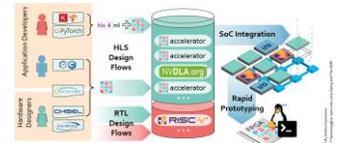
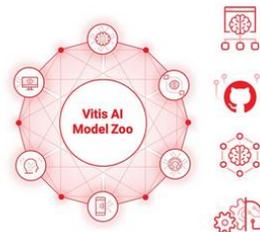
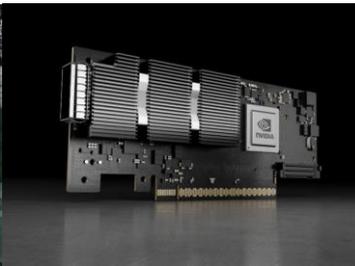
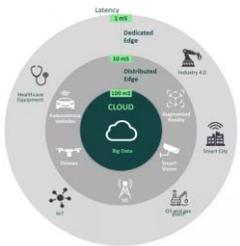


L'intelligence Artificielle dans les systèmes embarqués – Aspects Matériels –

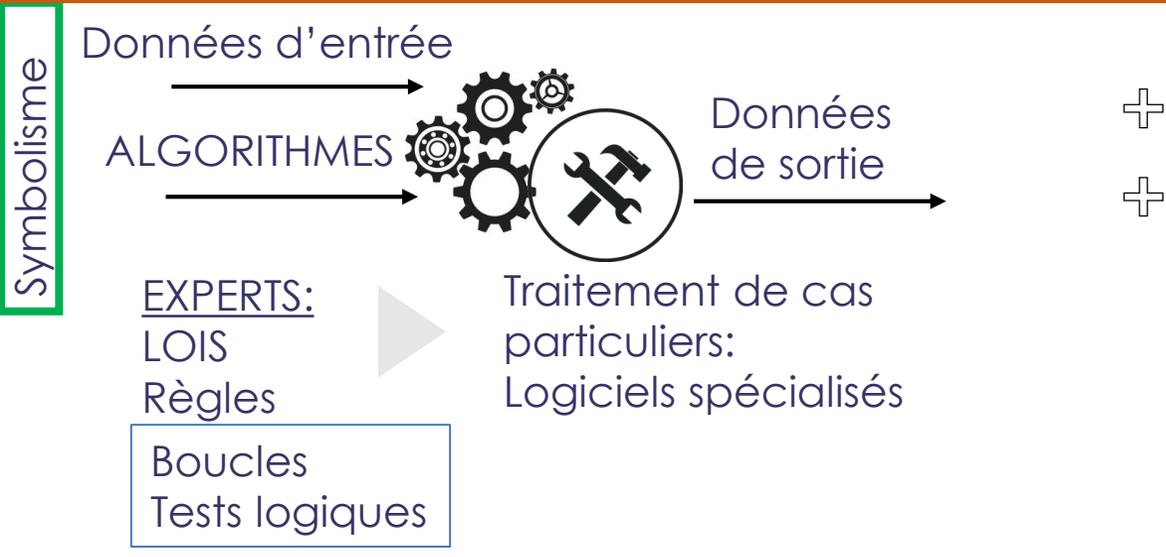


➔ Bilan Matériels
➔ Mise en œuvre

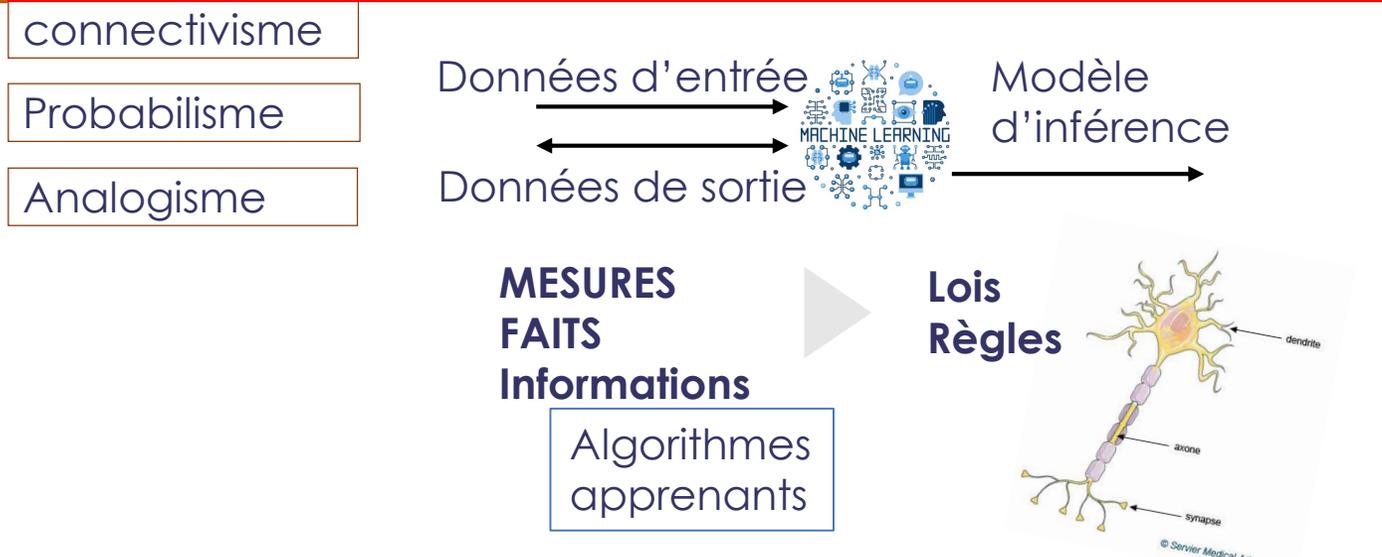


Accelerator generation with hls4ml 	Automatic integration in ESP 	Full-system RTL simulation 	Full-system test on FPGA
--	----------------------------------	--------------------------------	------------------------------

Approche déductive



Approche inférentielle (Inductive)



L'intelligence Artificielle dans les systèmes embarqués

Nouvelles architectures matérielles

Outils de générations de données

Réseaux de neurones et descente de gradient

CPU+ NNproc DSPs NMC

GPUs FPGAs

Frederic Druillolle LP21 Bd

SciPy

AI Custom Data and Models

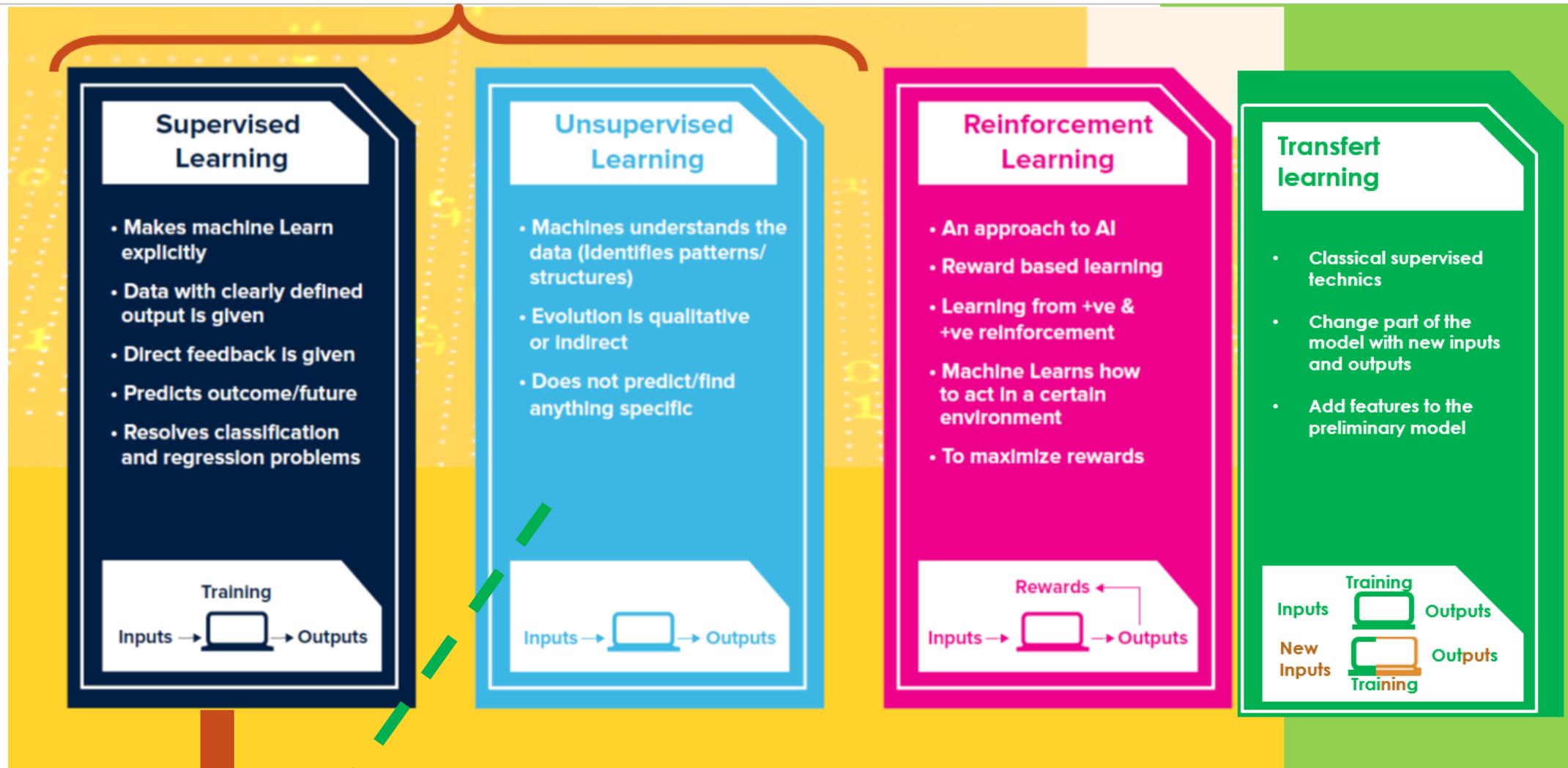
roboflow

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

- Input Cell
- BackFed Input Cell
- ▲ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- ▲ Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- ▲ Gated Memory Cell
- Kernel
- Convolution or Pool

