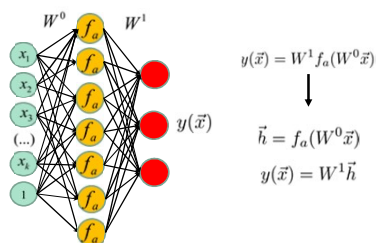


Réseaux de neurones IFT 780

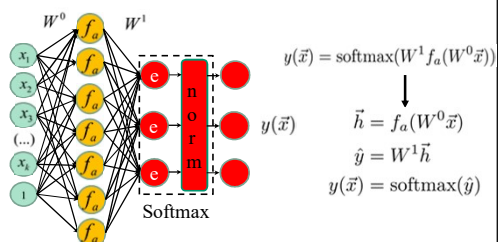
Réseaux récurrents
Par
Pierre-Marc Jodoin, Antoine Th  berge

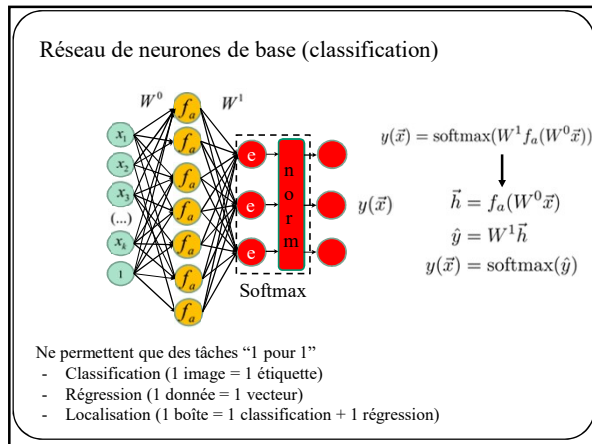
R  seau de neurones de base (r  gression)

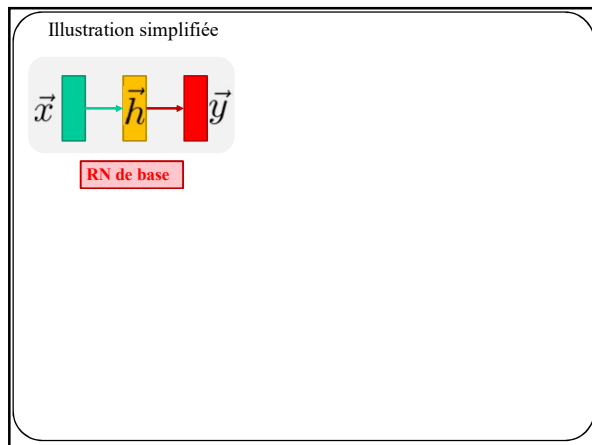


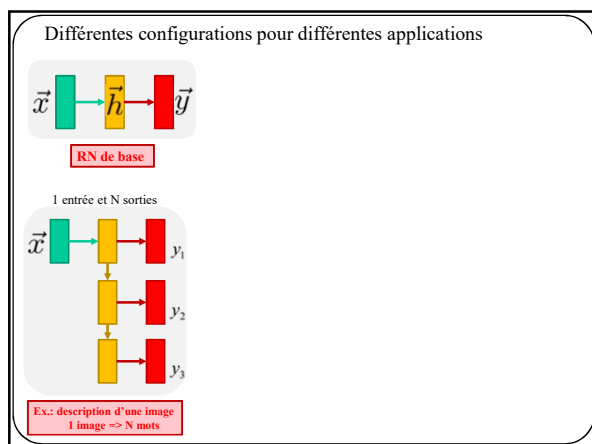
f_a : fonction d'activation

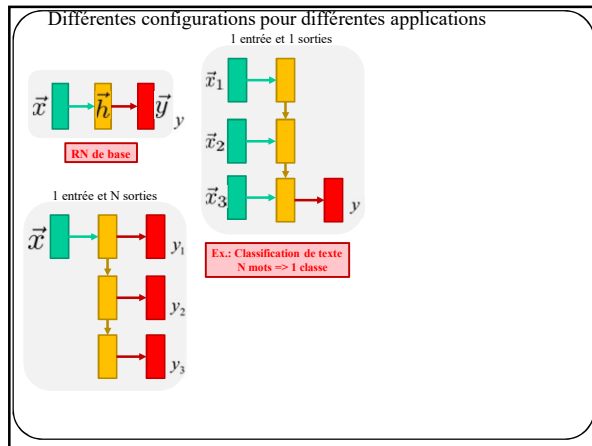
R  seau de neurones de base (classification)

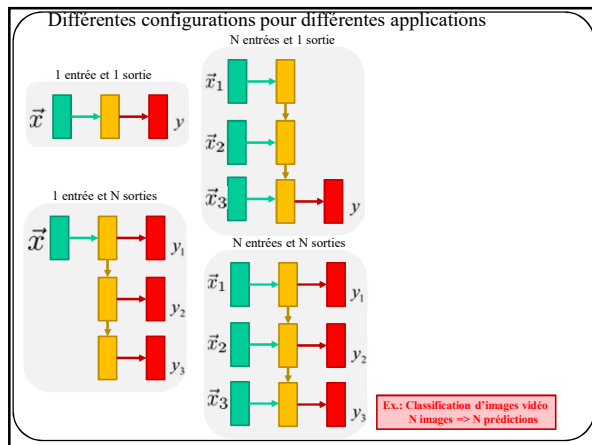


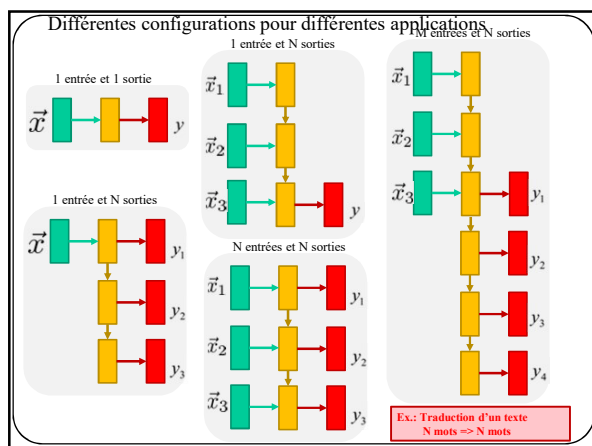




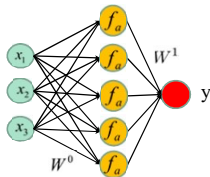




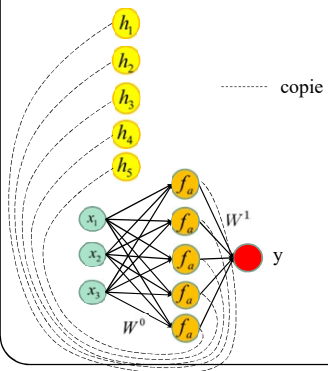




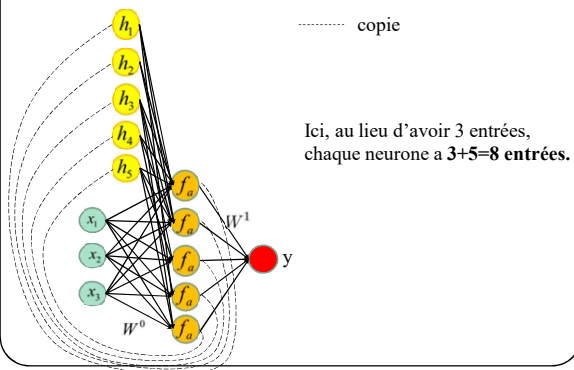
Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée

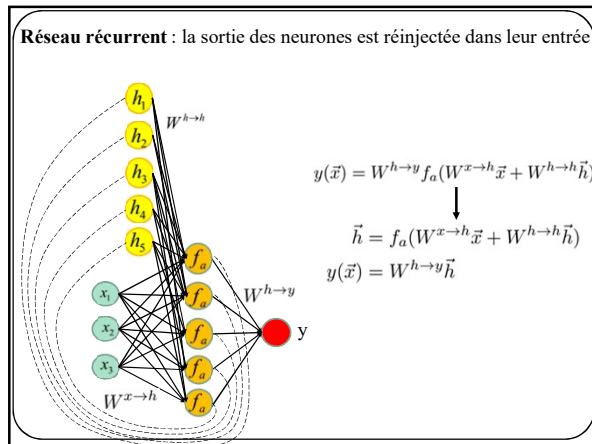


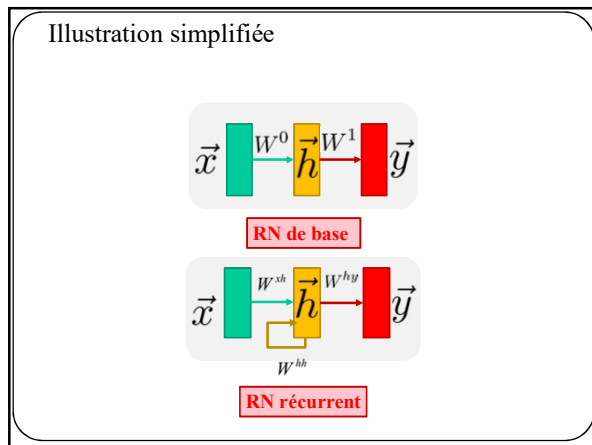
Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée

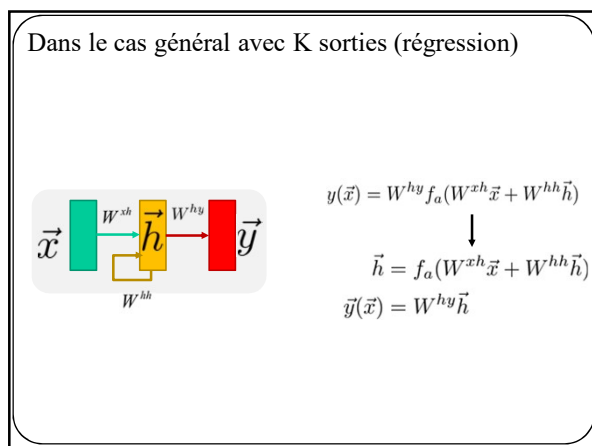


Réseau récurrent : la sortie des neurones est réinjectée dans leur entrée

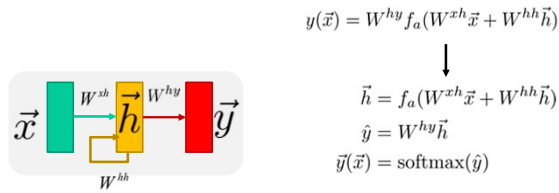






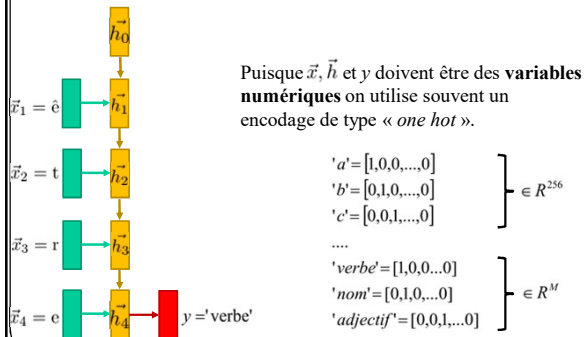


Dans le cas général avec K sorties (classification)



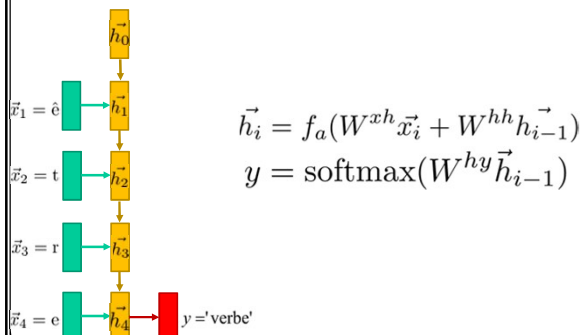
Exemple pour N entrées et 1 sortie:

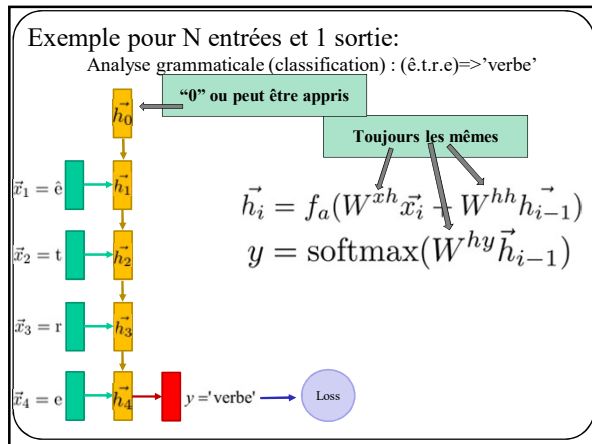
Analyse grammaticale (classification) : (ê.t.r.e) => 'verbe'

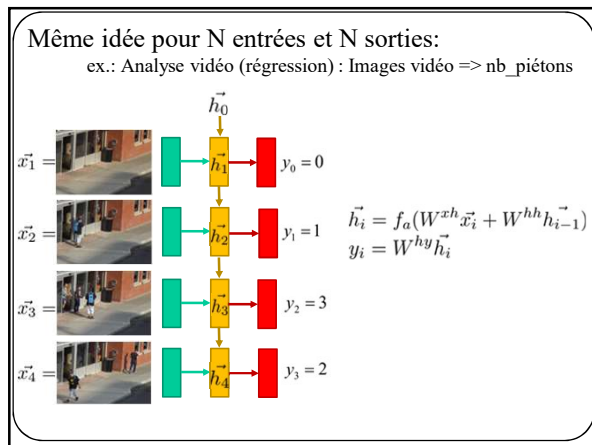


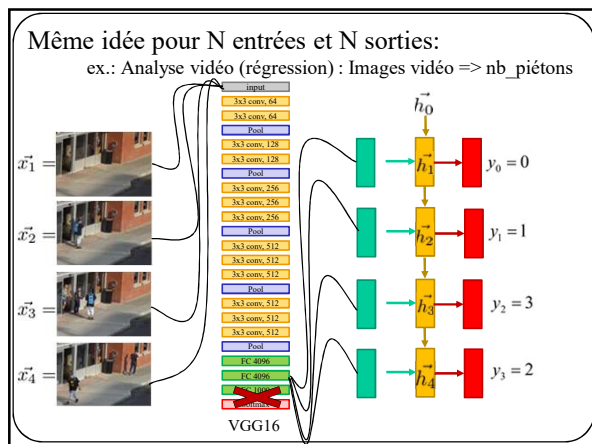
Exemple pour N entrées et 1 sortie:

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



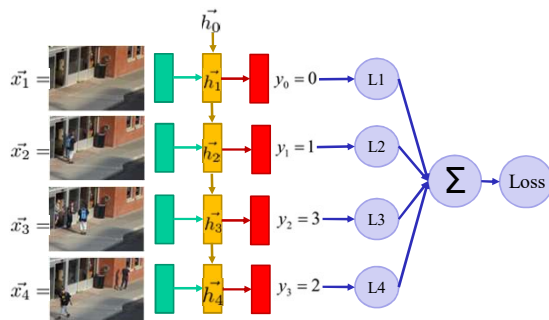






Même idée pour N entrées et N sorties:

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



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 ».

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

Alphabet : [a,e,m,s]

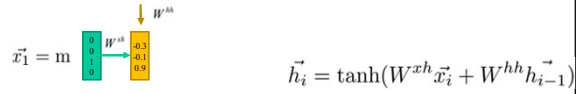
Entraîner un modèle à prédire les lettres du mot « masse ».

$\vec{x}_1 = m$ 

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

Alphabet :{a,e,m,s}

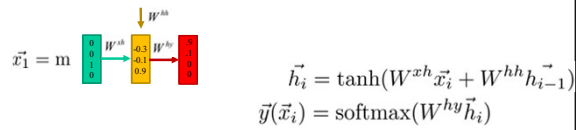
Entraîner un modèle à prédire les lettres du mot « masse ».



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

Alphabet :{a,e,m,s}

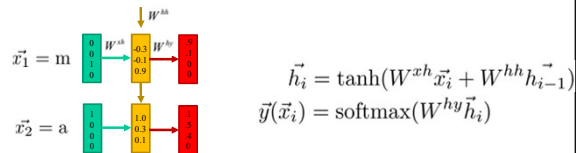
Entraîner un modèle à prédire les lettres du mot « masse ».



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

Alphabet :{a,e,m,s}

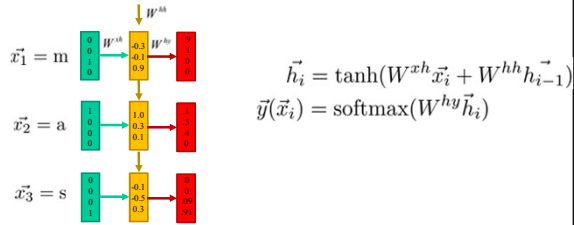
Entraîner un modèle à prédire les lettres du mot « masse ».



