

Réseaux de neurones

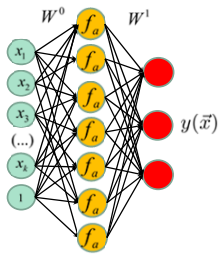
IFT 780

Réseaux récurrents

Par
Pierre-Marc Jodoin, Antoine Théberge

1

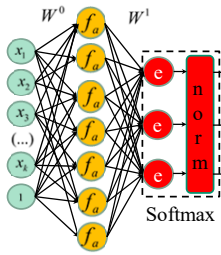
Réseau de neurones de base (régression)



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

2

Réseau de neurones de base (classification)

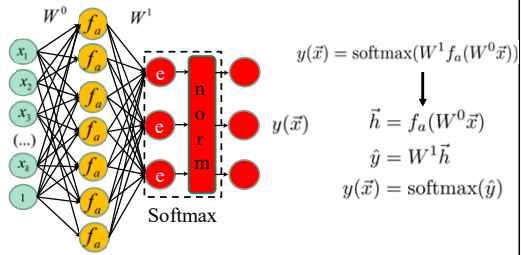


$$y(\vec{x}) = \text{softmax}(W^1 f_a(W^0 \vec{x}))$$
$$\vec{h} = f_a(W^0 \vec{x})$$
$$\hat{y} = W^1 \vec{h}$$
$$y(\vec{x}) = \text{softmax}(\hat{y})$$

3

1

Réseau de neurones de base (classification)

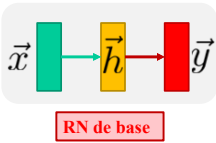


Ne permet que des tâches "1 pour 1"

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

4

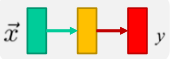
Illustration simplifiée



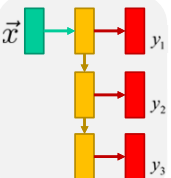
5

Différentes configurations pour différentes applications

1 entrée et 1 sortie

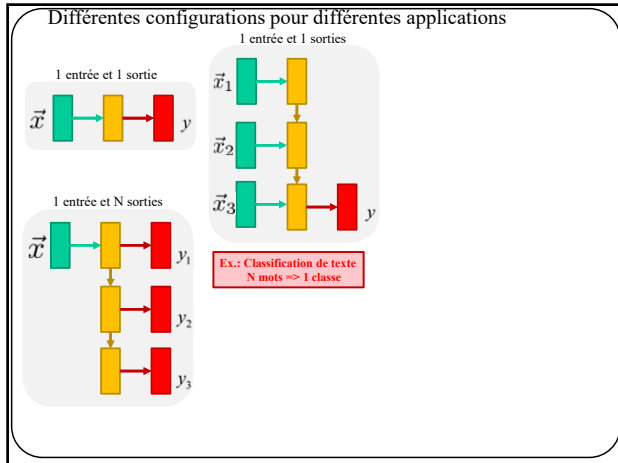


1 entrée et N sorties

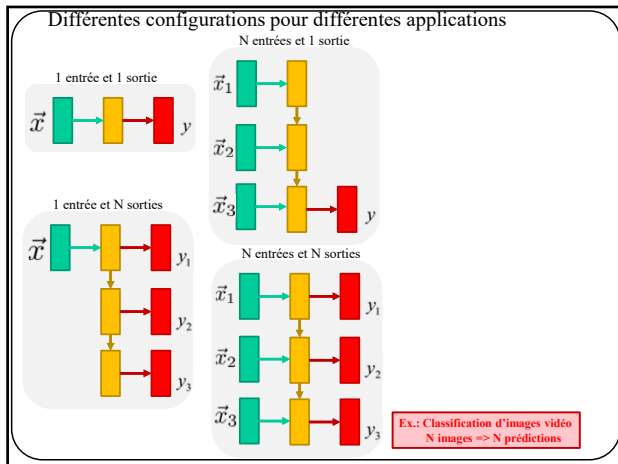


Ex.: description d'une image
1 image => N mots

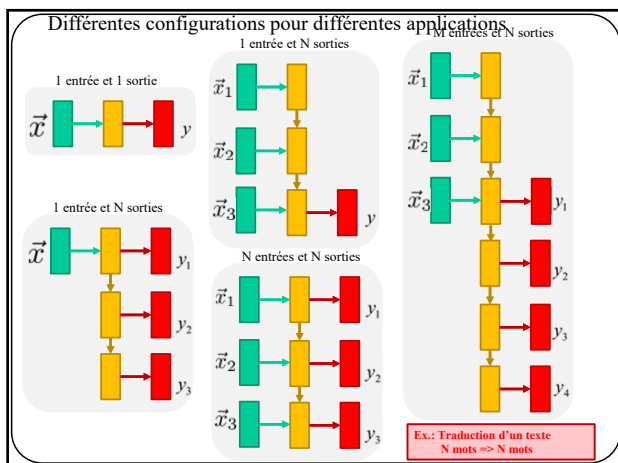
6



7

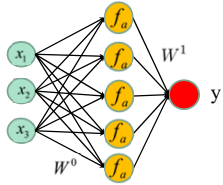


8



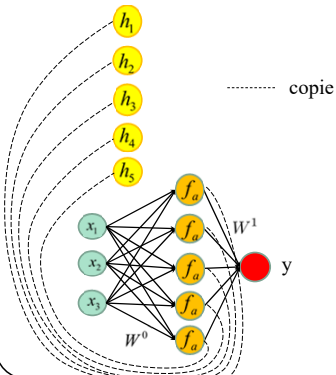
9

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



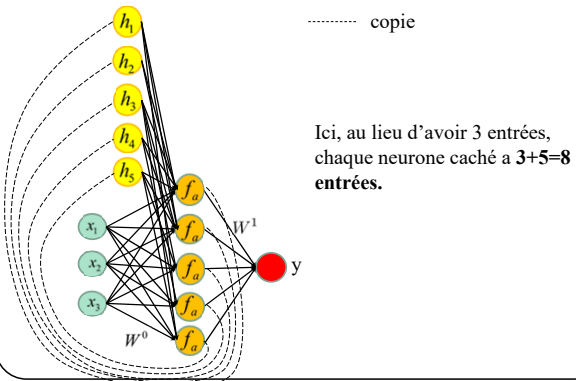
10

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

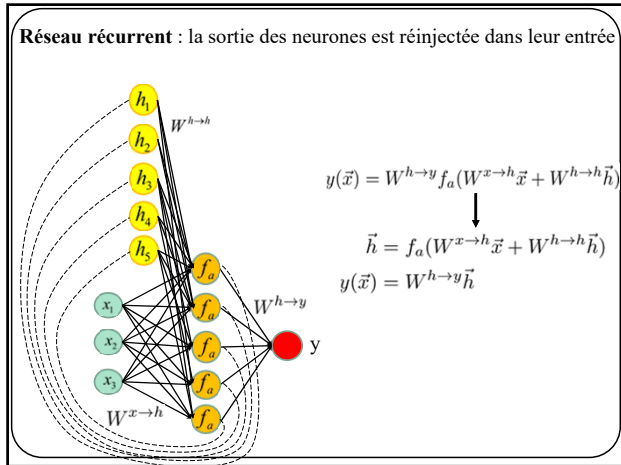


11

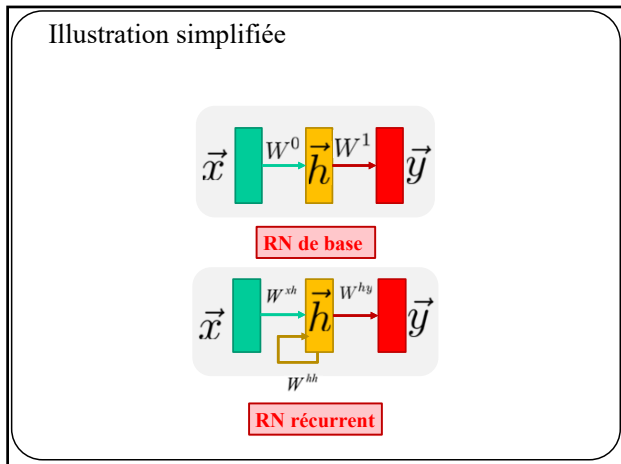
Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée



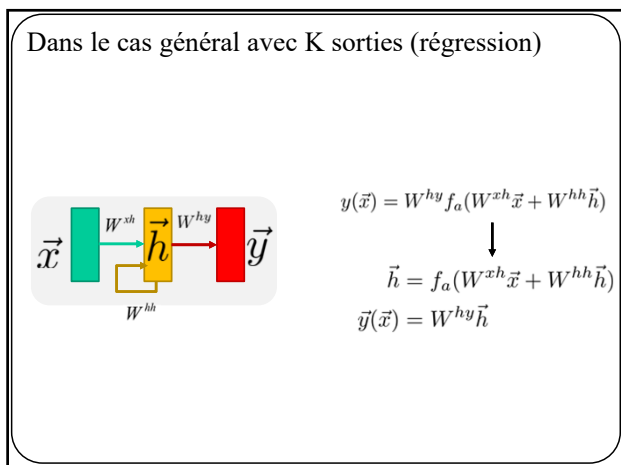
12



13

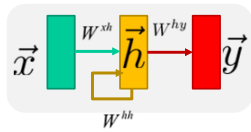


14



15

Dans le cas général avec K sorties (classification)



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

$$\downarrow$$

$$\vec{h} = f_a(W^{xh} \vec{x} + W^{hh} \vec{h})$$

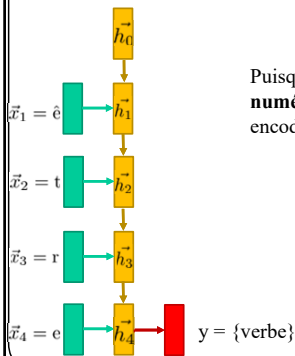
$$\hat{y} = W^{hy} \vec{h}$$

$$\vec{y}(\vec{x}) = \text{softmax}(\hat{y})$$

16

Exemple pour N entrées et 1 sortie:

Analyse grammaticale (classification) : (ê.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 ».

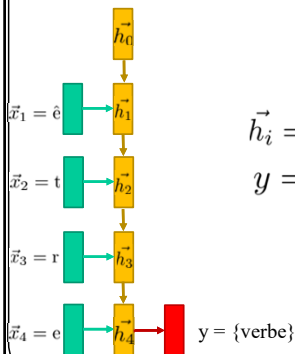
$$\left. \begin{array}{l} 'a' = [1, 0, 0, \dots, 0] \\ 'b' = [0, 1, 0, \dots, 0] \\ 'c' = [0, 0, 1, \dots, 0] \\ \dots \\ 'verbe' = [1, 0, 0, \dots, 0] \\ 'nom' = [0, 1, 0, \dots, 0] \\ 'adjectif' = [0, 0, 1, \dots, 0] \end{array} \right\} \in R^{256}$$

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

17

Exemple pour N entrées et 1 sortie:

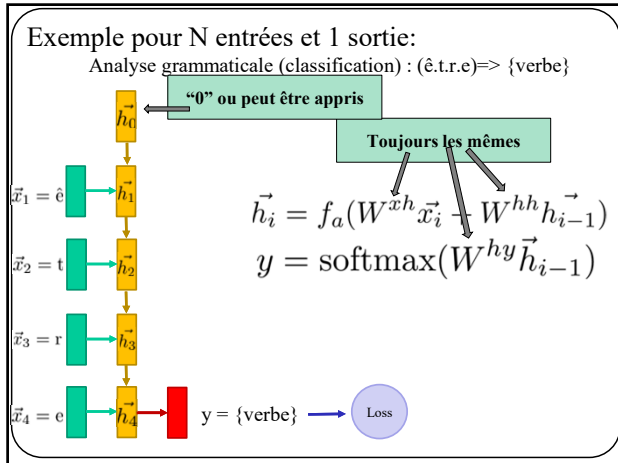
Analyse grammaticale (classification) : (ê.t.r.e) => {verbe}



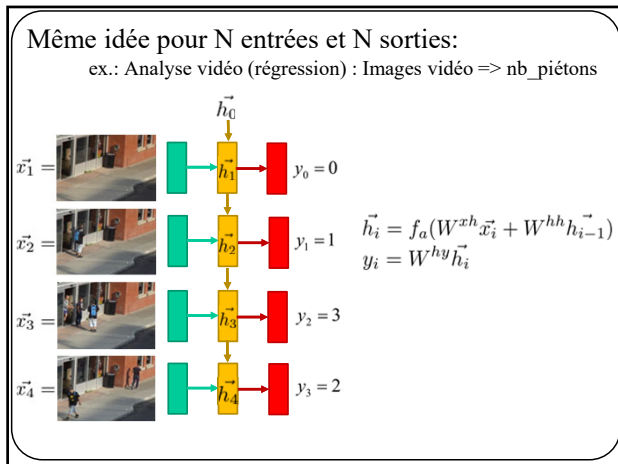
$$\vec{h}_i = f_a(W^{xh} \vec{x}_i + W^{hh} \vec{h}_{i-1})$$

$$y = \text{softmax}(W^{hy} \vec{h}_{i-1})$$

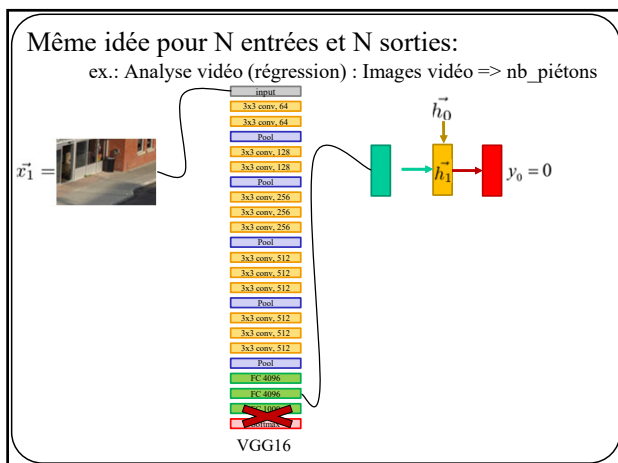
18



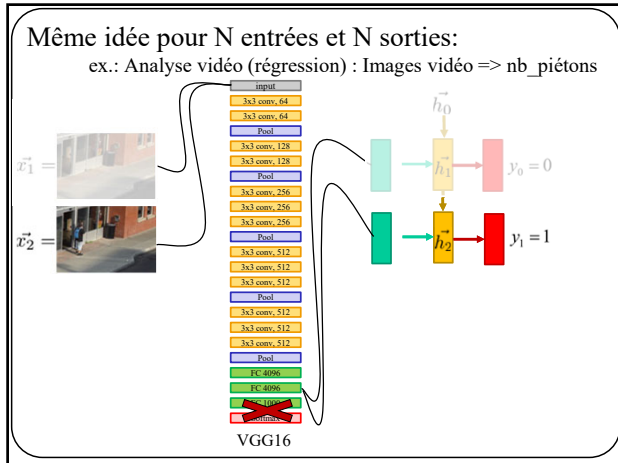
19



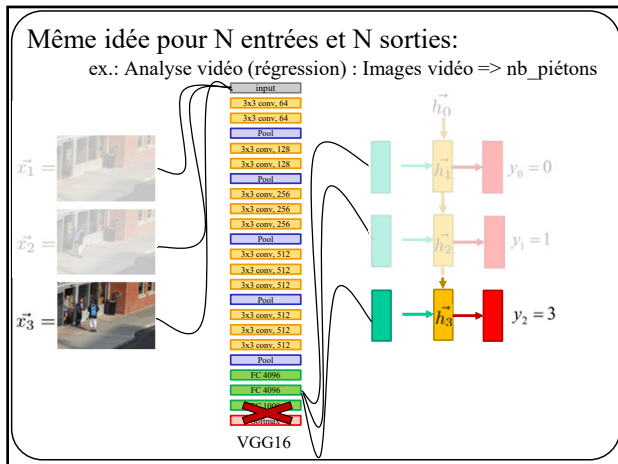
20



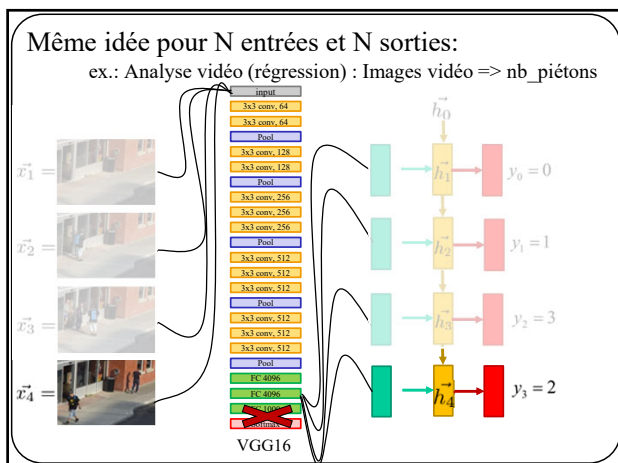
21



22



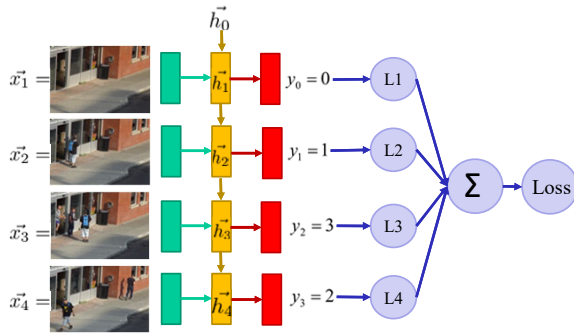
23



24

Même idée pour N entrées et N sorties:

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



25

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

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

Représentation « one hot » jouet:

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

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

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

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

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

26

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

Alphabet : [a,e,m,s]

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

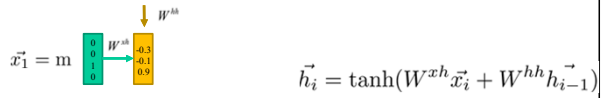
$\vec{x}_1 = m$

27

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

Alphabet :[a,e,m,s]