Les quatre catégories d'apprentissage profond



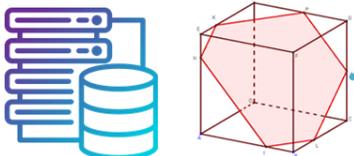
Système Embarqué

Algo Embarqué

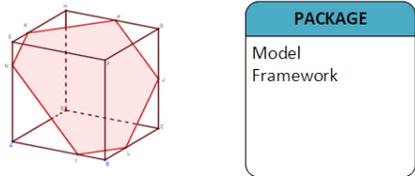
Méthodes

MLOps

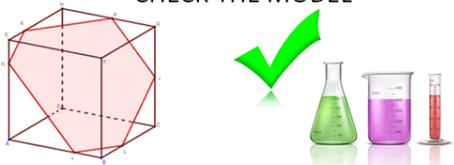
TRAIN THE MODEL



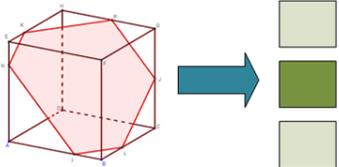
PACKAGE THE MODEL



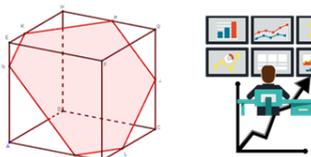
CHECK THE MODEL



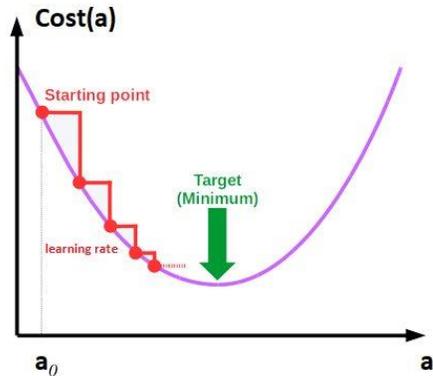
DEPLOY THE MODEL



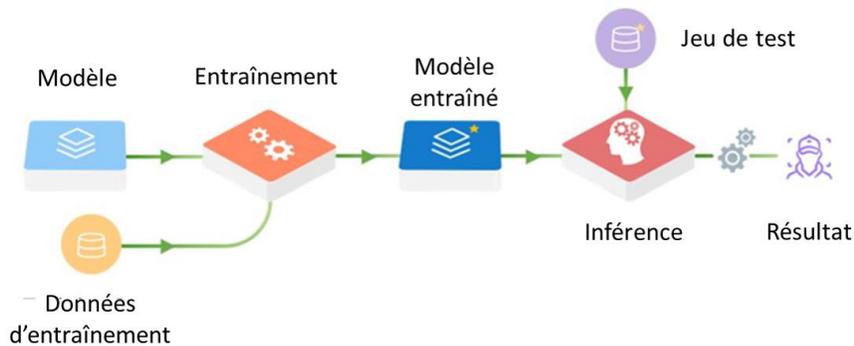
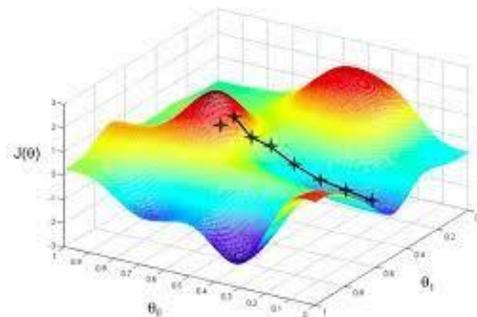
MONITOR THE MODEL



RETRAIN THE MODEL



Gradient Descent



IA cloud centralisé

Grandes compagnies

Gourmant en puissance électrique

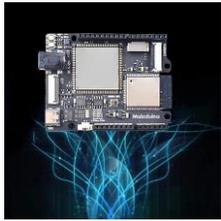
Gourmant en bande passante

Gourmant en quantité de données

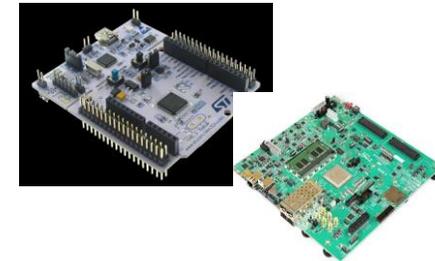
Migration

Réseau de composants (SE) pour de l'IA frontale

Kendryte K210



STM32 Cube AI



AMD-Xilinx



Intel



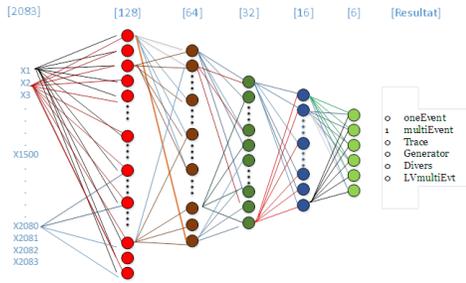
Digilent



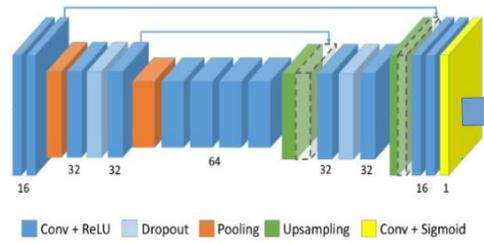
nvidia



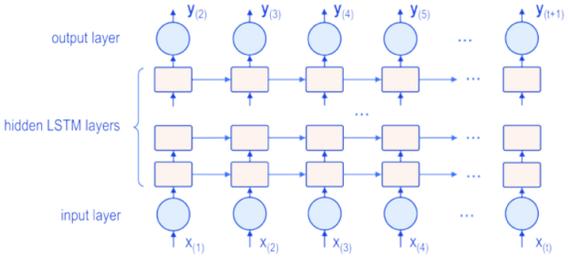
DEEP NEURAL NETWORK



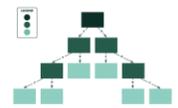
CONVOLUTIONAL NEURAL NETWORK



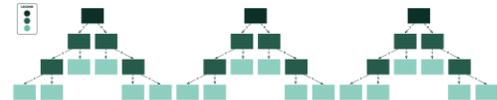
RECURRENT NEURAL NETWORK



DECISION TREE



RANDOM FOREST



Classique

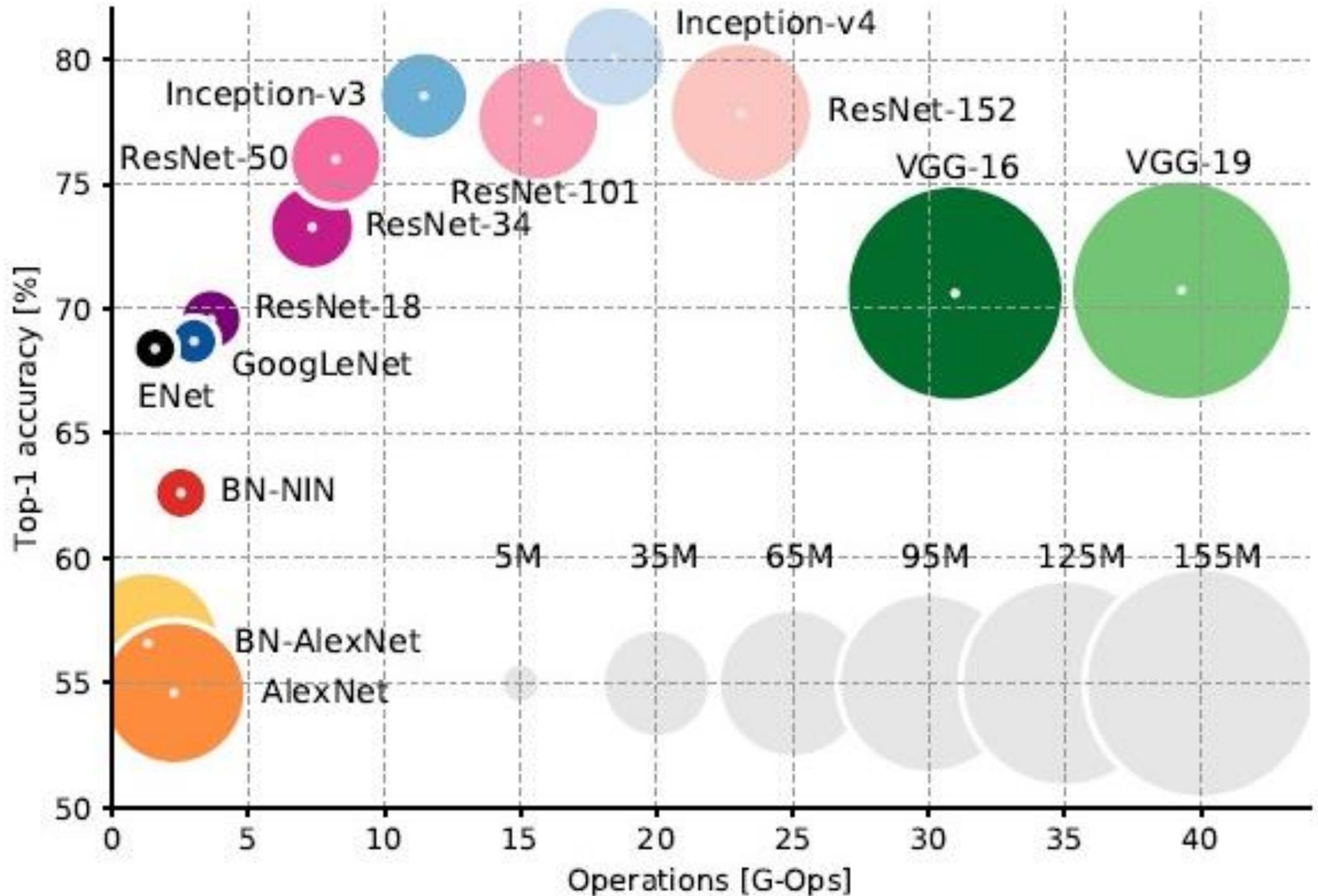
- Séries temporelles (LLM, ChatGPT)
- Données Audio
- Générateur d'information en apprentissage profond
- Apprentissage par renforcement
- Traitement des données avec graphes (chaînes)
- Traitement des images et des vidéos (computer Vision)
- Traitement du langage naturel
- Structuration des données