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

Alphabet :{a,e,m,s}

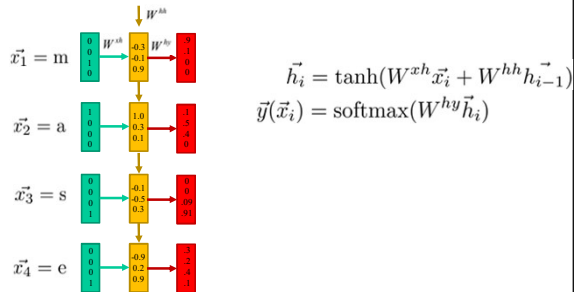
Entraîner un modèle à prédire les lettres du mot « masse ».



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

Alphabet :{a,e,m,s}

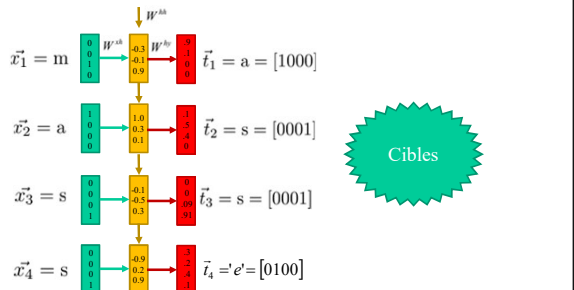
Entraîner un modèle à prédire les lettres du mot « masse ».

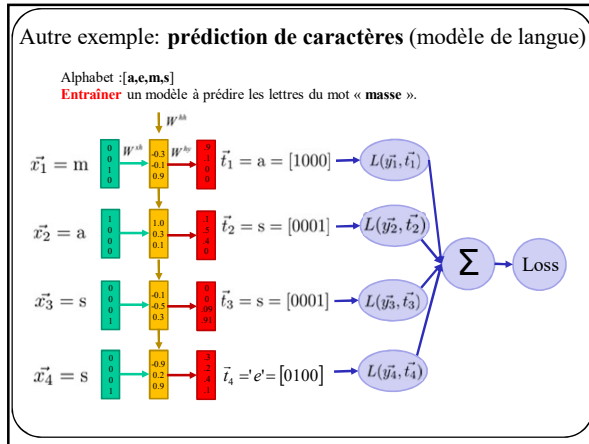


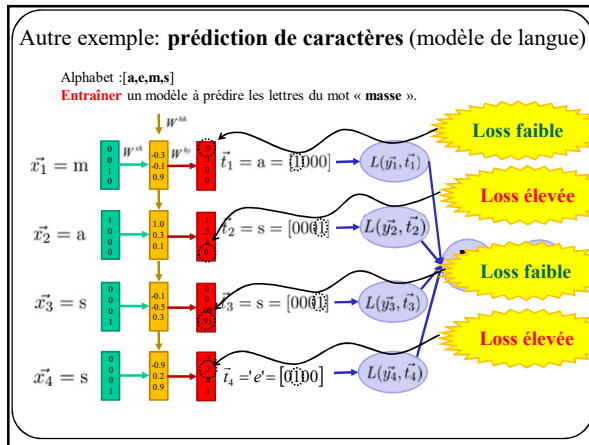
Autre exemple: **prédiction de caractères** (modèle de langue)

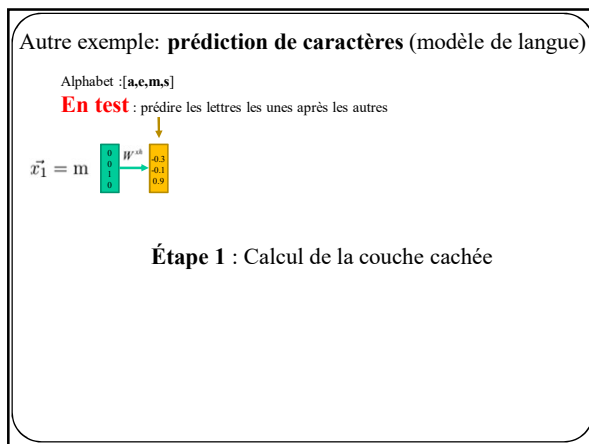
Alphabet :{a,e,m,s}

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









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

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

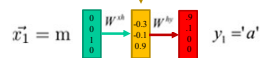


Étape 2 : Calcul de la sortie (softmax)

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

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

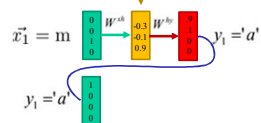


Étape 3 : Sélectionner le caractère le plus probable

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

Alphabet : [a,e,m,s]

En test : prédire les lettres les unes après les autres

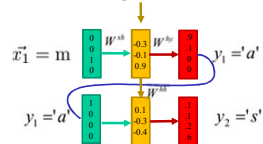


Étape 4 : Injecter le caractère prédit au début du réseau

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

Alphabet :{a,e,m,s}

En test : prédire les lettres les unes après les autres

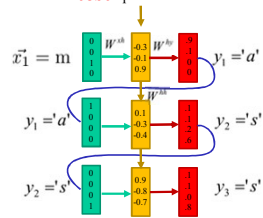


Et on recommence!

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

Alphabet :{a,e,m,s}

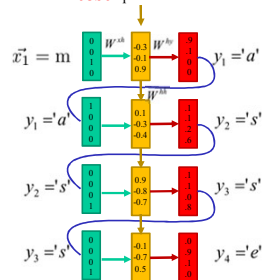
En test : prédire les lettres les unes après les autres



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

Alphabet :{a,e,m,s}

En test : prédire les lettres les unes après les autres



Code python: “mini-char-RNN” de A. Karpathy
<https://gist.github.com/karpathy/d4dee566867f8291f086>
 Un RNN en 112 lignes !

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

Code python: “mini-char-RNN” de A. Karpathy
<https://gist.github.com/karpathy/d4dee566867f8291f086>

THE SONNETS

by William Shakespeare

From latent creamers we desire increase,
That thereby beauty's rose might never die,
But as the ripen should by time decrease,
His tender hue might bear his memory;
But then, contracted to thine own bright eyes,
Feed'st thy light's flame with self-substantial fuel,
Making a famine where abundance lies,
Thyself thy foe, to thy sweet self too cruel;
Thou that art now the world's fresh ornament,
And only herald to the gaudy spring,
Within thine own bud hast hidden thy content,
And tender chaf'd male's waste in suggesting;
Pity the world, or else this glutton be,
To eat the world's due, by the grave and thee.

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

When forty winters shall besiege thy brow,
And dig deep trenches in thy beauty's field,
Thy youth's proud livery so gazed on now,
Will be a tattered weed of small worth held:
Then being asked, where all this beauty lies,
Where all the treasure of thy lusty days,
To say, within this bowditch solitary,
Where at all-making shamer, and thriftless pain,
How much more praise deserved thy beauty's use,
If thou couldst answer 'Tis this child of mine
Shall sum my beauty, and mine old excuse,
Proving I live by patterns young as thou:
Thou wast as new made when thou first didst old,
And see the blood made when thou first didst cold.

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

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

train more

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

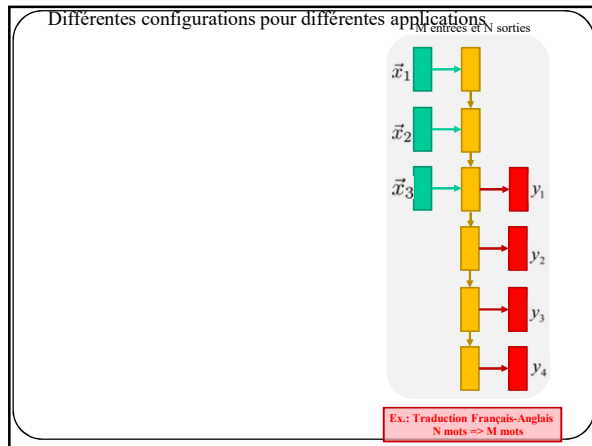
train more

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

train more

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

Crédit: A. Karpathy, CS231



Autre exemple: traduction

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

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

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

Autre exemple: traduction

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

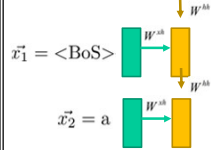
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

$x_1 = \text{<BoS>}$

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.

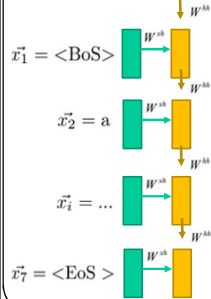
Traduire 'assez' -> 'enough'
 Alphabet fr : [**<BoS>**,a,e,s,z,<EoS>]
 Alphabet en : [**<BoS>**,e,g,h,u,o,u,<EoS>]



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.

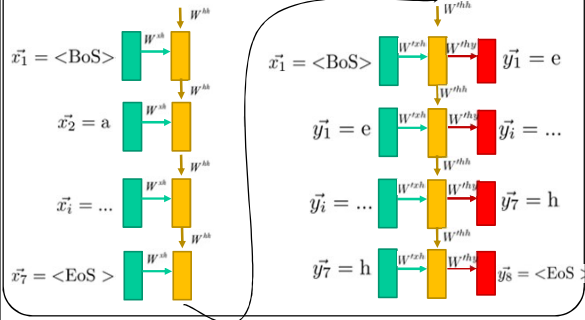
Traduire 'assez' -> 'enough'
 Alphabet fr : [**<BoS>**,a,e,s,z,<EoS>]
 Alphabet en : [**<BoS>**,e,g,h,u,o,u,<EoS>]

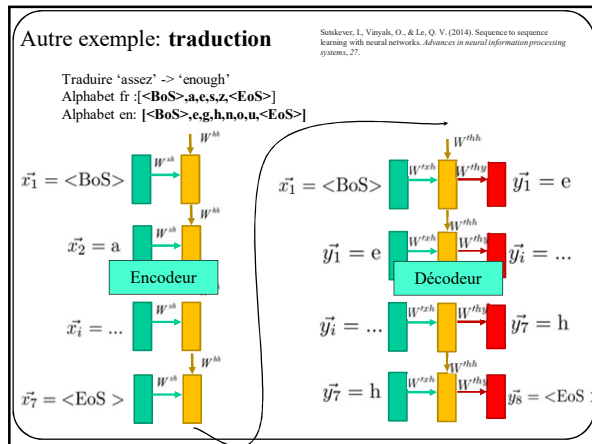


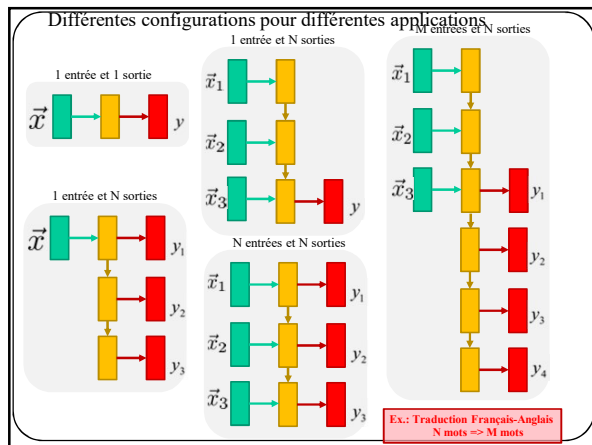
Autre exemple: traduction

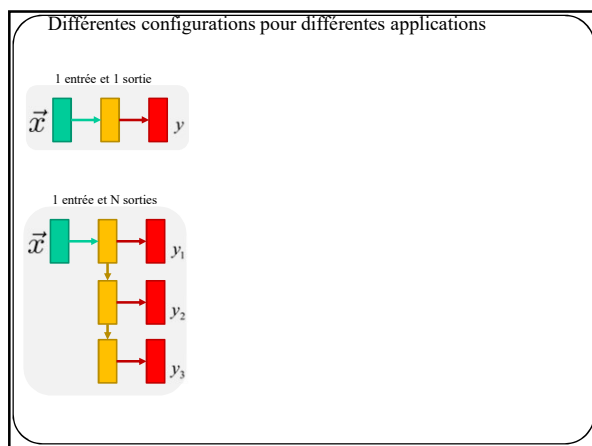
Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

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









Captioning




Diagram of VGG16 architecture for captioning:

- 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

Pré-entraîné sur ImageNet

VGG16

55

Captioning




Diagram of VGG16 architecture for captioning:

- 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~~

VGG16

Captioning




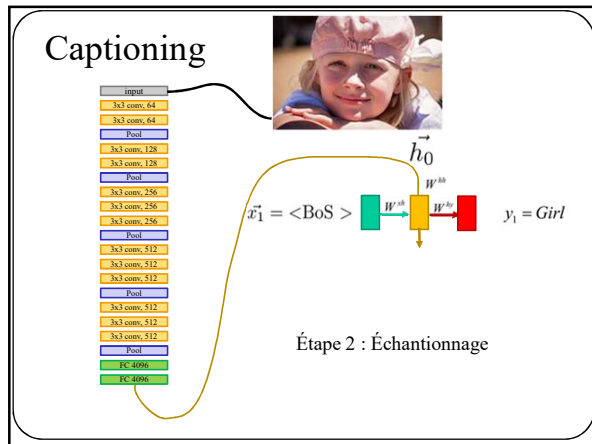
Diagram of VGG16 architecture for captioning:

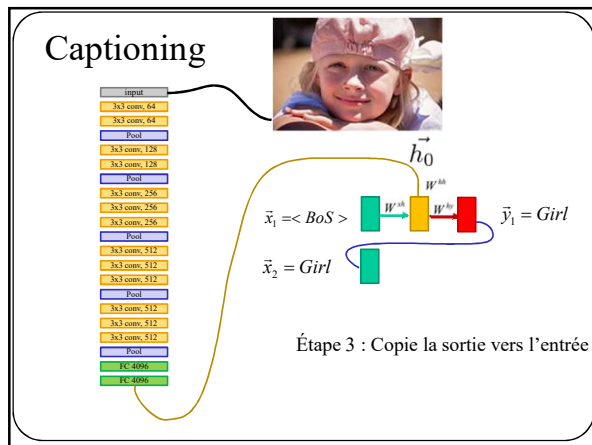
- input
- 3x3 conv, 64
- 3x3 conv, 64
- Pool
- 3x3 conv, 128
- 3x3 conv, 128
- Pool
- 3x3 conv, 256
- 3x3 conv, 256
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- 3x3 conv, 512
- 3x3 conv, 512
- Pool
- FC, 4096
- FC, 4096
- FC, 1000
- Softmax

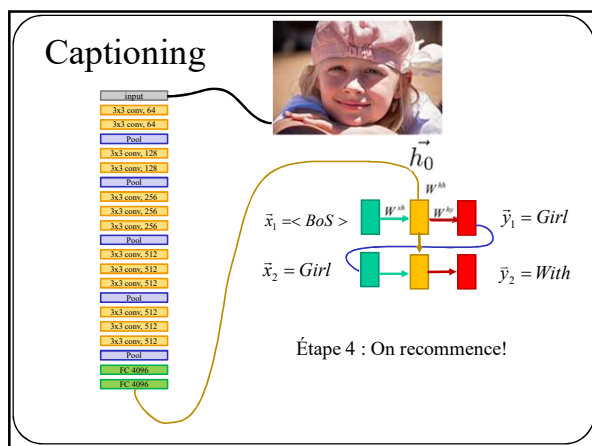
Diagram illustrating the first step of the captioning process:

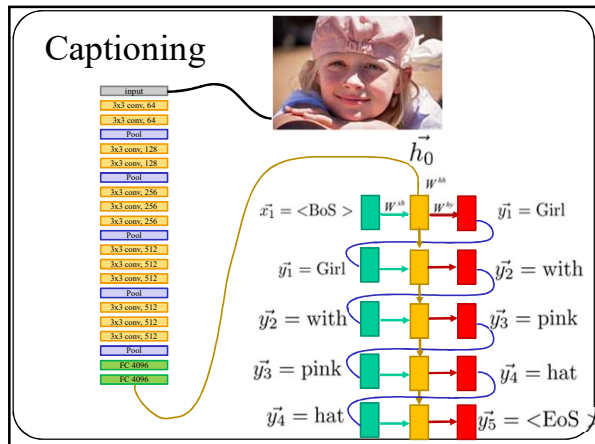
$\vec{x}_1 = \langle \text{BoS} \rangle$ (green box) $\xrightarrow{W^{ih}}$ \vec{h}_0 (yellow box) $\xrightarrow{W^{oh}}$ (red box)

Étape 1 : Init + Propagation avant













NeuralTalk and Walk

<https://vimeo.com/146492001>



64

Analyse de texte

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

Pour des **caractères**...

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

65

Analyse de texte

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

Pour des **mots**...

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

66

Prédiction sur des lettres vs. mots

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

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

Prédiction sur des fractions de mots

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

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

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.

