

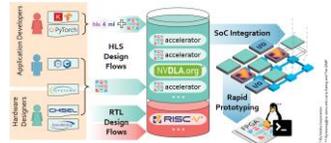
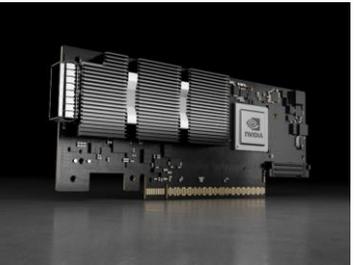
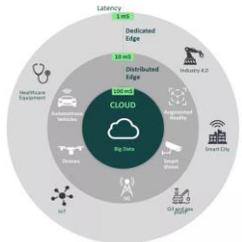
# L'intelligence Artificielle dans les systèmes embarqués

– Principes – Aspects Matériels – Mes recherches –  
– THINK –



Test of Hardware Inferring Neural network

Reseau DAQ – Hardware-Firmware-Software-

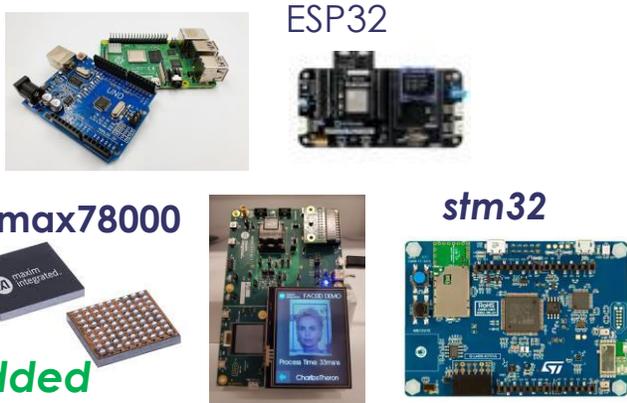


Accelerator generation with hls4ml 	Automatic integration in ESP 	Full-system RTL simulation 	Full-system test on FPGA 
--	----------------------------------	--------------------------------	------------------------------

# Conclusion Phase 1

1 Mparam

Edge AI



Embedded system



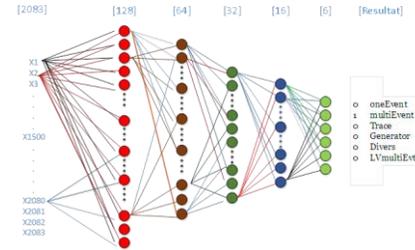
GPU



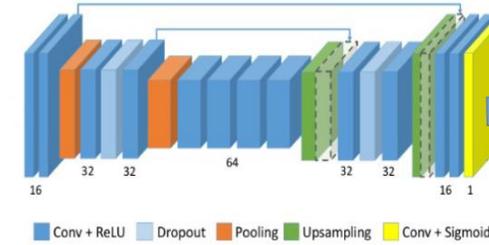
1 Gparam

Frederic Druillolle LP21 Bdx - IN2P3

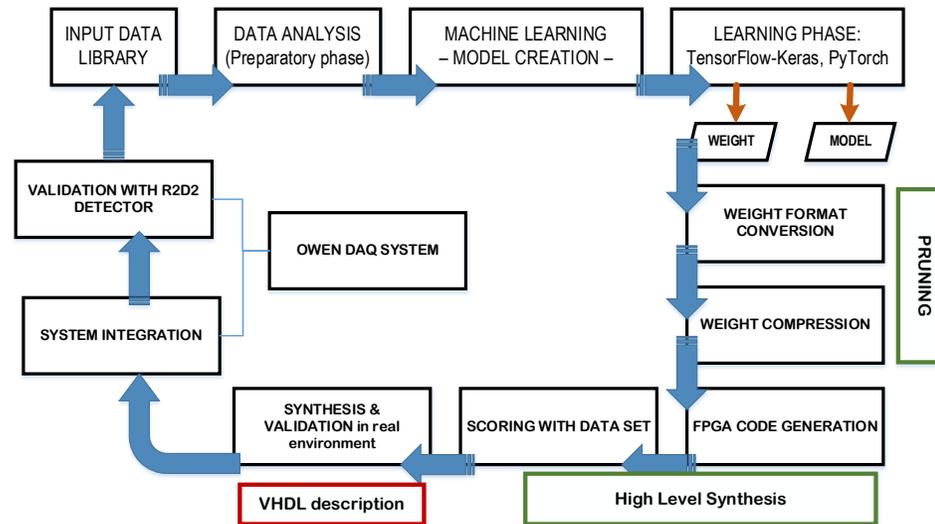
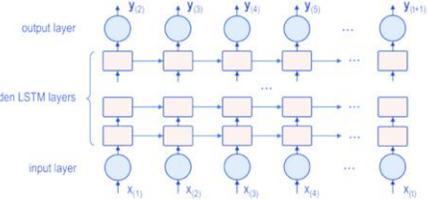
DEEP NEURAL NETWORK



CONVOLUTIONAL NEURAL NETWORK



RECURRENT NEURAL NETWORK



**Elagage (pruning) :**  
réduction du nombre de neurones

**Quantization:**  
Reduction de l'empreinte en mémoire (codage des coefficients et des données)

## IA embarqué : 2 solutions



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

ASIC  
Circuit neuromorphique  
FPGA

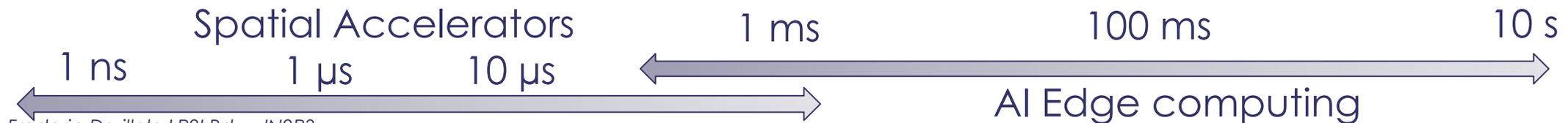
**2** En cours de développement



Reste basé sur de la programmation logicielle

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

**1** Pilote les développements actuels



## 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)

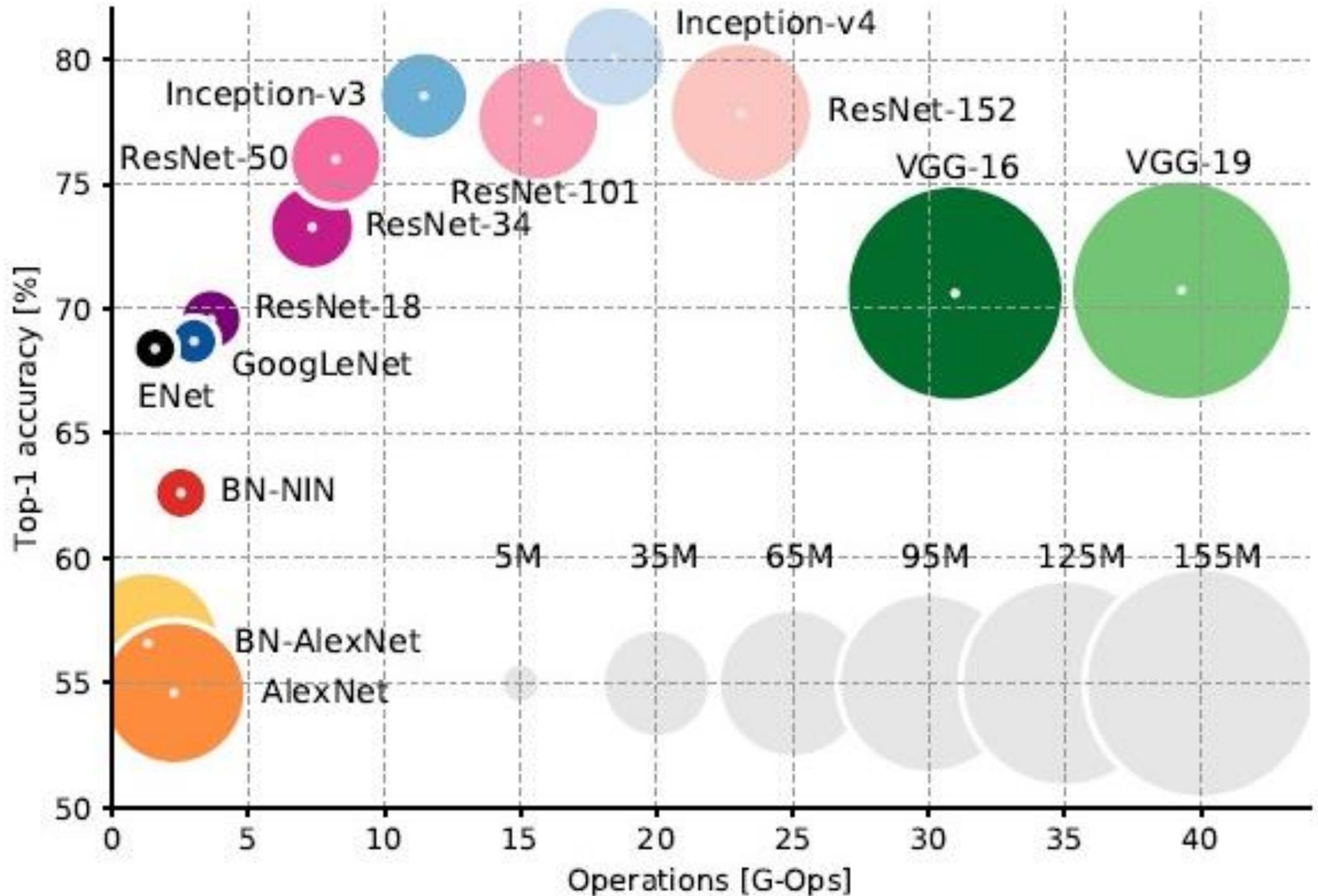
Création d'un framework dédié à Xilinx et Intel en VHDL

# Critères de comparaison des matériels pour l'IA

Critères	CPU	CPU+AI-CoPro	GPU	FPGA	ASIC
<b>Adaptabilité</b> (à une variété de modèles)	<b>Elevé</b>	<b>Elevé</b>	Moyen	<b>Faible</b>	Aucune
Puissance de calcul	<b>Faible</b>	Moyen-Elevé	<b>Elevé</b>	<b>Elevé</b>	Moyen
<b>Latence</b>	Moyen (ms)	Moyen (ms)	<b>Faible</b> ( $\mu$ s)	<b>Faible</b> (10 ns)	<b>Très faible</b>
<b>Flux d'entrée</b>	<b>Faible</b>	Moyen	<b>Elevé</b>	<b>Elevé</b>	<b>Elevé</b>
<b>Parallélisme</b>	<b>Faible</b>	Moyen	<b>Elevé</b>	<b>Elevé</b>	<b>Elevé</b>
Efficacité de la consommation électrique	Moyen	Moyen	Moyen	Moyen	<b>Elevé</b>
<b>Mise en Œuvre</b>	<b>Facile</b>	Moyen	Moyen	<b>Complexe</b>	<b>Très Complexe</b>
Encombrement/Densité des NN	Moyen	<b>Faible</b>	<b>Faible</b>	<b>Elevé</b>	<b>Elevé</b>