Embarqué

Off-line
On-line

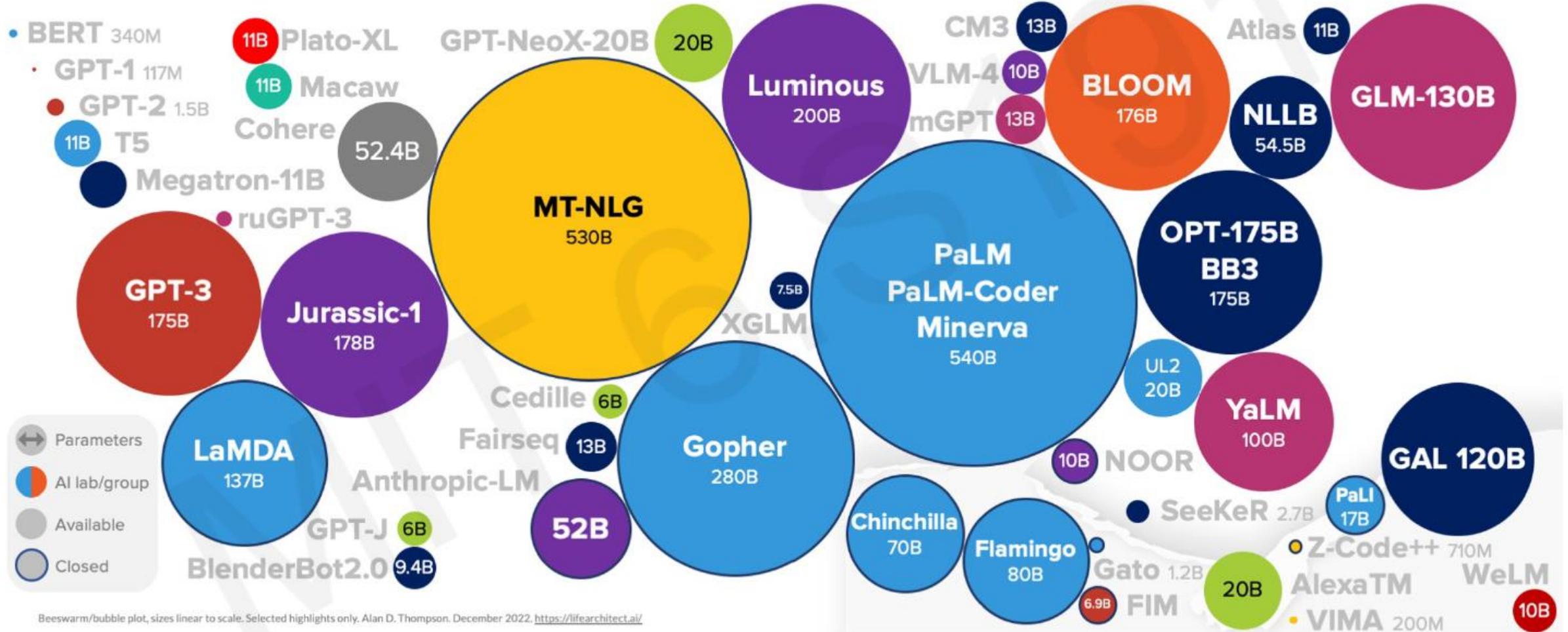
- Signal generation
- Design Optimisation
- Optimisation des Instruments
 - Signal recognition
 - Pile-Up recovery
 - Signal deconvolution
- Selection/ Classification/ Decision
 - Data selection
 - Data parameters prediction
 - Denoizing
- Data Compression
 - Reduced data format

Les autres réseaux : Performances (avant 2020)



