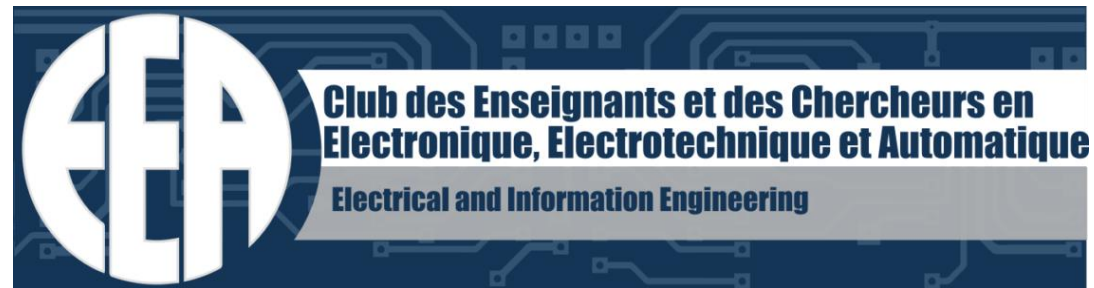


Architectures pour l'apprentissage profond

Bart.Lamiroy@univ-reims.fr

février 2024



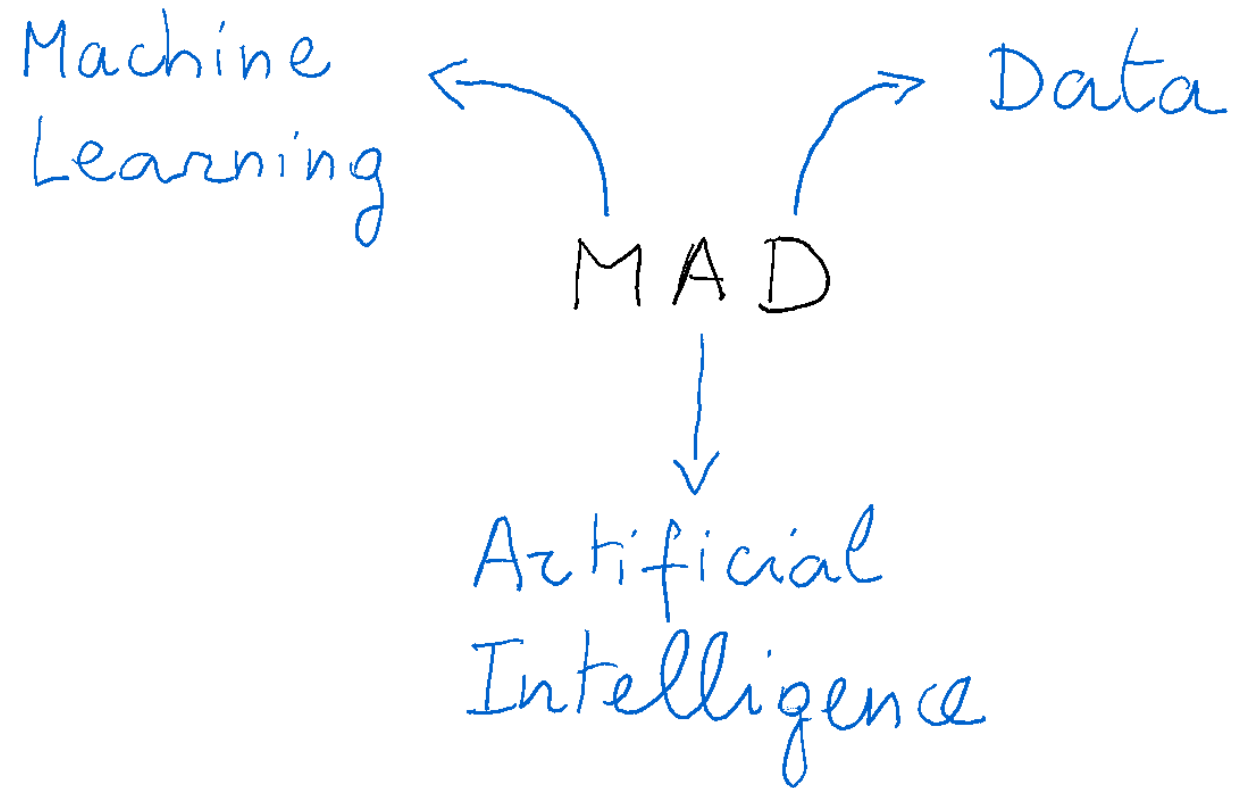
Intelligence Artificielle ?

- Evolution technologique
- Changement de culture
- Opportunité d'évolution



Source: giphy

MAD ...



Quelques mises en garde (voire conseils)

- Attention au vocabulaire anthropocentrique
« apprentissage », « intelligence », « décider », « voir », « choisir », « savoir », « comprendre » ...
- Quelqu'un qui se prétend « Expert en IA » est ~~très probablement~~ un mytho.
- Les projets de mise en œuvre MAD nécessitent **obligatoirement** des compétences transversales
(et donc des aller-retours entre experts en techniques de MAD et professionnels du domaine)
- Les moteurs actuels de l'adoption d'outils MAD et leur création viennent de « bricoleurs » éclairés ... n'hésitez pas d'essayer.

Objectifs généraux

Présentation des grandes familles de réseaux de neurones profonds

Comprendre

- Comment leur structure influe sur leurs capacités
- Comment ils sont entraînés
- Quels usages en découlent

Eclairage des LLM et agents conversationnels

Plan de l'exposé

- Vocabulaire et fondements scientifiques
- Quelques rappels sur le fonctionnement des réseaux de neurones
- Catégories de réseaux de neurones profonds
- Applications

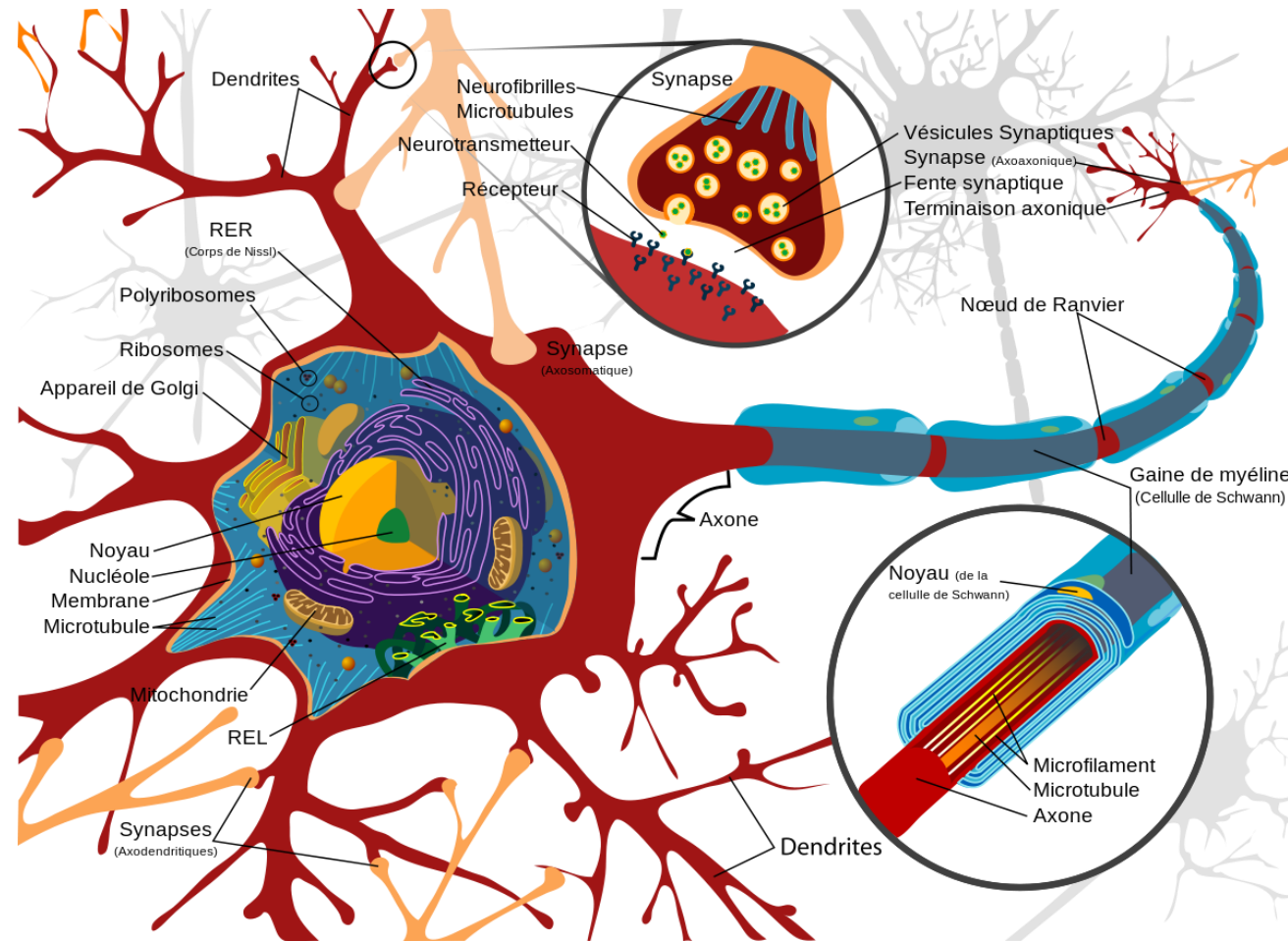
Quelques rappels (ou pas)

- Apprentissage supervisé
- Apprentissage auto-supervisé
- Apprentissage non supervisé

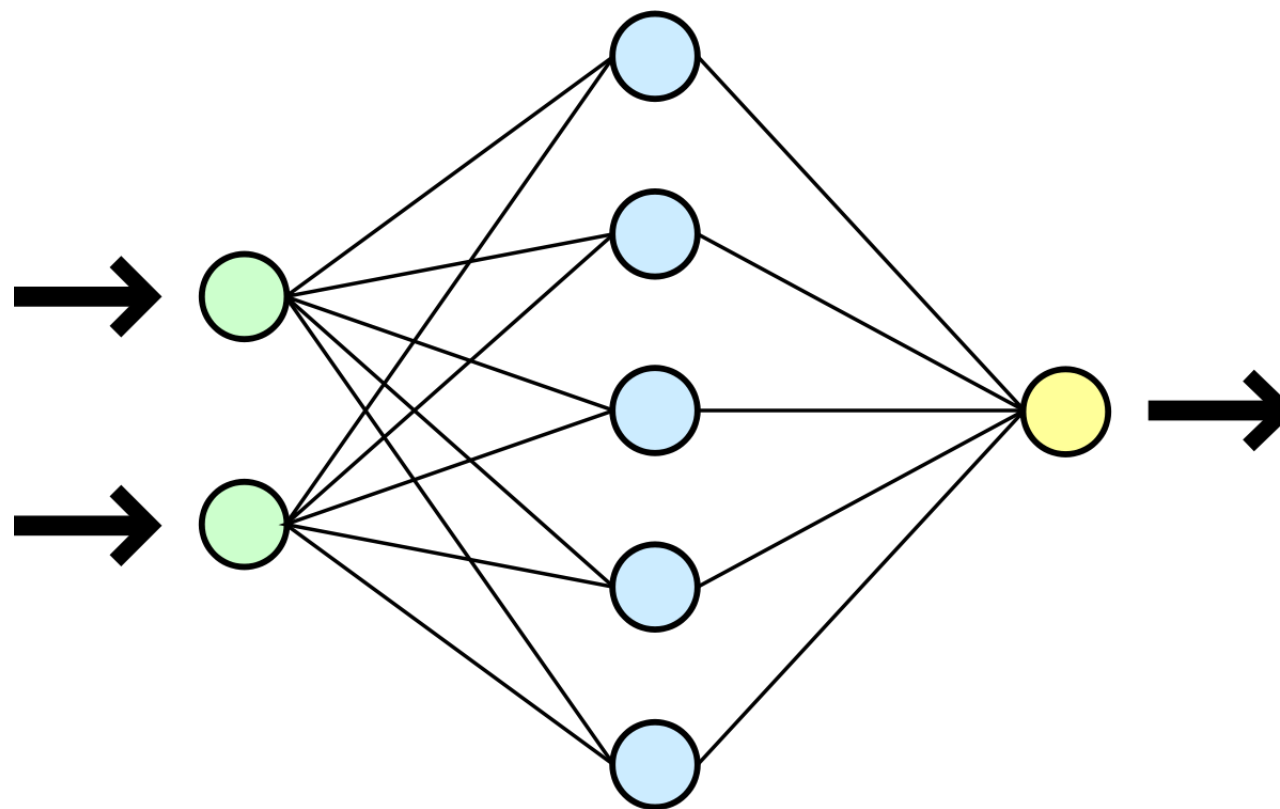
Réseaux de neurones artificiels

- Connexionnisme
- Fonction d'activation
- Fonction de coût
- Rétropropagation du gradient
- Théorème de l'approximateur universel

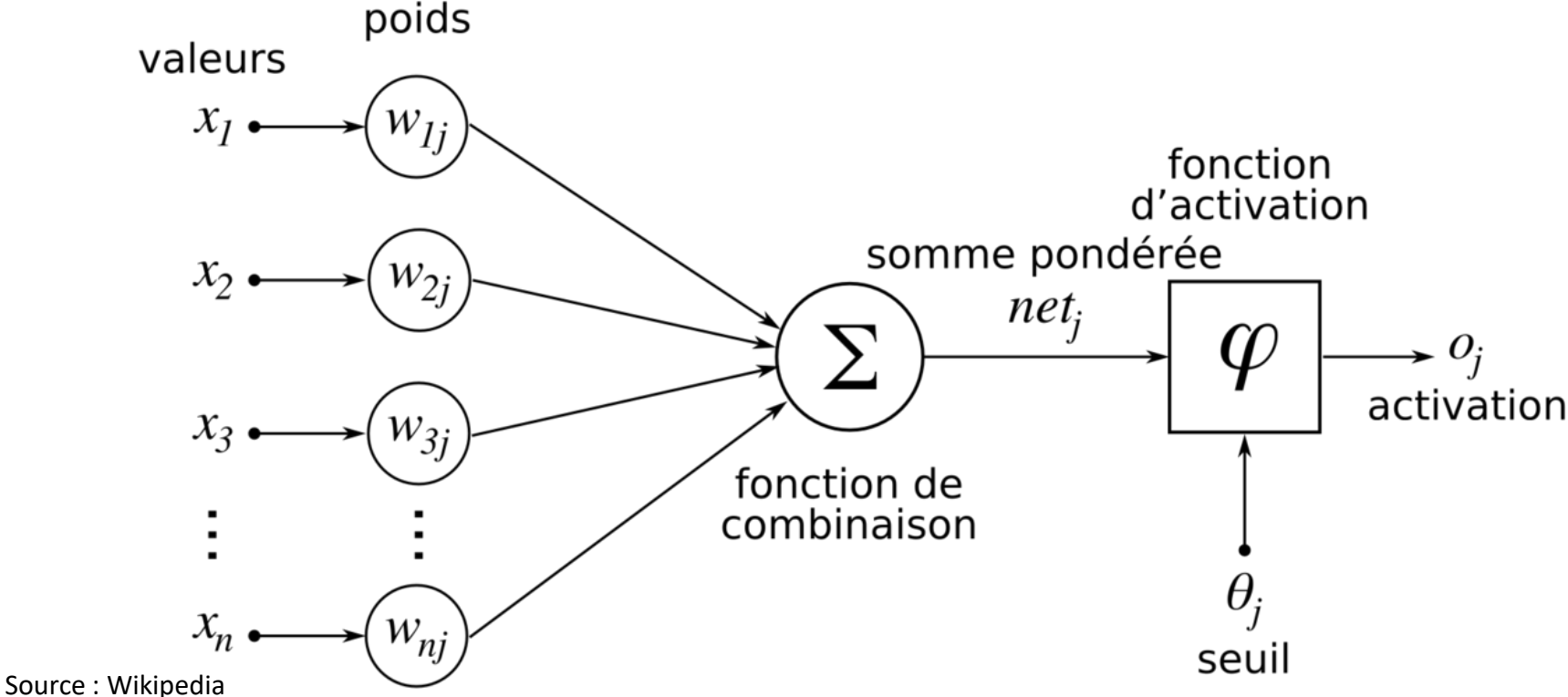
Le neurone biologique et son fonctionnement



Modèle simplifié : graphe orienté acyclique



Le neurone artificiel

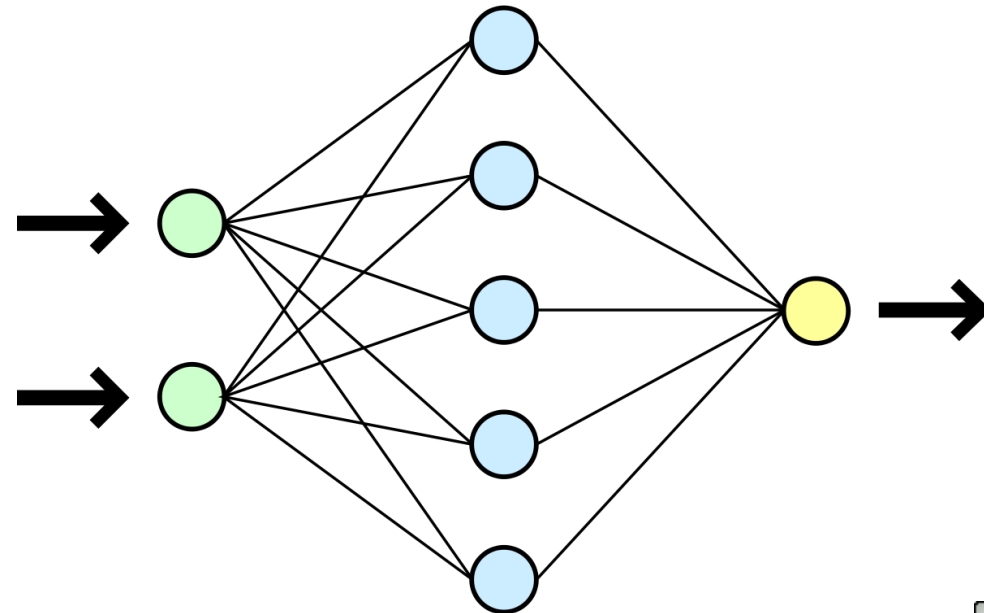


F. Rosenblatt (1958), "The perceptron: a probabilistic model for information storage and organization in the brain", repris dans J.A. Anderson & E. Rosenfeld (1988), Neurocomputing. Foundations of Research, MIT Press



D'un point de vue calculatoire

- Des produits scalaires
- Très parallélisable
- Opérations SIMD (Single Instruction, Multiple Data) très adaptées à des GPU
- Approximateur universel



Le théorème de l'approximateur universel

- Etant donnée ϕ une fonction strictement croissante bornée ;
- Soit I_m l'hypercube unite de dimension m et \mathcal{C} l'ensemble des fonctions continues sur I_m ;
- Alors pour tout ε , il existe $N \in \mathbb{N}$ tel que pour toute fonction $f \in \mathcal{C}$ il existe $v_i, b_i \in \mathbb{R}$ et $w_i \in \mathbb{R}^m$ tels que pour tout $x \in I_m$

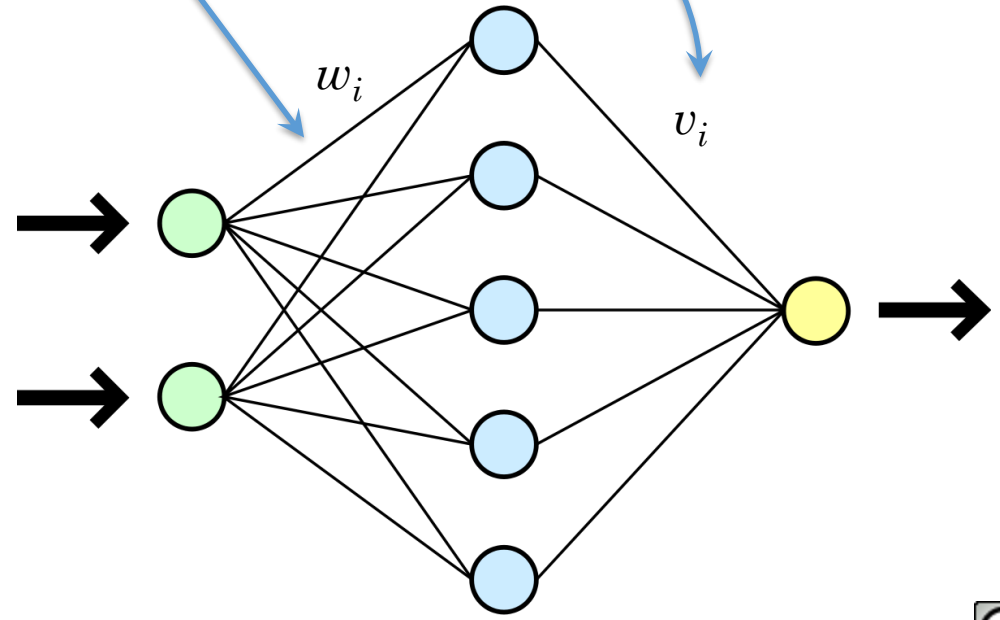
$$|F(x) - f(x)| < \varepsilon$$

$$F(x) = \sum_{i=1}^N v_i \phi(w_i^T x + b_i)$$

Le Perceptron

Formule de l'approximateur universel

$$F(x) = \sum_{i=1}^N v_i \varphi(w_i^T x + b_i)$$



K. Hornik (1991), "Approximation Capabilities of Multilayer Feedforward Networks", Neural Networks, Vol. 4, pp. 251-257, Pergamon Press, 1991.
https://web.njit.edu/~Eusman/courses/cs677_spring21/hornik-nn-1991.pdf

Le Perceptron

Défi : déterminer les paramètres $v_i, b_i \in \mathbb{R}$ et $w_i \in \mathbb{R}^m$

Solution : approche par *apprentissage supervisé* en formulant le problème comme une *minimisation par descente de gradient*.

