

## GLO-4030/7030 APPRENTISSAGE PAR RÉSEAUX DE NEURONES PROFONDS

## Réseaux Récurrents (RNN)

## Pourquoi RNN?

Traiter des données séquentielles

Image X vs. 
$$\{x^{(1)}, x^{(2)}, \dots, x^{(\tau)}\}$$

- séries temporelles
- séquences de pixels
- séquences de mots
- Souvent de longueur variable
- Pas clair d'avance où l'information pertinente est située

I went to Nepal in 2009

In 2009, I went to Nepal

I like whisky

```
vecteur
           MLP
(tokenization) [I, like, whisky]
```

```
vecteur \{ h_1 \}
                MLP
         [I, like, whisky]
```

```
vecteur { h<sub>1</sub> h<sub>2</sub> h<sub>3</sub> }

↑

MLP

↑

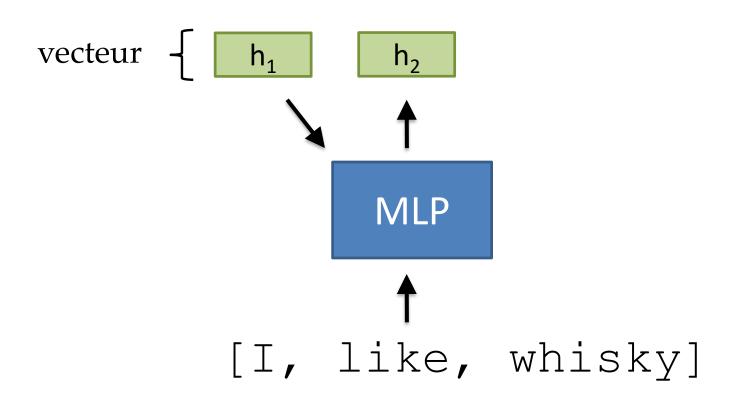
[I, like, whisky]
```

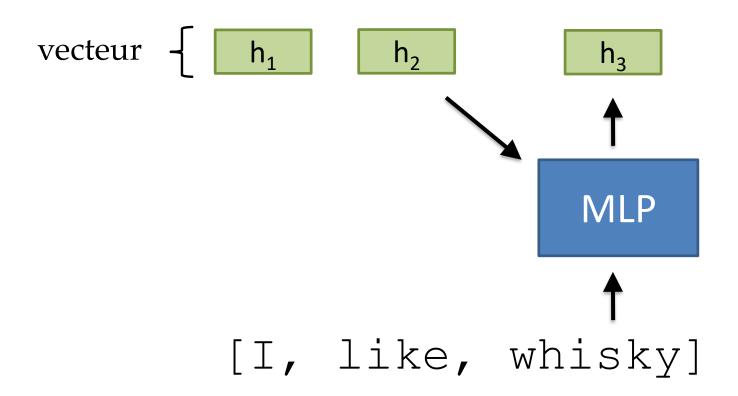
h<sub>1</sub> h<sub>2</sub>

Pas de relation entre les h

[I, like, whisky]

```
vecteurs -
        MLP
        [I, like, whisky]
```





 $h_1$   $h_2$ 

Maintenant  $h_3$  pourra contenir de l'information sur toute la séquence

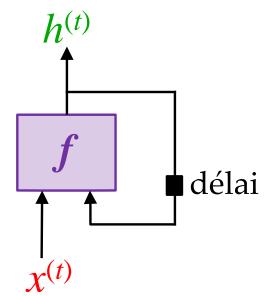
[I, like, whisky]

## Modélisation séquentielle idéale

- 1. Être capable de traiter les séquences de longueur variable
- 2. Garder la trace des dépendances à longterme
- 3. Conserver l'information sur l'ordre
- 4. Partager les paramètres le long de la séquence