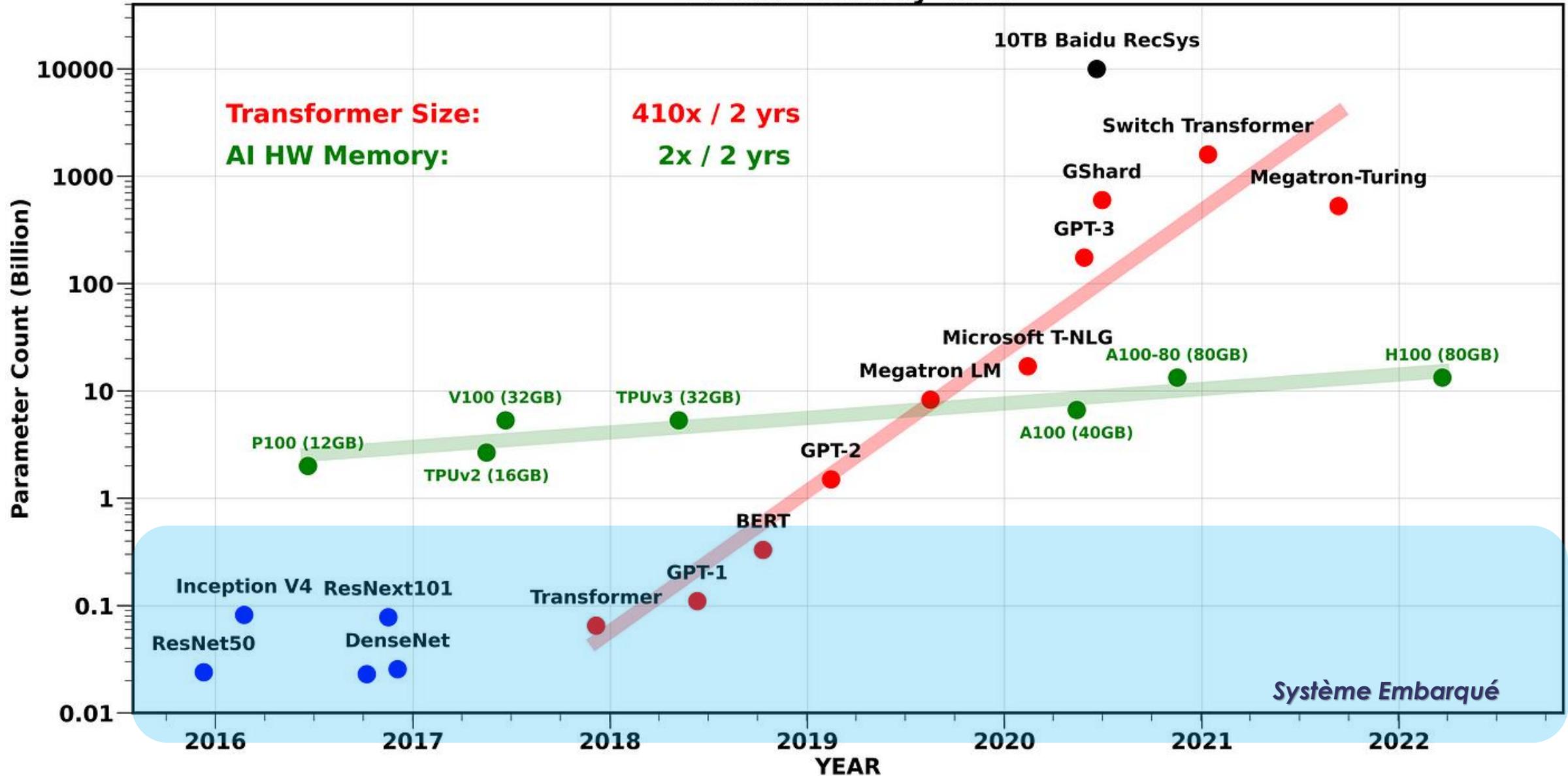
Modern Era of Statistics

Language Models size – up to Dec, 2022

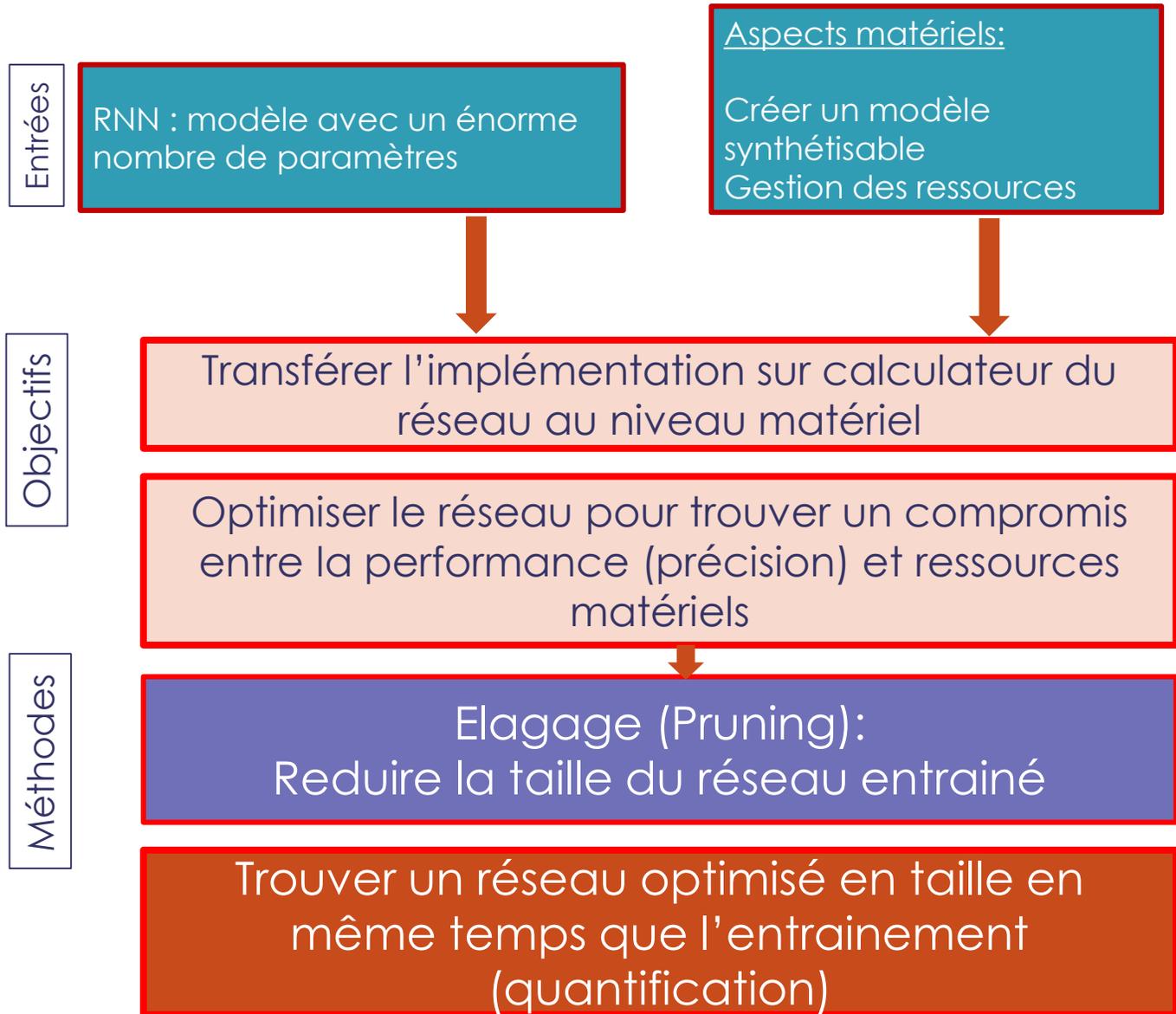
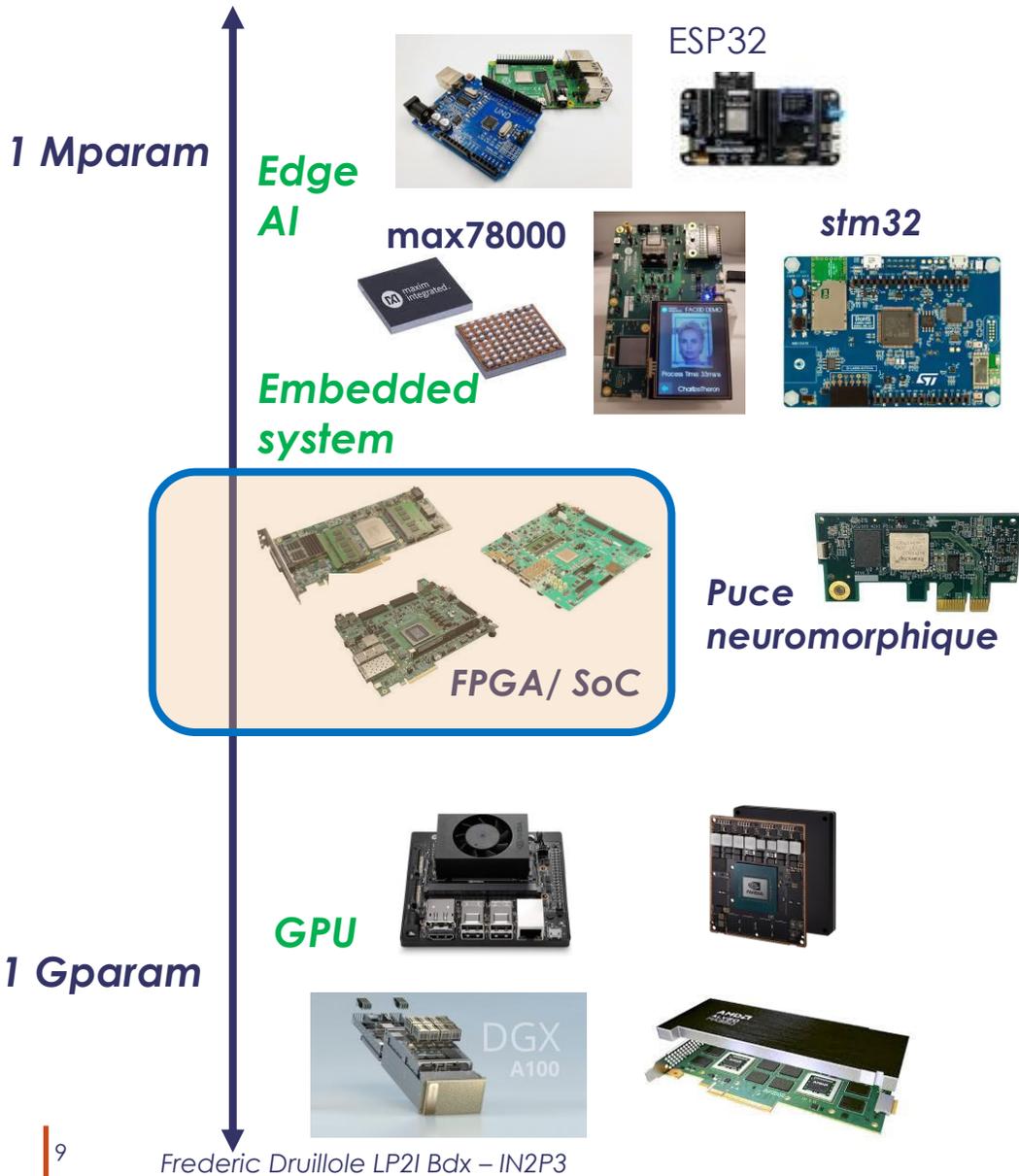


Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. Alan D. Thompson. December 2022. <https://lifearchitect.ai/>

AI and Memory Wall



Choisir la plateforme et Optimiser le modèle



IA embarqué : 2 solutions



Basé sur des fonctions matérielles dédiées aux calculs sans logiciel (calcul matriciel)