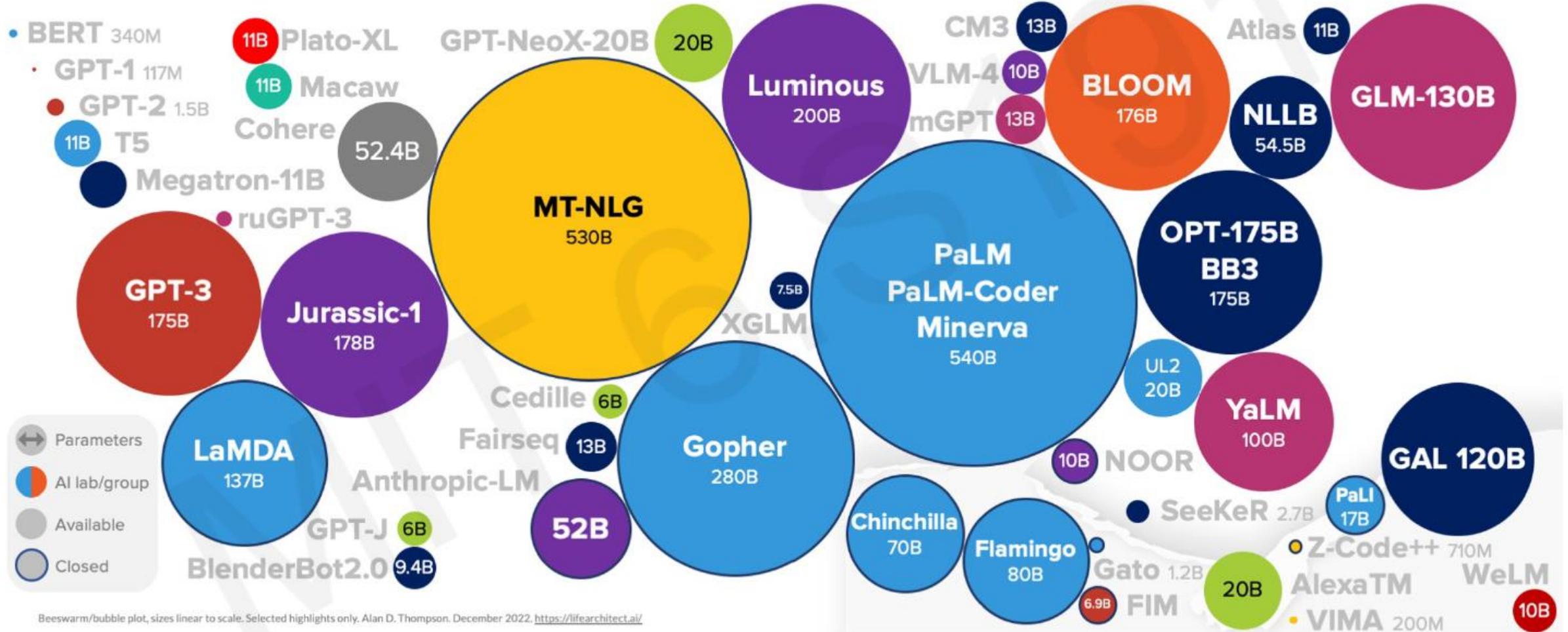


## Les autres réseaux : Performances (avant 2020)



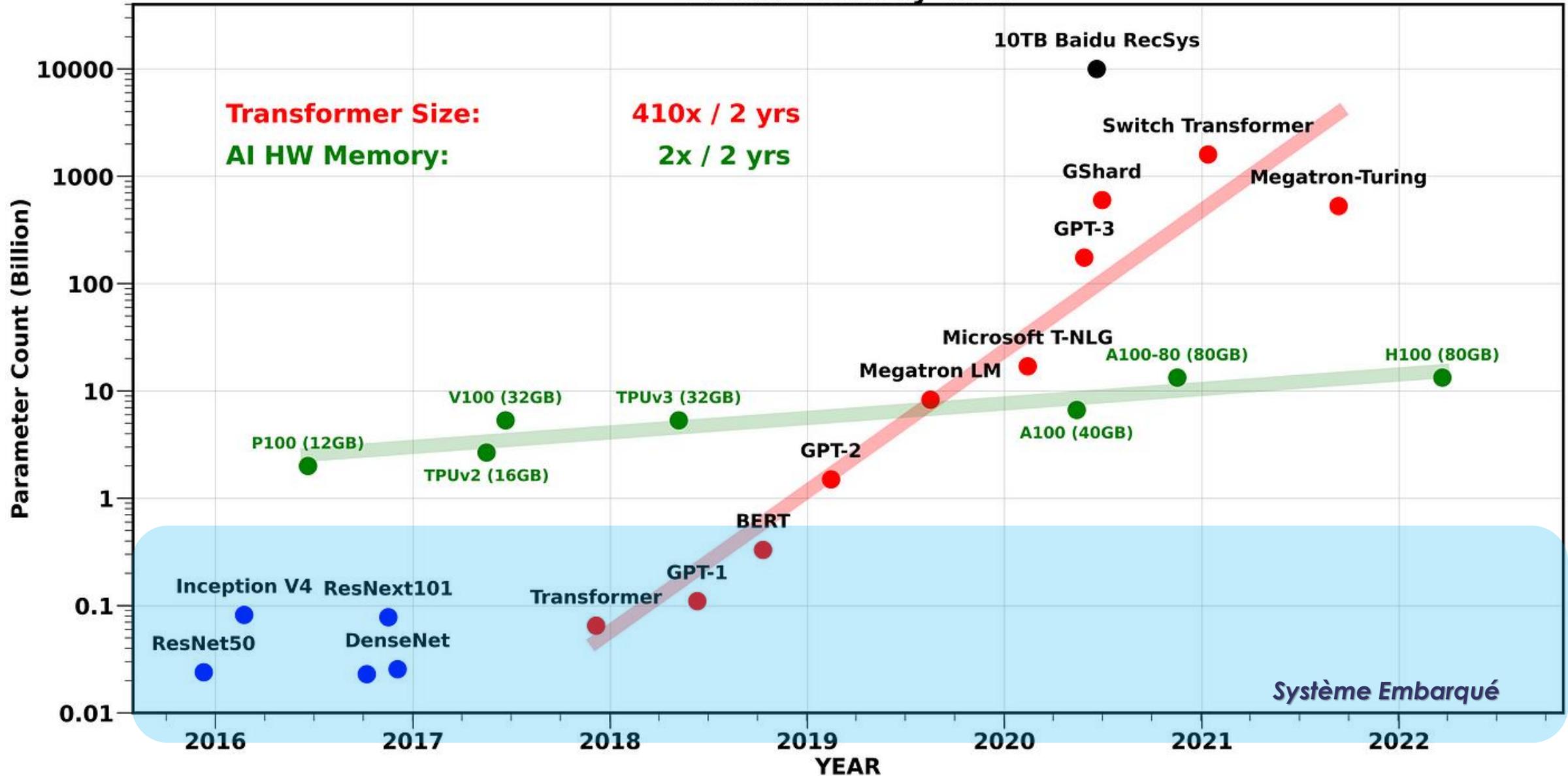
# Modern Era of Statistics

Language Models size – up to Dec, 2022

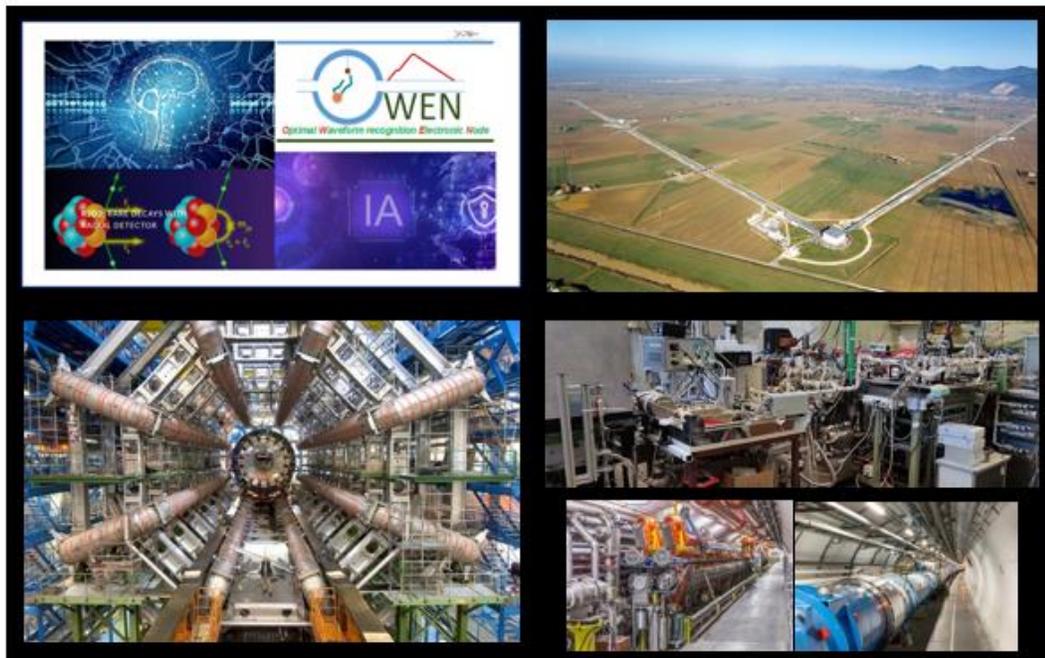


Beeswarm/bubble plot, sizes linear to scale. Selected highlights only. Alan D. Thompson. December 2022. <https://lifearchitect.ai/>

# AI and Memory Wall



*L'ambition de THINK est de comprendre et évaluer la possibilité d'implanter des algorithmes de réseaux de neurones relativement en amont pour nos futurs besoins en instrumentation.*



F. Magniette, M. Melenec



G. Aad, R. Bertrand



S. Viret, G. Galbit, Q. David (2025), X. Chen



F. Druillolle, B. Mansoux, J. Domange



S. Geiger, W. Perrin, M. Fozé, J. Wurtz

# Conclusion Phase 1

1 Mparam

**Edge AI**

ESP32 *Think2-1erAnnée*

max78000

stm32

**Embedded system**

**FPGA/ SoC**

**Puce neuromorphique**

→ Plateformes utilisées

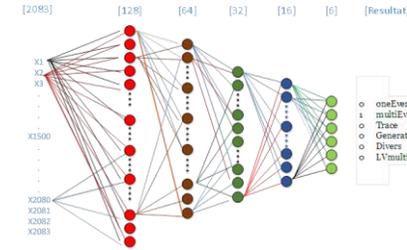
1 Gparam

**GPU**

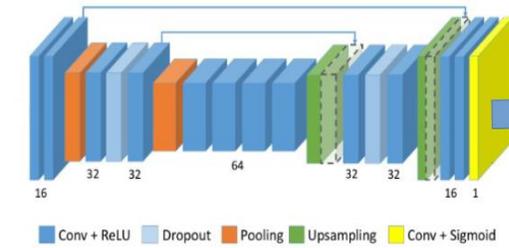
DGX A100

Frederic Druillolle LP21 Bdx - IN2P3

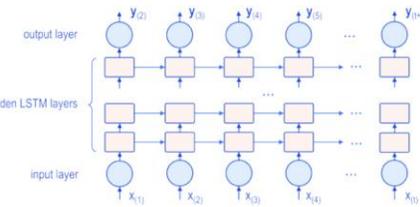
DEEP NEURAL NETWORK



CONVOLUTIONAL NEURAL NETWORK

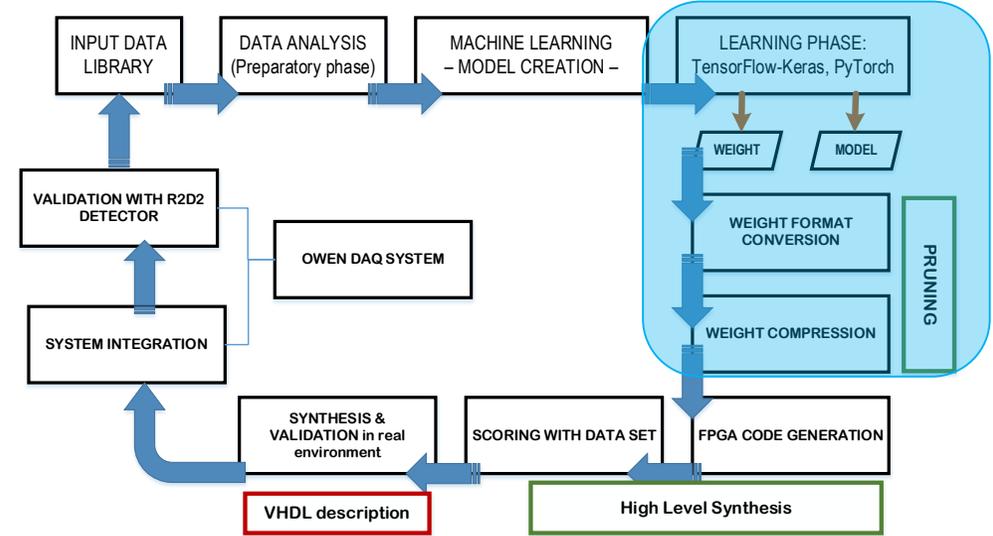


RECURRENT NEURAL NETWORK



Think Phase 2

→ Optimisation des modèles



Thèmes	Responsable	Labos
Les particules de hautes énergies	George Aad	CPPM Presentation AG DI2I 2024
La physique des neutrinos	Frédéric Druillole	LP2I Bordeaux Presentation AG DI2I 2025
Les ondes gravitationnelles	Sébastien Viret	IP2I Lyon Presentation AG DI2I 2025
L'instrument optimisé	Jocelyn Domange	LP2I Bordeaux



Statut THINK2 Presentation Journées des Metiers de l'électronique Juin 2025 Strasbourg

