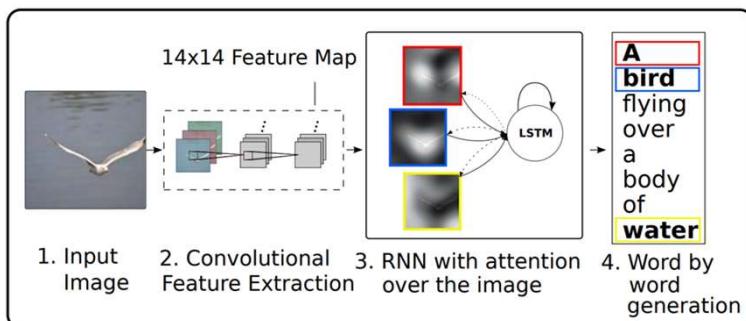
Seq2Seq avec attention



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

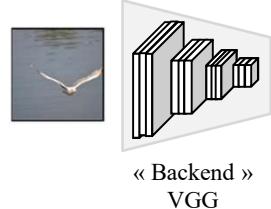
Modèles d'attention pour la description d'images

Réseau récurrent pour du captioning capable de « concentrer son attention » sur les zones de l'image associés aux mots.



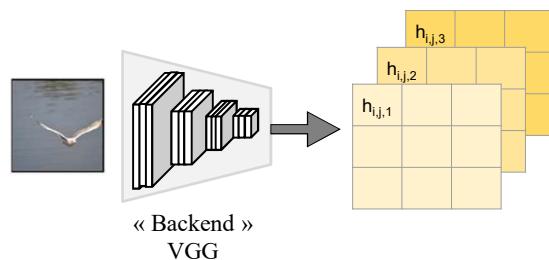
Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., ... & Bengio, Y. (2015, June). Show, attend and tell: Neural image caption generation with visual attention. In International conference on machine learning (pp. 2048-2057). PMLR.

Modèles d'attention pour la description d'images (version simplifiée)

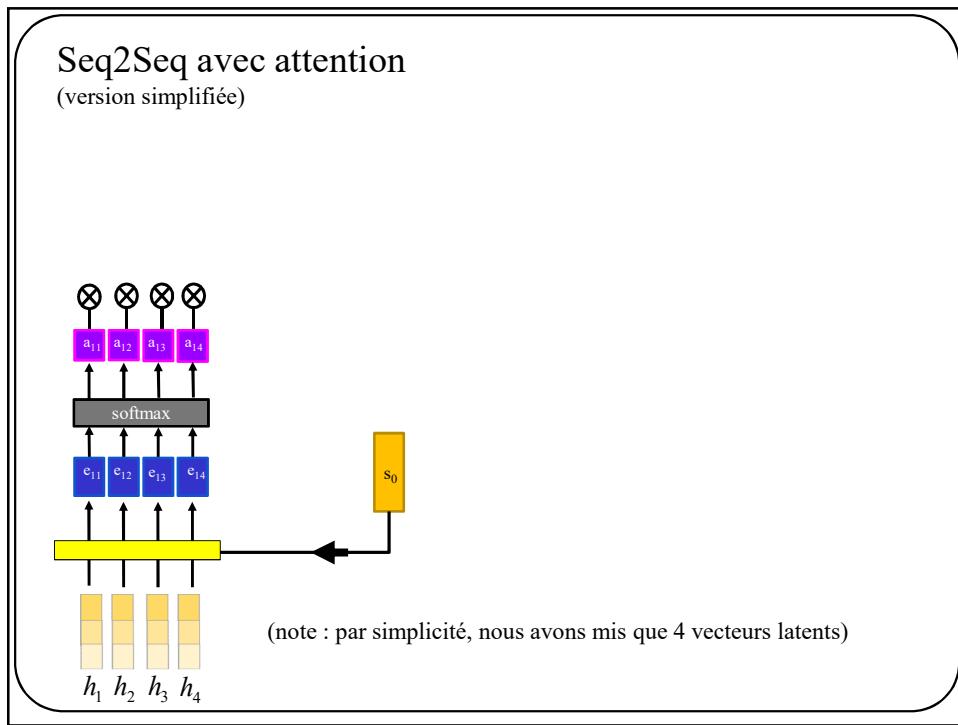
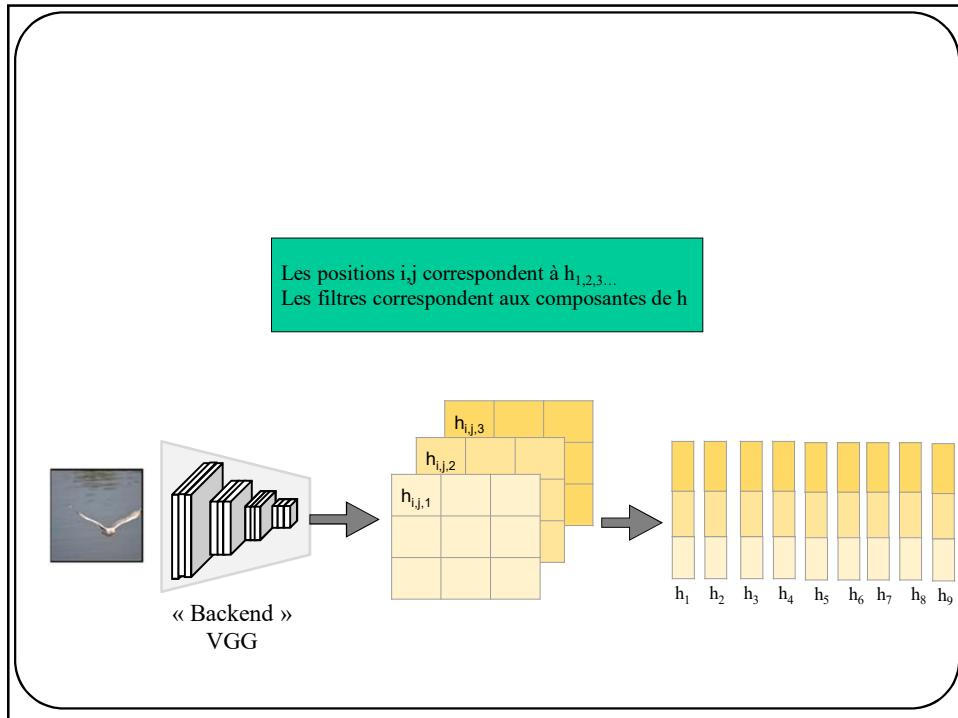


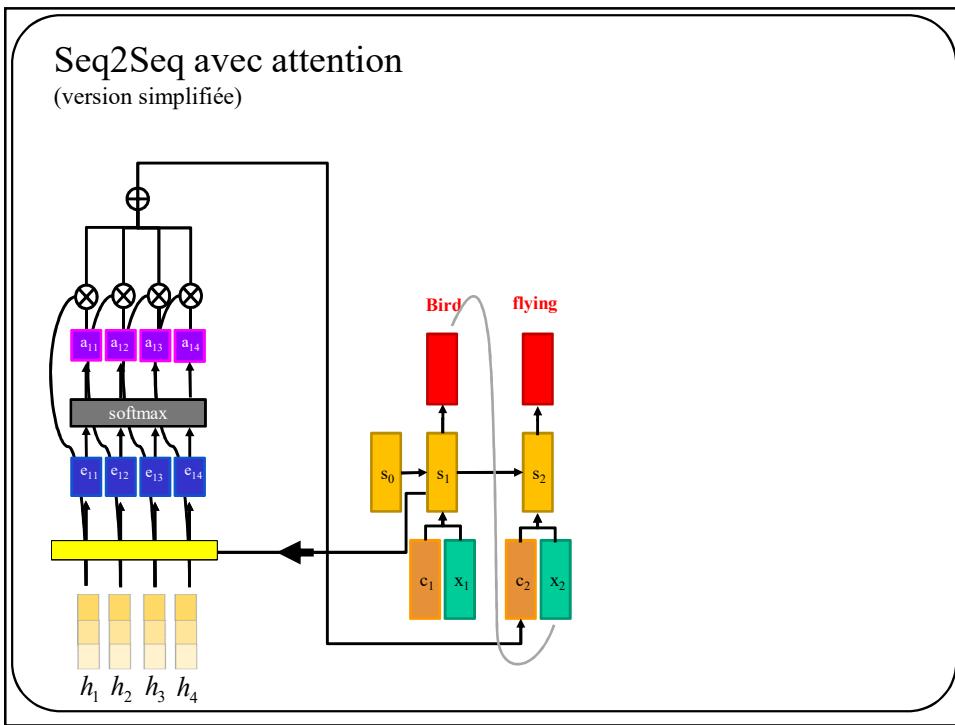
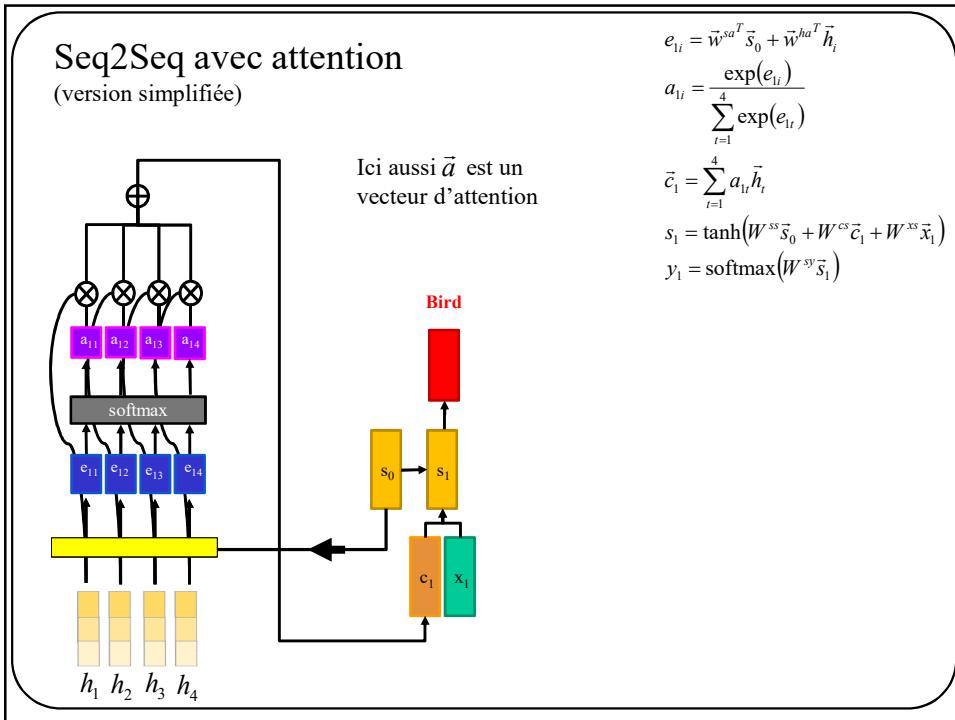
« Backend »
VGG

