

Faire de la science à partir des données ?

L'avènement du deep learning

Journée Science Ouverte CNRS 2022 - Mercredi 30 novembre 2022

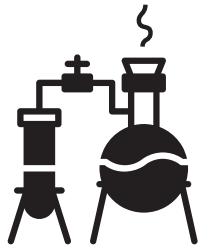
La science ouverte et les données de la recherche

Jean-Luc.Parouty@cnrs.fr



Scientific paradigms

1st paradigm



Experimental science

2nd paradigm

$$i\hbar \frac{d}{dt} |\Psi(t)\rangle = \hat{H} |\Psi(t)\rangle$$

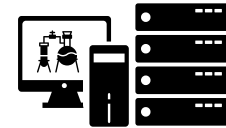
$$\nabla \times H = J + \frac{\partial D}{\partial t}$$

$$F = G \cdot \frac{m_1 \cdot m_2}{r^2}$$

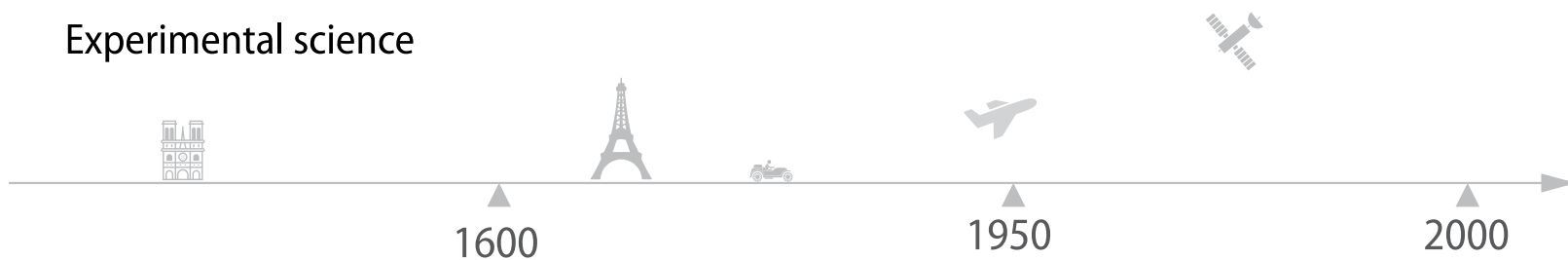
Theoretical science

3rd paradigm

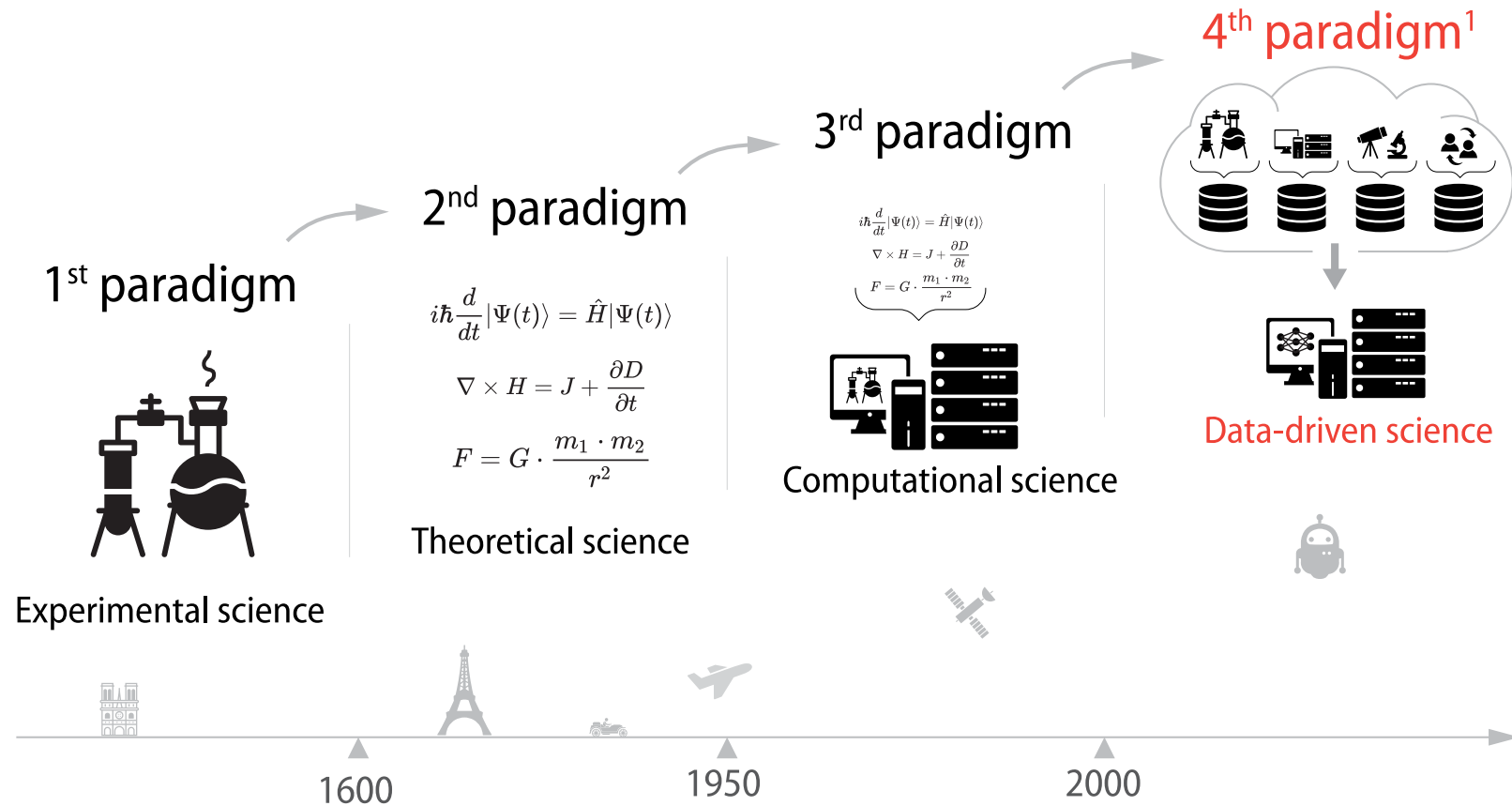
$$i\hbar \frac{d}{dt} |\Psi(t)\rangle = \hat{H} |\Psi(t)\rangle$$
$$\nabla \times H = J + \frac{\partial D}{\partial t}$$
$$F = G \cdot \frac{m_1 \cdot m_2}{r^2}$$



Computational science

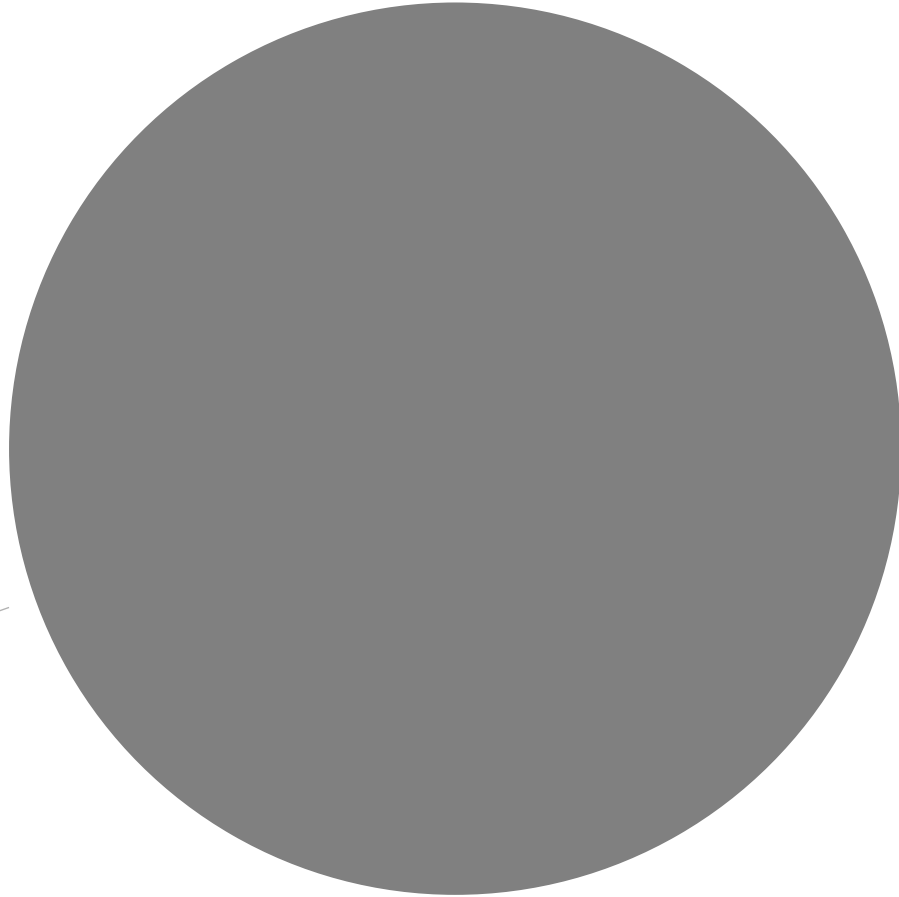


Scientific paradigms

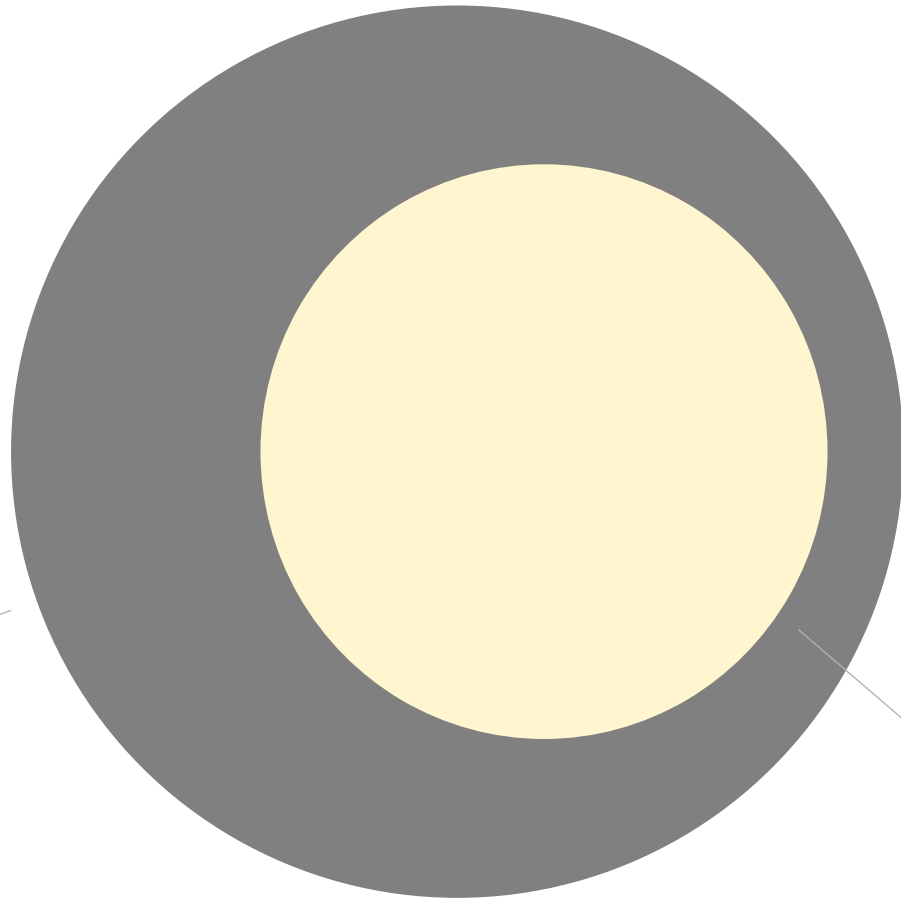


¹ Jim Gray, 2007 [GRAY]