### Veille technologique

- μProcesseurs et coProcesseur IA
- RISC-V et IA

### Méthode d'optimisation des réseaux par optimisation bayésienne:

- Tutorial fait par Frédéric Magniette
- Test applicatif au CPPM et IP2I

### Stratégies d'entraînement

### Travail sur le LAr d'Atlas

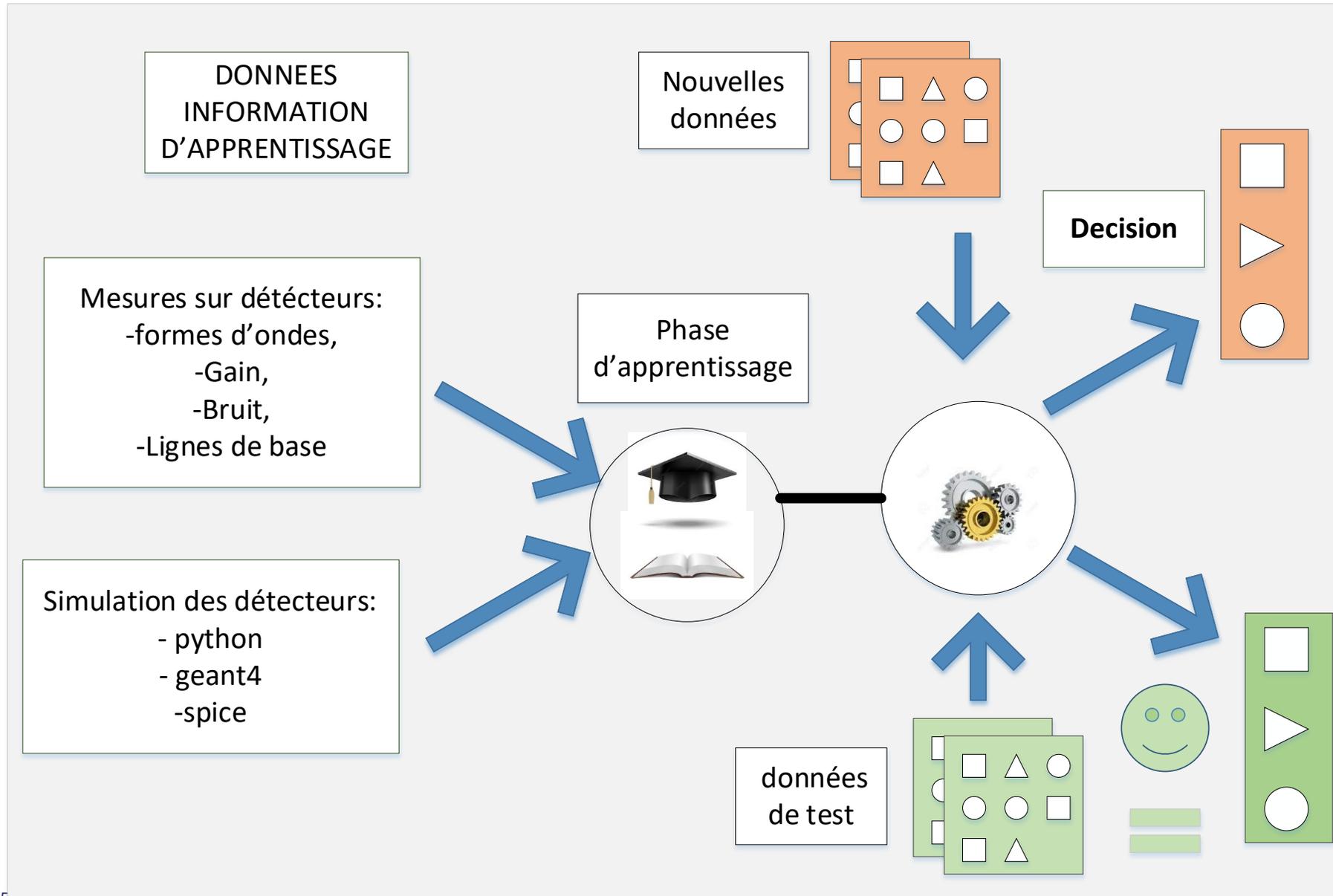
- Amélioration des mesures avec presence de PileUp
- G. Aad et R. Bertrand
- Concours IR CPPM

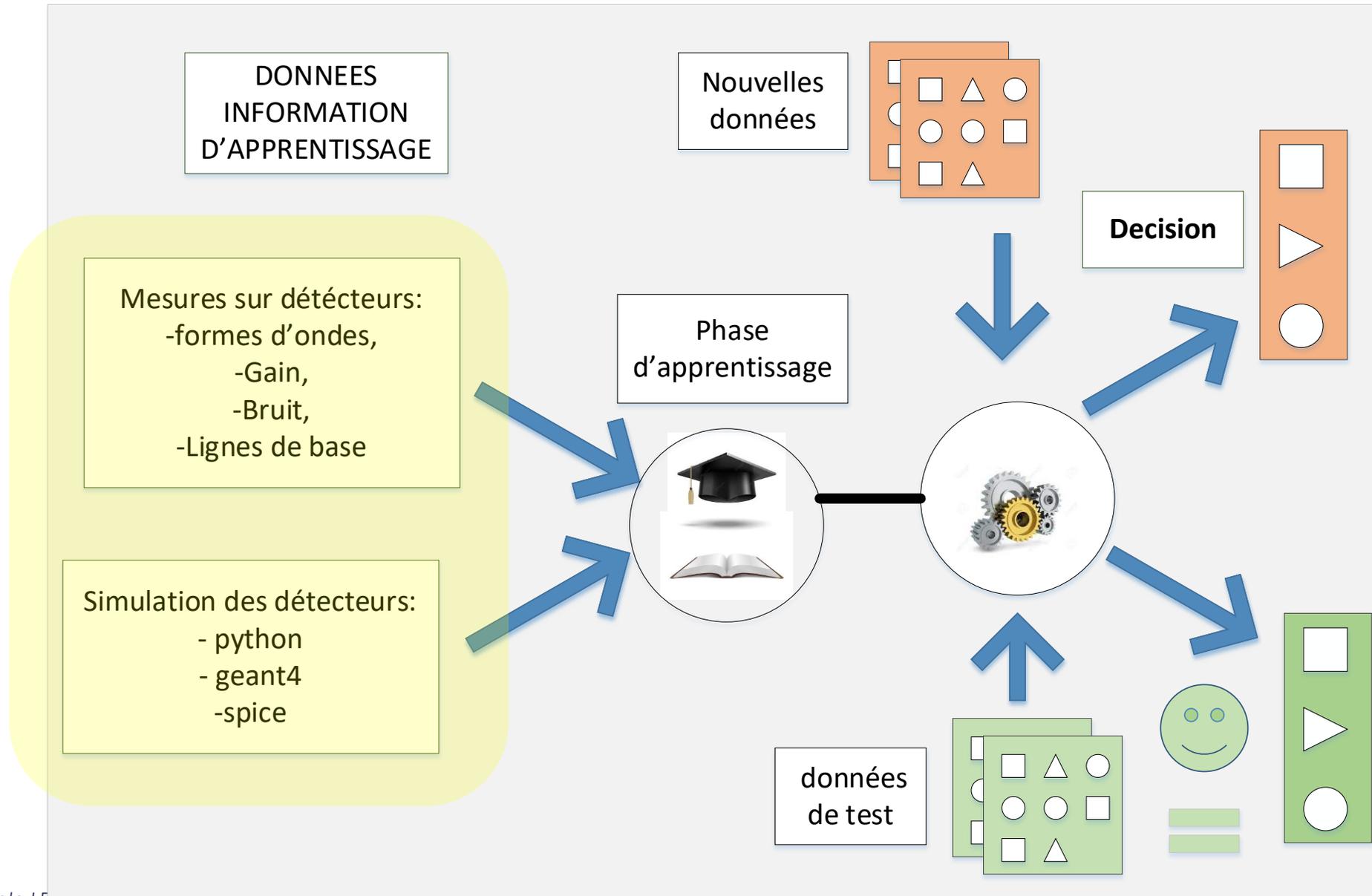
### Travail sur OG

- Optimisation du VHDL par python (G. Galbit DI2I)
- Modèle Liquid Neural Network (S. Viret)

## Actions entreprises dans THINK2

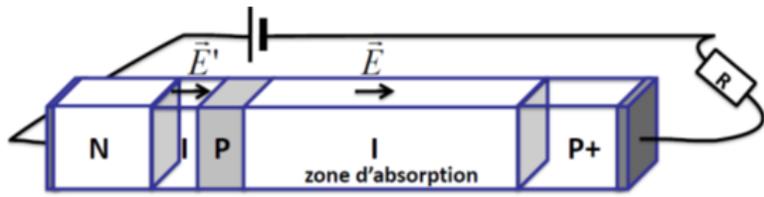
- ➔ Travail sur R2D2 (Neutrino)
  - Modèle ANN avec des mesures
  - Génération de signaux par VAE et GAN
  - Création de nouveaux modèles (détection d'anomalie...)
- ➔ Travail sur Juno (Muons):
  - utiliser un réseau de neurones pour le trigger de niveau 2 du Top Tracker
  - Détecter des muons
- ➔ Création d'un espace GitLab
  - Stockage des données d'apprentissage
  - Stockage des tutoriaux de modèles (Jupyter lab)
- ➔ Diffusion et acculturation
  - Demos
  - Présentation
- ➔ Modélisation détecteurs (J. Domange)
  - Données d'apprentissage
  - Prise de données réelles (Chambre à Brouillard, APD, SiPM)





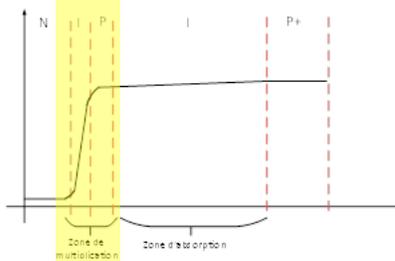
# Création de données pour l'apprentissage: Modélisation de détecteur APD en python

→ Multiplication des électrons par avalanche du signal de ionisation dû au photon

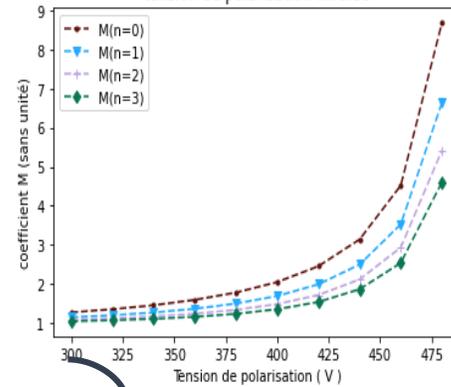


Coefficient de multiplication empirique

$$M = \frac{1}{1 - \left(\frac{V_{bias}}{V_{claquage}}\right)^n}$$

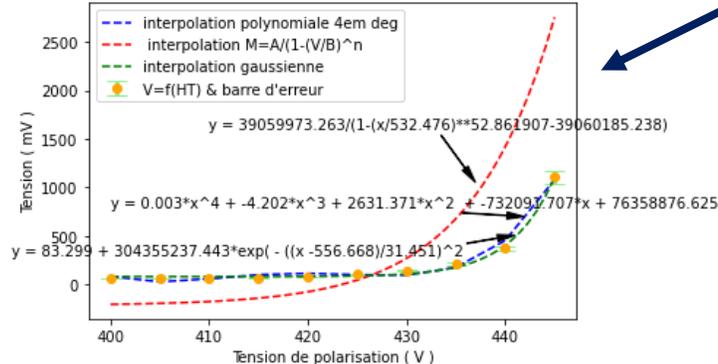


Coefficient M en fonction de la tension de polarisation inverse

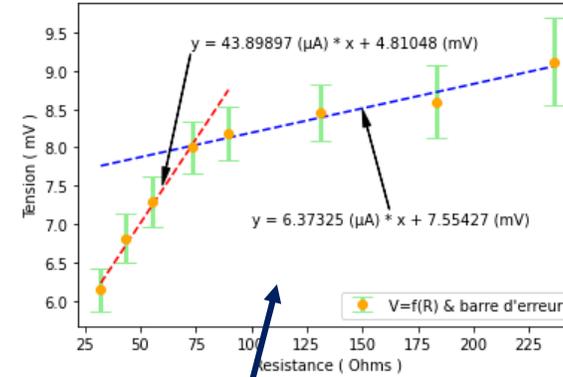


Mesures

Tension de charge de 2000 photons (400nm) suivant la valeur de la resistance de charge avec une Haute Tension = 440 V



Tension de charge de 2000 photons (420nm) suivant la valeur de la resistance de charge avec une Haute Tension = 440 V