## Idée générale

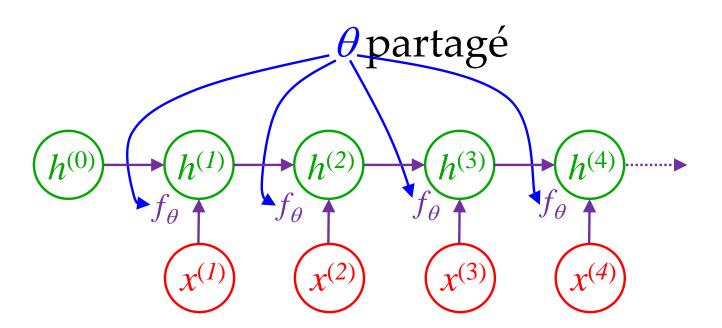
$$h^{(t)} = f(h^{(t-1)}, x^{(t)}, \theta)$$



- Mêmes paramètres  $\theta$  (weight sharing)
- Limite le pouvoir de représentation
  - régularisation
- Relation *f* stationnaire : ne change pas selon *t* 
  - p. e. règle grammaire indépendante de la position
- Lien avec systèmes dynamiques (GMC, GEL)

## Graphe calcul déroulé

$$h^{(t)} = f(h^{(t-1)}, x^{(t)}, \theta)$$



### Variable cachée h

- Résumé sémantique (avec perte) de la séquence (passée, si causal)
- En lien direct avec la tâche:
  - p. e. si on cherche des dates, des mots comme mercredi vont influencer h plus que Québec
  - backprop fera le travail de trouver la fonction
     f favorisant cette représentation
- Taille de *h* influencera la quantité d'information pouvant y être stockée
  - pourra difficilement résumer À la recherche du temps perdu de M. Proust (4 215 pages)
  - généralisation plus difficile si *h* est grand

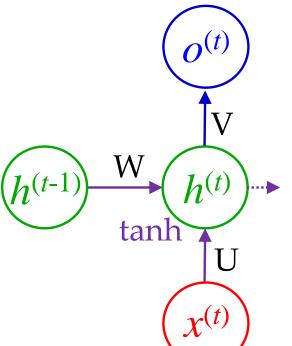
## RNN universel (vanille)

- Utilise des fonctions affines
- tanh comme non-linéarité

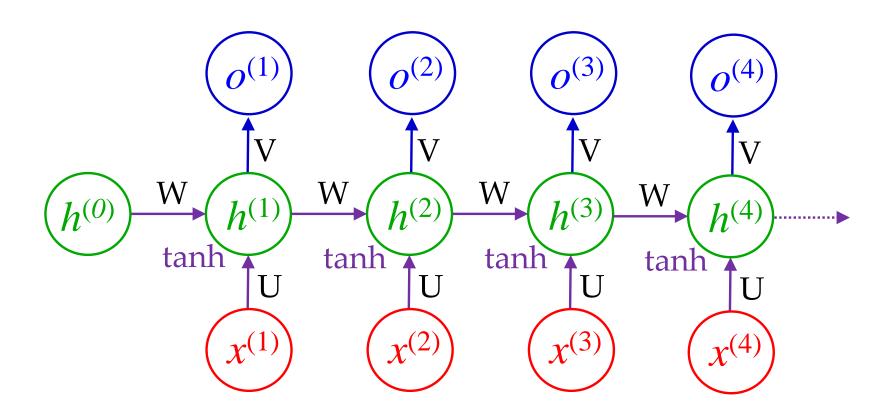
$$h^{(t)} = tanh(Wh^{(t-1)} + Ux^{(t)} + b)$$

$$o^{(t)} = Vh^{(t)} + c$$

- Peut accomplir autant qu'une machine de Turing
- Variante assez commune
- Défaut : on ne peut pas paralléliser forward/backward pass
  - doit faire la séquence au complet en sériel



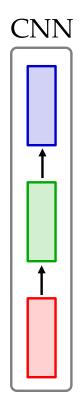
## RNN vanille déroulé



## Pourquoi tanh?

- Non-linéaire
- Toujours dérivable
- Sortie [-1,1] (enlever/ajouter)
- Symétrique
- Pas de biais systématique
  - sigmoïde va de [0,1], induit biais
- Autres?