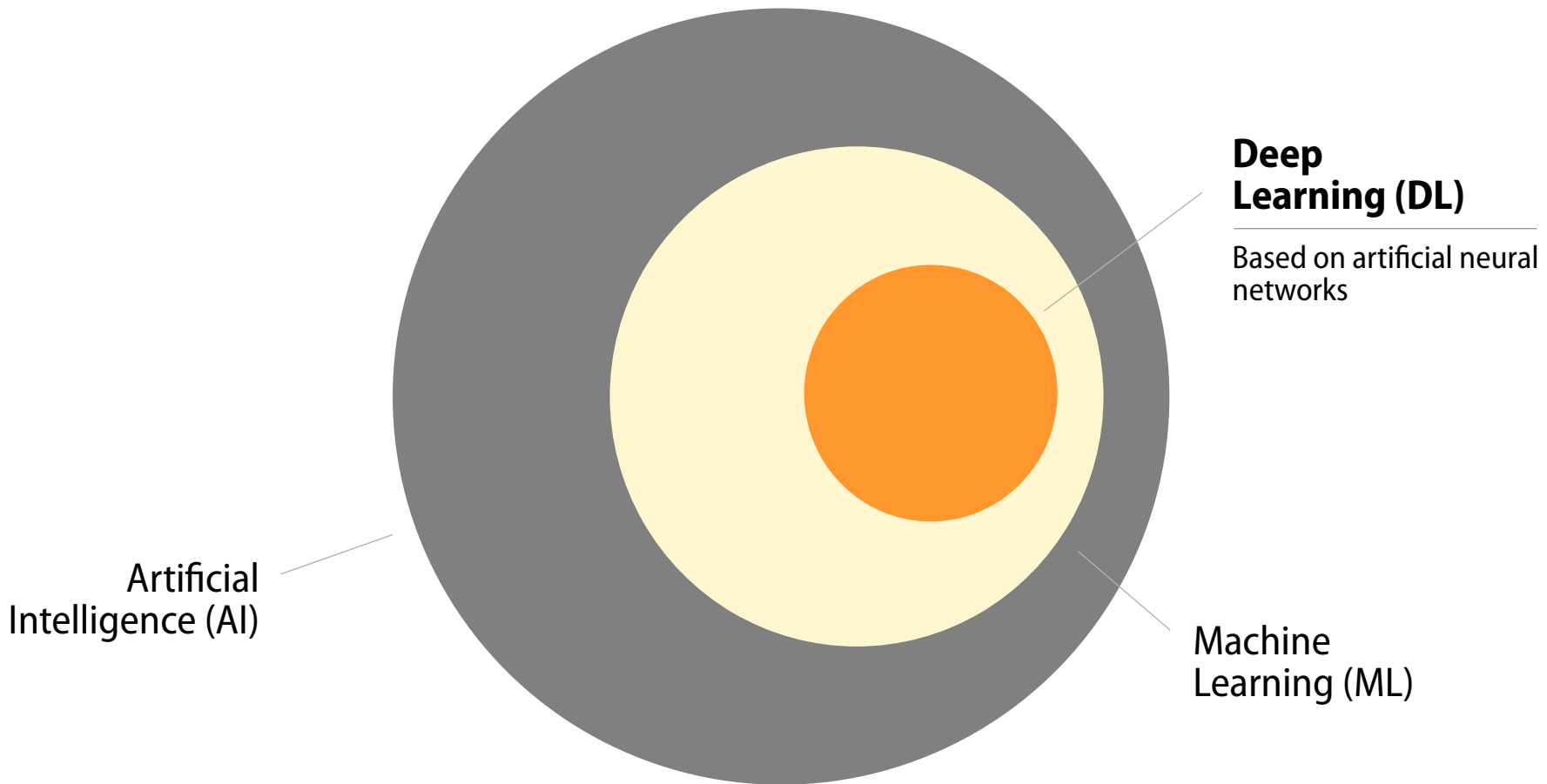
Artificial
Intelligence (AI)

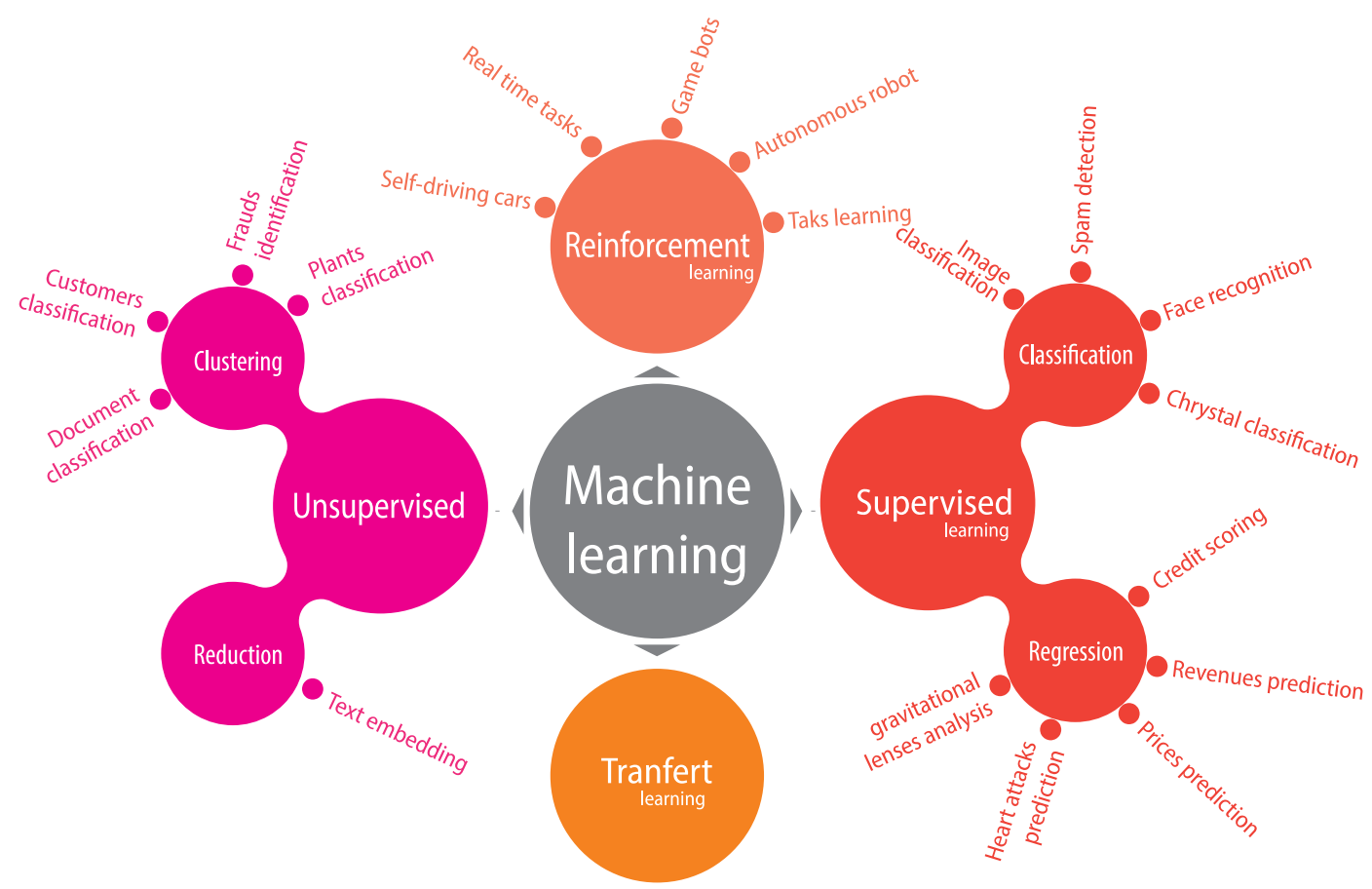


Artificial
Intelligence (AI)

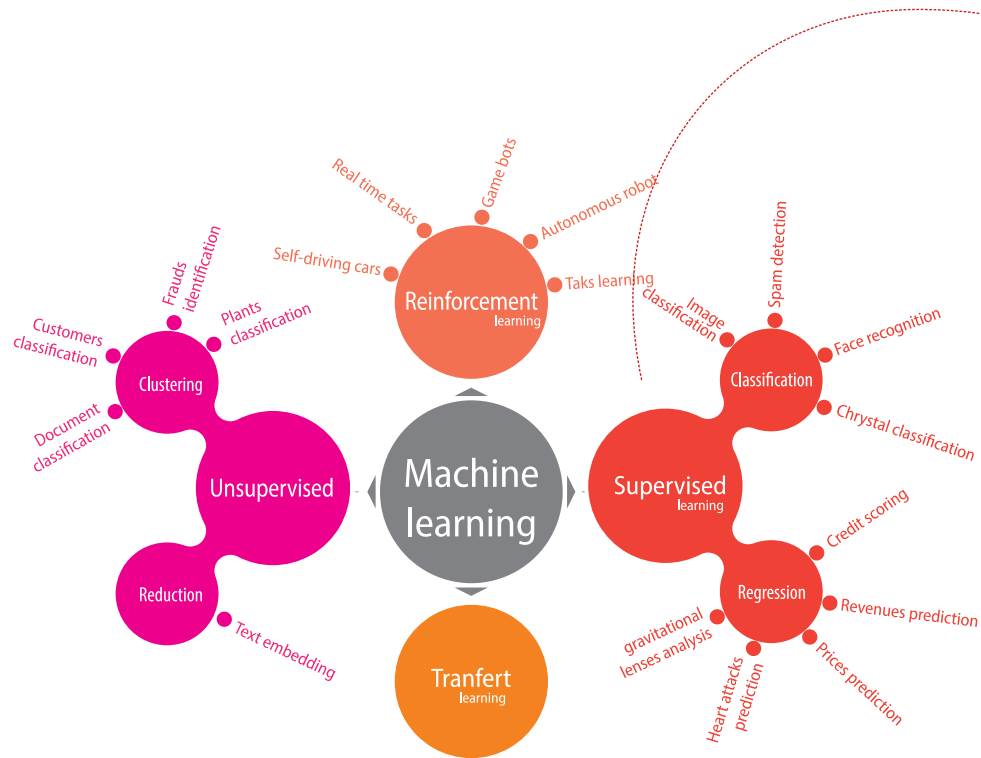


Machine
Learning (ML)

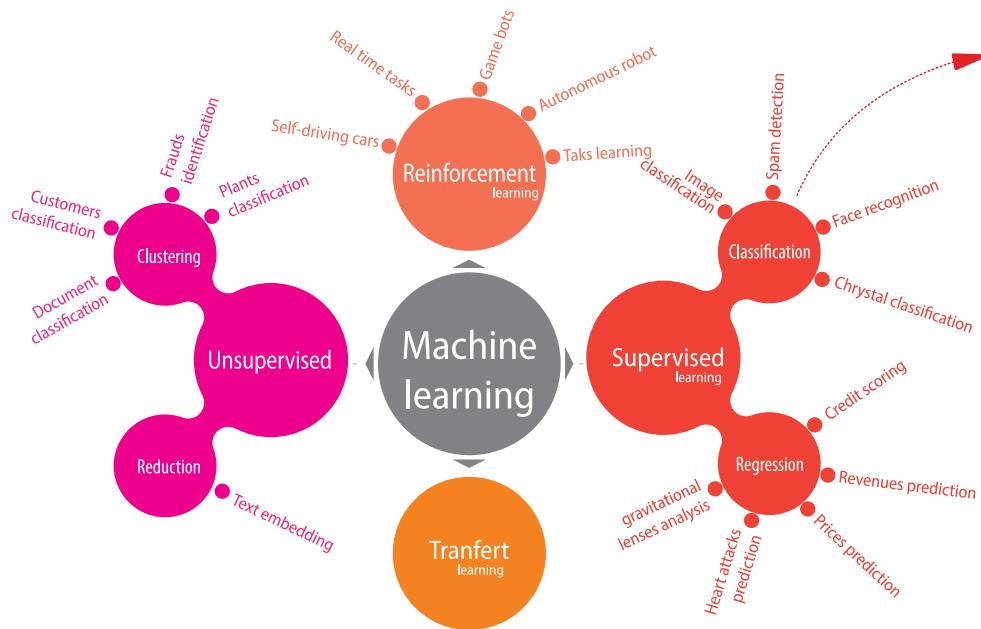




Supervised learning



Learning from examples



Classification :

Predict qualitative informations



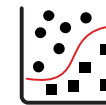
This is a cat

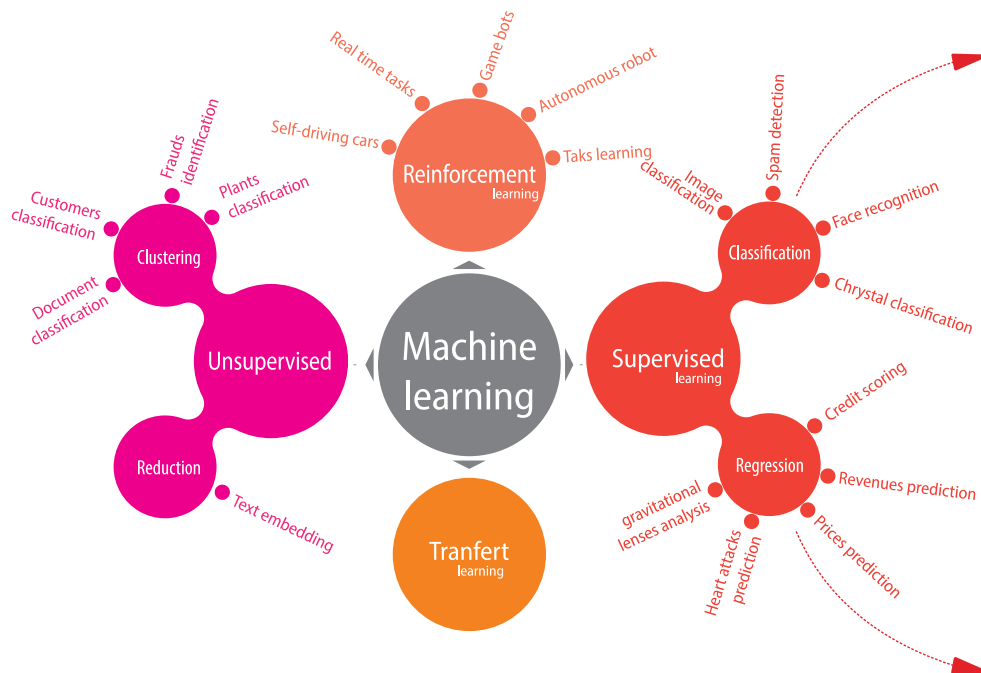


This is a rabbit



Tell me,
what is it ?





