

# AVANCÉES EN THÉORIE PAC-BAYÉSIENNE

## DE BORNES EN GÉNÉRALISATION À DES ALGORITHMES D'APPRENTISSAGE SUPERVISÉ ET DE TRANSFERT

Soutenance HDR — 7 avril 2025

**Emilie MORVANT**

Université Jean Monnet - Laboratoire Hubert Curien

|             |                            |              |                          |
|-------------|----------------------------|--------------|--------------------------|
| Rapporteurs | <b>Marianne CLAUSEL</b>    | Professeure  | Université de Lorraine   |
|             | <b>Colin DE LA HIGUERA</b> | Professeur   | Université de Nantes     |
|             | <b>Liva RALAIVOLA</b>      | VP Recherche | Criteo AI Lab            |
| Examineurs  | <b>Stéphane CHRÉTIEN</b>   | Professeur   | Aix-Marseille Université |
|             | <b>Rémi EMONET</b>         | MCF HDR      | Université Lyon 2        |
|             | <b>François JACQUENET</b>  | Professeur   | Université Jean Monnet   |
| Tuteur      | <b>François JACQUENET</b>  | Professeur   | Université Jean Monnet   |

# Parcours



# Parcours



Responsable Licence 2 d'informatique

Responsable  
L1 d'info

Monitorat

Décharge jeune chercheur

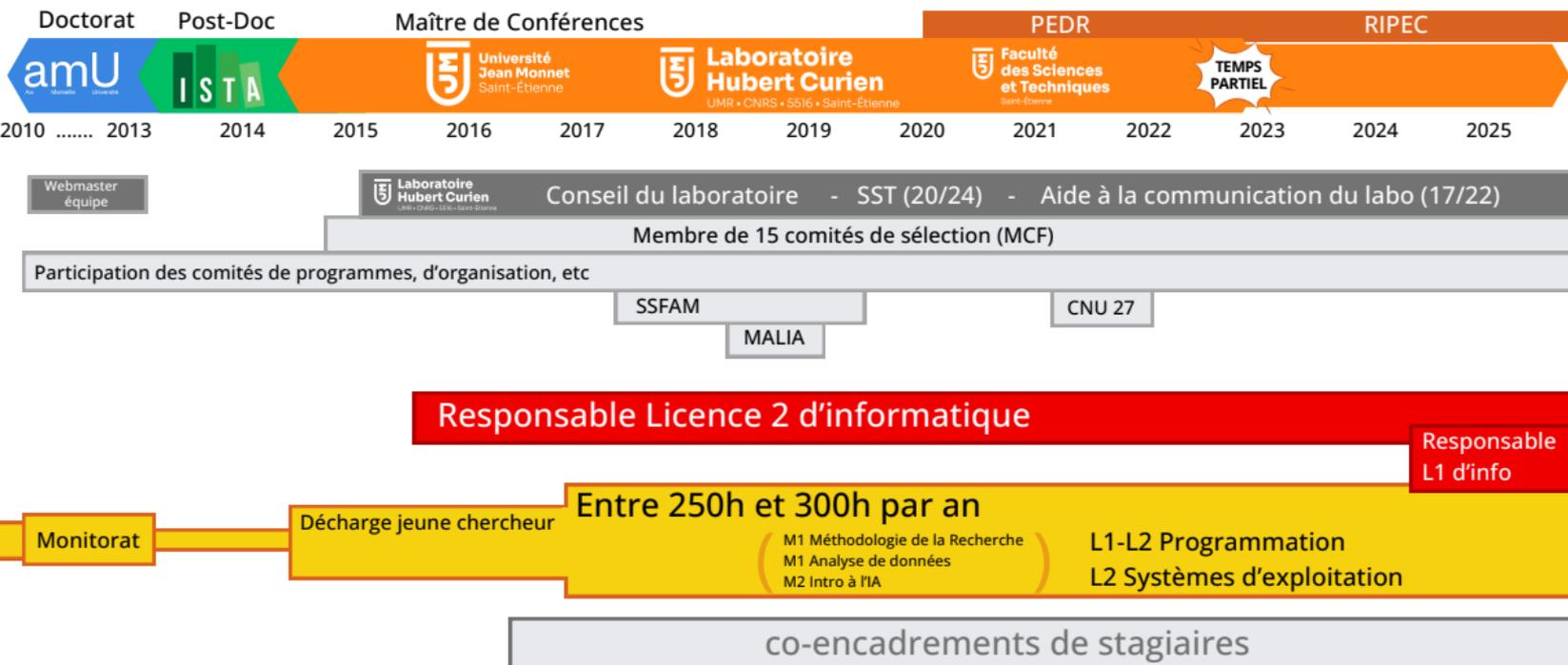
Entre 250h et 300h par an

( M1 Méthodologie de la Recherche  
M1 Analyse de données  
M2 Intro à l'IA )