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

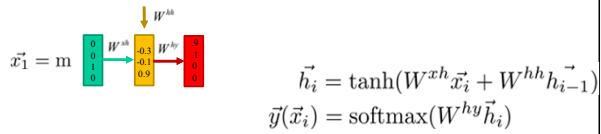


28

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

Alphabet :[a,e,m,s]

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

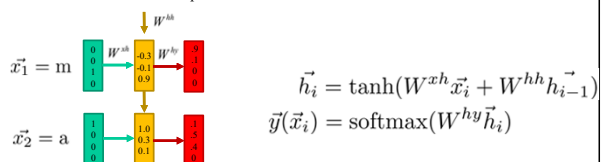


29

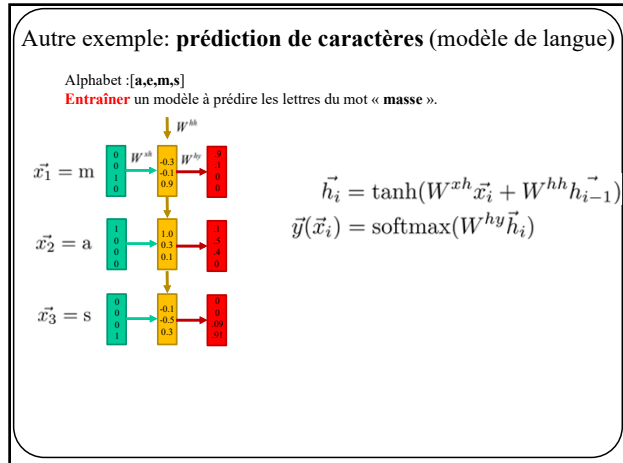
Autre exemple: **prédiction de caractères** (modèle de langue)

Alphabet :[a,e,m,s]

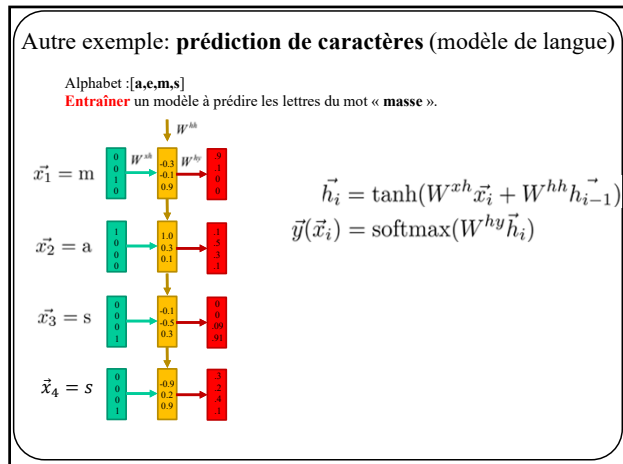
Entraîner un modèle à prédire les lettres du mot « masse ».



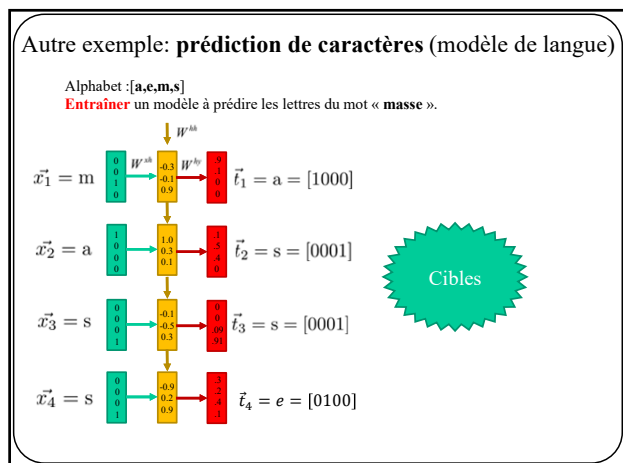
30



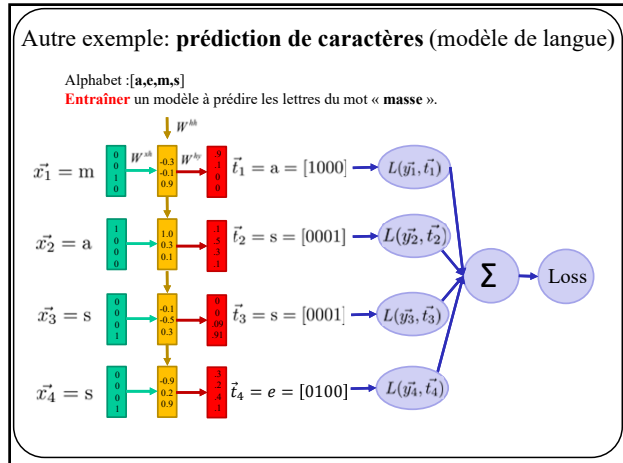
31



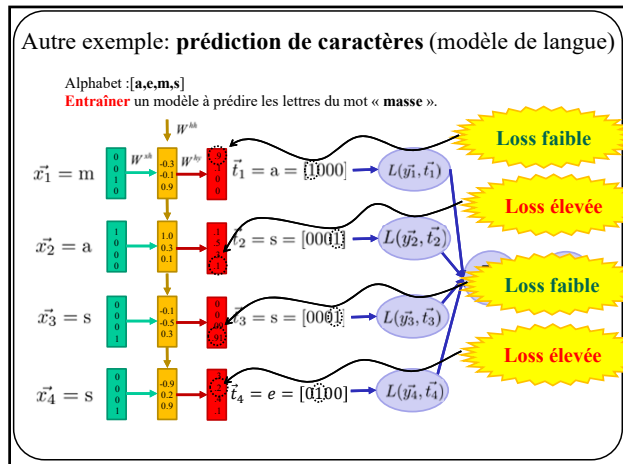
32



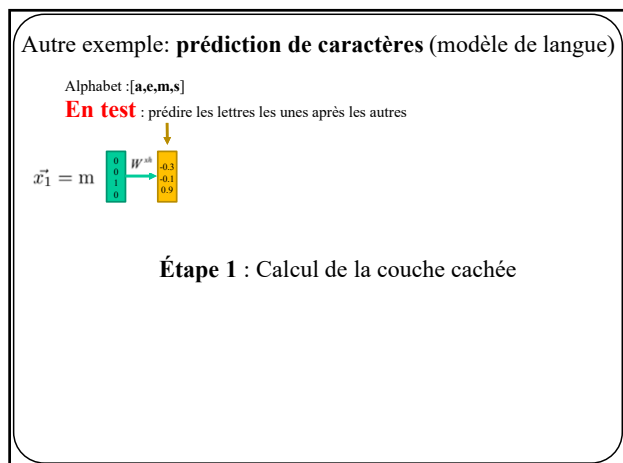
33



34



35

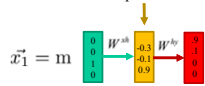


36

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

Alphabet :[a,e,m,s]

En test : prédire les lettres les unes après les autres



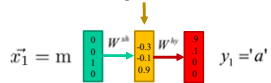
Étape 2 : Calcul de la sortie (softmax)

37

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

Alphabet :[a,e,m,s]

En test : prédire les lettres les unes après les autres



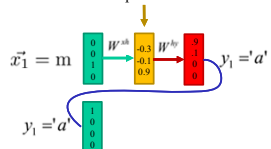
Étape 3 : Sélectionner le caractère le plus probable

38

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

Alphabet :[a,e,m,s]

En test : prédire les lettres les unes après les autres

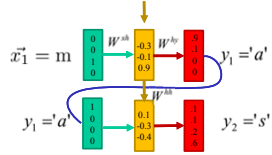


Étape 4 : Injecter le caractère prédit au début du réseau

39

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

Alphabet :[a,e,m,s]

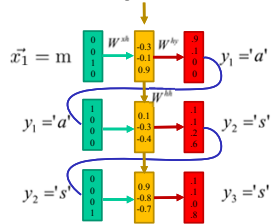
En test : prédire les lettres les unes après les autres

Et on recommence!

40

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

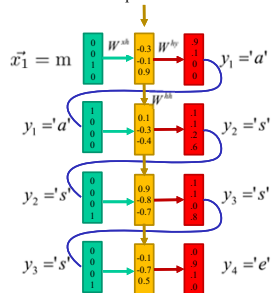
Alphabet :[a,e,m,s]

En test : prédire les lettres les unes après les autres

41

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

Alphabet :[a,e,m,s]

En test : prédire les lettres les unes après les autres

42

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 newer 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 news 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 age, not a man and in fire,
To show the raining 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 goodness and our finest,
We spare with hours, but cut thy council I am great.
Murdered and by thy master's ready there
My power to give thou but so much as hell:
Some service is the noble's hand here,
Would show him to bar wise.

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

Crédit: A. Karpathy, CS231

46

Entraînement sur le code source de Linux en C++

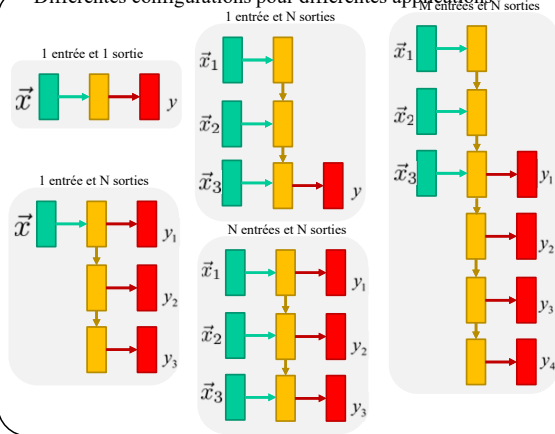
```
static void do_command(struct seq_file *s, void *)
{
    int column = 32 << (cmd[] & 0x00);
    if (state)
        cmd = (int)(int_state * (int)(int_state_flags & Cmd) 7 2 * 1);
    else
        seq = 1;
    for (i = 0; i < 16; i++) {
        if (i & 0x0001) {
            pipe = (int)(int_state * UNKNTREAD_UNCKA) +
                ((count & 0x00000000ffffff) & 0x00000000) << 8;
            if (count == 0)
                subpid_ppc_md_kexec_handle, 0x20000000;
            pipe_set_bytes(i, 0);
        }
        /* From our user pages pointer to place names if all dash */
        subsystem_info = &of_changepage_PAGE_SIZE;
        sub_control(effect, id, substate);
        /* Now we want to deliberately put it to derive */
        control_check_polarity(content, val, 0);
        for (i = 0; i < UNKNTREAD_UNCKA; i++)
            seq_puts(s, "policy");
    }
}
```

```
static void read_stat_HPC_read_mostly_offset(struct seq_operations, \
    p0M[1]);
static void
seq_operations(long seq)
{
    struct seq_operations
    P0T_PARAM_AND(1, seq) = seq_state_state();
    seq_get_stat(read_mostly_offset, current_state_stat);
    read_mostly_offset = seq_full; long
}
```

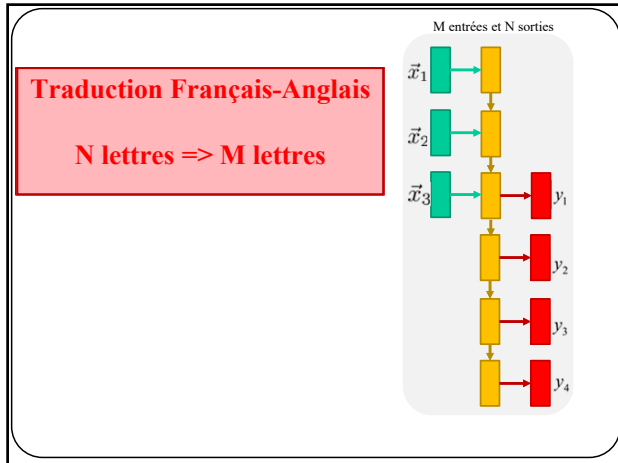
Crédit: A. Karpathy, CS231

47

Différentes configurations pour différentes applications



48



49

Autre exemple: traduction

Traduire 'assez' -> 'enough'

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

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

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

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

50

Autre exemple: traduction

Traduire 'assez' -> 'enough'

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

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

$\vec{x}_1 = \text{<BoS>}$

51

Autre exemple: traduction

Traduire 'assez' -> 'enough'
 Alphabet fr : [**<BoS>**,a,e,s,z,<EoS>]
 Alphabet en : [**<BoS>**,e,g,h,n,o,u,<EoS>]

$\vec{x}_1 = \langle \text{BoS} \rangle$
 $\vec{x}_2 = a$

52

Autre exemple: traduction

Traduire 'assez' -> 'enough'
 Alphabet fr : [**<BoS>**,a,e,s,z,<EoS>]
 Alphabet en : [**<BoS>**,e,g,h,n,o,u,<EoS>]

$\vec{x}_1 = \langle \text{BoS} \rangle$
 $\vec{x}_2 = a$
 $\vec{x}_i = \dots$
 $\vec{x}_7 = \langle \text{EoS} \rangle$

53

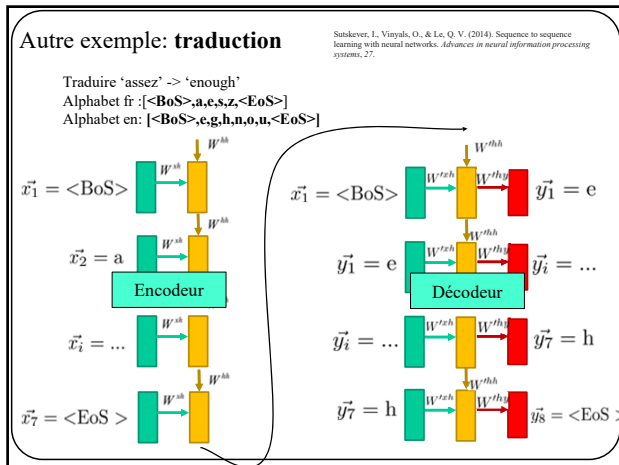
Autre exemple: traduction

Traduire 'assez' -> 'enough'
 Alphabet fr : [**<BoS>**,a,e,s,z,<EoS>]
 Alphabet en : [**<BoS>**,e,g,h,n,o,u,<EoS>]

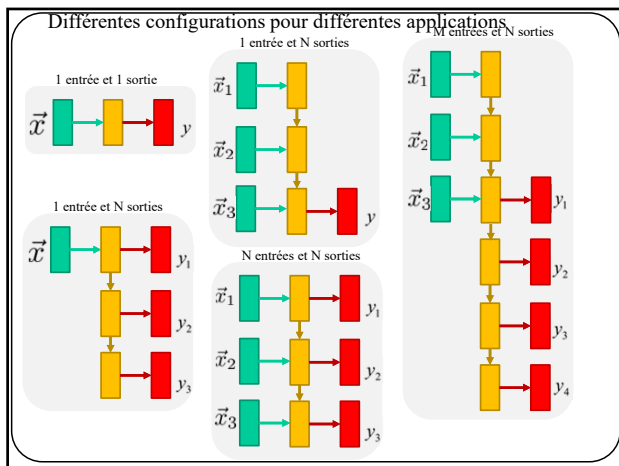
$\vec{x}_1 = \langle \text{BoS} \rangle$
 $\vec{x}_2 = a$
 $\vec{x}_i = \dots$
 $\vec{x}_7 = \langle \text{EoS} \rangle$

$\vec{y}_1 = e$
 $\vec{y}_i = \dots$
 $\vec{y}_7 = h$
 $\vec{y}_8 = \langle \text{EoS} \rangle$

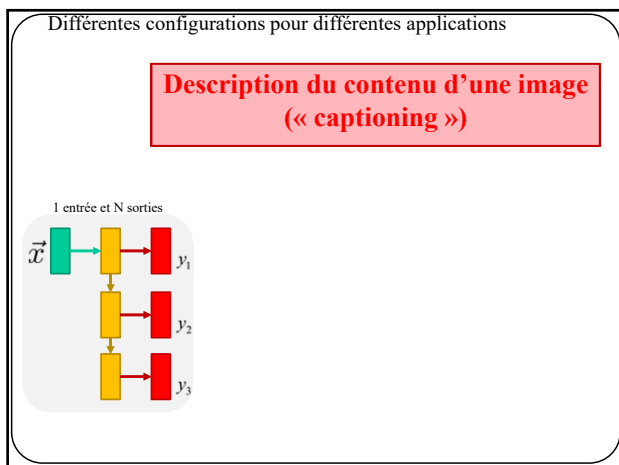
54



55



56



57

The diagram illustrates the VGG16 architecture, which is a deep convolutional neural network. It starts with an input layer, followed by two main stages of convolution and pooling. The first stage consists of two layers of 3x3 convolutions with 64 filters, followed by a max pooling layer. The second stage consists of two layers of 3x3 convolutions with 128 filters, followed by a max pooling layer. This is followed by three more layers of 3x3 convolutions with 256 filters, then three more layers of 3x3 convolutions with 512 filters. The final layers are two fully connected layers of 4096 units each, followed by a Softmax layer for classification. The diagram also includes a small image of a child's face, which is the input to the network.