Senrich, 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'(?!\b)%s\b' % bigram)
    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,
        'h o w s </w>': 4, 'h o w 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)
```

l' → f'
l e → l o
l o w → l o w
e f' → e f'

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.

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

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

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

Word2Vec s'appuie sur 2 idées fondamentales

Un 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

1 ligne = code
pour 1 mot

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!)

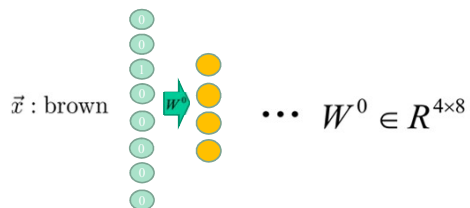
Ex: sélectionner le code de « brown »

$$\begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} \text{Dictionnaire} \\ \text{(matrice d'encodage)} \\ \begin{array}{cccc} 2 & 3 & 4 & 5 \\ -1 & -3 & -2 & 2 \\ 11 & 6 & 4 & -3 \\ -4 & 8 & -4 & 4 \\ 24 & -6 & 42 & 17 \\ 91 & 13 & 14 & -5 \\ 0 & 36 & 4 & 56 \\ -1 & 0 & 1 & 35 \end{array} \end{pmatrix} = \begin{pmatrix} 11 & 6 & 4 & -3 \end{pmatrix}$$

Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

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



Word2Vec s'appuie sur 2 idées fondamentales

Idée 1: Dictionnaire = matrice d'encodage

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

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



Word2Vec s'appuie sur 2 idées fondamentales

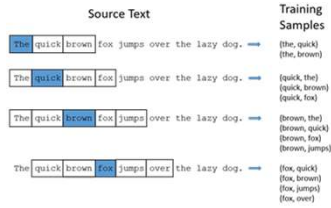
Idée 1: Dictionnaire = matrice d'encodage

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



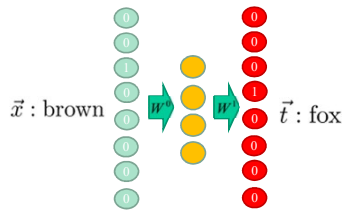
Word2Vec s'appuie sur 2 idées fondamentales

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



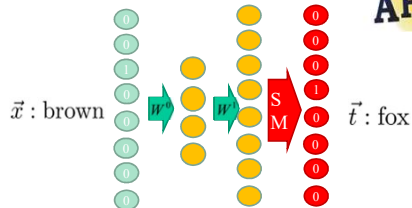
Basé sur un corpus de texte, on va créer des **millions de paires de mots**

Word2Vec [Mikolov et al. '13]



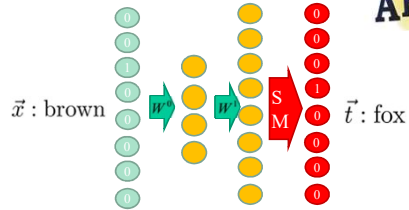
Entraîner un réseau de neurones
à reproduire le 2^e mot partant du 1^{er}

Word2Vec [Mikolov et al. '13]



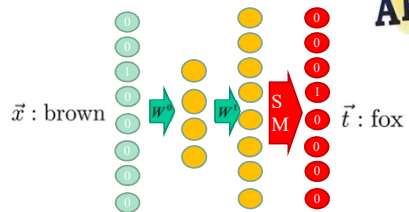
Puisque la sortie est de type « one-hot »
on utilise un softmax

Word2Vec [Mikolov et al. '13]



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

Word2Vec [Mikolov et al. '13]



Lorsqu'entraîné, utiliser W^0
comme dictionnaire

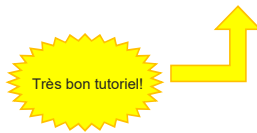
Word2Vec [Mikolov et al. '13]

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/>

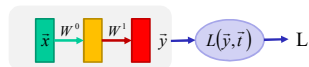


T.Mikolov et al. (2013). "Efficient Estimation of Word Representations in Vector Space", in ICLR 2013

Comment entraîner un RNN?

Histoire de gradients

RN de classification avec entropie croisée

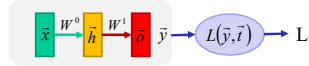


$$\tilde{y}(\tilde{x}) = S_M \left(W^1 \tanh \left(W^0 \tilde{x} \right) \right)$$

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

Histoire de gradients

Simple RN de classification avec entropie croisée



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

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

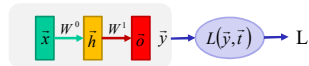
$$\bar{y} = S_M(\bar{o})$$

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

↓
Propagation
avant

Histoire de gradients

Simple RN de classification avec entropie croisée



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

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

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

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

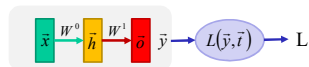
Pour entraîner le réseau
il faut calculer

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

↓

Histoire de gradients

Simple RN de classification avec entropie croisée



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

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

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

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

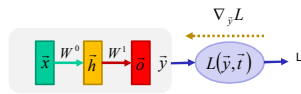
Dérivée en chaîne

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

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

↓

Histoire de gradients



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

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

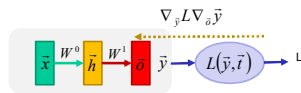
$$\bar{y} = S_M(\bar{o})$$

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

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

Rétro-propagation

Histoire de gradients



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

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

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

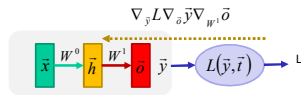
$$L = L_{CE}(\bar{y}, \bar{t})$$

$$\nabla_{\bar{o}} \bar{y} = \text{Id} \bar{y}^T - \bar{y}^T \bar{y}$$

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

Rétro-propagation

Histoire de gradients



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

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

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

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

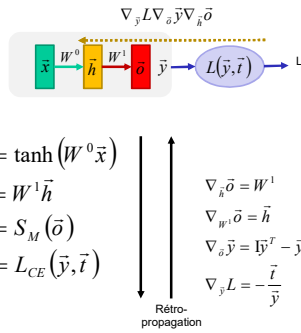
$$\nabla_{\bar{w}^1} \bar{o} = \bar{h}$$

$$\nabla_{\bar{o}} \bar{y} = \text{Id} \bar{y}^T - \bar{y}^T \bar{y}$$

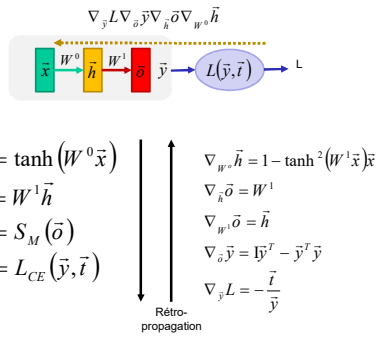
$$\nabla_{\bar{y}} L = -\frac{\bar{t}}{\bar{y}}$$

Rétro-propagation

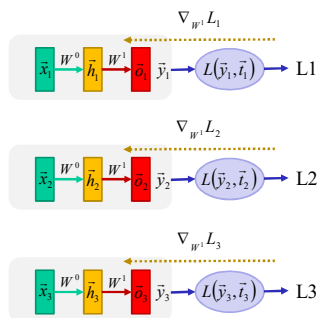
Histoire de gradients



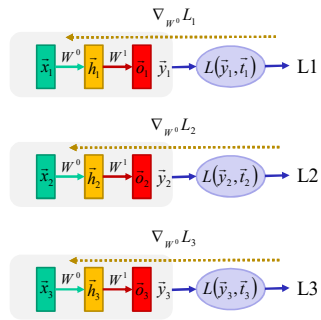
Histoire de gradients



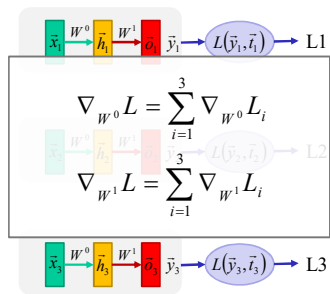
Ex.: 3 données, 3 rétro-propagations

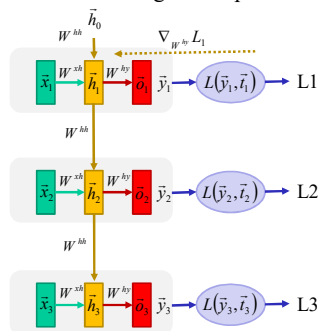


Ex.: 3 données, 3 rétro-propagations

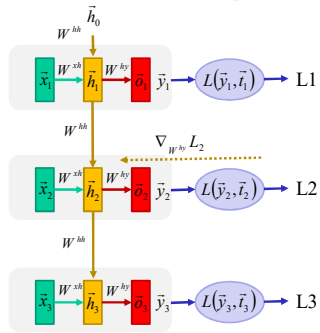


3 rétro-propagations

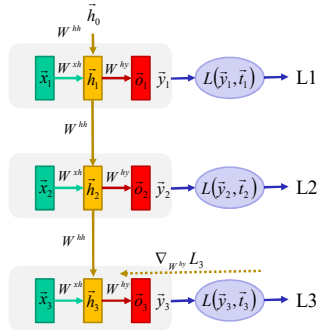


Réseau récurrent: gradient pour W^{hy} 

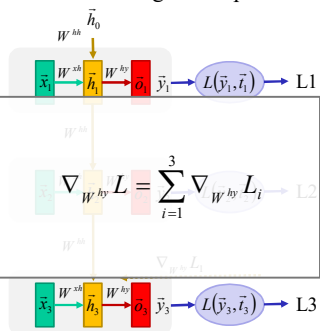
Réseau récurrent: gradient pour W^{hy}



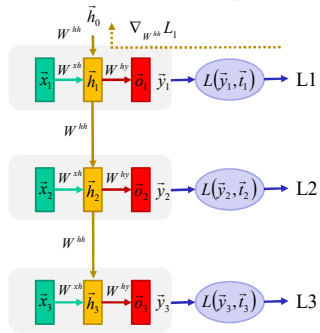
Réseau récurrent: gradient pour W^{hy}



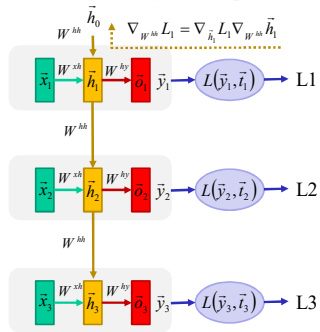
Réseau récurrent: gradient pour W^{hy}



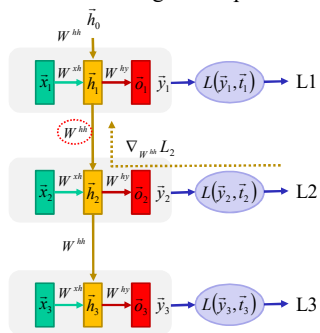
Réseau récurrent: gradient pour W^{hh}



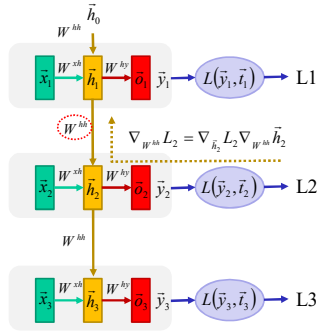
Réseau récurrent: gradient pour W^{hh}



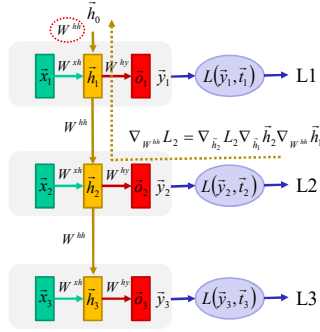
Réseau récurrent: gradient pour W^{hh}



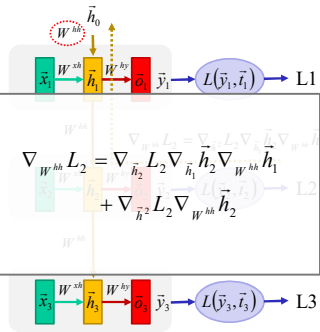
Réseau récurrent: gradient pour W^{hh}



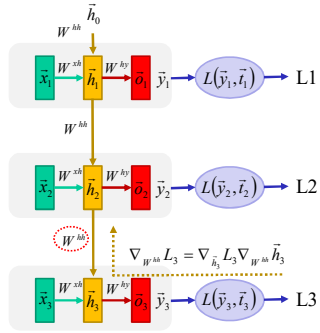
Réseau récurrent: gradient pour W^{hh}



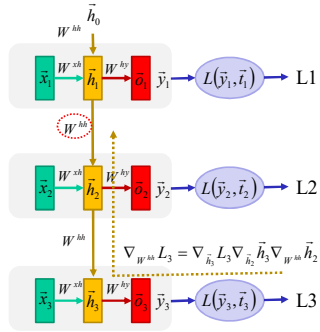
Réseau récurrent: gradient pour W^{hh}



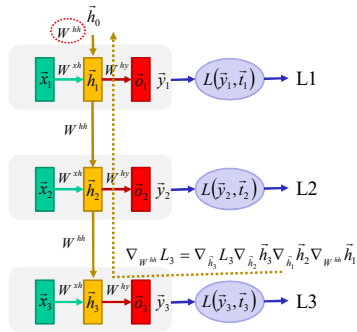
Réseau récurrent: gradient pour W^{hh}



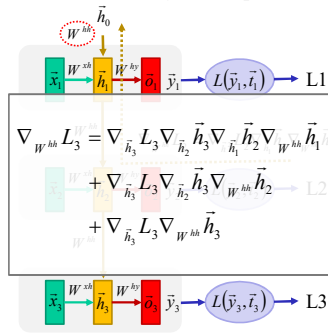
Réseau récurrent: gradient pour W^{hh}



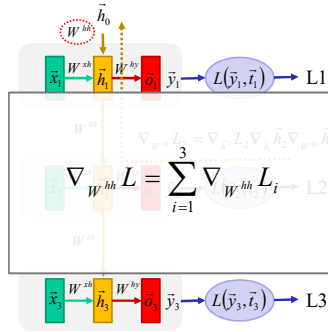
Réseau récurrent: gradient pour W^{hh}



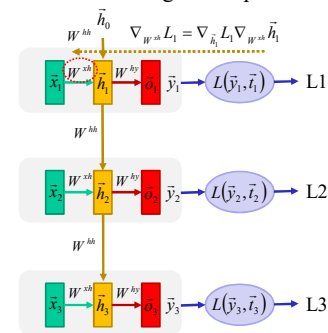
Réseau récurrent: gradient pour W^{hh}

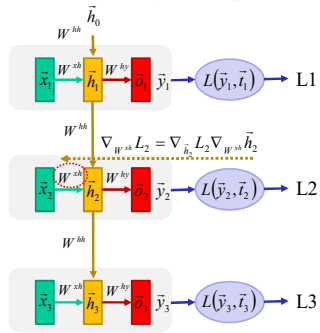


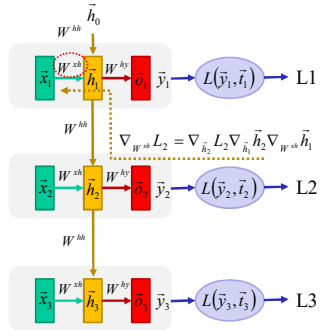
Réseau récurrent: gradient pour W^{hh}

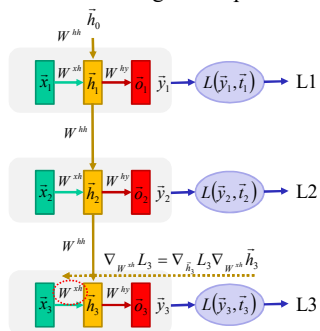


Réseau récurrent: gradient pour W^{xh}

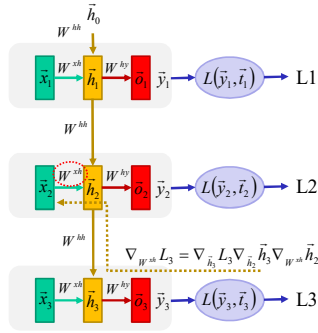


Réseau récurrent: gradient pour W^{xh} 

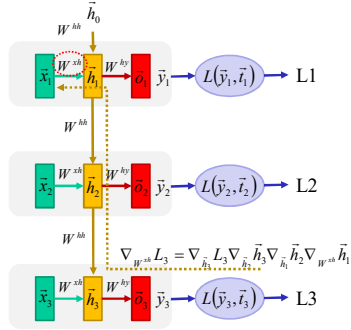
Réseau récurrent: gradient pour W^{xh} 

Réseau récurrent: gradient pour W^{xh} 

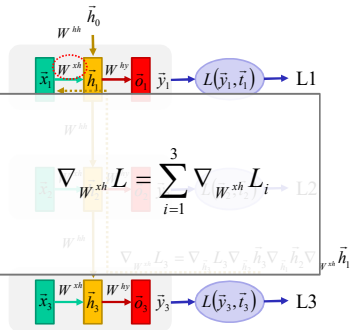
Réseau récurrent: gradient pour W^{xh}



Réseau récurrent: gradient pour W^{xh}



Réseau récurrent: gradient pour W^{xh}



Réseau récurrent: calcul du gradient

Moins difficile qu'il n'y paraît.