Classification :

Predict qualitative informations



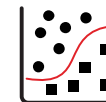
This is a cat



This is a rabbit



Tell me,
what is it ?



Régression :

Predict quantitative informations



150 K€



400 K€



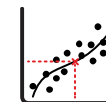
120 K€



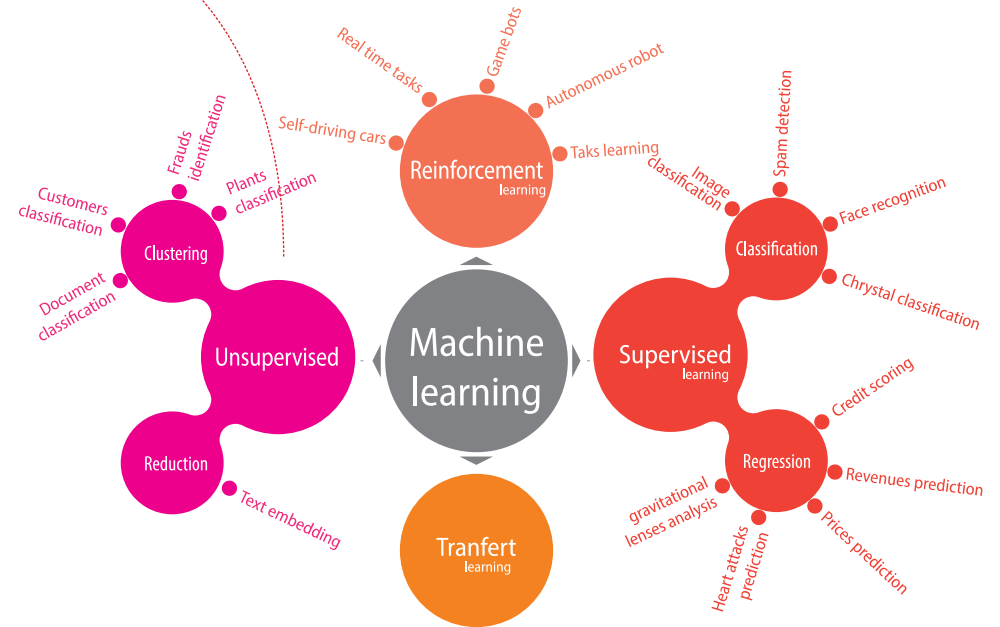
100 K€



Tell me,
what's the
price ?

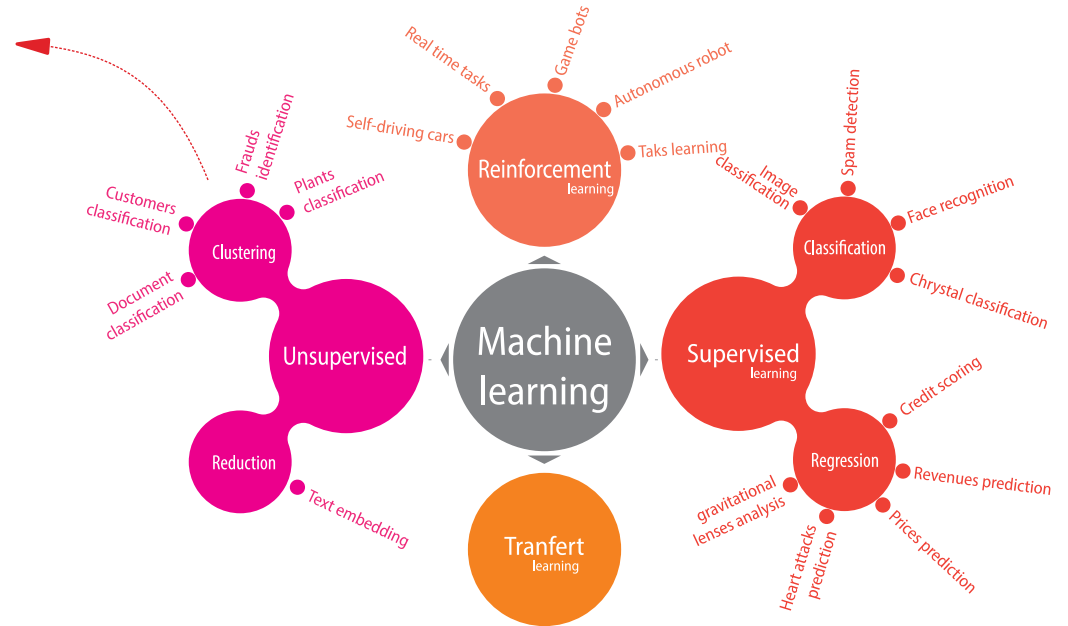
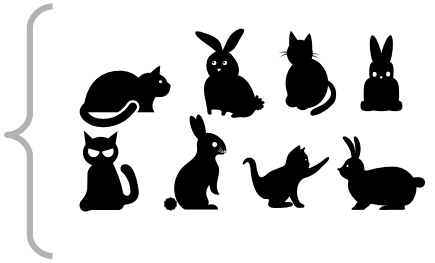


Learning from data alone



Clustering:
Finding Common Relationships

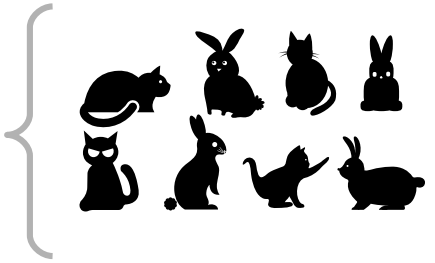
What is the relationship between these data?



Clustering:
Finding Common Relationships



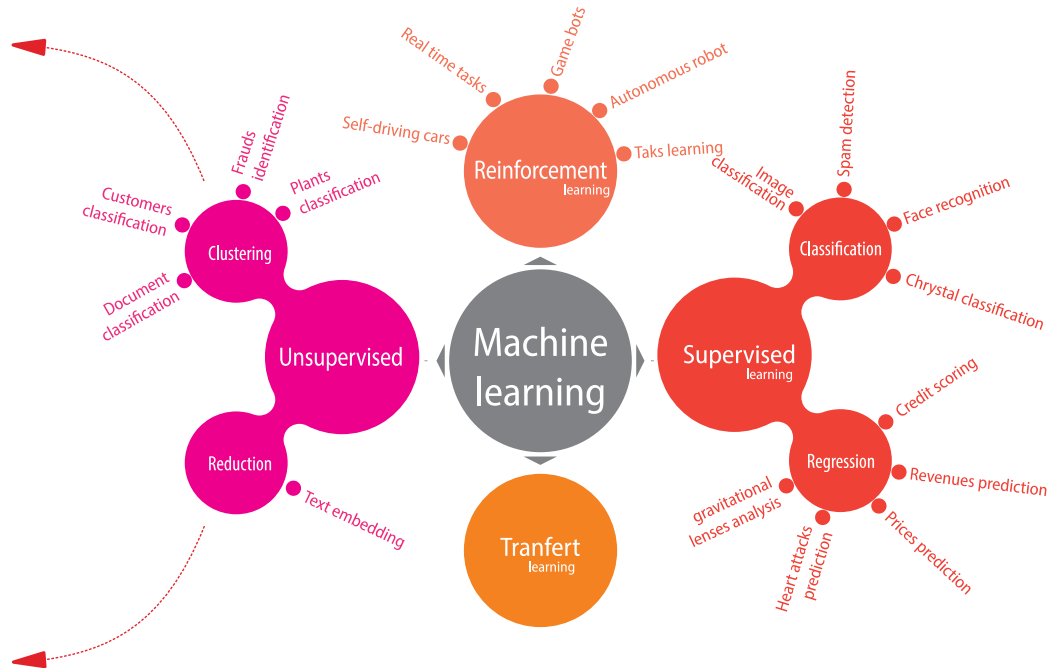
What is the relationship between these data?



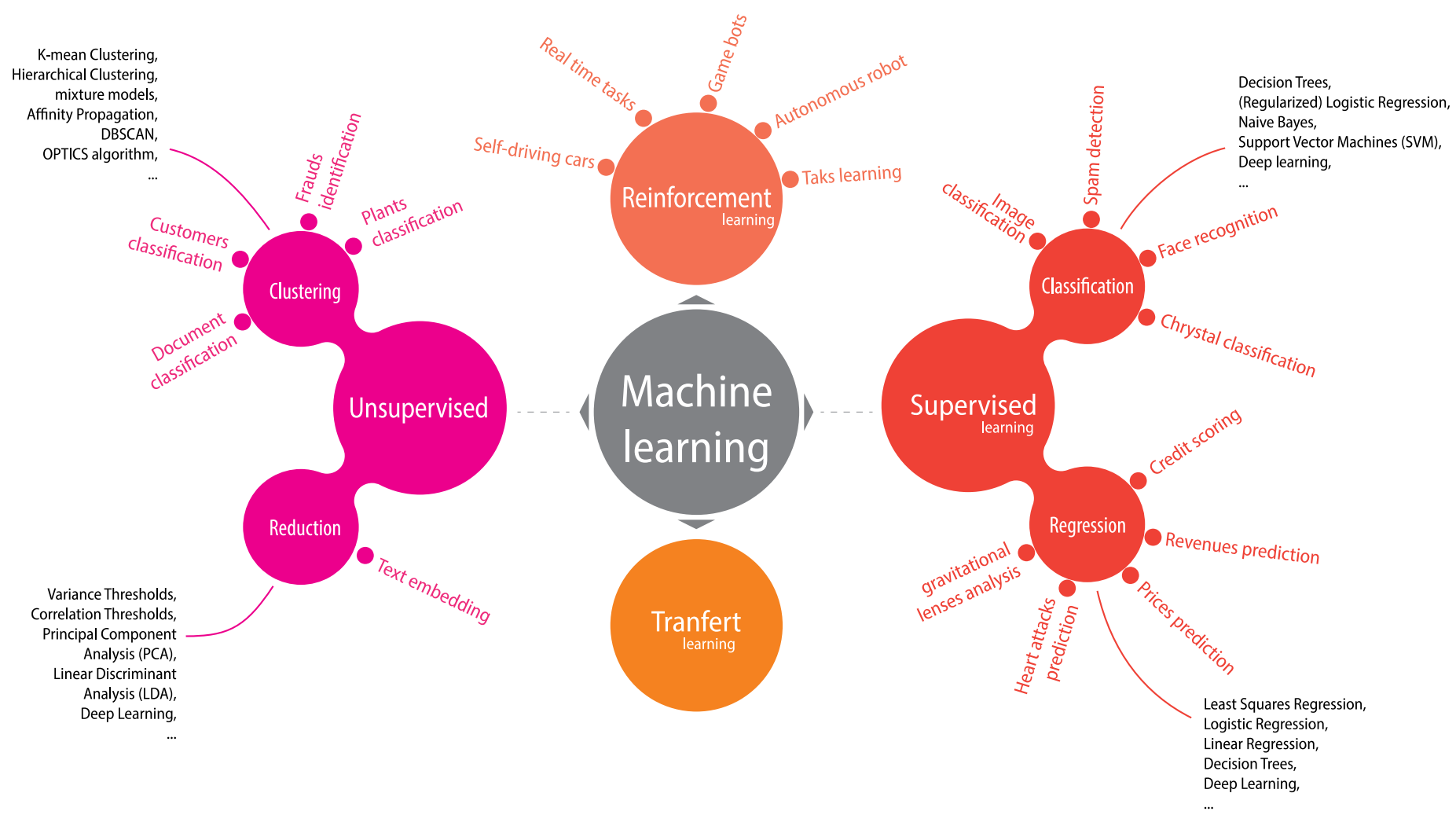
Reduction:
Reduce the number of dimensions



Simplify while keeping meaning



[*-learning]



[intelligence]



[intelligence]

« Capacité de percevoir ou d'inférer l'information, et de la conserver comme une connaissance à appliquer à des comportements adaptatifs dans un environnement ou un contexte donné »

*« Ability to perceive or infer information, and to retain it as knowledge to be applied towards adaptive behaviors within an environment or context »**



[intelligence]

« Ensemble des **fonctions** mentales ayant pour objet la connaissance **conceptuelle** et **rationnelle** »*

« Set of mental functions aimed at conceptual and rational knowledge »

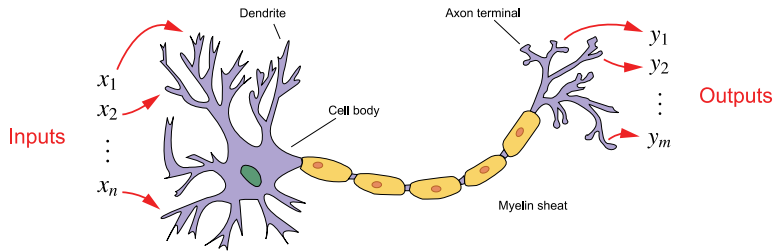
Modelling the brain :

« Penser s'apparente à un calcul massivement parallèle de **fonctions élémentaires**.