L1-L2 Programmation  
L2 Systèmes d'exploitation

co-encadrements de stagiaires

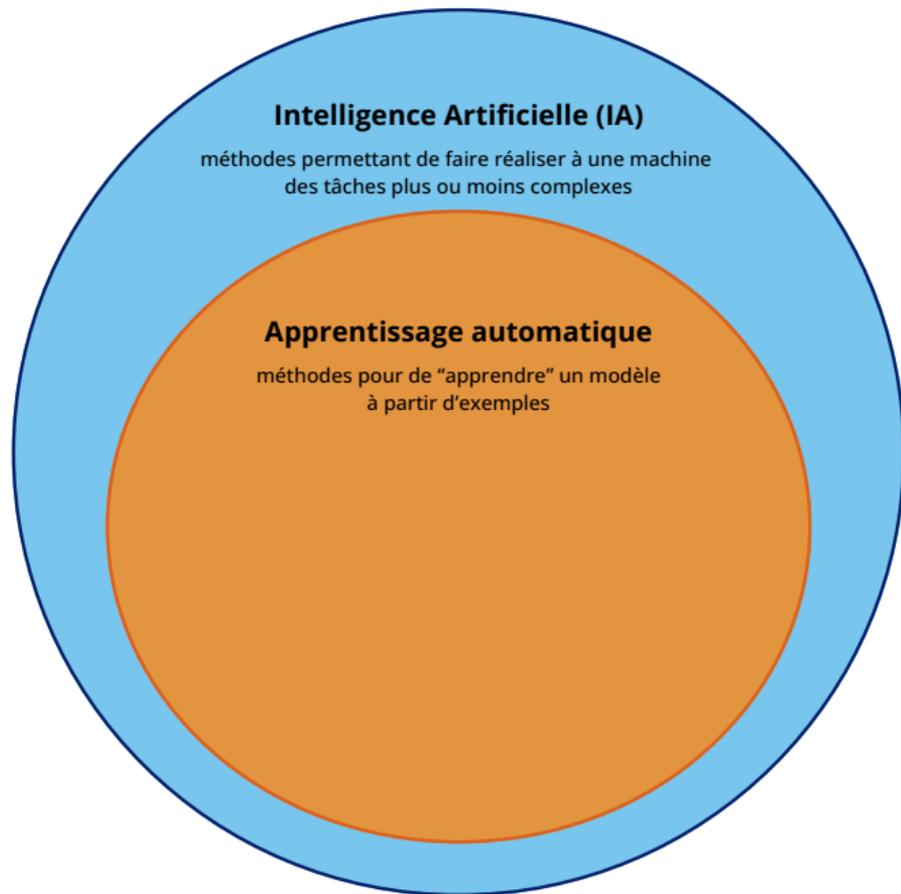
# Parcours



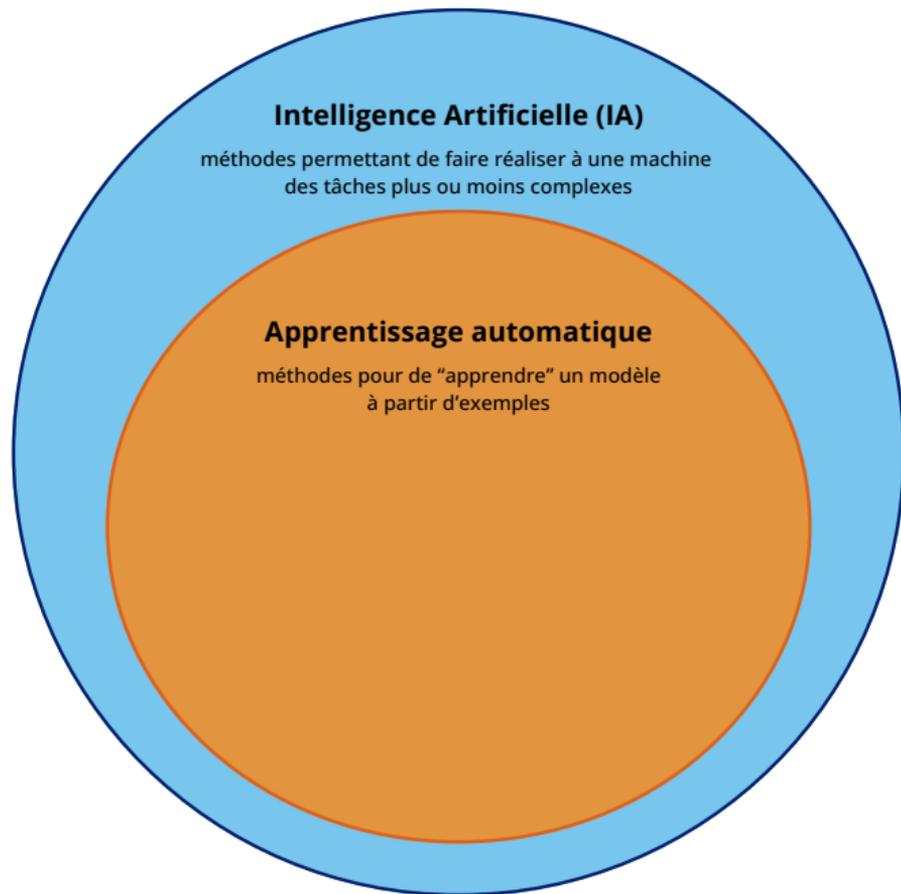
## Intelligence Artificielle (IA)