ASIC
Circuit neuromorphique
FPGA

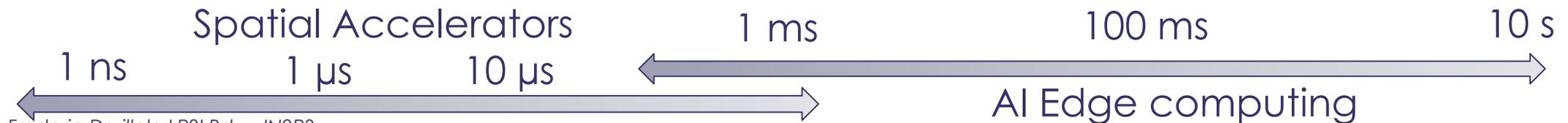
2 En cours de développement



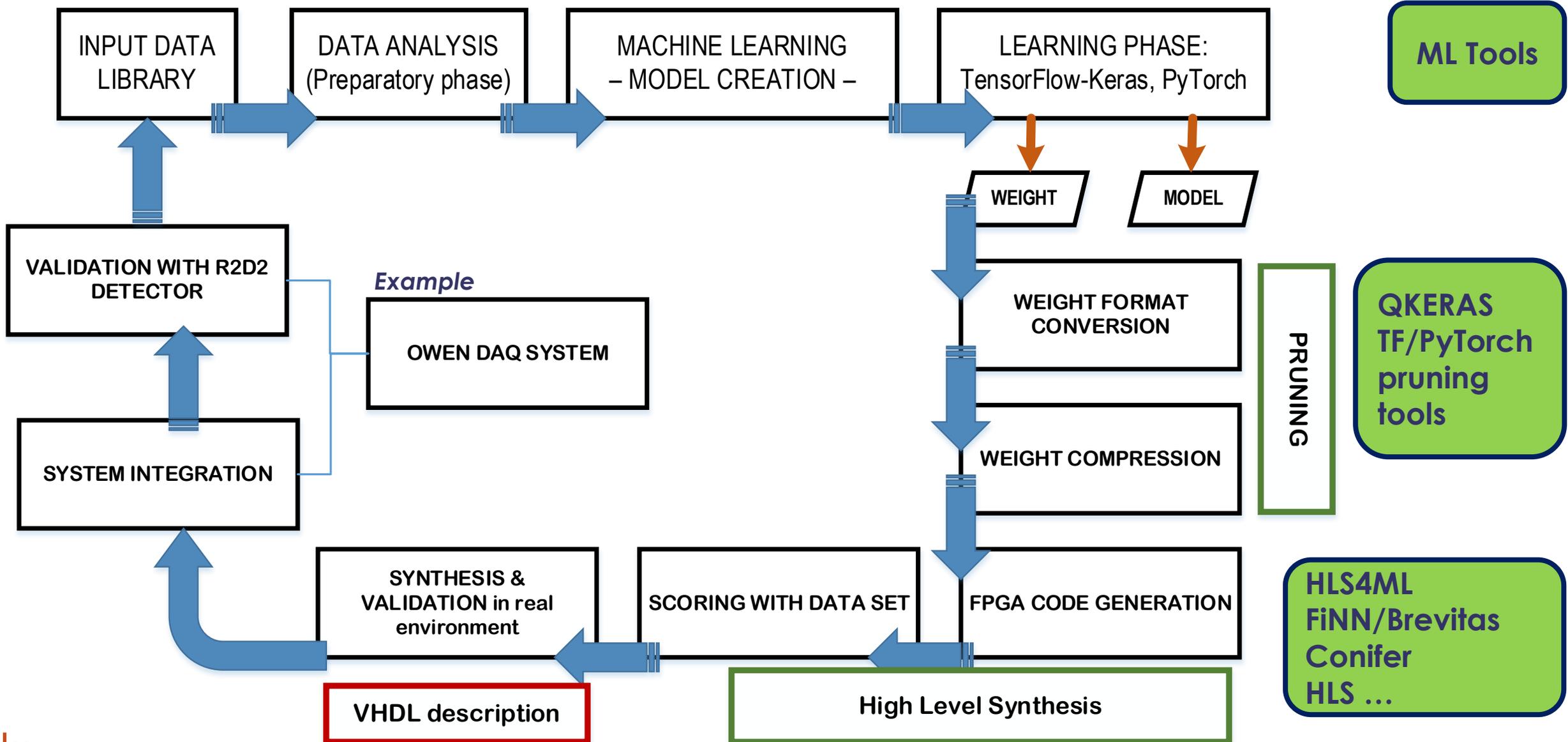
Reste basé sur de la programmation logicielle

MMPA
GPU
FPGA
PU redesigned pour l'IA (TPU, KPU...)

1 Pilote les développements actuels



→ Outils différents suivant la cible

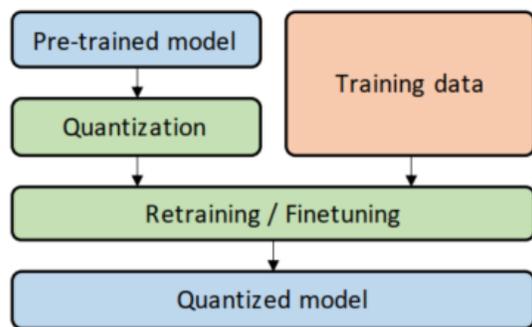
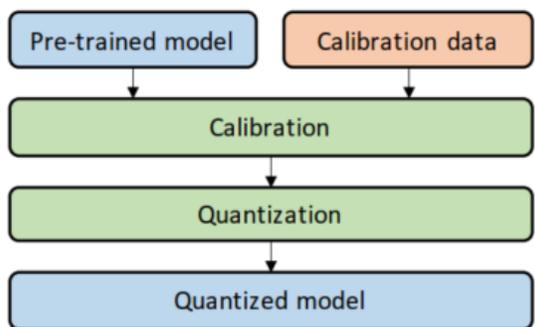


QUANTIZATION

Floating Point \rightarrow integer Fixed Point
32bits \rightarrow 16,8,4,2 bits

PTQ: **Post-Training Quantization**

QAT: **Quantization Aware Training**



Conversion des données et du modèle après apprentissage.

Réentraînement du modèle en ayant convertit la précision des données et du modèle

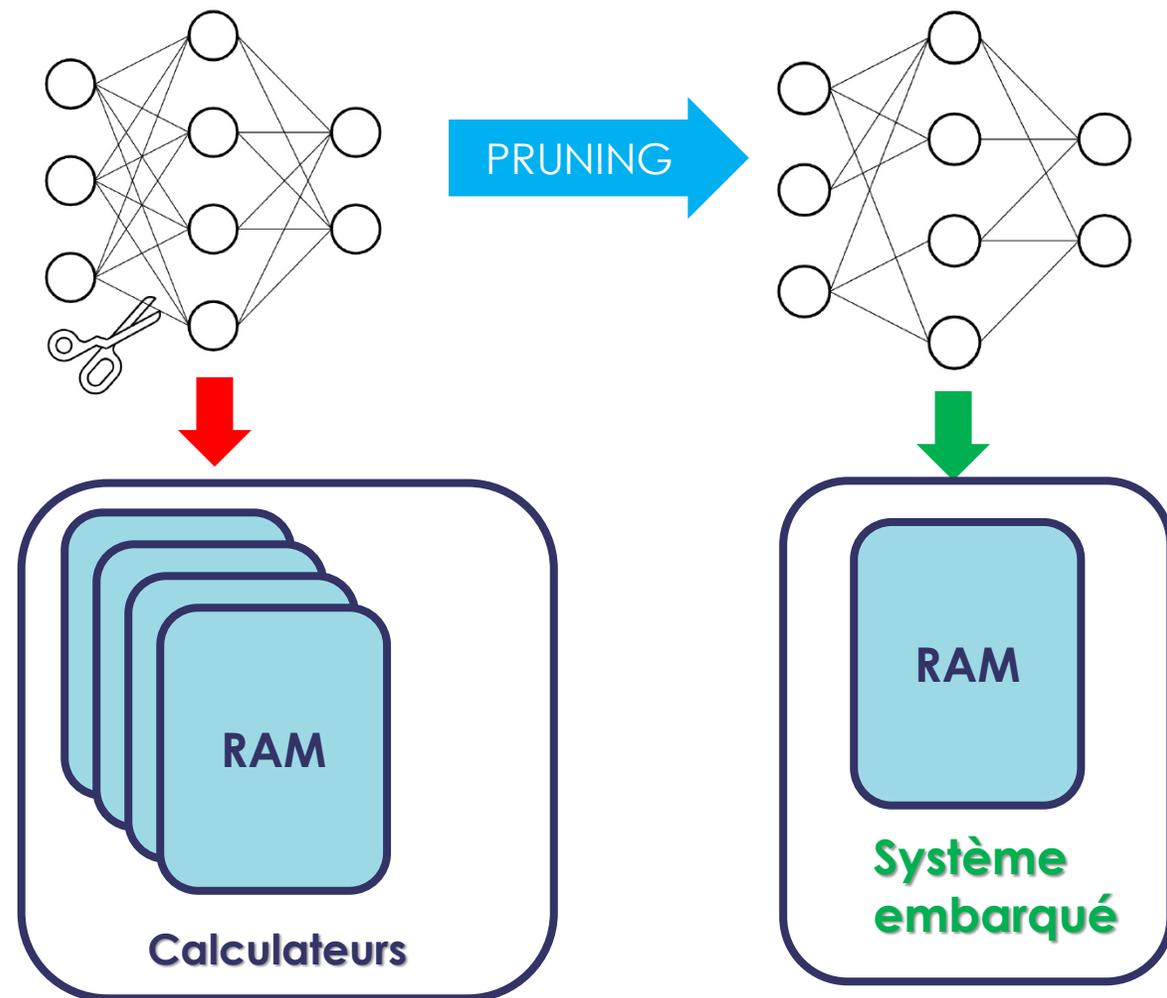
FP32

INT8 <3,5>

 Perte de précision \rightarrow Justesse du modèle

PRUNING

On conserve les poids au dessus d'une certaine valeur :
Elagage



Travail sur la précision du modèle

AI Edge computing

■ **Ax8i (max78000)**

■ **ST Microelectronics**

■ CubeIDE + Cube-AI

■ **VITIS AI (AMD)**

■ **OpenVino (Intel)**

■ **N2D2 – DeepGreen
Eclipse Aidge**

Solution
Européenne
Open

■ **Edge Impulse
(société)**

Spatial Accelerators

■ **VIVADO HLS**

■ HLS4ML

■ FiNN (Xilinx R&D)

■ ChaiDNN

Consortium ?