L'information est un **signal** avant d'être un code »¹

Connectionnism

Modelling the brain
Modéliser le cerveau



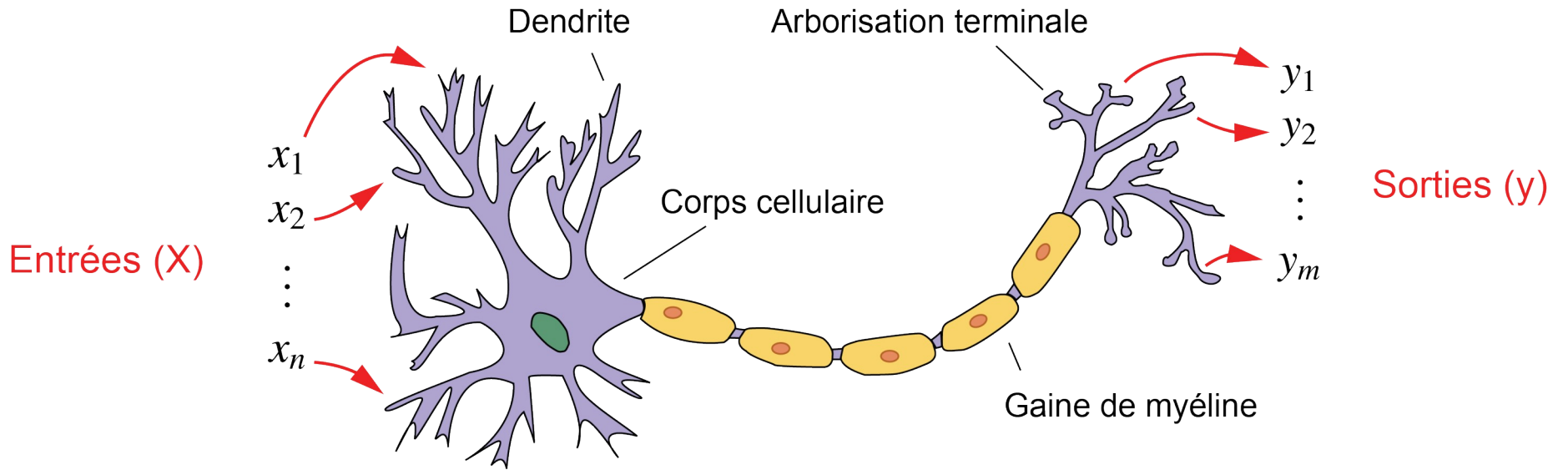
VS

Symbolic

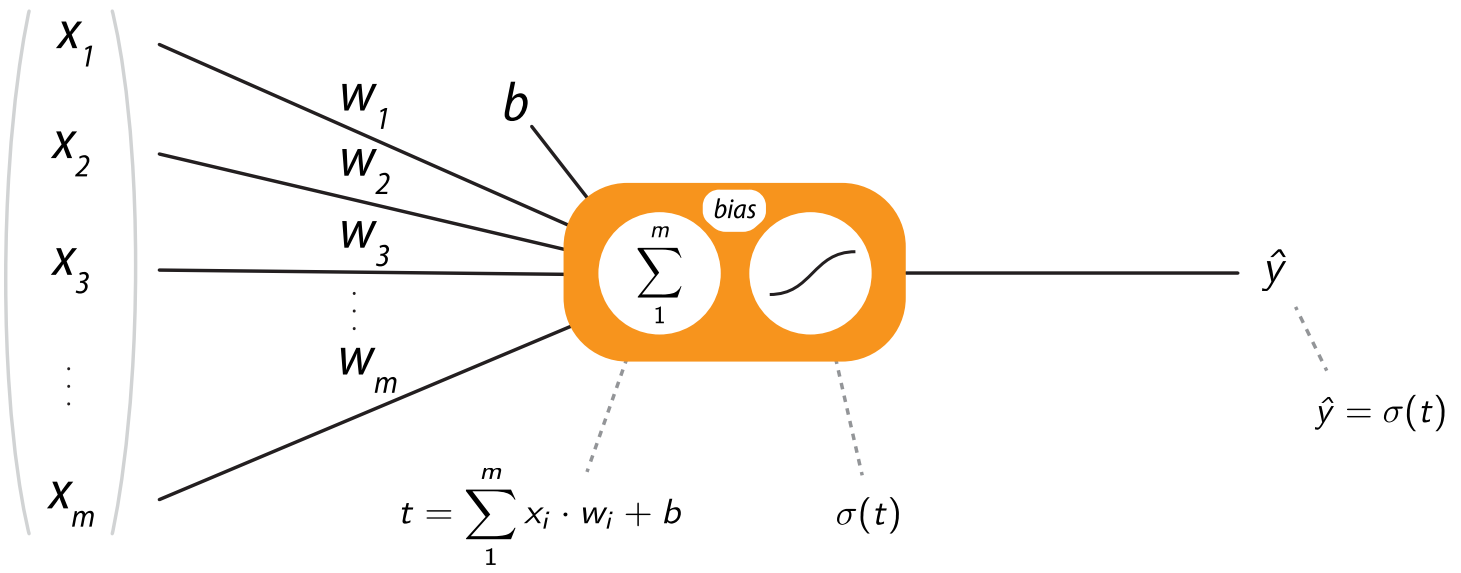
Making a mind
Forger une opinion

Tout [homme] est [mortel]
 [Socrate] est un [homme]
 Donc [Socrate] est [mortel]

¹ Dominique Cardon, Jean-Philippe Cointet, Antoine Mazieres (2018) [LRDN]



$$\hat{y} = \sigma(\Theta^T \cdot X + b)$$



Input
 X

Bias / Weight
 Θ, b

Activation function
 $\sigma(t)$

Output
 \hat{y}

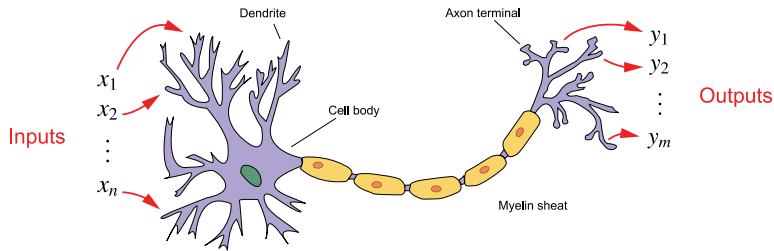
Modelling the brain :

« Penser s'apparente à un calcul massivement parallèle de **fonctions élémentaires**.

L'information est un **signal** avant d'être un code »¹

Connectionnism

Modelling the brain
Modéliser le cerveau



Making a mind :

« Penser, c'est calculer des **symboles** qui ont à la fois une réalité matérielle et une valeur sémantique de représentation »¹

L'information est une donnée symbolique de **haut niveau**.

Symbolic

Making a mind
Forger une opinion

Tout [homme] est [mortel]
[Socrate] est un [homme]
Donc [Socrate] est [mortel]

VS

¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

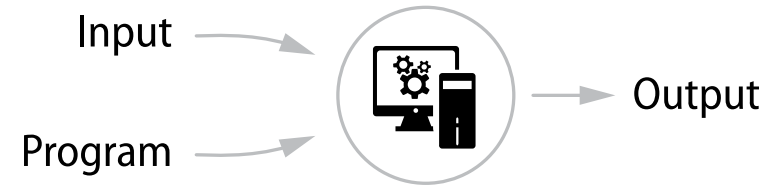
Inductive approach



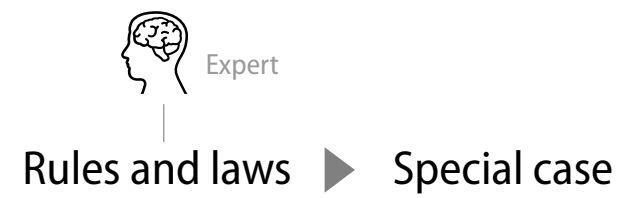
Connectionnism

vs

Deductive approach

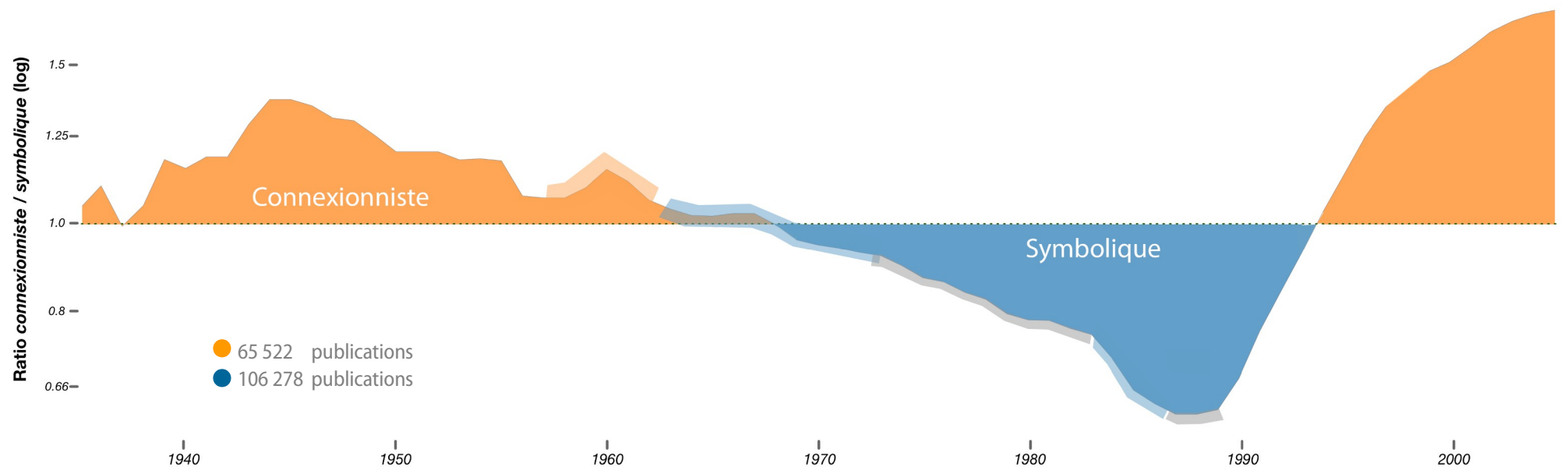


Symbolic



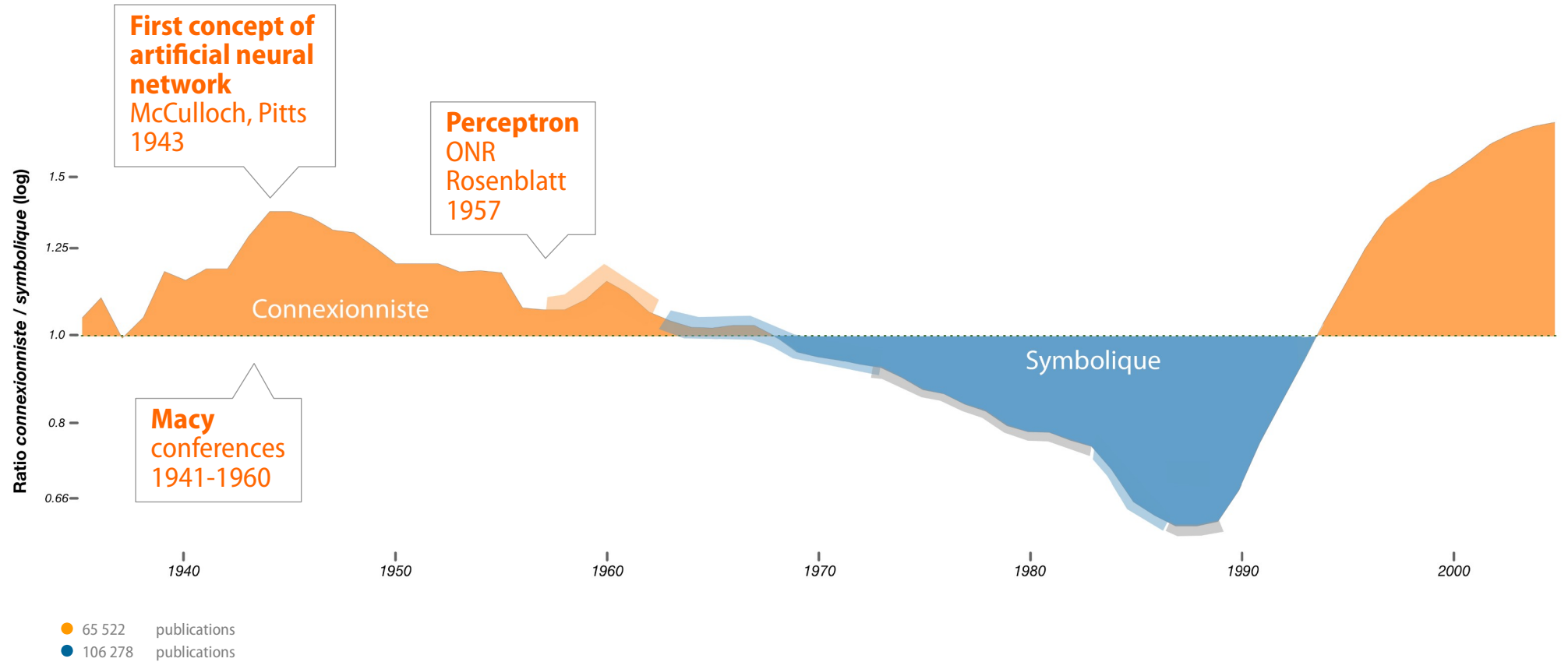
Evolution of the academic influence of connexionist and symbolic approaches¹