méthodes permettant de faire réaliser à une machine  
des tâches plus ou moins complexes

# Ma trajectoire de recherche

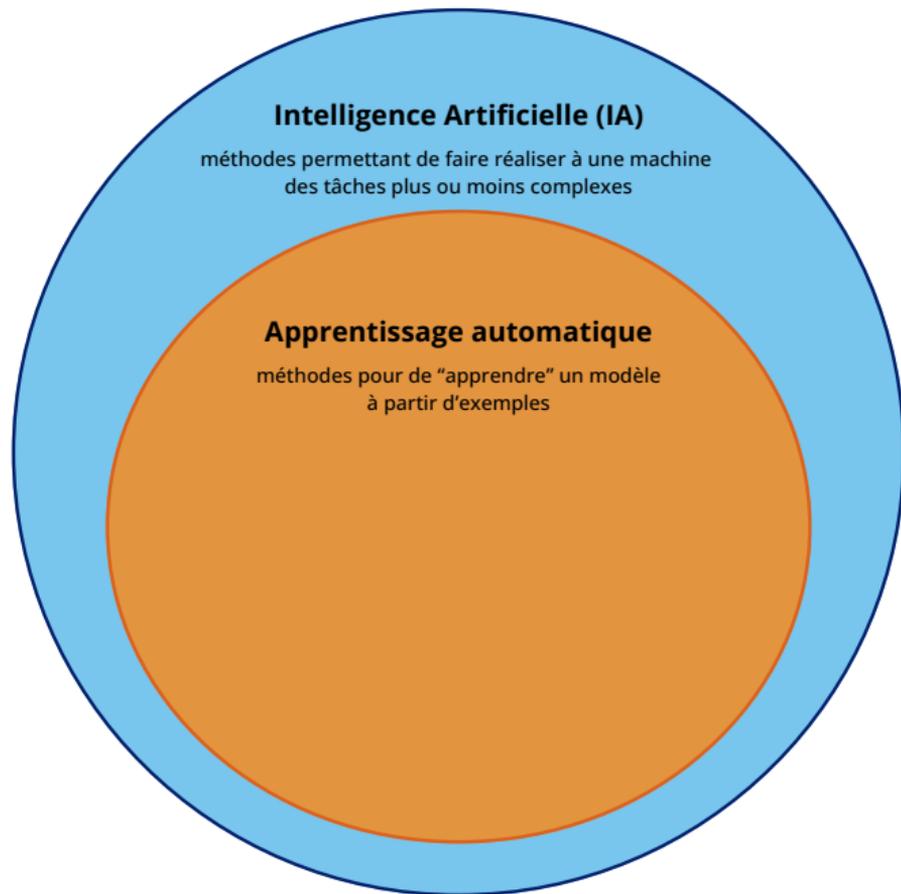


# Ma trajectoire de recherche



## ENJEUX SOCIÉTAUX EN