■ **HADDOC2**

■ **DNNWeaver**

■ ...

Framework servant à réaliser l'accélération des calculs des modèles IA pour l'inférence.

Utilise des modèles éprouvés

Difficile de créer des inférences sur des composants à faibles ressources (multi-voies, proche des capteurs)

Création d'un framework dédié à Xilinx et Intel en VHDL

Critères de comparaison des matériels

Critères	CPU	CPU+AI-CoPro	GPU	FPGA	ASIC
Adaptabilité (à une variété de modèles)	Elevé	Elevé	Moyen	Faible	Aucune
Puissance de calcul	Faible	Moyen-Elevé	Elevé	Elevé	Moyen
Latence	Moyen (ms)	Moyen (ms)	Faible (μ s)	Faible (10 ns)	Très faible
Flux d'entrée	Faible	Moyen	Elevé	Elevé	Elevé
Parrallelisme	Faible	Moyen	Elevé	Elevé	Elevé
Efficacité de la consommation électrique	Moyen	Moyen	Moyen	Moyen	Elevé
Mise en Oeuvre	Facile	Moyen	Moyen	Complexe	Très Complexe

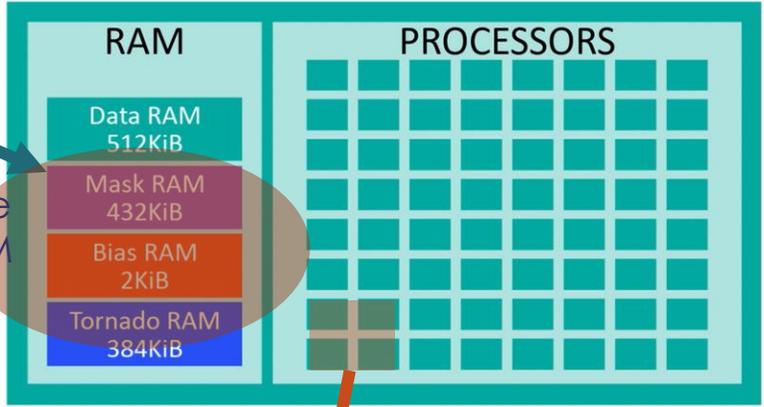
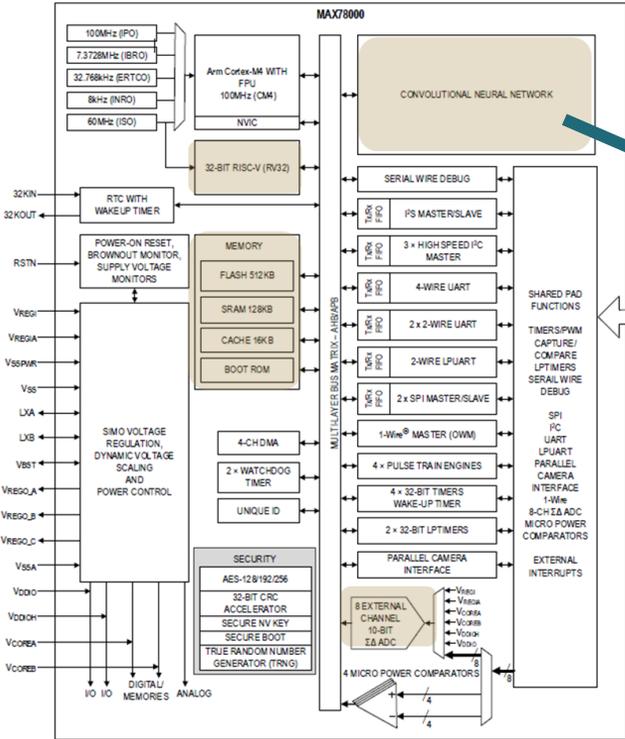
Savoir-Faire



Architecture Max78000



MAX78000 CNN Accelerator

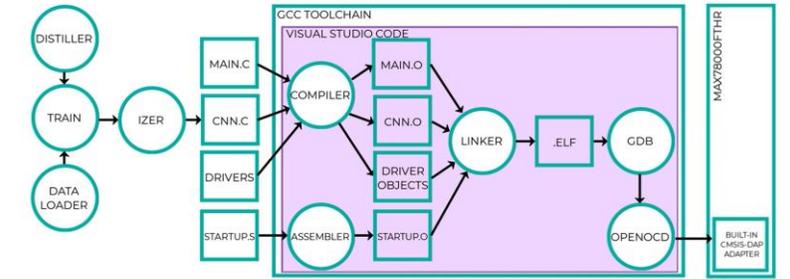


Modèle en RAM

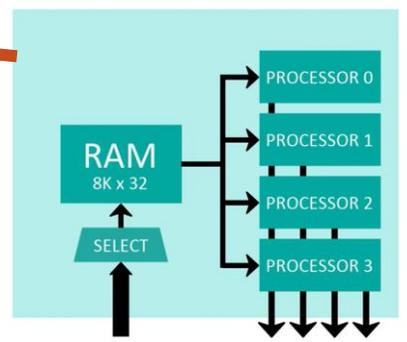
Génération du code C++ cnn.h et main.c
 → Chargement en mémoire du modèle, machine d'état pour réaliser l'inférence



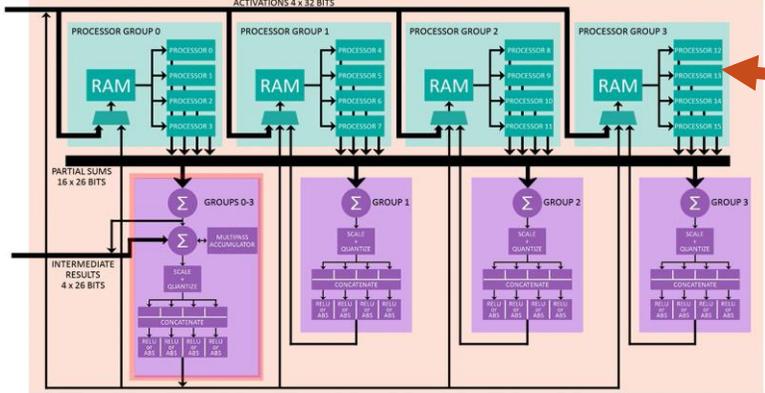
A Choice of Wrappers



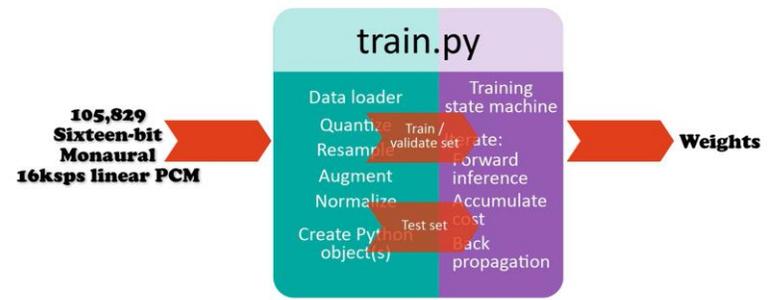
The Processor Group



The Processor Quadrant



Training KWS70v3



Données Y et X

Modèle PyTorch

Apprentissage classique

Architecture
Poids du réseau

Outil sous Linux

On réécrit le modèle avec la bibliothèque ai8x

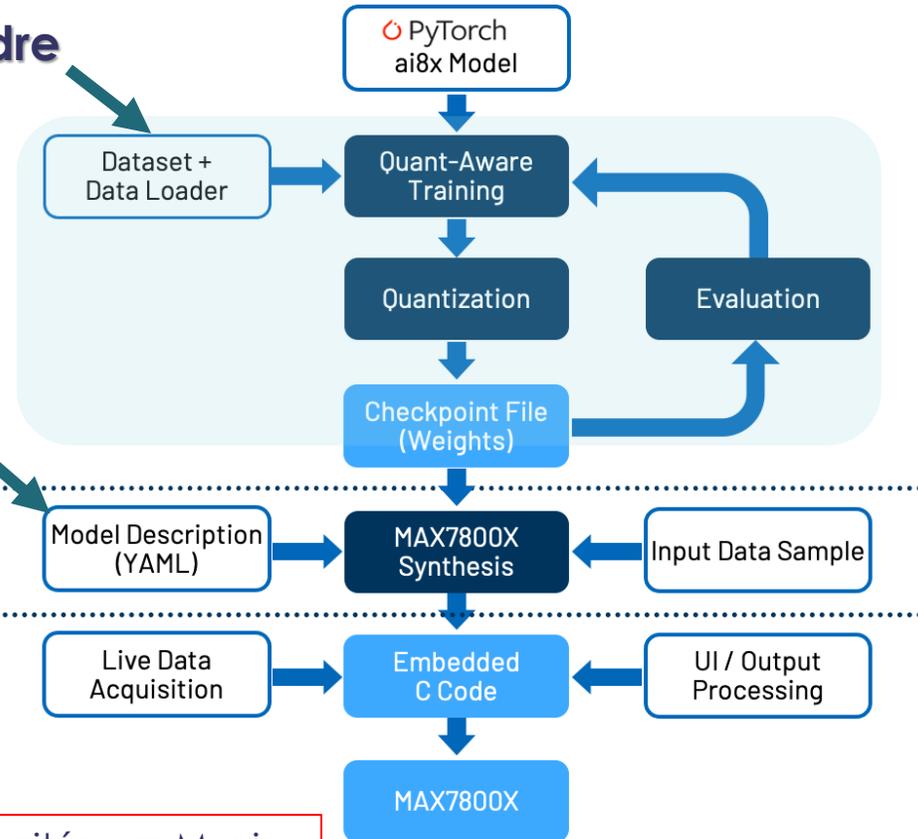
Fonctions à apprendre

① Training
ai8x-training

Syntaxe à apprendre

② Synthesis
ai8x-synthesis

③ Deployment
Embedded SDK



Difficile d'adapter les outils à des modèles non traités par Maxim
Reconstruire les scripts Python
Retro-engineering à réaliser
Validation de la chaine à refaire