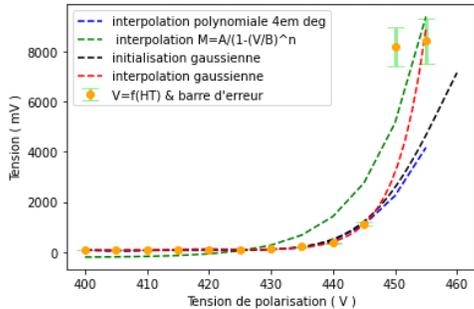


$$I_{ph} = M(V_b) \cdot R(\lambda) \cdot \frac{N \cdot q}{h \cdot \nu} \cdot e^{-\frac{t}{\tau}}$$

Limitation de la BP par Rcharge:

$$T = R_{charge} \times C_{APD}$$

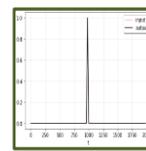
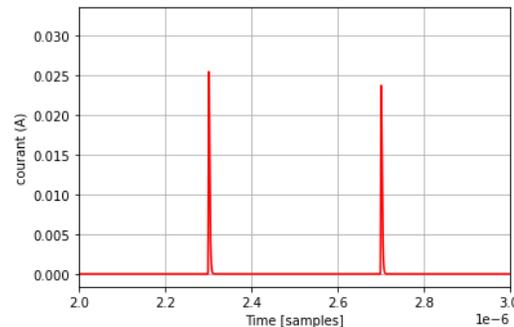
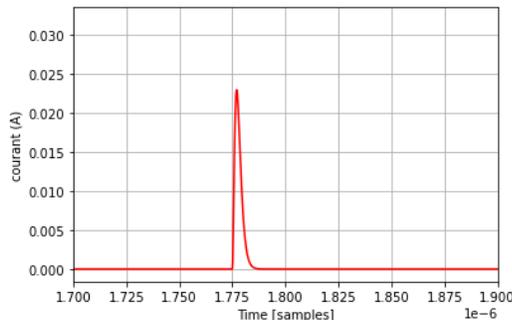
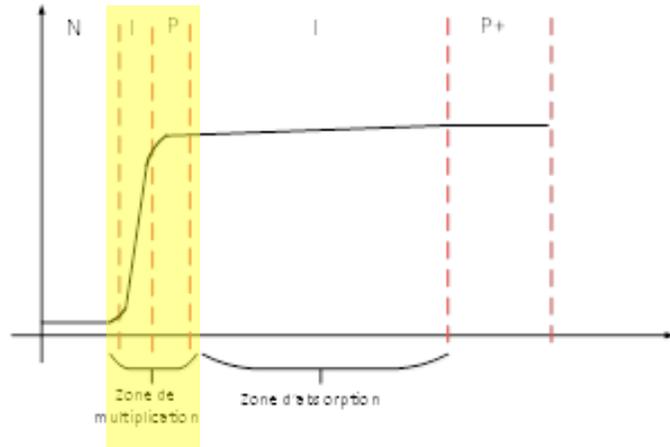
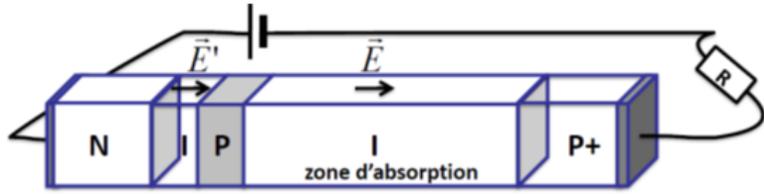
Tension de charge de 2000 photons (400nm) suivant la valeur de la resistance de charge avec une Haute Tension = 440 V



Interpolation de  $M(V_{bias})$  par la montée d'une fonction gaussienne marche mieux.

# Création de données pour l'apprentissage: Modélisation de détecteur APD en python

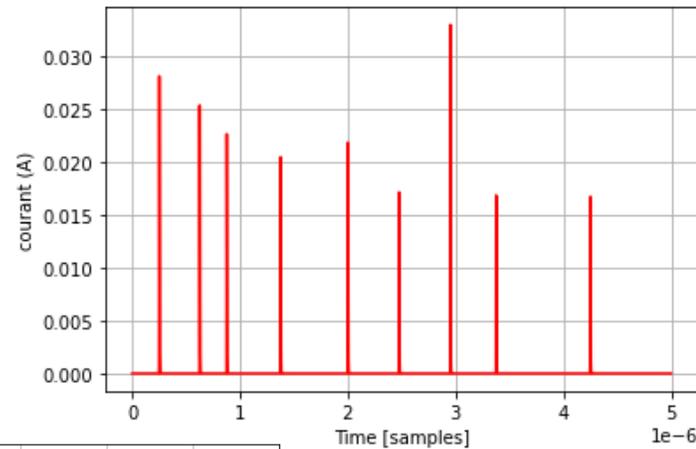
→ Multiplication des électrons par avalanche du signal de ionisation dû au photon



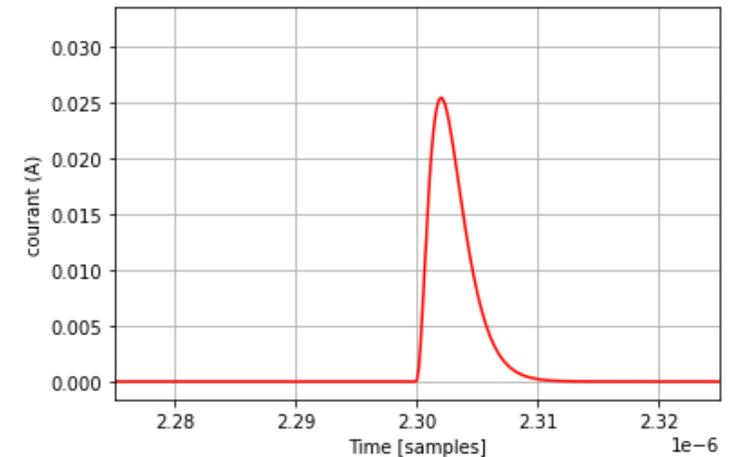
$$q \cdot R \cdot M(V_b) \cdot QE(\lambda) \cdot \frac{E_{\text{photon}}}{E_i \cdot \Delta t}$$

Modèle complet

$$\frac{12}{(1 + 1e^{-9p})^3} \cdot \text{Offset} \rightarrow V_{\text{ph}}$$

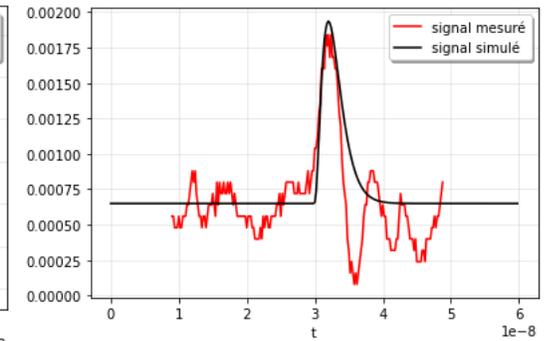
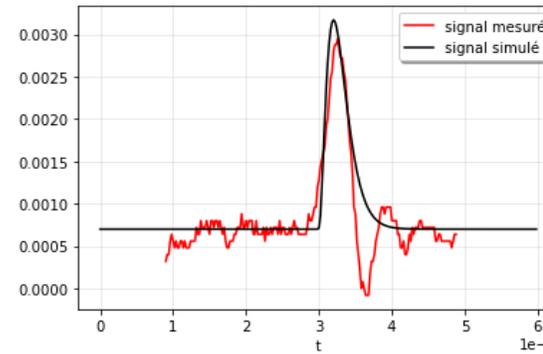
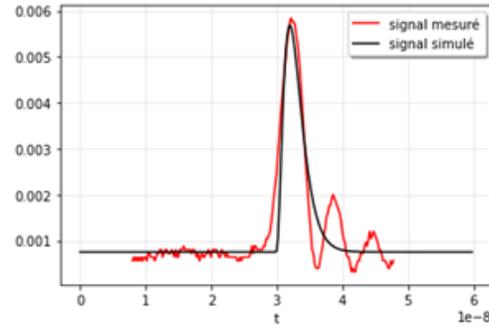
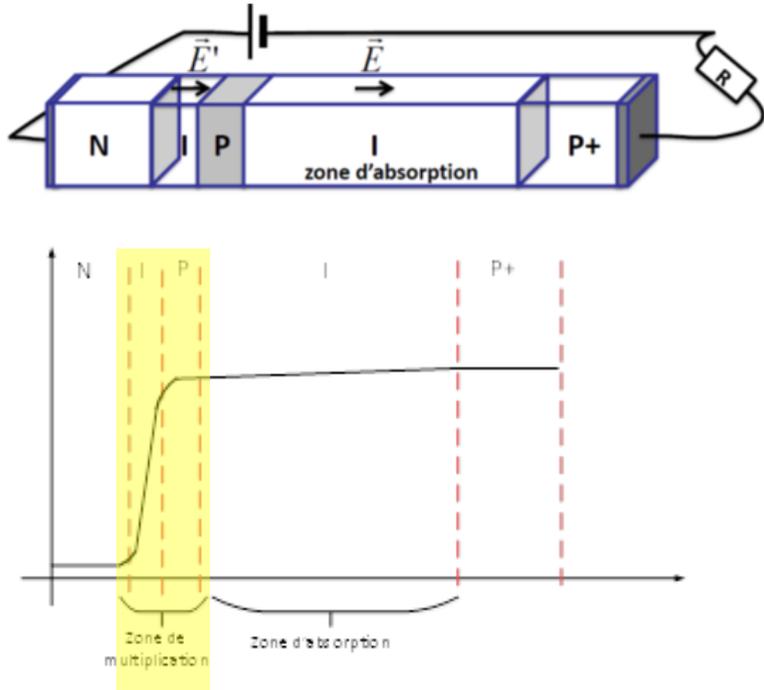


Temps par tirage poissonnien  
Amplitude suivant Energie

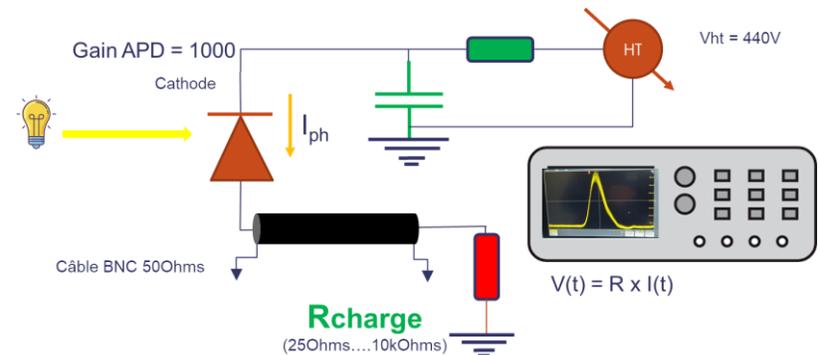
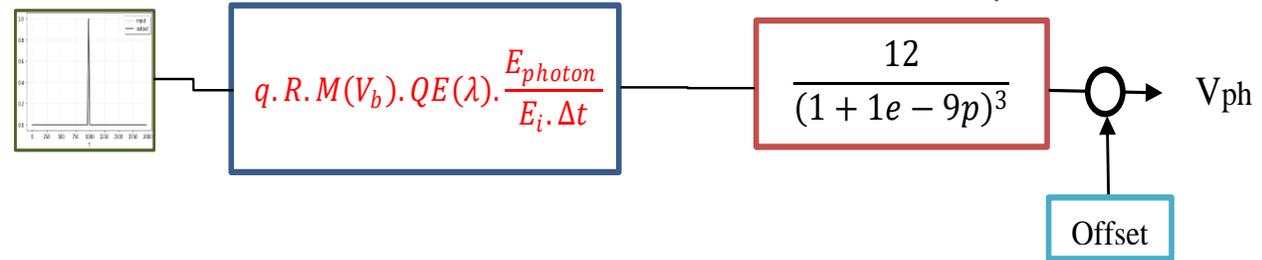


# Création de données pour l'apprentissage: Modélisation de détecteur APD en python

→ Multiplication des électrons par avalanche du signal de ionisation dû au photon