Diagram illustrating the VGG16 architecture (Réseau VGG pré-entraîné sur ImageNet):

- Input
- 3x3 conv, 64
- 3x3 conv, 64
- Pool
- 3x3 conv, 128
- 3x3 conv, 128
- Pool
- 3x3 conv, 256
- 3x3 conv, 256
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- fc 4096
- fc 4096
- fc 1000
- Softmax

58

Captioning




Diagram illustrating the VGG16 architecture for image captioning. The input image is processed through a series of layers:

- Input
- 3x3 conv, 64
- 3x3 conv, 64
- Pool
- 3x3 conv, 128
- 3x3 conv, 128
- Pool
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- fc, 4096
- fc, 4096
- ~~fc, 1000~~

VGG16

59

Captioning

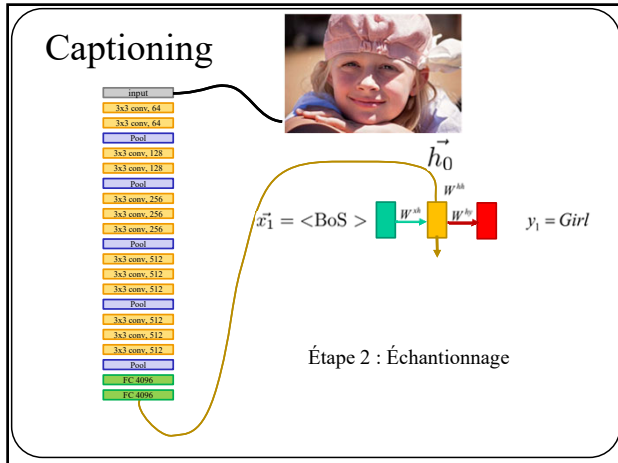
Diagram illustrating the architecture for Captioning, showing a sequence of layers:

- input
- 3x3 conv, 64
- 3x3 conv, 64
- Pool
- 3x3 conv, 128
- 3x3 conv, 128
- Pool
- 3x3 conv, 256
- 3x3 conv, 256
- 3x3 conv, 256
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- FC 4096
- FC 4096

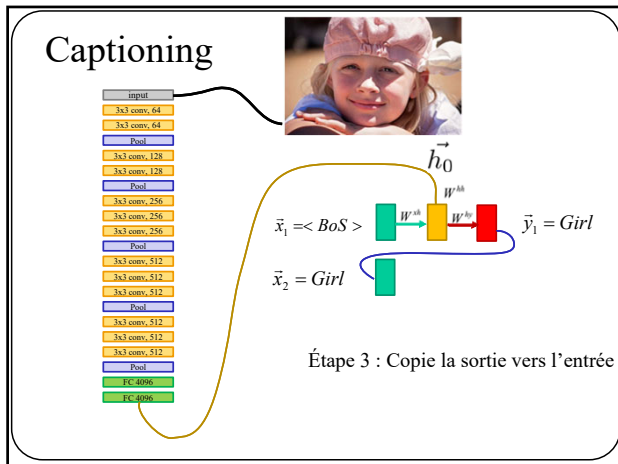
The diagram also shows a yellow arrow indicating the forward pass (Init + Propagation avant) and a red arrow indicating the backward pass (Init + Propagation arrière).

Étape 1 : Init + Propagation avant

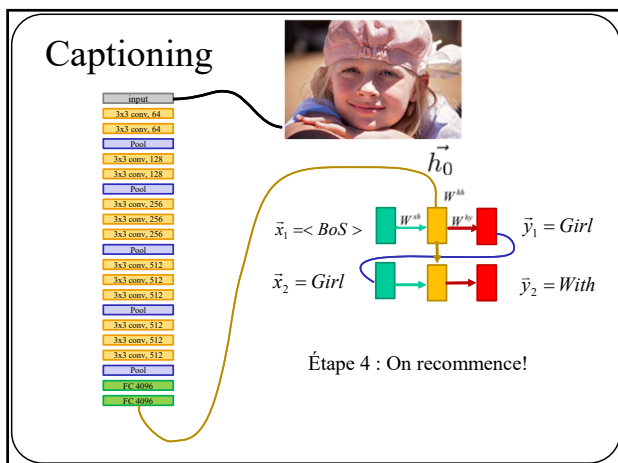
60



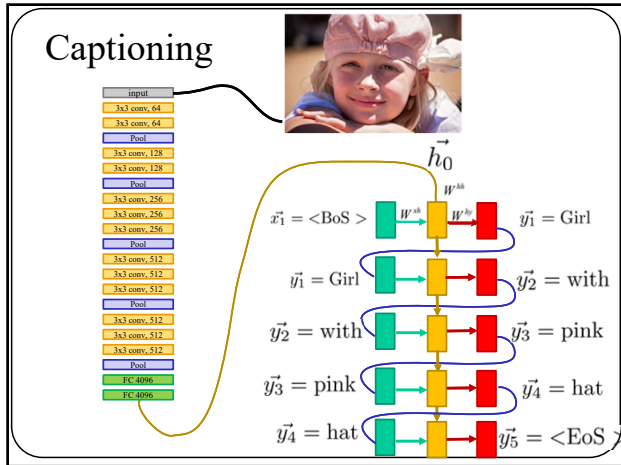
61



62



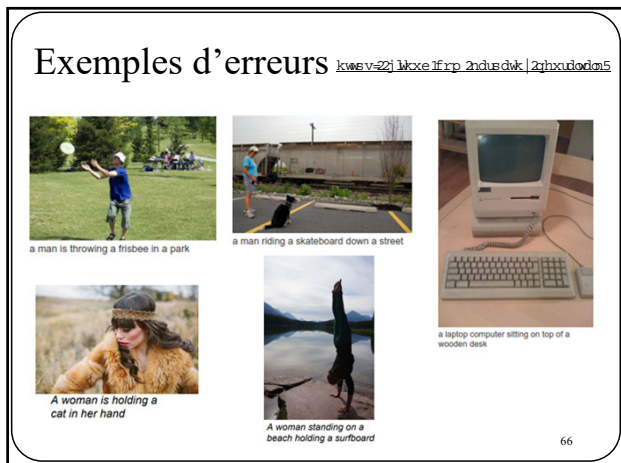
63



64



65



66

NeuralTalk and Walk

<https://arxiv.org/abs/1507.04808>



67

67

Analyse de texte

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

Pour des **caractères**...

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

68

68

Analyse de texte

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

Pour des **mots**...

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

69

69

Prédiction sur des lettres vs. mots

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

...

...

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

...

Prédiction sur des lettres

Prédiction sur des mots

70

Prédiction sur des lettres vs. mots

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

...

...

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

...

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

71

Limites des Jetons « one-hot »

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

- 1- Peu efficace en mémoire lorsque non compressés
ex.: 10,000 bits pour encoder le mot « je » dans une langue à 10,000 mots!
- 2- Pas de distance sémantique entre les Jetons:

Ex.

$\text{distance}[\text{one-hot}(\text{'bon'}), \text{one-hot}(\text{'bien'})] = \text{distance}[\text{one-hot}(\text{'bon'}), \text{one-hot}(\text{'trottoir'})]$

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

$\text{distance}[\text{code}(\text{'bon'}), \text{code}(\text{'bien'})] \ll \text{distance}[\text{code}(\text{'bon'}), \text{code}(\text{'trottoir'})]$
 $\text{distance}[\text{code}(\text{'Jean'}), \text{code}(\text{'Chantal'})] \ll \text{distance}[\text{code}(\text{'bon'}), \text{code}(\text{'trottoir'})]$
 $\text{distance}[\text{code}(\text{'Inde'}), \text{code}(\text{'Liban'})] \ll \text{distance}[\text{code}(\text{'bon'}), \text{code}(\text{'trottoir'})]$

74

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

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

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

1 ligne = code
pour 1 mot

75

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

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

Ex: sélectionner le jeton de « brown »

$$\begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 & 3 & 4 & 5 \\ -1 & -3 & -2 & 2 \\ 11 & 6 & 4 & -3 \\ -4 & 8 & -4 & 4 \\ 24 & -6 & 42 & 17 \\ 91 & 13 & 14 & -5 \\ 0 & 36 & 4 & 56 \\ -1 & 0 & 1 & 35 \end{pmatrix} = \begin{pmatrix} 11 & 6 & 4 & -3 \end{pmatrix}$$

76

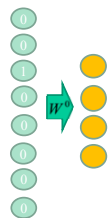
Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire jeton = matrice d'encodage

Première couche d'un réseau de neurones

=
matrice d'encodage

\vec{x} : brown



W^0

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



77

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

Première couche d'un réseau de neurones
=
matrice d'encodage



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

78

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

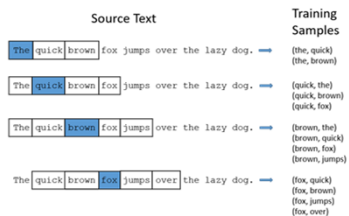


On pourra donc utiliser un réseau de neurones
pour calculer le contenu du dictionnaire

79

Word2Vec s'appuie sur 2 idées fondamentales

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



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

80

Word2Vec [Mikolov et al. '13]

$\vec{x} : \text{brown}$ $\vec{t} : \text{fox}$

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

81

Word2Vec [Mikolov et al. '13]

$\vec{x} : \text{brown}$ $\vec{t} : \text{fox}$

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

82

Word2Vec [Mikolov et al. '13]

$\vec{x} : \text{brown}$ $\vec{t} : \text{fox}$

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

83

Word2Vec [Mikolov et al. '13]

\vec{x} : brown W^0 \vec{t} : fox

Lorsqu'entraîné, utiliser W^0 comme dictionnaire

84

Word2Vec [Mikolov et al. '13]

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

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

85

Distance sémantique entre deux mots = distance entre leur jeton

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

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

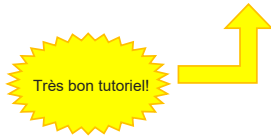
Ahmia, Oussama & Béchet, Nicolas & Marteau, Pierre-Francois. *Two Multilingual Corpora Extracted from the Tenders Electronic Daily for Machine Learning and Machine Translation Applications*. in LREC 2018

86

86

Word2Vec

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



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

87

Comment entraîner un RNN?

88

Histoire de gradients

RN de classification avec entropie croisée



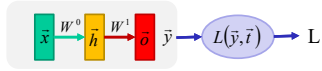
$$\hat{y}(\vec{x}) = S_M(W^1 \tanh(W^0 \vec{x}))$$

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

89

Histoire de gradients

Simple RN de classification avec entropie croisée



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

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

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

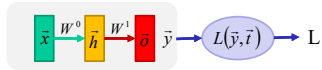
$$L = L_{CE}(\bar{y}, \bar{t})$$

Propagation
avant

90

Histoire de gradients

Simple RN de classification avec entropie croisée



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

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

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

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

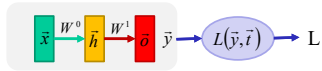
Pour entraîner le réseau
il faut calculer