hls4ml pipeline

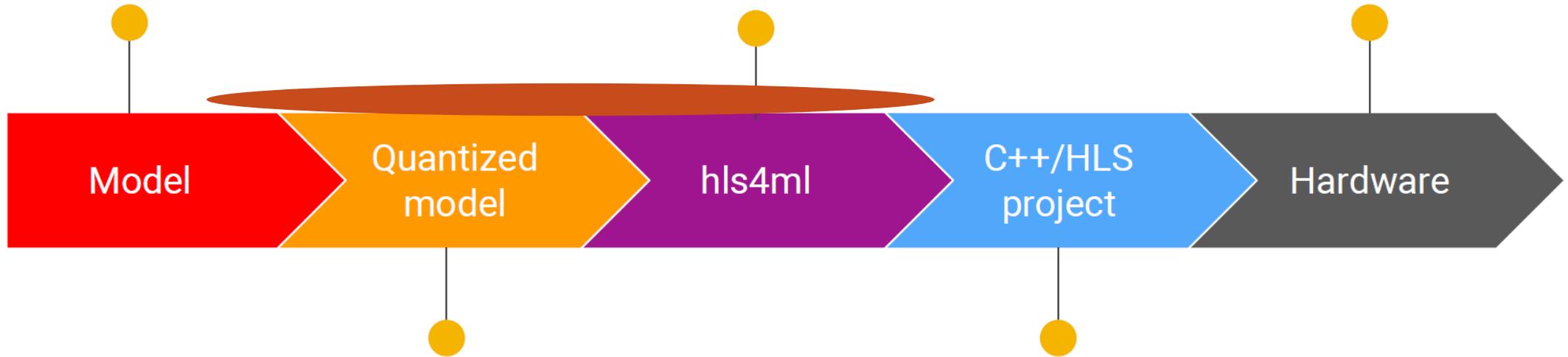
Supported DL frameworks:

-  Keras
-  PyTorch
-  ONNX

Model conversion,
optimization, profiling &
tuning



Xilinx FPGAs, Intel/Altera
FPGAs, Intel x86 CPUs



Quantization and pruning
techniques:

- [QKeras + AutoQ](#) (Keras)
- [Brevitas](#) (PyTorch)



• **Zynq + HLS4ML: io_parallel / io_stream / Reusefactor / Resource / Latency**

	ARTY-Z7 CH1	ARTY Z7 CH3	ARTY Z7 CH3 Optimisé	ZCU102 CH3	CH4	ZCU102 CH4 Qbit<16,2>	ZCU102 CH5	ZCU102 CH7 Optim Ressource	ZCU102 CH7 Optim Latence	ARTY Z7 CH7 Optim HLS Stream
Nombre de cellules	16	25801	25801	25801	658951	658951	156951	156951	156951	156951
Horloge de reference	4,166ns	9,408ns	9,408ns	4,396ns		4,369ns	4,369ns	4,369ns	4,028ns	9,410ns
Temps de latence	70ns	19,804us	19,804us	2,510us		5,510us	5,510us	10,010us	10,015us	141us
BRAM	0%	41%	8%	6%		36%	18%	9%	15%	72%
DSP48E	21%	115%	11%	10%		59%	59%	12%	24%	6%
FF	2%	85%	28%	10%		28%	22%	32%	53%	167%
LUT	1%	280%	68%	46%		103%	113%	81%	104%	423%
URAM	0%	0%	0%	0%		0%	0%	0%	0%	0%



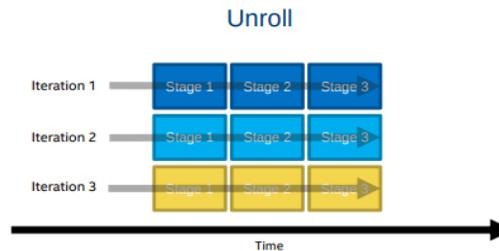
ARTY Z7
ZCU104



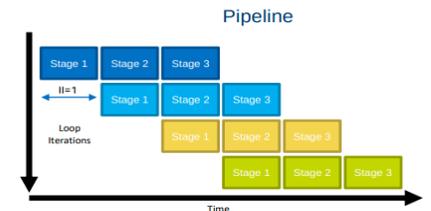
Tensorflow/Keras
+ HLS4ML
VITIS HLS



*depending on VITIS-HLS #pragma
The way HLS handles vector/matrix before DSP*

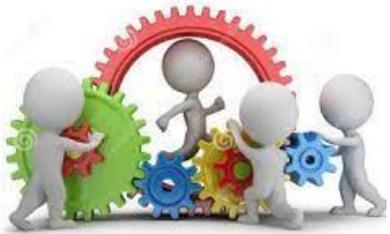


Question of HLS optimisation



What if directly written in VHDL ?

Conclusion : challenges du ML en Système embarqué



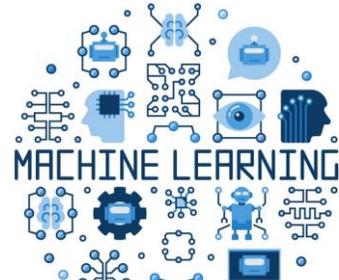
Roles & competencies

- **Data Physicist**
- **System Engineering team**
- **ML Engineer**
- **Software Engineer**
- **Hardware Engineer**
- **Infra & Security teams**

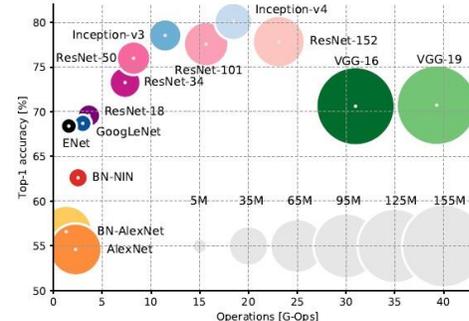


Tools

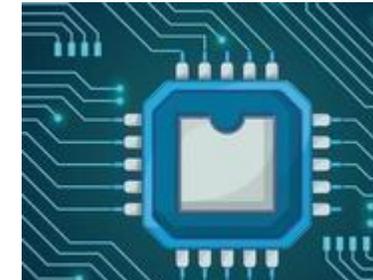
- **ML Tools:**
 - **TF-KERAS, PyTorch ...**
- **HLS4ML (Xilinx...)**
- **HLS**
- **Brevitas & FiNN(Xilinx)**
- **CONIFER (LLR)**
- **N2D2 (CEA)**
- **VHDL**
- **C/C++**
- **SDK dédié**