Ration of publications between connexionists and symbolists



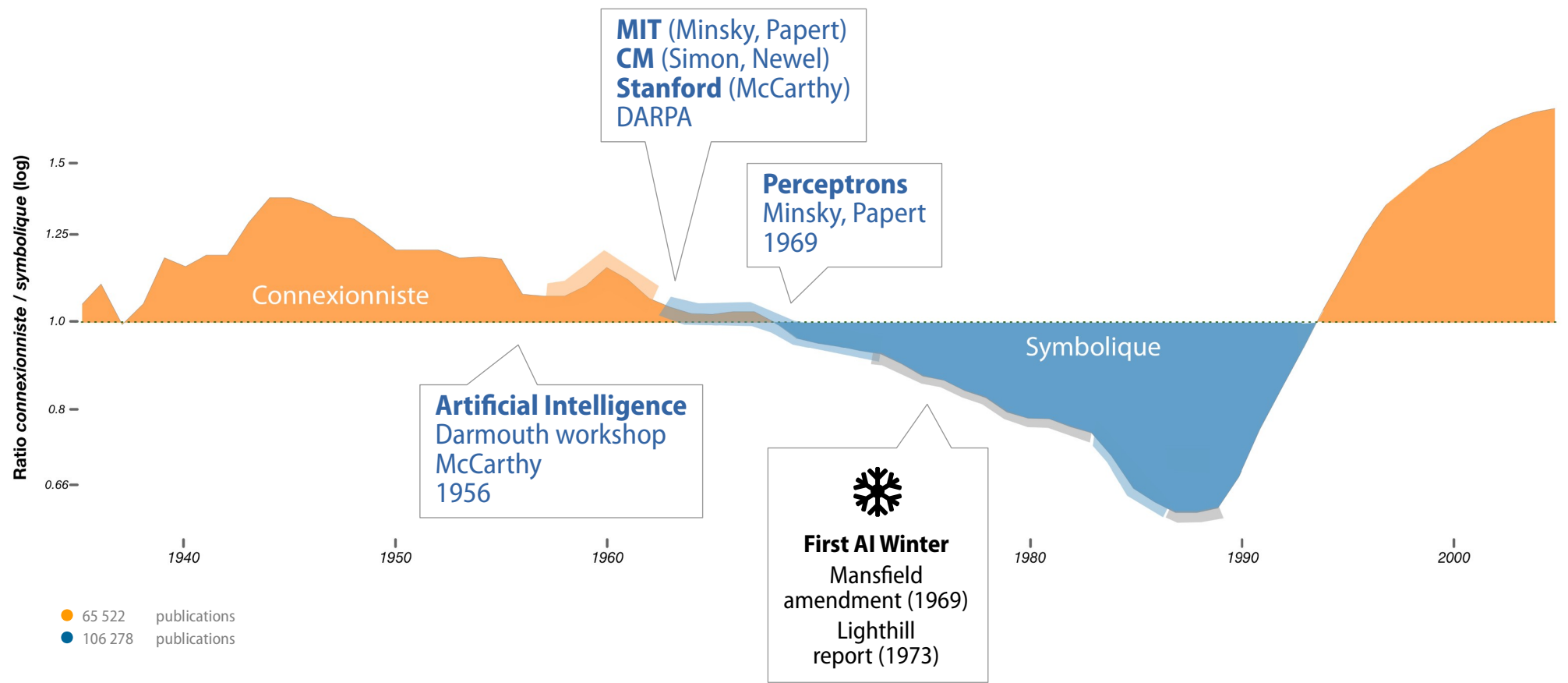
¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

Evolution of the academic influence of connexionist and symbolic approaches¹



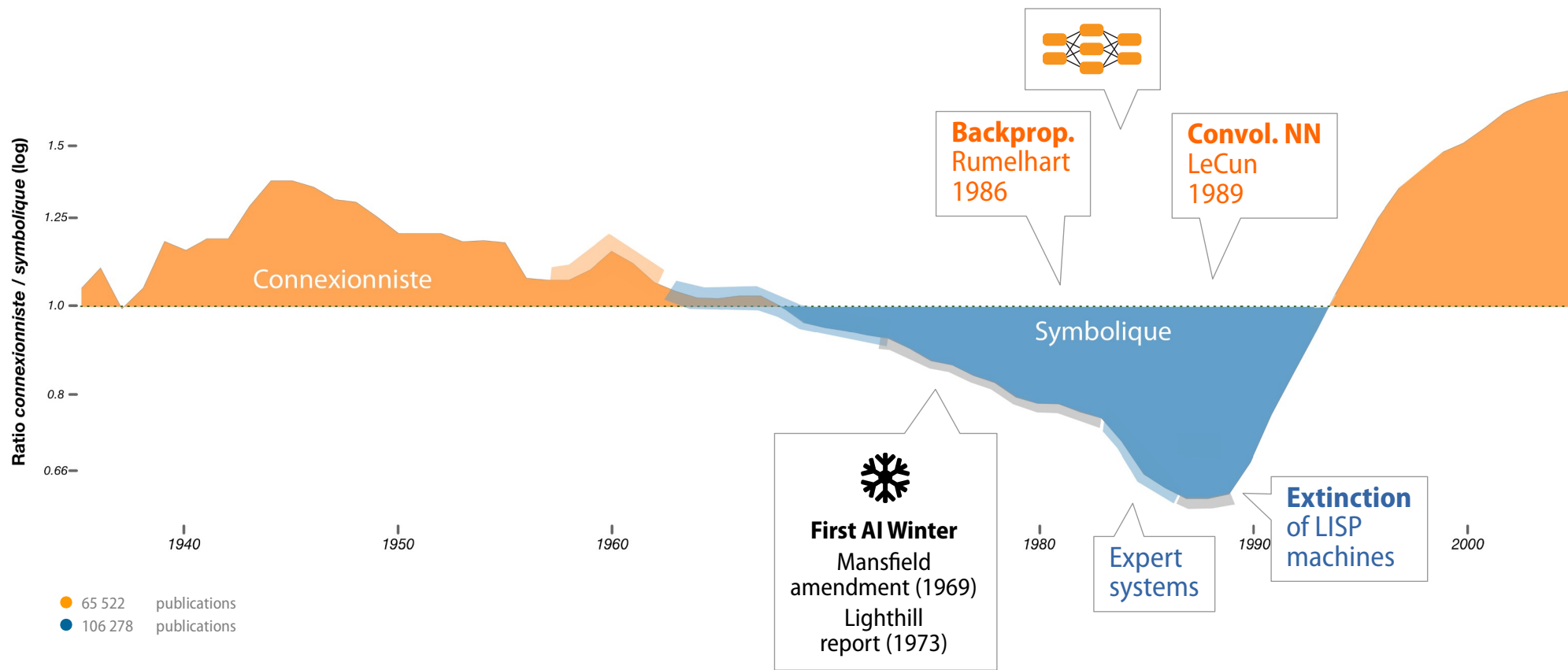
¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

Evolution of the academic influence of connexionist and symbolic approaches¹



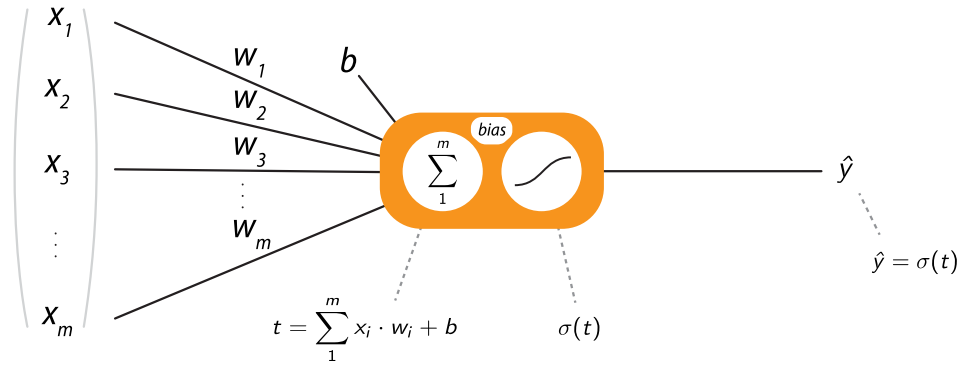
¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