```

35 # backward pass: compute gradients going backwards
36 dwh, dwh, dwh = np.zeros_like(wh), np.zeros_like(wh), np.zeros_like(wh)
37 dth, dth = np.zeros_like(th), np.zeros_like(th)
38 dnext = np.zeros_like(x0)
39 for t in reversed(range(len(inputs))):
40     dy = np.copy(y[t])
41     dy[targets[t]] = 1 # backprop into y, see http://cs231n.github.io/neural-networks-case-study/grad-if-confused-here
42     dwh = np.dot(dy, wh[t])
43     dth = dy
44     dh = np.dot(dwh, T) + dnext # backprop into h
45     dthraw = (1 - h[t]) * h[t] * dh # backprop through tanh nonlinearity
46     dth = np.dot(dthraw, x[t-1:T])
47     dwh = np.dot(dthraw, h[t-1:T])
48     dnext = np.dot(wh, dthraw)
49     for dthraw in [dth, dwh, dth, dth, dth]:
50         np.clip(dthraw, -5, 5, out=dthraw) # clip to mitigate exploding gradients
51     return loss, dwh, dwh, dwh, dth, dth, h[inputs-1:]

```

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













Les réseaux récurrents ont un
inconvenient majeur:

difficile à établir des
relations à longue distance

130

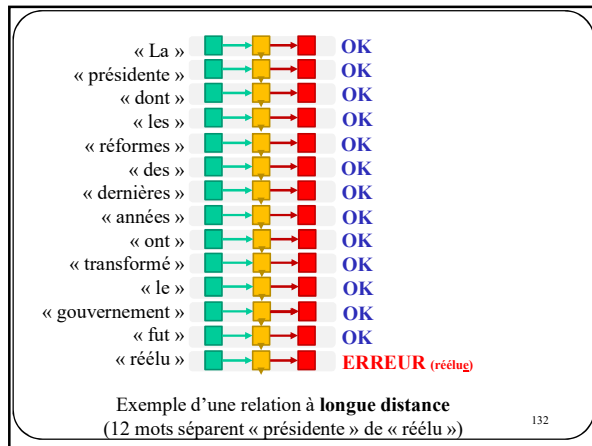
Exemples: analyse grammaticale

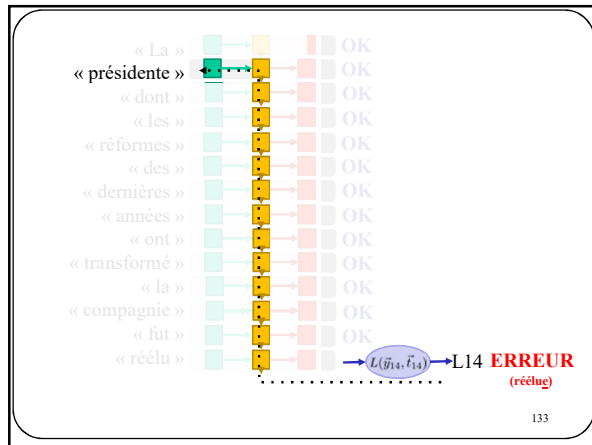
Entraîner un réseau à détecter des erreurs grammaticales

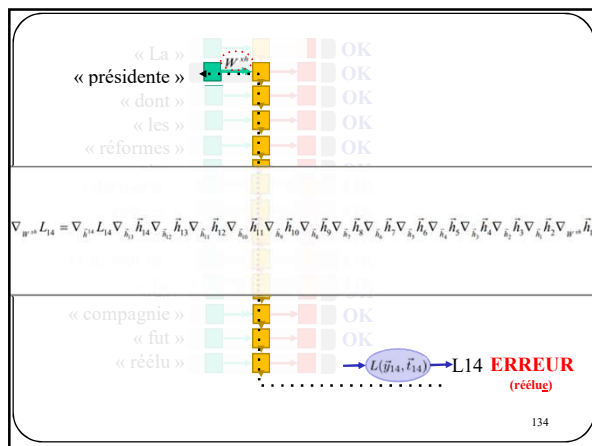
« La »    OK
« présidente »    OK
« fut »    OK
« réélu »    **ERREUR (réélu)**

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

131







« La »

« présidente »

« dont »

« la »

« transformée »

« la »

« compagnie »

« fut »

« réélu »

$\nabla_{W^{sh}} L_{14} = \nabla_{h^{14}} L_{14} \left[\prod_{j=14}^2 \nabla_{h_{j-1}} h_j \right] \nabla_{W^{sh}} h_1$

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

$L(\tilde{y}_{14}, t_{14}) \rightarrow L14$ **ERREUR (réélu)**

135

« La »

« présidente »

« dont »

« la »

« transformée »

« la »

« compagnie »

« fut »

« réélu »

$\nabla_{W^{sh}} L_{14} = \nabla_{h^{14}} L_{14} \left[\prod_{j=14}^2 \nabla_{h_{j-1}} h_j \right] \nabla_{W^{sh}} h_1$

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

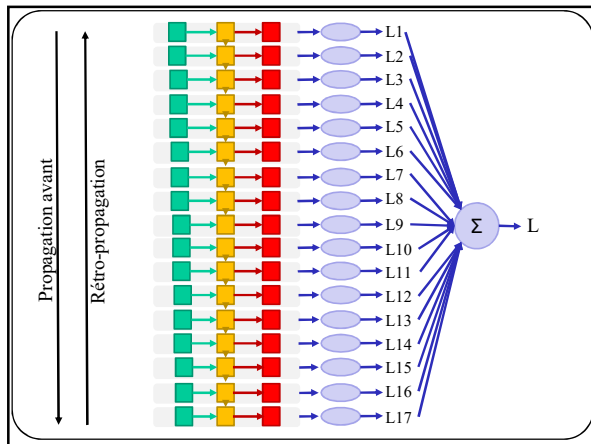
$L(\tilde{y}_{14}, t_{14}) \rightarrow L14$ **ERREUR (réélu)**

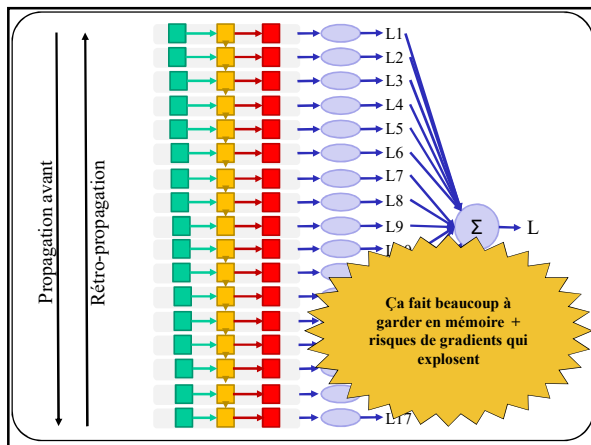
136

Problème connexe

Gestion de la mémoire

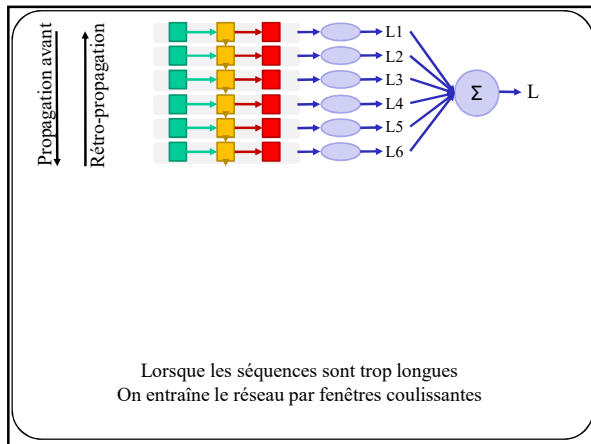
137

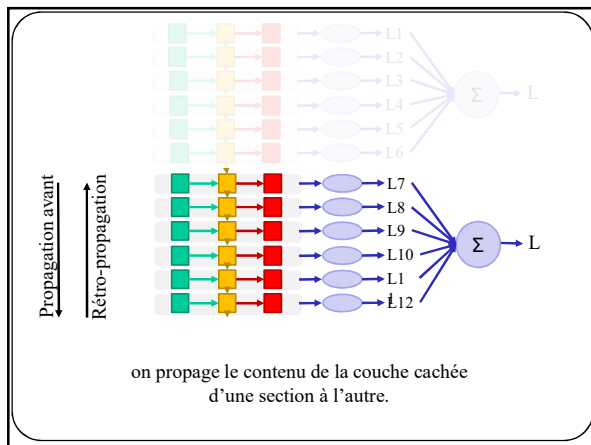


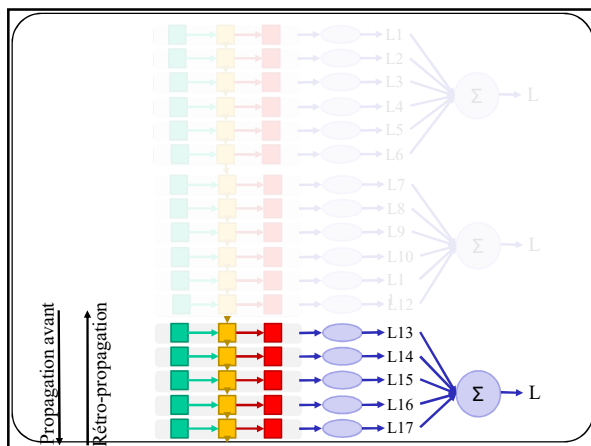


Solution pour la gestion de la
mémoire
Fenêtres coulissantes

140





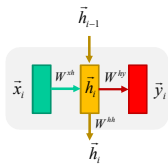


Solution à la disparition du gradient:

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

144

Illustration + formulation d'un RNN



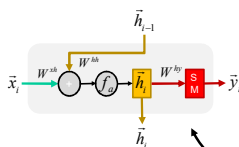
$$\tilde{h}_t = f_a(W^{xh}\tilde{x}_t + W^{hh}\tilde{h}_{t-1})$$

$$\hat{y}_t = W^{hy}\tilde{h}_t$$

$$\tilde{y}_t = \text{SMAX}(\hat{y}_t)$$

145

Autre illustration du même RNN



Identique

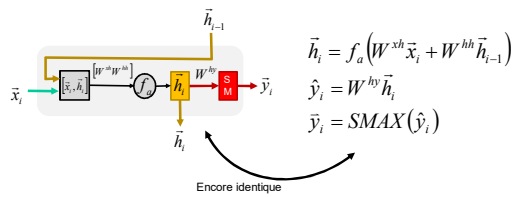
$$\tilde{h}_t = f_a(W^{xh}\tilde{x}_t + W^{hh}\tilde{h}_{t-1})$$

$$\hat{y}_t = W^{hy}\tilde{h}_t$$

$$\tilde{y}_t = \text{SMAX}(\hat{y}_t)$$

146

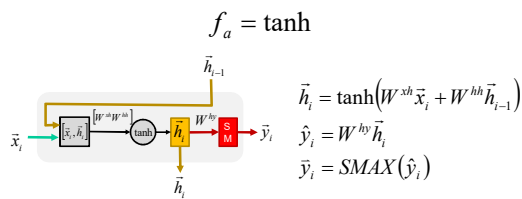
Autre illustration du même RNN



147

GRU (Gated Recurrent Unit)

Modif 1

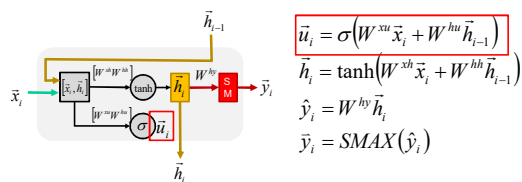


148

GRU (Gated Recurrent Unit)

Modif 2

Update gate

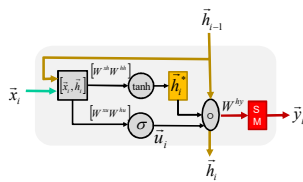
 $\sigma = \text{sigmoid}$ 

149

GRU (Gated Recurrent Unit)

Modif 2

Update gate



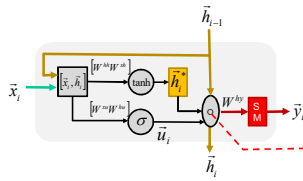
$$\begin{aligned}\bar{u}_i &= \sigma(W^{xu} \bar{x}_i + W^{hu} \bar{h}_{i-1}) \\ \bar{h}_i^* &= \tanh(W^{xh} \bar{x}_i + W^{hh} \bar{h}_{i-1}) \\ \bar{h}_i &= \bar{u}_i \circ \bar{h}_i^* + (1 - \bar{u}_i) \circ \bar{h}_{i-1} \\ \hat{y}_i &= W^{hy} \bar{h}_i \\ \bar{y}_i &= \text{SMAX}(\hat{y}_i)\end{aligned}$$

150

GRU (Gated Recurrent Unit)

Modif 2

« element-wise product »
ou
« Hadamard product »



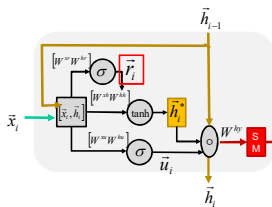
$$(a, b, c) \circ (x, y, z) = (ax, by, cz)$$

151

GRU (Gated Recurrent Unit)

Modif 3

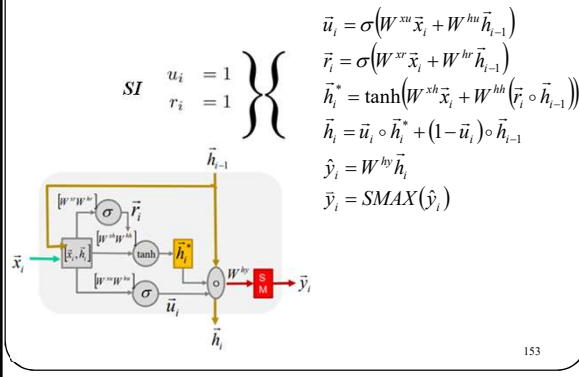
Reset gate



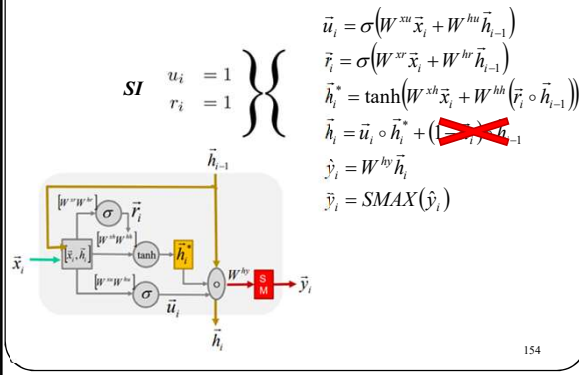