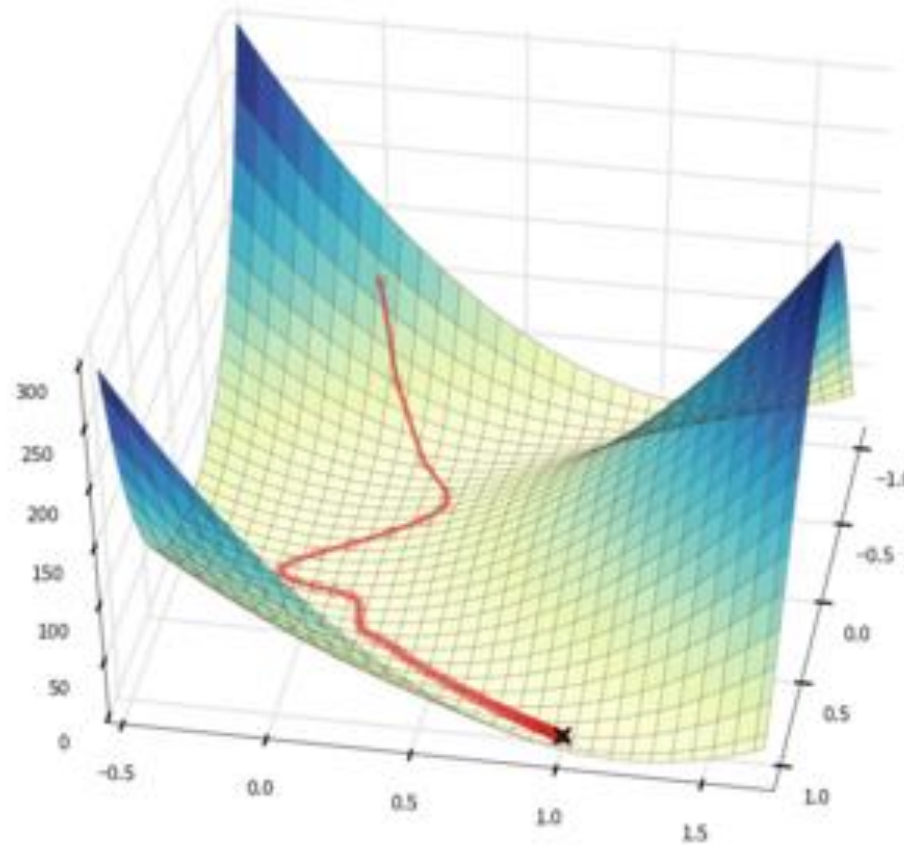
Petit souci : on sait que N existe, mais on ne sait pas comment le calculer.

$$F(x) = \sum_{i=1}^N v_i \varphi(w_i^T x + b_i)$$

Apprentissage supervisé et calcul de gradient

- Couples (entrée \mathbf{x} , sortie attendue $f(\mathbf{x})$)
- Présenter l'entrée et observer la sortie $F(\mathbf{x})$
- La sortie $F(\mathbf{x})$ est une combinaison de fonctions linéaires et la fonction d'activation, paramétrées par les poids des « synapses » \mathbf{b}_i , \mathbf{v}_i , \mathbf{w}_i .
- Il est possible de reformuler le paramétrage du réseau comme une minimisation de l'erreur par descente du gradient sur les poids \mathbf{b}_i , \mathbf{v}_i , \mathbf{w}_i .

Descente de gradient



Calcul du gradient

$$E = \frac{1}{2} \sum_{i=1}^p \|\mathbf{o}_i - \mathbf{t}_i\|^2.$$

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_\ell} \right).$$

Chaque poids w_i est ajusté avec l'incrément Δw_i

$$\Delta w_i = -\gamma \frac{\partial E}{\partial w_i} \quad i = 1, \dots, \ell,$$

Calcul du gradient

Erreur observée
(fonction de coût)

$$E = \frac{1}{2} \sum_{i=1}^p \| \underbrace{o_i}_{\text{Sortie observée}} - \underbrace{t_i}_{\text{Sortie attendue}} \|^2.$$

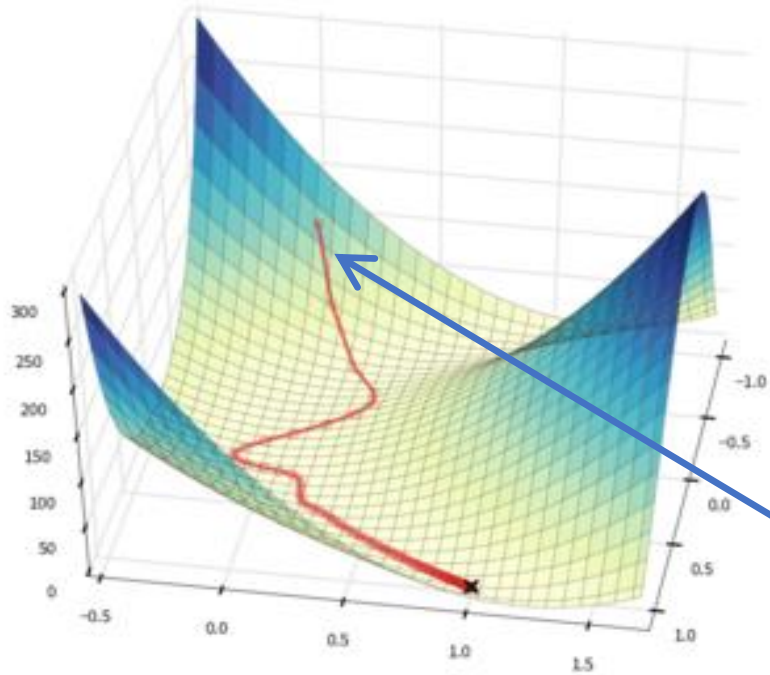
$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_\ell} \right).$$

Poids du
reseau de
neurones

Chaque poids w_i est ajusté avec l'incrément Δw_i

$$\Delta w_i = -\gamma \frac{\partial E}{\partial w_i} \quad i = 1, \dots, \ell,$$

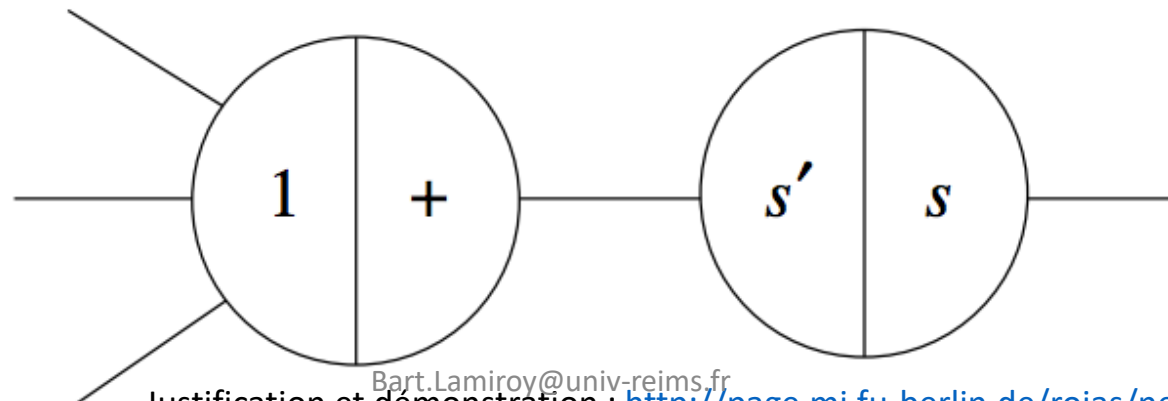
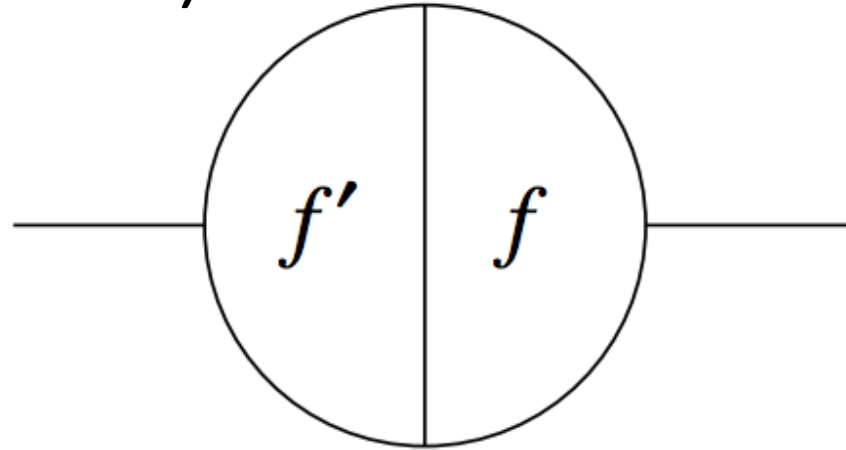
Calcul du gradient



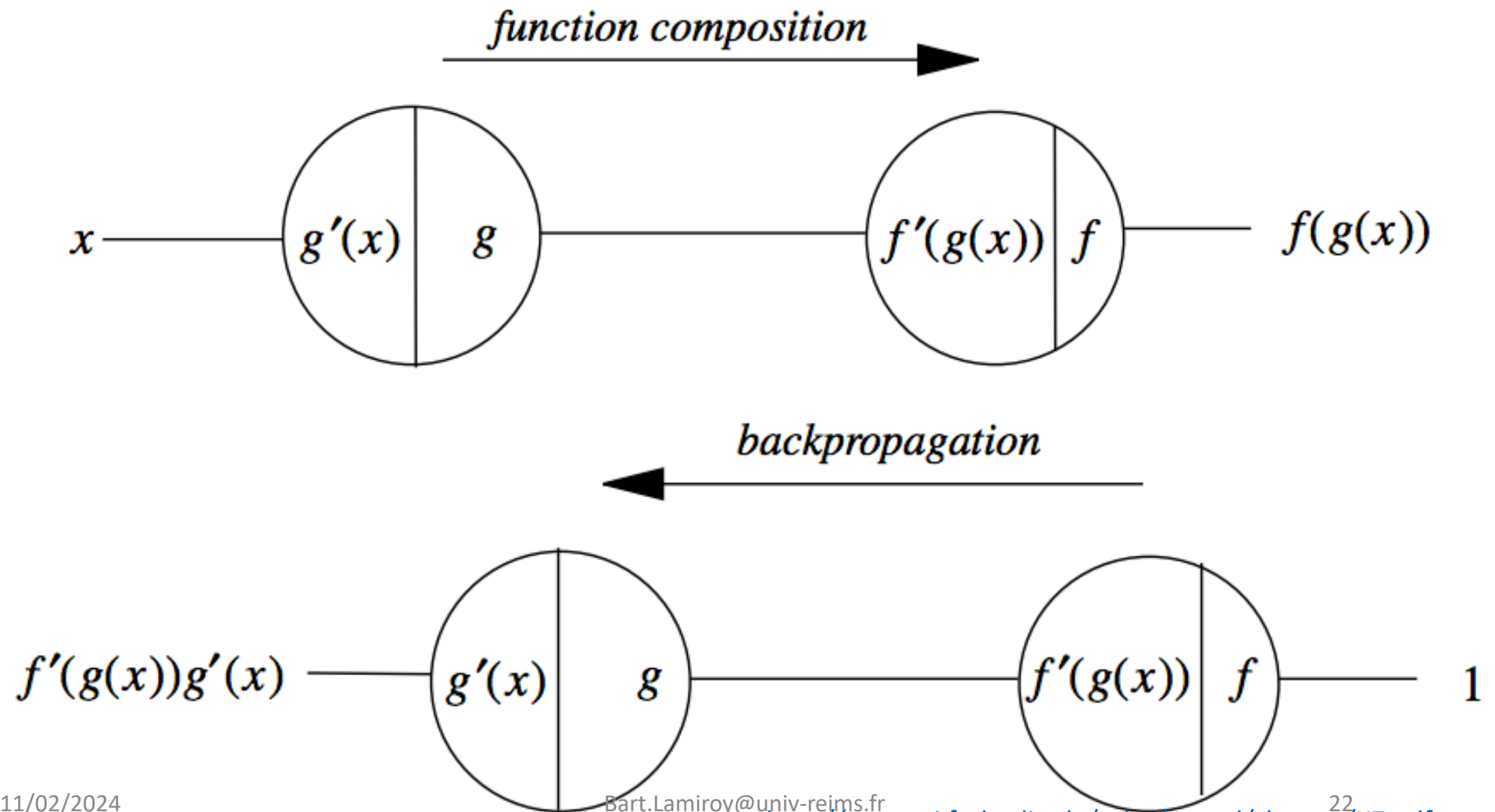
- Learning rate
- Decay
- Momentum

$$\Delta w_i = -\gamma \frac{\partial E}{\partial w_i} \quad i = 1, \dots, \ell,$$

Algorithme de rétropropagation du gradient (cas à une variable)



Algorithme de rétropropagation



Résultat

Il est assez trivial de démontrer que

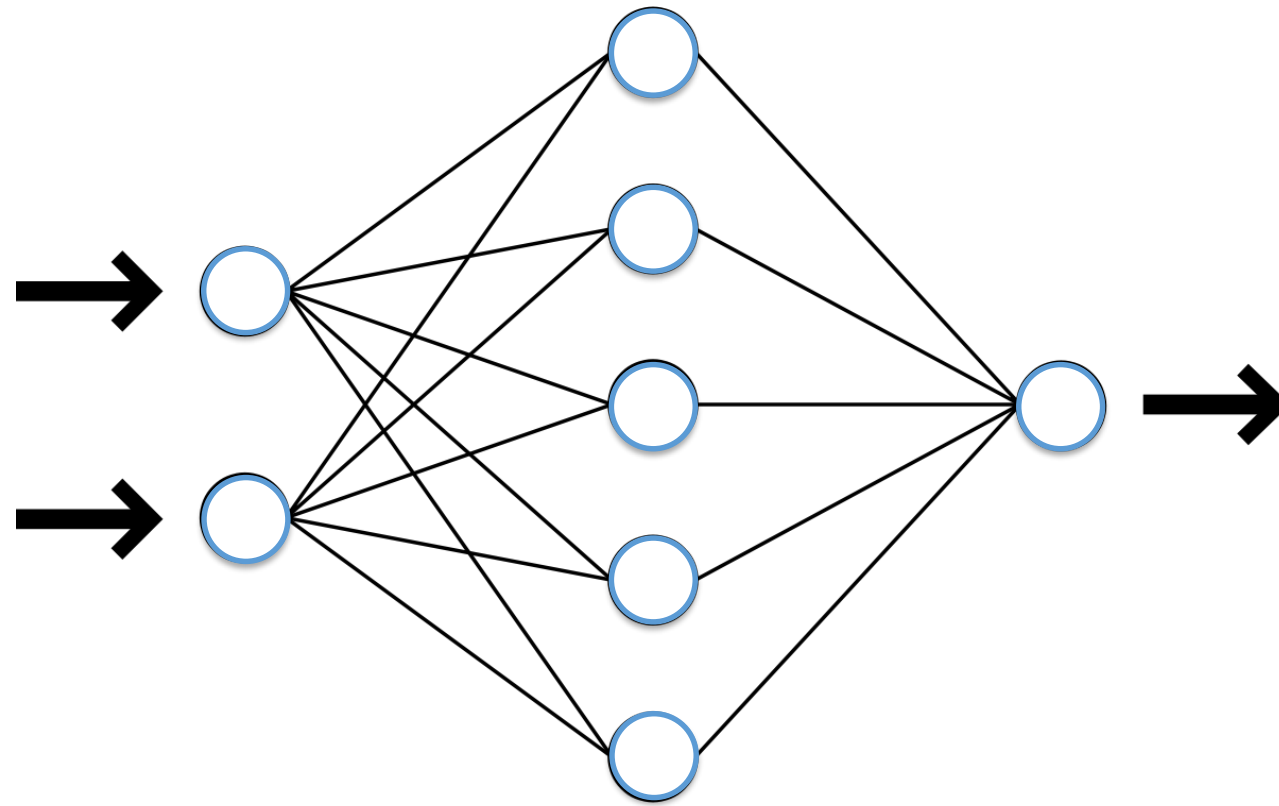
$$\frac{\partial E}{\partial w_{ij}} = o_i \delta_j. \quad \Delta w_{ij} = -\gamma o_i \delta_j.$$

Où δ_j est la valeur de la fonction de coût rétropropagée au nœud j et où o_i est la valeur calculée au nœud i lors du calcul *feed-forward*.