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

91

Histoire de gradients

Simple RN de classification avec entropie croisée



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

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

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

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

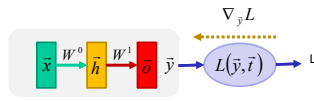
Dérivée en chaîne

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

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

92

Histoire de gradients



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

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

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

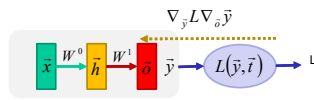
$$L = L_{CE}(\tilde{y}, \tilde{t})$$

Rétro-propagation

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

93

Histoire de gradients



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

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

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

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

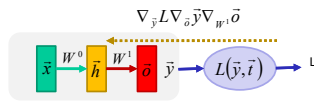
Rétro-propagation

$$\nabla_{\tilde{o}} \tilde{y} = \text{I}\tilde{y}^T - \tilde{y}^T \tilde{y}$$

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

94

Histoire de gradients



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

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

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

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

Rétro-propagation

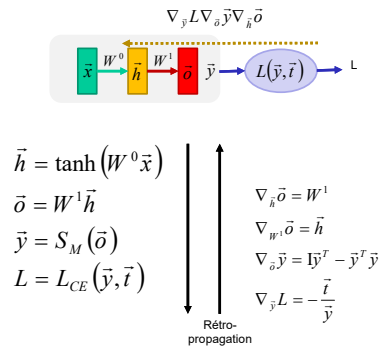
$$\nabla_{w^1} \tilde{o} = \tilde{h}$$

$$\nabla_{\tilde{o}} \tilde{y} = \text{I}\tilde{y}^T - \tilde{y}^T \tilde{y}$$

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

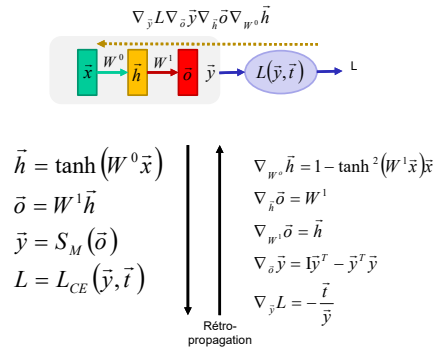
95

Histoire de gradients



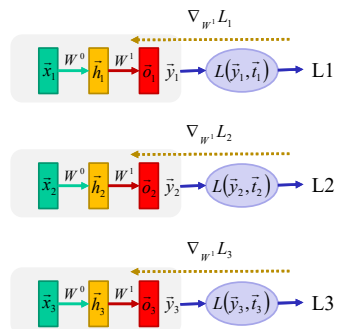
96

Histoire de gradients



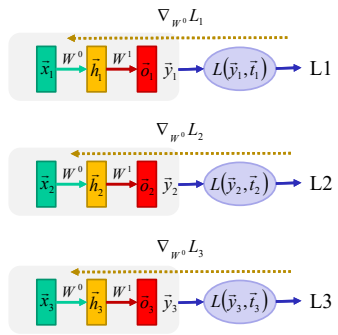
97

Ex.: 3 données, 3 rétro-propagations



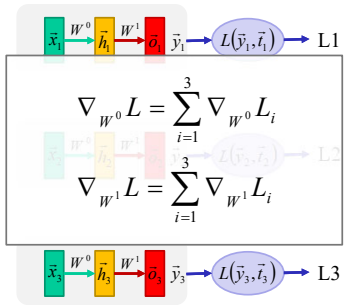
98

Ex.: 3 données, 3 rétro-propagations

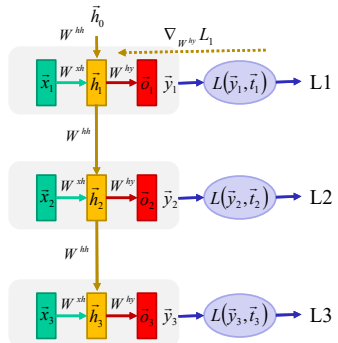


99

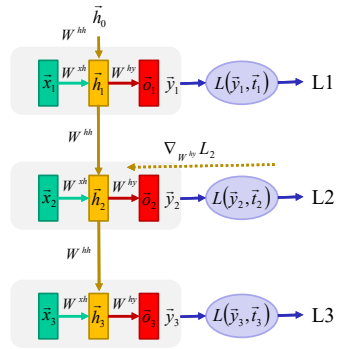
3 rétro-propagations



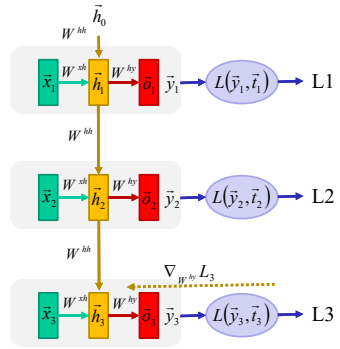
100

Réseau récurrent: gradient pour W^{hy} 

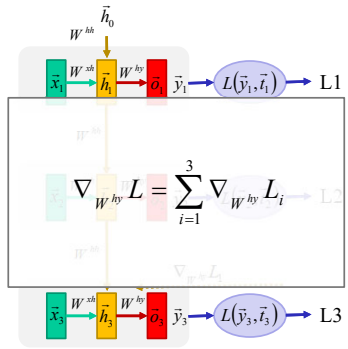
101

Réseau récurrent: gradient pour W^{hy} 

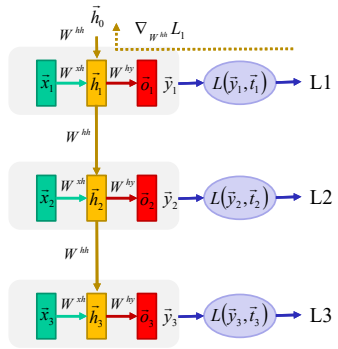
102

Réseau récurrent: gradient pour W^{hy} 

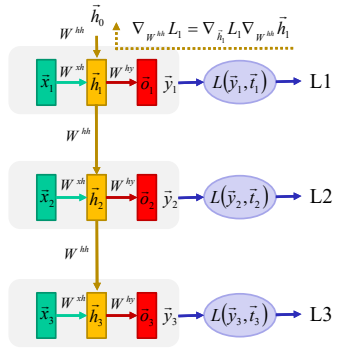
103

Réseau récurrent: gradient pour W^{hy} 

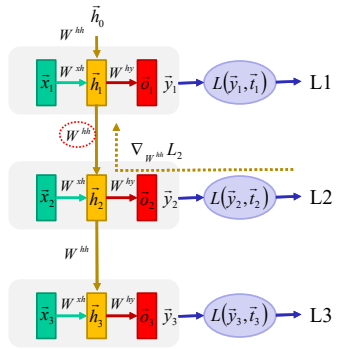
104

Réseau récurrent: gradient pour W^{hh} 

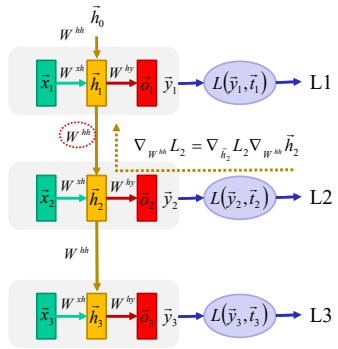
105

Réseau récurrent: gradient pour W^{hh} 

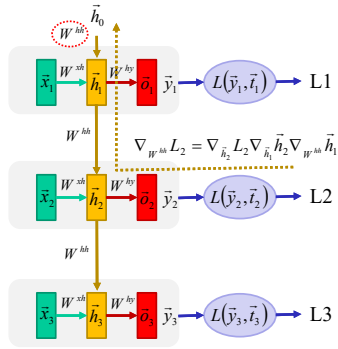
106

Réseau récurrent: gradient pour W^{hh} 

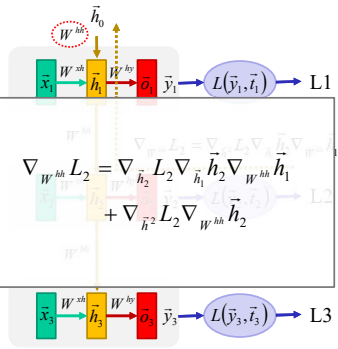
107

Réseau récurrent: gradient pour W^{hh} 

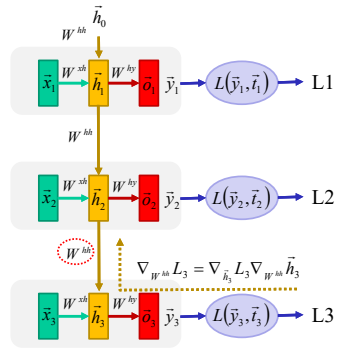
108

Réseau récurrent: gradient pour W^{hh} 

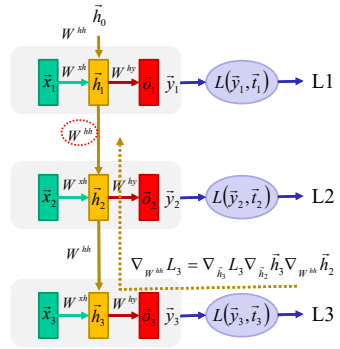
109

Réseau récurrent: gradient pour W^{hh} 

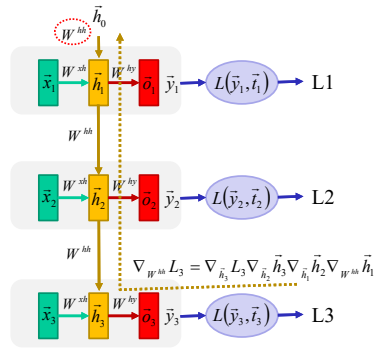
110

Réseau récurrent: gradient pour W^{hh} 

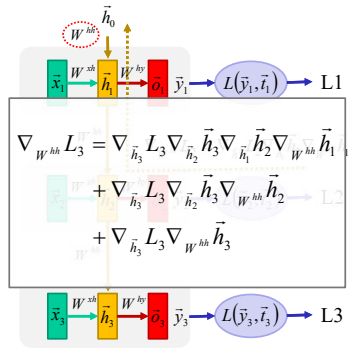
111

Réseau récurrent: gradient pour W^{hh} 

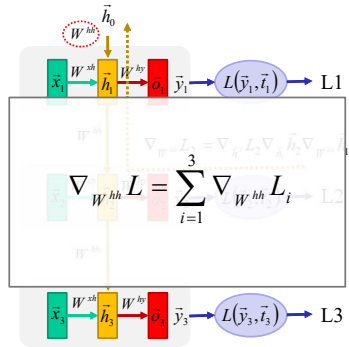
112

Réseau récurrent: gradient pour W^{hh} 

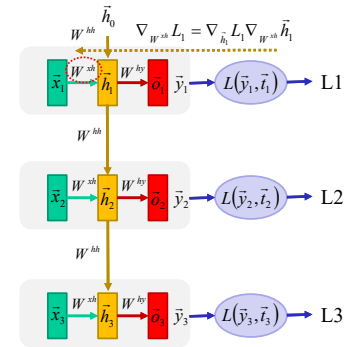
113

Réseau récurrent: gradient pour W^{hh} 

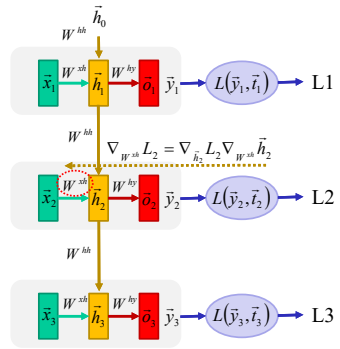
114

Réseau récurrent: gradient pour W^{hh} 

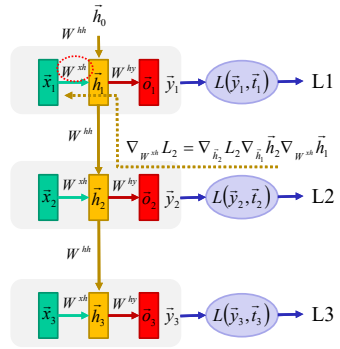
116

Réseau récurrent: gradient pour W^{xh} 

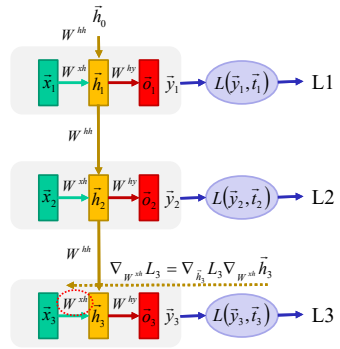
118

Réseau récurrent: gradient pour W^{xh} 

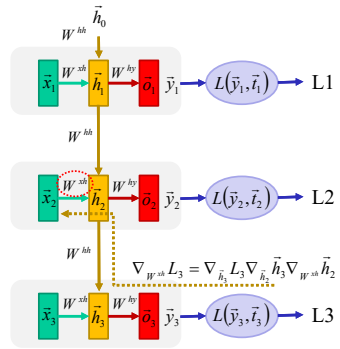
119

Réseau récurrent: gradient pour W^{xh} 

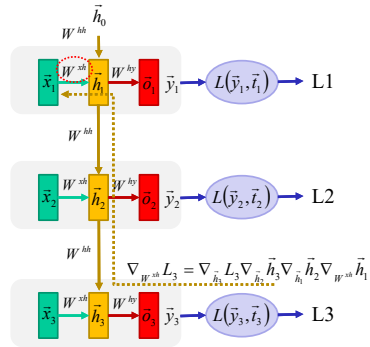
120

Réseau récurrent: gradient pour W^{xh} 

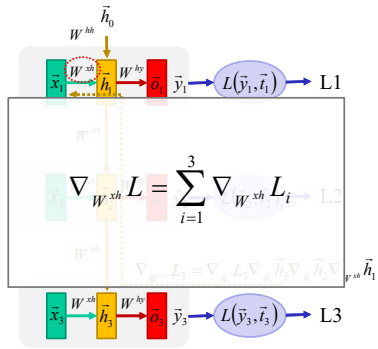
121

Réseau récurrent: gradient pour W^{xh} 

122

Réseau récurrent: gradient pour W^{xh} 

123

Réseau récurrent: gradient pour W^{xh} 

124

Réseau récurrent: calcul du gradient

Moins difficile qu'il n'y paraît.

```

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

```

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

125

Les réseaux récurrents ont un
inconvenient majeur:













difficile à établir des
relations à longue distance

126

126

Exemples: analyse grammaticale

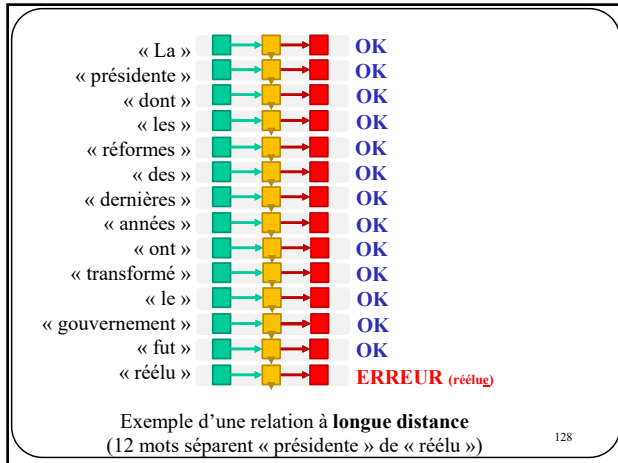
Entraîner un réseau à détecter des erreurs grammaticales

« La »				OK
« présidente »				OK
« fut »				OK
« réélu »				ERREUR (réélu)

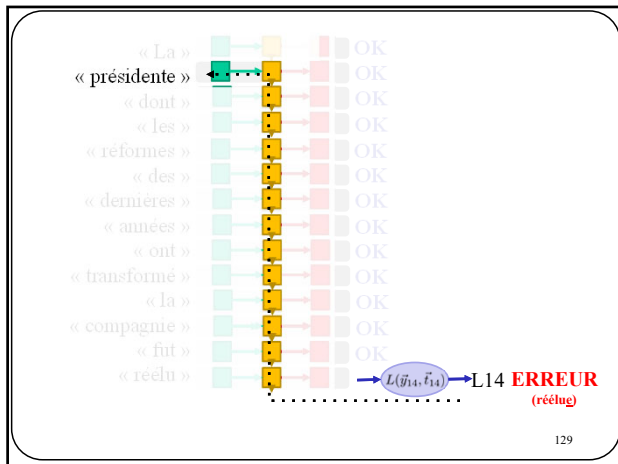
Exemple d'une relation à **courte distance**
(1 mot sépare « présidente » de « réélu »)

127

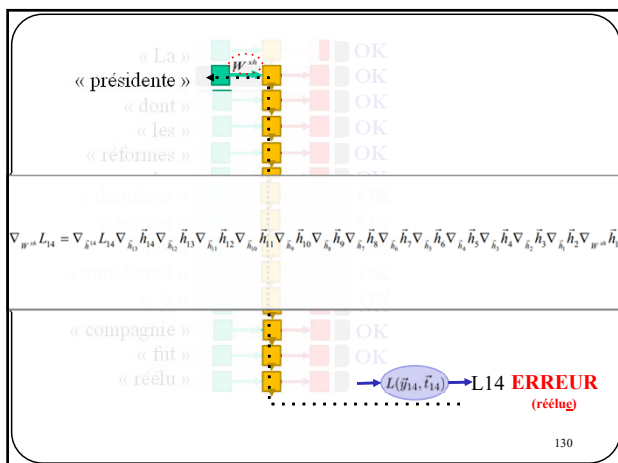
127



128



129



130

« La » OK
« présidente » OK
« dont » OK
« l' » OK
« réforme » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« réélu » OK

« réélu » → L14 **ERREUR** (réélu)

Disparition du gradient
si $-1 < \nabla_{h_{j-1}} h_j < 1$

131

« La » OK
« présidente » OK
« dont » OK
« l' » OK
« réforme » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« » OK
« réélu » OK

« réélu » → L14 **ERREUR** (réélu)

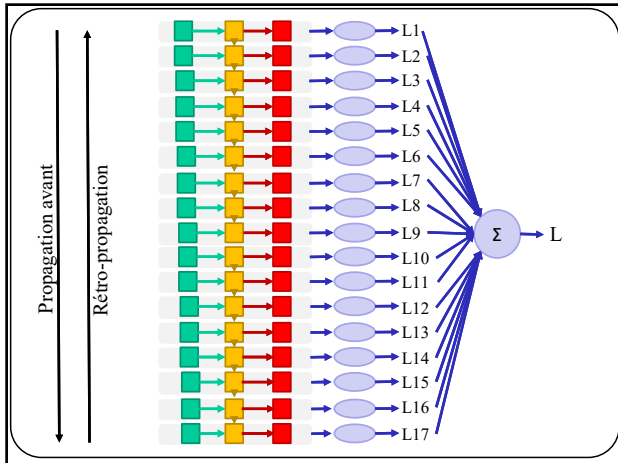
Explosion du gradient
si $-1 > \nabla_{h_{j-1}} h_j > 1$

132

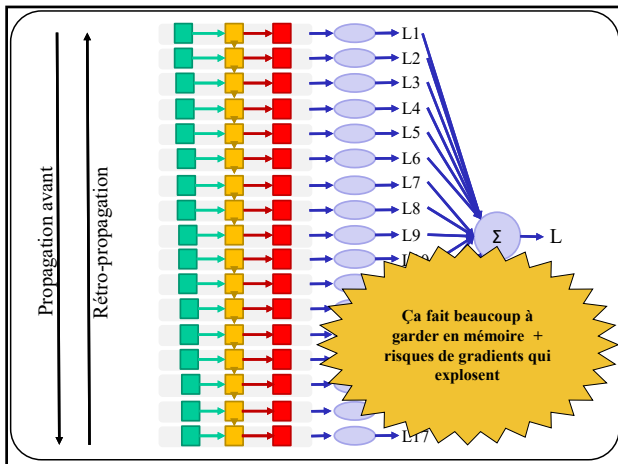
Problème connexe

Gestion de la mémoire

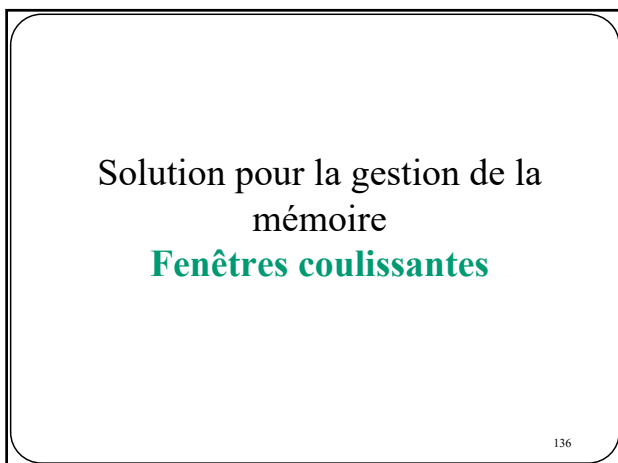
133



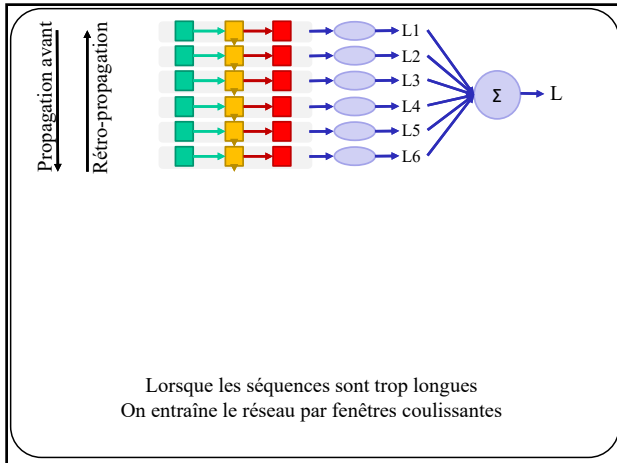
134



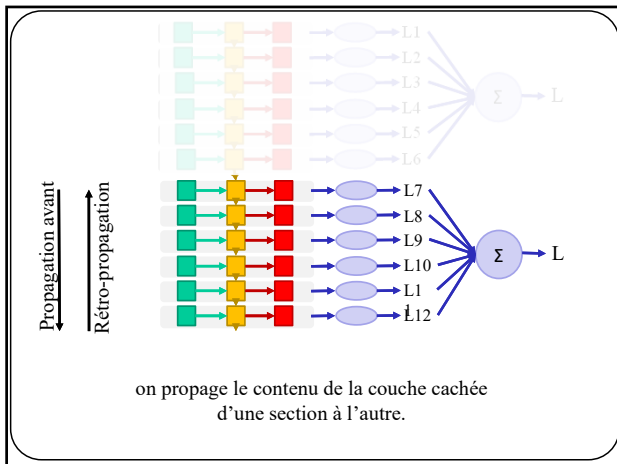
135



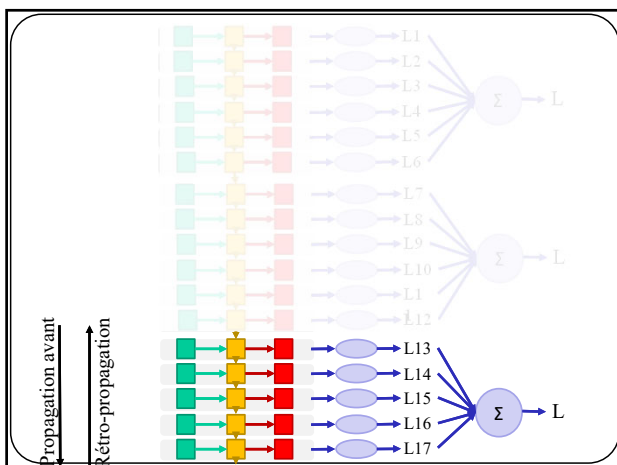
136



137



138



139

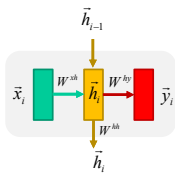
Solution à la disparition du gradient:

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

140

140

Illustration + formulation d'un RNN

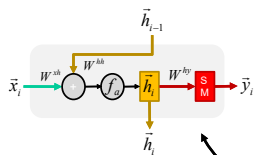


$$\begin{aligned}\vec{h}_i &= f_a(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \hat{y}_i &= W^{hy}\vec{h}_i \\ \bar{y}_i &= SMAX(\hat{y}_i)\end{aligned}$$

141

141

Autre illustration du même RNN



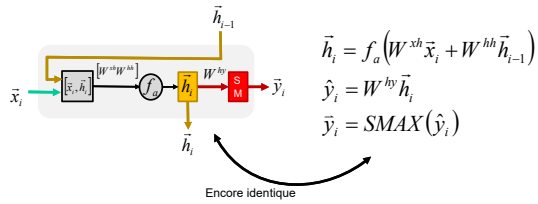
$$\begin{aligned}\vec{h}_i &= f_a(W^{xh}\vec{x}_i + W^{hh}\vec{h}_{i-1}) \\ \hat{y}_i &= W^{hy}\vec{h}_i \\ \bar{y}_i &= SMAX(\hat{y}_i)\end{aligned}$$

Identique

142

142

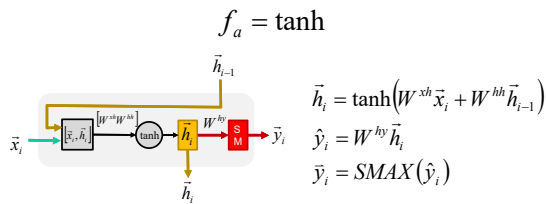
Autre illustration du même RNN



143

GRU (Gated Recurrent Unit)

Modif 1

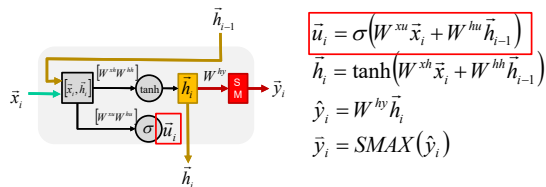


144

GRU (Gated Recurrent Unit)