$$\begin{aligned}\bar{u}_i &= \sigma(W^{xu} \bar{x}_i + W^{hu} \bar{h}_{i-1}) \\ \bar{r}_i &= \sigma(W^{xr} \bar{x}_i + W^{hr} \bar{h}_{i-1}) \\ \bar{h}_i^* &= \tanh(W^{xh} \bar{x}_i + W^{hh} (\bar{r}_i \circ \bar{h}_{i-1})) \\ \bar{h}_i &= \bar{u}_i \circ \bar{h}_i^* + (1 - \bar{u}_i) \circ \bar{h}_{i-1}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

152

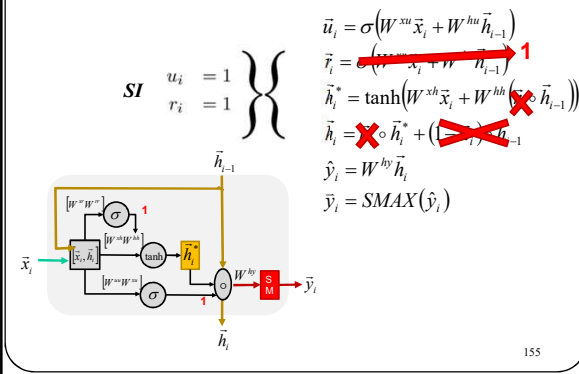
Comprendre les gates



Comprendre les gates

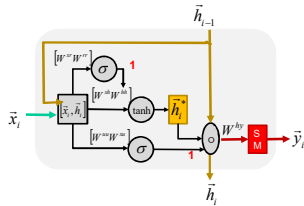


Comprendre les gates



Comprendre les gates

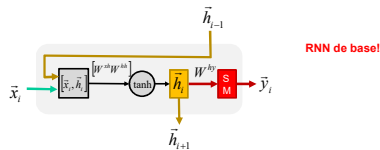
$$SI \quad \left. \begin{array}{l} \tilde{u}_i = 1 \\ \tilde{r}_i = 1 \end{array} \right\} \left\{ \begin{array}{l} \tilde{h}_i^* = \tanh(W^{xh}\tilde{x}_i + W^{hh}\tilde{h}_{i-1}) \\ \tilde{h}_i = \tilde{h}_i^* \\ \hat{y}_i = W^{hy}\tilde{h}_i \\ \tilde{y}_i = SMAX(\hat{y}_i) \end{array} \right.$$



156

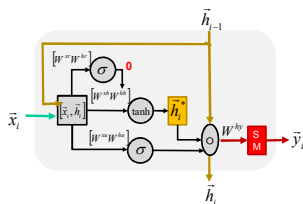
Comprendre les gates

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



157

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



158

$$SI \quad \left. \begin{array}{l} \tilde{u}_i \neq 1 \\ \tilde{r}_i = 0 \end{array} \right\}$$

$$\begin{aligned} \tilde{u}_i &= \sigma(W^{uu} \tilde{x}_i + W^{hu} \tilde{h}_{i-1}) \\ \tilde{h}_i^* &= \tanh(W^{sh} \tilde{x}_i) \\ \tilde{h}_i &= \tilde{u}_i \circ \tilde{h}_i^* + (1 - \tilde{u}_i) \circ \tilde{h}_{i-1} \\ \hat{y}_i &= W^{hy} \tilde{h}_i \\ \tilde{y}_i &= SMAX(\hat{y}_i) \end{aligned}$$

\tilde{h}_{i-1}^* ne dépend que de \tilde{x}_i et non de \tilde{h}_{i-1}

159

$$SI \quad \left. \begin{array}{l} \tilde{u}_i = 0 \\ \tilde{r}_i = 0 \end{array} \right\}$$

~~$$\begin{aligned} \tilde{u}_i &= \sigma(W^{uu} \tilde{x}_i + W^{hu} \tilde{h}_{i-1}) \\ \tilde{h}_i^* &= \tanh(W^{sh} \tilde{x}_i + W^{hs} \tilde{h}_{i-1}) \\ \tilde{h}_i &= \tilde{u}_i \circ \tilde{h}_i^* + (1 - \tilde{u}_i) \circ \tilde{h}_{i-1} \\ \hat{y}_i &= W^{hy} \tilde{h}_i \\ \tilde{y}_i &= SMAX(\hat{y}_i) \end{aligned}$$~~

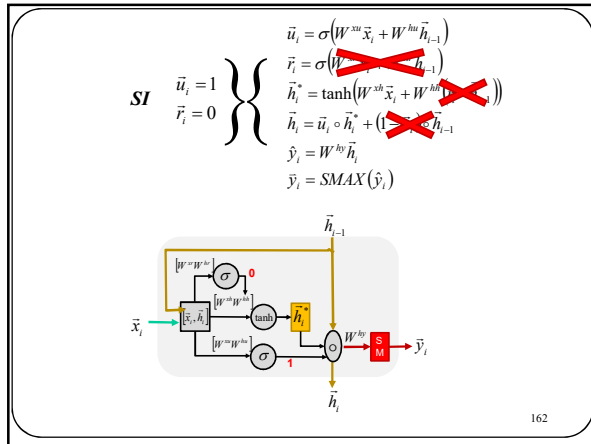
160

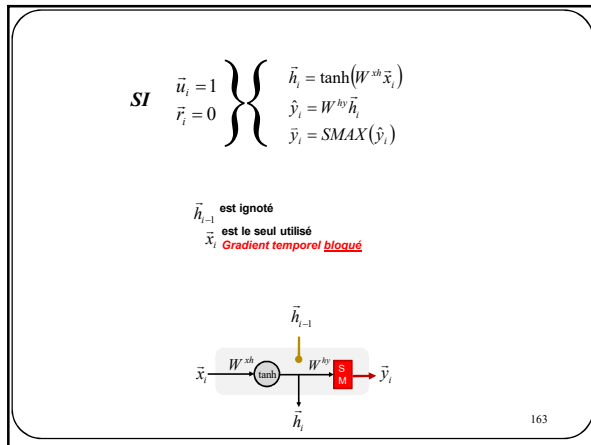
$$SI \quad \left. \begin{array}{l} \tilde{u}_i = 0 \\ \tilde{r}_i = 0 \end{array} \right\}$$

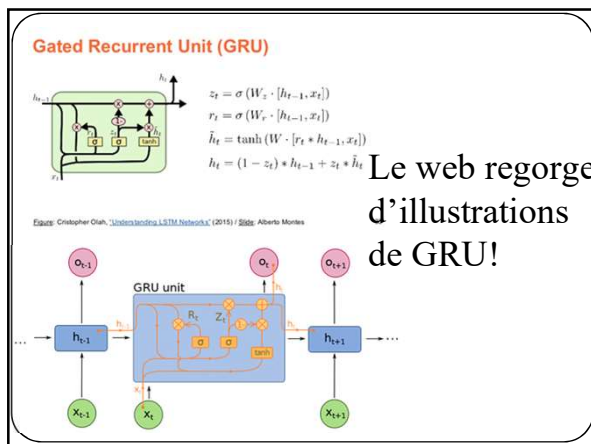
$$\begin{aligned} \tilde{h}_i &= \tilde{h}_{i-1} \\ \hat{y}_i &= W^{cy} \tilde{h}_i \\ \tilde{y}_i &= SMAX(\hat{y}_i) \end{aligned}$$

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

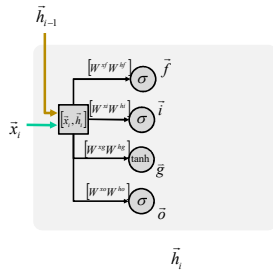
161







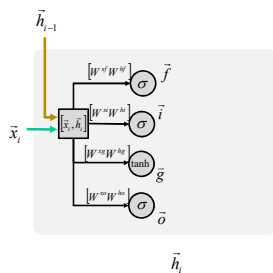
LSTM (Long Short Term Memory)



\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

LSTM (Long Short Term Memory)

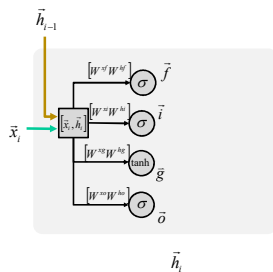


\tilde{f} : forget gate
 \tilde{i} : input gate
 \tilde{g} : gate gate

Modèle récurrent
 le plus utilisé.
 À bien
 comprendre !

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

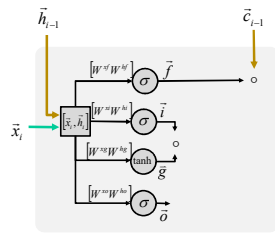
LSTM (Long Short Term Memory)



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

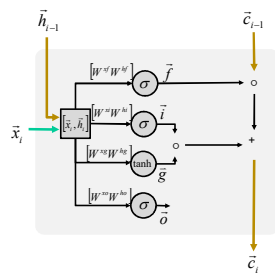
167

LSTM (Long Short Term Memory)



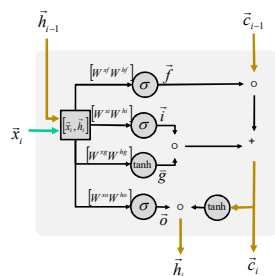
168

LSTM (Long Short Term Memory)

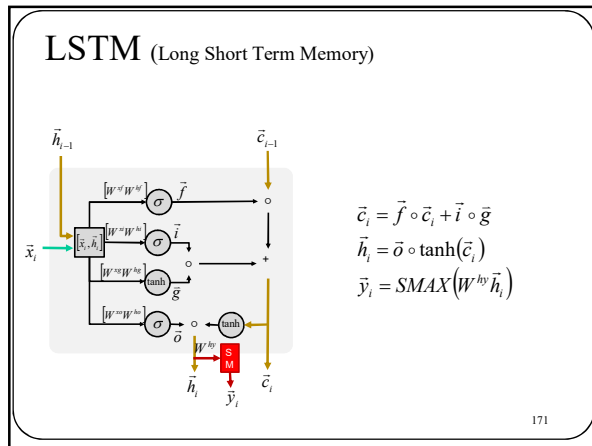


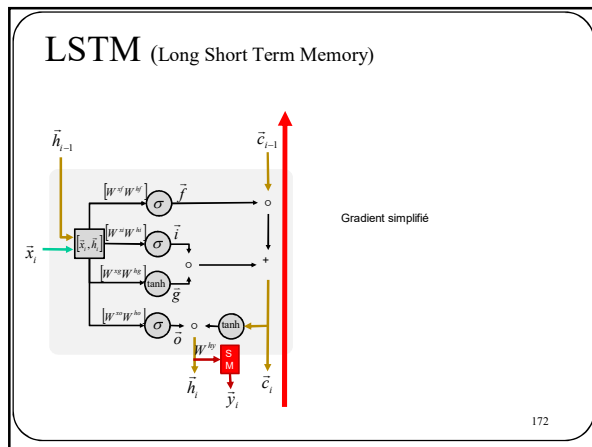
169

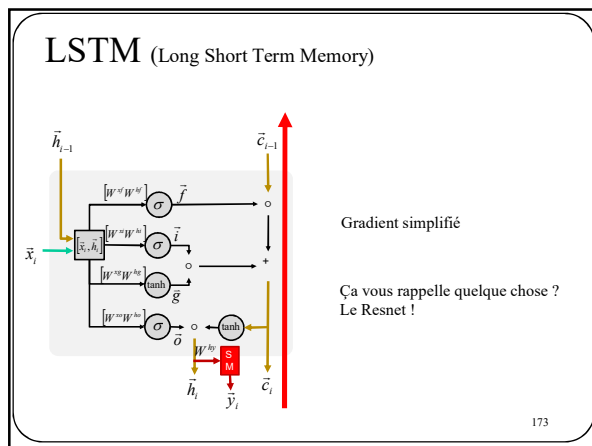
LSTM (Long Short Term Memory)

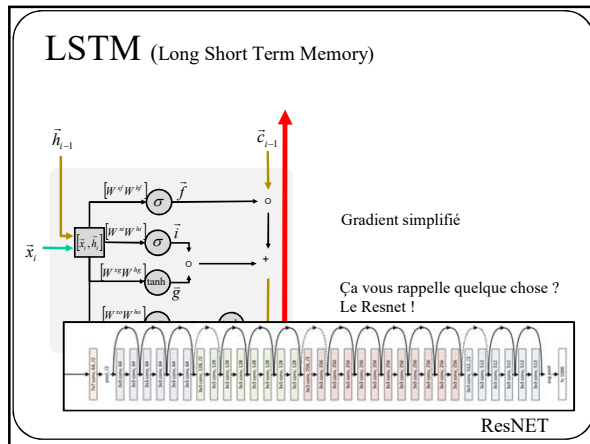


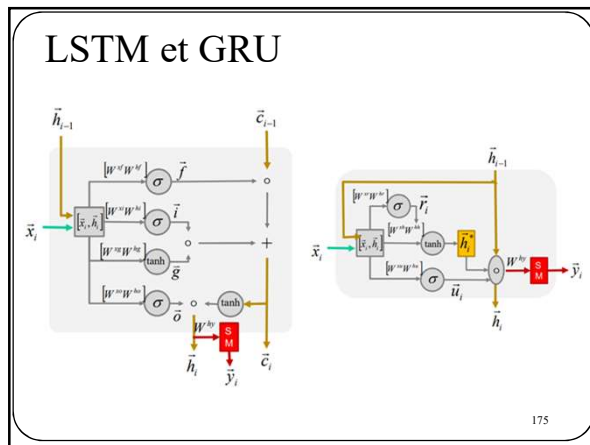
170

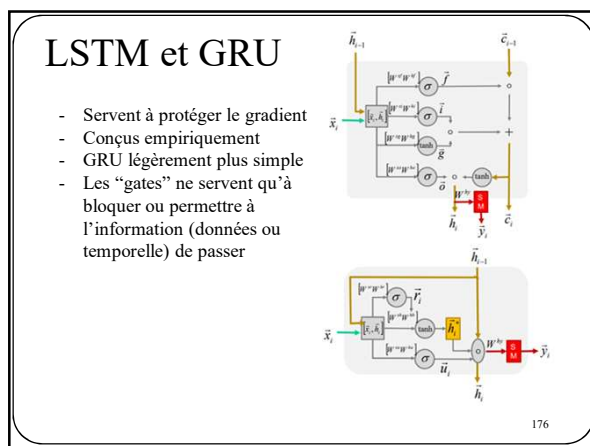




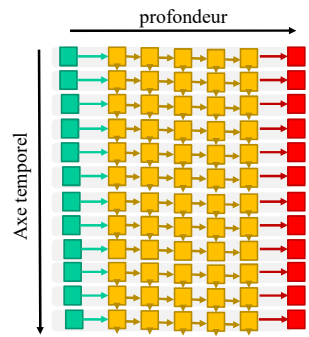








RNN multi-couches

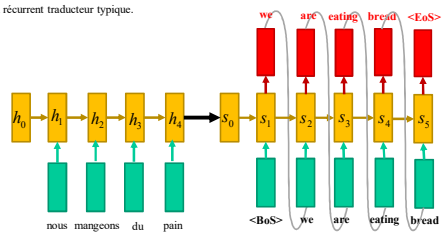


Modèles d'attention

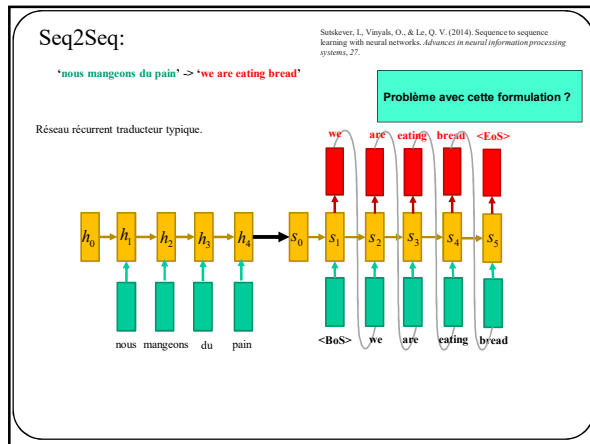
Seq2Seq:

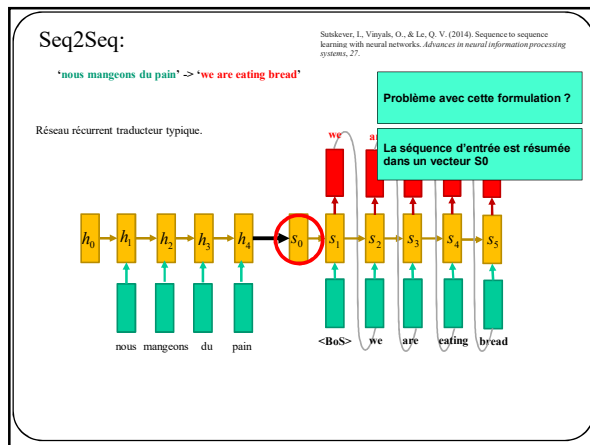
'nous mangeons du pain' → 'we are eating bread'

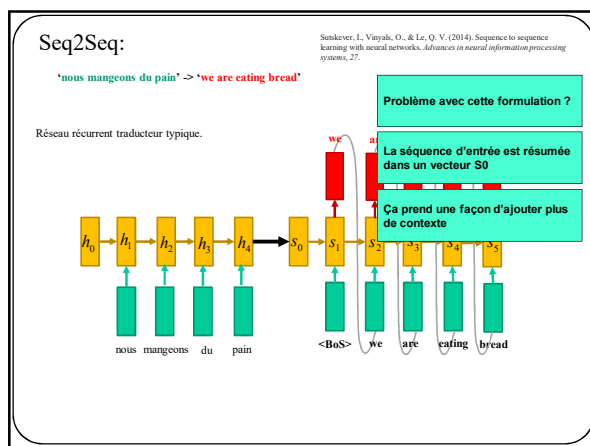
Réseau récurrent traducteur typique.



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







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:1409.0327.

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

nous mangeons du pain

$s_0 = \text{init value}$

Seq2Seq avec attention
(version simplifiée)

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

$$e_{li} = W^{ev} \vec{s}_0 + W^{eh} \tilde{h}_i \in R$$

nous mangeons du pain

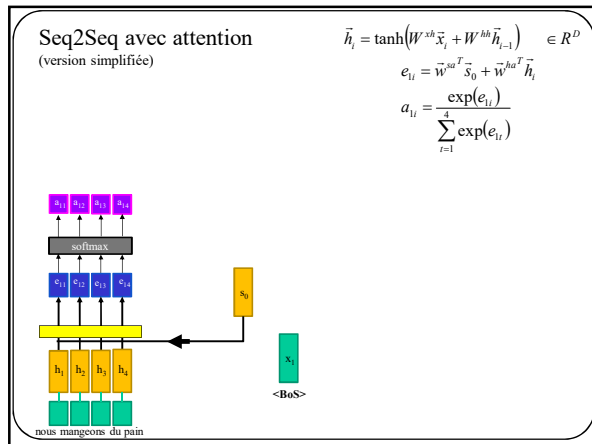
Seq2Seq avec attention
(version simplifiée)

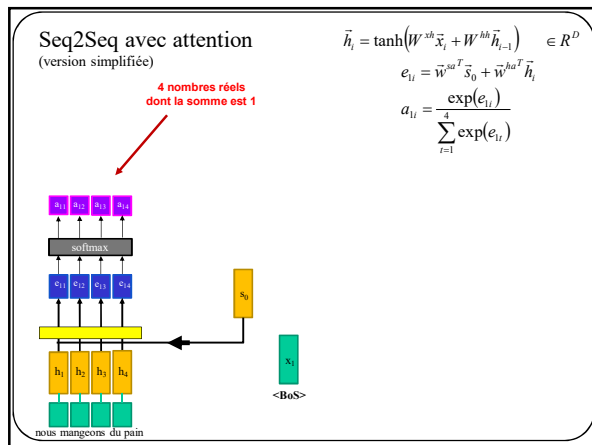