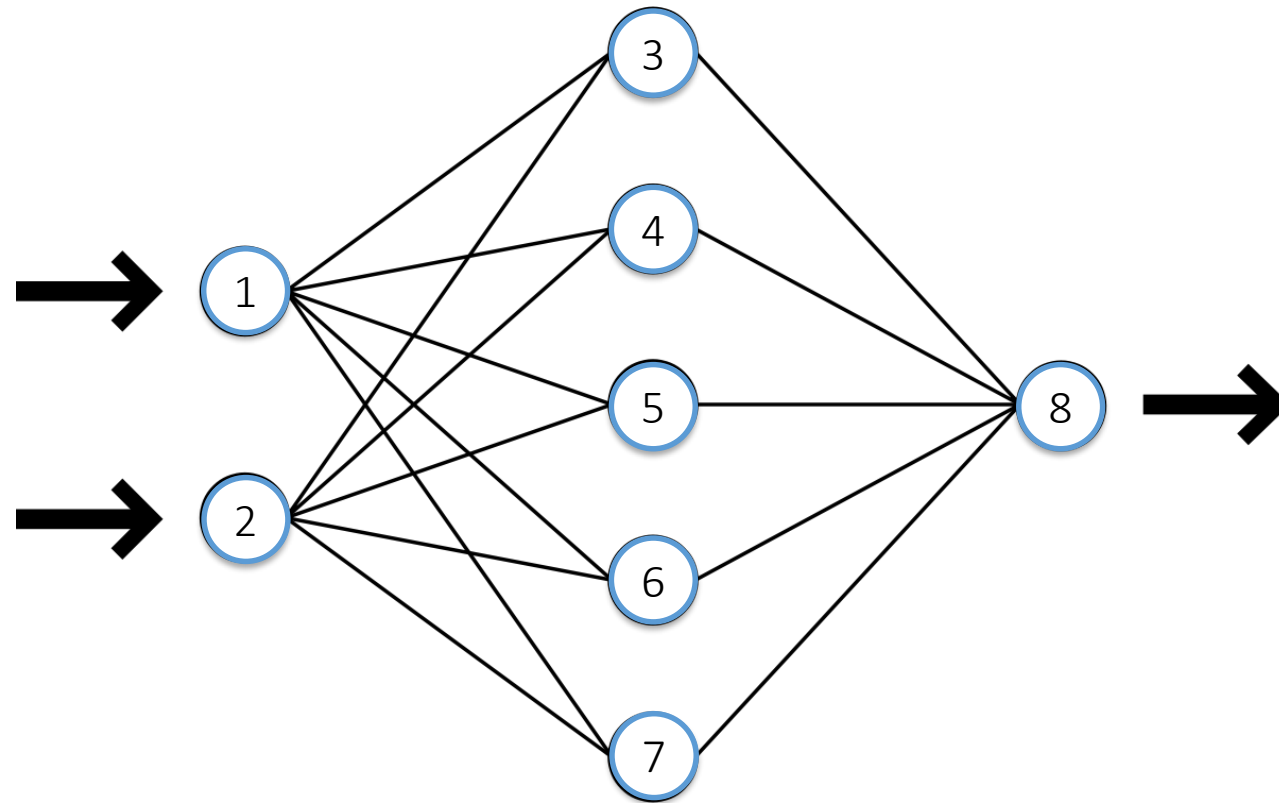
Algorithme final d'apprentissage supervisé

- Recueillir “suffisamment” de données de type $(\mathbf{x}, \mathbf{f}(\mathbf{x}))$
- Initialiser les poids \mathbf{w}_{ij} du réseau de neurones à des valeurs aléatoires
- De façon itérative :
 - Exécuter le réseau en *feed-forward* avec une valeur \mathbf{x}_k et mémoriser tous les états o_i des neurones intermédiaires.
 - Calculer l'erreur E entre la sortie du réseau $F(\mathbf{x})$ et la valeur $\mathbf{f}(\mathbf{x})$ attendue
 - Rétropropager l'erreur à travers le réseau et mémoriser les valeurs de δ_j
 - Ajuster les poids du réseau $\Delta \mathbf{w}_{ij} = -\gamma o_i \delta_j$
 - Continuer jusqu'à convergence

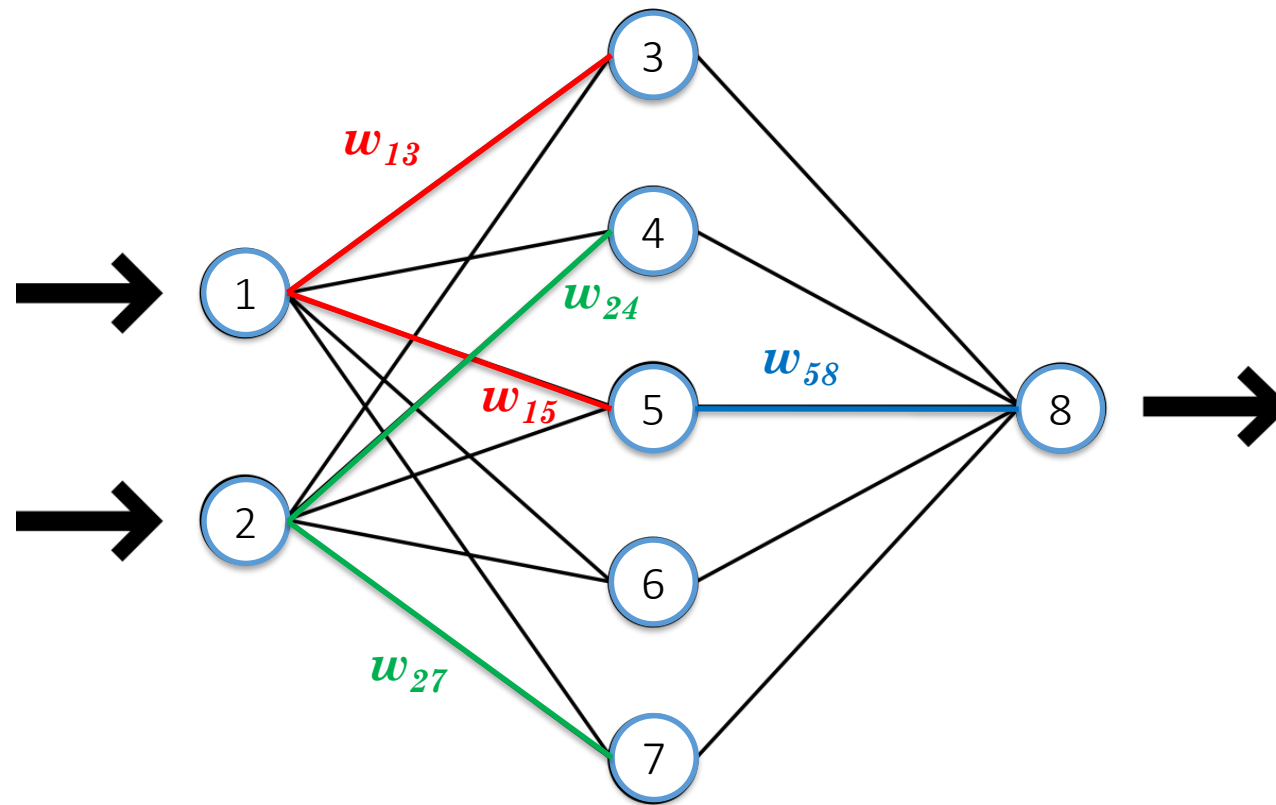
Feed-Forward Phase



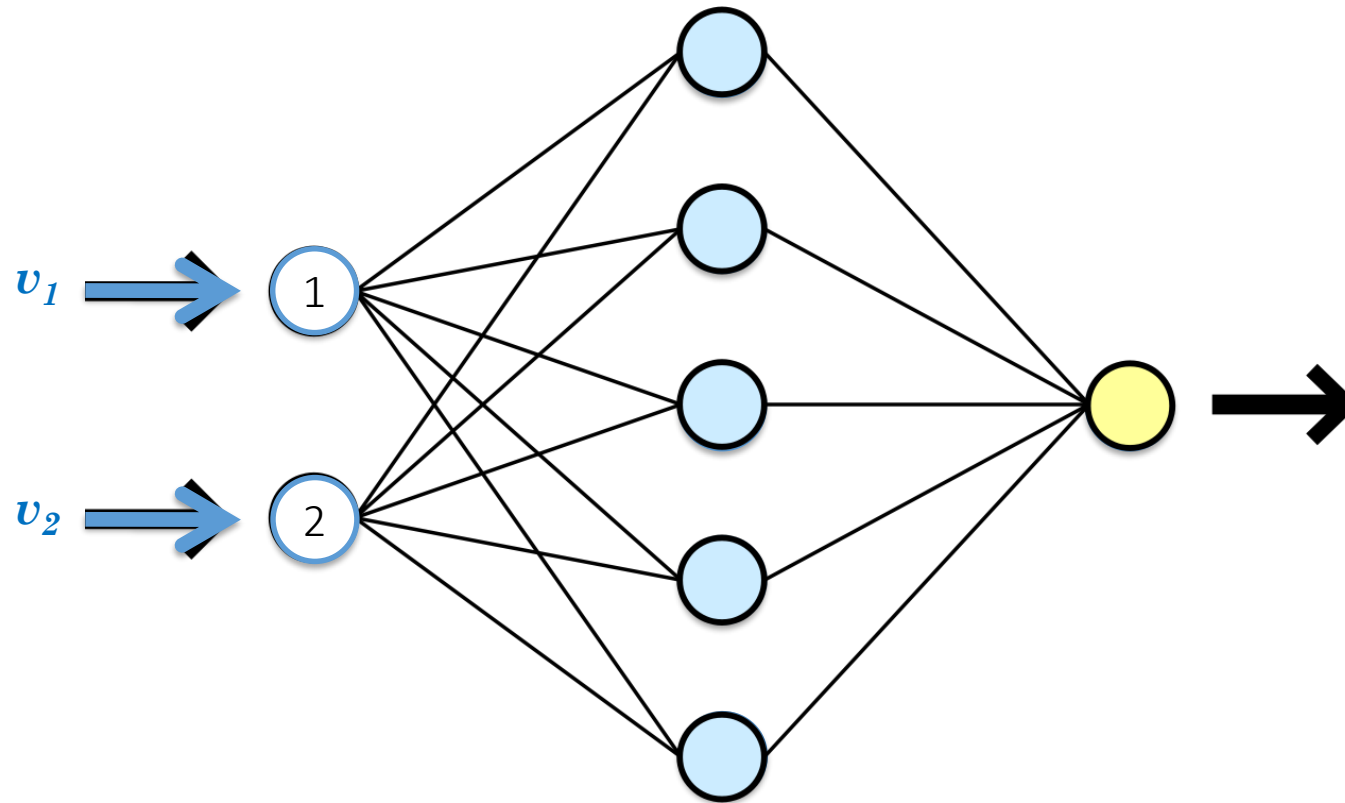
Feed-Forward Phase (notations)



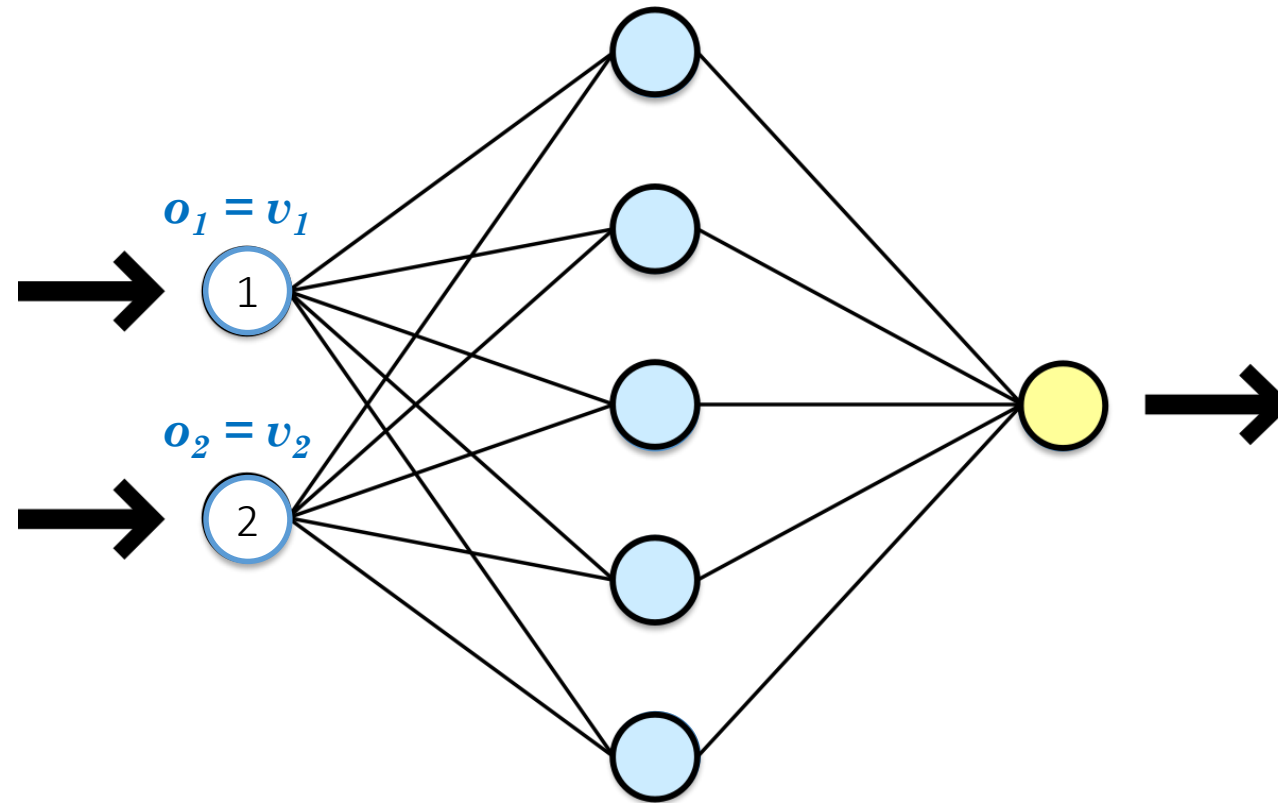
Feed-Forward Phase (notations)



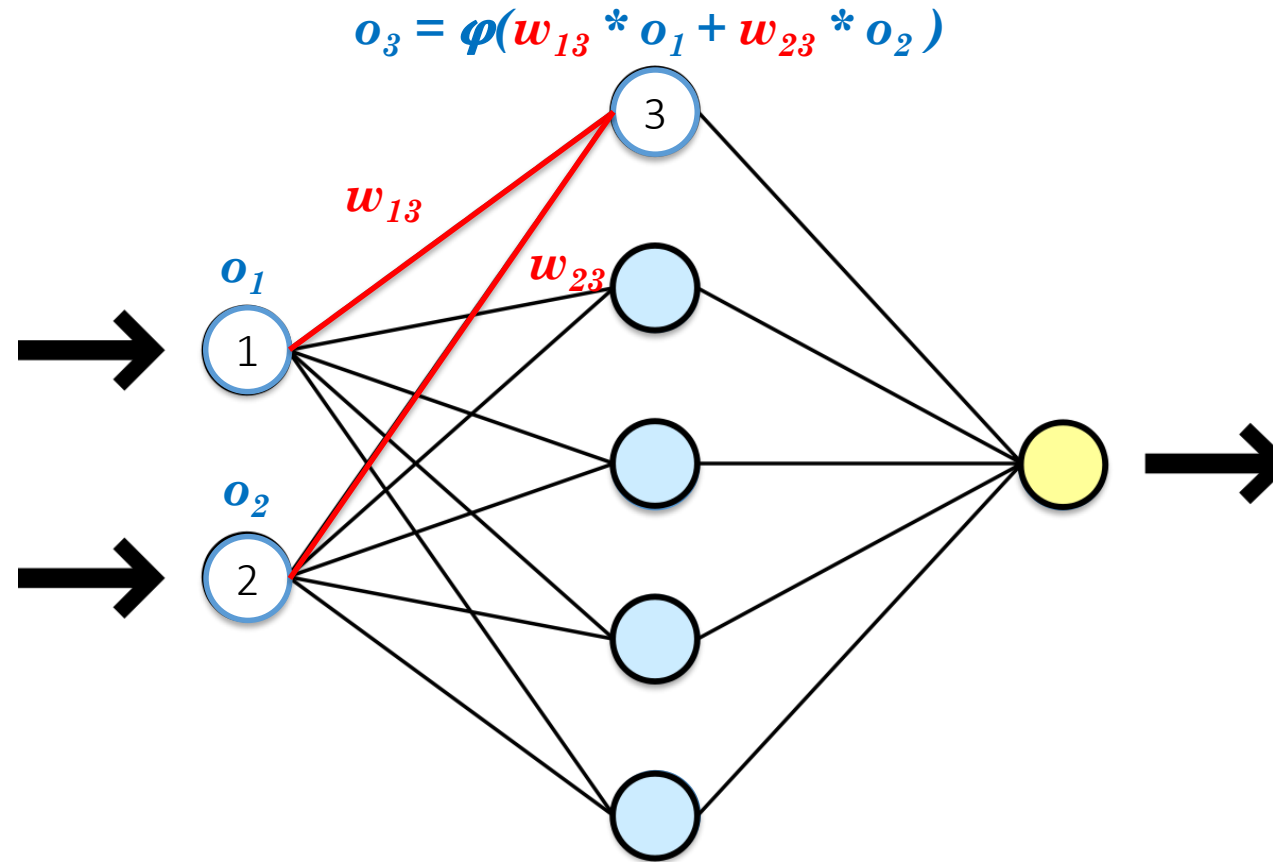
Feed-Forward Phase (Step 0)



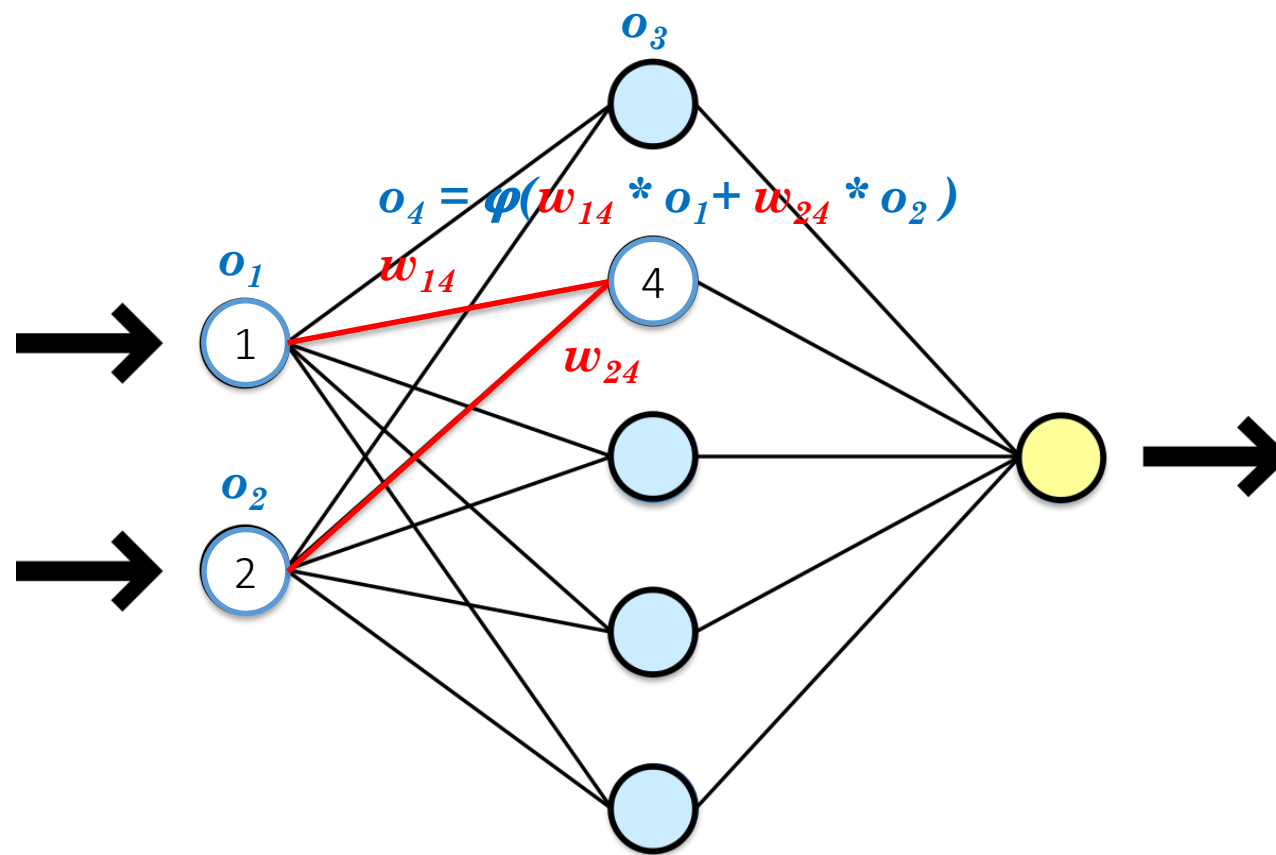
Feed-Forward Phase (Step 1)



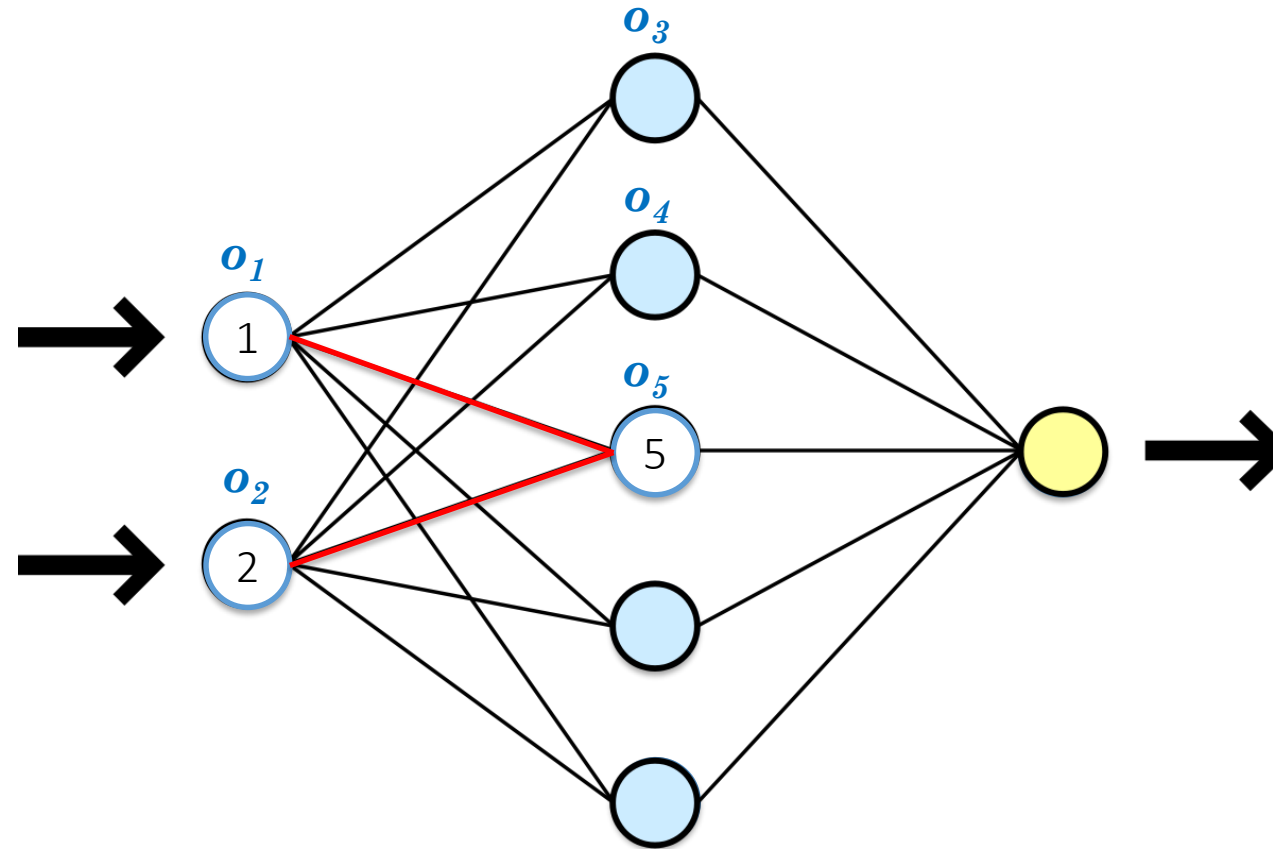
Feed-Forward Phase (Step 2)