Artefacts & ML zoology



- **Model**
- **Code source...**



Digital hardware technologies

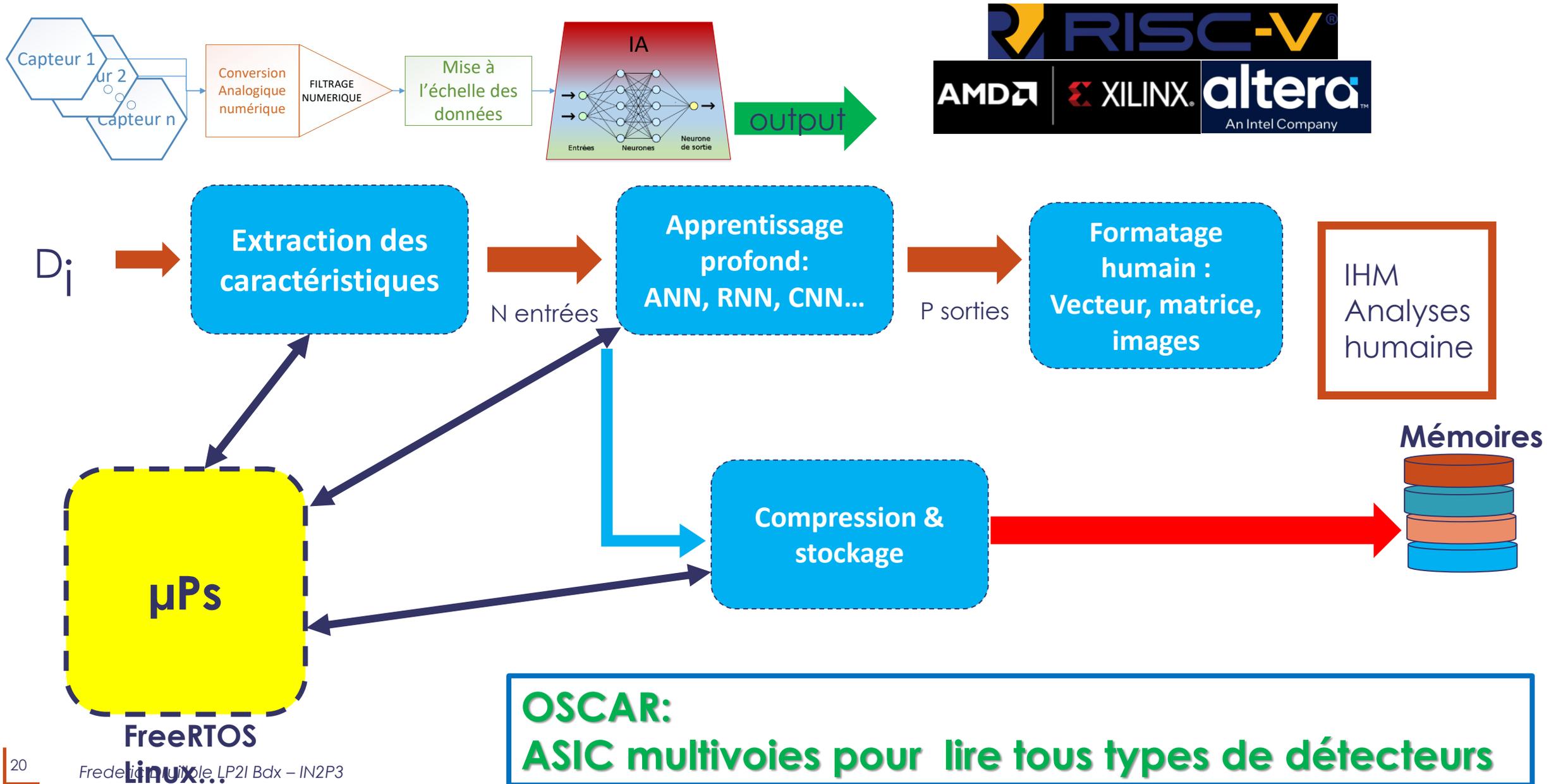
- **CPU + CoPU**
- **FPGA SOM**
- **SNN**
- **MPPA**
- **GPU**
- **...**



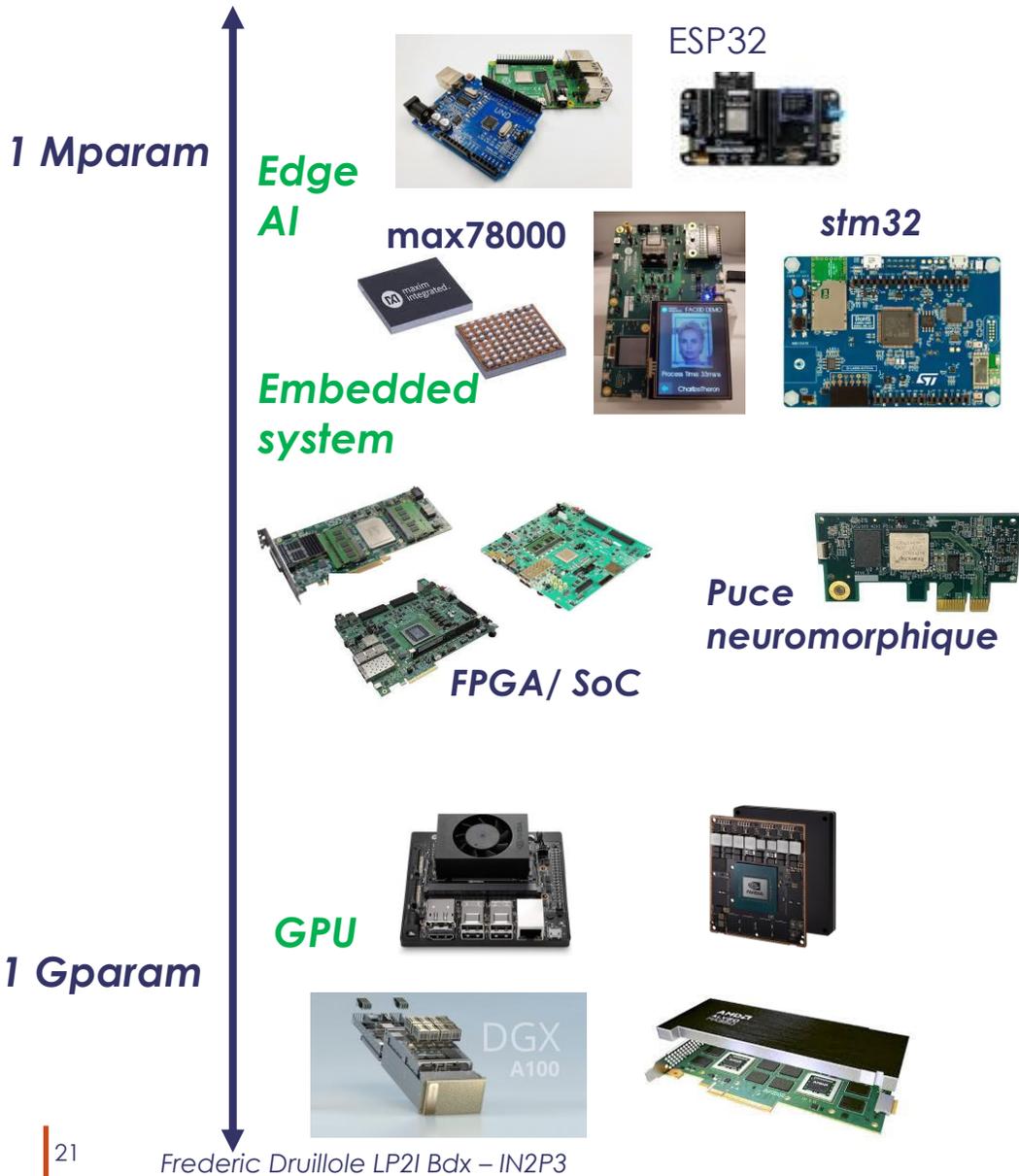
Deployment & Operational AI

- **GitLab/Git**
- **Training Service skew**
- **Model Monitoring**
- **Responsible AI**
- **...**

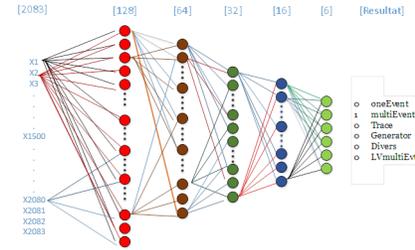
OSCAR: L'architecture numérique générale / idéale



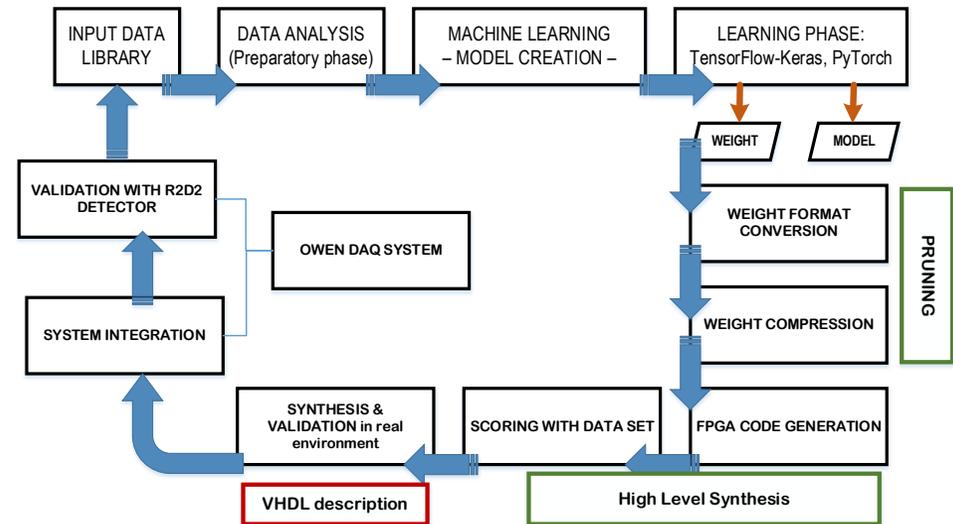
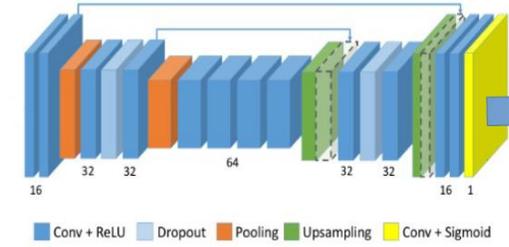
Conclusion



DEEP NEURAL NETWORK



CONVOLUTIONAL NEURAL NETWORK



Quelques conseils

- Ne travaillez qu'avec les modèles dont vous comprenez réellement le fonctionnement mathématique.
- Privilégiez les modèles les plus rapides et les plus simples (DNN et CNN)
- Donnez-vous une estimation de la latence du modèle.
- Tenez compte des transformations des données pour optimiser le modèle (opération supplémentaire)
- Créer des modèles autour de 100k paramètres maximum.