Modif 2

Update gate

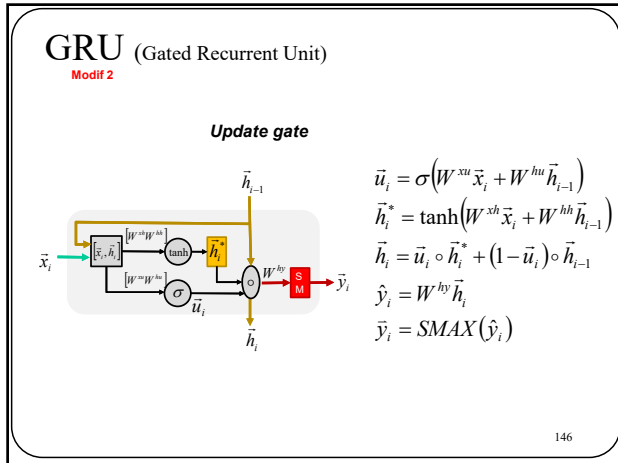
 $\sigma = \text{sigmoid}$ 

145

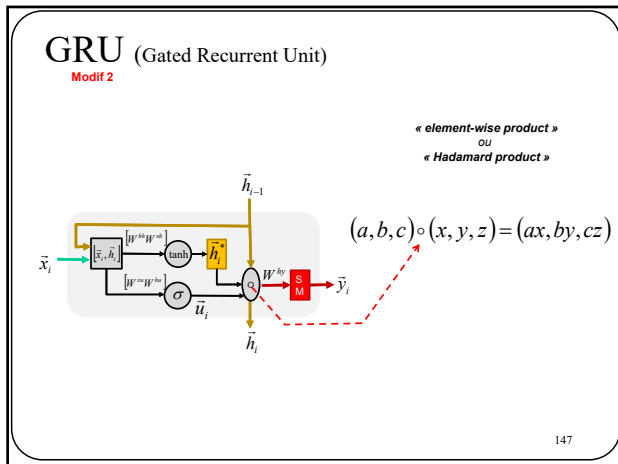
143

144

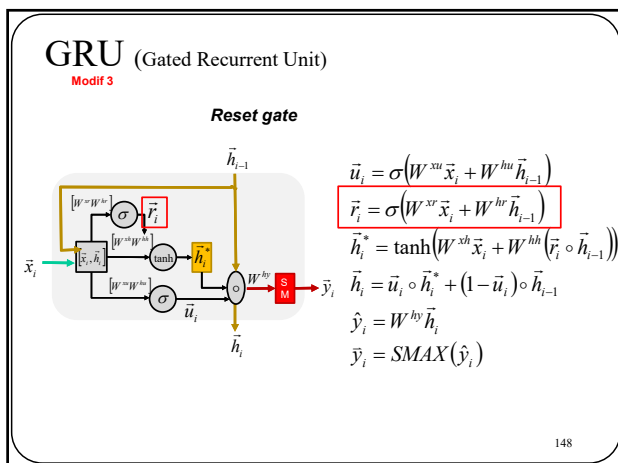
145



146



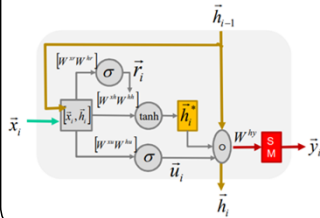
147



148

Comprendre les *gates*

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

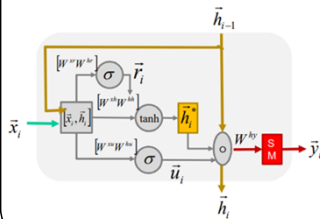


149

149

Comprendre les *gates*

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

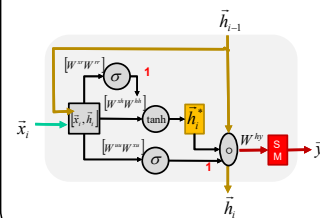


150

150

Comprendre les *gates*

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

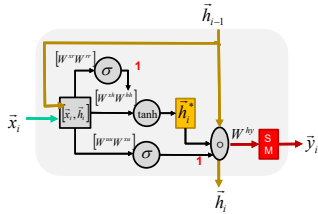


151

151

Comprendre les *gates*

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

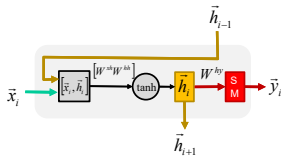


152

152

Comprendre les *gates*

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

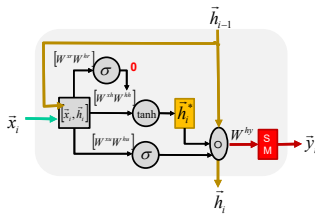


RNN de base!

153

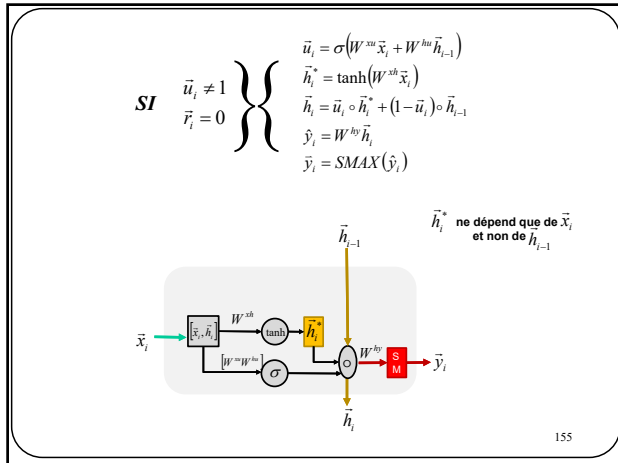
153

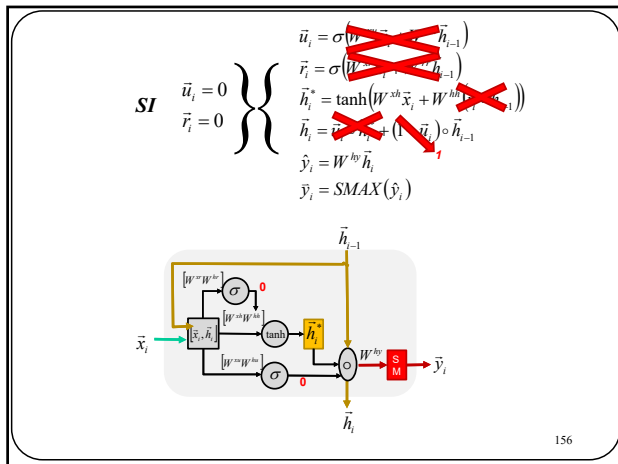
$$SI \quad \left. \begin{array}{l} \tilde{u}_i \neq 1 \\ \tilde{r}_i = 0 \end{array} \right\} \left\{ \begin{array}{l} \tilde{u}_i = \sigma(W^{xu}\tilde{x}_i + W^{hu}\tilde{h}_{i-1}) \\ \tilde{r}_i = \sigma(W^{xr}\tilde{x}_i + W^{hr}\tilde{h}_{i-1}) \\ \tilde{h}_i^* = \tanh(W^{xh}\tilde{x}_i + W^{hh}\tilde{h}_{i-1})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

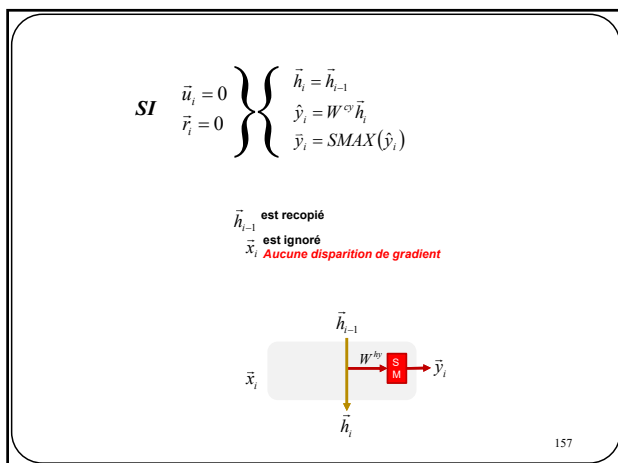


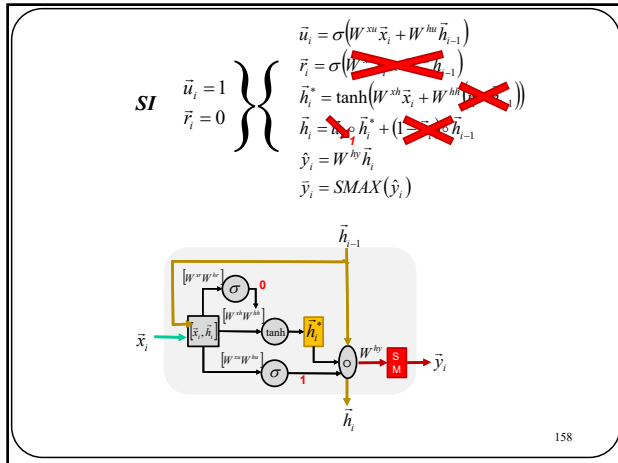
154

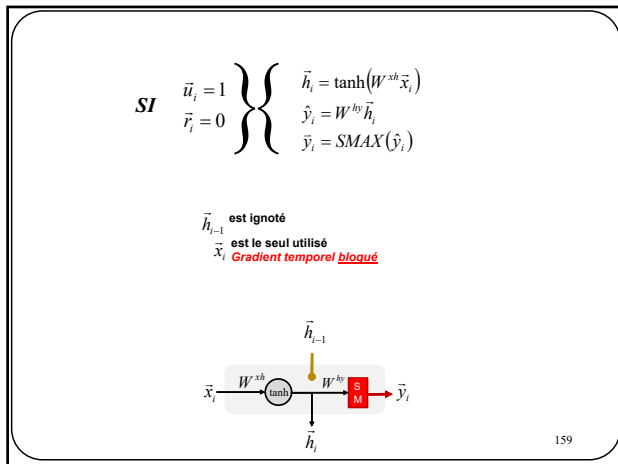
154

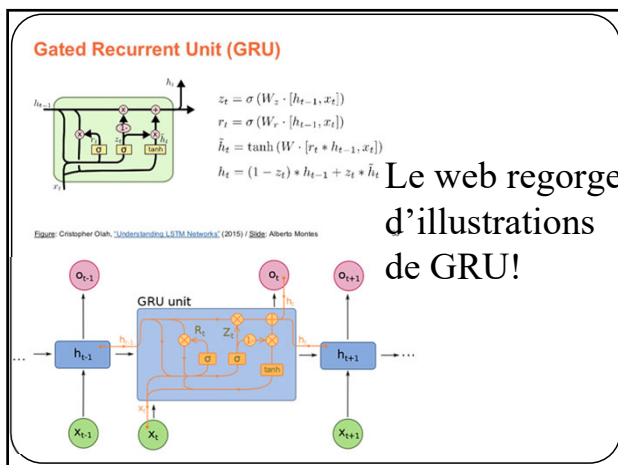




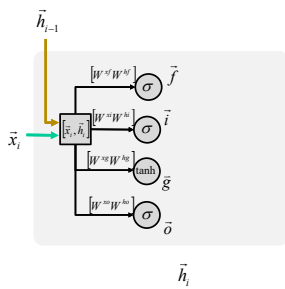








LSTM (Long Short Term Memory)

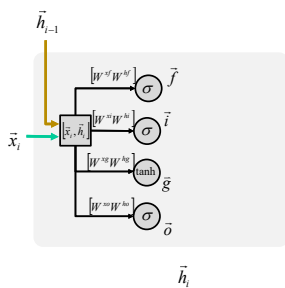


\tilde{f} : forget gate
 \tilde{i} : input gate
 \tilde{g} : gate gate
 \tilde{o} : output gate

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

161

LSTM (Long Short Term Memory)



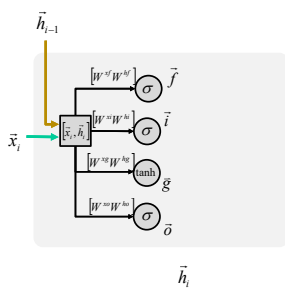
\tilde{f} : forget gate
 \tilde{i} : input gate
 \tilde{g} : gate gate

Modèle récurrent
le plus utilisé.
À bien
comprendre !

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

162

LSTM (Long Short Term Memory)

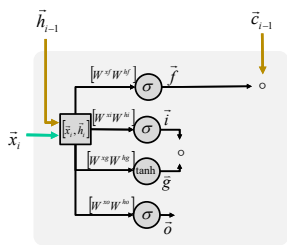


$$\begin{aligned}
 \tilde{f} &= \sigma(W^{xf} \tilde{x}_i + W^{hf} \tilde{h}_{i-1}) \\
 \tilde{i} &= \sigma(W^{xi} \tilde{x}_i + W^{hi} \tilde{h}_{i-1}) \\
 \tilde{g} &= \tanh(W^{xg} \tilde{x}_i + W^{hg} \tilde{h}_{i-1}) \\
 \tilde{o} &= \sigma(W^{xo} \tilde{x}_i + W^{ho} \tilde{h}_{i-1})
 \end{aligned}$$

163

163

LSTM (Long Short Term Memory)

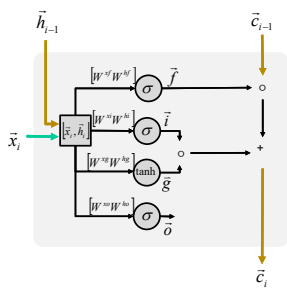


$$\begin{aligned}\tilde{f} &\circ \tilde{c}_i \\ \tilde{i} &\circ \tilde{g}\end{aligned}$$

164

164

LSTM (Long Short Term Memory)

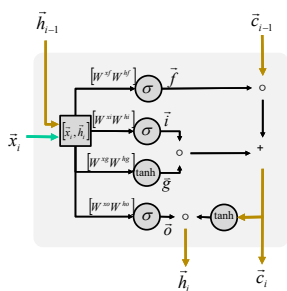


$$\tilde{c}_i = \tilde{f} \circ \tilde{c}_i + \tilde{i} \circ \tilde{g}$$

165

165

LSTM (Long Short Term Memory)

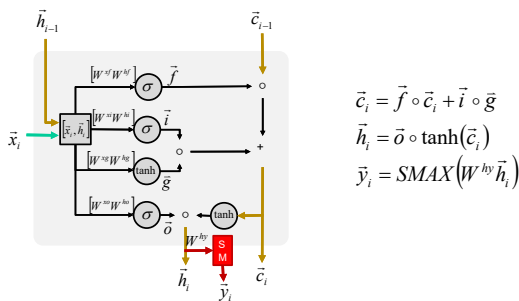


$$\begin{aligned}\tilde{c}_i &= \tilde{f} \circ \tilde{c}_i + \tilde{i} \circ \tilde{g} \\ \tilde{h}_i &= \tilde{o} \circ \tanh(\tilde{c}_i)\end{aligned}$$

166

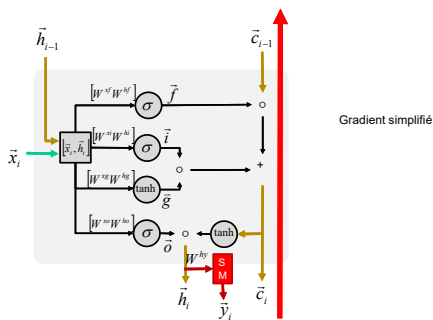
166

LSTM (Long Short Term Memory)



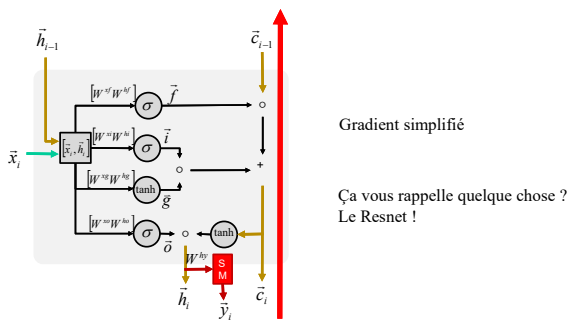
167

LSTM (Long Short Term Memory)



168

LSTM (Long Short Term Memory)

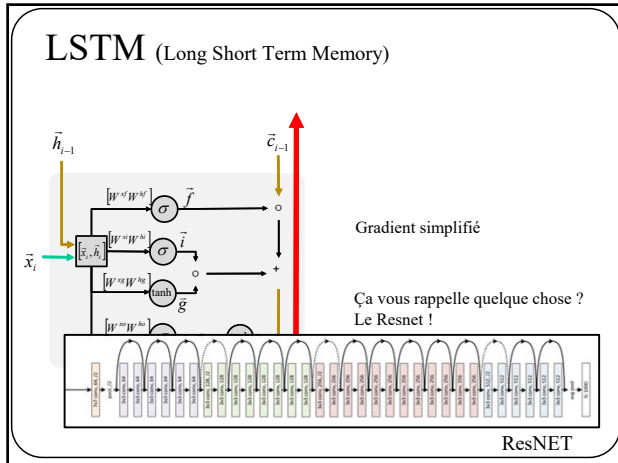


169

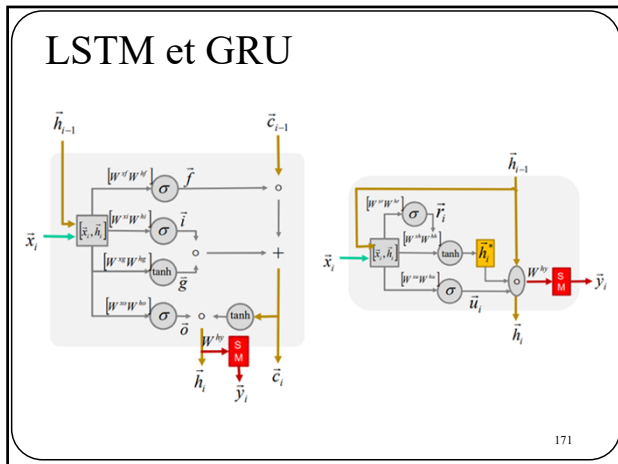
167

168

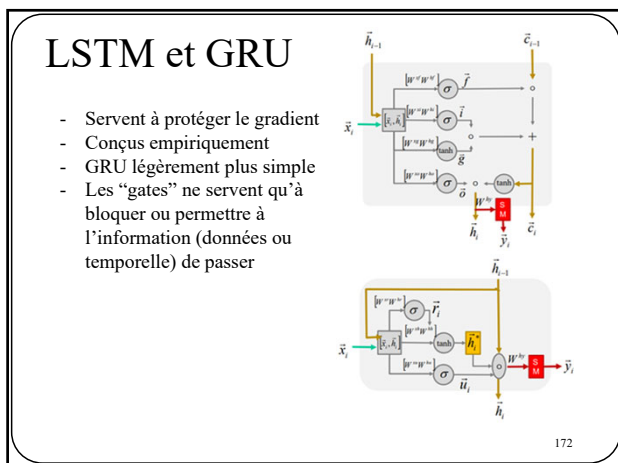
169



170

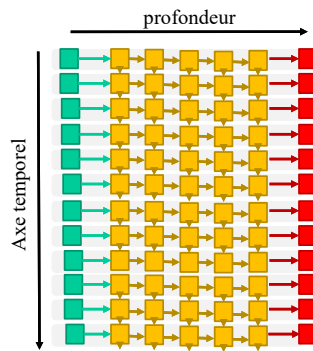


171



172

RNN multi-couches



173

Modèles d'attention

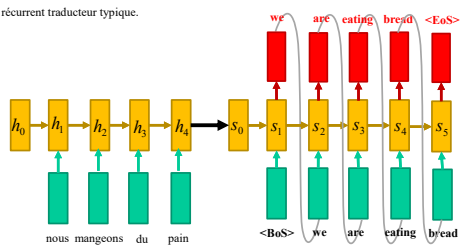
174

Seq2Seq:

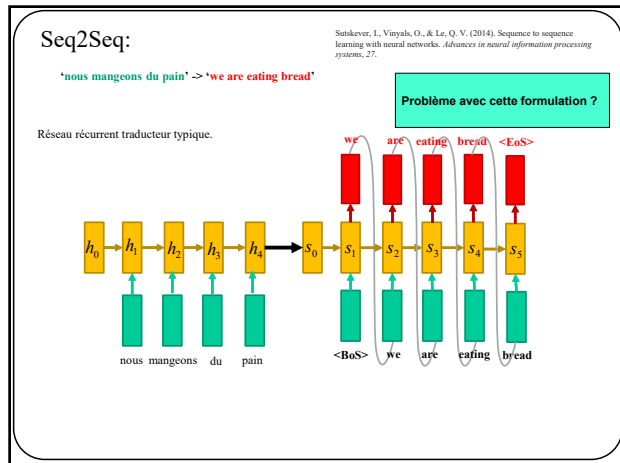
'nous mangeons du pain' -> 'we are eating bread'

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

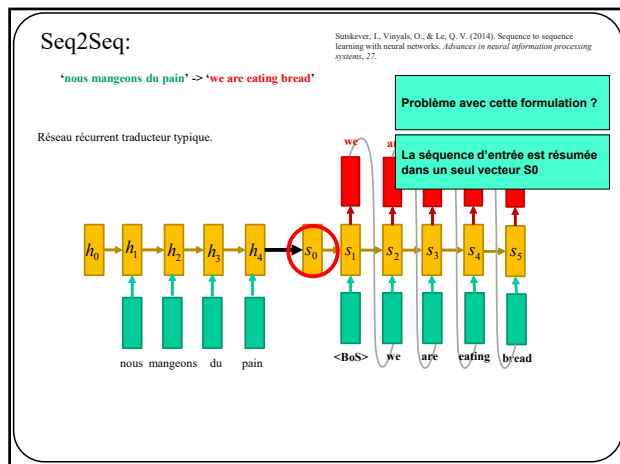
Réseau récurrent traducteur typique.



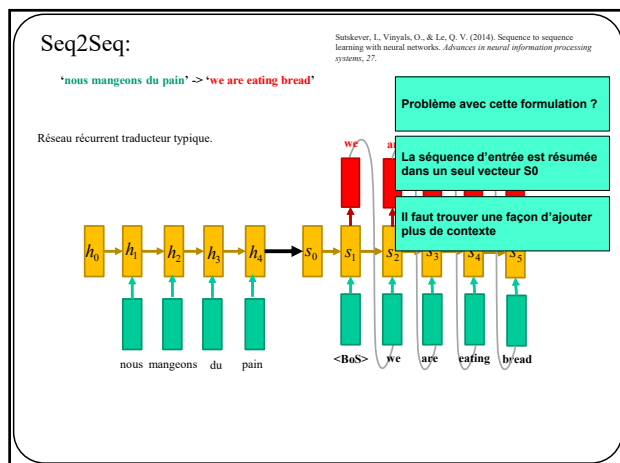
175



176



177



178

Seq2Seq avec attention

(version simplifiée)

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

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

nous mangeons du pain

s_0 = init value

179

Seq2Seq avec attention

(version simplifiée)

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

$$e_{li} = \tilde{w}^{saT} \tilde{s}_0 + \tilde{w}^{haT} \tilde{h}_i \in R$$

nous mangeons du pain

180

Seq2Seq avec attention

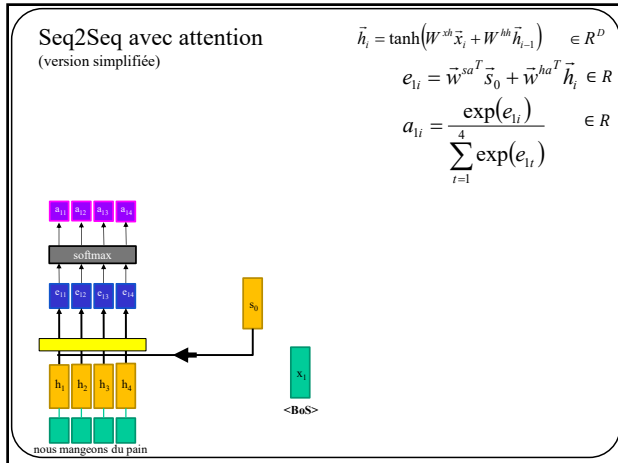
(version simplifiée)

$$\tilde{h}_i = \tanh(W^{sh} \tilde{x}_i + W^{hs} \tilde{h}_{i-1}) \in R^D$$

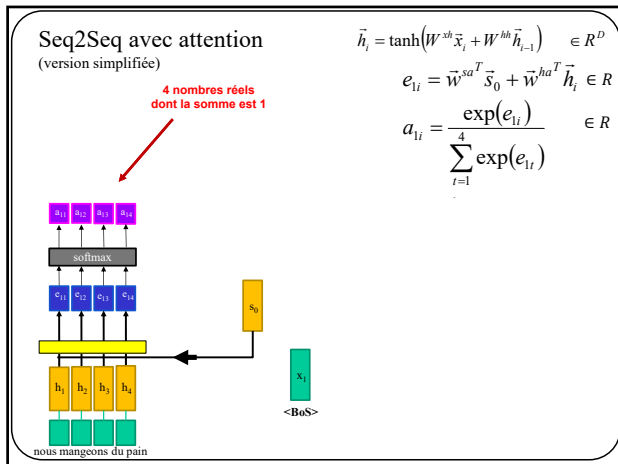
$$e_{li} = \bar{W}^{saT} \bar{s}_0 + \bar{W}^{haT} \tilde{h}_i \in R$$

The diagram illustrates a Seq2Seq model with attention. At the bottom, the input sequence "nous mangeons du pain" is shown in green boxes, which are mapped to hidden states h_1, h_2, h_3, h_4 in yellow boxes. These states are fed into a yellow rectangular block representing the encoder. Above this block, four blue boxes represent the context vectors $e_{l1}, e_{l2}, e_{l3}, e_{l4}$. A red arrow points from the text "4 nombres réels" to these blue boxes. To the right, a green box represents the source state s_0 , which is also fed into the yellow block. The output of the encoder block is a black arrow pointing to the right, representing the hidden state \tilde{h}_i , which is then used to calculate the attention weights e_{li} .

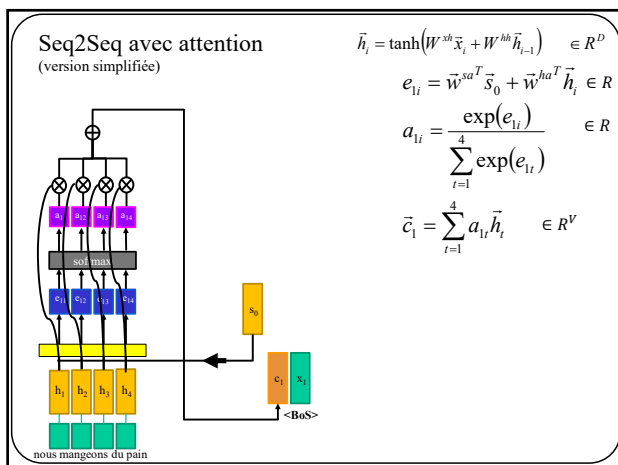
181



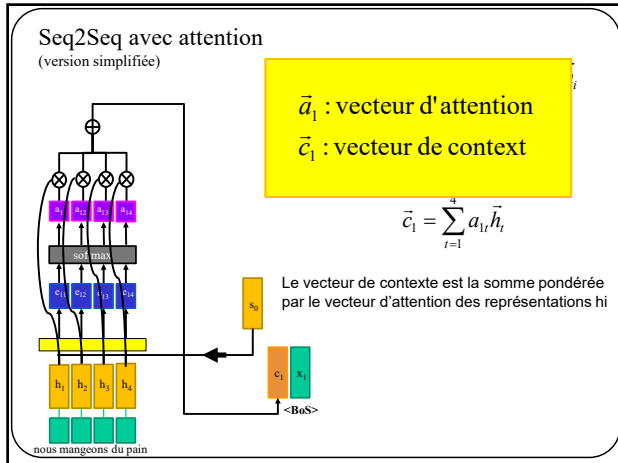
182



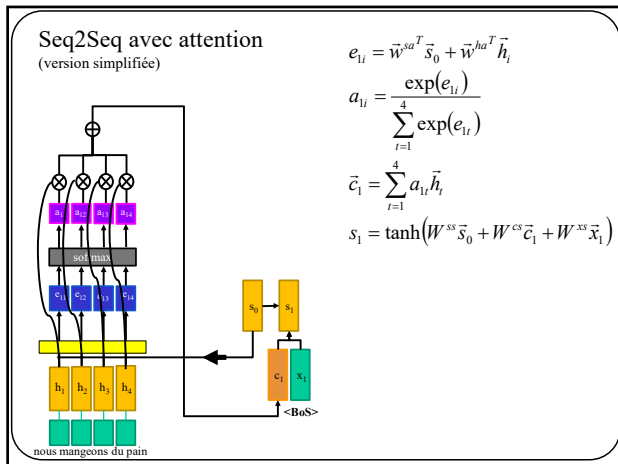
183



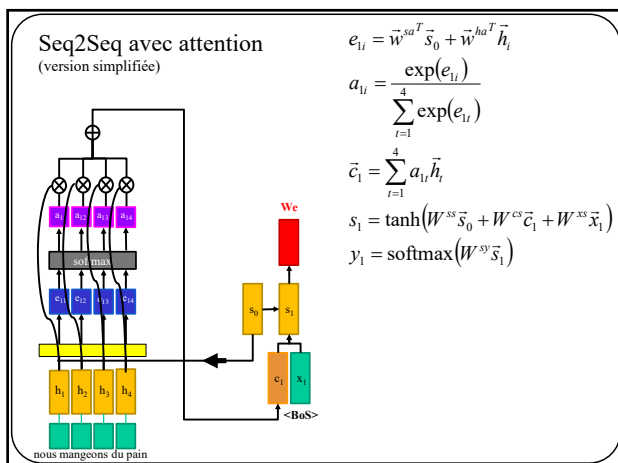
184



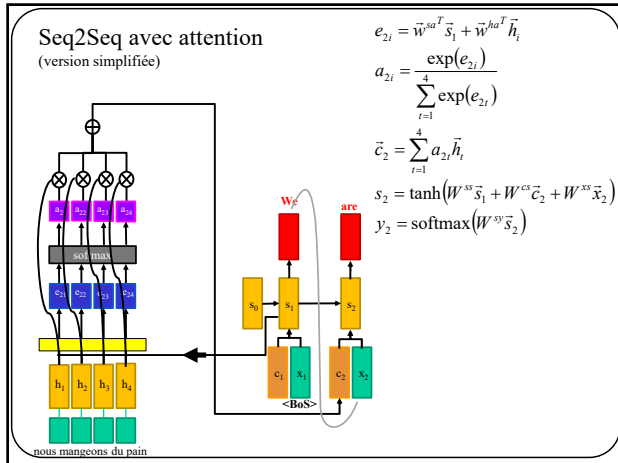
185



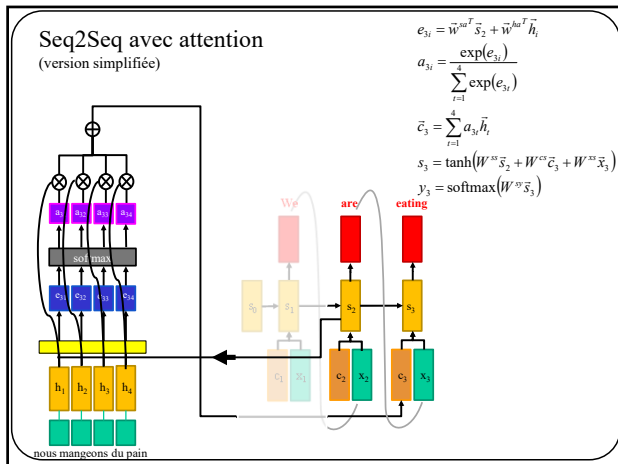
186



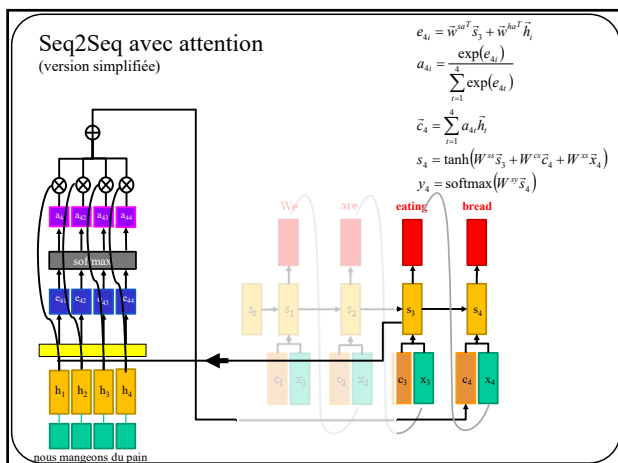
187



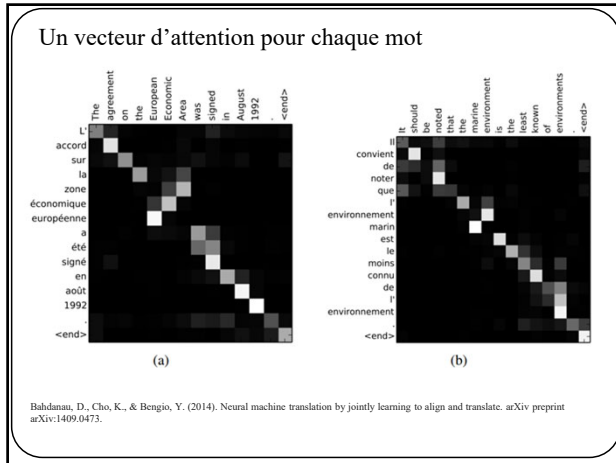
188



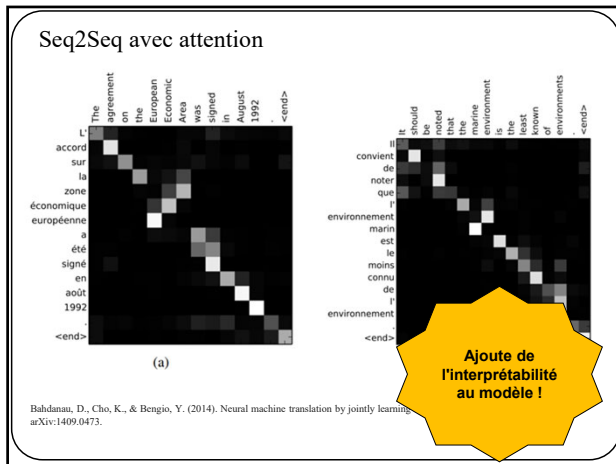
189



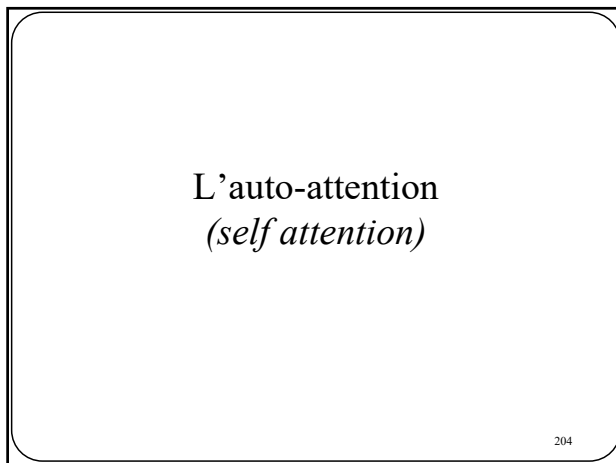
190



191



192



204

Revenons à la base : multiplication matricielle

Considérons les 4 matrices suivantes

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

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

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

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

205

205

Revenons à la base : multiplication matricielle

Leur multiplication donne:

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

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

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

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

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

$$W^v X = V = \begin{pmatrix} V^x_{11} & V^x_{12} & V^x_{13} & V^x_{14} \\ V^x_{21} & V^x_{22} & V^x_{23} & V^x_{24} \end{pmatrix} \in \mathbb{R}^{2 \times 4}$$