$$\tilde{h}_i = \tanh(W^{vh} \vec{x}_i + W^{hh} \tilde{h}_{i-1}) \in R^D$$

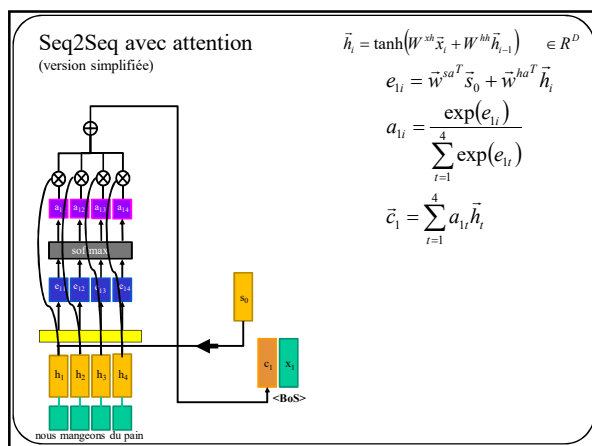
$$e_{li} = W^{ev} \vec{s}_0 + W^{eh} \tilde{h}_i \in R$$

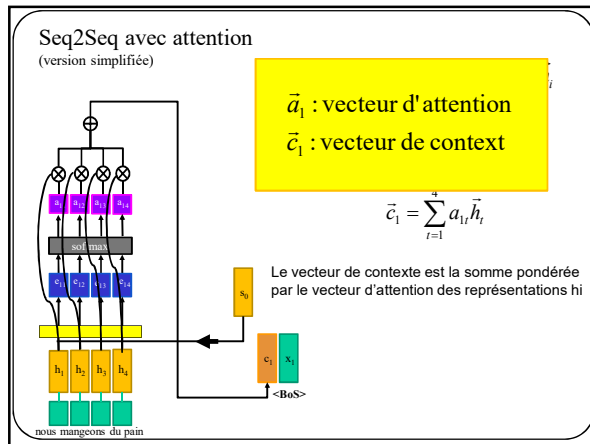
4 nombres réels

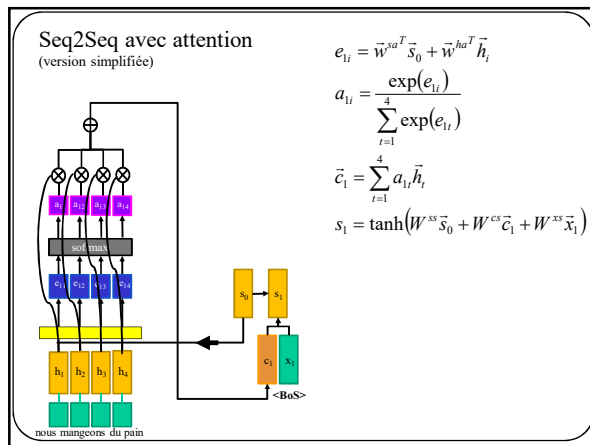
nous mangeons du pain

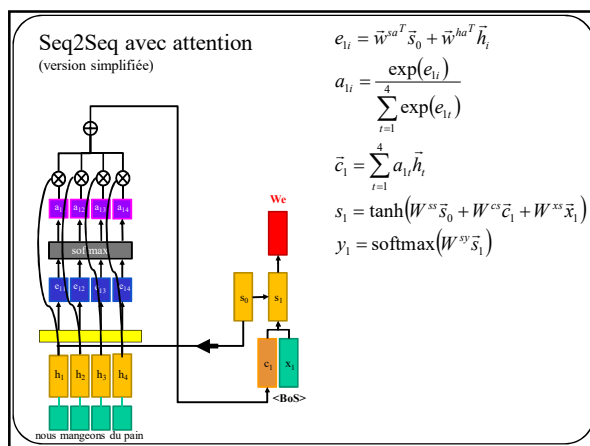


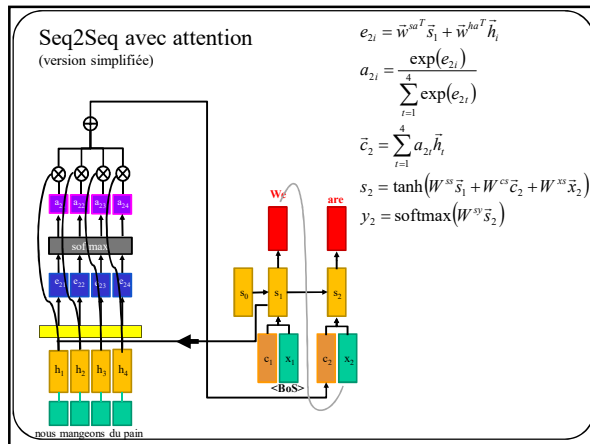


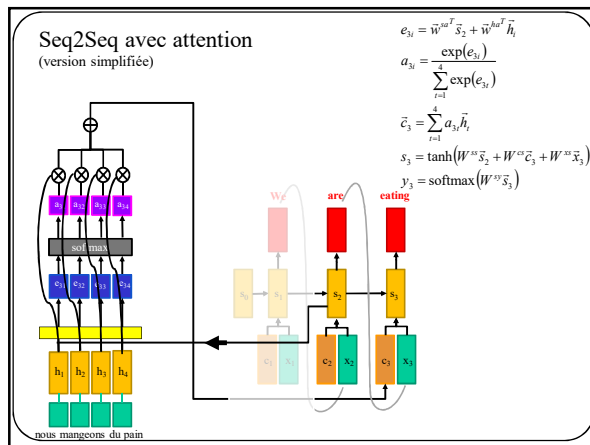


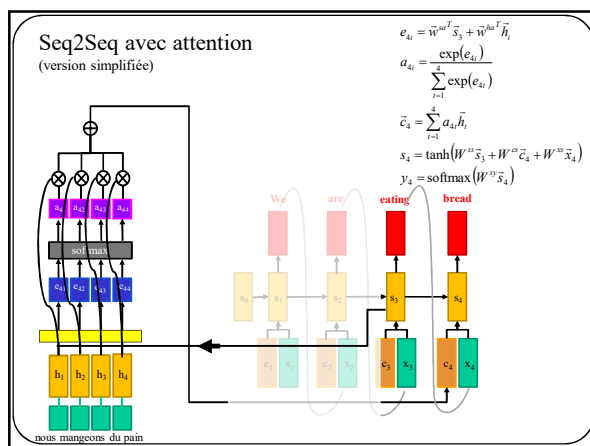




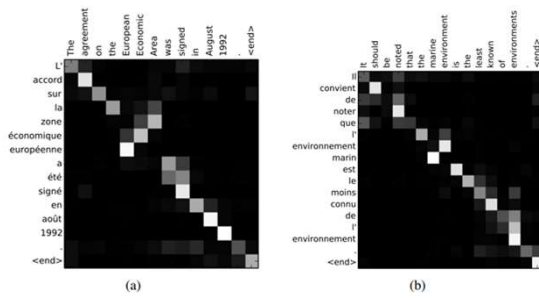






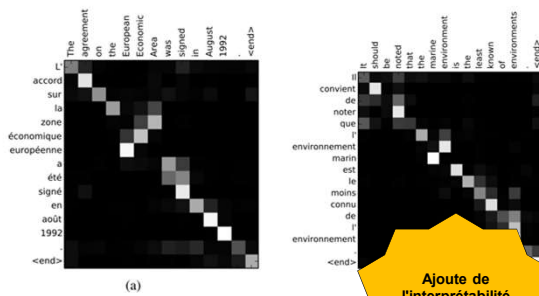


Un vecteur d'attention pour chaque mot



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

Seq2Seq avec attention

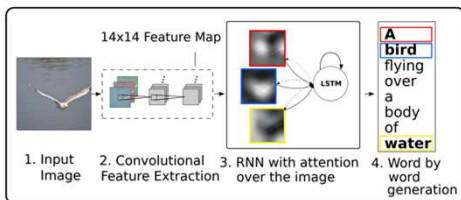


Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473.

Ajoute de l'interprétabilité au modèle !

Modèles d'attention pour la description d'images

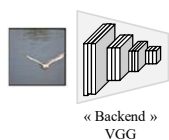
Réseau récurrent pour du captioning capable de « concentrer son attention » sur les zones de l'image associés aux mots.



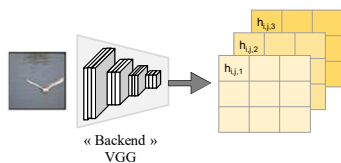
Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., ... & Bengio, Y. (2015, June). Show, attend and tell: Neural image caption generation with visual attention. In International conference on machine learning (pp. 2048-2057). PMLR.

197

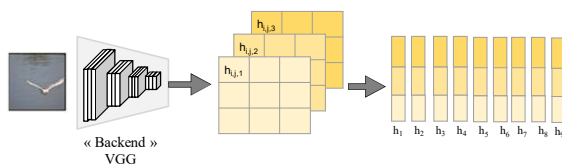
Modèles d'attention pour la description d'images (version simplifiée)

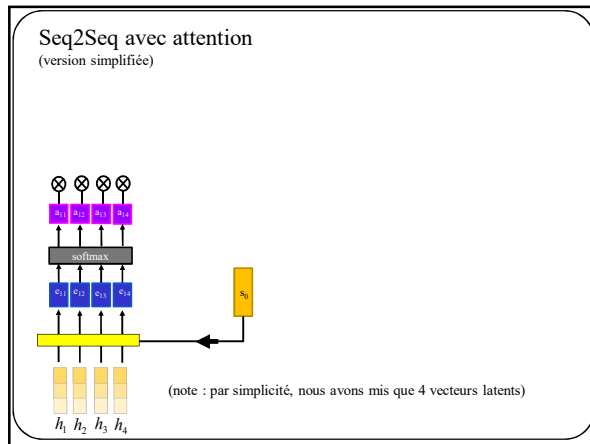


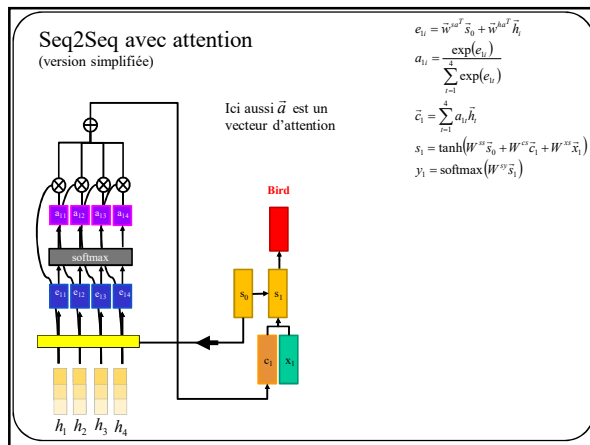
198

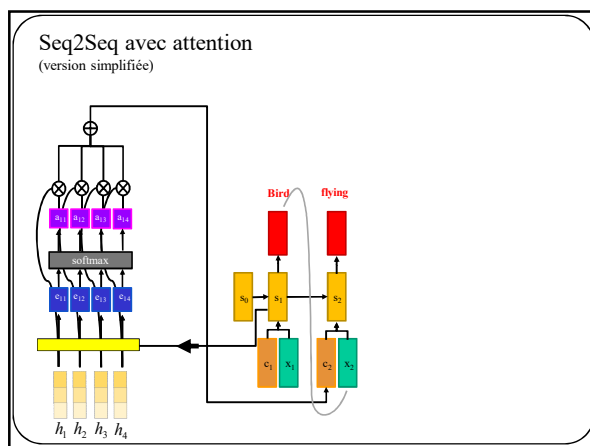


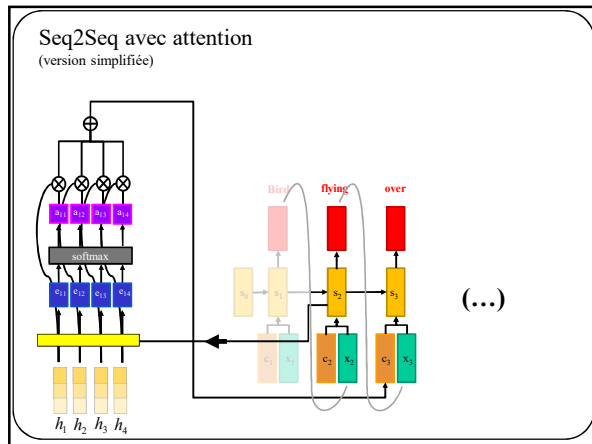
Les positions ij correspondent à $h_{i,2,3}$...
Les filtres correspondent aux composantes de h

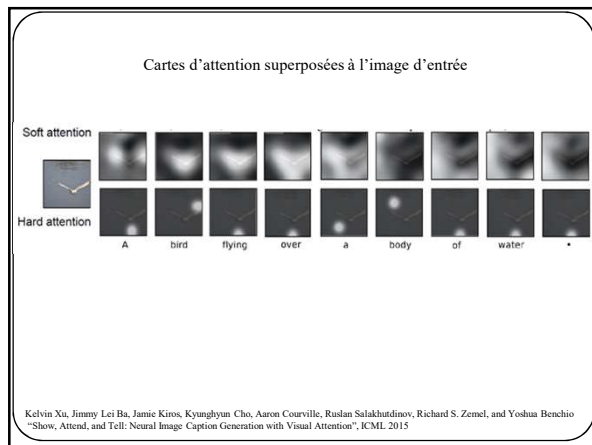


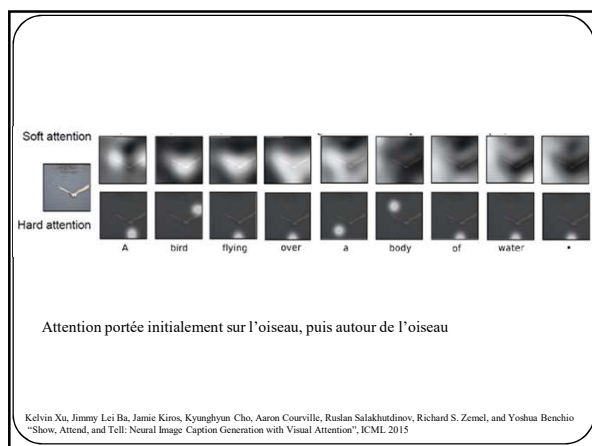














L'auto-attention
(*self attention*)

208

Revenons à la base : multiplication matricielle

Considérons les 4 matrices suivantes

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

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

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

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

209

Revenons à la base : multiplication matricielle

Leur multiplication donne:

$$\begin{aligned}
 X &= \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix} \in \mathbb{R}^{3 \times 4} \\
 W^q X &= Q = \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} & Q_{14} \\ Q_{21} & Q_{22} & Q_{23} & Q_{24} \\ Q_{31} & Q_{32} & Q_{33} & Q_{34} \end{pmatrix} \in \mathbb{R}^{3 \times 4} \\
 W^q &= \begin{pmatrix} W^q_{11} & W^q_{12} & W^q_{13} \\ W^q_{21} & W^q_{22} & W^q_{23} \\ W^q_{31} & W^q_{32} & W^q_{33} \end{pmatrix} \in \mathbb{R}^{3 \times 3} \\
 W^k X &= K = \begin{pmatrix} K^k_{11} & K^k_{12} & K^k_{13} & K^k_{14} \\ K^k_{21} & K^k_{22} & K^k_{23} & K^k_{24} \\ K^k_{31} & K^k_{32} & K^k_{33} & K^k_{34} \end{pmatrix} \in \mathbb{R}^{3 \times 4} \\
 W^k &= \begin{pmatrix} W^k_{11} & W^k_{12} & W^k_{13} \\ W^k_{21} & W^k_{22} & W^k_{23} \\ W^k_{31} & W^k_{32} & W^k_{33} \end{pmatrix} \in \mathbb{R}^{3 \times 3} \\
 W^V X &= V = \begin{pmatrix} V^V_{11} & V^V_{12} & V^V_{13} & V^V_{14} \\ V^V_{21} & V^V_{22} & V^V_{23} & V^V_{24} \end{pmatrix} \in \mathbb{R}^{2 \times 4} \\
 W^V &= \begin{pmatrix} W^V_{11} & W^V_{12} & W^V_{13} \\ W^V_{21} & W^V_{22} & W^V_{23} \end{pmatrix} \in \mathbb{R}^{2 \times 3}
 \end{aligned}$$

210

Auto attention

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

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

↑ Nous ↑ mangeons ↑ du ↑ pain

Dans cet exemple, 4 mots en entrée donc 4 colonnes dans X
 Les vecteurs 3D peuvent être obtenus par **Word2Vec**

211

Auto attention

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

↑ Nous ↑ mangeons ↑ du ↑ pain

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

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

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

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

212

Auto attention

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

↑ Nous ↑ mangeons ↑ du ↑ pain

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

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

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