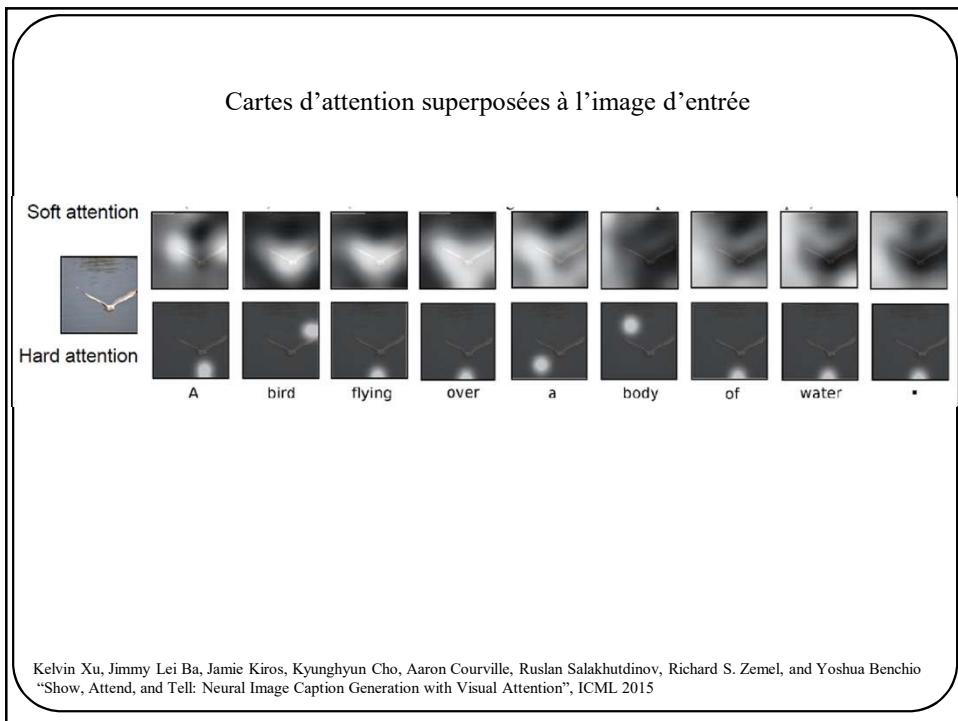
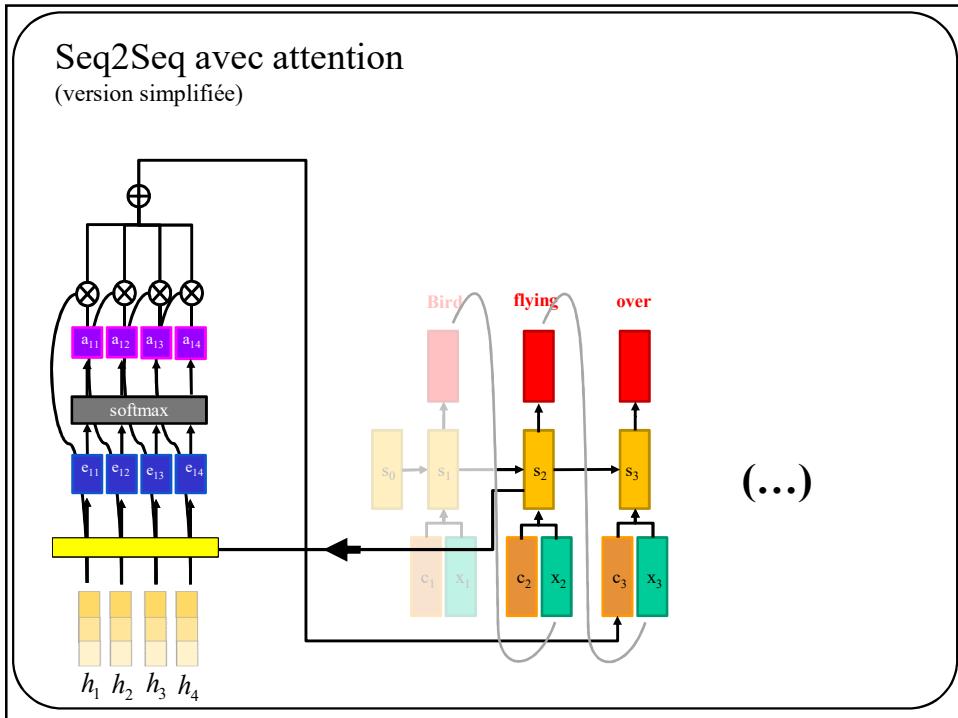
198

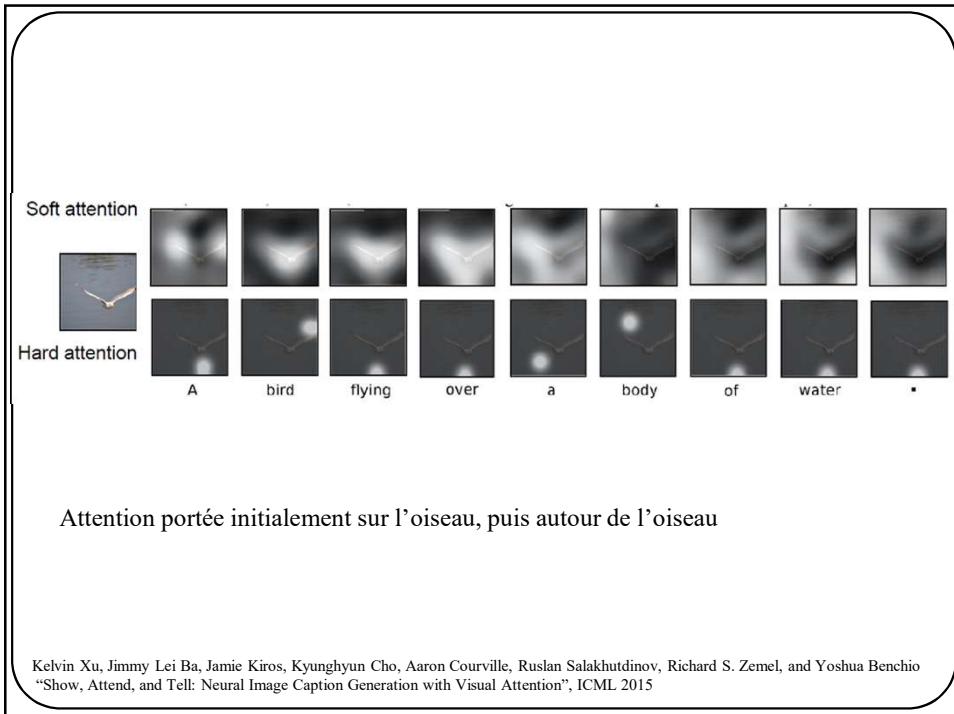


« Backend »
VGG









L'auto-attention (*self attention*)

208

Revenons à la base : multiplication matricielle

Considérons les 4 matrices suivantes

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

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

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

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

209

Revenons à la base : multiplication matricielle

Leur multiplication donne:

$$\begin{aligned}
 X &= \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in R^{3x4} & W^q X = Q &= \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} & Q_{14} \\ Q_{21} & Q_{22} & Q_{23} & Q_{24} \\ Q_{31} & Q_{32} & Q_{33} & Q_{34} \end{pmatrix} \in R^{3x4} \\
 W^q &= \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in R^{3x3} & W^K X = K &= \begin{pmatrix} K^x_{11} & K^x_{12} & K^x_{13} & K^x_{14} \\ K^x_{21} & K^x_{22} & K^x_{23} & K^x_{24} \\ K^x_{31} & K^x_{32} & K^x_{33} & K^x_{34} \end{pmatrix} \in R^{3x4} \\
 W^k &= \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in R^{3x3} & W^V X = V &= \begin{pmatrix} V_{11}^x & V_{12}^x & V_{13}^x & V_{14}^x \\ V_{21}^x & V_{22}^x & V_{23}^x & V_{24}^x \end{pmatrix} \in R^{2x4} \\
 W^V &= \begin{pmatrix} W^V_{11} & W^V_{12} & W^V_{13} \\ W^V_{21} & W^V_{22} & W^V_{23} \end{pmatrix} \in R^{2x3}
 \end{aligned}$$

210

Auto attention

X est une matrice de données pour laquelle chaque colonne i correspond à un vecteur en entrée \vec{x}_i

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

Dans cet exemple, 4 mots en entrée donc 4 colonnes dans X
Les vecteurs 3D peuvent être obtenus par [Word2Vec](#)

211

Auto attention

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

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

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

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

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

212

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)

213

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

214

$$V: \begin{pmatrix} V_{11} & V_{12} & V_{13} & V_{14} \\ V_{21} & V_{22} & V_{23} & V_{24} \end{pmatrix}$$

Q,K,V contiennent une transformation linéaire de la matrice d'entrée X. Chaque mot **Xi** a été transformé en des vecteurs **Qi**, **Ki** et **Vi**

$$W^v X$$

$$K: \begin{pmatrix} K_{11} & K_{12} & K_{13} & K_{14} \\ K_{21} & K_{22} & K_{23} & K_{24} \\ K_{31} & K_{32} & K_{33} & K_{34} \end{pmatrix}$$

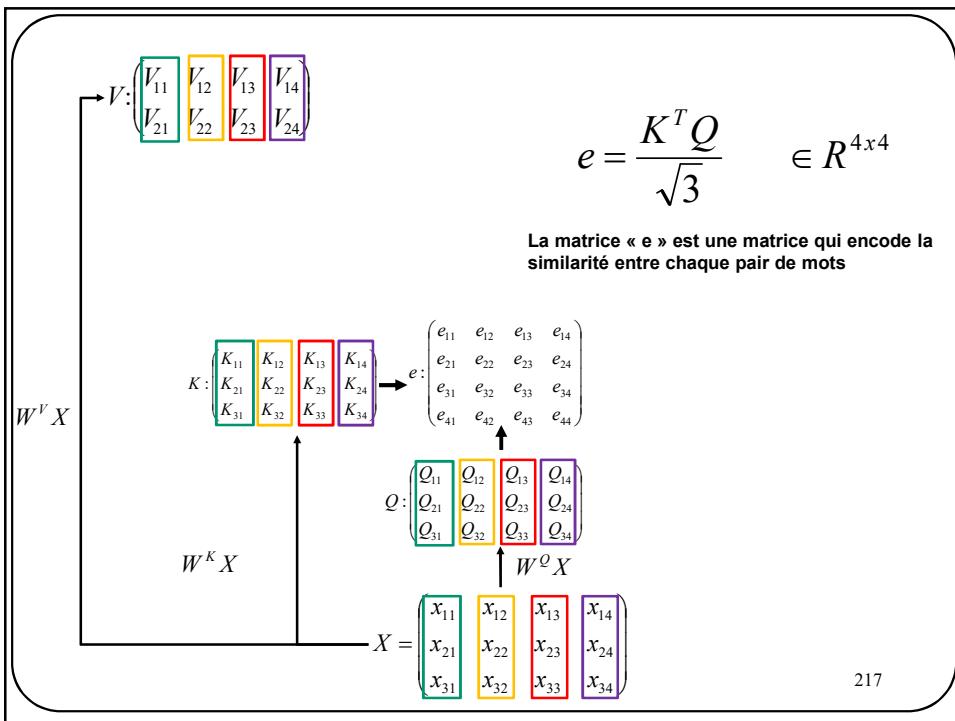
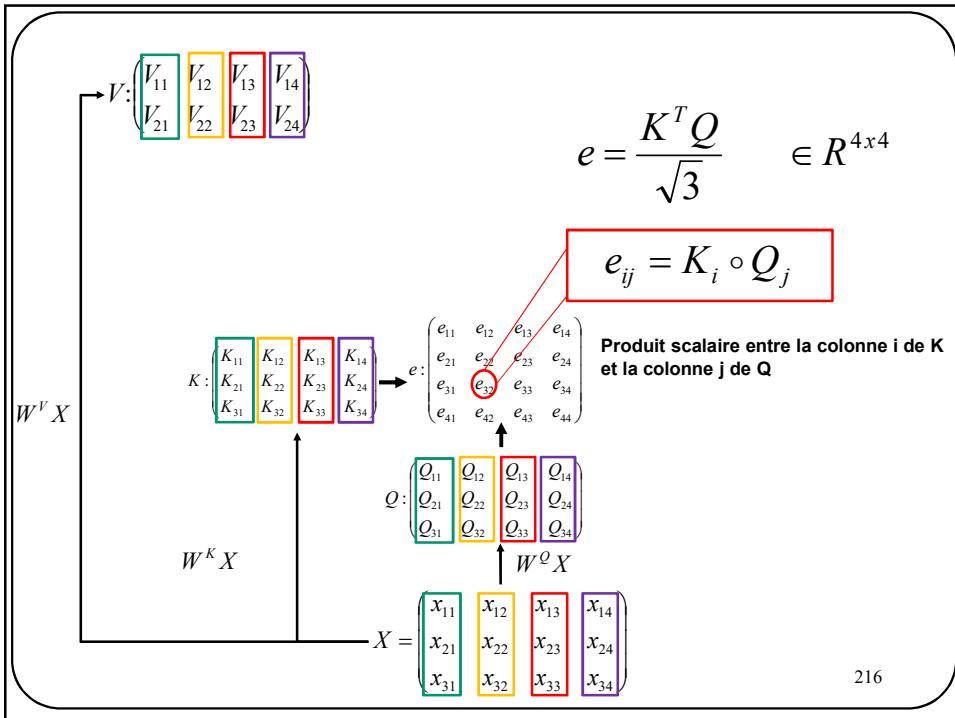
$$Q: \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} & Q_{14} \\ Q_{21} & Q_{22} & Q_{23} & Q_{24} \\ Q_{31} & Q_{32} & Q_{33} & Q_{34} \end{pmatrix}$$