Evolution of the academic influence of connexionist and symbolic approaches¹

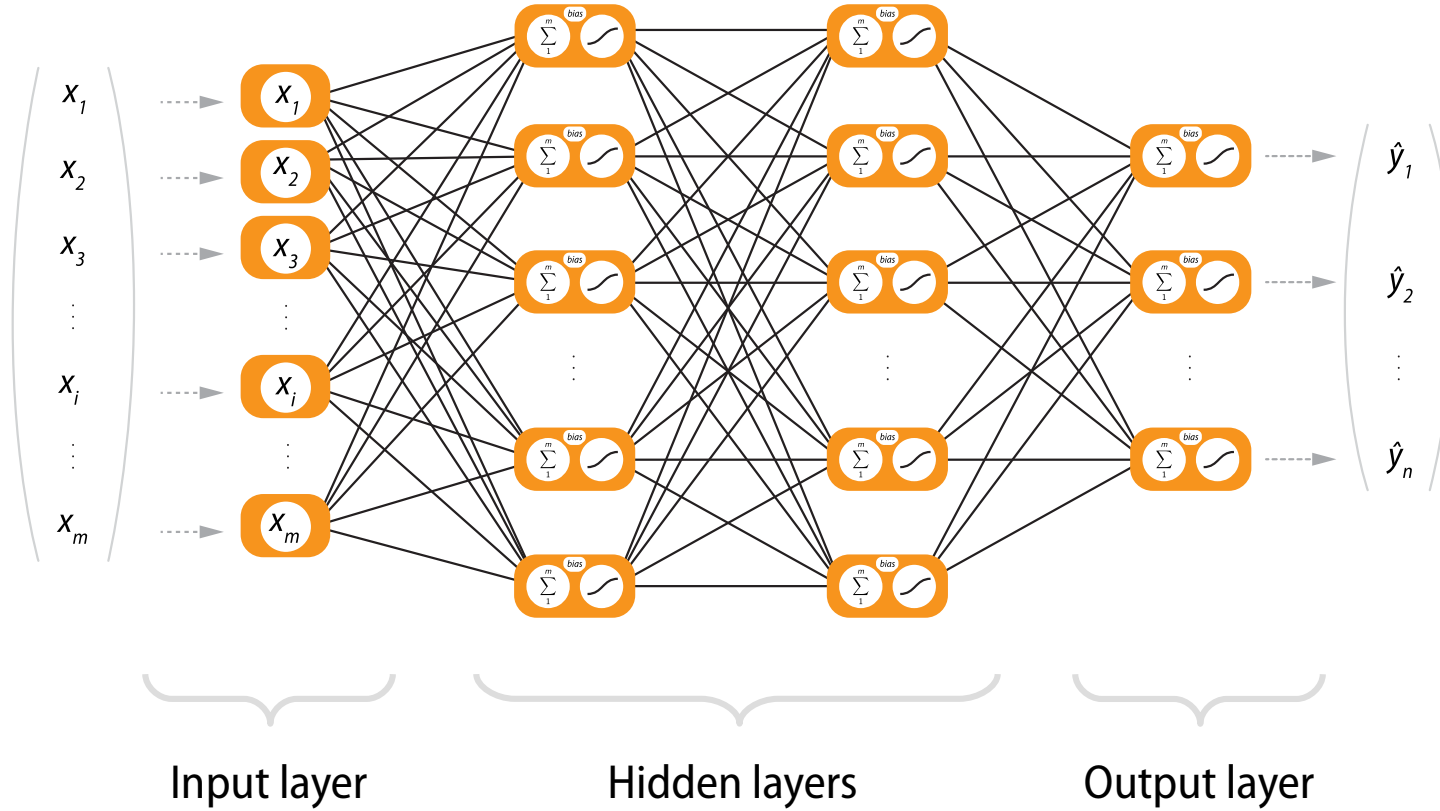


¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

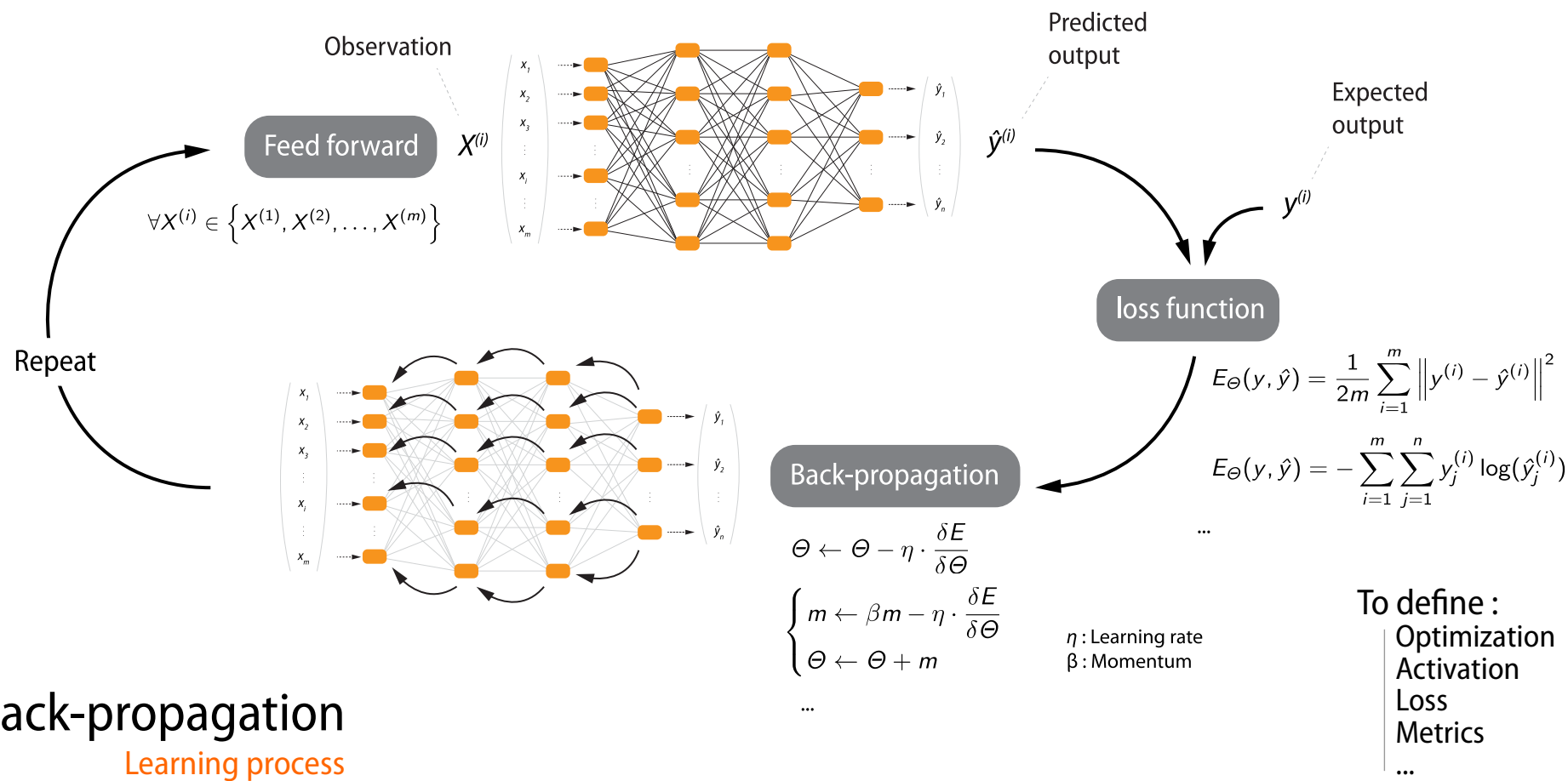
Deep Neural Networks



Deep Neural Networks

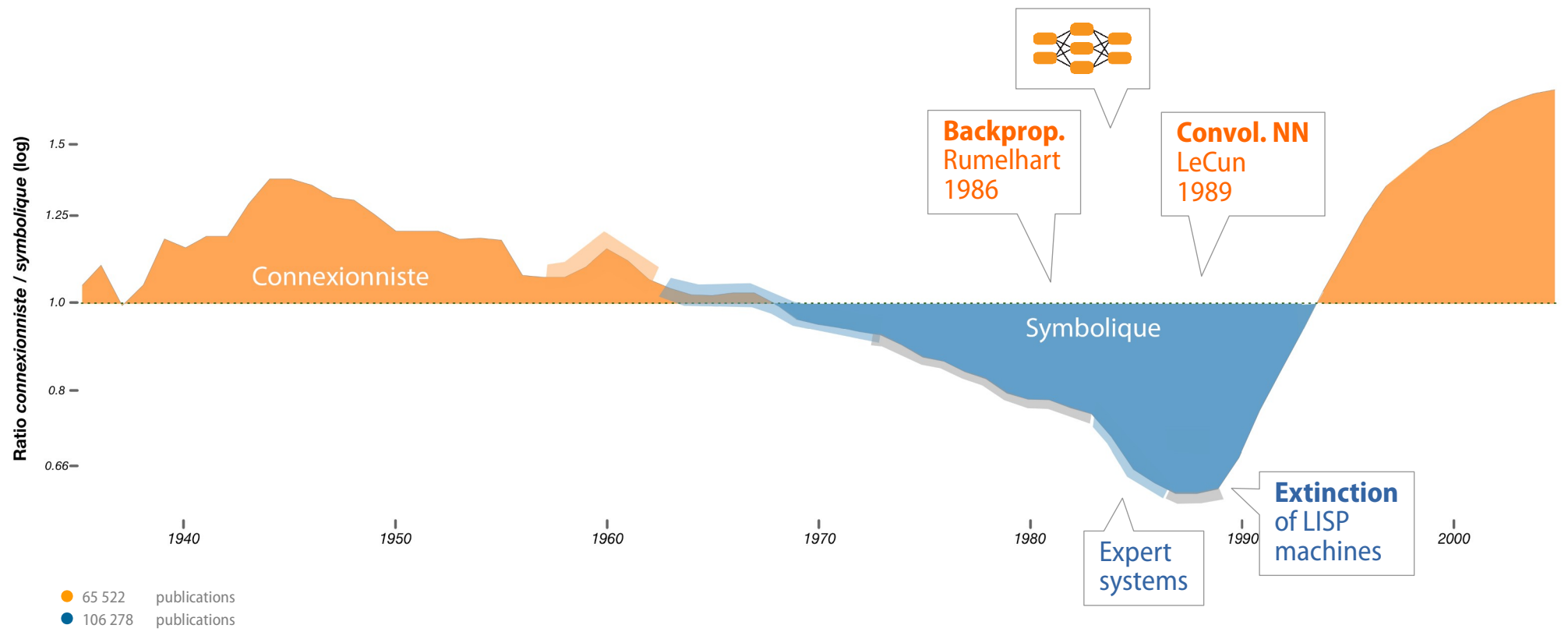


Deep Neural Networks



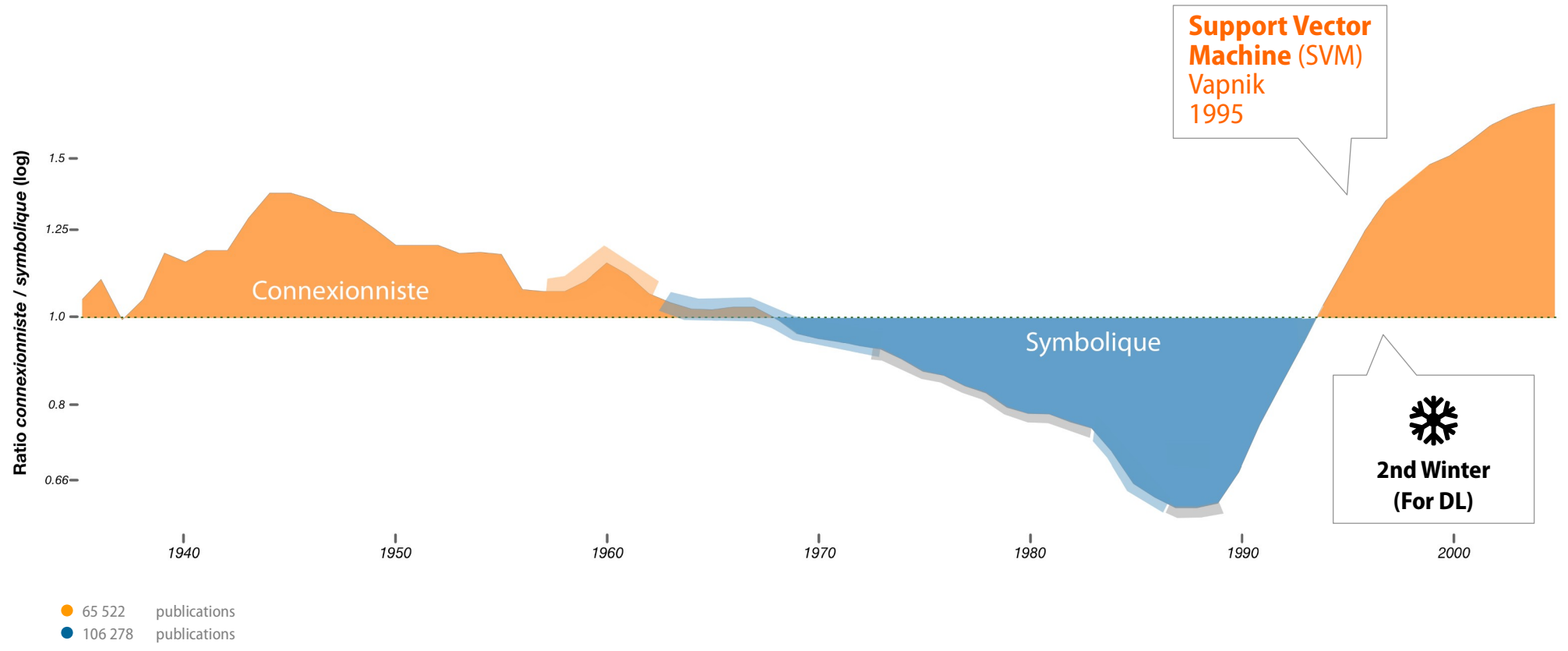
Back-propagation
Learning process

Evolution of the academic influence of connexionist and symbolic approaches¹



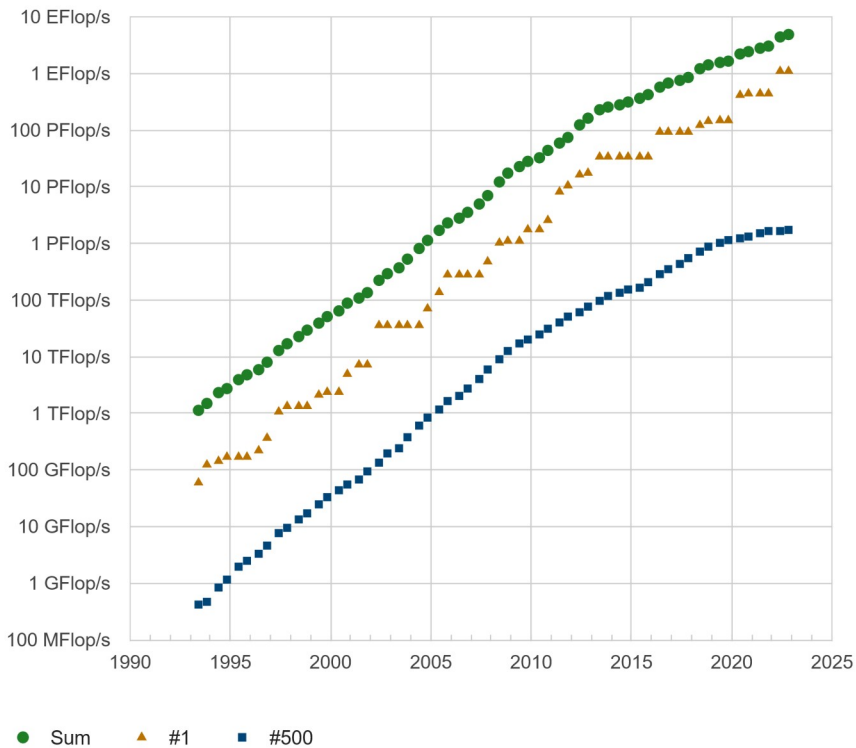
¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