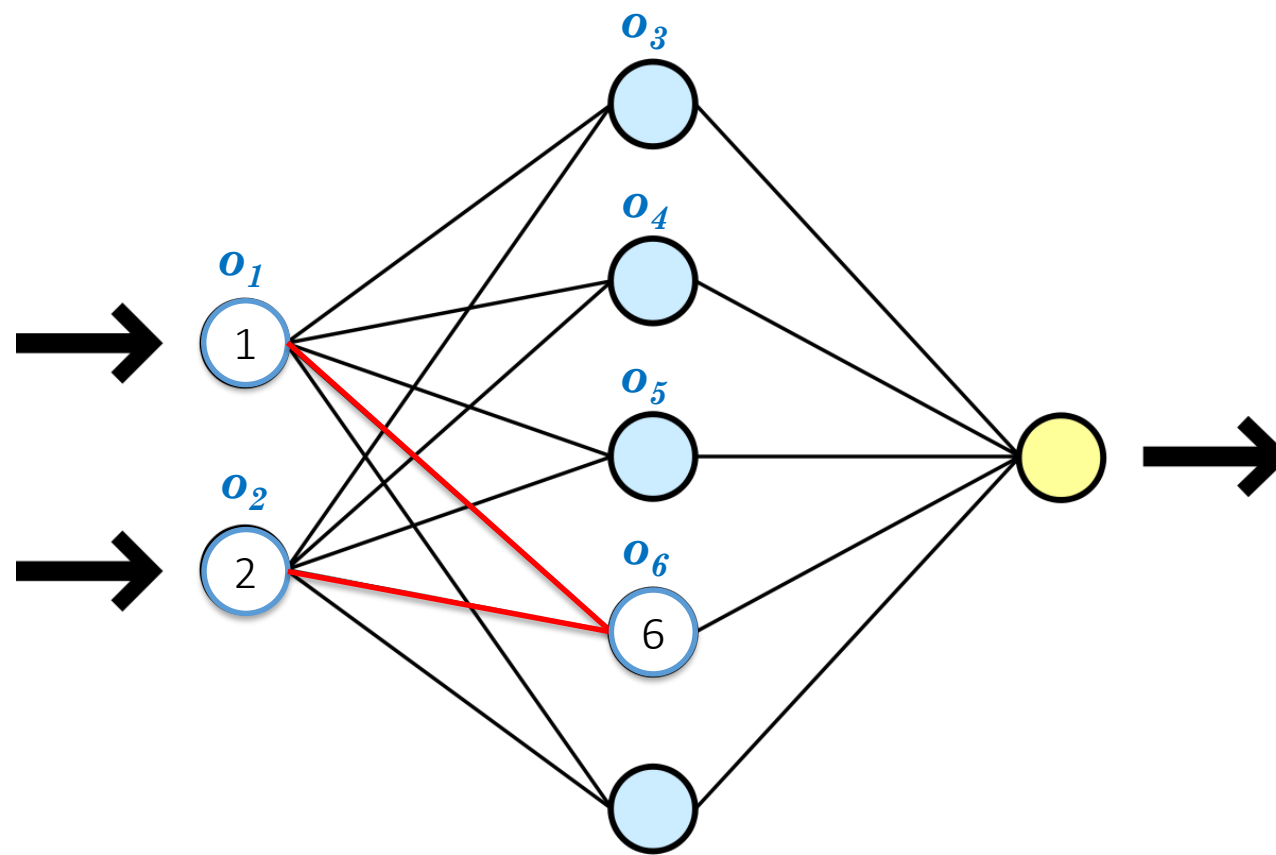
Feed-Forward Phase (Step 2)



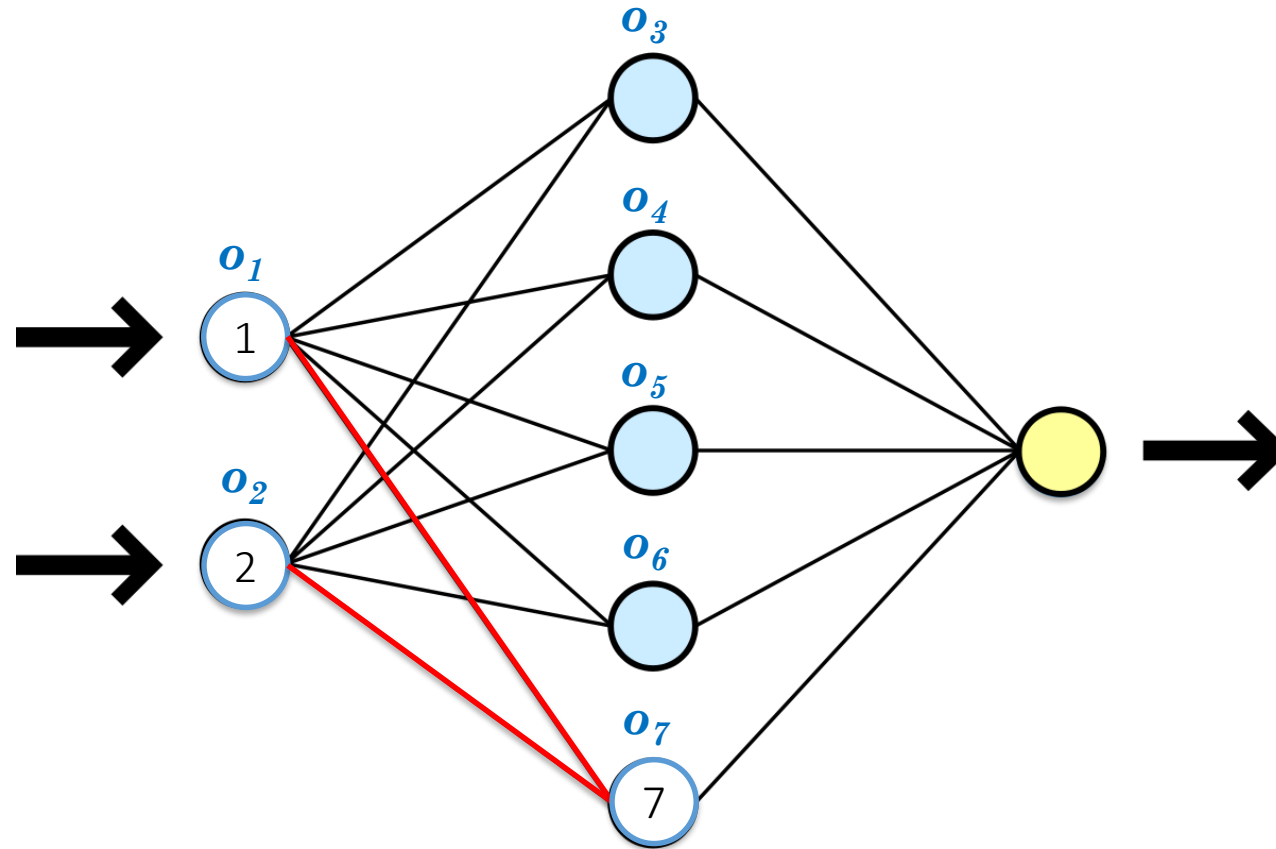
Feed-Forward Phase (Step 2)



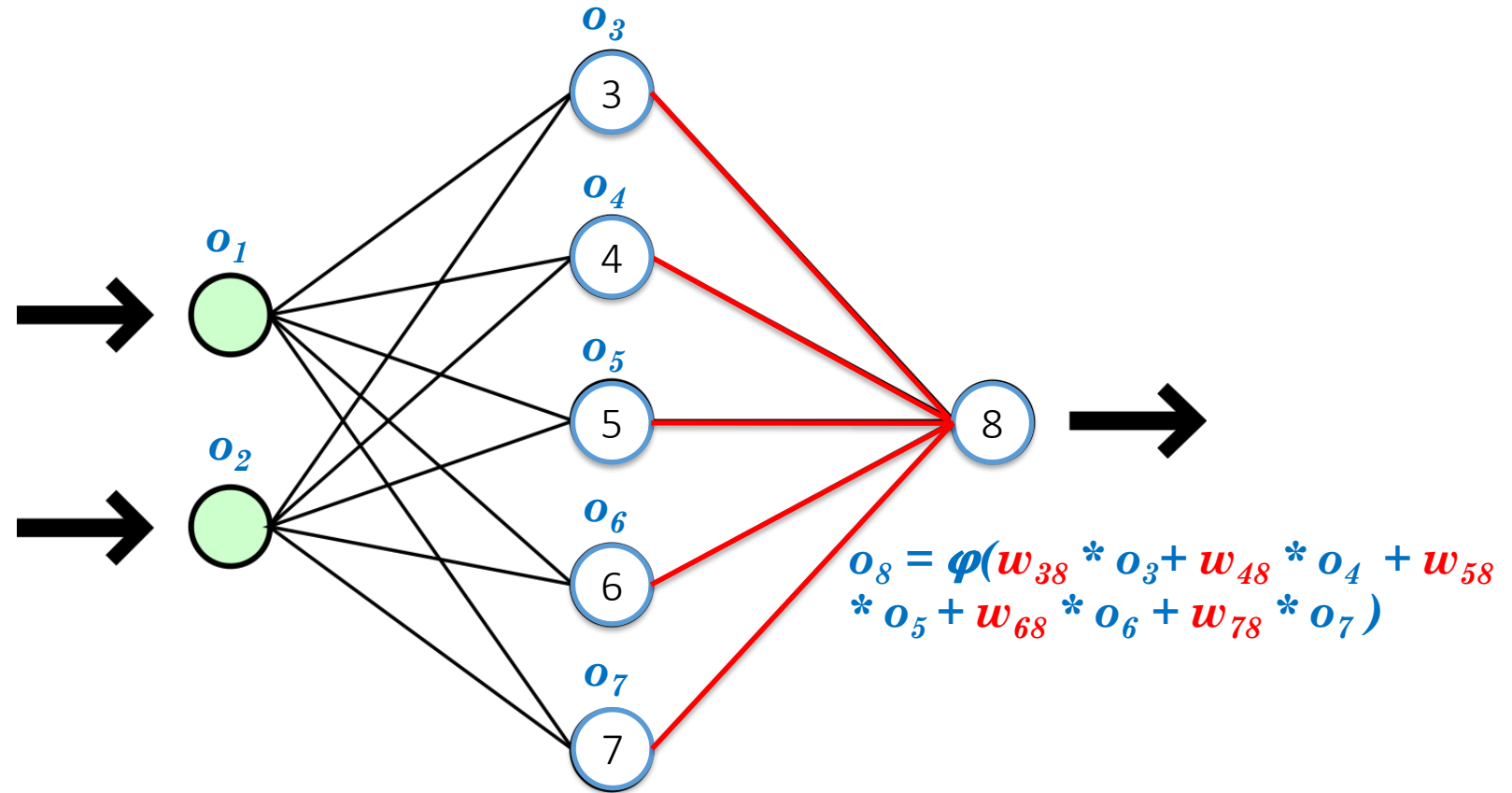
Feed-Forward Phase (Step 2)



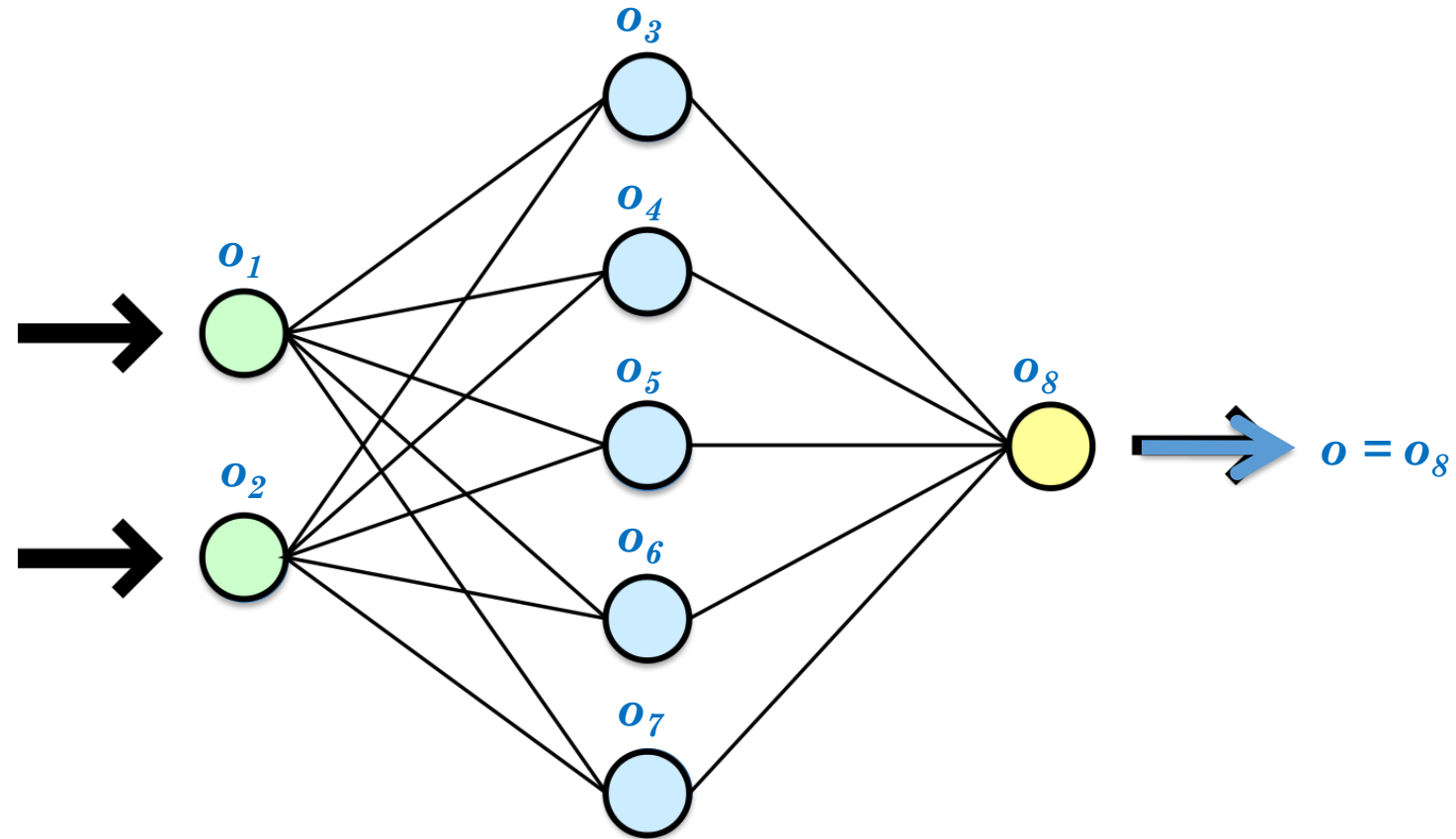
Feed-Forward Phase (Step 2)



Feed-Forward Phase (Step 3)



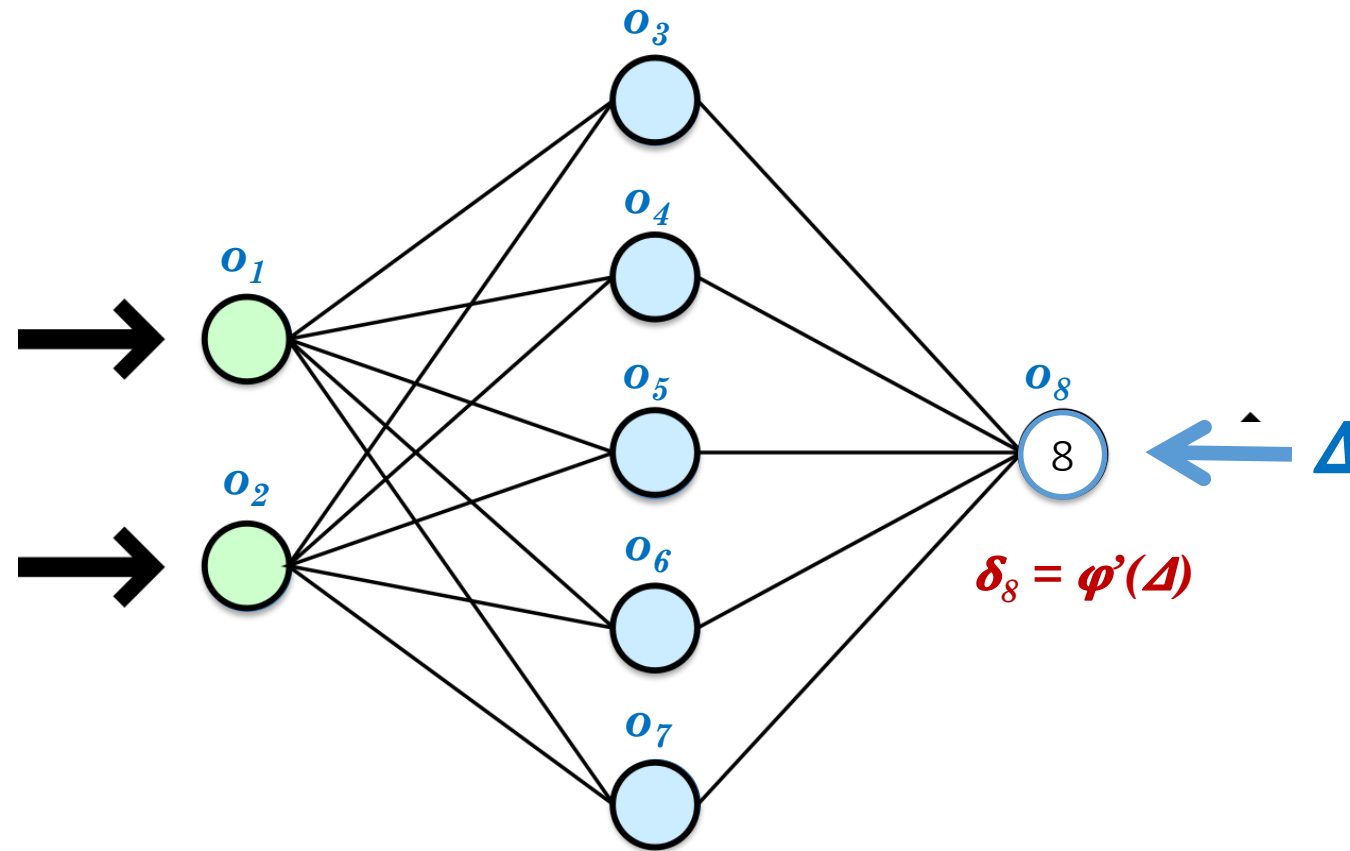
End of *Feed-Forward* Phase



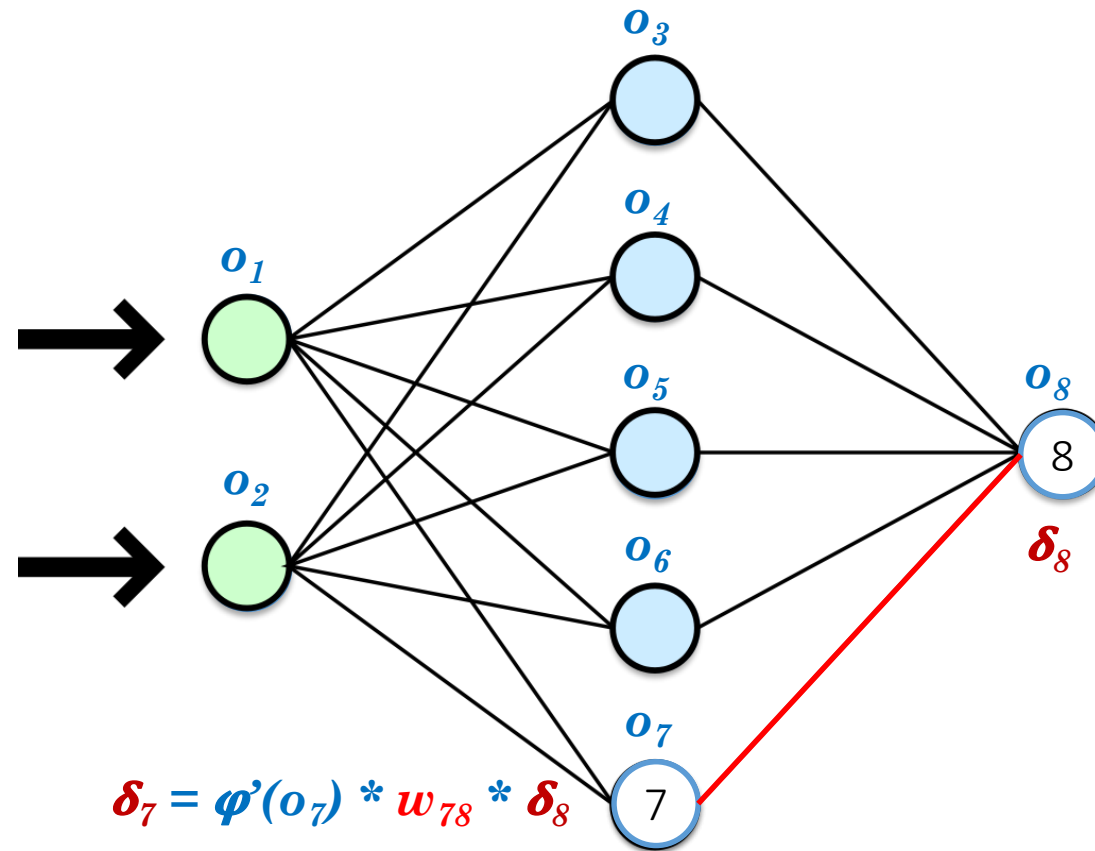
Analyse de l'erreur observée

- o est le résultat obtenu à partir de l'entrée $[v1, v2]$
- Le résultat attendu était t
- On calcule donc $\Delta = |o - t|$