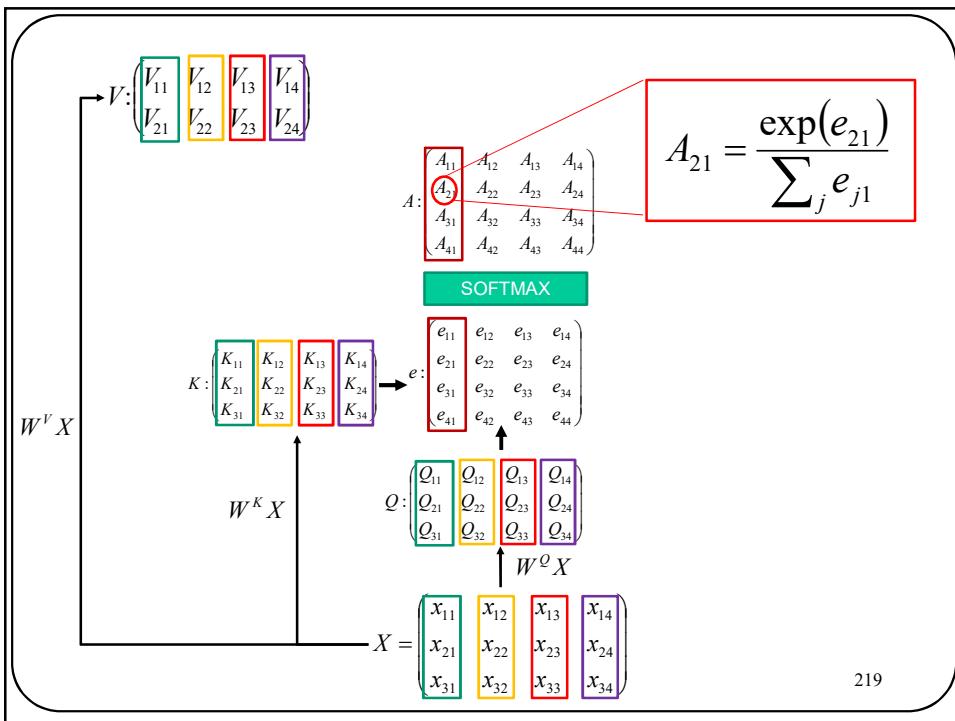
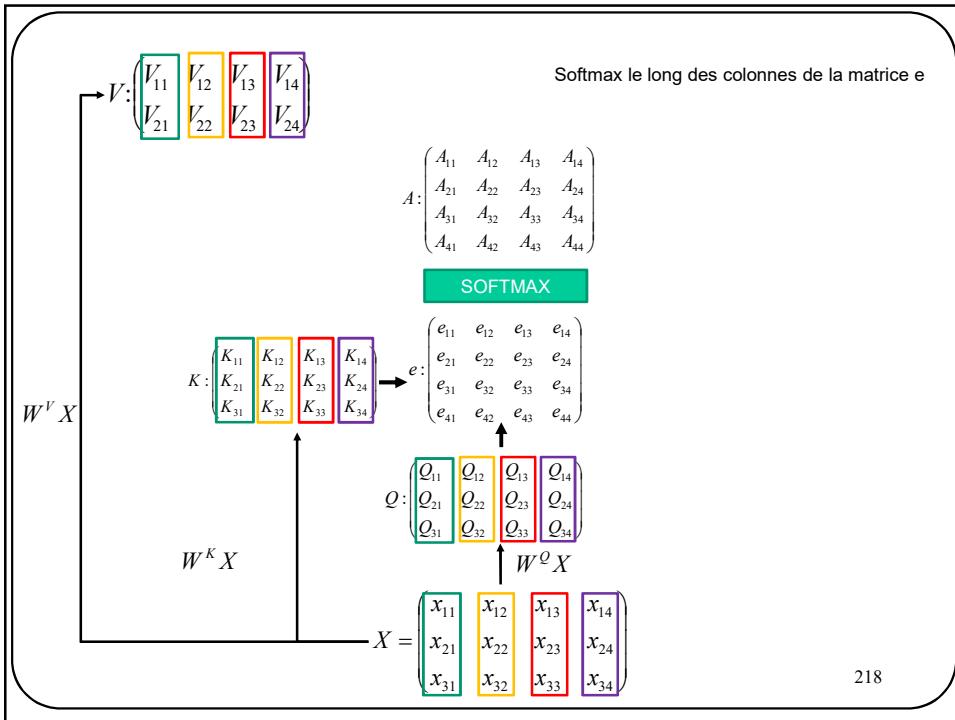
$$W^K X$$

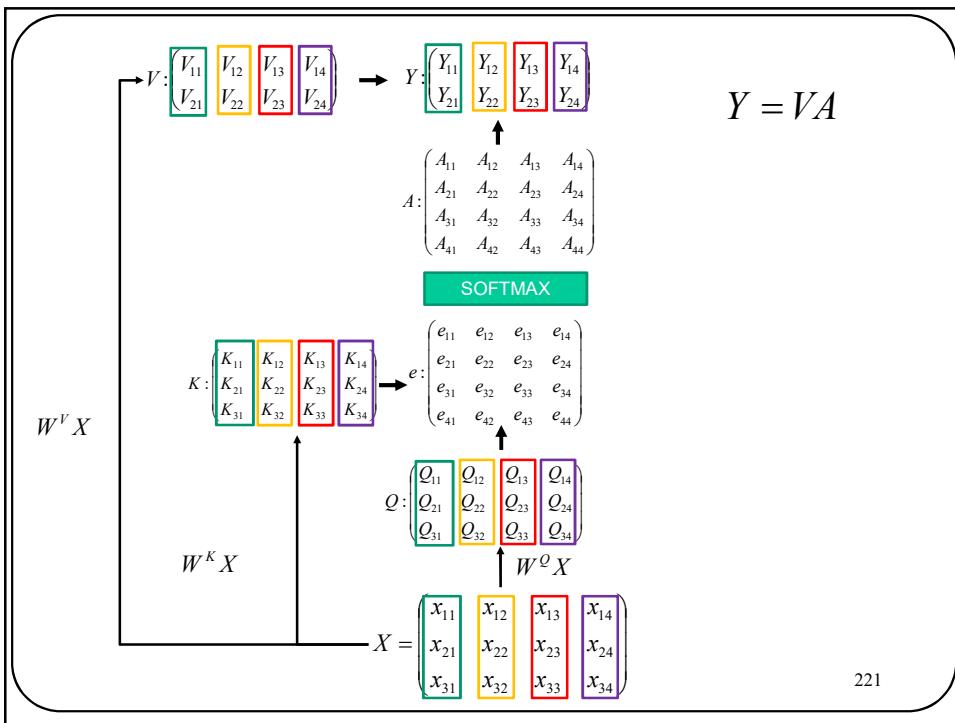
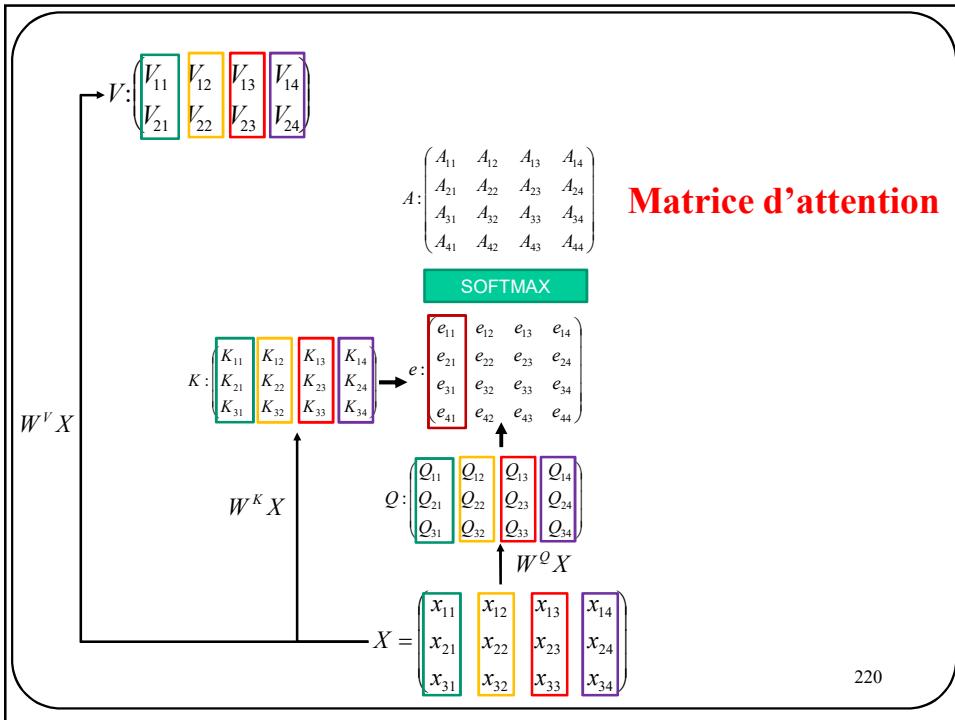
$$W^Q X$$

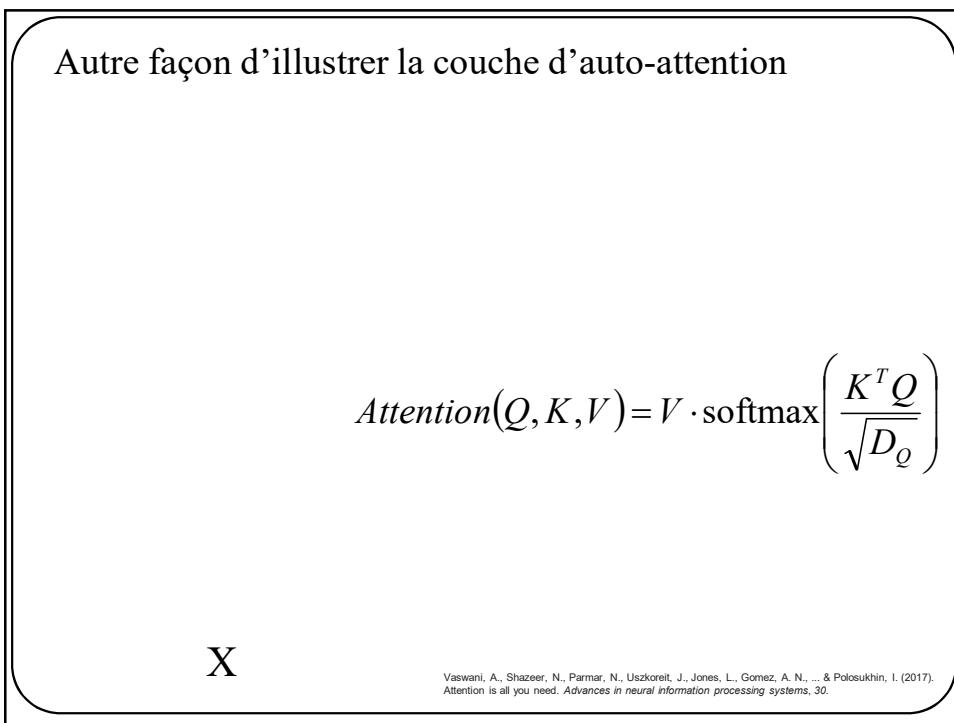
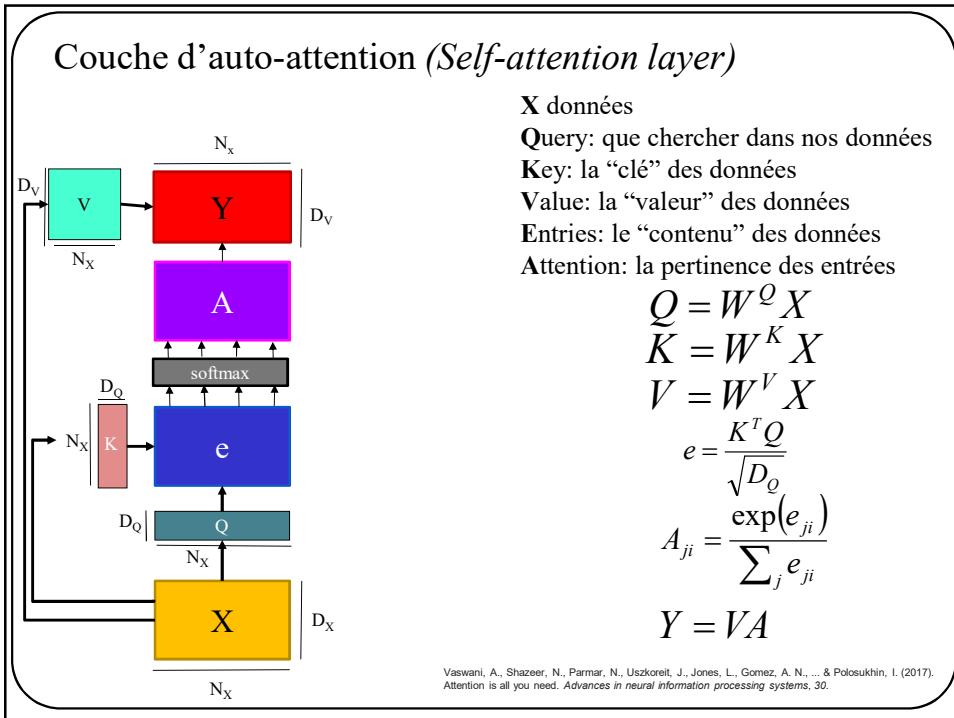
$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix}$$

215



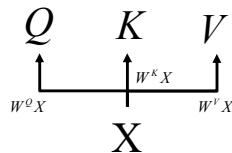






Autre façon d'illustrer la couche d'auto-attention

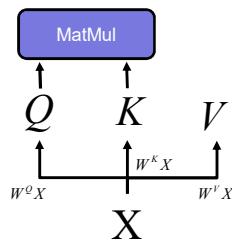
$$\text{Attention}(Q, K, V) = V \cdot \text{softmax} \left(\frac{K^T Q}{\sqrt{D_Q}} \right)$$



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.

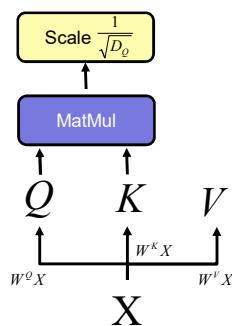
Autre façon d'illustrer la couche d'auto-attention

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



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.