## Topologie Feedforward



adapté de cs231n

## Topologie RNN

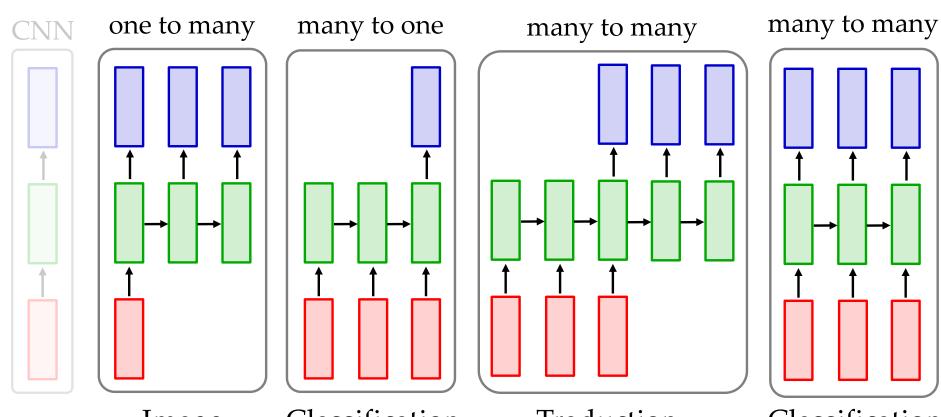


Image captioning

Classification de sentiment (texte)

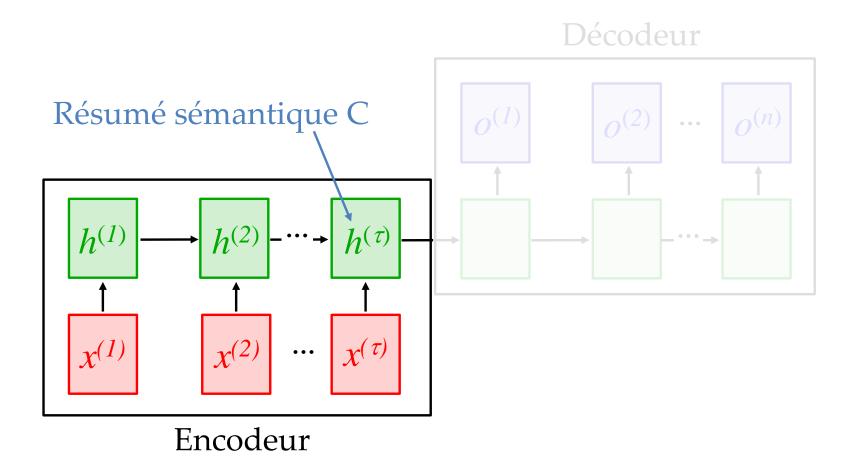
Traduction,
Réponse aux
questions
(tailles entrée/sortie
variables)