Evolution of the academic influence of connexionist and symbolic approaches¹

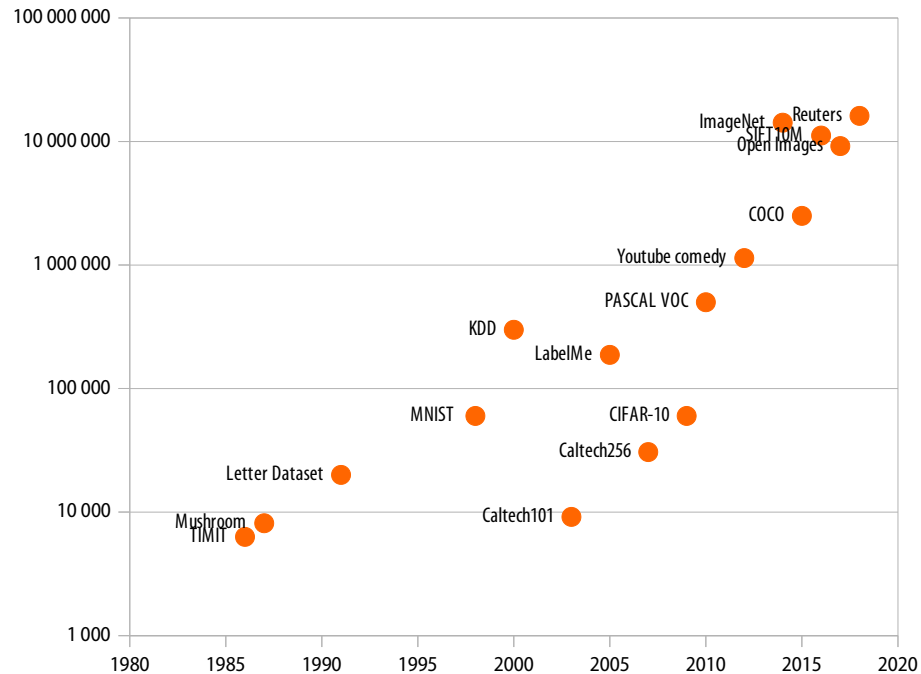


¹ D Cardon, JP Cointet, A Mazieres, 2018 [LRDN]

Performance Development¹



Datasets for machine-learning²



Laboratoire
Cas particulier

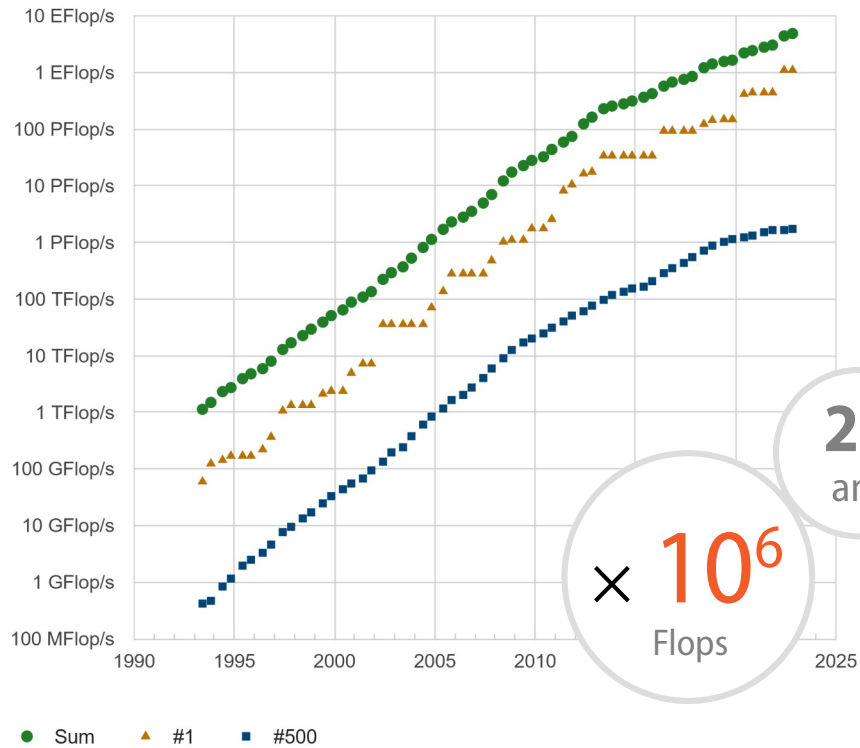


Monde réel

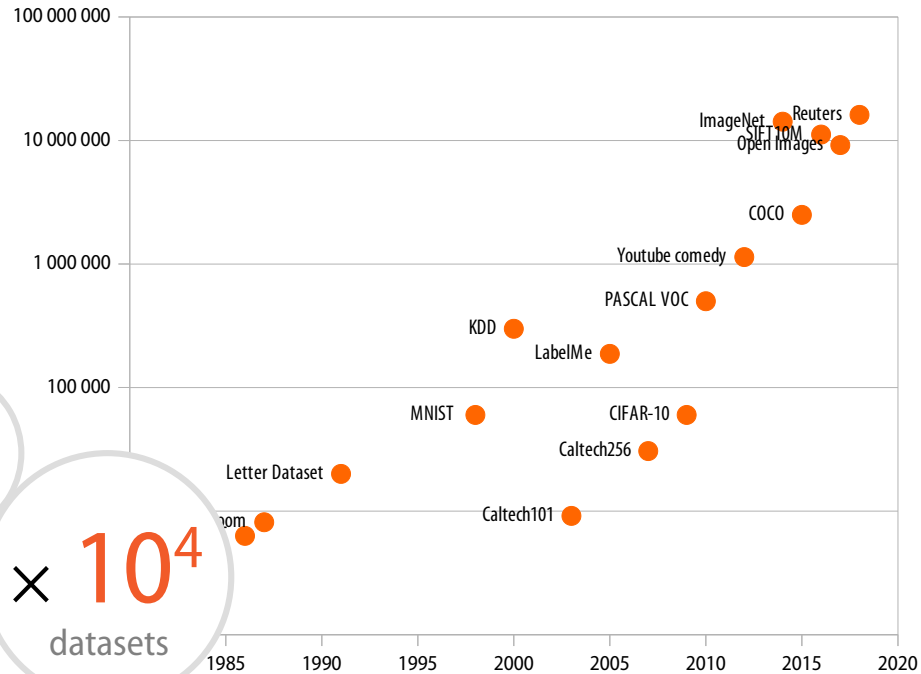
¹ TOP500 List [TOP500]


² Wikipedia [WKP1]

Performance Development¹



Datasets for machine-learning²



Laboratoire Cas particulier  Monde réel

× 10⁶
Flops

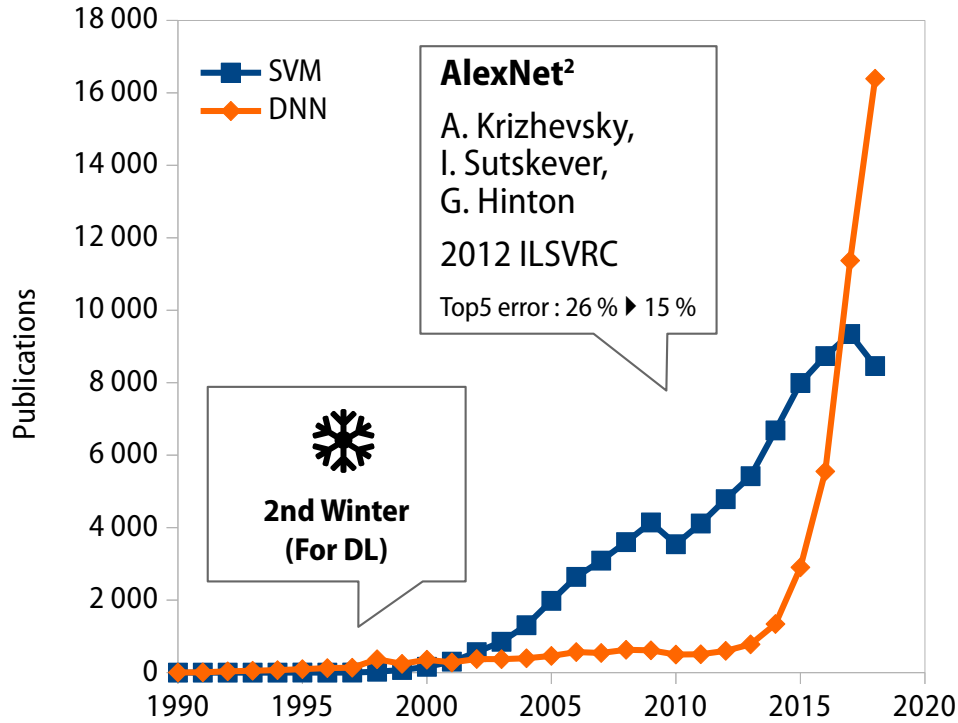
25 ans

× 10⁴
datasets

¹ TOP500 List [TOP500]

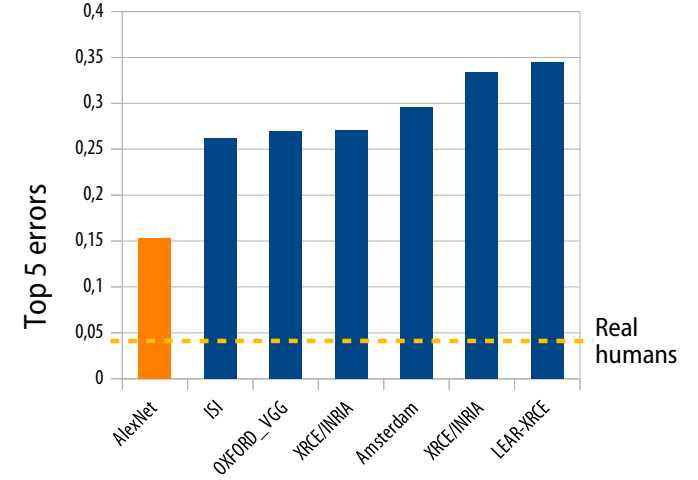
² Wikipedia [WKP1]

Publications SVM vs DNN¹



DNN

Images classification Top 5 error at ILSVRC 2012^{3,4}



Without mathematical guarantee, DNN have proven to be more effective in the face of the **complexity of the real world!**

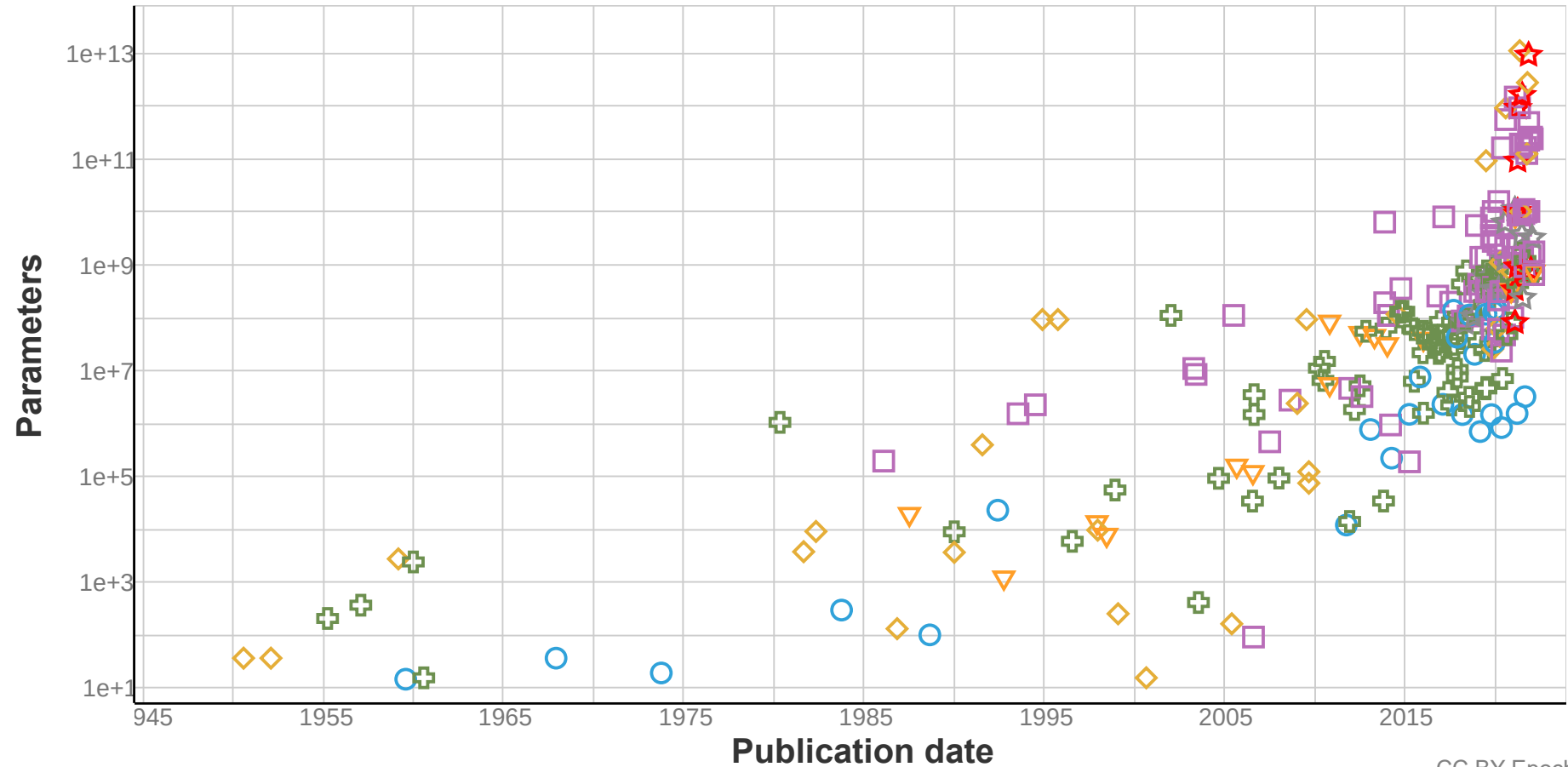
¹ Web of Science [WOS1][WOS2]

² AlexNet [ALEX]

³ ImageNet Large Scale Visual Recognition [ILSVRC]

⁴ Similar evolution in Natural language processing, translation, board games, etc.
 See : DeepL.com, AlphaGo, AlphaZero, ...