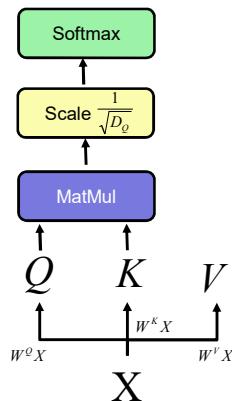
Autre façon d'illustrer la couche d'auto-attention



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

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.

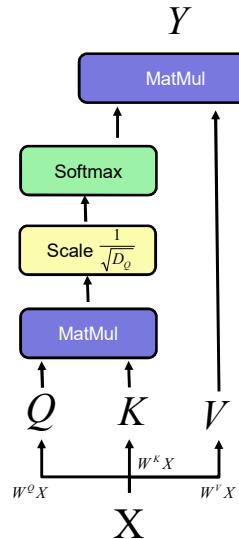
Autre façon d'illustrer la couche d'auto-attention



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

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.

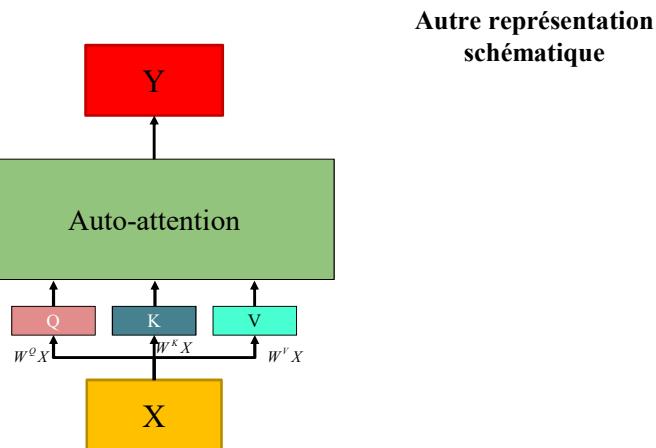
Autre façon d'illustrer la couche d'auto-attention



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

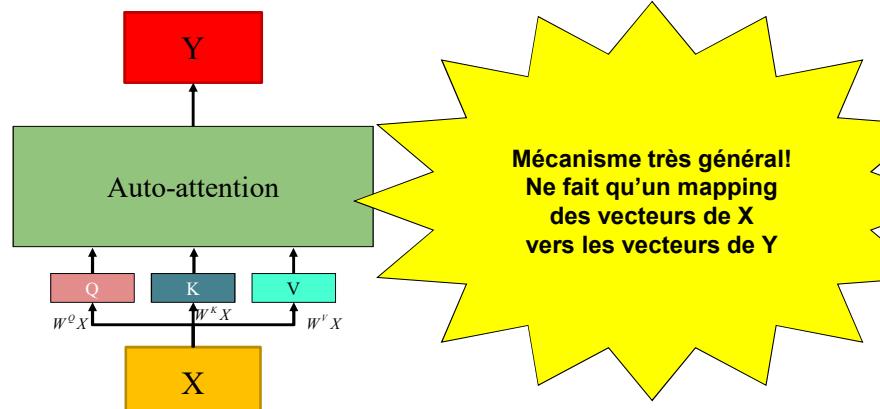
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

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



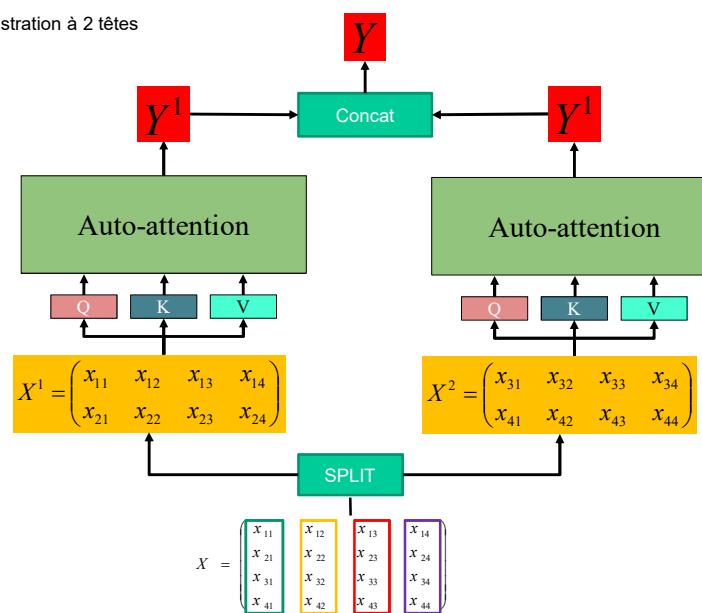
Couche d'auto-attention (*Self-attention layer*)

Autre représentation schématique

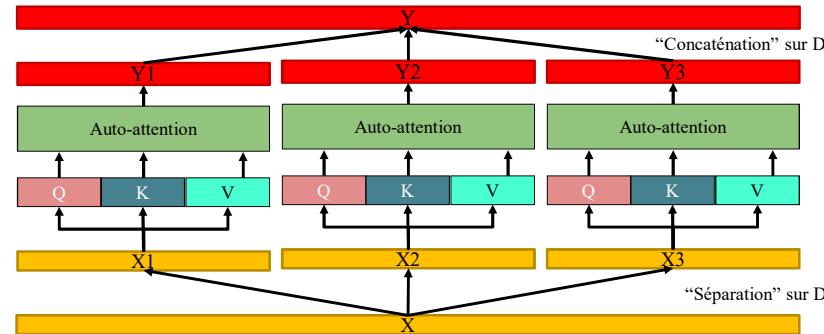


Auto-attention multi-têtes (*Multi-head Self-attention*)

Illustration à 2 têtes



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.

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.

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.

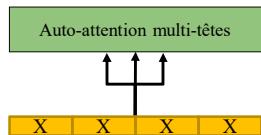
Transformer (Attention is all you need)



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.

Transformer (*Attention is all you need*)

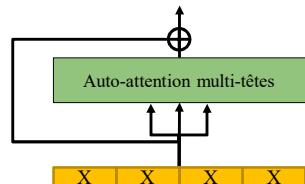
- Auto-attention multi-têtes sur les dimensions de X



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.

Transformer (*Attention is all you need*)

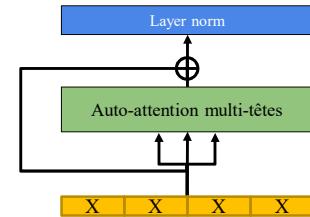
- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle



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.

Transformer (*Attention is all you need*)

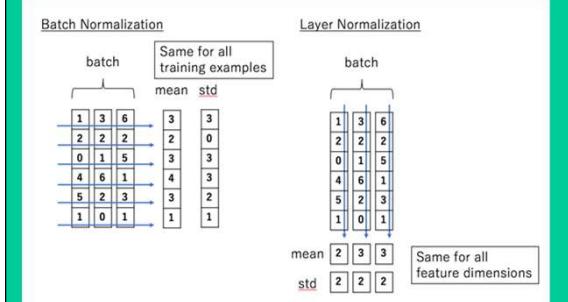
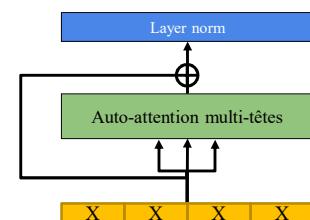
- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”



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.

Transformer (*Attention is all you need*)

$$y = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

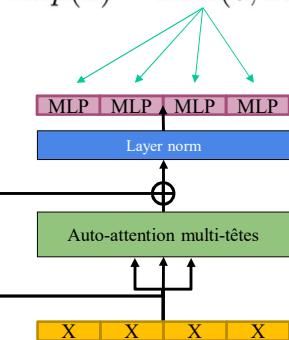


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

Transformer (*Attention is all you need*)

$$mlp(x) = \max(0, xW^1)W^2$$

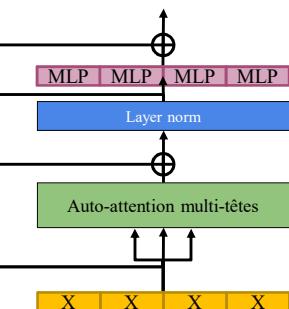
- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”
- MLP par entrée



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.

Transformer (*Attention is all you need*)

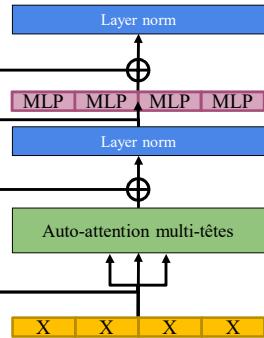
- Intra-attention multi-tête sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”
- MLP par entrée



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.

Transformer (*Attention is all you need*)

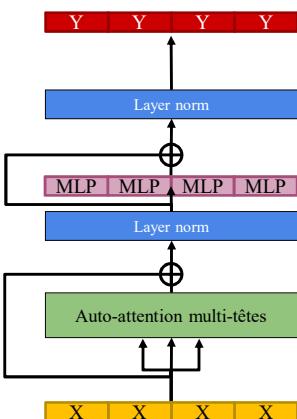
- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”
- MLP par entrée



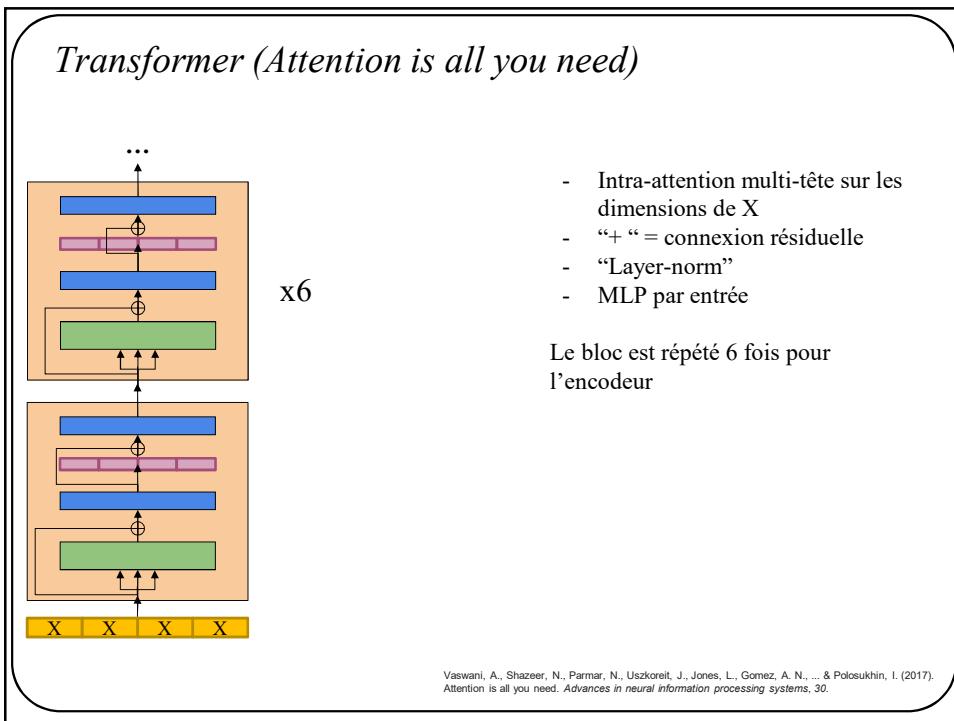
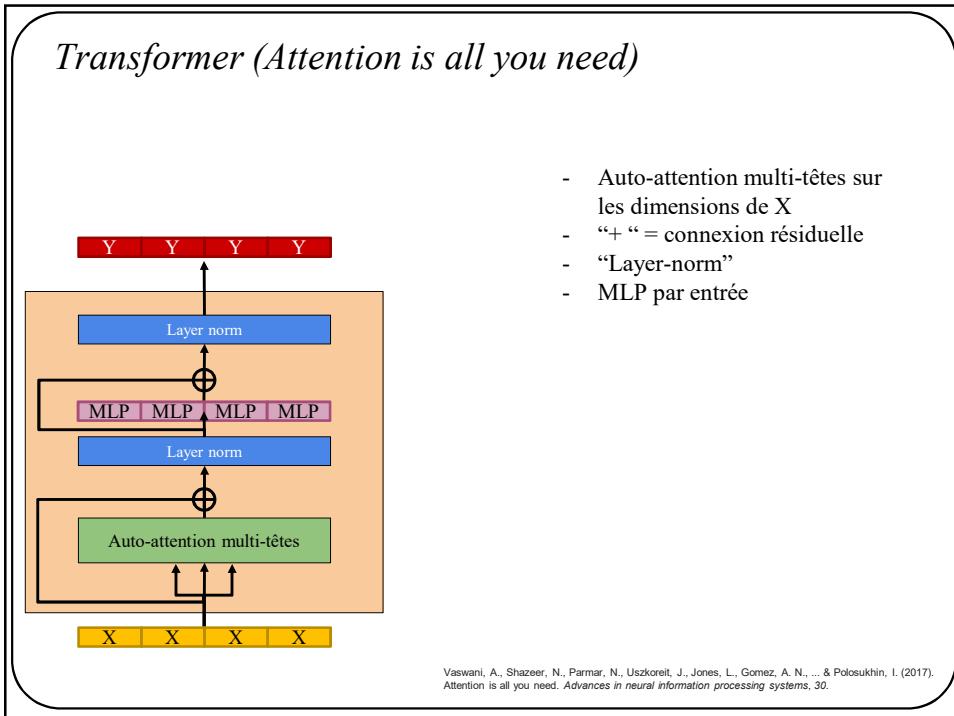
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