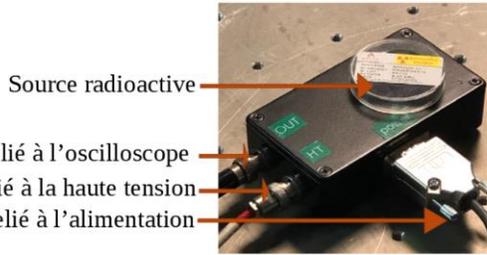
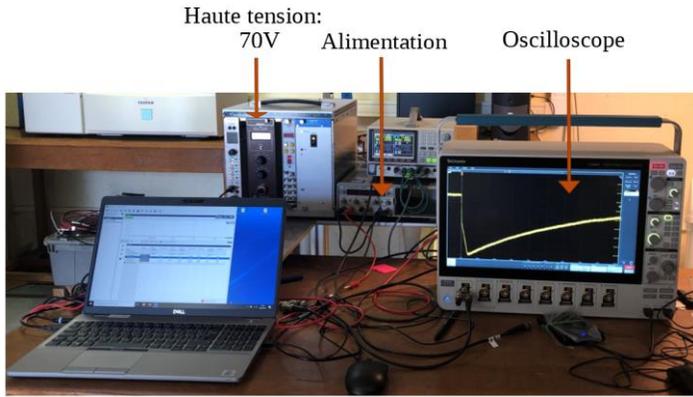


Modèle complet



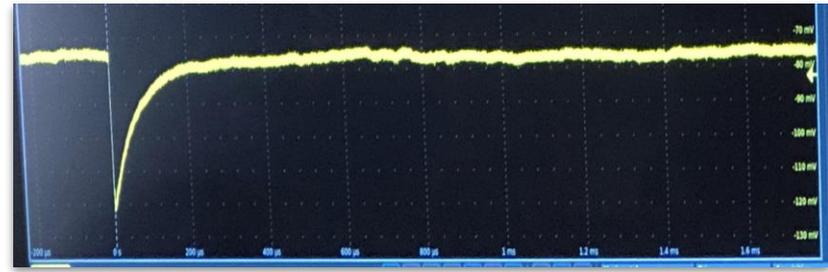
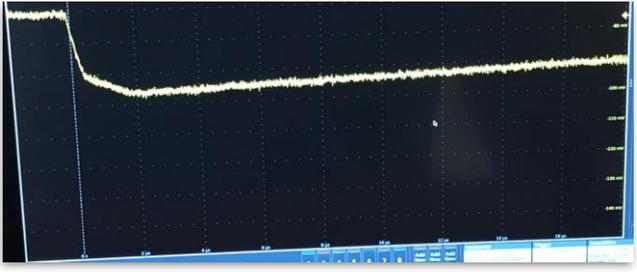
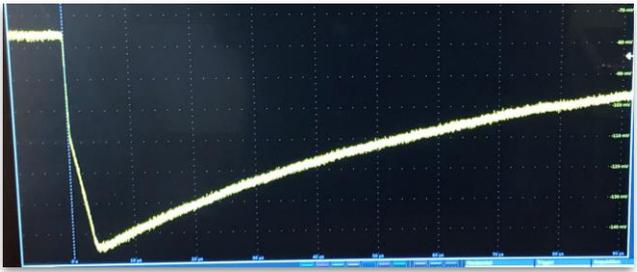
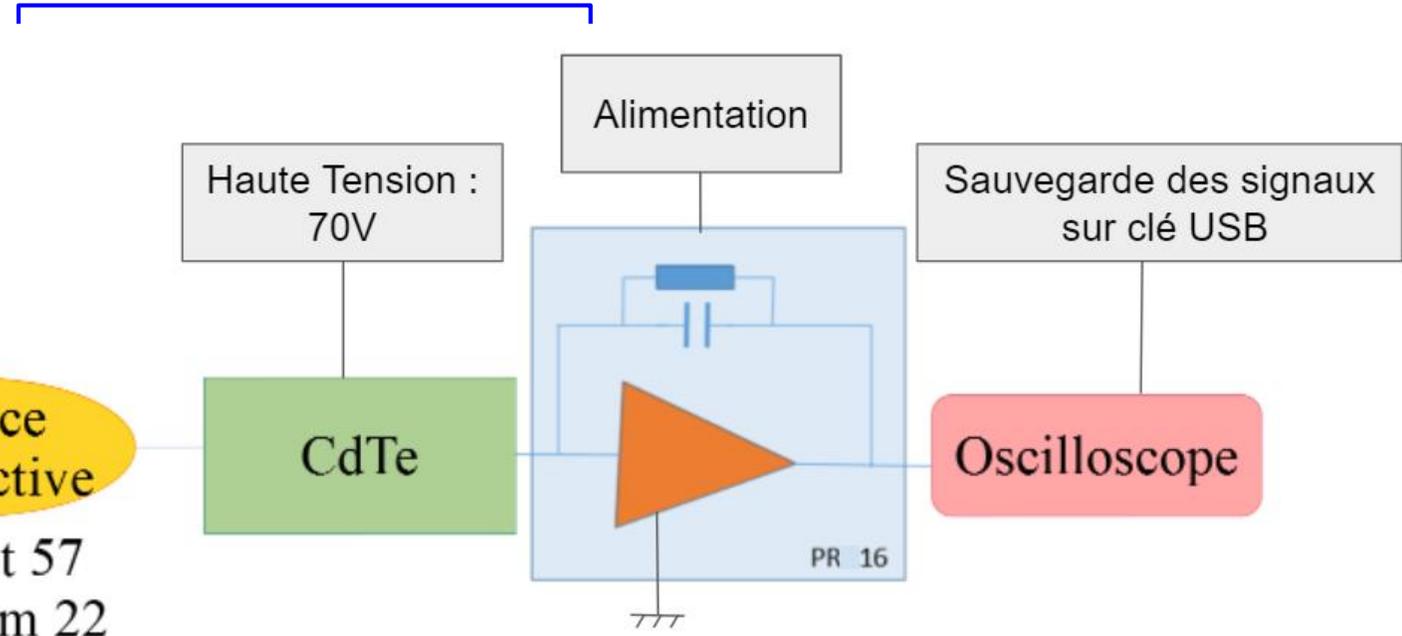
# Création de données pour l'apprentissage: Modélisation de détecteur CdTe en python

## Schéma du montage :



Cd  
PR1

Source Radioactive  
– Cobalt 57  
– Sodium 22



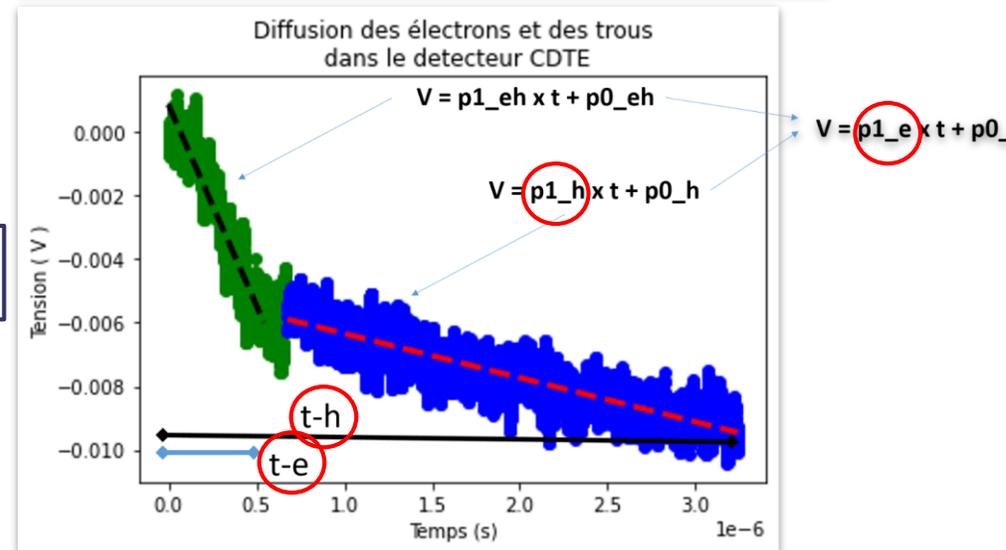
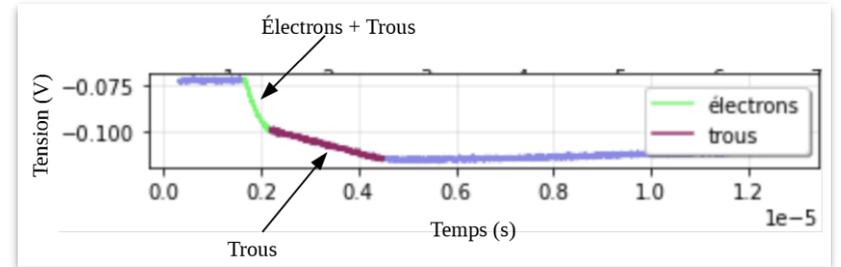
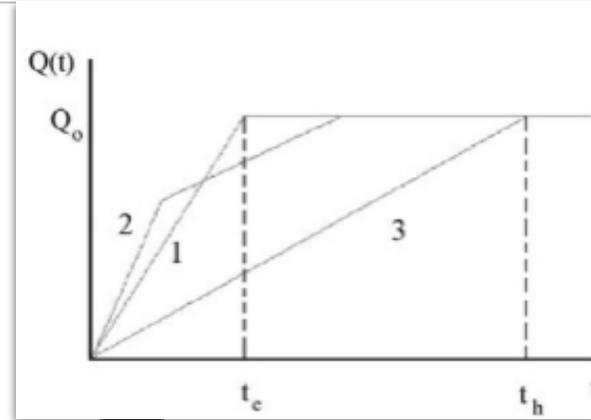
## Modèle théorique : Modèle de Hecht

$$Q(t) = \frac{E_p \cdot q}{E_i \cdot L} \cdot \left[ \lambda_e \left( 1 - e^{-\frac{-(L-x_0)}{\lambda_e}} \right) + \lambda_h \left( 1 - e^{-\frac{-x_0}{\lambda_h}} \right) \right]$$

Charge déposée par le photon détecté par unité de longueur

Partie électrons

Partie trous



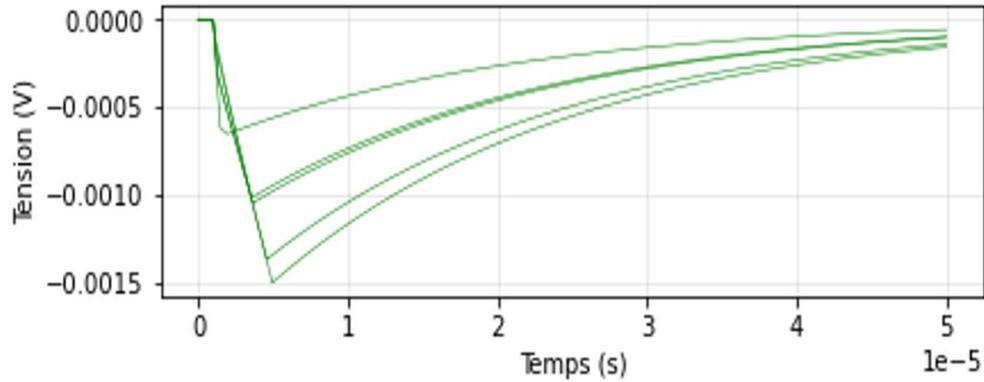
- E<sub>p</sub> : Énergie émise par la source radioactive (eV)
- E<sub>i</sub> : Énergie de ionisation du détecteur (4,43eV)
- λ<sub>e</sub> : Longueur de diffusion moyenne des électrons (m)
- λ<sub>h</sub> : Longueur de diffusion moyenne des trous (m)
- τ<sub>e</sub> : Temps de piégeage des électrons (s)
- τ<sub>h</sub> : Temps de piégeage des trous (s)
- μ<sub>e</sub> : Mobilité des électrons (m<sup>2</sup>/Vs)
- μ<sub>h</sub> : Mobilité des trous (m<sup>2</sup>/Vs)
- L : Longueur du détecteur (2mm)
- x<sub>0</sub> : Profondeur d'interaction de la particule (m)
- q : Charge de la particule (C)

$$\lambda_e = \mu_e \cdot \tau_e \cdot E_{\text{électrique}}$$

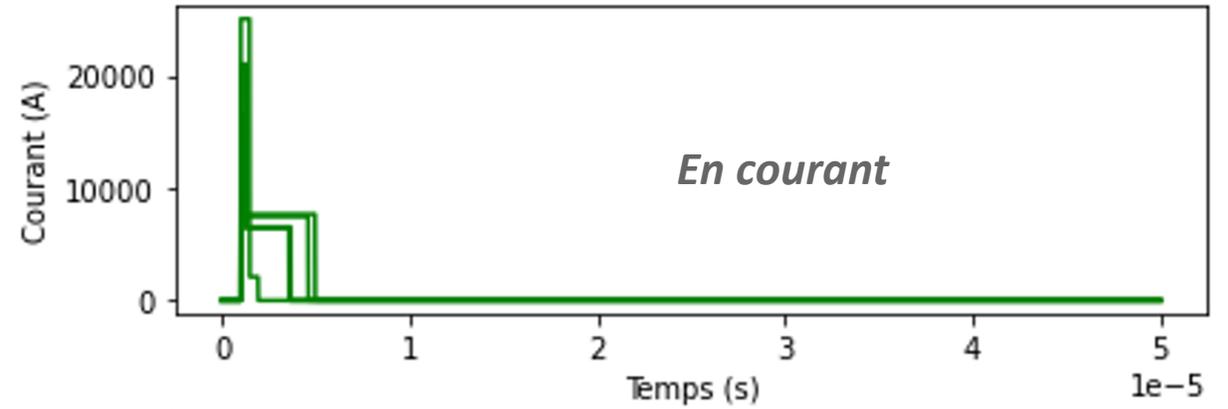
$$\lambda_h = \mu_h \cdot \tau_h \cdot E_{\text{électrique}}$$