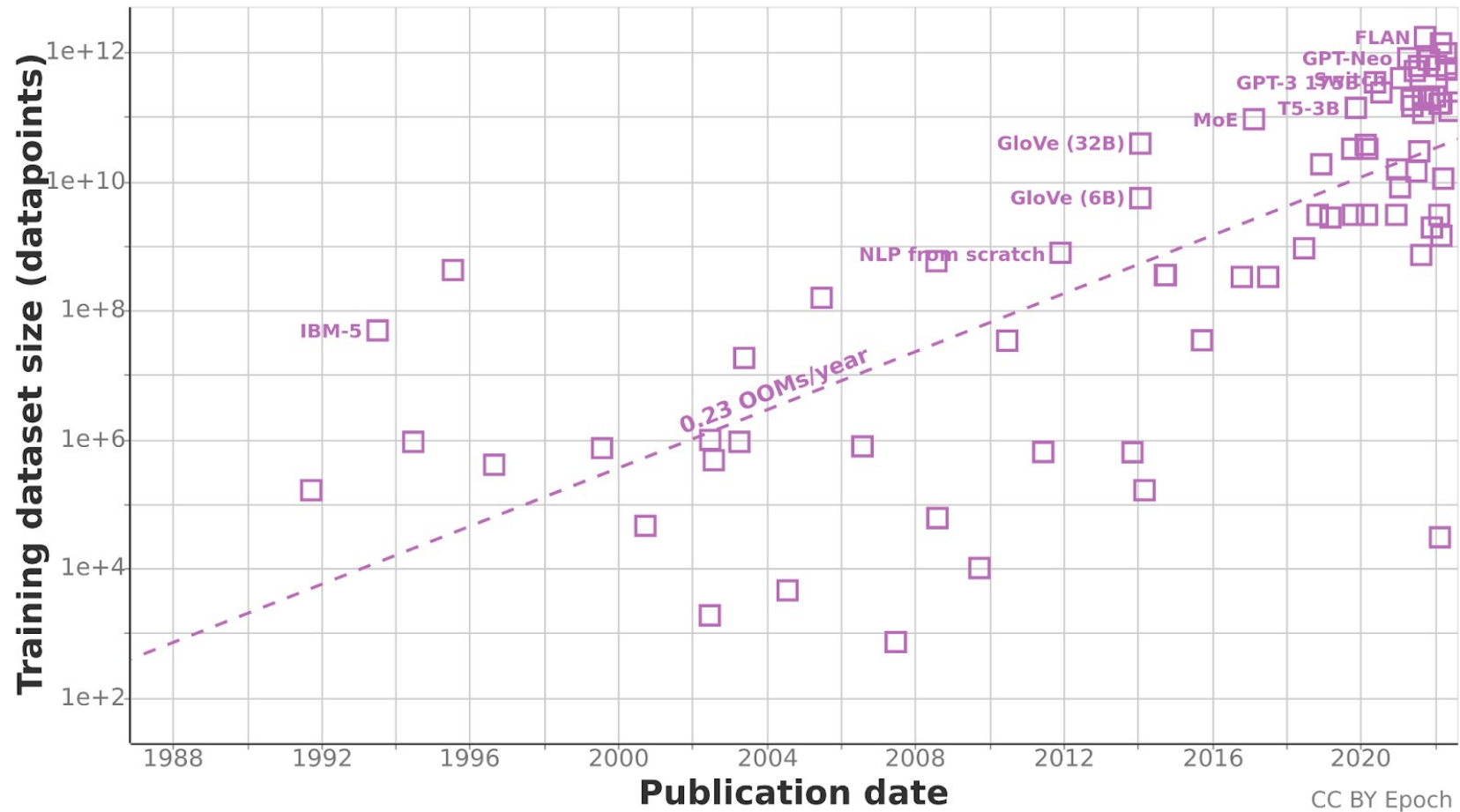
Machine Learning Model Sizes



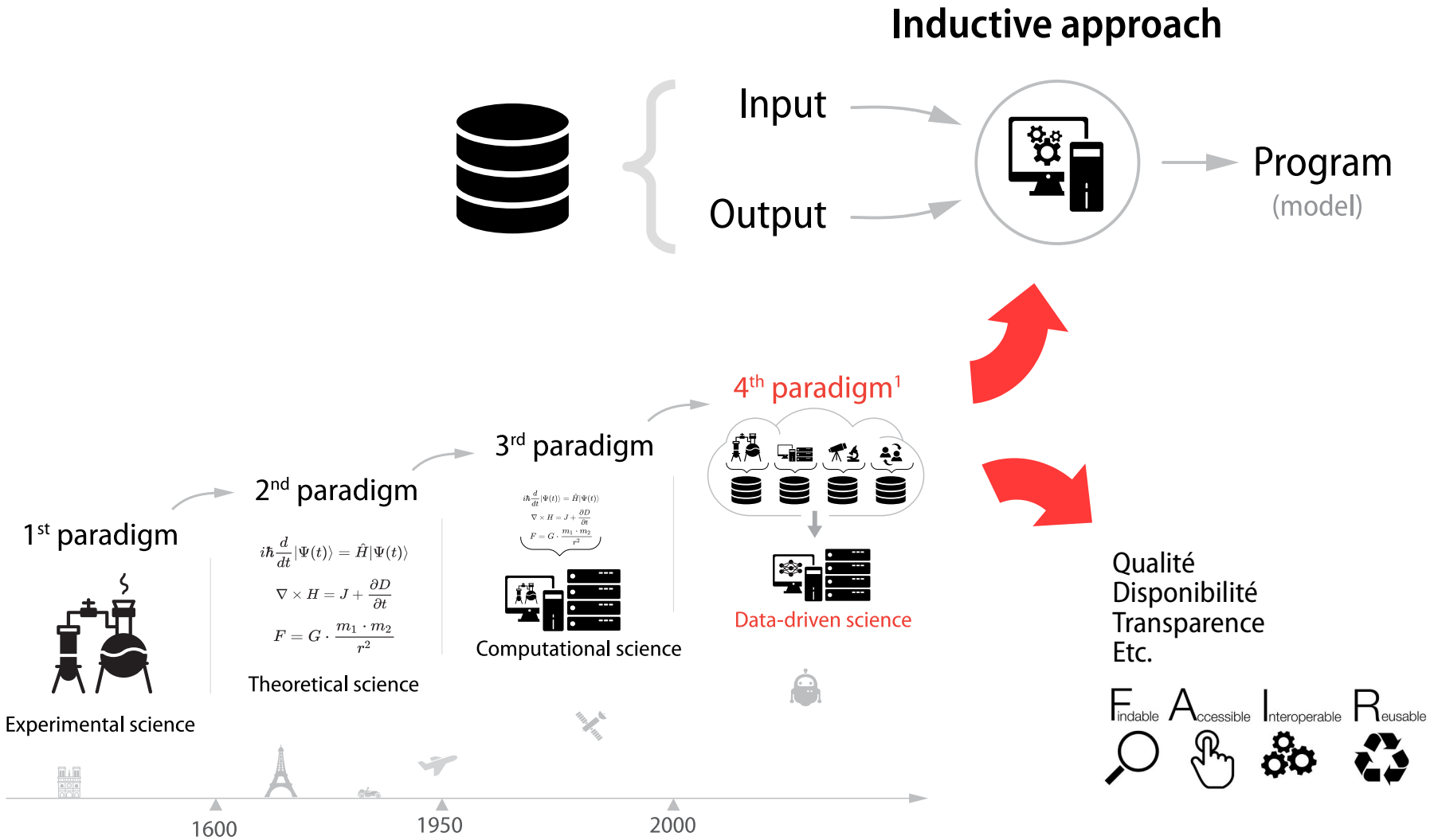
CC BY Epoch

Machine Learning Model Sizes and the Parameter Gap
Pablo Villalobos, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, Anson Ho, Marius Hobbhahn (2022)
<https://doi.org/10.48550/arXiv.2207.02852>
<https://epochai.org/blog/machine-learning-model-sizes-and-the-parameter-gap>

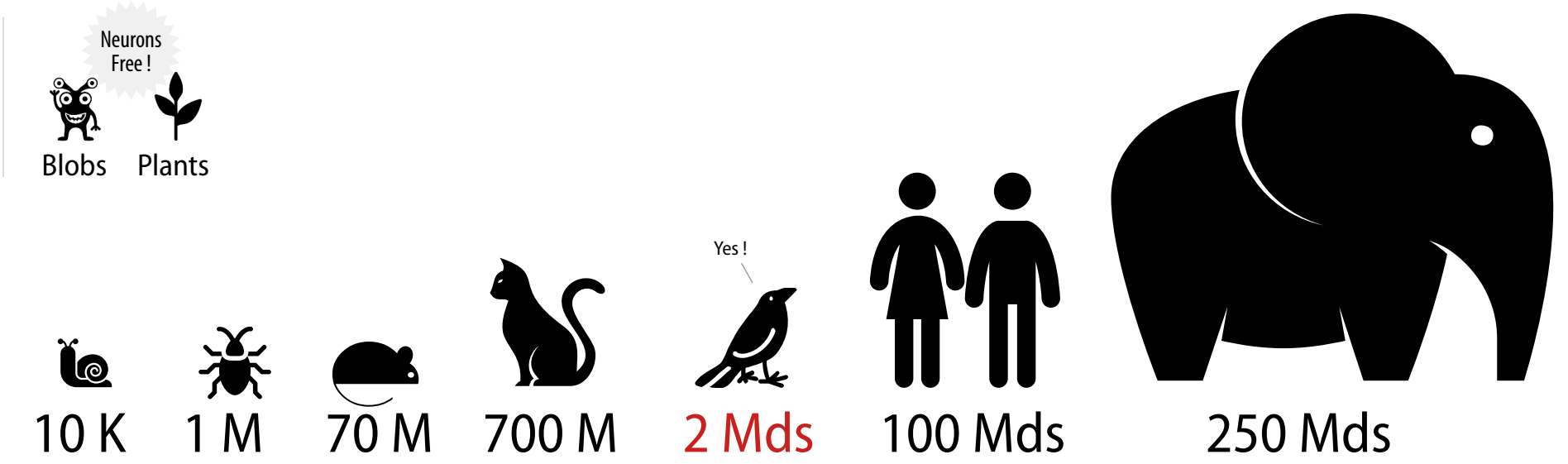
Datasets Sizes



CC BY Epoch



Some brain sizes, in number of neurons...



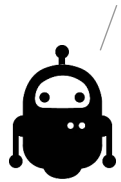


Yes!



Cette fois, c'est bien fini ;-)

Merci beaucoup !



Références

- [JGRAY] Gray, J. (2001), from « The Fourth Paradigm: Data-Intensive Scientific Discovery » Tony Hey, Stewart Tansley, Kristin Tolle (2009). Published by Microsoft Research.
ISBN: 978-0-9825442-0-4
- [FROS] Rosenblatt, Frank. (1958). « The perceptron: A probabilistic model for information storage and organization in the brain. » *Psychological Review*, 65(6), 386-408.
- [MIPA] Minsky, Marvin; Papert, Seymour. (1969). « Perceptrons : An Introduction to Computational Geometry », MIT Press
- [DRUM] Rumelhart, David E.; Hinton, Geoffrey E.; Williams, Ronald J. (1986). « Learning representations by back-propagating errors ». *Nature*. 323 (6088): 533–536. doi:10.1038/323533a0.
- [YLEC1] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, L. D. Jackel, « Backpropagation Applied to Handwritten Zip Code Recognition », AT&T Bell Laboratories
- [LRDN] Dominique Cardon, Jean-Philippe Cointet, Antoine Mazieres. (2018). « La revanche des neurones », *Réseaux, La Découverte*, 5 (211), <10.3917/res.211.0173>. <hal-01925644>
- [TOP500] Statistics on top 500 high-performance computers. (2018) « Exponential growth of supercomputing power as recorded by the TOP500 list ». <https://www.top500.org>
- [WKP1] Wikipedia/en. (2018) « List of datasets for machine-learning research ». <https://en.wikipedia.org>
- [WOS1] Core database : TS=("support vector machine*" OR ("SVM" AND "classification") OR ("SVM" AND "regression") OR ("SVM" AND "classifier") OR "support vector network*" OR ("SVM" AND "kernel trick*"))
- [WOS2] Core database : TS=("deep learning" OR "deep neural network*" OR ("DNN" AND "neural network*") OR "convolutional neural network*" OR ("CNN" AND "neural network*") OR "recurrent neural network*" OR ("LSTM" AND "neural network*") OR ("RNN*" AND "neural network*"))

Illustrations :

Potato plant From *Die Giftpflanzen Deutschlands*, Peter Esser, 1910, via [iconspng.com](https://www.iconspng.com)

Neuron Wikimedia Commons, the free media repository.

Icons thenounproject.com