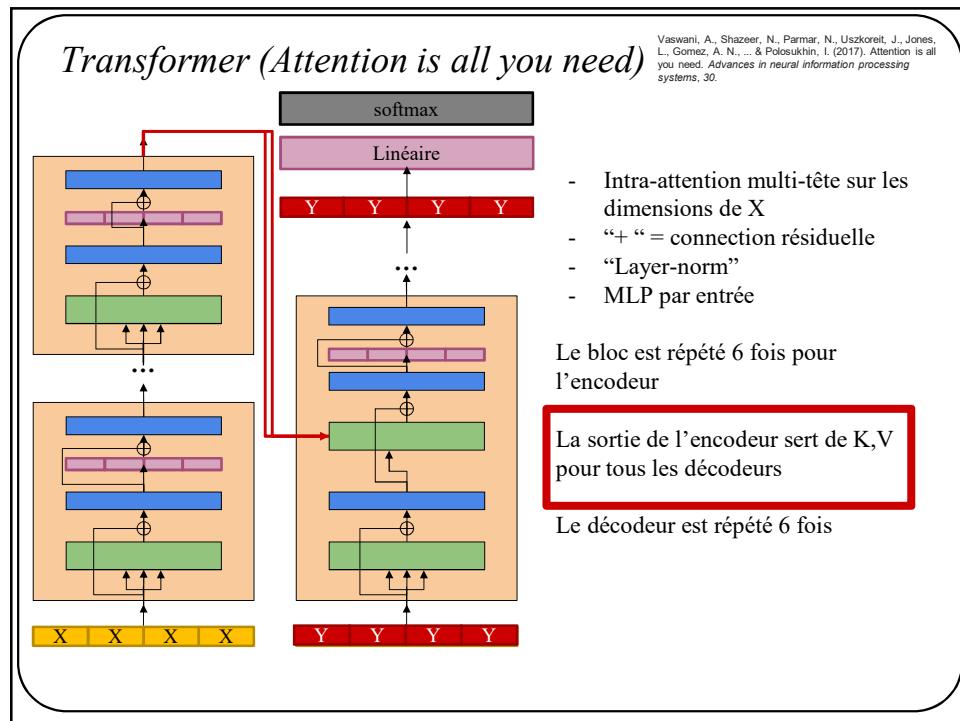
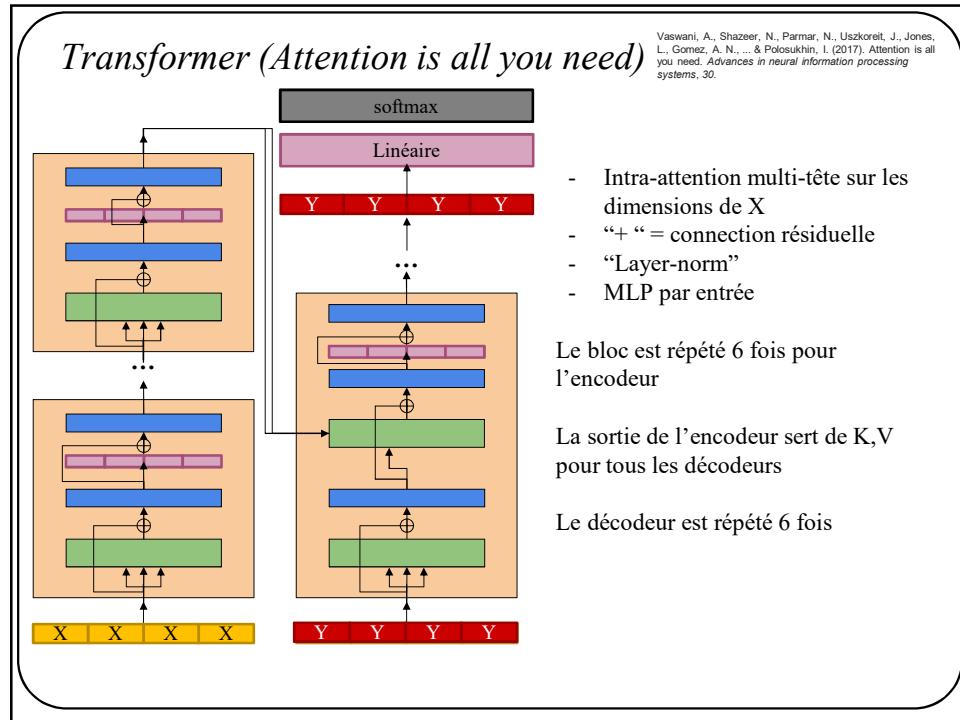
Transformer (*Attention is all you need*)

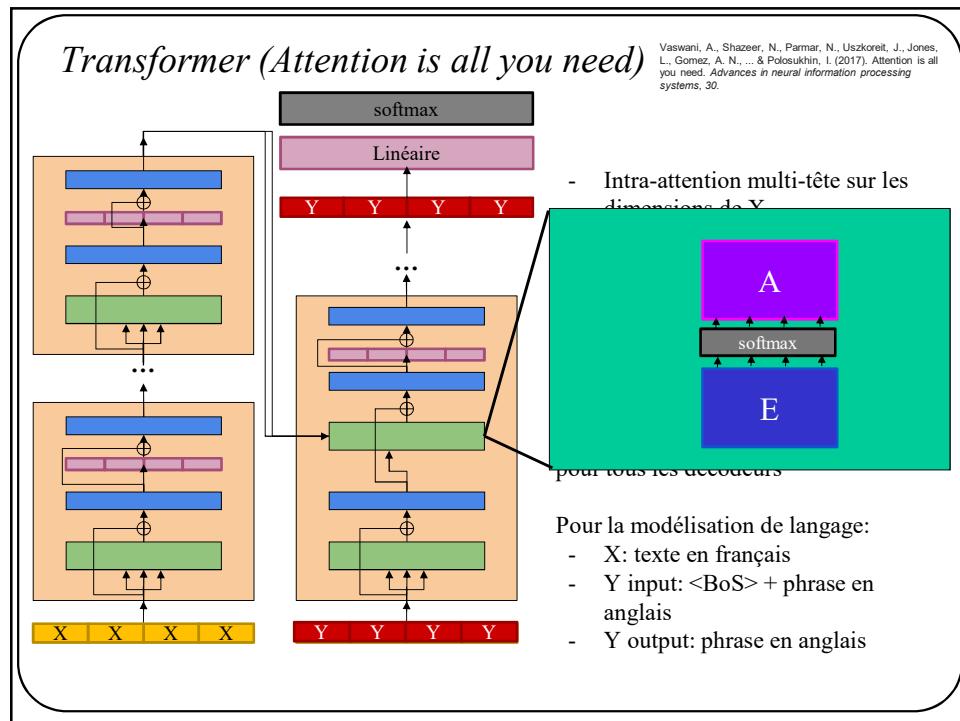
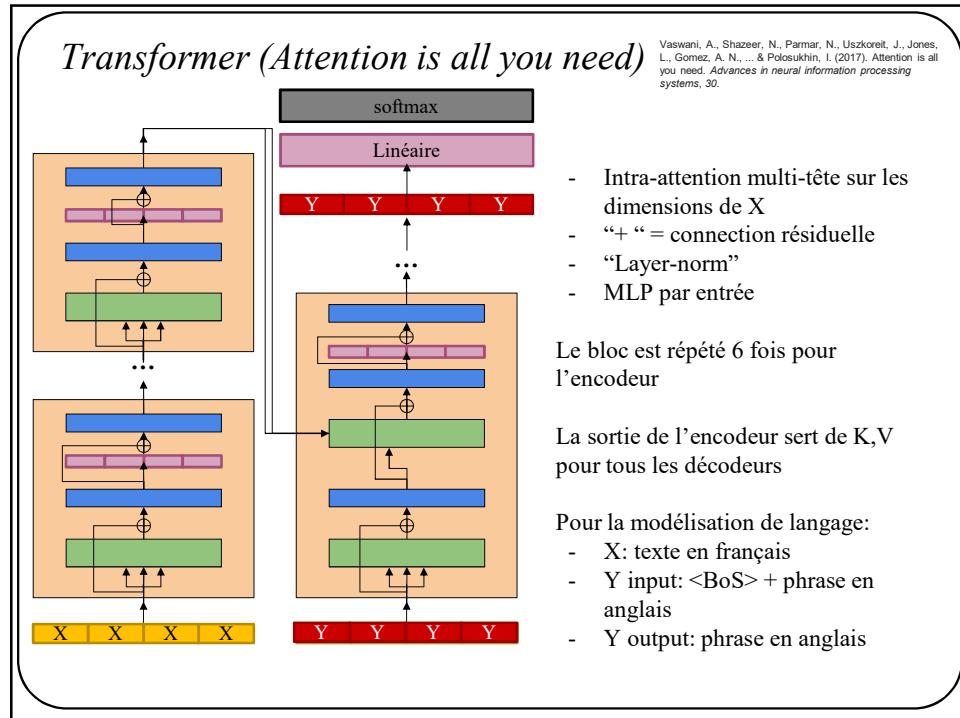
- Auto-attention multi-têtes sur les dimensions de X
- “+” = connexion résiduelle
- “Layer-norm”
- MLP par entrée

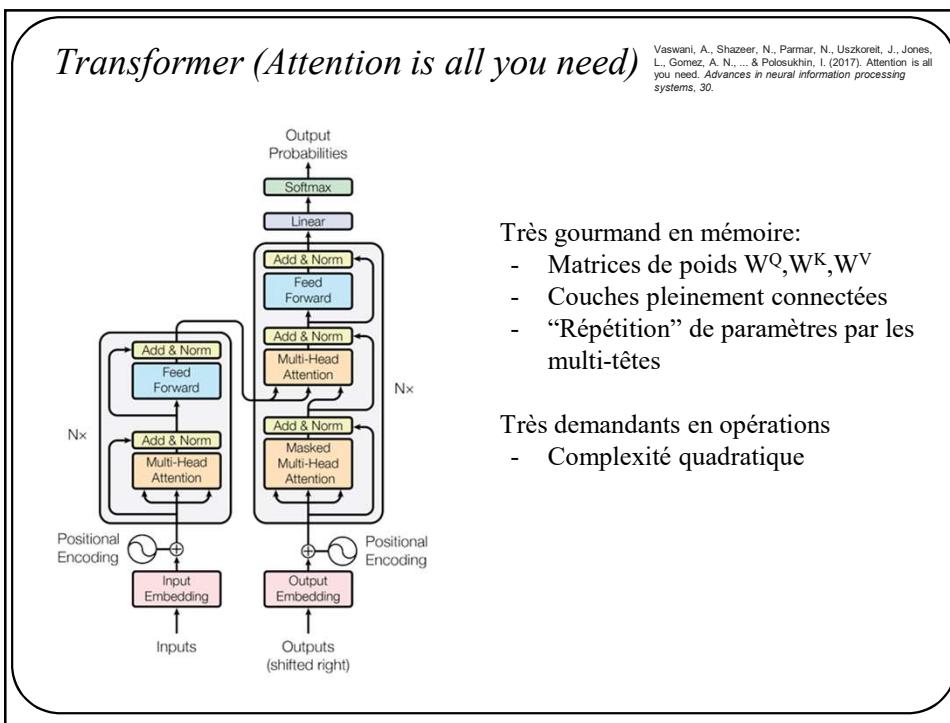
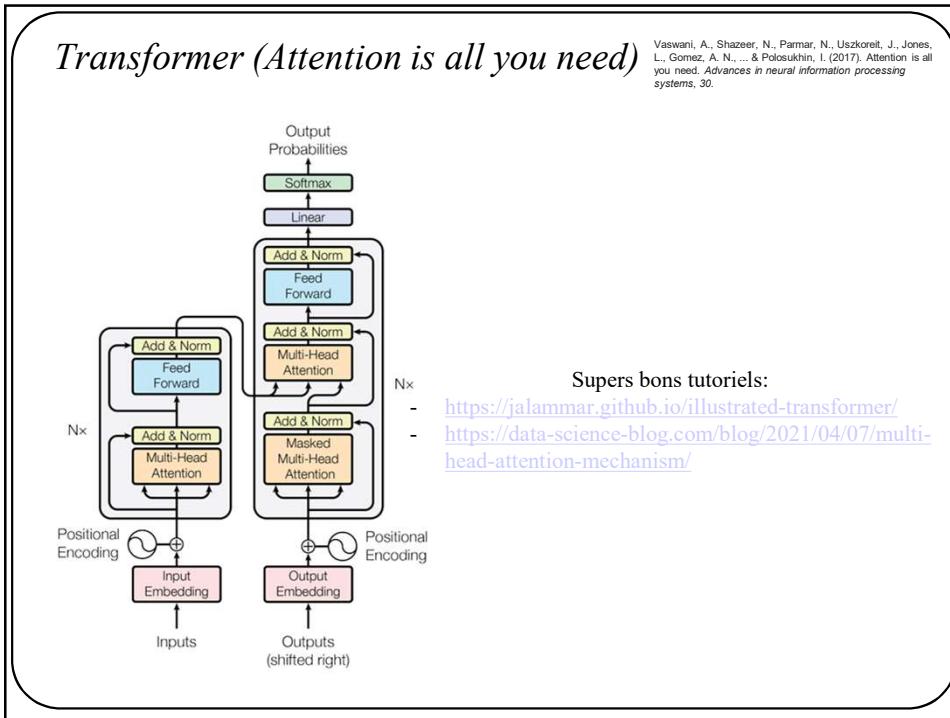