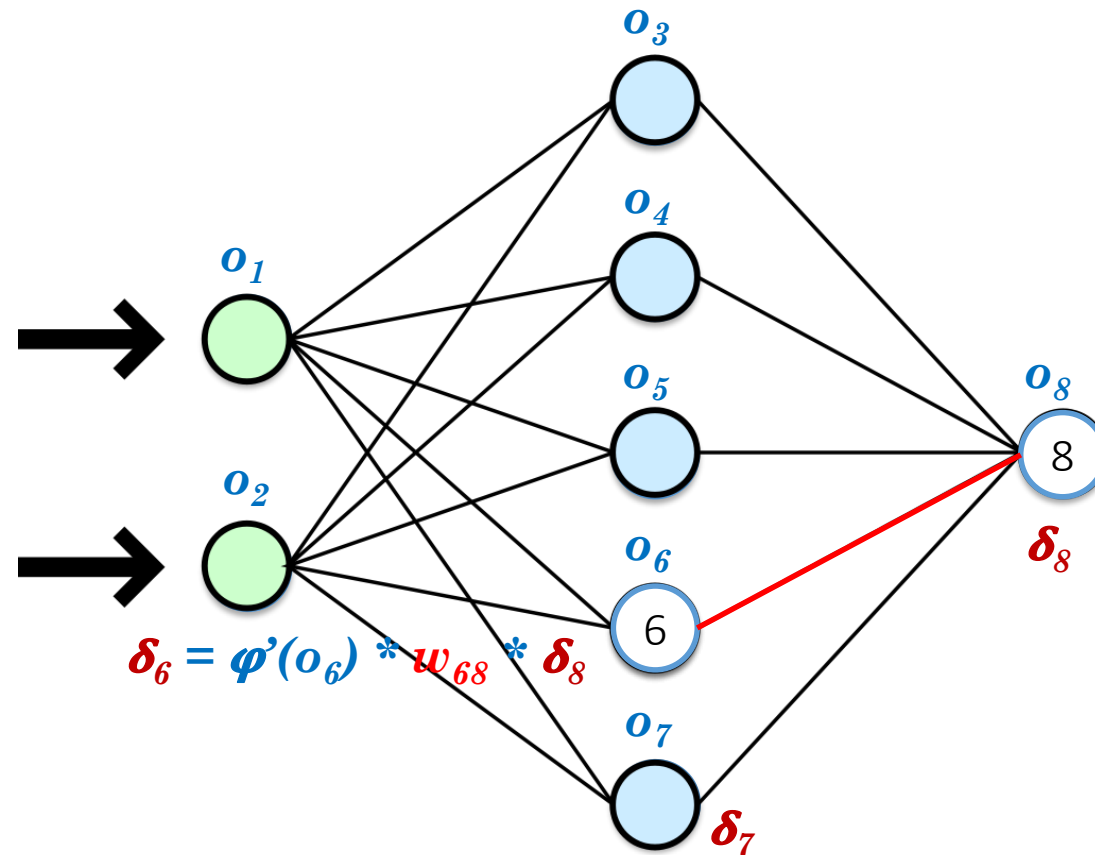
Phase de rétropropagation



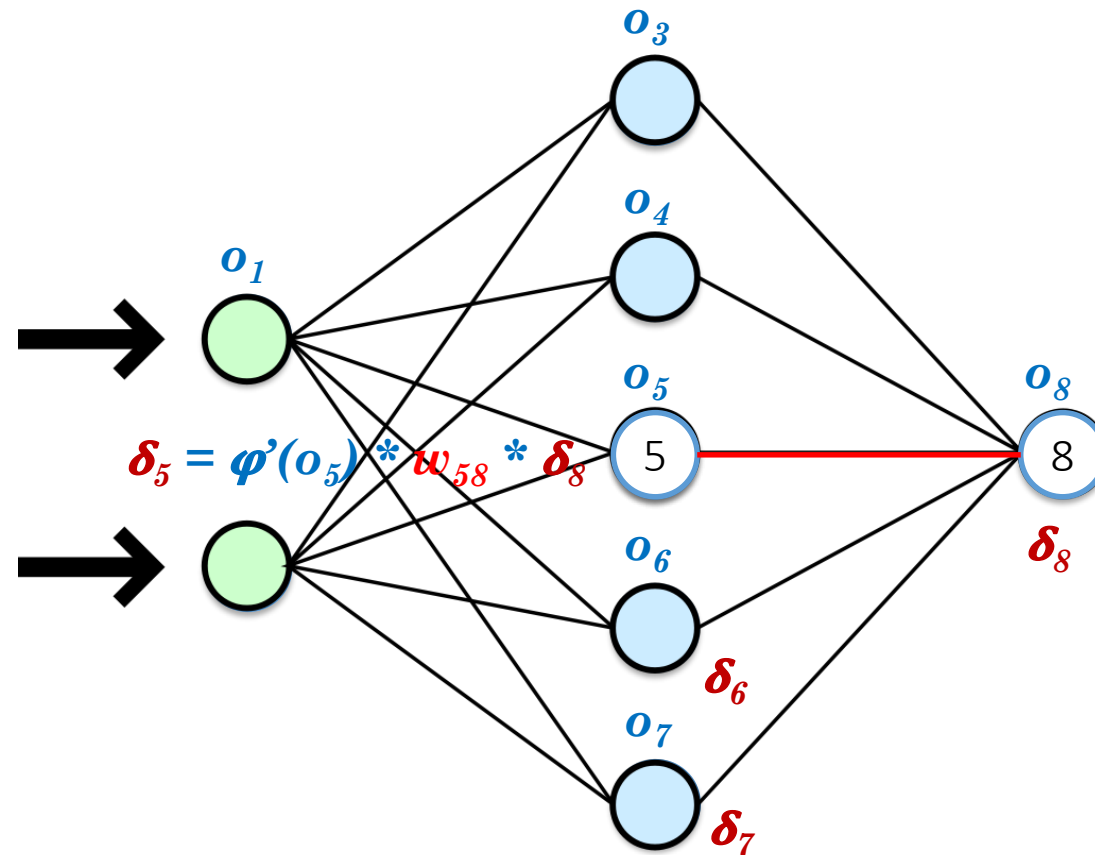
Phase de rétropropagation



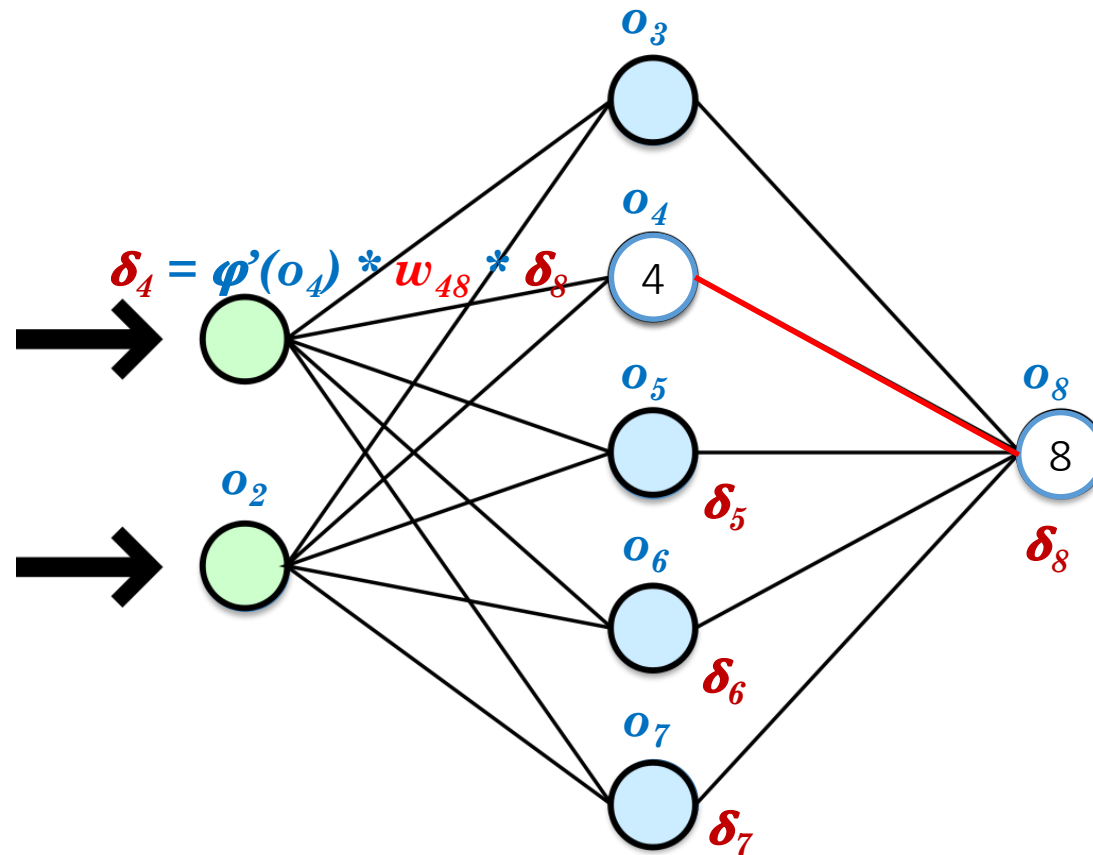
Phase de rétropropagation



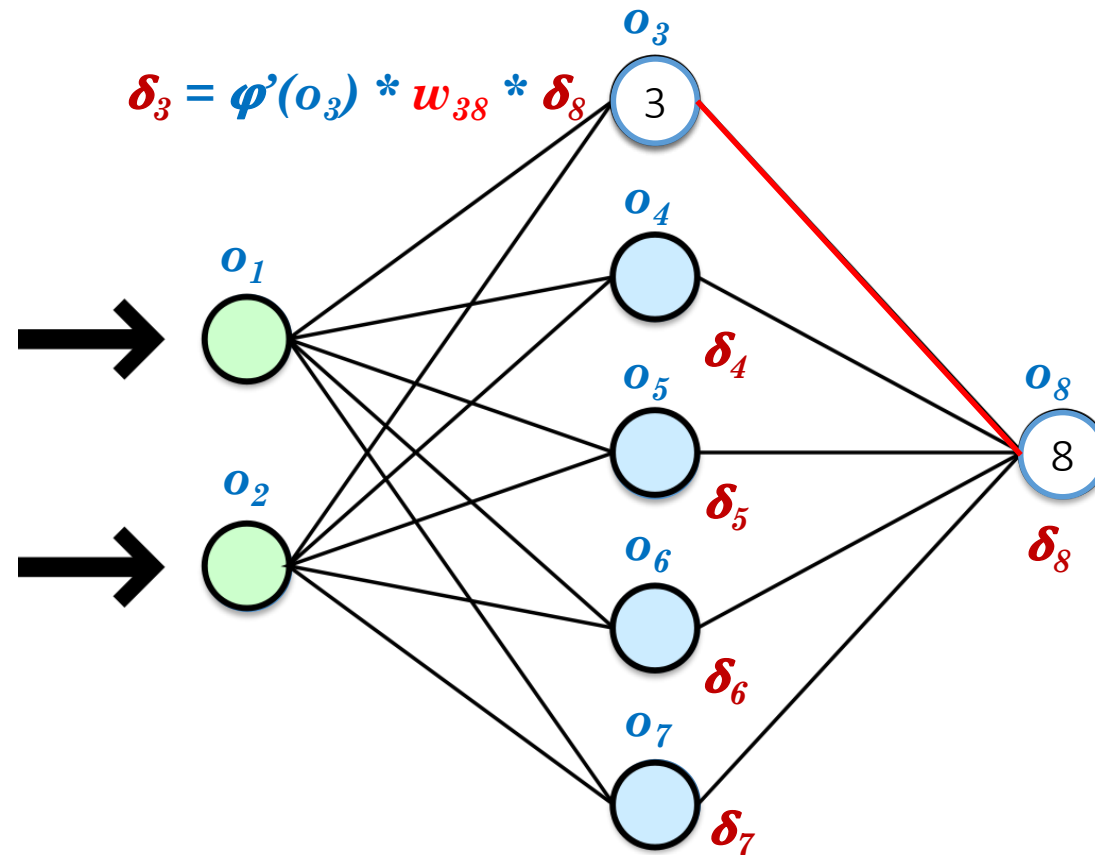
Phase de rétropropagation



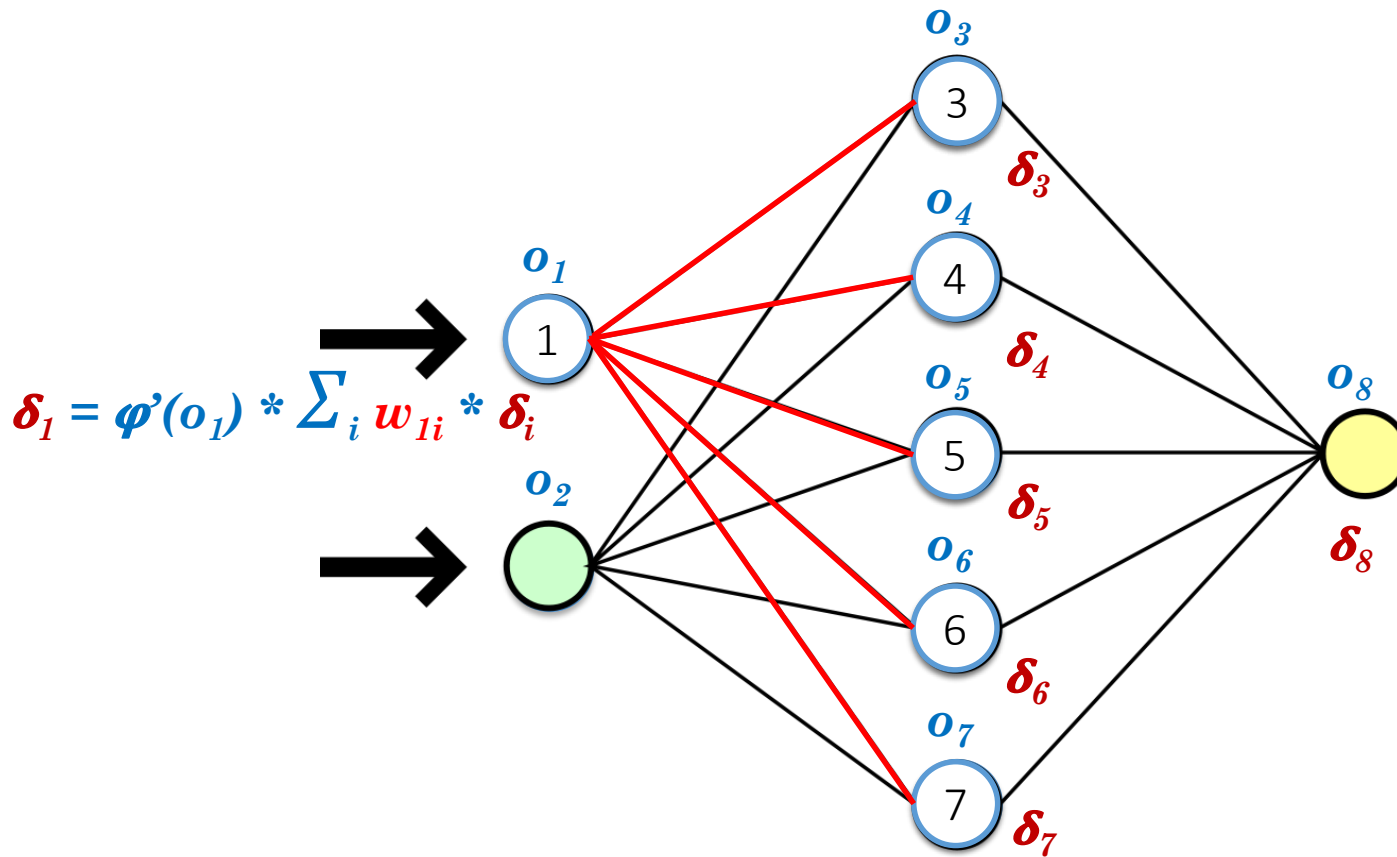
Phase de rétropropagation



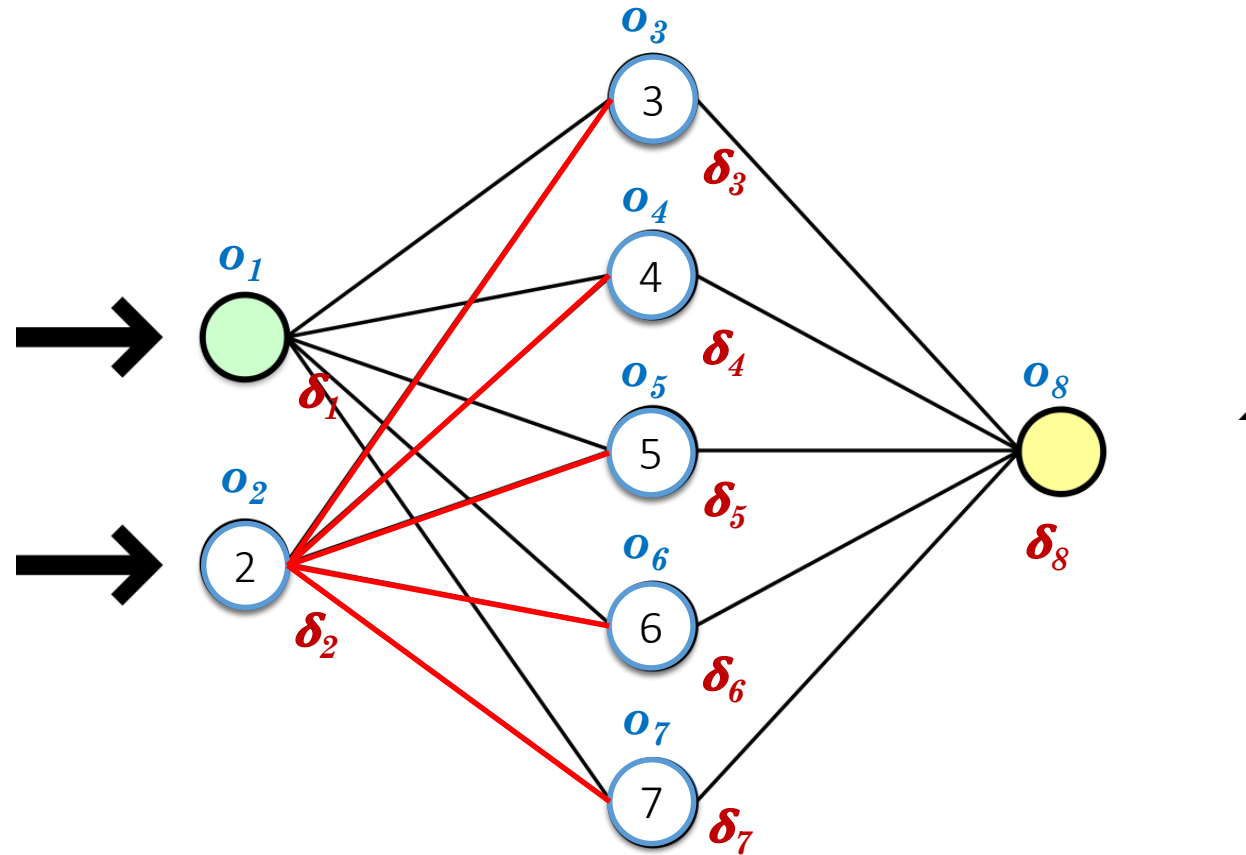
Phase de rétropropagation



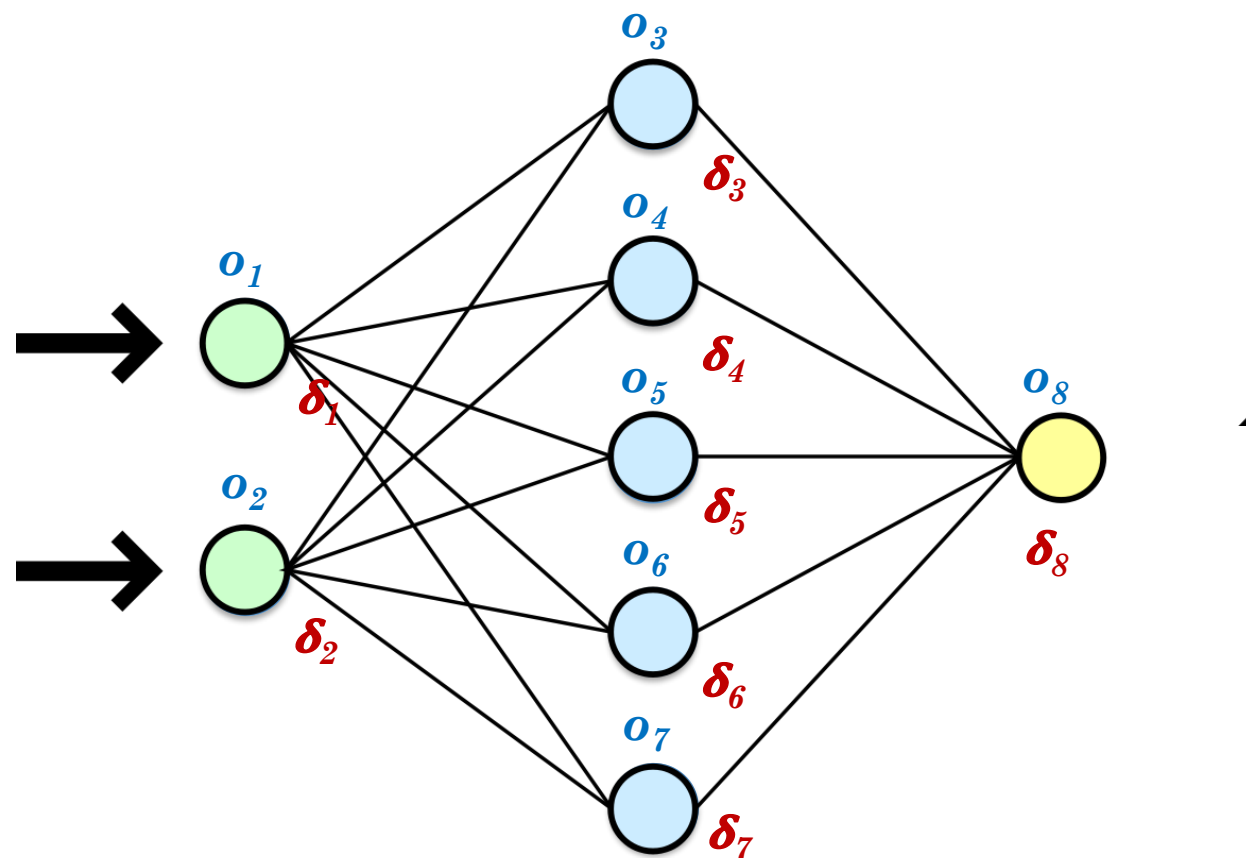
Phase de rétropropagation



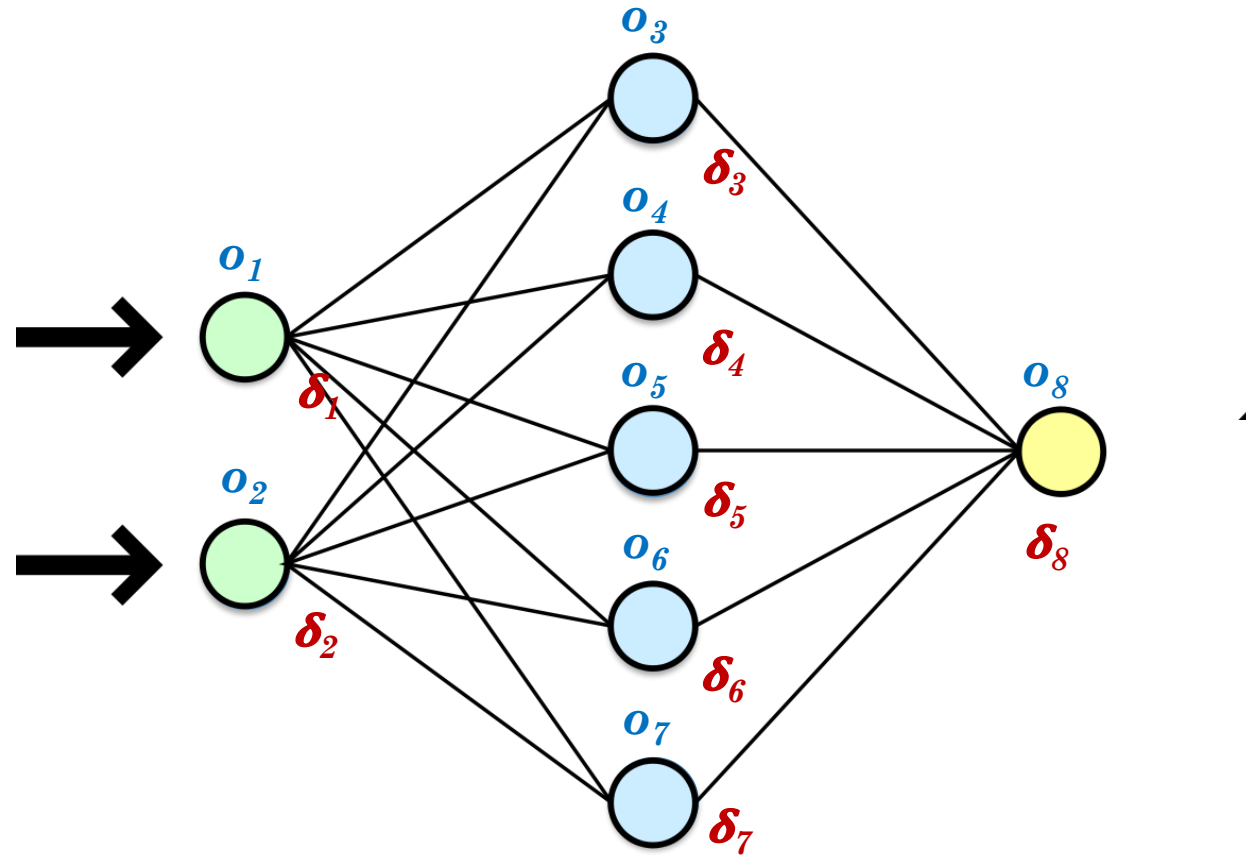
Phase de rétropropagation



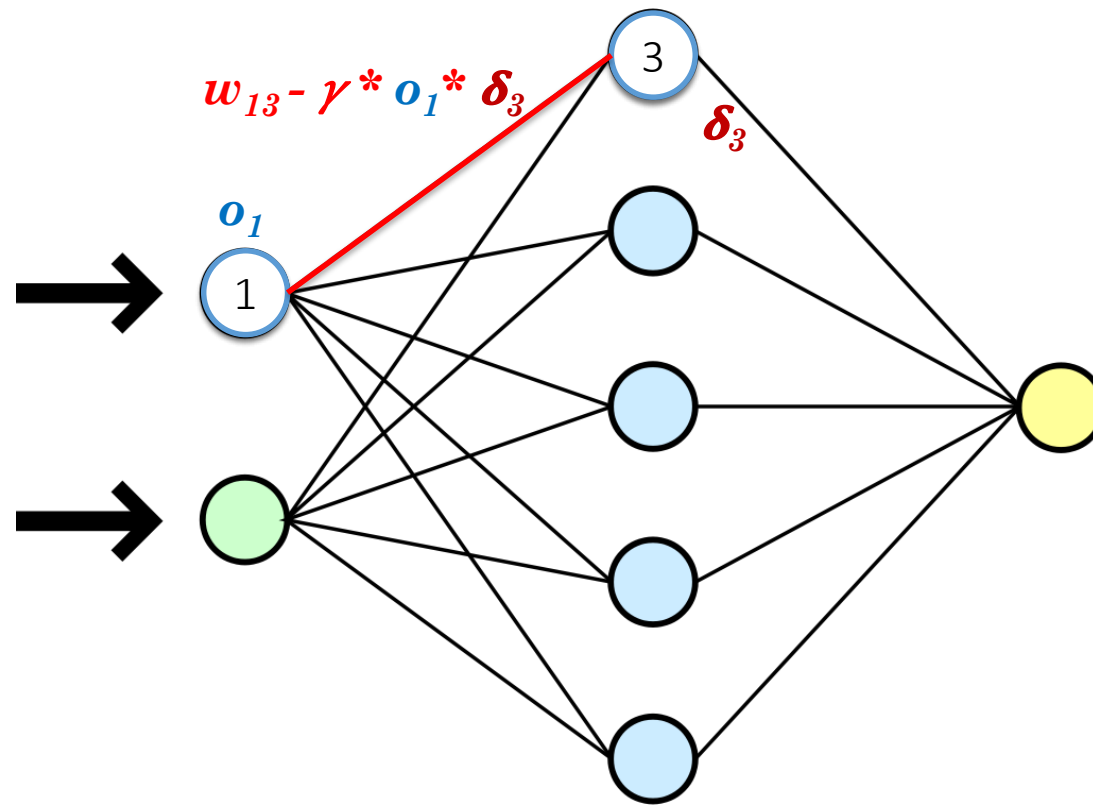
Résultat de la rétropropagation