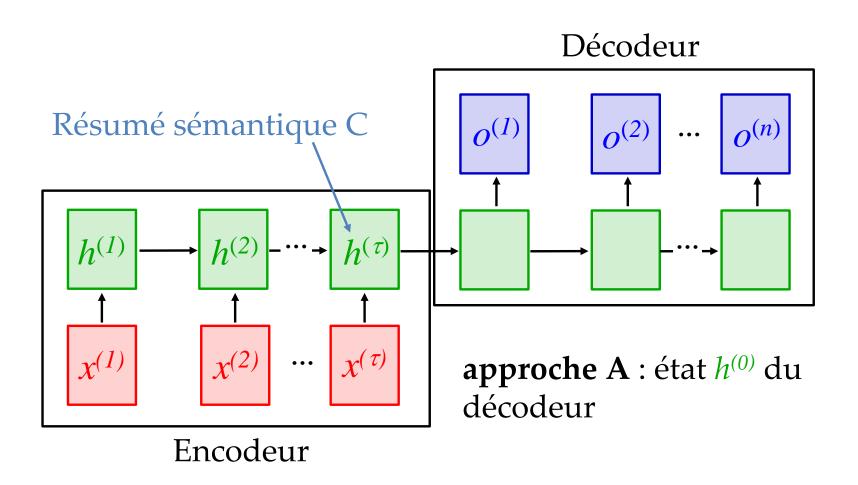
Classification de trames vidéos

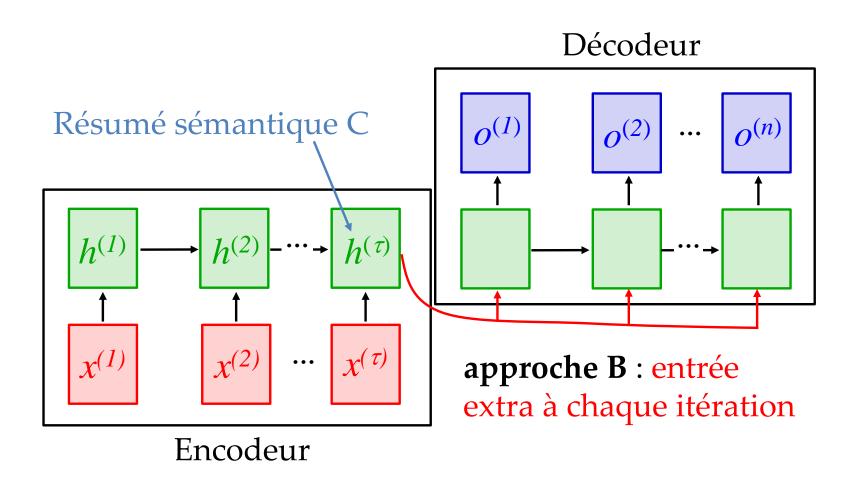
Architecture many to many

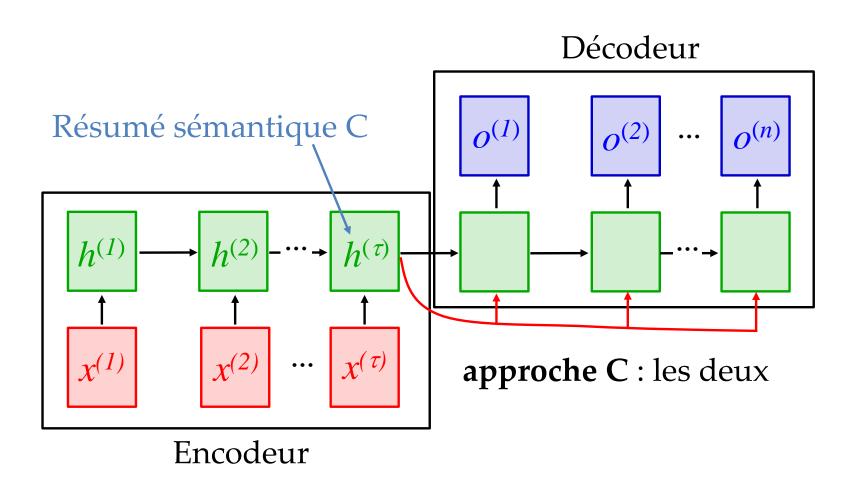
contexte

• Généré une séquence à partir d'un résumé C









# Exemple de génération avec RNN

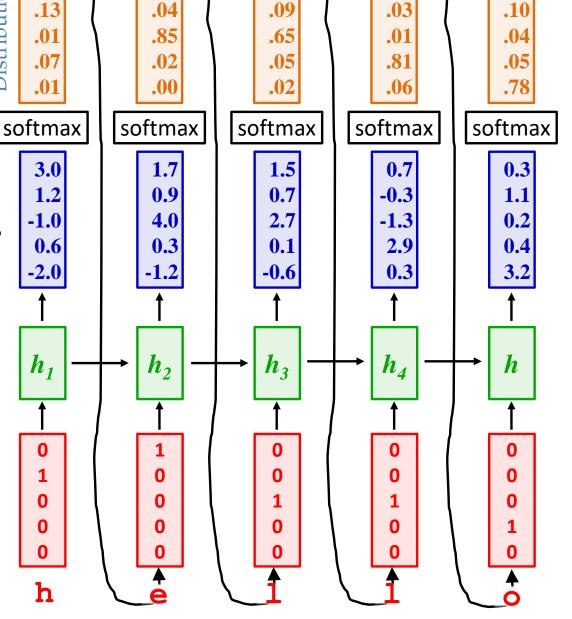
pige e

.78

pige 1

.09

- Réseau entraîné à prédire des caractères {e,h,l,o,<end>}
- Entraîné sur hello



pige 1

.19

pige o pige <end>

.04

27

.09

## Exemple: entraîné sur Shakespeare

### Sortie:

- Réseau RNN à trois couches
- 512 neurones cachées par couche
- Entraîné sur 4.4
   Mo de données
   texte

### **PANDARUS:**

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