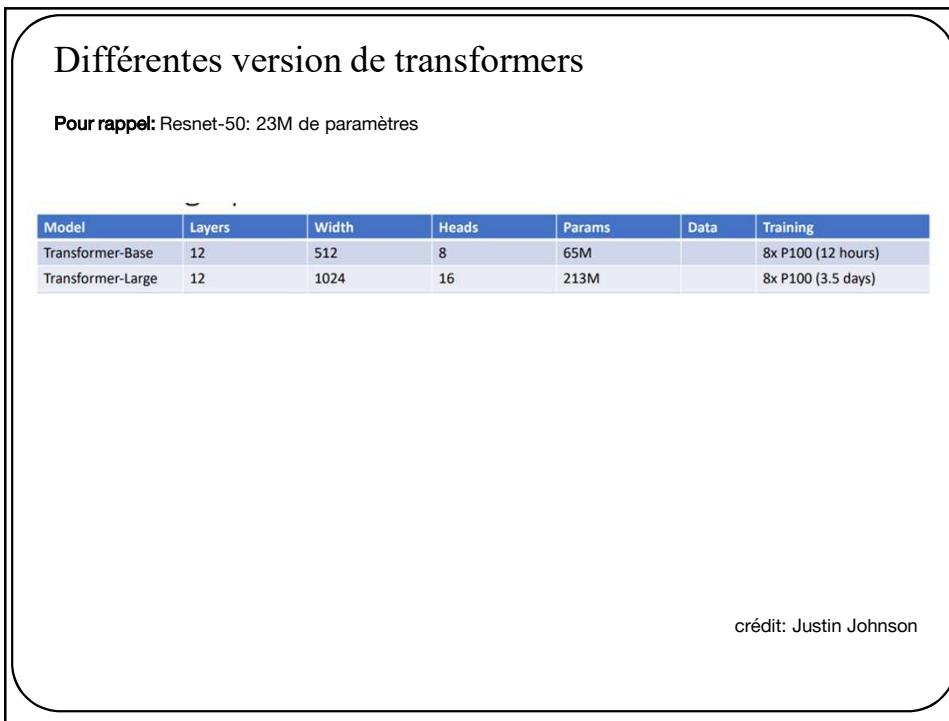
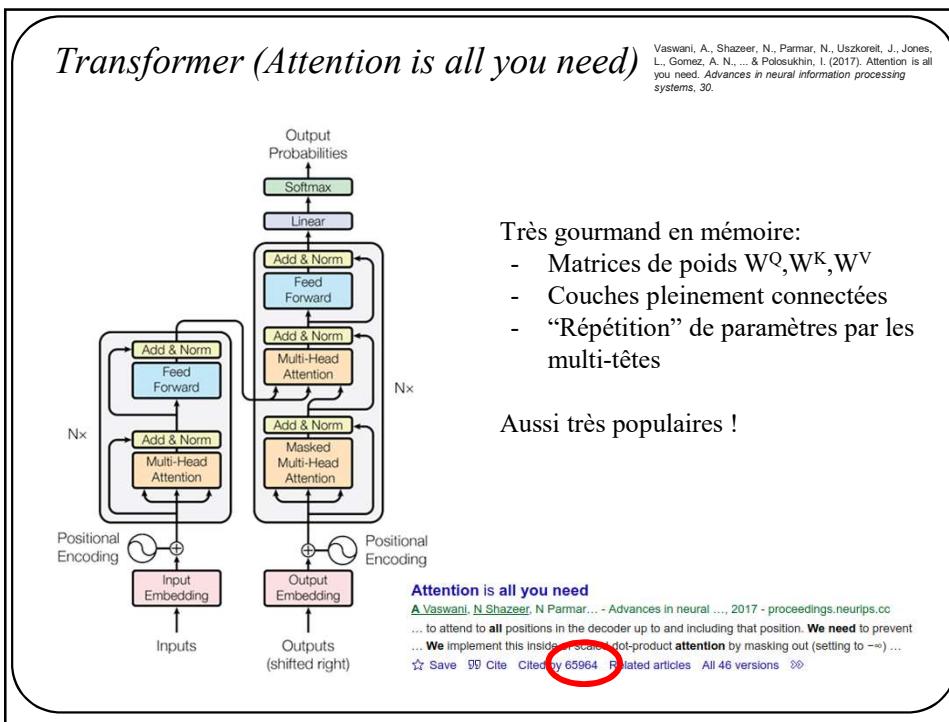
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.











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

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

Transformers

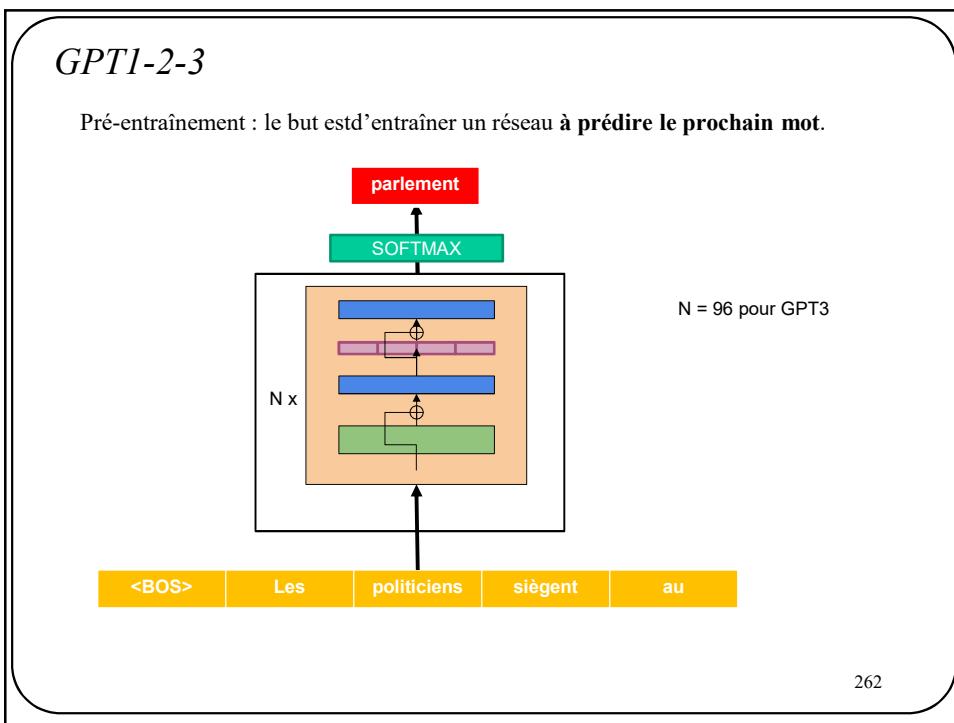
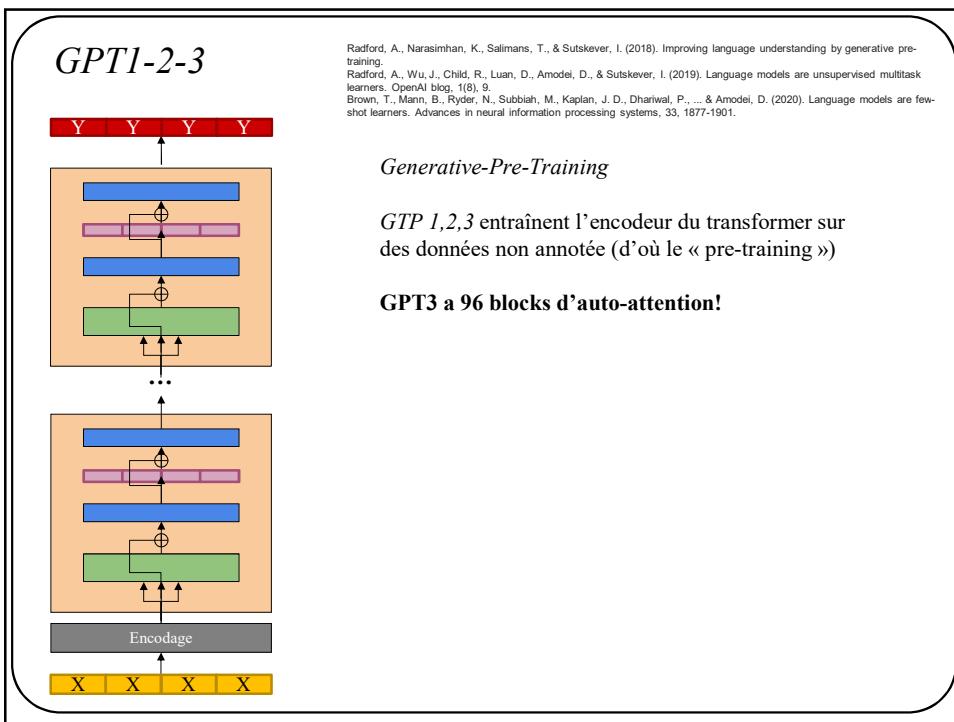
Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160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

Transformers

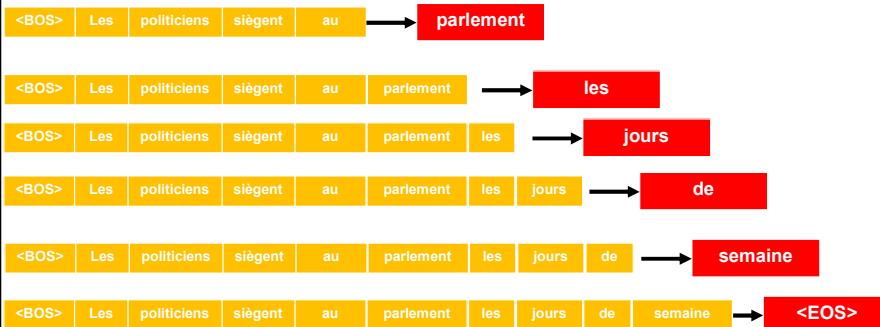
Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	213M		8x P100 (3.5 days)
BERT-Base	12	768	12	110M	13 GB	
BERT-Large	24	1024	16	340M	13 GB	
XLNet-Large	24	1024	16	~340M	126 GB	512x TPU-v3 (2.5 days)
RoBERTa	24	1024	16	355M	160 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	
Megatron-LM	40	1536	16	1.2B	174 GB	64x V100 GPU
Megatron-LM	54	1920	20	2.5B	174 GB	128x V100 GPU
Megatron-LM	64	2304	24	4.2B	174 GB	256x V100 GPU (10 days)
Megatron-LM	72	3072	32	8.3B	174 GB	512x V100 GPU (9 days)

crédit: Justin Johnson



GPT1-2-3

Et on recommence:



263

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>

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 et 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>

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-3.27B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3.67B	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.

GPTI-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.

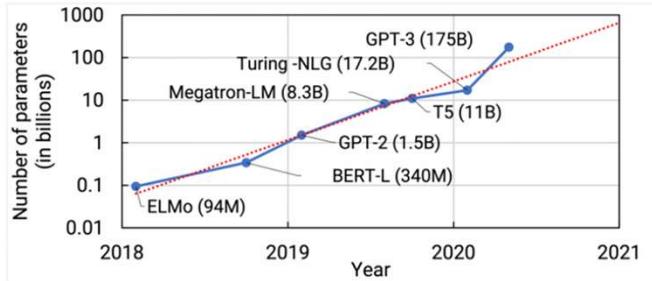


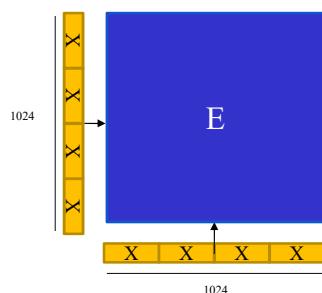
Figure 1. Trend of state-of-the-art NLP model sizes with time.