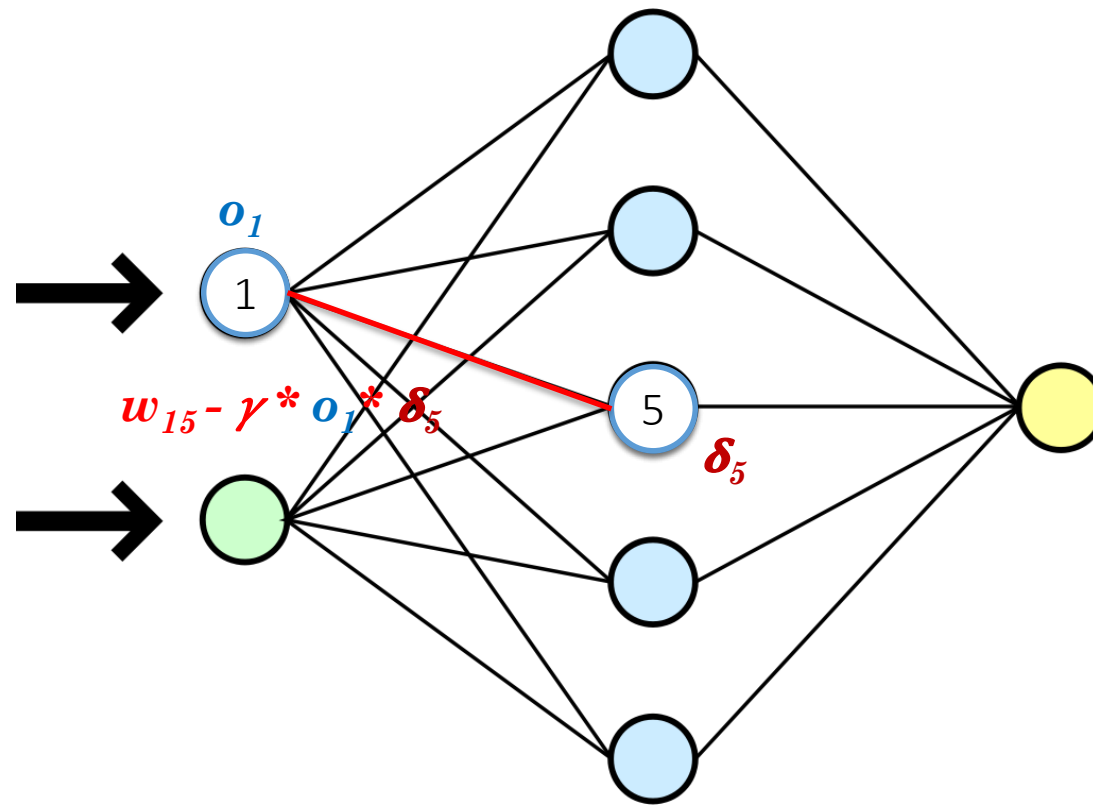
Mise à jour des poids



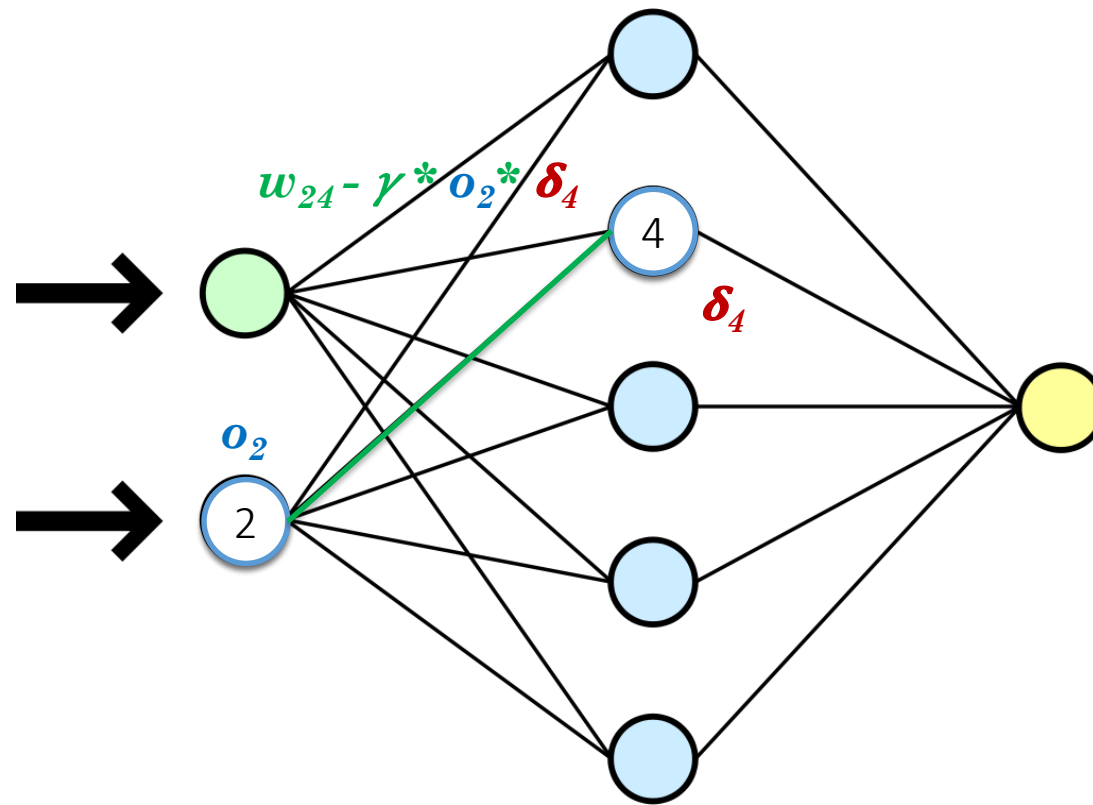
Mise à jour des poids



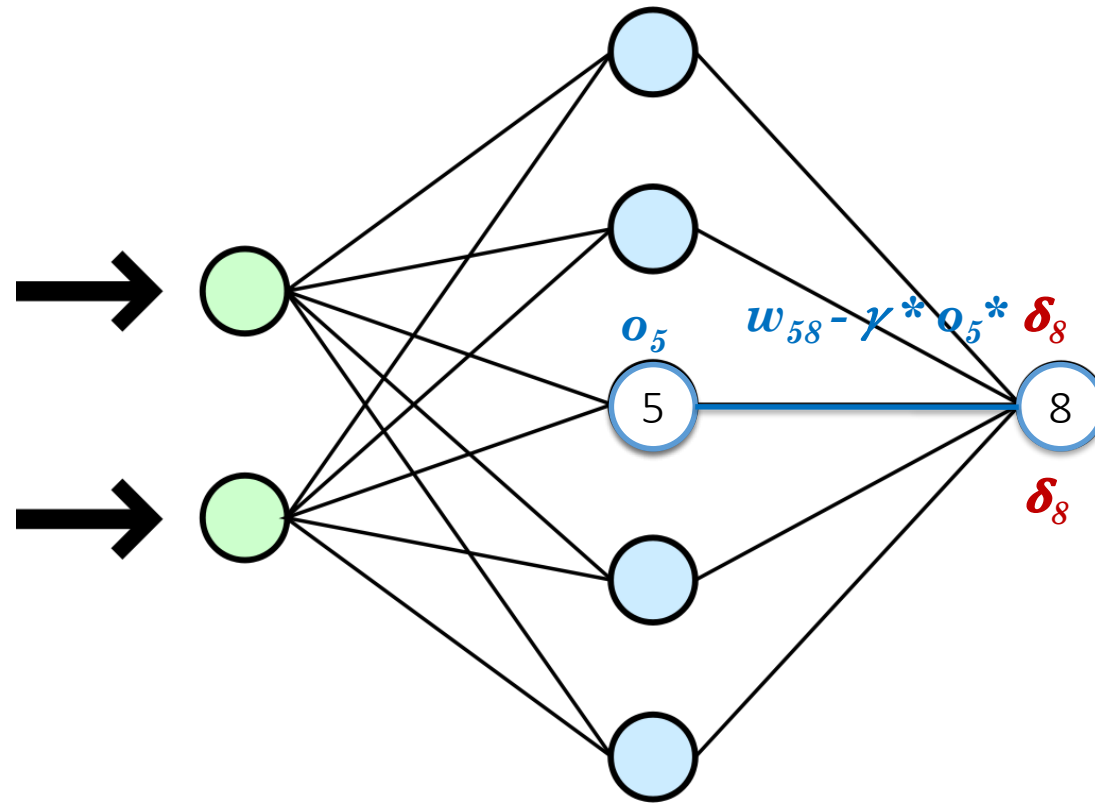
Mise à jour des poids



Mise à jour des poids



Mise à jour des poids



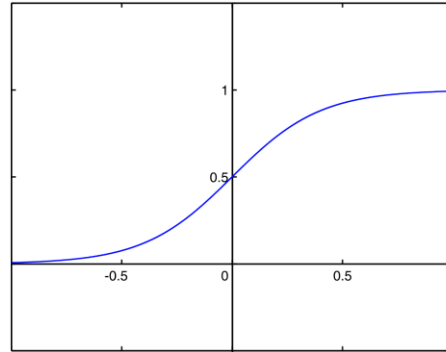
Conclusions intermédiaires

- Paradigme générique de “programmation” de fonctions
- Pas besoin de modèle, uniquement de données

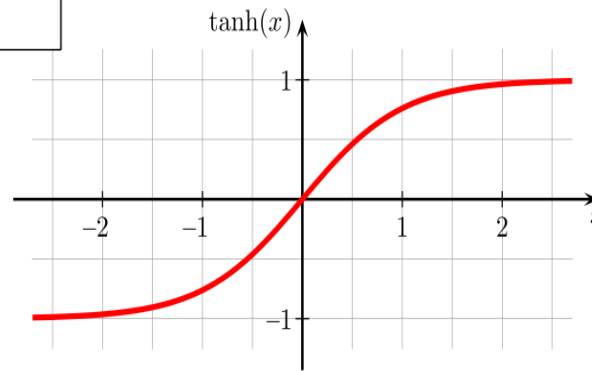
- Généralisation à des architectures de réseaux plus complexes
- Problème des choix des hyper-paramètres du modèle
- Problème du biais des données choisies