### Second Senator:

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

### **DUKE VINCENTIO:**

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

### Second Lord:

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

### Clown:

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

### VTOLA:

I'll drink it.

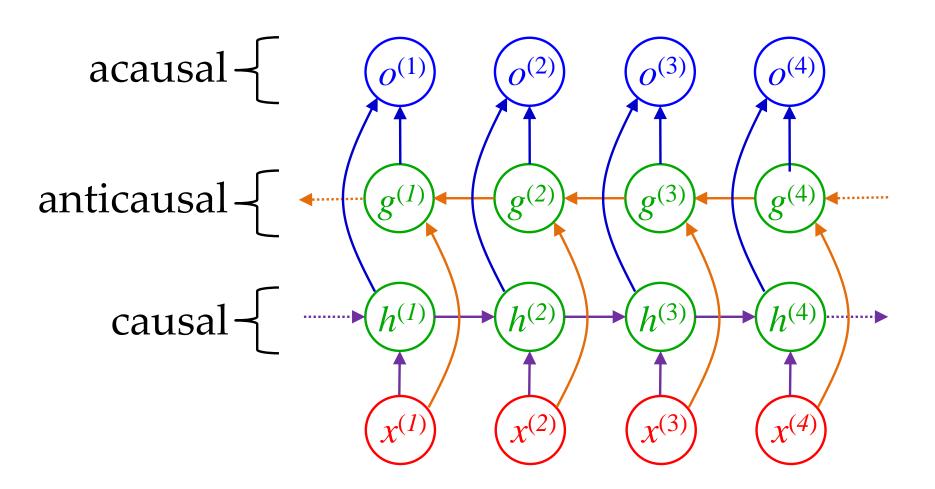
## Longueur de sortie o(t)

- Lors de la génération, on doit s'avoir quand arrêter d'échantillonner le RNN
- 3 stratégies :
  - 1. Symbole spécial (**<END>**)
  - 2. Sortie supplémentaire 0-1 (via sigmoïde), qui prédit la fin
  - 3. Sortie qui prédit  $\tau$  directement (régression)