<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"

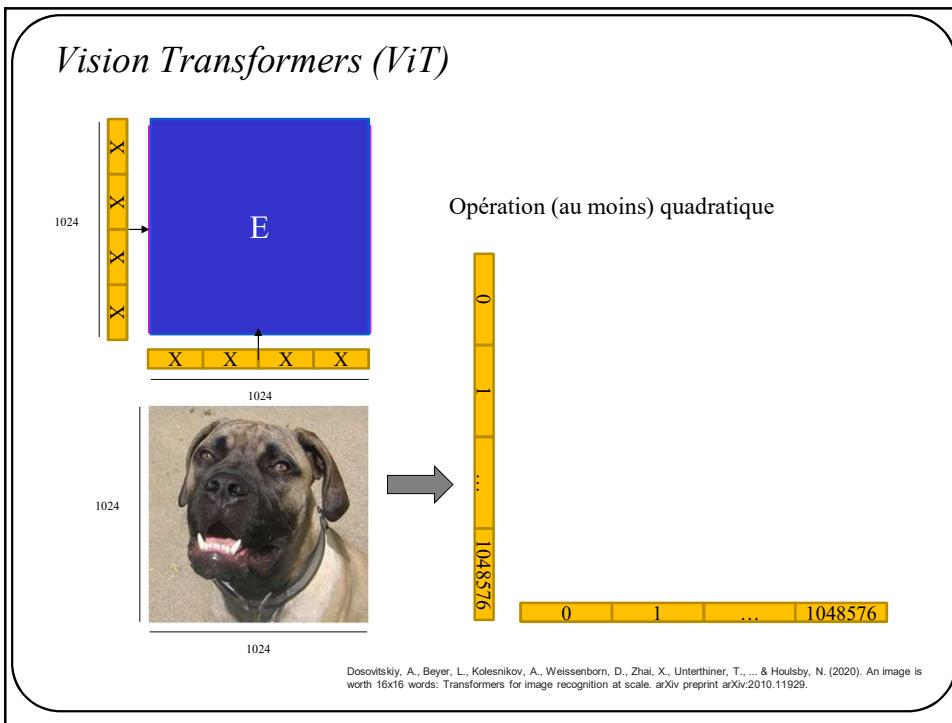
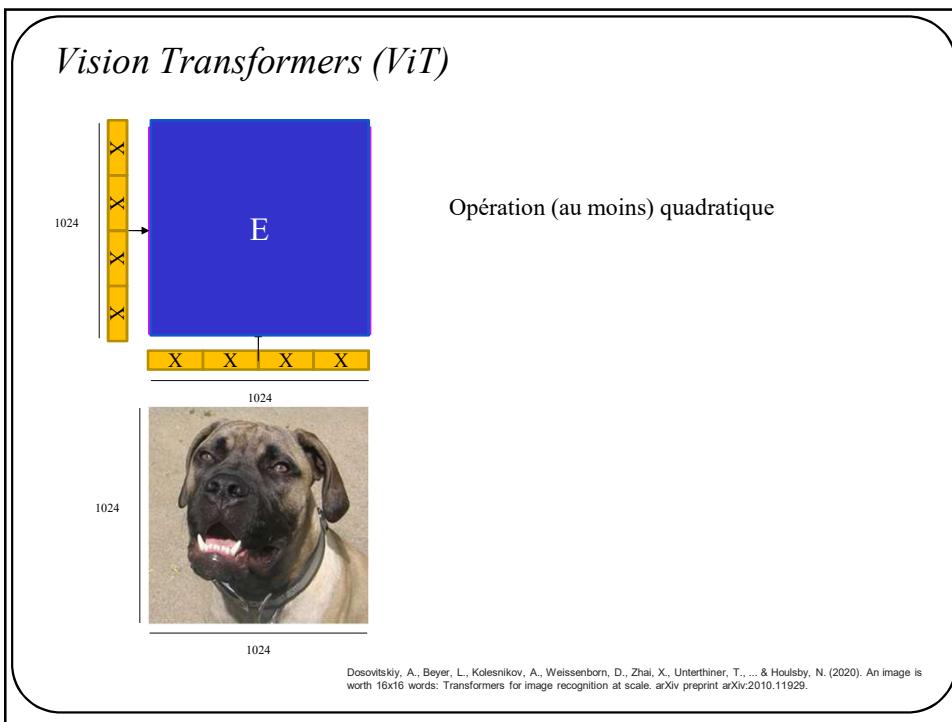
<https://bdtechtalks.com/2020/09/21/gpt-3-economy-business-model/>

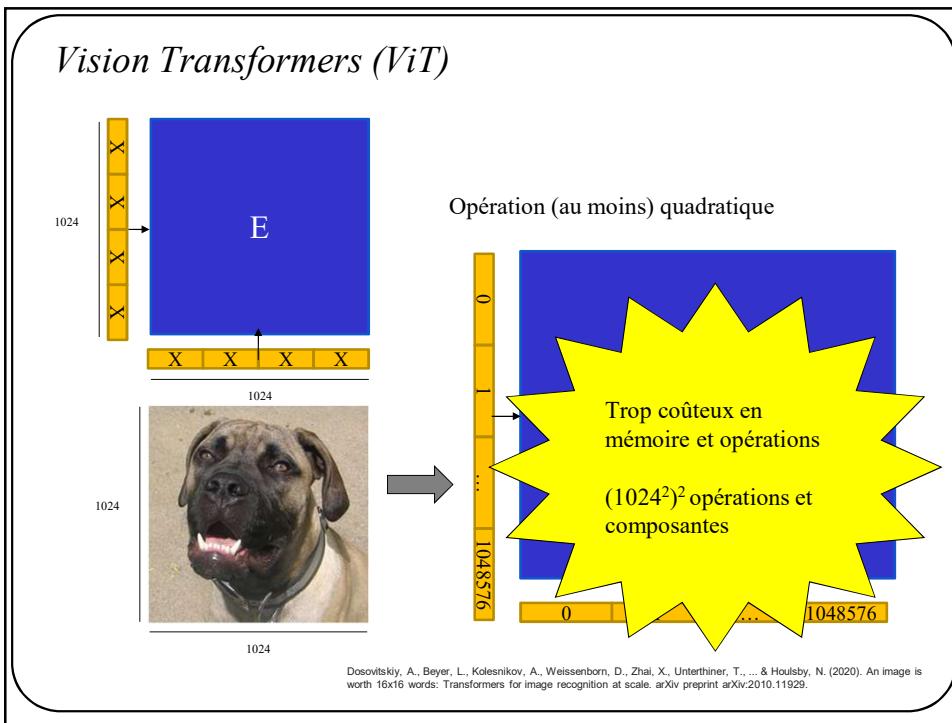
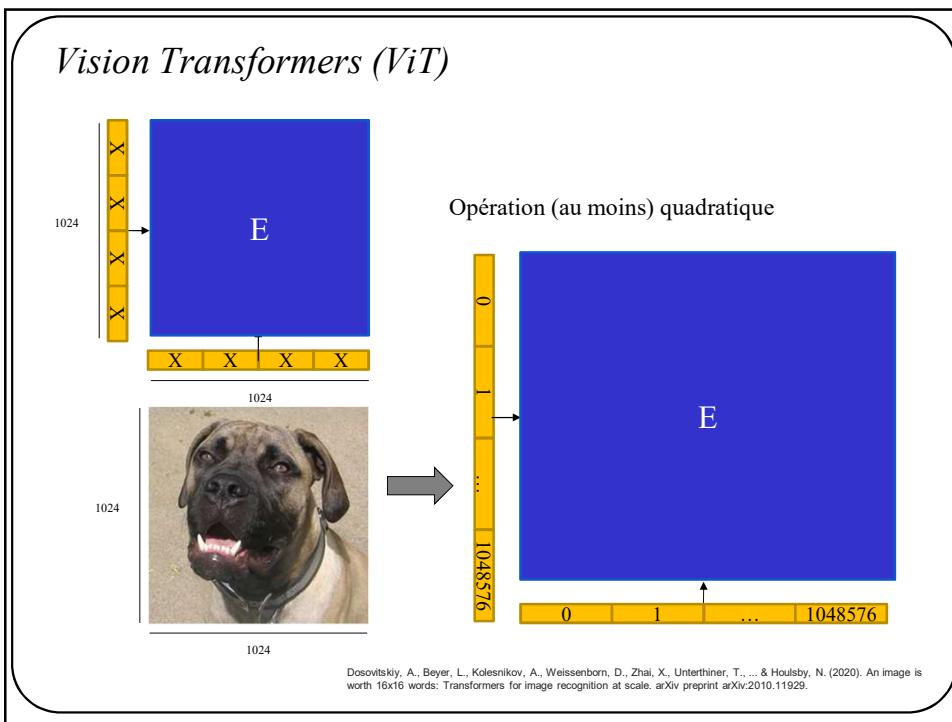
Vision Transformers (ViT)

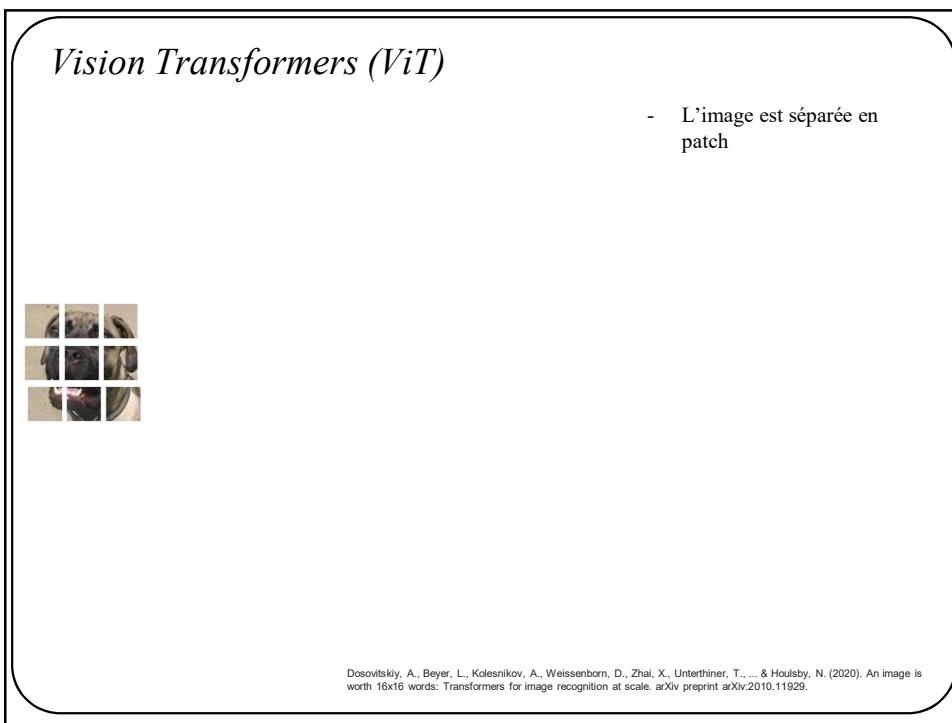
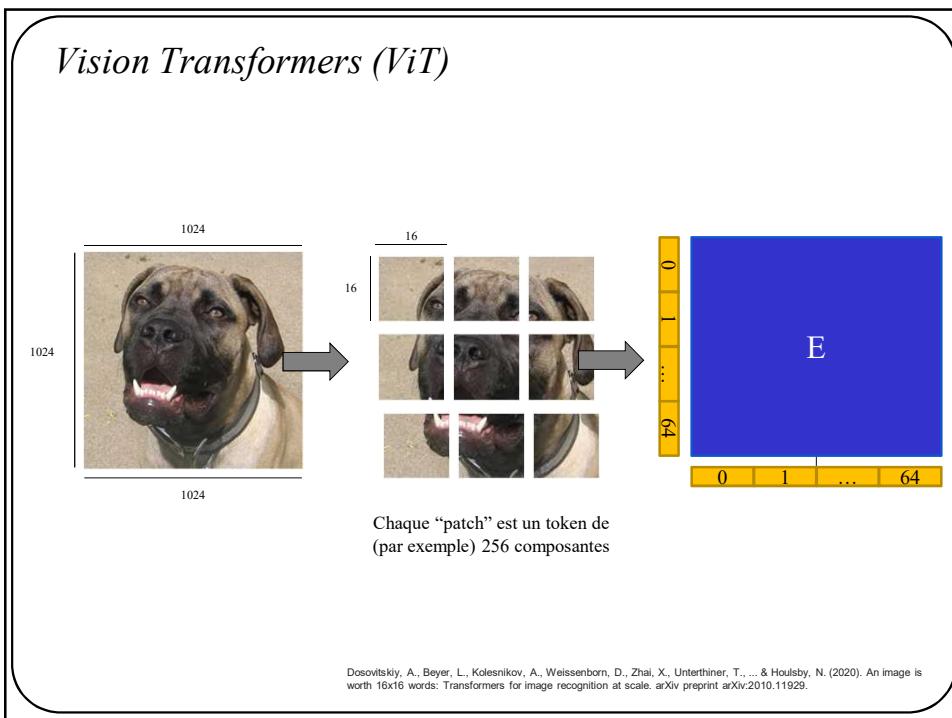


Opération (au moins) quadratique =
 1024^2 opérations et composantes à garder en mémoire

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.

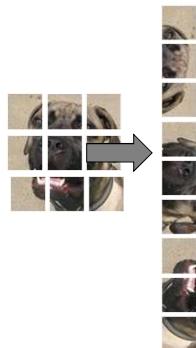






Vision Transformers (ViT)

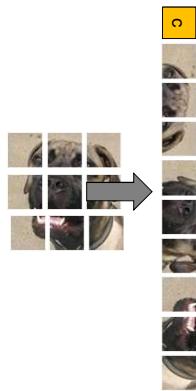
- L'image est séparée en patch
- Chaque patch est linéarisée



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.