Rapide focus sur les fonctions d'activation

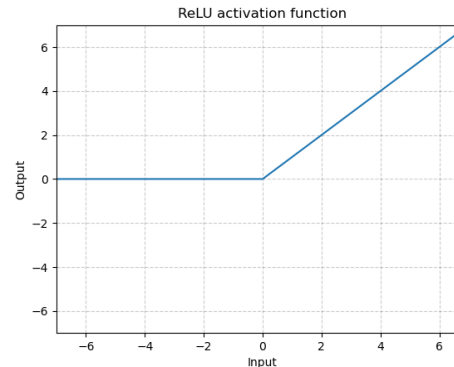
- Sigmoide



- Tangente hyperbolique



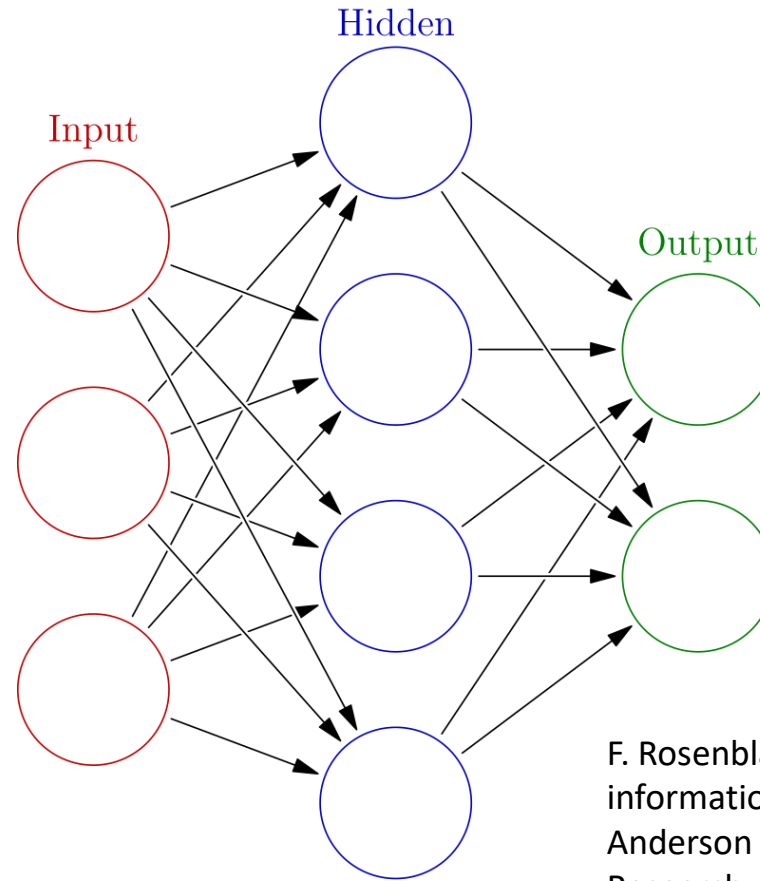
- ReLU (Rectified Linear Unit)



Types de réseaux

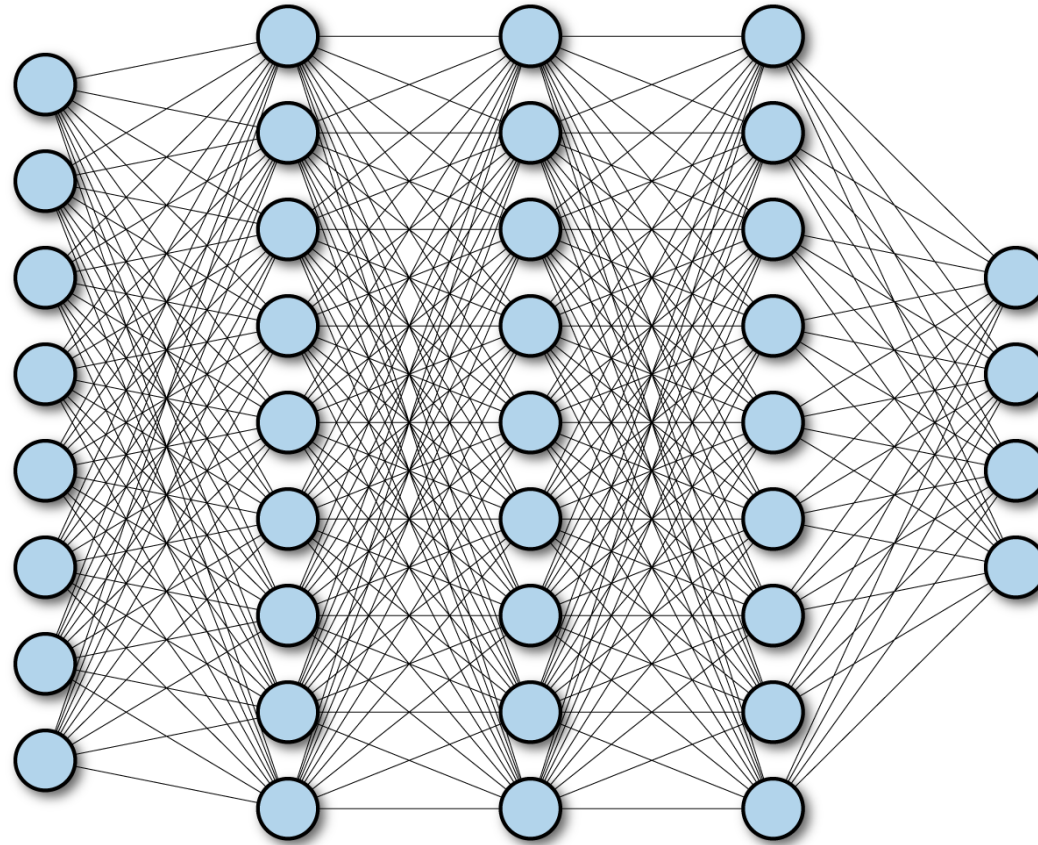
- *Hidden Layer Perceptron*
- *FCNN – Fully Connected Neural Network*
- *CNN – Convolutional Neural Network*
- *RNN – Recurrent Neural Network*
 - *LSTM – Long Short-Term Memory*
 - *GRU – Gated Recurrent Unit*
 - *CRNN – Convolutional Recurrent Neural Network*
- *AutoEncoder*
- *Transformer*
- *GAN – Generative Adversial Network*

Hidden Layer Perceptron



F. Rosenblatt (1958), "The perceptron: a probabilistic model for information storage and organization in the brain", repris dans J.A. Anderson & E. Rosenfeld (1988), Neurocomputing. Foundations of Research, MIT Press

FCNN – Fully Connected Neural Network

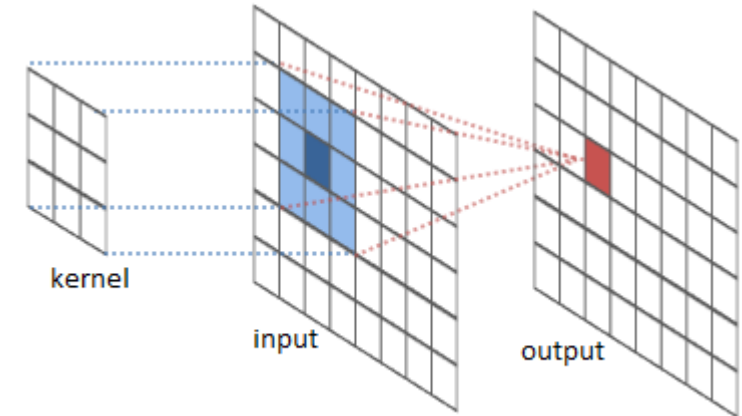


https://www.oreilly.com/library/view/tensorflow-for-deep/9781491980446/assets/tfdl_0402.png

CNN – Convolutional Neural Network

- Convolution
- *Pooling* (sous-échantillonnage)
- Aggrégation

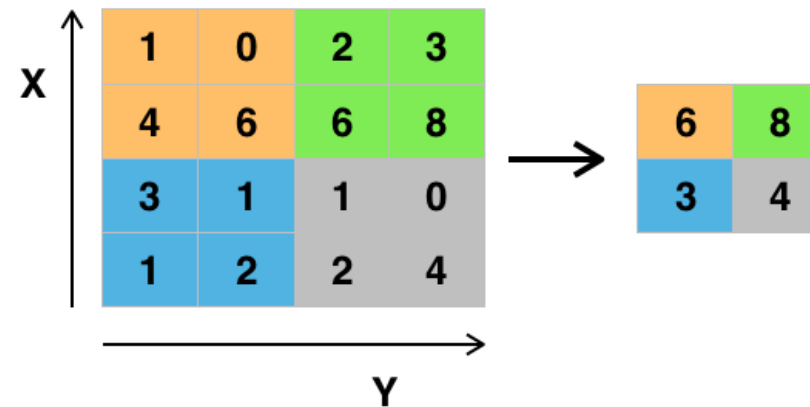
$$(f * g)(t) := \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau.$$



<https://mathinfo.alwaysdata.net/wp-content/uploads/2016/11/RiverTrain-ImageConvDiagram.png>

Max Pooling

Single depth slice



By Aphex34 - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=45673581>

CNN – Convolutional Neural Network

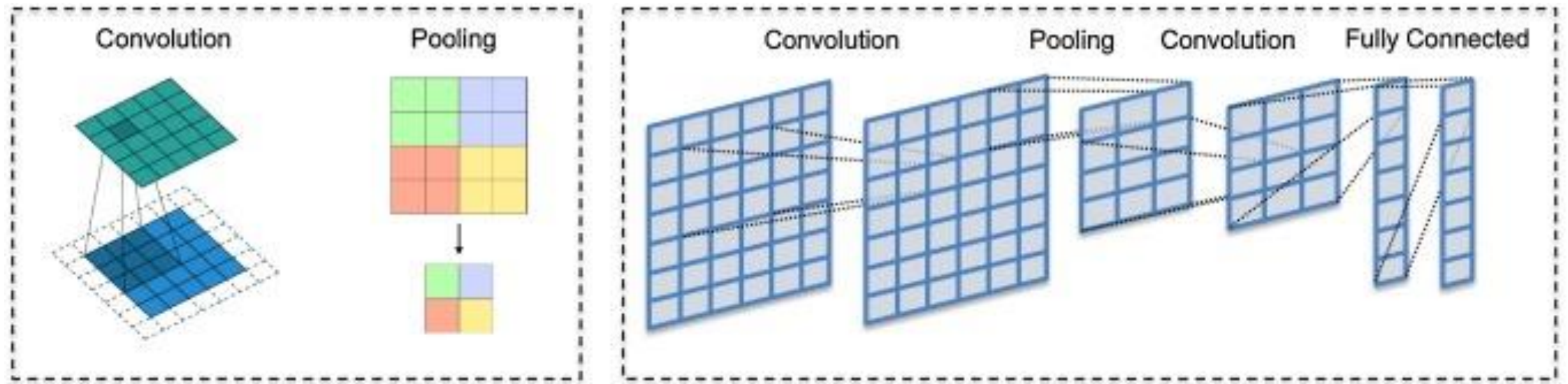
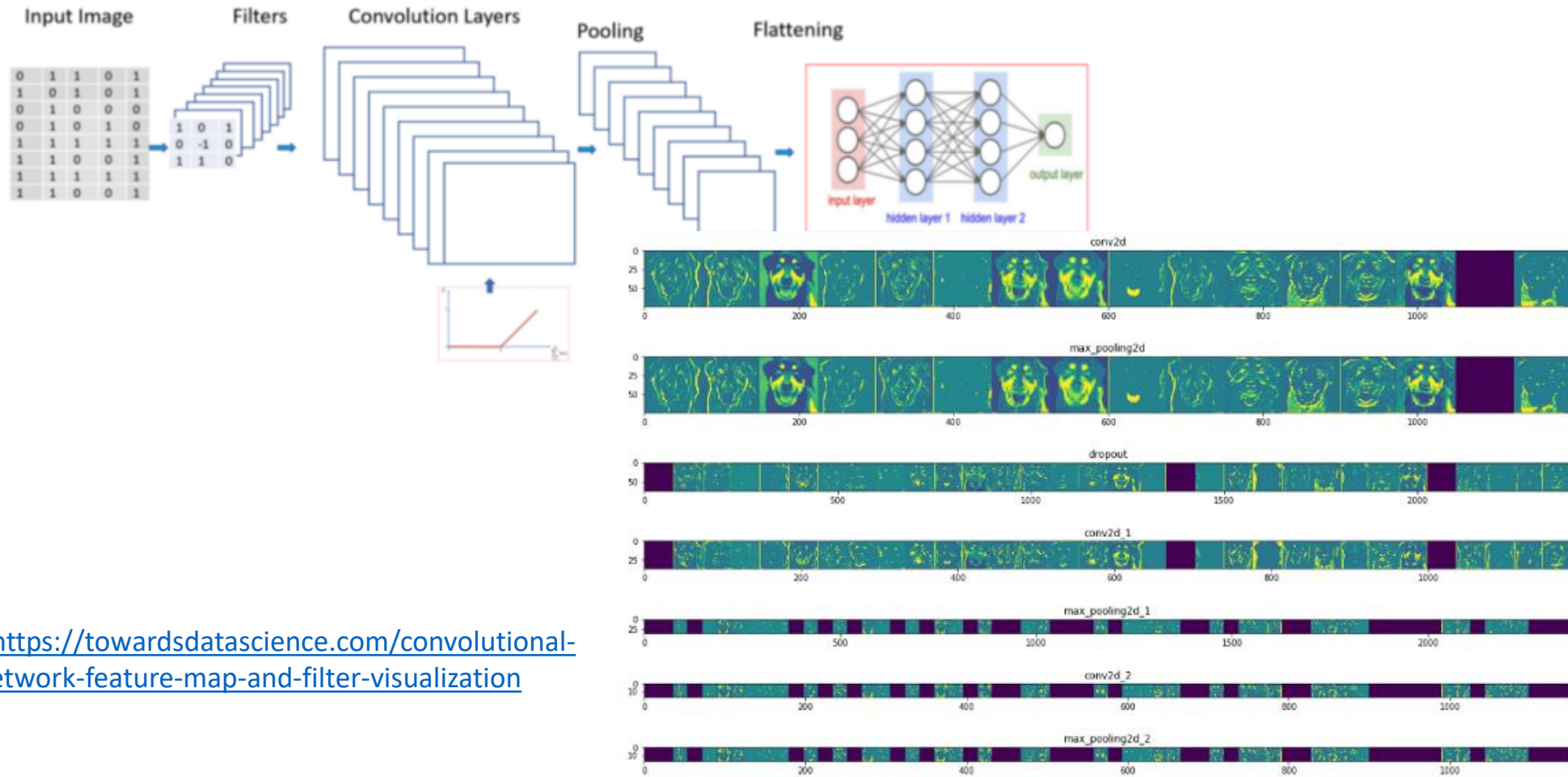


Figure 6 from Andreas Maier, Christopher Syben, Tobias Lasser, Christian Riess. "A gentle introduction to deep learning in medical image processing". Zeitschrift für Medizinische Physik Volume 29, Issue 2, May 2019, Pages 86-101.

<https://www.sciencedirect.com/science/article/pii/S093938891830120X#fig0010>

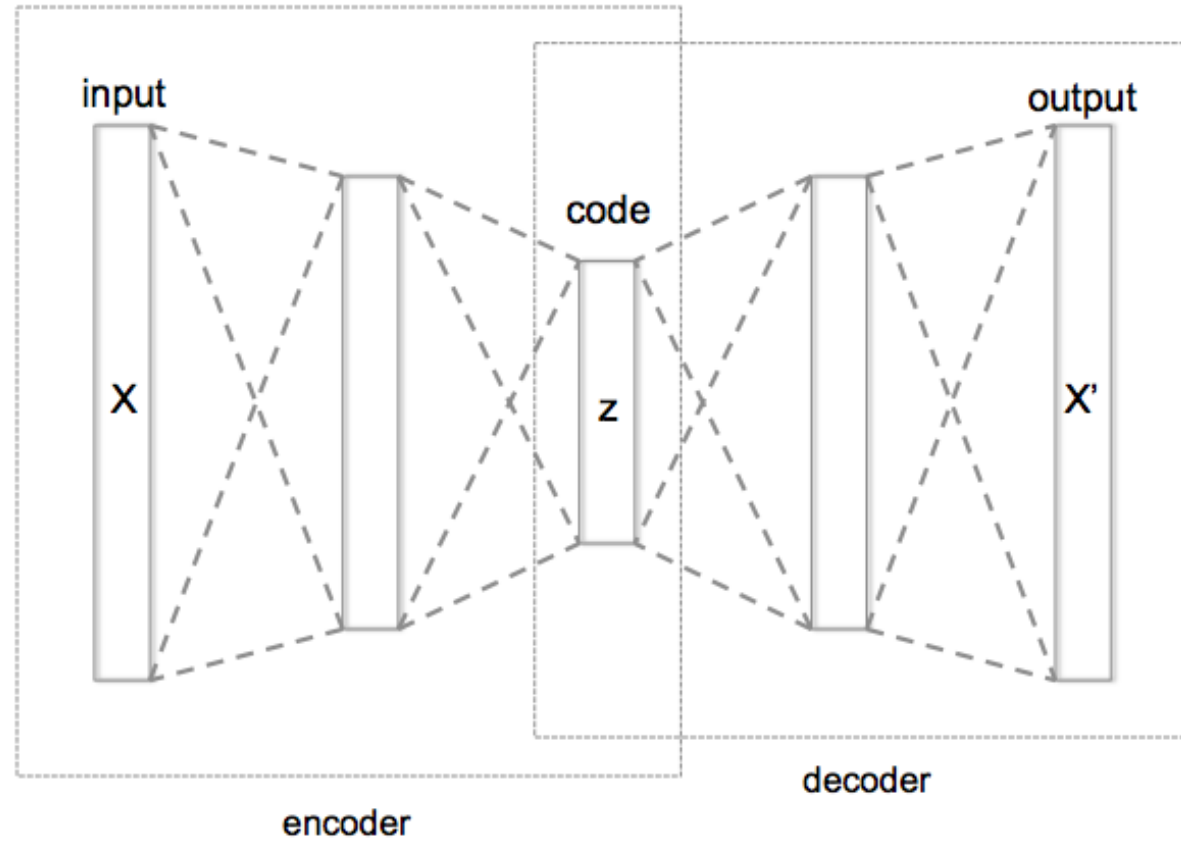
CNN end-to-end



Source : <https://towardsdatascience.com/convolutional-neural-network-feature-map-and-filter-visualization>

Autoencodeur

Espace de représentation latent
« plongement » (*embedding*)



Par Chervinskii — Travail personnel, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=4555552>

Autoencoder Types

Regularization

- Sparse (SAE)
- Denoising
- Variational (VAE)