## RNN bi-directionnel

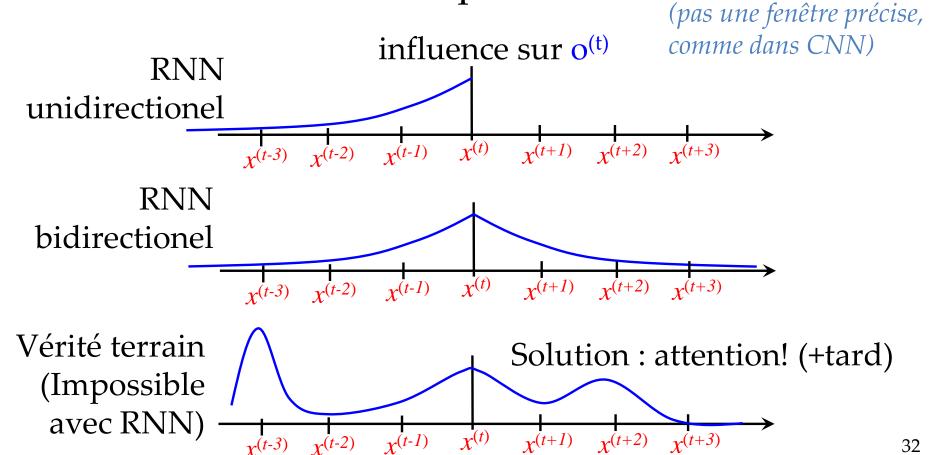
- Sortie  $o^{(t)}$  peut dépendre de toute la séquence (1 à  $\tau$ )
- Information pertinente parfois après une entrée *x* 
  - ordre des mots dans une langue
    - adjectif avant ou après un mot
    - langue SVO, SOV, V2, etc...
  - reconnaissance de la voix
    - coarticulation
  - bio-informatique

## RNN bi-directionnel



## Longue portée

- Influence à longue portée difficile dans RNN
- CNN: champ récepteur croissant en profondeur
- RNN : décroissance exponentielle de l'influence



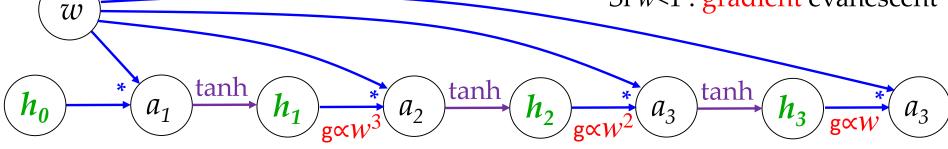
## Gradient et entrainement

## Exploding/vanishing gradient

- Poids W partagés
- Exemple simplifié:

Si *w*>1 : gradient explose

Si *w*<1 : gradient évanescent



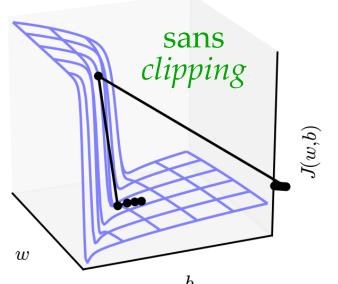
Pour un réseau récurrent linéaire (simplification) :  $h^{(t)} = W^{T} h^{(t-1)}$ 

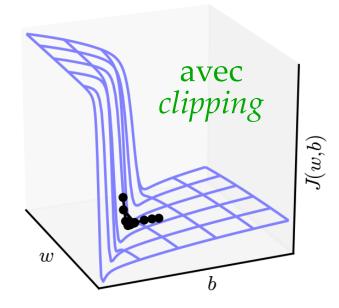
Décomposition en éléments propres  $W = Q\Lambda Q^T$ 

$$h^{(t)} = Q^{T} \Lambda^{t} Q h^{(0)}$$
  $\Longrightarrow$  si valeur propre  $\lambda > 1$ : vecteur propre explose si valeur propre  $\lambda < 1$ : vecteur propre évanescent

## Gradient clipping pour entraînement RNN

• Ravins typique dans les RNN:





• Solution, clipper:

la norme du gradient

$$if \|\vec{g}\| > v$$

$$\vec{g} \frac{v}{\|\vec{g}\|}$$

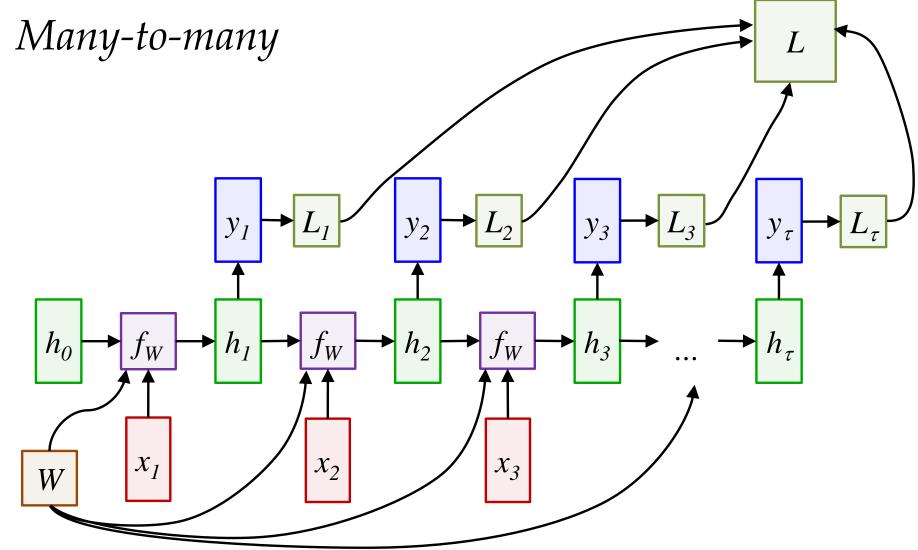
les entrées du gradient, individuellement

$$\vec{g} = \begin{bmatrix} 88493.4 \\ -0.3 \\ ... \\ -9948423 \end{bmatrix} = \begin{bmatrix} v \\ -0.3 \\ ... \\ -v \end{bmatrix}$$

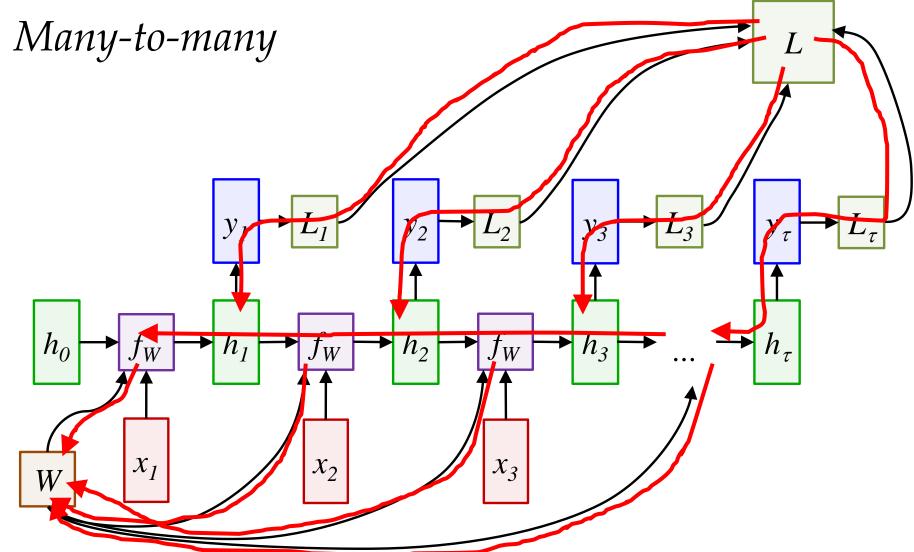
(similaire comme scaling par paramètre, lors de l'optimisation)