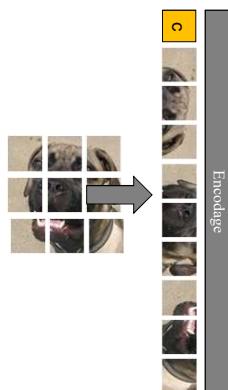
Vision Transformers (ViT)

- 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.

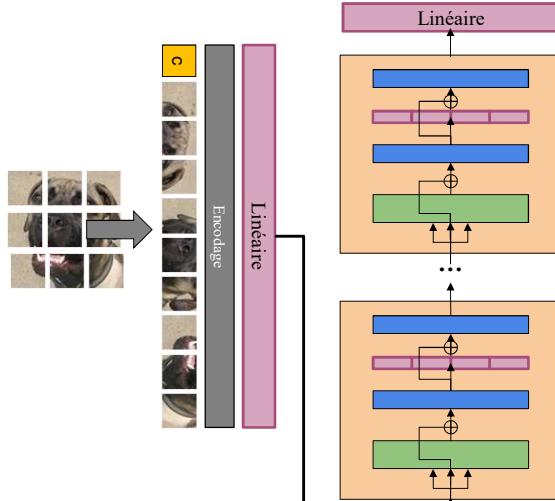
Vision Transformers (ViT)



- 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

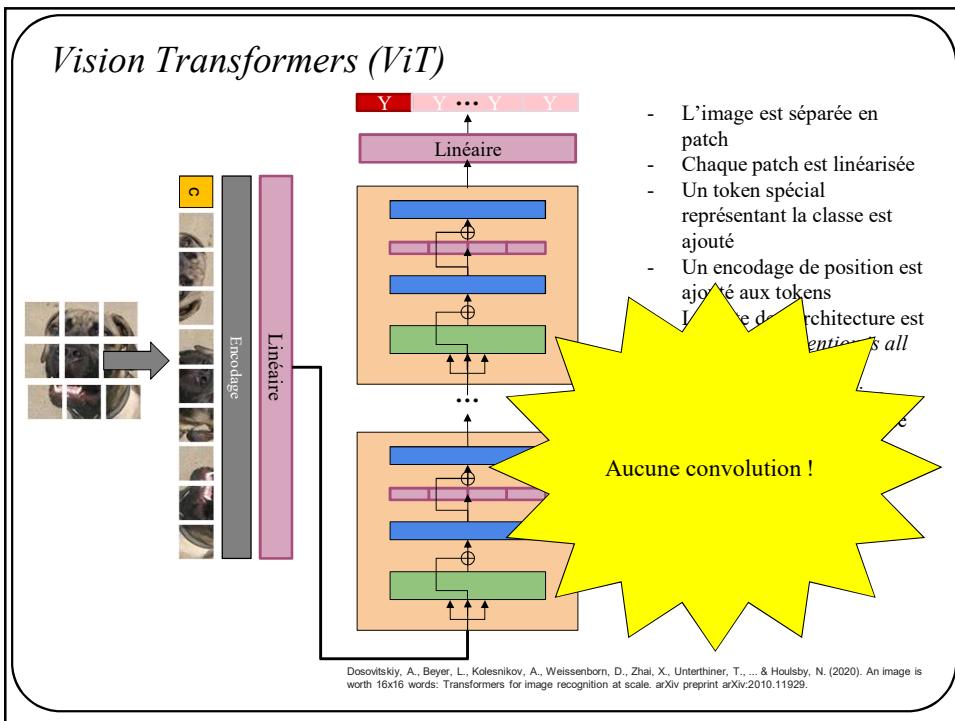
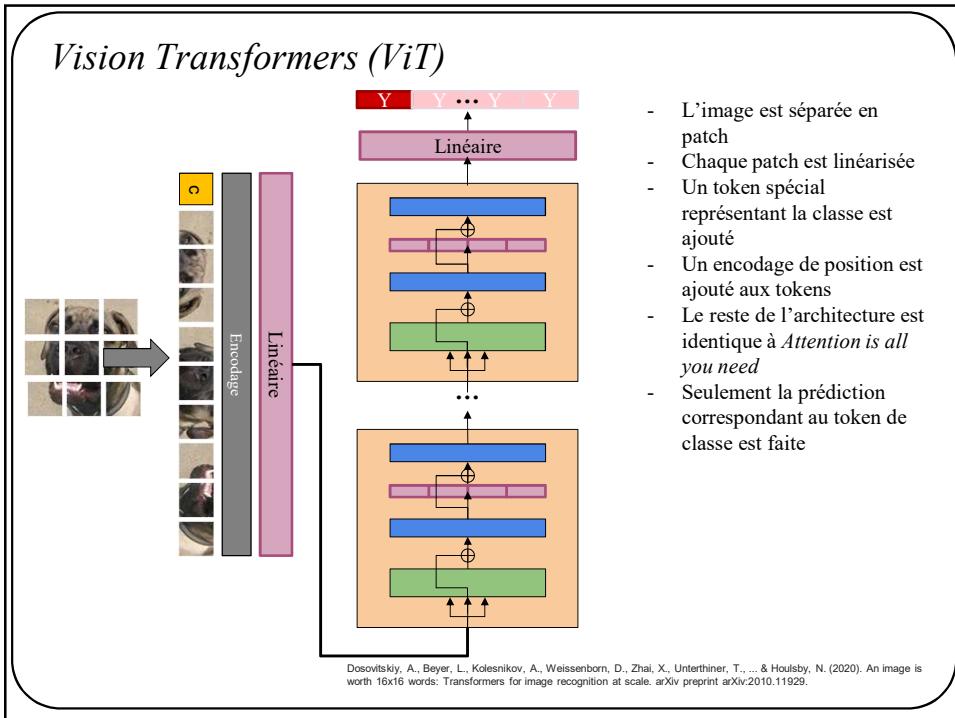
Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

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.



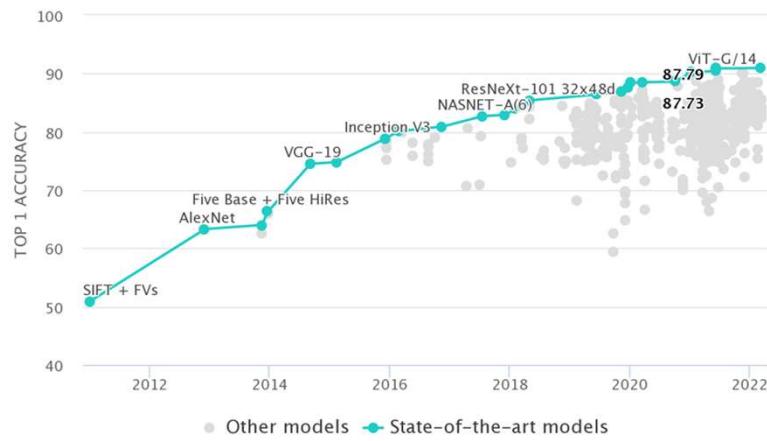
Vision Transformers (ViT)

	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 \pm 0.04	87.76 \pm 0.03	85.30 \pm 0.02	87.54 \pm 0.02	88.4 / 88.5*
ImageNet ReL	90.72 \pm 0.05	90.54 \pm 0.03	88.62 \pm 0.05	90.54	90.55
CIFAR-10	99.50 \pm 0.06	99.42 \pm 0.03	99.15 \pm 0.03	99.37 \pm 0.06	—
CIFAR-100	94.55 \pm 0.04	93.90 \pm 0.05	93.25 \pm 0.05	93.51 \pm 0.08	—
Oxford-IIIT Pets	97.56 \pm 0.03	97.32 \pm 0.11	94.67 \pm 0.15	96.62 \pm 0.23	—
Oxford Flowers-102	99.68 \pm 0.02	99.74 \pm 0.00	99.61 \pm 0.02	99.63 \pm 0.03	—
VTAB (19 tasks)	77.63 \pm 0.23	76.28 \pm 0.46	72.72 \pm 0.21	76.29 \pm 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.

Vision Transformers (ViT)



Les vision transformers dominent la classification depuis leur arrivée

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

Vision Transformers (ViT)

Rank	Model	Top 1 Accuracy	Top 3 Accuracy	Number of params	Extra Training Data	Paper	Code	Result	Year	Tips ↗	
1	Model soups (ViT-G14)	90.94%		1843M	✓	Model soups: averaging weights of multiple ViT models improves accuracy without increasing inference time.	View	Run	2022	Transformer PyTorch	
2	CoAtNet-7	90.88%		2440M	✓	CoAtNet: Mixing Convolution and Attention for All Data Sizes	View	Run	2021	Transformer PyTorch	
3	VIT-G/14	90.45%		1843M	✓	Scaling Vision Transformers	View	Run	2021	Transformer PyTorch	
4	CoAtNet-6	90.45%		1470M	✓	CoAtNet: Mixing Convolution and Attention for All Data Sizes	View	Run	2021	Transformer PyTorch	
5	V-MoE-15B (Every2)	90.35%		14700M	✓	Scaling Vision with Sparse Mixture of Experts	View	Run	2021	Transformer	
6	Meta Pseudo Labels (EfficientNet-L2)	90.2%	98.8%	480M	✓	Meta Pseudo Labels	View	Run	2021	EfficientNet PyTorch	
7	SwinV2-G	90.17%			✓	Swin Transformer V2: Scaling Up Capacity and Resolution	View	Run	2021	Transformer	
8	Florence-CoSwin-H	90.05%	99.02%		✓	Florence: A New Foundation Model for Computer Vision	View	Run	2021	Transformer	
9	Meta Pseudo Labels (EfficientNet-B6-Wide)	90%	98.7%	390M	✓	Meta Pseudo Labels	View	Run	2021	EfficientNet PyTorch	
10	NFNet-F4+	89.2%		527M	✓	High-Performance Large-Scale Image Recognition Without Normalization	View	Run	2021	PyTorch	

Les vision transformers dominent la classification depuis leur arrivée

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

Sommaire

A diagram illustrating a sequence-to-sequence model architecture. An input vector \vec{x} is processed by a hidden state \vec{h} , which then produces an output vector \vec{y} . A feedback loop from \vec{y} back to \vec{h} is shown.

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

↓ train more

"Mont thithey" fomescerlundi
Keushey, Thom here
sheuke, ammerenith ol sivh I lalterthend Bleipile shawy fil on aseterlome
coanigennc Phe lism thond hon at. Meidimorion in ther thize."

↓ train more

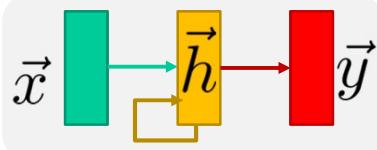
Aftair fall unsuch that the hall for Prince Welonski's that me of
her hearly, and behr to so aravage living were to it beloge, pavu say falling misfor
how, and Gogition is so overelical and ofter.

↓ 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 asking his soul came to the packs and drove up his father-in-law women.

Modélisation de langage

Sommaire

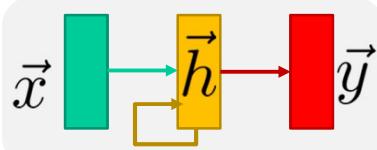


a group of people playing a game with nintendo wii controllers

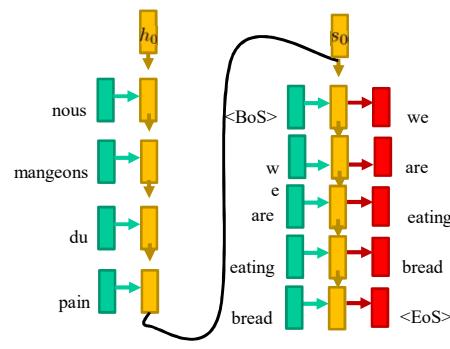
Description d'images

- 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

Sommaire

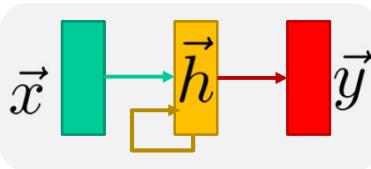
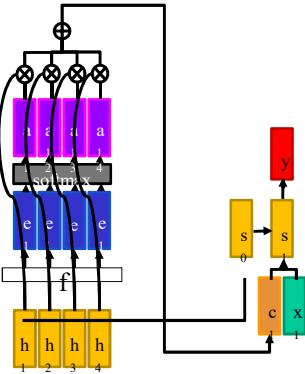


- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique



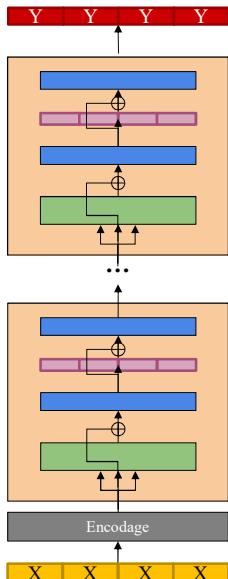
Traduction

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

Sommaire



- Un *Transformer* sont un modèle extrêmement puissant pour les tâches liées au *langage naturel*
- Les *transformers* n'utilisent *que* l'attention (pas un modèle récurrent)
- Les *transformers* sont demandant en ressources
- Ceux-ci ne sont pas limités au langage!