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

206

206

Auto attention

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

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

↑ Nous ↑ mangeons ↑ du ↑ pain

Dans cet exemple, 4 mots en entrée donc 4 colonnes dans X
Les jetons peuvent être obtenus par **Word2Vec**

207

207

Auto attention

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

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

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

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

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

↑ Nous ↑ mangeons ↑ du ↑ pain

208

208

Auto attention

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

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

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

Matrices de paramètres appris par rétropropagation

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

209

209

Auto attention

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

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

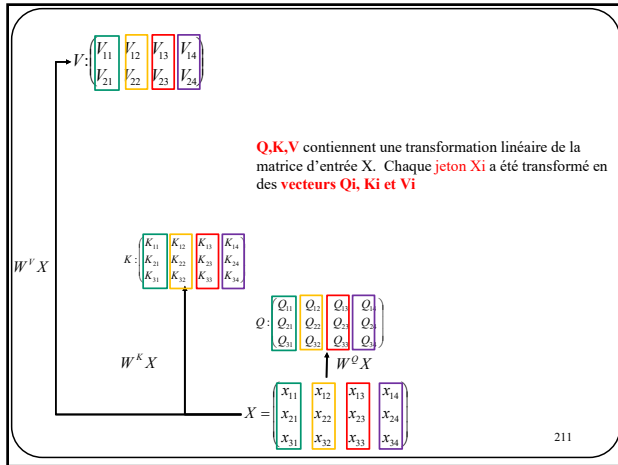
$$W^v = \begin{pmatrix} W^v_{11} & W^v_{12} & W^v_{13} \\ W^v_{21} & W^v_{22} & W^v_{23} \end{pmatrix} \in \mathbb{R}^{2 \times 3}$$

Matrices de paramètres appris par rétropropagation

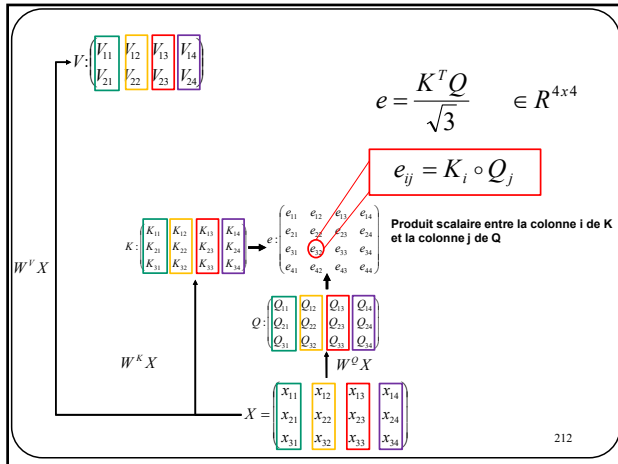
Pour ces 3 matrices, le nombre de ligne (3,3,2) est arbitraire

210

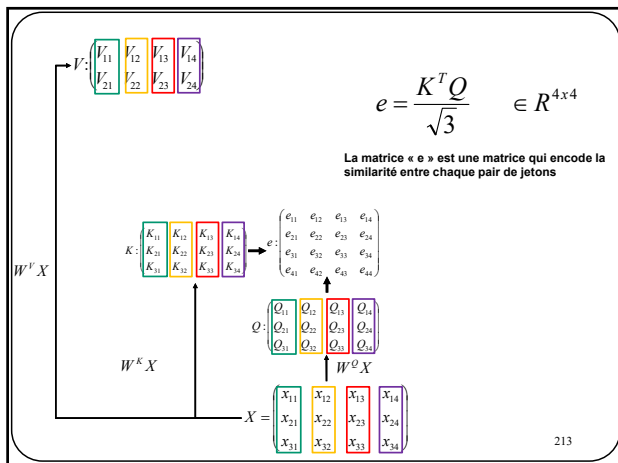
210



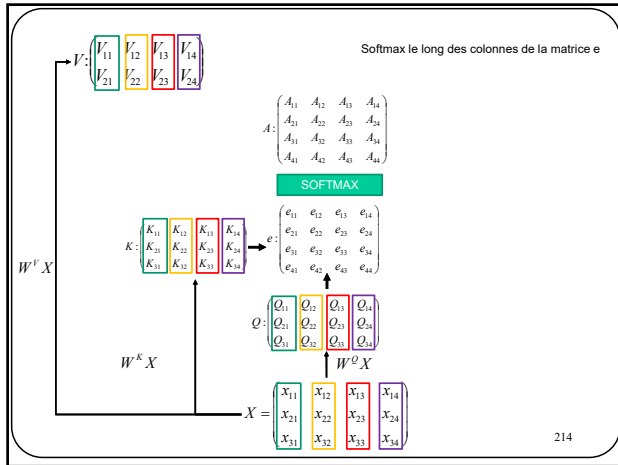
211



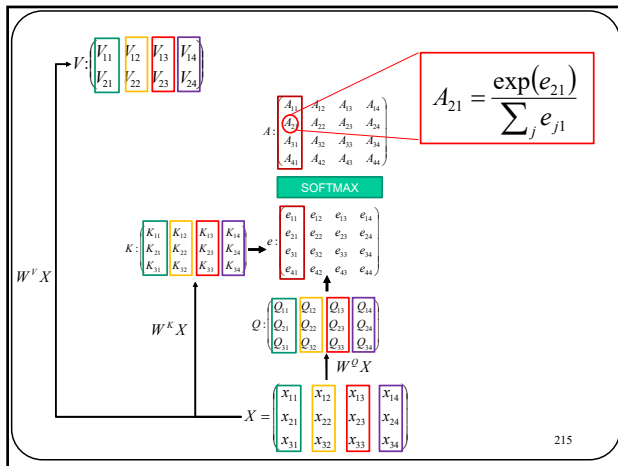
212



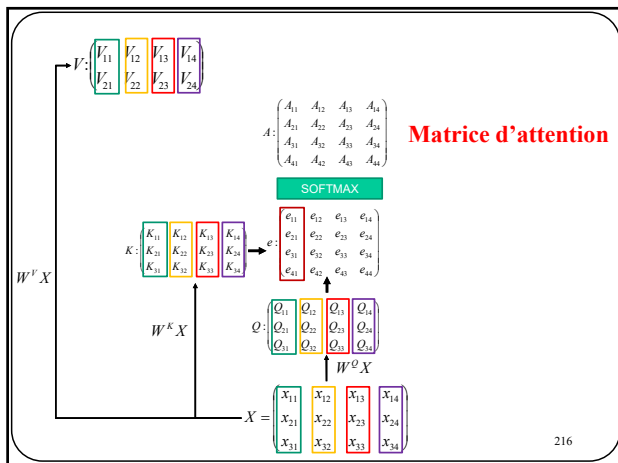
213



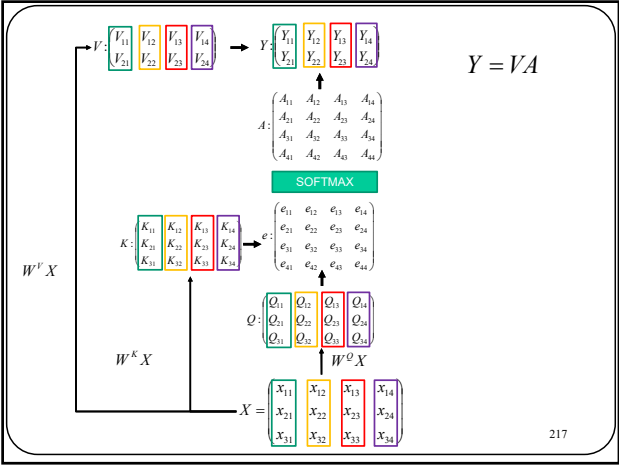
214



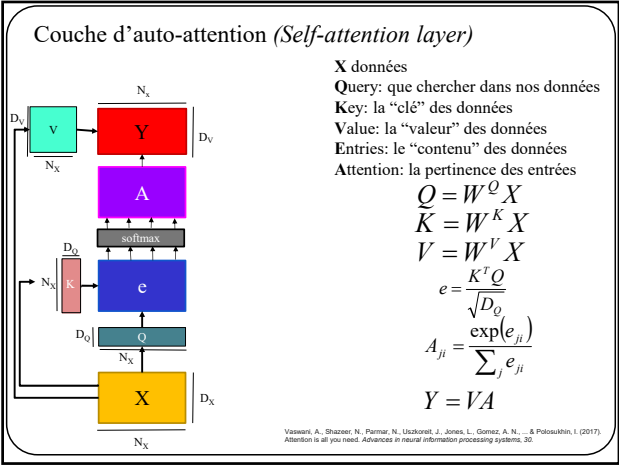
215

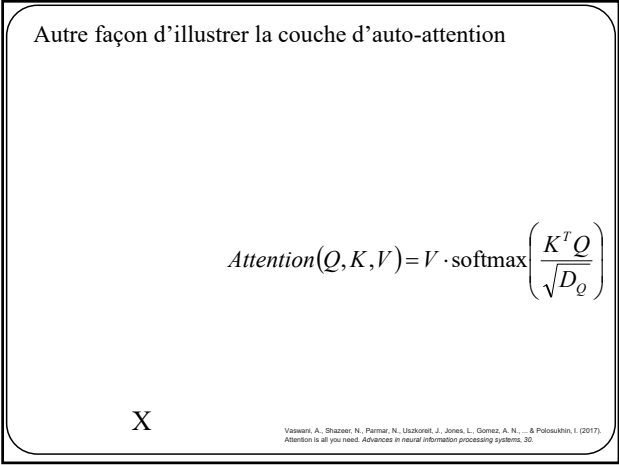


216



217





219

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

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

The diagram shows an input vector X at the bottom. Three arrows point upwards from X to vectors Q , K , and V . The arrows from X to Q and K are labeled $W^Q X$ and $W^K X$ respectively. The arrow from X to V is labeled $W^V X$.

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.

220

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

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

The diagram shows an input vector X at the bottom. Three arrows point upwards from X to vectors Q , K , and V . The arrows from X to Q and K are labeled $W^Q X$ and $W^K X$ respectively. The arrow from X to V is labeled $W^V X$. A blue box labeled "MatMul" is positioned above Q and K , with arrows pointing to them from below.

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

221

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

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

The diagram shows an input vector X at the bottom. Three arrows point upwards from X to vectors Q , K , and V . The arrows from X to Q and K are labeled $W^Q X$ and $W^K X$ respectively. The arrow from X to V is labeled $W^V X$. A blue box labeled "MatMul" is positioned above Q and K , with arrows pointing to them from below. A yellow box labeled "Scale $\frac{1}{\sqrt{D_Q}}$ " is positioned above the "MatMul" box, with an arrow pointing to it from below.

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

222

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

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

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

223

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

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

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

224

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

Autre représentation schématique

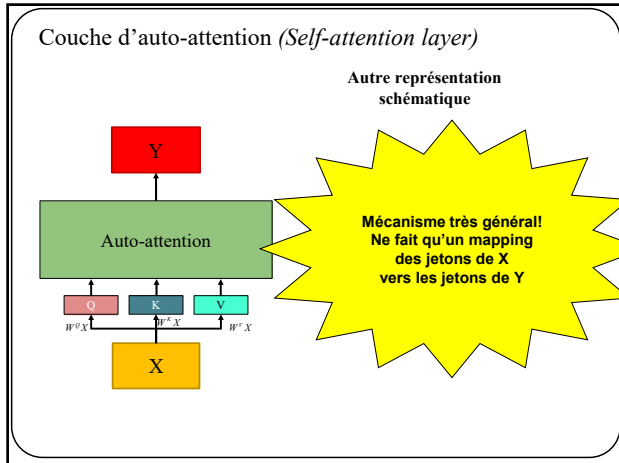
Y

Auto-attention

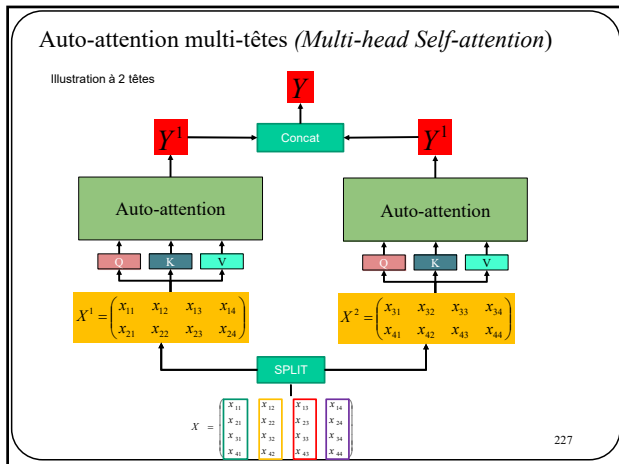
Q K V

X

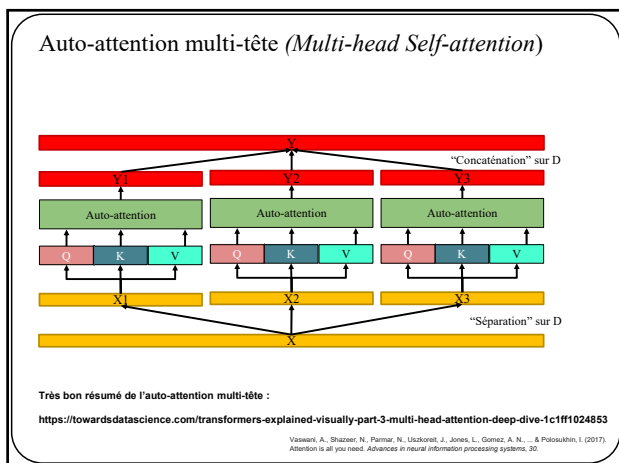
225



226



227



228

L'apothéose des réseaux de neurones

Transformer
(Attention is all you need)

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

229


Transformer

Implique aucune notion de récurrence

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

230

Transformer (Attention is all you need)

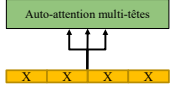


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

231

Transformer (Attention is all you need)

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

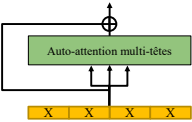


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

232

Transformer (Attention is all you need)

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

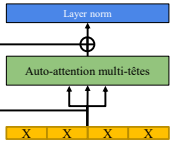


Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

233

Transformer (Attention is all you need)

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



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

234

Transformer (Attention is all you need)

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

Batch Normalization

batch				mean		std	
1	3	4	2	3			
2	2	2	2	0			
0	1	5	3	3			
4	6	1	4	3			
5	2	3	3	2			
1	0	1	1	1			
				2	3	3	
				2	2	2	

Layer Normalization

batch				mean		std	
1	3	4	2				
2	2	2	2				
0	1	5	3				
4	6	1	4				
5	2	3	3				
1	0	1	1				
				2	3	3	
				2	2	2	

Same for all training examples

Same for all feature dimensions

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.

235

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 pour chaque jeton

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.

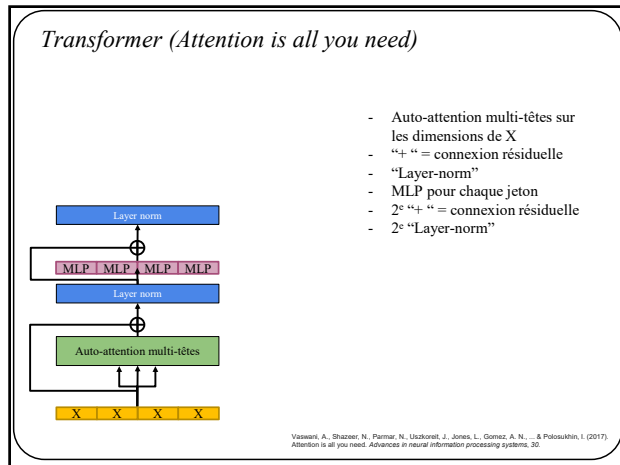
236

Transformer (Attention is all you need)

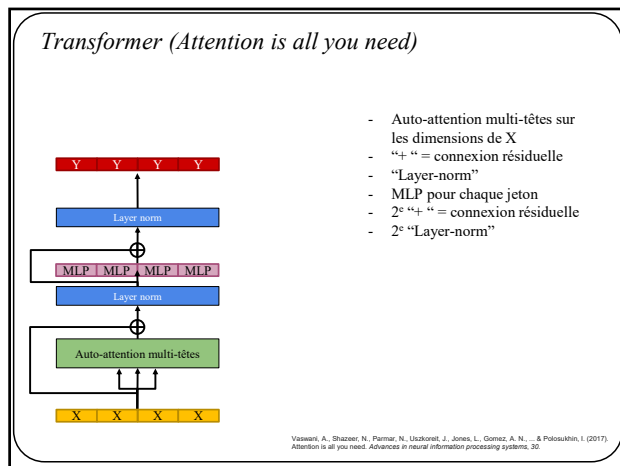
- Intra-attention multi-tête sur les dimensions de X
- "+" = connexion résiduelle
- "Layer-norm"
- MLP pour chaque jeton
- 2e "+" = connexion résiduelle