Si NaN, bouger au hasard d'une magnitude v

## Calcul du gradient sur graphe déroulé



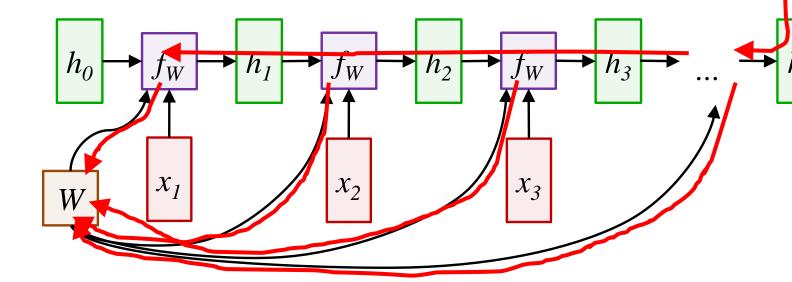
Calcul du gradient sur graphe déroulé



## Calcul du gradient sur graphe déroulé

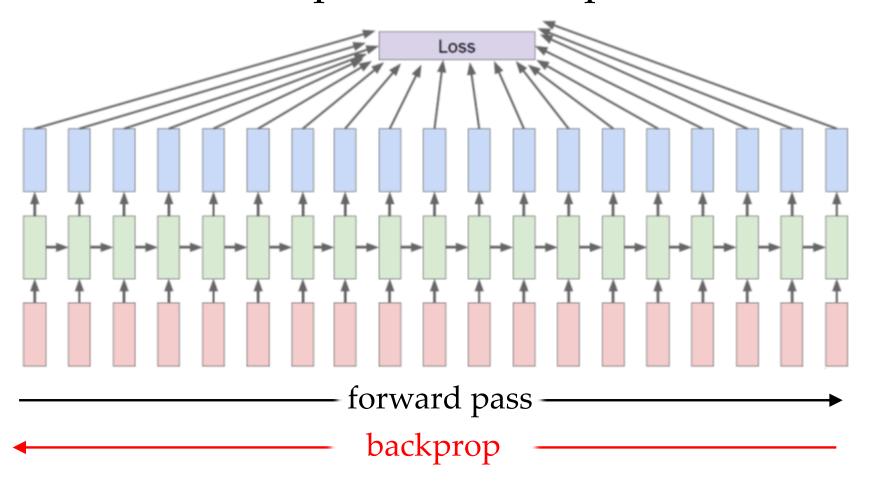
Many-to-one

Moins d'entrées du **gradient** dans le graphe + vanishing gradient : entraînement plus difficile.

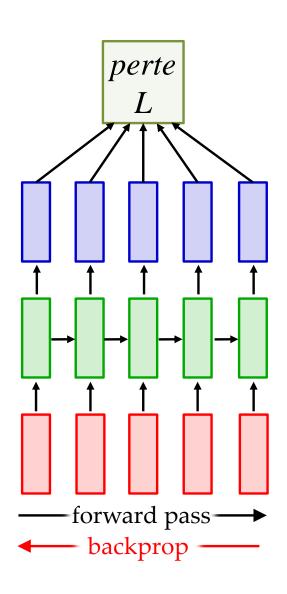


## Backprop through time (BPTT)

Calcule la séquence au complet

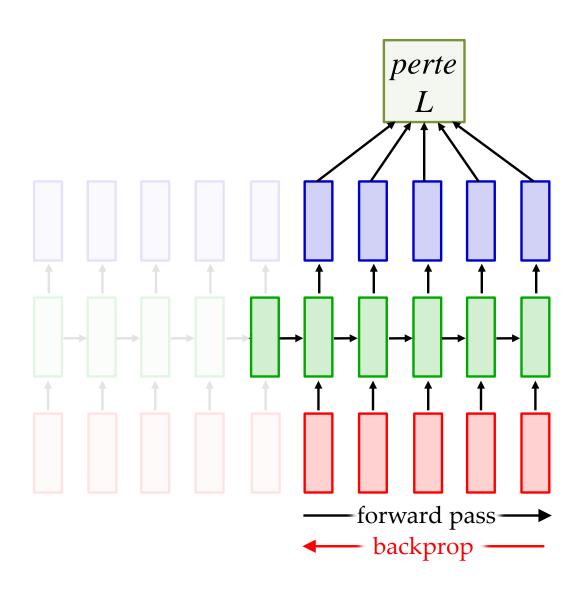


## Truncated BPTT



Effectue BPTT sur des segments de la séquence

## **Truncated BPTT**



## Truncated BPTT