- Résultats : Exemple modélisation de Hecht pour  $E = 122\text{keV}$  à profondeur de 1 mm :

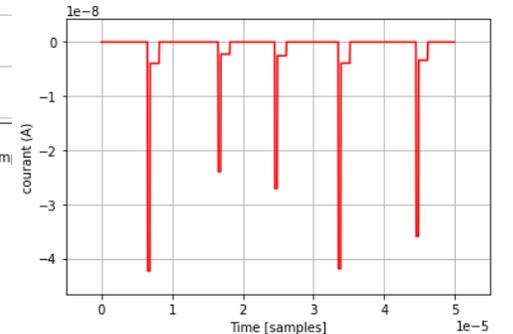
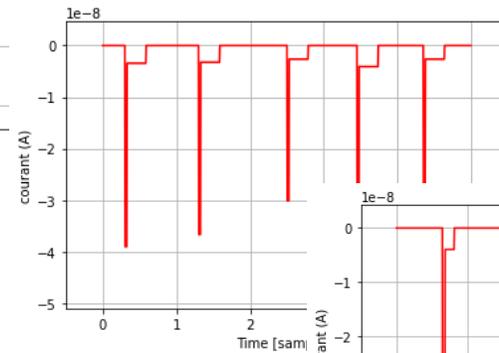
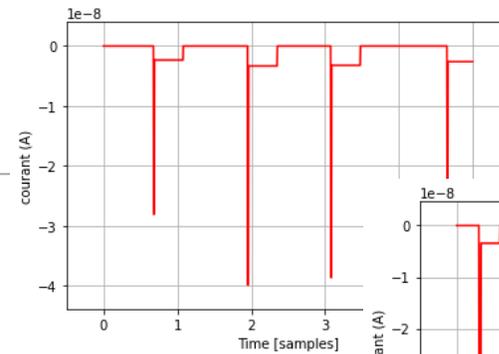
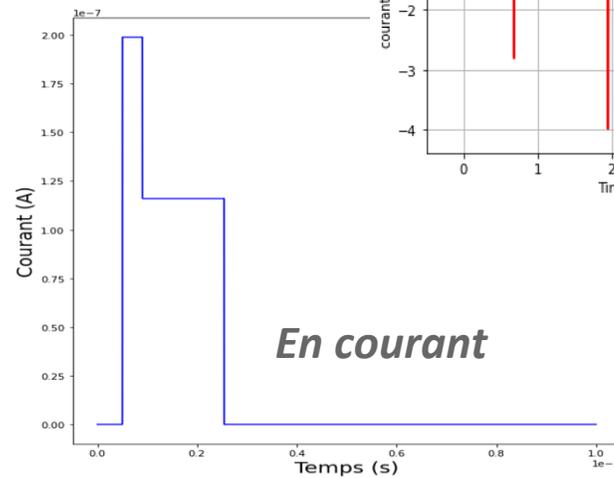
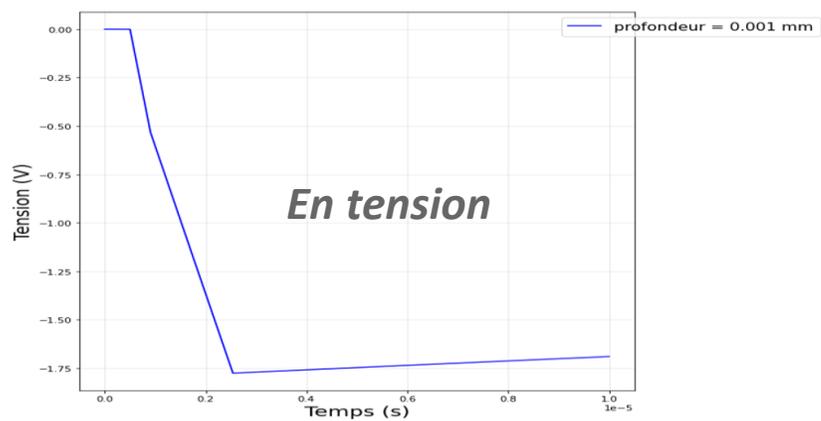
– Avec valeurs théoriques :



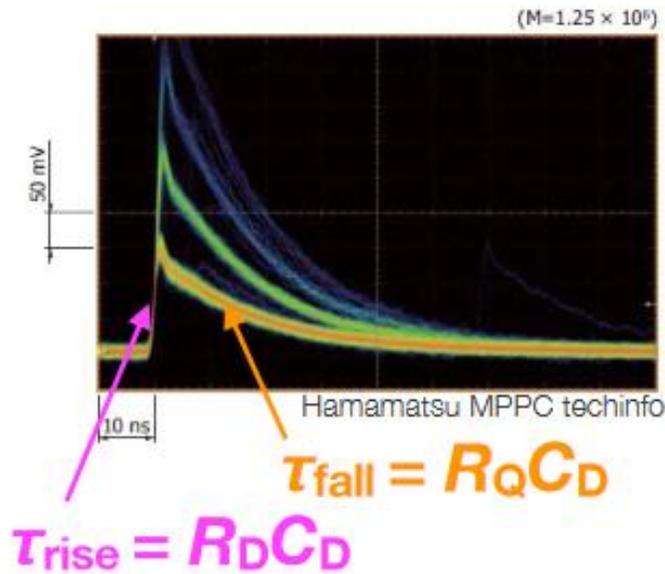
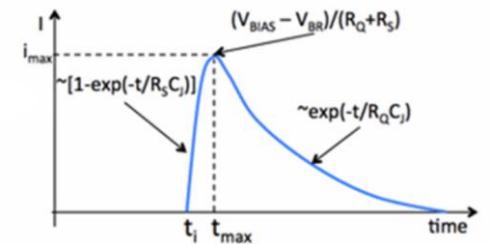
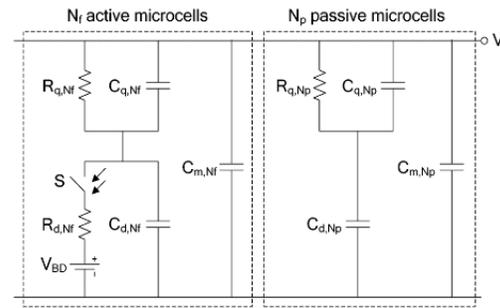
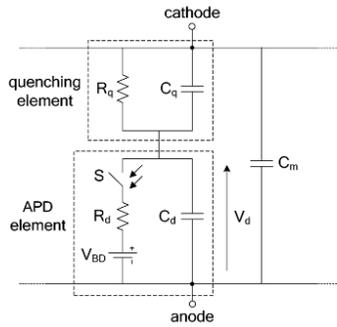
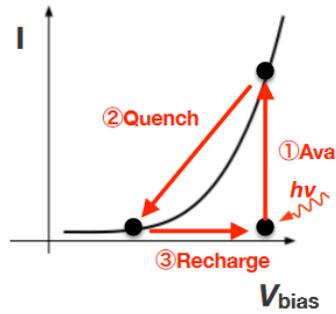
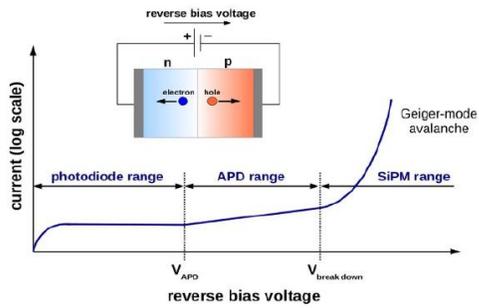
*En tension*



– Avec valeurs expérimentales :



# Création de données pour l'apprentissage: Modélisation de détecteur SiPM en python



```

sipm_prop.setRq(300)
sipm_prop.setRq(300e3)
sipm_prop.setCd(200e-15)
sipm_prop.setCq(20e-15)
sipm_prop.setCm(5e-15)
sipm_prop.setVov(2)
sipm_prop.setPropSimu('NbreCellules', 100)
sipm_prop.setPropSimu('NbreCelluleDeclenchee', 1)
sipm_prop.setPropSimu('NbreCellules', 1000)
sipm_prop.setPropSimu('NbreCelluleDeclenchee', 1)
sipm_prop.setPropSimu('courantLimite', 100e-6)
sipm_prop.setProperty('temps_de_descente_lent', 3e-9)
sipm_prop.setProperty('tempsRecuperation', 3e-10)
    
```

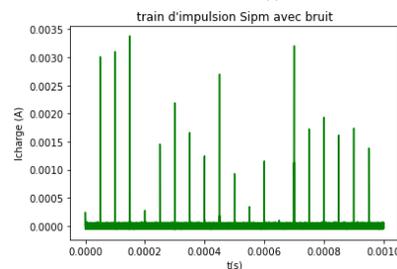
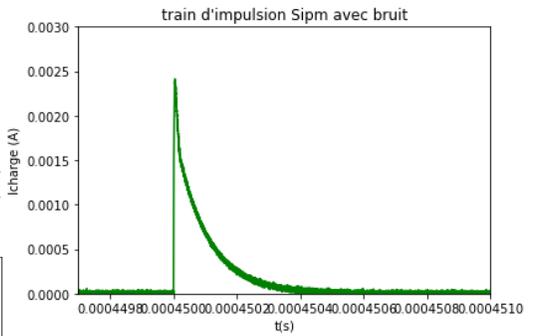
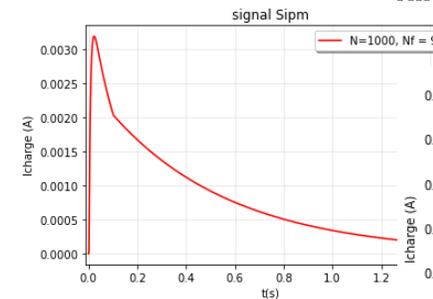
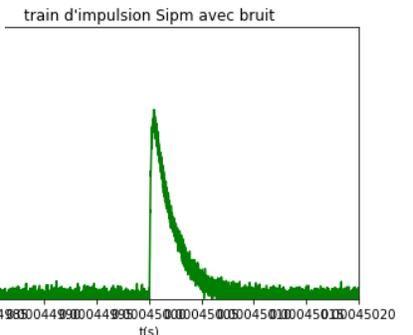
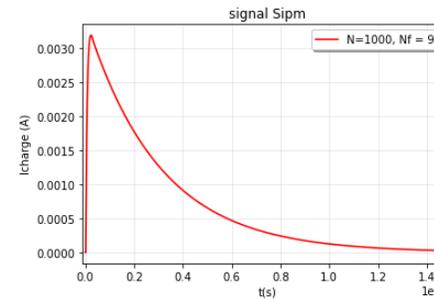
Calcul:

- $C_{eq} = f(C_d, C_q, C_m)$
- $H_o, -H_{RL}$
- $I_{RL}, -I_{RC}$
- $I_{peak}$

Modèle Signal Analogique

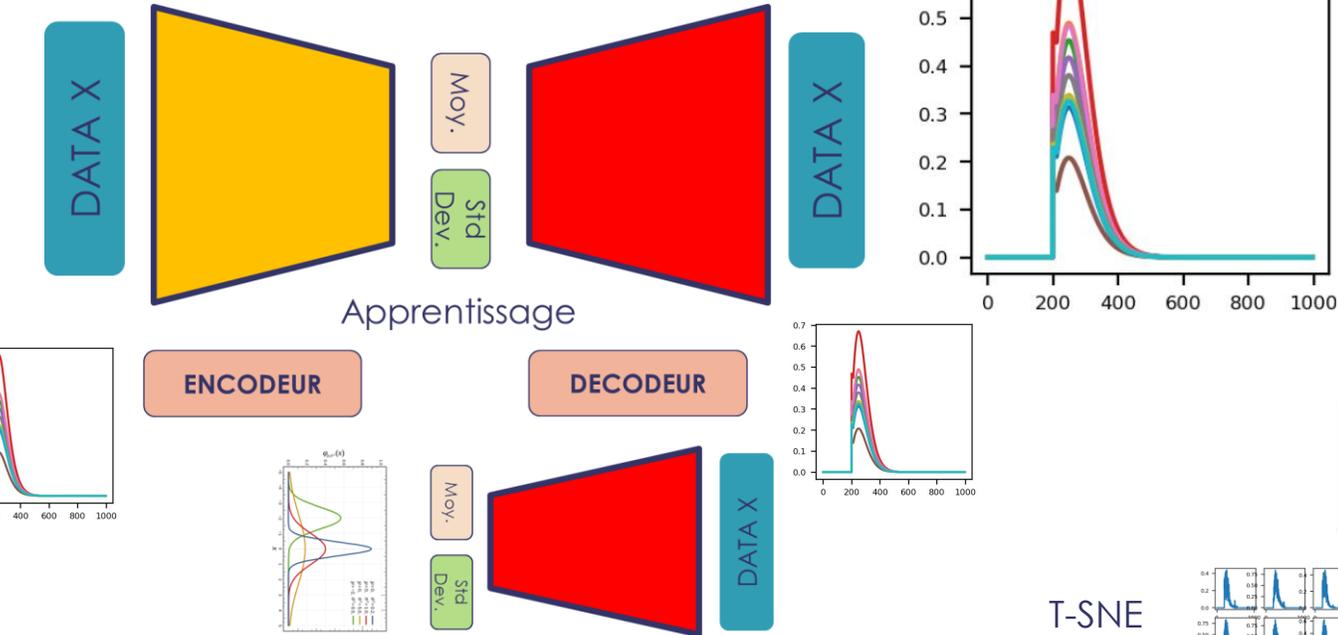
```

t1 = math.sqrt(2 * Rq * Cd)
tmax = math.sqrt(2 * Rq * Cd)
Ipeak = (Vbias - Vbr) / (Rq + Rs)
Idec = Ipeak * exp(-t / (Rq * Cd))
    
```

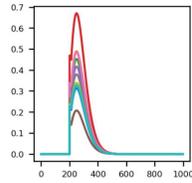
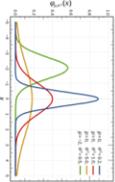
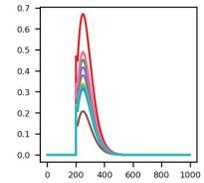
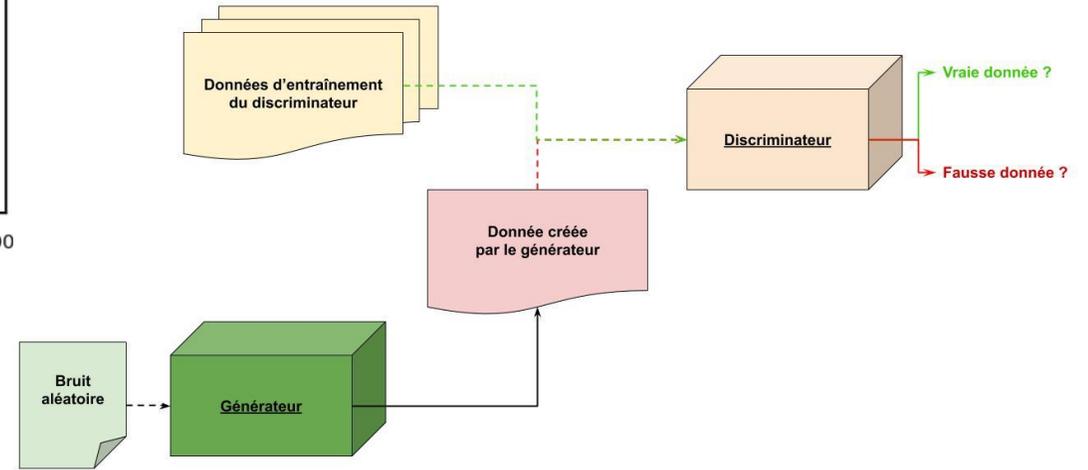


# Pipeline de données: Génération de données

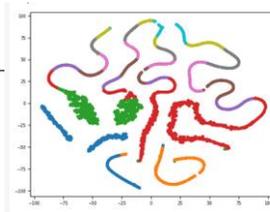
## Variational Auto-Encoder



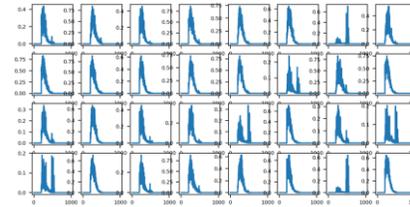
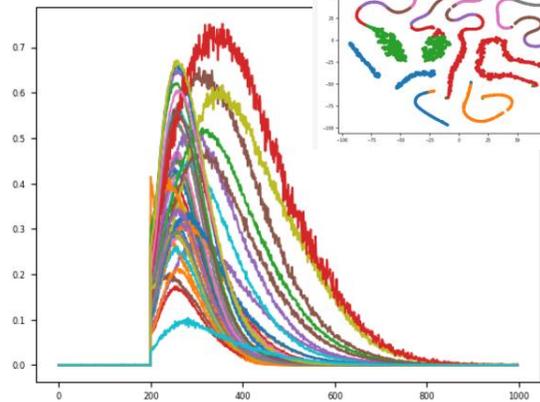
## Generative Adversal network



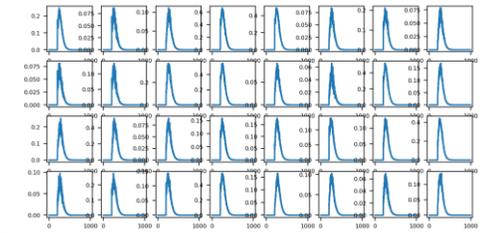
T-SNE



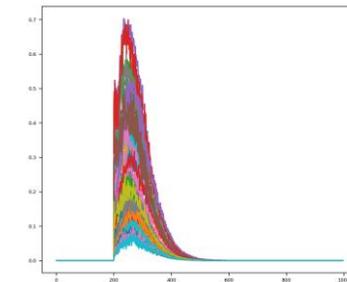
ANN



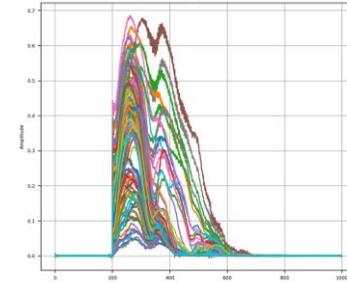
Learning phase



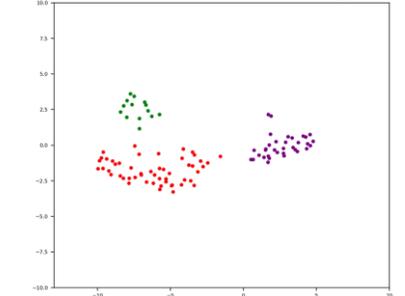
ANN



CNN



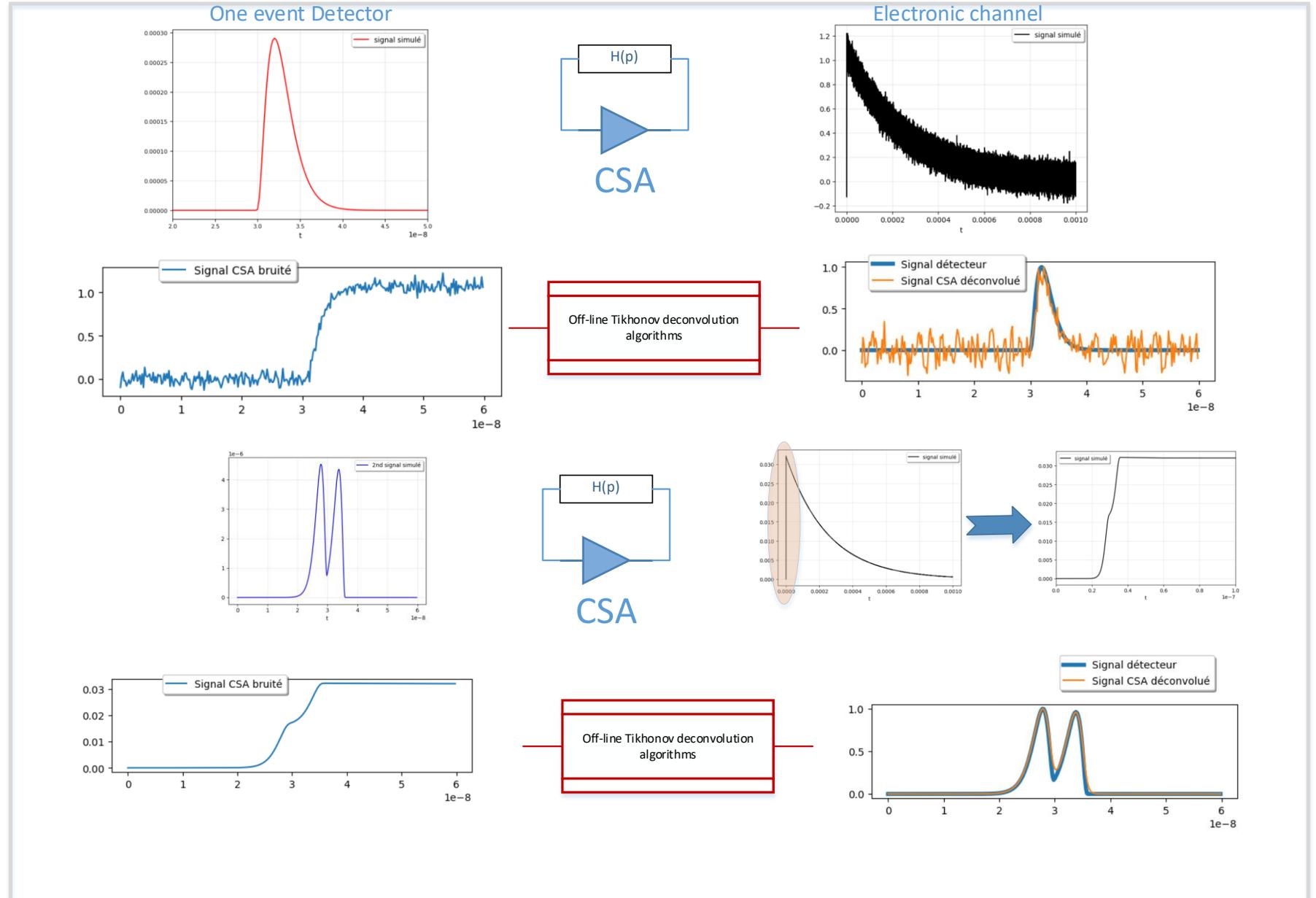
T-SNE



# Déconvolution par algo de Tikhonov (Off-Line)

Problème mal posé au sens d'Hadamard:

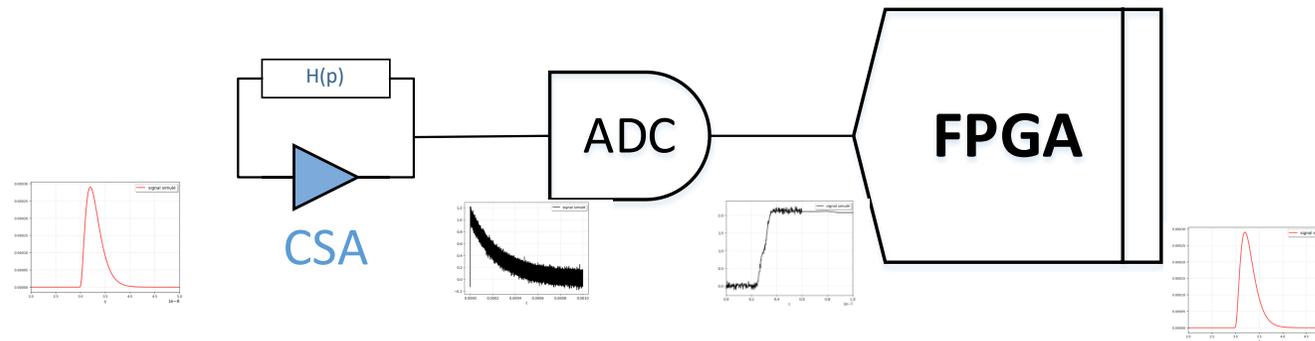
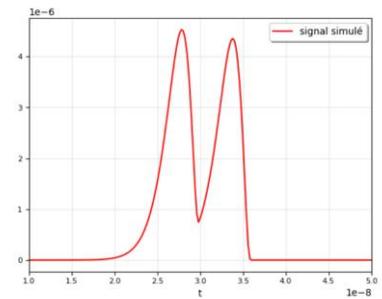
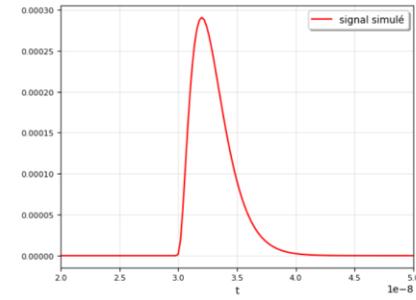
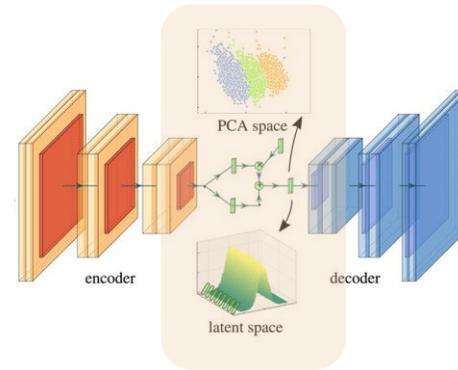
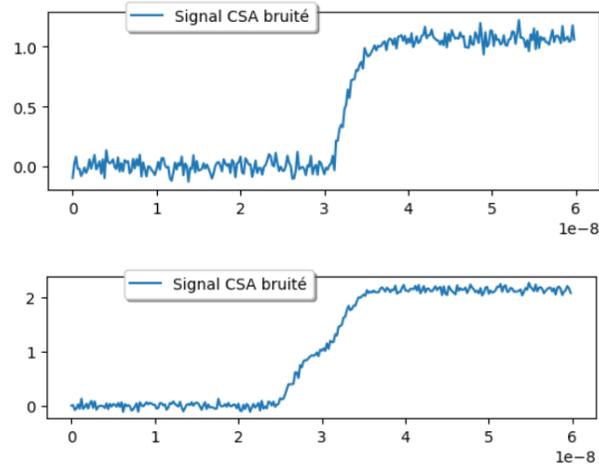
→ Algorithme Off-Line



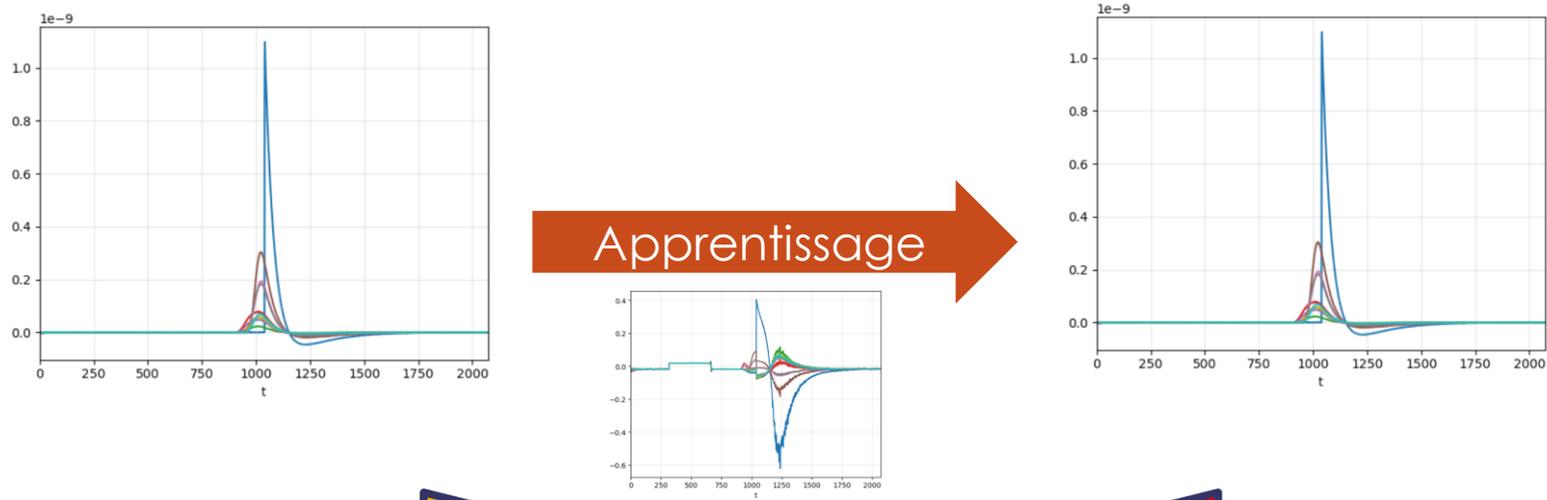
Problème mal posé  
au sens  
d'Hadamard:

→ Algorithme On-  
Line

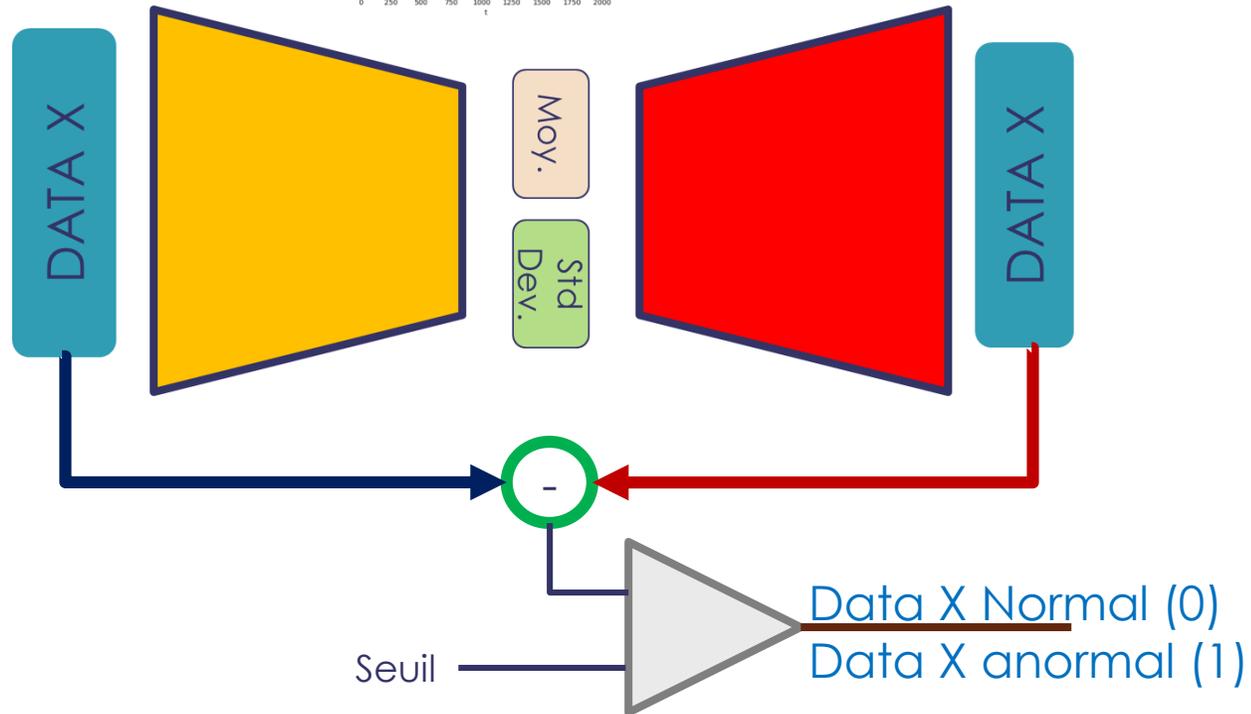
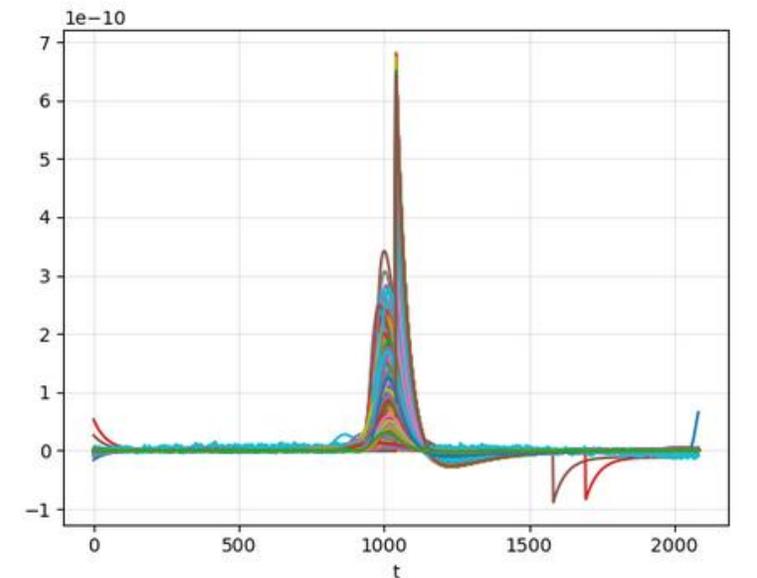
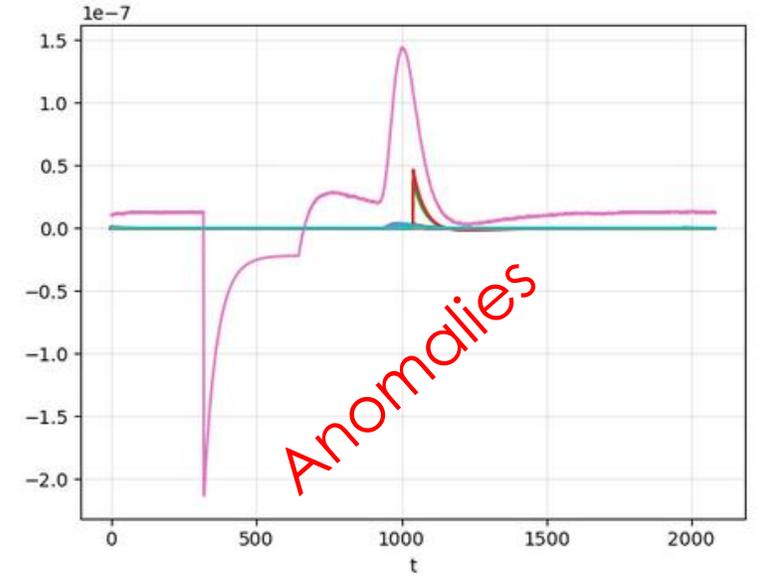
## Auto-encoders denoising and deconvolution



# Détection d'anomalie



## Inférence avec les vrai mesures



## Conclusion

- ▶ J'ai créé des modèles de détecteurs pour préparer les données d'apprentissage
  - ▶ APD
  - ▶ CdTe
  - ▶ SiPM
- ▶ J'ai testé des modèles d'IA pour:
  - ▶ Générer des données par IA
  - ▶ Classifier des signaux
  - ▶ Détecter des signaux anormaux
- ▶ J'ai étudié le principe de la déconvolution par IA

- Développement des applications SID (système intelligent pour détecteur )
- Mettre en œuvre les codes, les données pour une application réelle:
  - R2D2: signaux double bêta
  - Optimisation de spectre par déconvolution
  - Amélioration du rapport signal sur bruit sur SiPM et APD.

- ❑ **Utiliser de l'IA pour le traitement des signaux de détecteurs, c'est:**
  - ❑ **Connaitre la physique de création des signaux dans le détecteurs**
  - ❑ **Modéliser et simuler la création de signaux en sortie de détecteurs**
  - ❑ **Ajouter la chaîne électronique complète**
  - ❑ **Effectuer des vraies mesures pour confirmer les formes des signaux.**

- ❑ On peut appliquer les réseaux de neurones pour:
  - ❑ Générer de plus grandes quantités de données
  - ❑ Avoir plus de diversités dans les données générées
  - ❑ Créer des fonctions (approximateur universel) pour:
    - ❑ Trouver les données anormales
    - ❑ Catégoriser les événements
    - ❑ Élaborer de nouvelles fonctions de transfert à partir des données

- ❑ Inférer des réseaux de neurones au plus près des détecteurs nécessite:
  - ❑ Une méthode d'optimisation des modèles (elagage, quantisation)
  - ❑ Un savoir faire en IA (experiences/ exploration)
  - ❑ Une connaissance des composants et de leur SDK