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

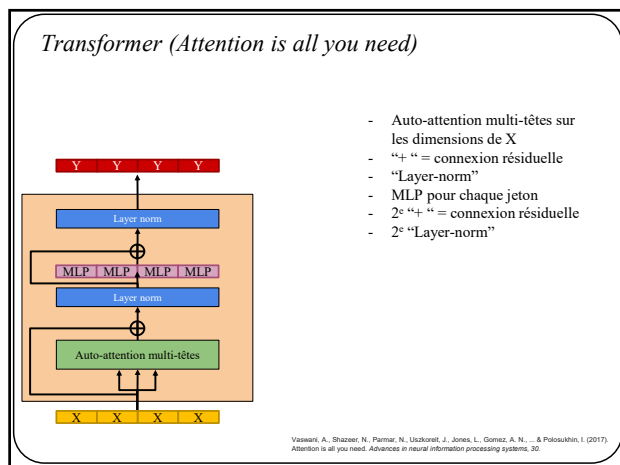
237



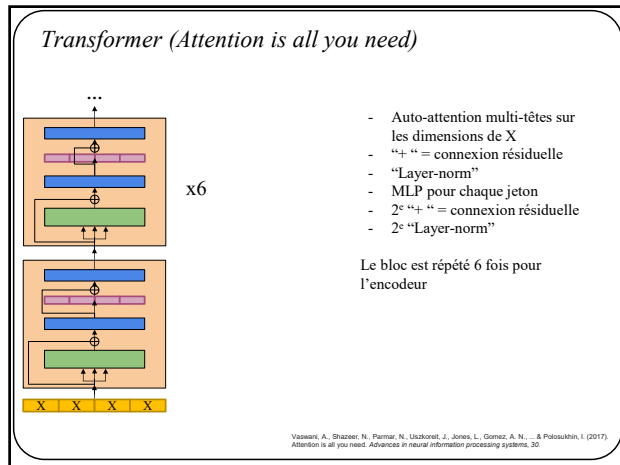
238



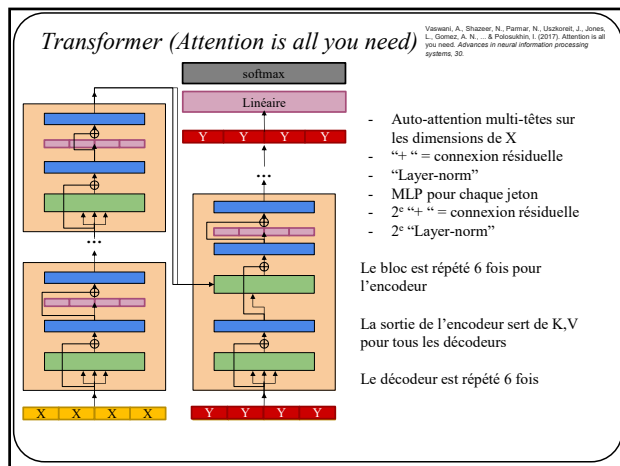
239



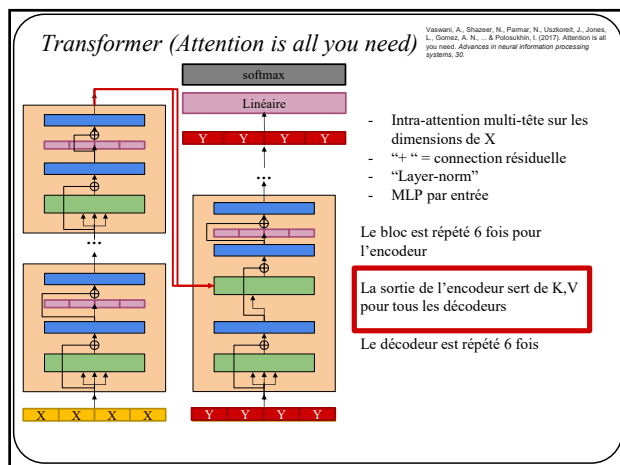
240



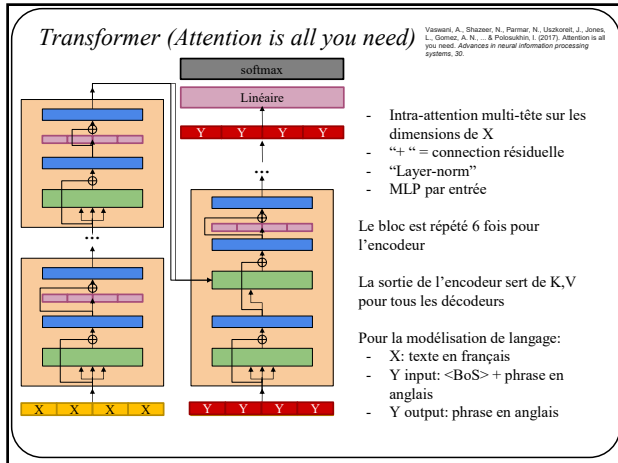
241



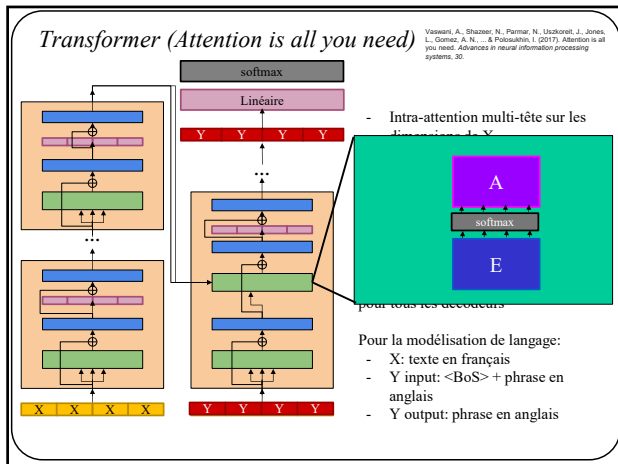
242



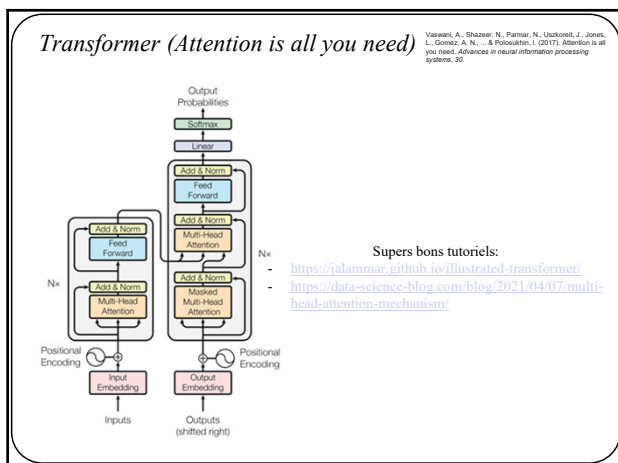
243



244



245



249

Transformer (Attention is all you need)

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

Très gourmand en mémoire:

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

Très demandant en opérations

- Complexité quadratique

250

Transformer (Attention is all you need)

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

Très gourmand en mémoire:

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

Aussi très populaires !

Attention is all you need

▲ Vaswani, N. Shazeer, N. Parmar, ... - Advances in neural ..., 2017 - proceedings.neurips.cc ... to attend to all positions in the decoder up to and including that position. We need to prevent ... We implement this inside dot-product attention by masking out (setting to -∞) ...

☆ Save 00 Cite 98797 Related articles All 62 versions 00

251

Différentes version de transformers

Source: OpenAI GPT-3

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)

Full name: GPT-3

252

Transformers

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

Full Stack ML Engineer

253

Transformers

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

Full Stack ML Engineer

254

Transformers

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

Full Stack ML Engineer

255

Transformers

Model	Layers	Width	Heads	Params	Data	Training
Transformer-Base	12	512	8	65M		8x P100 (12 hours)
Transformer-Large	12	1024	16	219M		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)

fußg le Mccvñtß Mrc kqvzq

256

GPT1-2-3

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. OpenAI blog, 1(8), 9.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. Advances in neural information processing systems, 33, 1877-1901.

Generative-Pre-Training

GTP 1,2,3 entraînent l'encodeur du transformer sur des données non annotée (d'où le « pre-training »)

GPT3 a 96 blocks d'auto-attention!

257

GPT1-2-3

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

parlement

SOFTMAX

N x

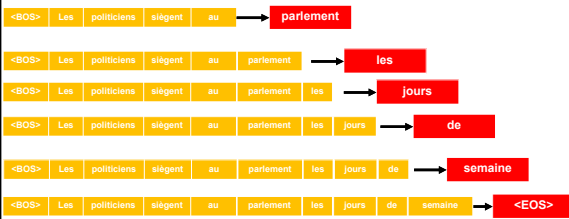
<BOS> Les politiciens siègent au

N = 96 pour GPT3

258

GPT1-2-3

Et on recommence:



259

259

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>

260

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, headache decision: how do I want to spend my post - masters year?

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

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

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

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

261

GPT-1-2-3

Rafford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.
 Rafford, 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_{tokens}	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	36	1280	20	64	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	48	1600	28	64	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	60	2048	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	80	2560	40	128	2M	1.2×10^{-4}
GPT-3 175B	175B	96	5120	60	128	2M	1.0×10^{-4}
GPT-3 175B or "GPT-3"	175B	96	5120	60	128	3.2M	0.6×10^{-4}

Table 2.1: Size, architecture, 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.

262

GPT-1-2-3

Rafford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training.
 Rafford, 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.

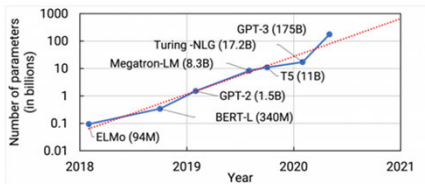


Figure 1. Trend of state-of-the-art NLP model sizes with time.

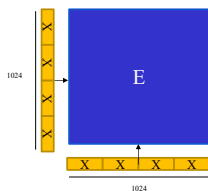
<https://developer.mda.com/blog/scaling-language-model-training-to-a-trillion-parameters-using-megatron/>

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

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

263

Vision Transformers (ViT)



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

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

264

Vision Transformers (ViT)

Opération (au moins) quadratique

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissert, 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.

265

Vision Transformers (ViT)

Opération (au moins) quadratique

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissert, 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.

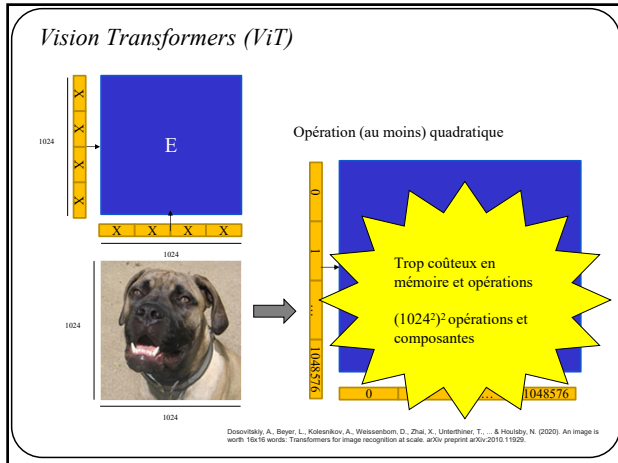
266

Vision Transformers (ViT)

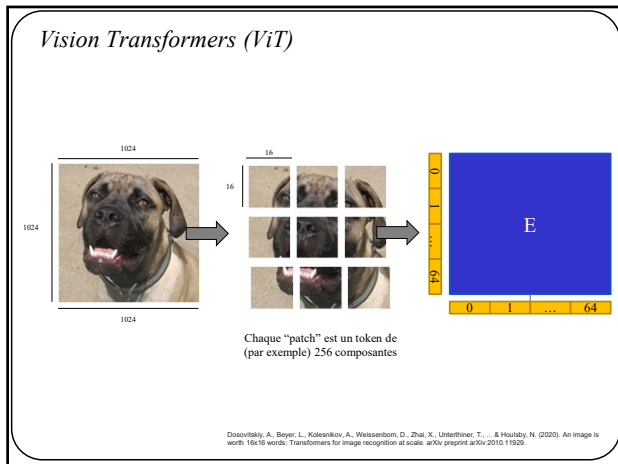
Opération (au moins) quadratique

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissert, 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.

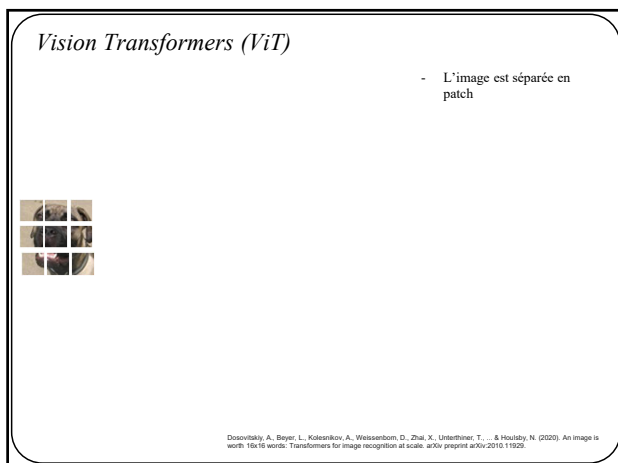
267



268




269



270

Vision Transformers (ViT)

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




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

271

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., Weissert, D., Zhai, X., Unterthiner, T., ... & Houtisby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

272

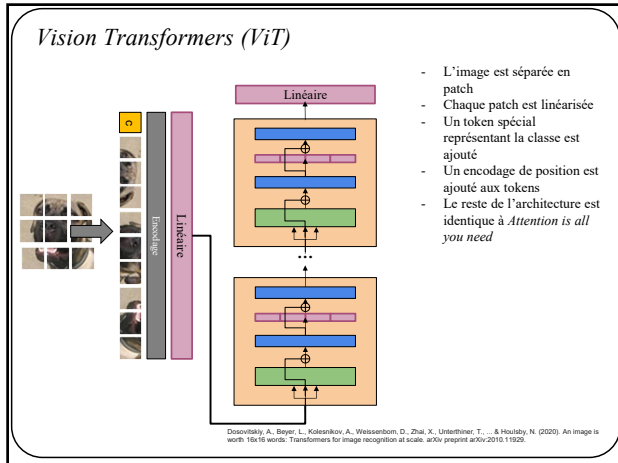
Vision Transformers (ViT)

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

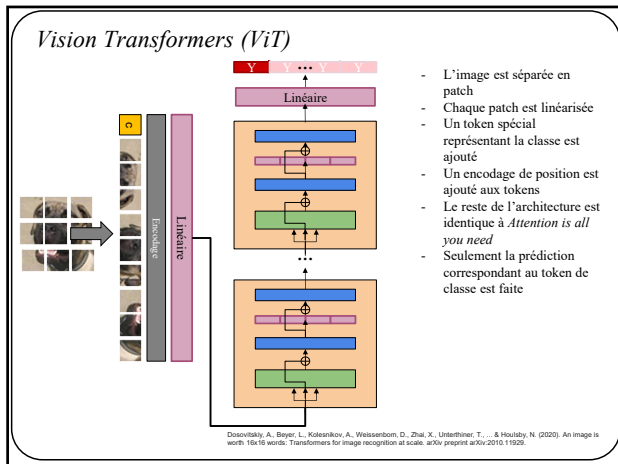


Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissert, D., Zhai, X., Unterthiner, T., ... & Houtisby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

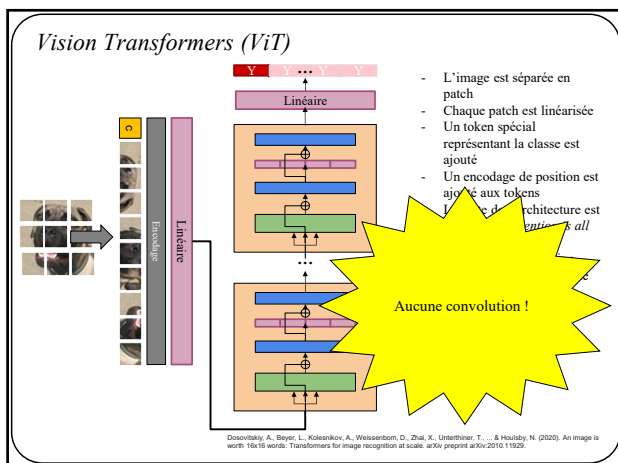
273



274



275



276

Vision Transformers (ViT)

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

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

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

277

Vision Transformers (ViT)

Les vision transformers dominent la classification depuis leur arrivée

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

278

Vision Transformers (ViT)

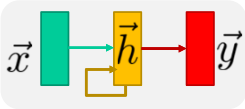
Rank	Model	Top-1 Accuracy	Top-5 Accuracy	Params (M)	Epochs	Year	Type	Link
1	DeiT-S/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
2	DeiT-T/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
3	DeiT-L/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
4	DeiT-XL/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
5	DeiT-XXL/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
6	DeiT-S/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
7	DeiT-T/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
8	DeiT-L/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
9	DeiT-XL/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link
10	DeiT-XXL/36	90.8%	98.2%	1.7B	300k	2021	ViT	Link

Les vision transformers dominent la classification depuis leur arrivée

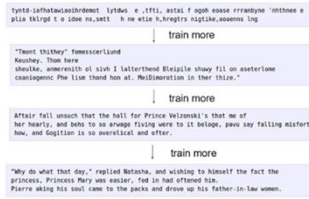
<https://paperswithcode.com/sota/image-classification-on-imagenet>

279

Sommaire



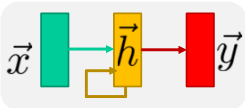
- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique



Modélisation de langage

280

Sommaire



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

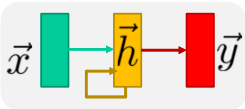


a group of people playing a game with nintendo wii controllers

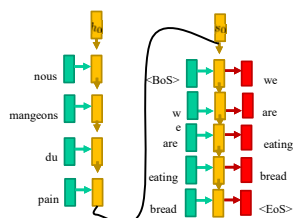
Description d'images

281

Sommaire



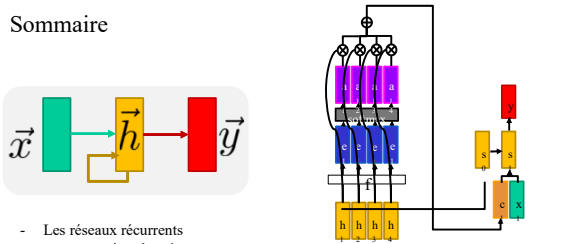
- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique



Traduction

282

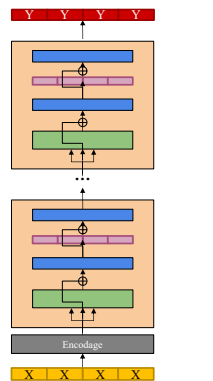
Sommaire



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

283

Sommaire



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

284