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)

213

Auto attention

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

↑ Nous ↑ mangeons ↑ du ↑ pain

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

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

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

Matrices de paramètres appris par rétropropagation
Pour ces 3 matrices, le nombre de ligne (3,3,2) est arbitraire

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 x_i a été transformé en des **vecteurs Q_i, K_i et V_i**

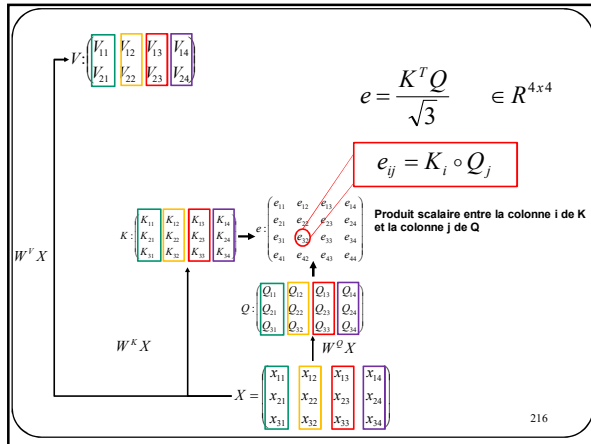
$K: \begin{pmatrix} K_{11} & K_{12} & K_{13} & K_{14} \\ K_{21} & K_{22} & K_{23} & K_{24} \\ K_{31} & K_{32} & K_{33} & K_{34} \end{pmatrix}$

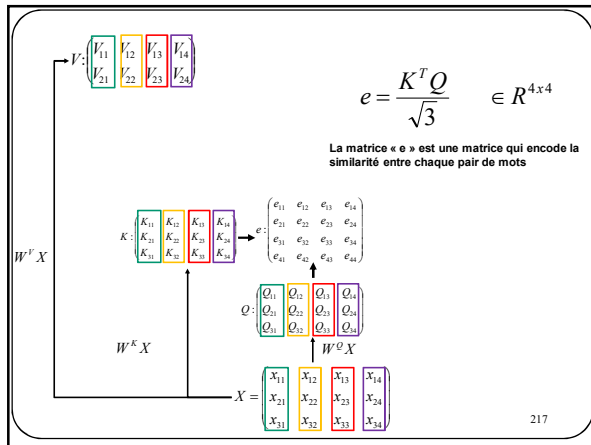
$Q: \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} & Q_{14} \\ Q_{21} & Q_{22} & Q_{23} & Q_{24} \\ Q_{31} & Q_{32} & Q_{33} & Q_{34} \end{pmatrix}$

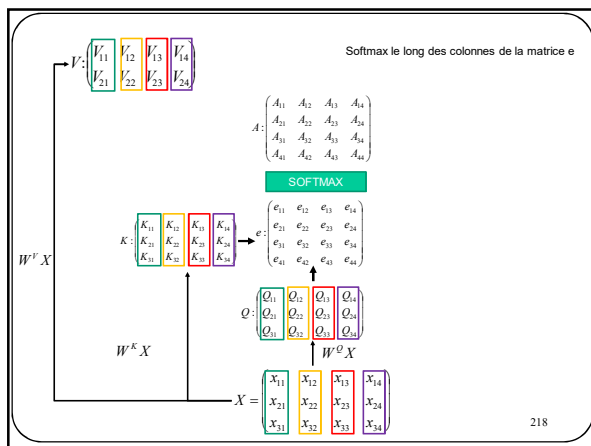
$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{pmatrix}$

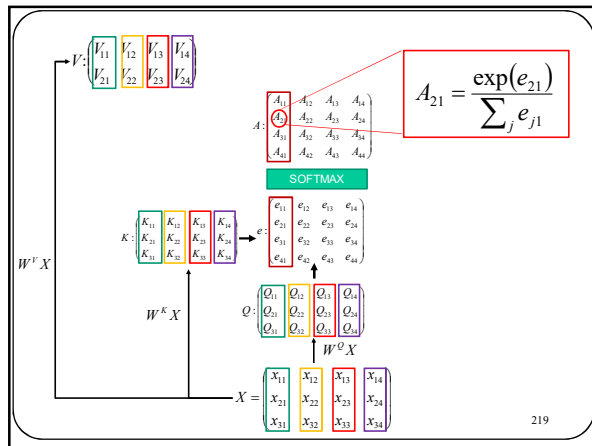
$W^q X$ $W^k X$ $W^v X$

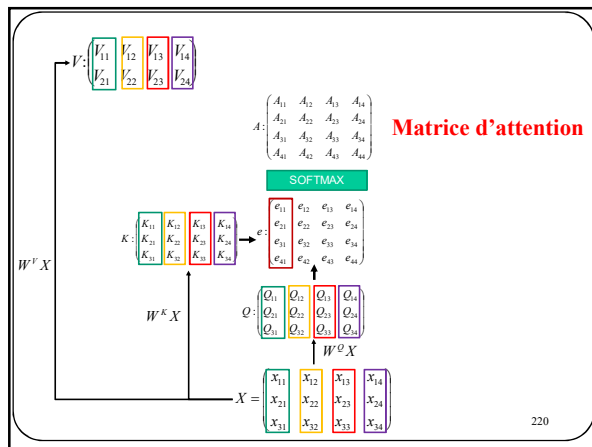
215

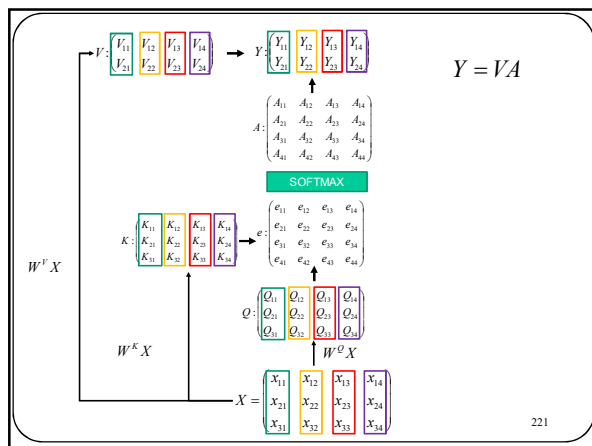


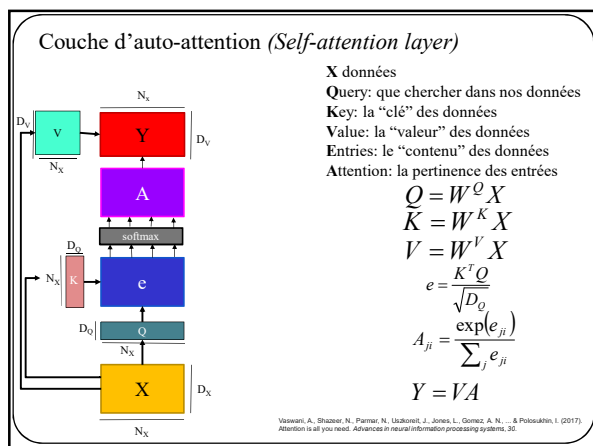


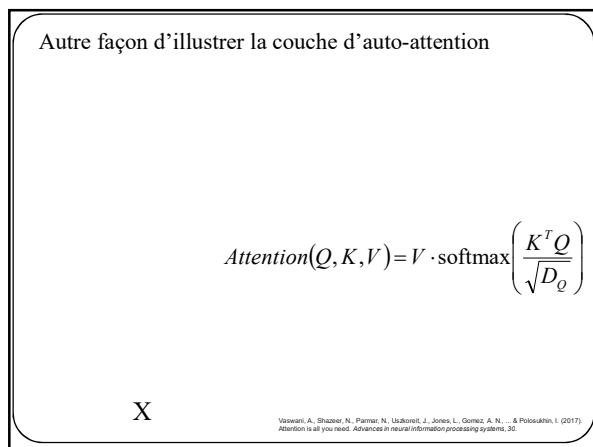


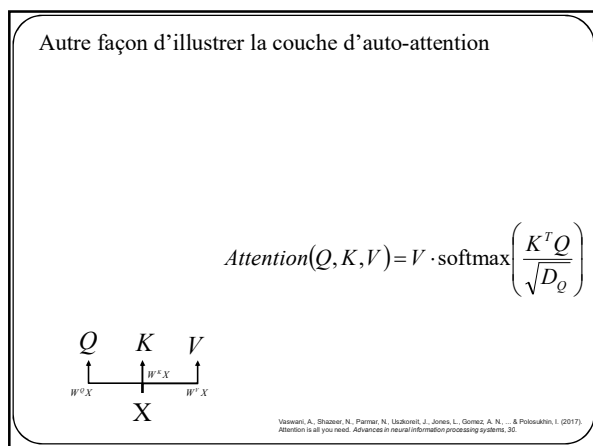




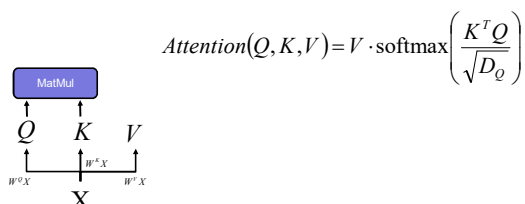




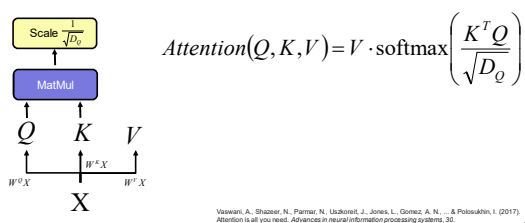




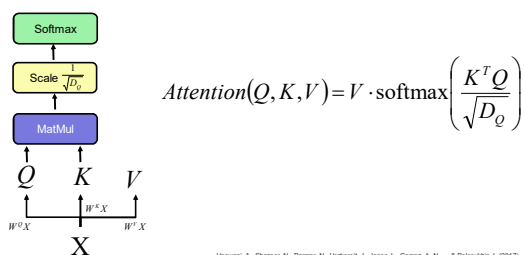
Autre façon d'illustrer la couche d'auto-attention

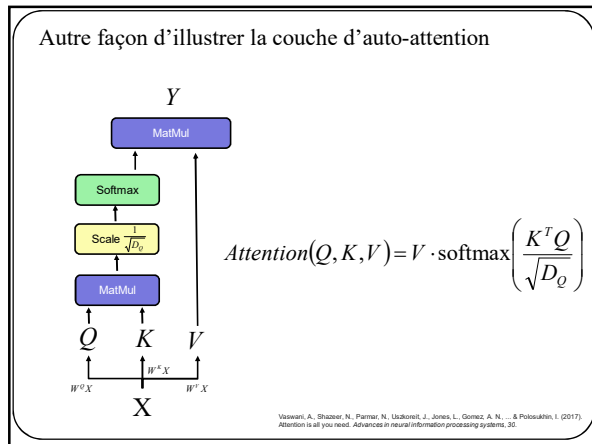


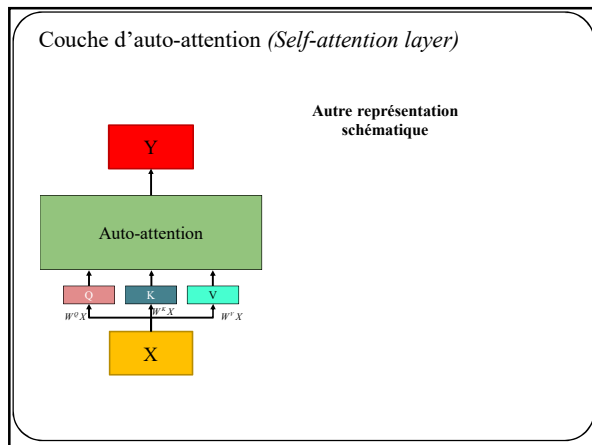
Autre façon d'illustrer la couche d'auto-attention

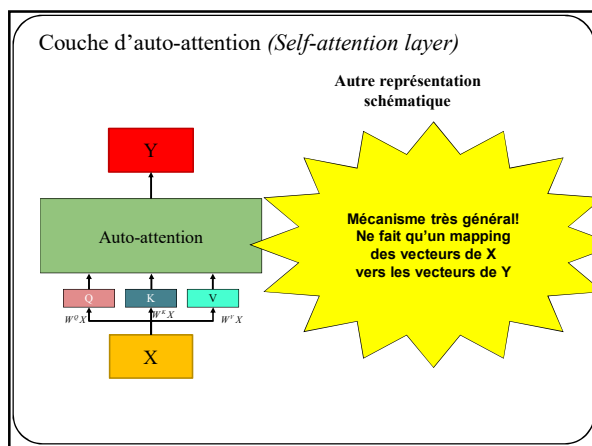


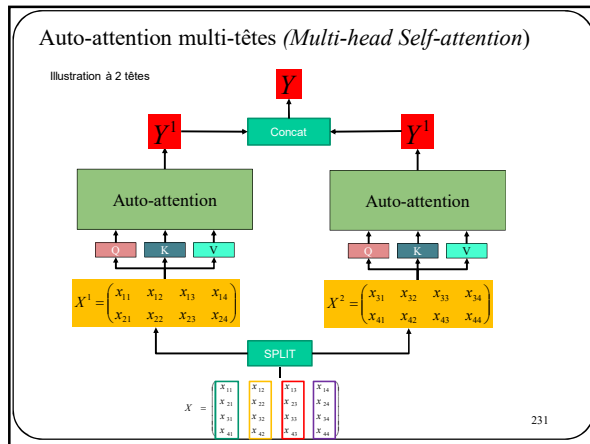
Autre façon d'illustrer la couche d'auto-attention

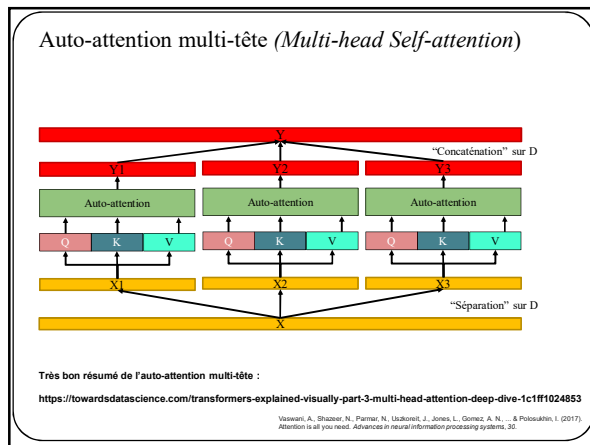


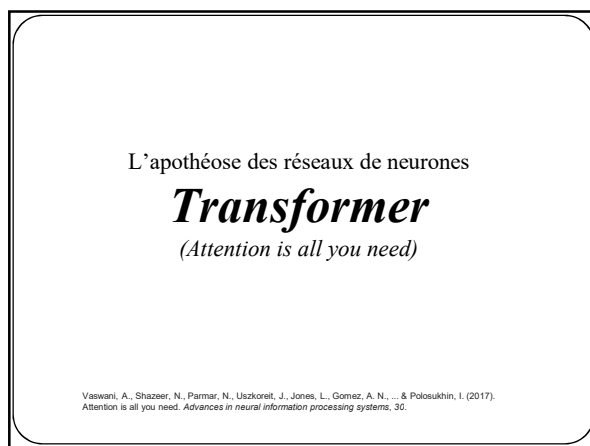













Transformer

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.

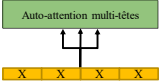
Transformer (Attention is all you need)



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)

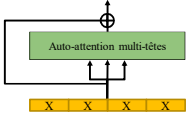
- Auto-attention multi-têtes sur les dimensions de X



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)

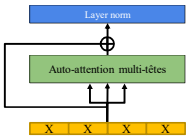
- Auto-attention multi-têtes sur les dimensions de X
- "+ " = connexion résiduelle



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)

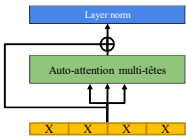
- Auto-attention multi-têtes sur les dimensions de X
- "+ " = connexion résiduelle
- "Layer-norm"



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)

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



Batch Normalization

batch

1	3	4
2	1	2
6	1	5
4	6	1
5	2	3
1	0	1

mean 2 3 3
std 2 2 2

Same for all training examples

Layer Normalization

batch

1	3	4
2	2	2
0	1	6
4	6	1
5	2	3
1	0	1

mean 2 3 3
std 2 2 2

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.

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)

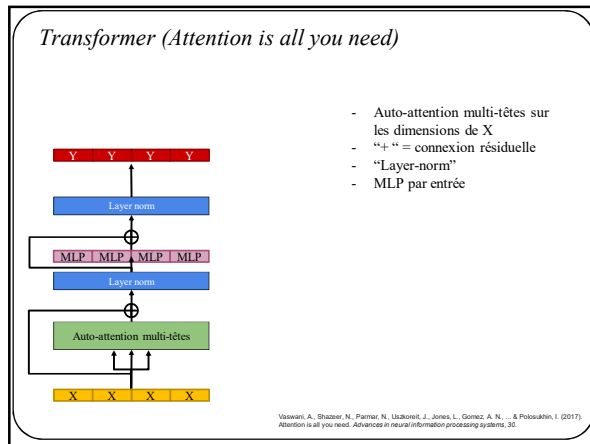
- Intra-attention multi-tête sur les dimensions de X
- "+" = connexion résiduelle
- "Layer-norm"
- MLP par entrée

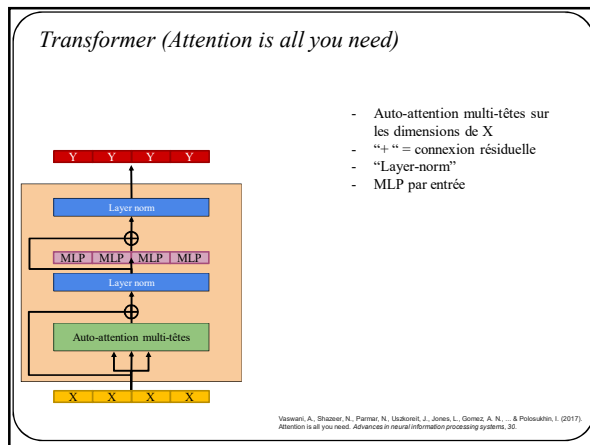
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

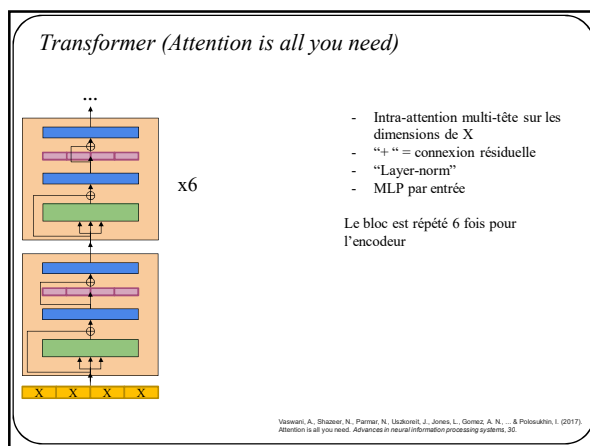
Transformer (Attention is all you need)

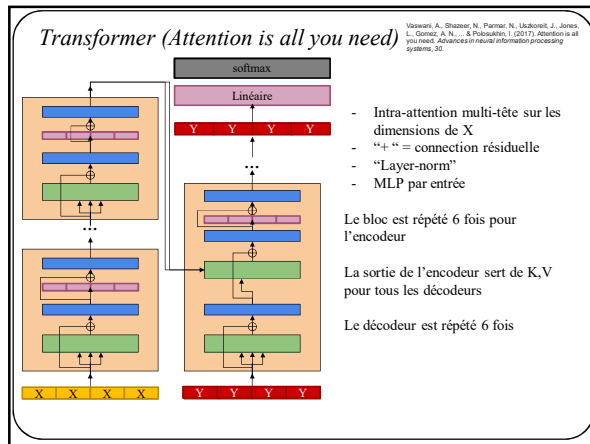
- Auto-attention multi-têtes sur les dimensions de X
- "+" = connexion résiduelle
- "Layer-norm"
- MLP par entrée

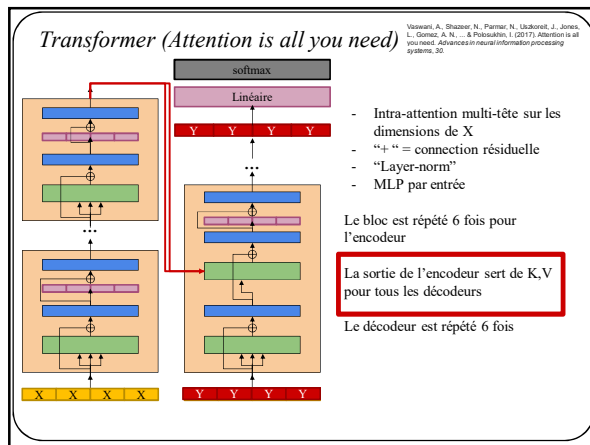
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

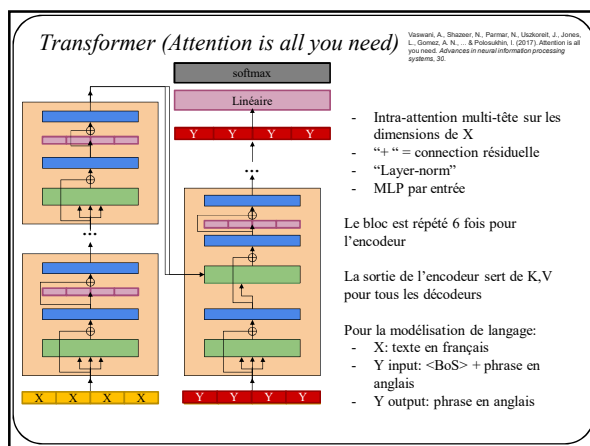


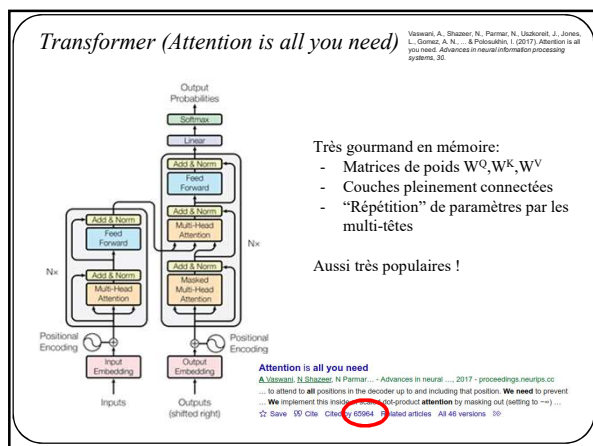












Différentes version de transformers

Pour rappel: Resnet-50: 23M de paramètres

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

crédit: Justin Johnson

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	1.3B	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

Rafferty, A., Naveen, K., Salinas, T., & Subbier, I. (2018). Improving language understanding by generative pre-training.
Rafferty, A., Wu, J., Choi, R., Lavin, D., Amodei, D., & Subbier, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 9.