Exemple d'utilisation des plongements :

reine = $(\alpha_0, \alpha_1 \dots \alpha_n)$

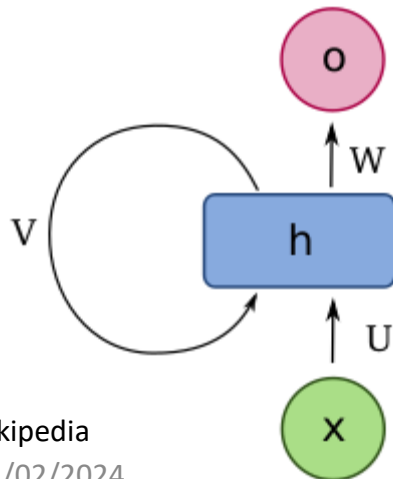
femme = $(\beta_0, \beta_1 \dots \beta_n)$

bâtiment = $(\gamma_0, \gamma_1 \dots \gamma_n)$

reine - femme + bâtiment = château

Recurrent Neural Networks

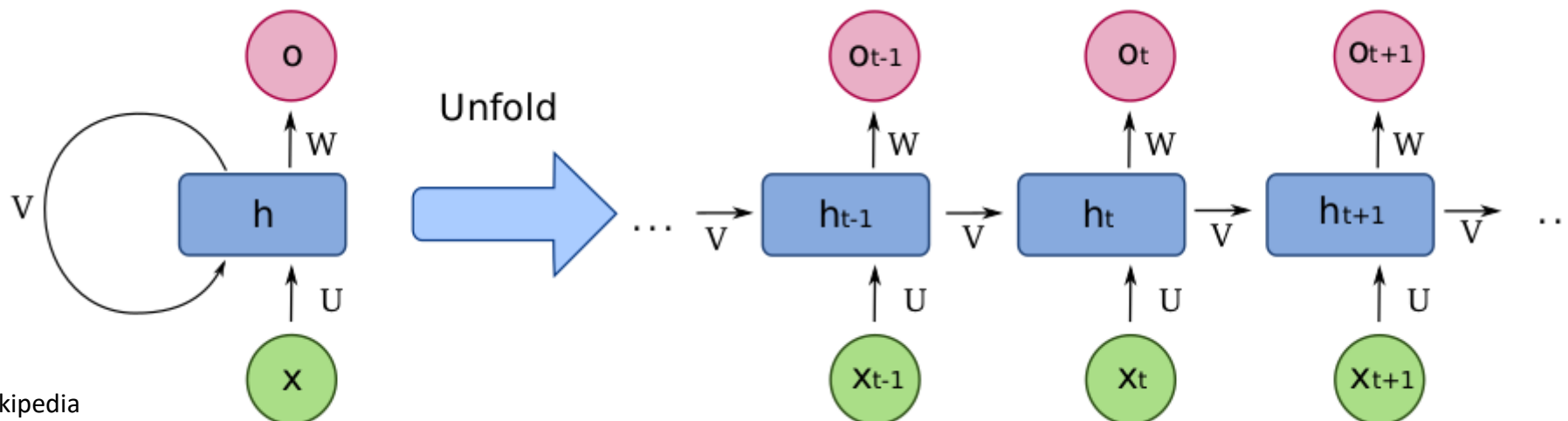
- Introduction d'un temps de latence/mémoire dans la construction des réseaux (RNN : Recurrent Neural Networks)
- RNN initialement conçus pour traiter des **séries temporelles**
- Principe d'une mémoire (*hidden state*) h_t qui dépend de l'entrée x_t au temps t et de la mémoire h_{t-1} au temps $t-1$



$$h_0 = [0, 0 \dots 0]$$
$$h_t = V h_{t-1} + U x_t$$
$$o_t = W h_t$$

Recurrent Neural Networks

- Introduction d'un temps de latence/mémoire dans la construction des réseaux (RNN : Recurrent Neural Networks)
- RNN initialement conçus pour traiter des séries temporelles
- Principe d'une mémoire (*hidden state*) h_t qui dépend de l'entrée x_t au temps t et de la mémoire h_{t-1} au temps $t-1$



Source: Wikipedia

11/02/2024

Bart.Lamiroy@univ-reims.fr

65

Extensions des RNN

- **LSTM**

Long Short-Term Memory, S. Hochreiter, J. Schmidhuber, Neural Computation
Volume 9, Issue 8, November 15, 1997, p.1735-1780,
<https://doi.org/10.1162/neco.1997.9.8.1735>

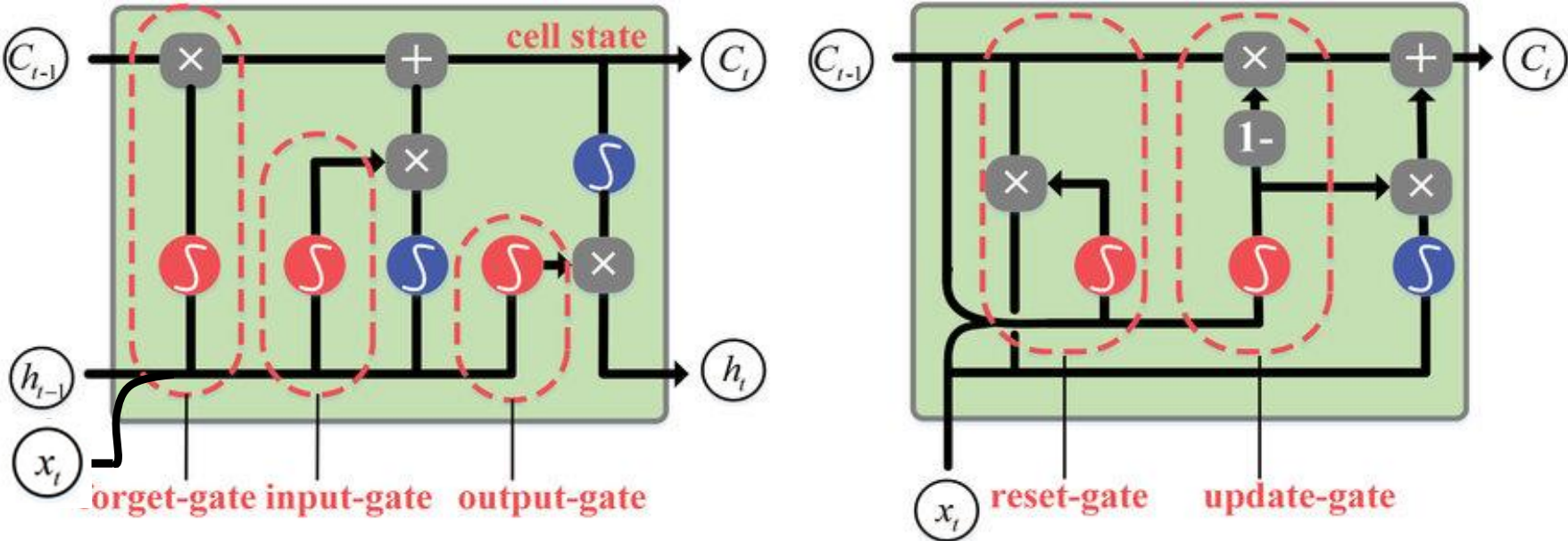
- **GRU – Gated Recurrent Unit**

Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation, K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio, 2014. arXiv:1406.1078

- **CRNN –Convolutional Recurrent Neural Networks**

<https://stanford.edu/~shervine/teaching/cs-230/cheatsheet-recurrent-neural-networks>

LSTM vs. GRU



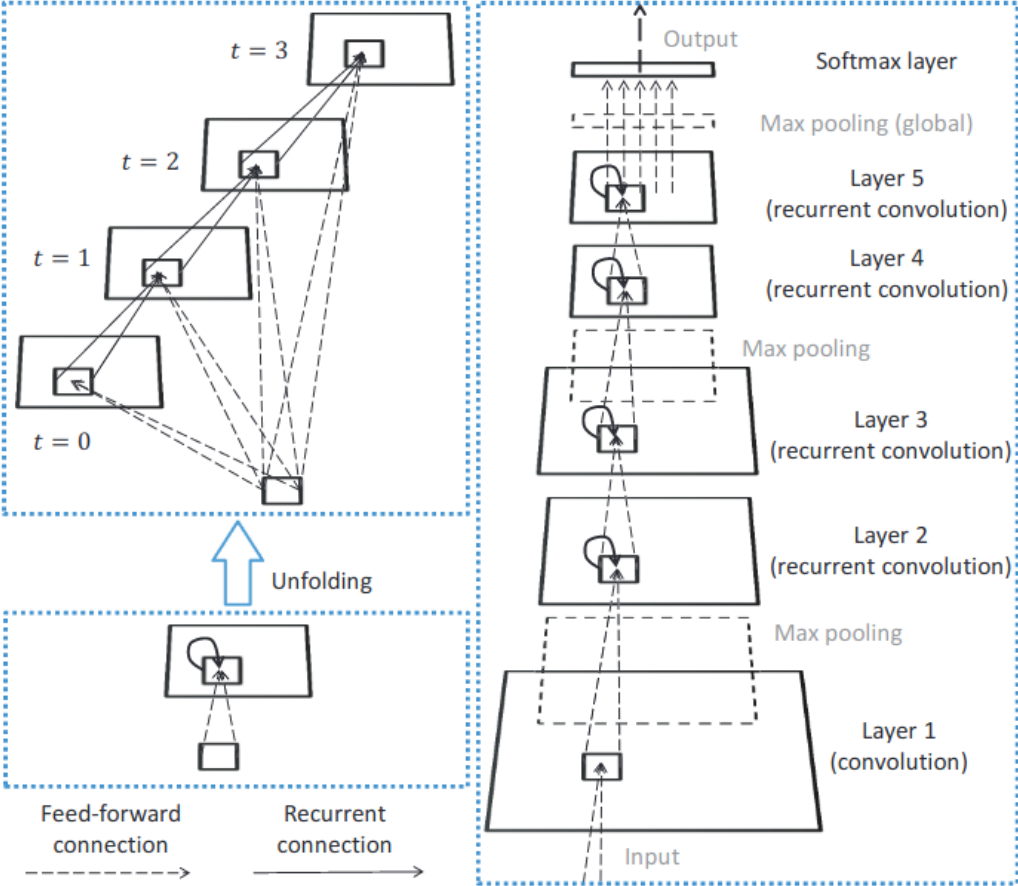
(a) LSTM simplified structure

(b) GRU simplified structure

	sigmoid		pointwise multiplication		output = 1-input		previous cell state		previous hidden state		vector concatenation
	tanh		pointwise addition		current cell state		current hidden state		input		

Zhao, Junhui & Nie, Yiwen & Shanjin, Ni & Sun, Xiaoke. (2020). Traffic Data Imputation and Prediction: An Efficient Realization of Deep Learning. IEEE Access. PP. 1-1. 10.1109/ACCESS.2020.2978530.

Recurrent Convolutional Networks

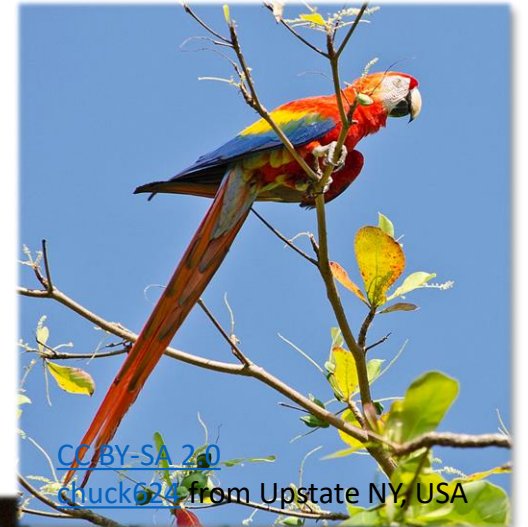


Source : [Recurrent Convolutional Neural Network for Object Recognition](#), M. Liang, X. Hu, CVPR 2015



ChatGPT

Transformers et mécanisme d'attention



Droits appartenant aux propriétaires respectifs

Pourquoi « transformers »

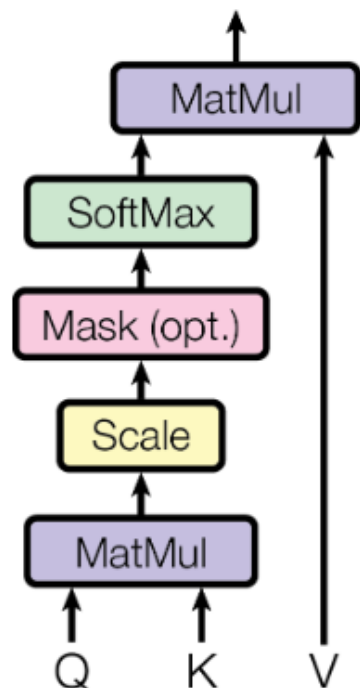
- Transformation de séquences d'entrée en séquences de sortie (*seq2seq*)
- Mécanisme d'auto-régression
 - Principe de régression type série temporelle (prédire le « futur » en observant le « passé »)
 - Enrichissement la séquence initiale (le « passé ») par les prédictions successives à chaque étape.

“Attention is all you Need”, A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I Polosukhin
Advances in Neural Information Processing Systems 30 (NIPS 2017)

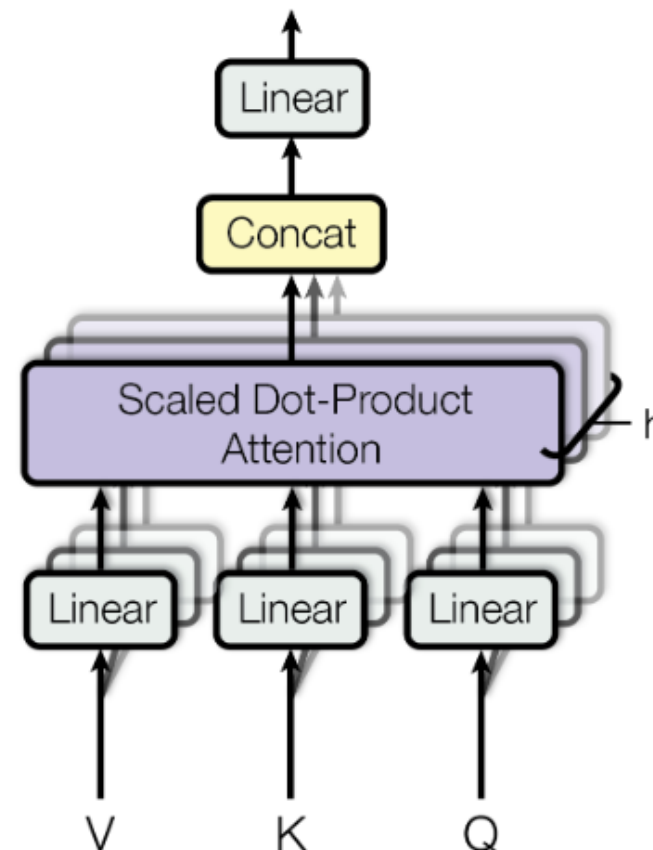
https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

Mécanisme d'attention

Scaled Dot-Product Attention



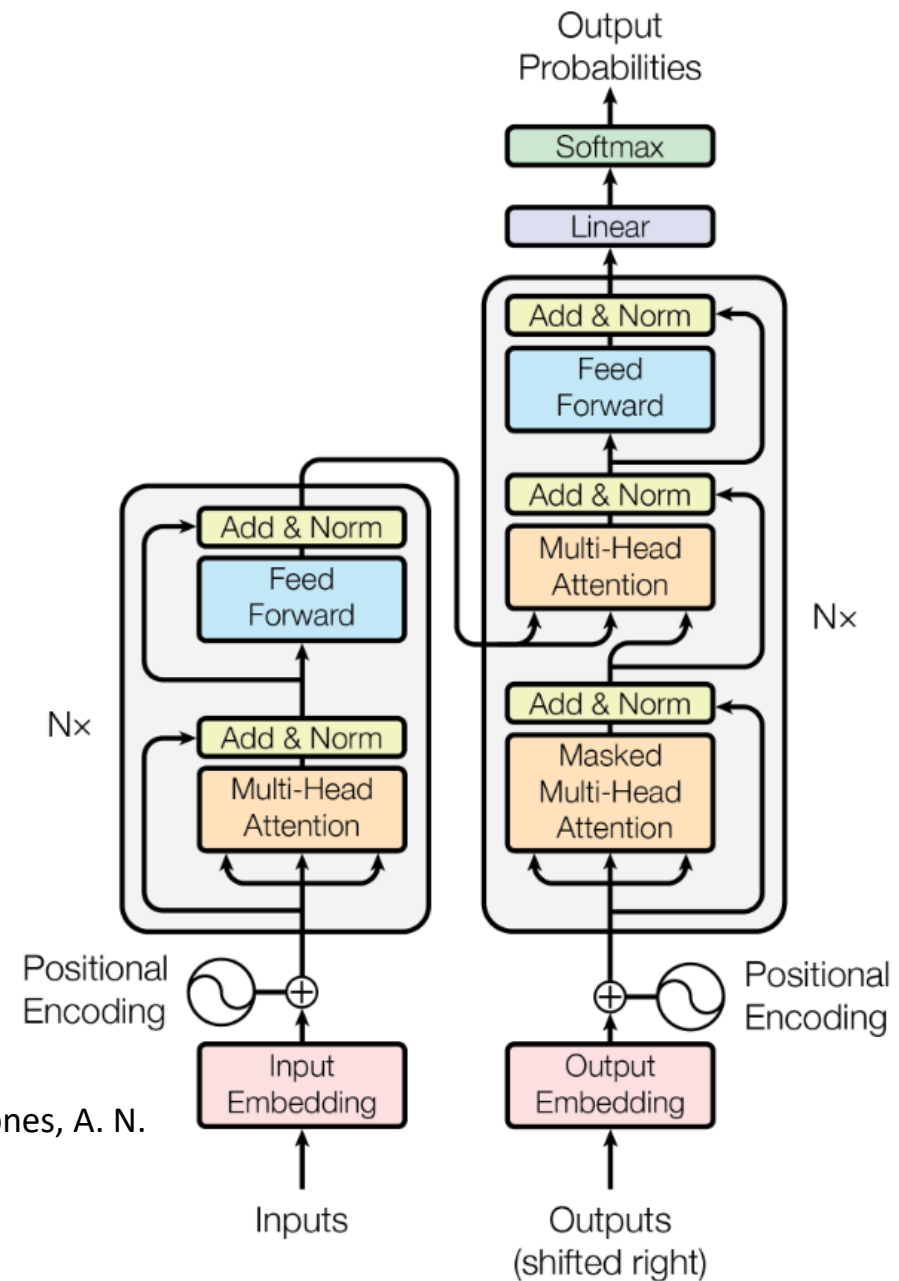
Multi-Head Attention



“Attention is all you Need”, A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I Polosukhin

Advances in Neural Information Processing Systems 30 (NIPS 2017)

Transformer



“Attention is all you Need”, A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I Polosukhin
Advances in Neural Information Processing Systems 30 (NIPS 2017)

Résultat : LLM

- LLM = Large Language Model
- Apprentissage auto-supervisé
- Principe de base : seq2seq autorégressif
 - Entrée = séquence
 - Sortie = le mot le plus probable qui complète la séquence
- Perroquet stochastique

« *On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?* 🦜 »

E. M. Bender, T. Gebru, A. McMillan-Major, M. Mitchell, FAccT '21: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, March 2021, Pages 610–623

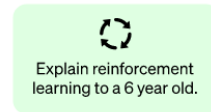
<https://dl.acm.org/doi/10.1145/3442188.3445922>

ChatGPT : Entraînement humain

Step 1

Collect demonstration data and train a supervised policy.

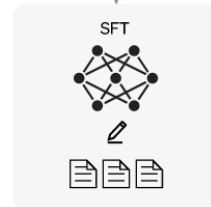
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



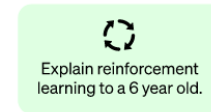
This data is used to fine-tune GPT-3.5 with supervised learning.



Step 2

Collect comparison data and train a reward model.

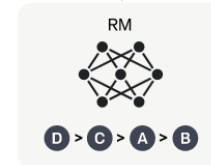
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

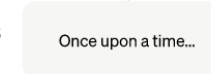
A new prompt is sampled from the dataset.



The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



Source : https://images.openai.com/blob/cf717bdb-0c8c-428a-b82b-3c3add87a600/ChatGPT_Diagram.svg?width=10&height=10&quality=50

Métrique chinchilla

$$\begin{cases} C = C_0 N D \\ L = \frac{A}{N^\alpha} + \frac{B}{D^\beta} + L_0 \end{cases}$$

where the variables are

- C is the cost of training the model, in **FLOPs**.
- N is the number of parameters in the model.
- D is the number of tokens in the training set.
- L is the average negative log-likelihood loss per token (**nats/token**), achieved by the trained LLM on the test dataset.
 - L_0 represents the loss of an ideal generative process on the test data
 - $\frac{A}{N^\alpha}$ captures the fact that a Transformer language model with N parameters underperforms the ideal generative process
 - $\frac{B}{D^\beta}$ captures the fact that the model trained on D tokens underperforms the ideal generative process

A retenir (pour des LLM)

On peut atteindre la même performance (**L constant**) sur un modèle plus complexe (**augmenter N d'un facteur α**) à condition d'augmenter la quantité de données d'entraînement de la même proportion (**augmenter D d'un facteur α**) et d'augmenter le coût de calcul façon quadratique (**augmenter C d'un facteur α^2**).

Training Compute-Optimal Large Language Models
Jordan Hoffmann et al. DeepMind
<https://arxiv.org/pdf/2203.15556.pdf>

On pourrait continuer longtemps ...

- *Hallucinations as a feature*

“One of the sort of non-obvious things is that a lot of value from these systems is heavily related to the fact that they do hallucinate. If you want to look something up in a database, we already have good stuff for that.”

— Sam Altman, CEO OpenAI, Dreamforce 2023, San Francisco

- Modèles multi-langues
- Modèles multi-modaux
- Modèles par diffusion

...

... très longtemps

- L'important est-ce l'algorithme ou la data ?
- Comment évaluer les algorithmes
- Comment corriger les algorithmes
- ...



Source : <http://www.signwithrobert.com/>

Questions ?