Brenin, T., Mann, B., Pyle, N., Subbier, M., Kaplan, J. D., Dhurand, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in neural information processing systems*, 33, 1877-1901.

Generative-Pre-Training

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

GPT3 a 96 blocks d'auto-attention!

GPT1-2-3

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

N = 96 pour GPT3

262

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 (necery, 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, 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>

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

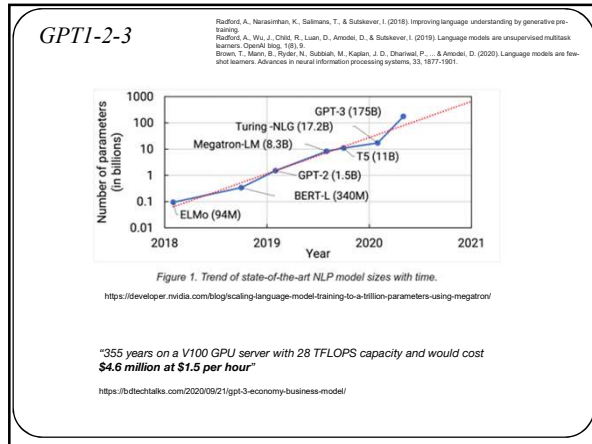
Parameters	Layers	d_{model}
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600

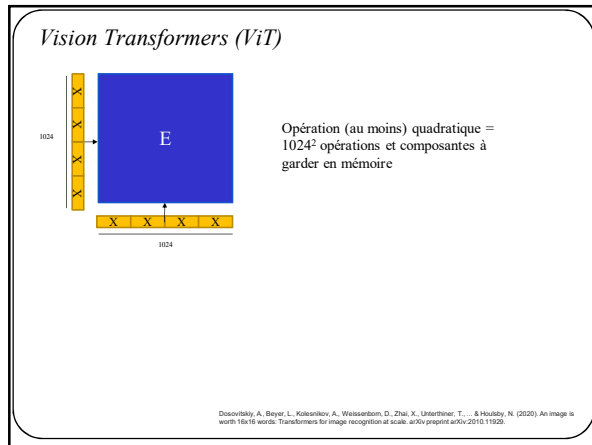
Table 2. Architecture hyperparameters for the 4 model sizes.

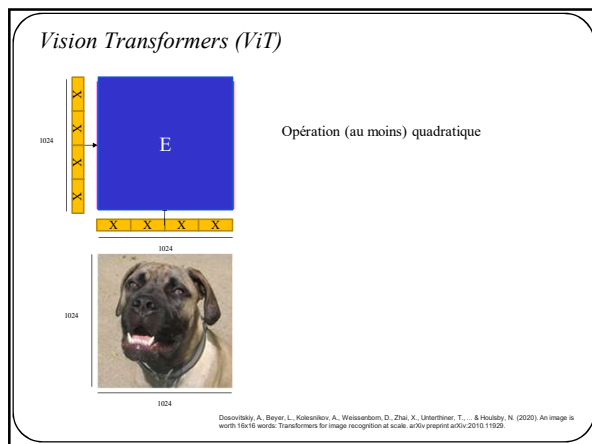
GPT-3

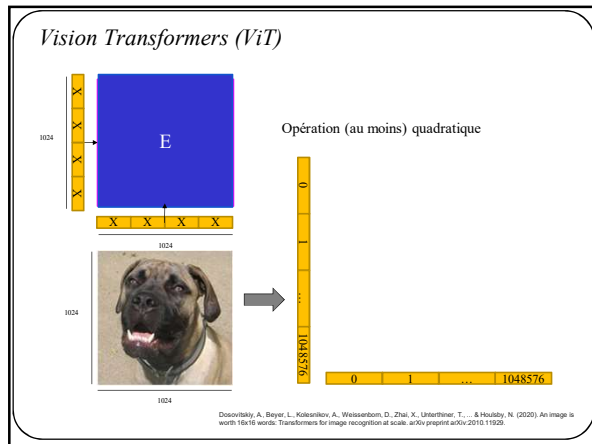
Model Name	n_{tokens}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	34	1156	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	44	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	52	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	52	4096	32	128	2M	1.2×10^{-4}
GPT-3 1.3B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 1.75B or "GPT-3"	175.0B	46	12288	96	128	3.2M	0.6×10^{-4}

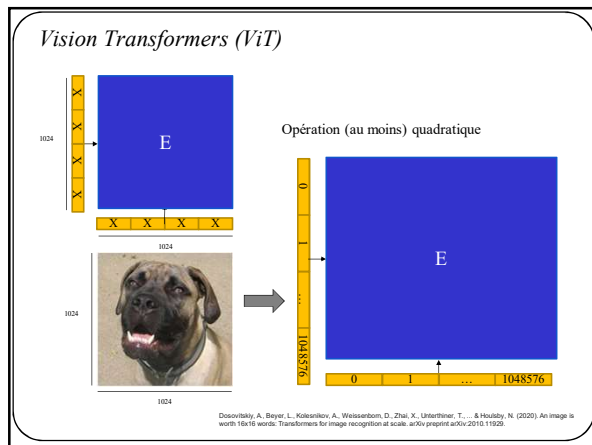
Table 3.1: Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

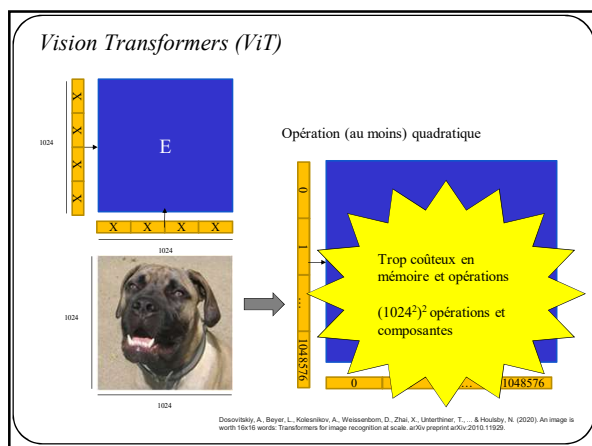












Vision Transformers (ViT)

Chaque "patch" est un token de (par exemple) 256 composantes

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.

Vision Transformers (ViT)

- L'image est séparée en patch

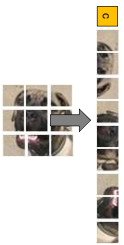
Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissert, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

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., ... & 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., 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.

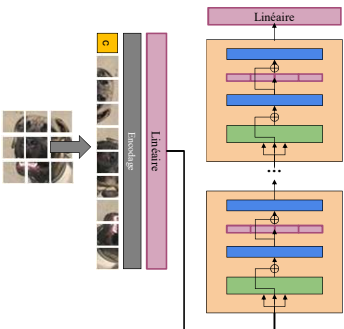
Vision Transformers (ViT)



- 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

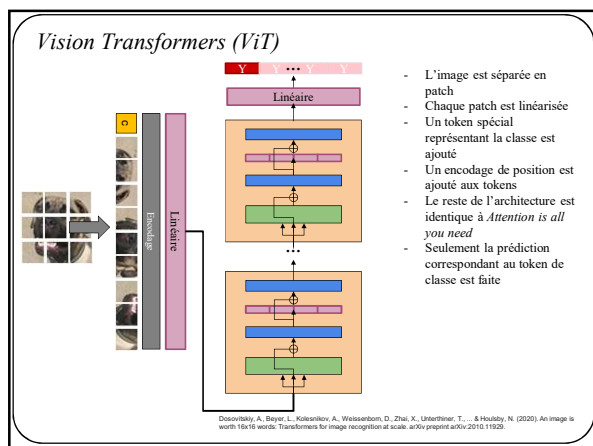
Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissert, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.

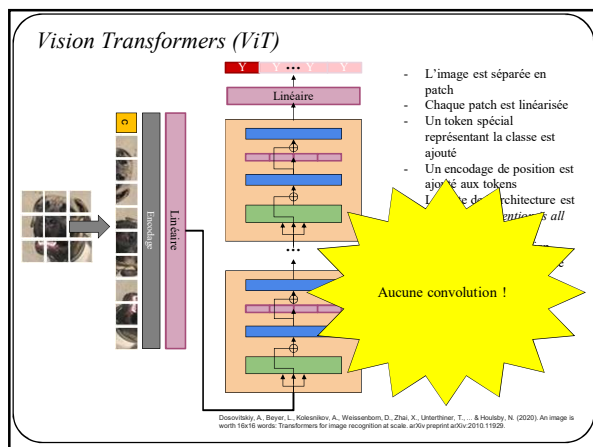
Vision Transformers (ViT)



- 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., 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.



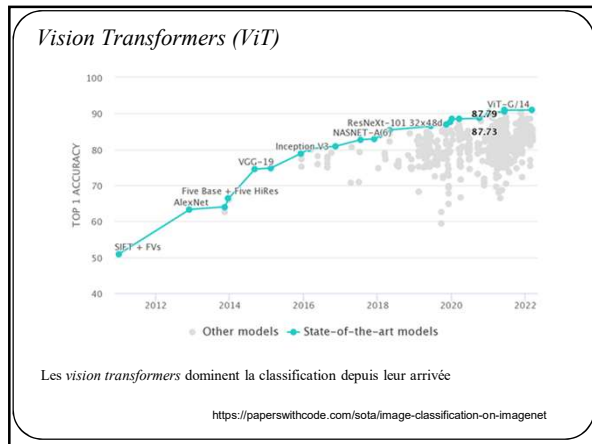


Vision Transformers (ViT)

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

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

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



Vision Transformers (ViT)

Rank	Model	Top-1 Accuracy	Top-5 Accuracy	Params	FLOPs	Open Source	Year
1	ViT-G/14	87.73%	97.90%	140M	~1.3T	Yes	2021
2	DeiT-S/16	87.00%	97.80%	170M	~0.8T	Yes	2021
3	DeiT-T/16	86.70%	97.70%	300M	~1.5T	Yes	2021
4	DeiT-B/16	86.50%	97.60%	600M	~3.0T	Yes	2021
5	ViT-L/14	86.30%	97.50%	300M	~1.3T	Yes	2021
6	ViT-H/14	86.10%	97.40%	600M	~3.0T	Yes	2021
7	DeiT-XL/16	85.90%	97.30%	1.2B	~6.0T	Yes	2021
8	DeiT-M/16	85.70%	97.20%	300M	~1.3T	Yes	2021
9	DeiT-S/16	85.50%	97.10%	170M	~0.8T	Yes	2021
10	DeiT-T/16	85.30%	97.00%	300M	~1.5T	Yes	2021

Les vision transformers dominent la classification depuis leur arrivée

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

Sommaire

train more

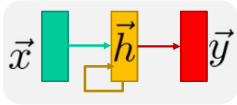
train more

train more

- 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

Sommaire



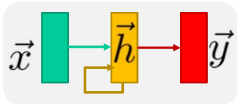
- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique



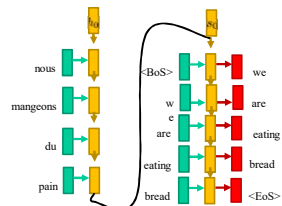
a group of people playing a game with nintendo wii controllers

Description d'images

Sommaire

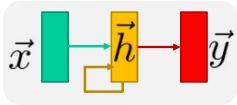


- Les réseaux récurrents peuvent traiter des séquences
- Ils ne requièrent que de légères modifications à des réseaux pleinement connectés
- Ils sont instables sur de longues séquences
- LSTM/GRU sont utilisés en pratique

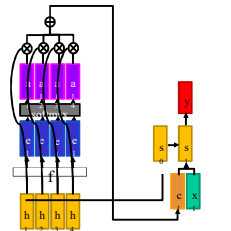


Traduction

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!