IA

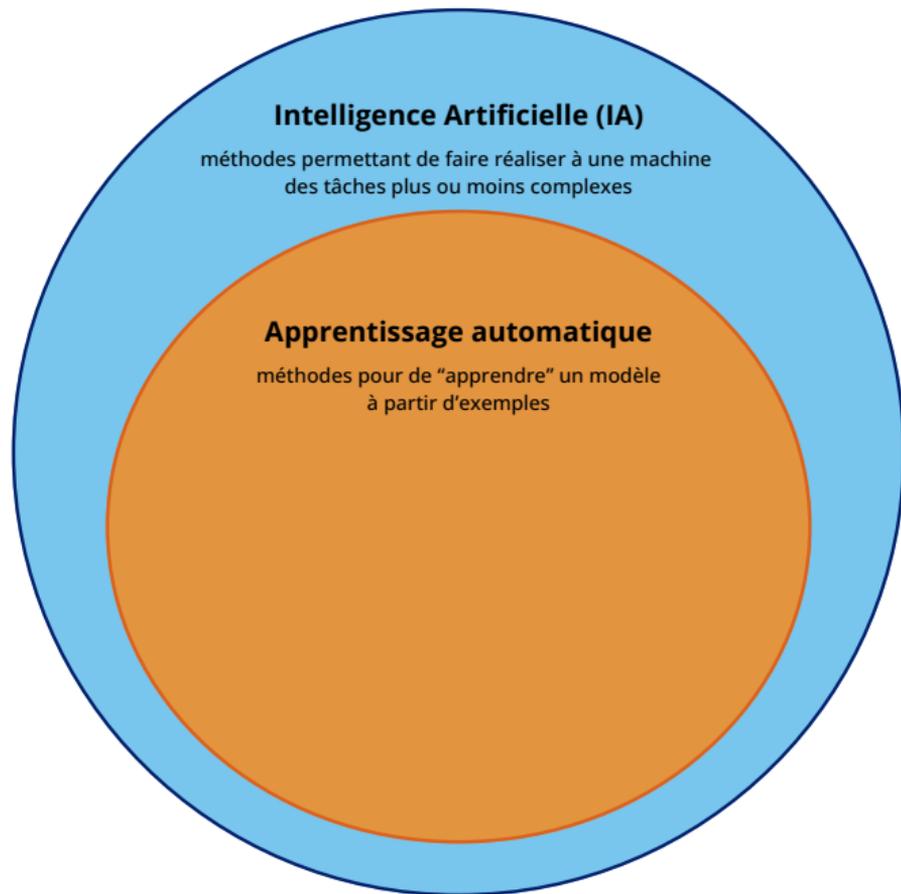
# Ma trajectoire de recherche



## ENJEUX SOCIÉTAUX EN IA

- **Confiance**

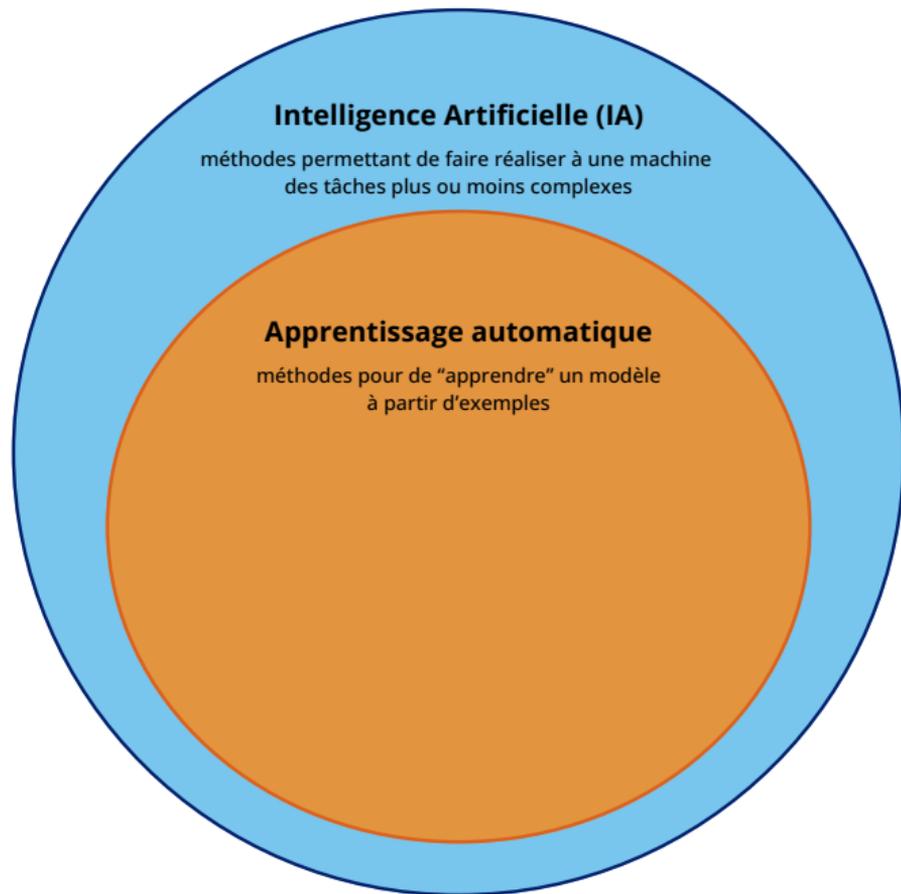
# Ma trajectoire de recherche



## ENJEUX SOCIÉTAUX EN IA

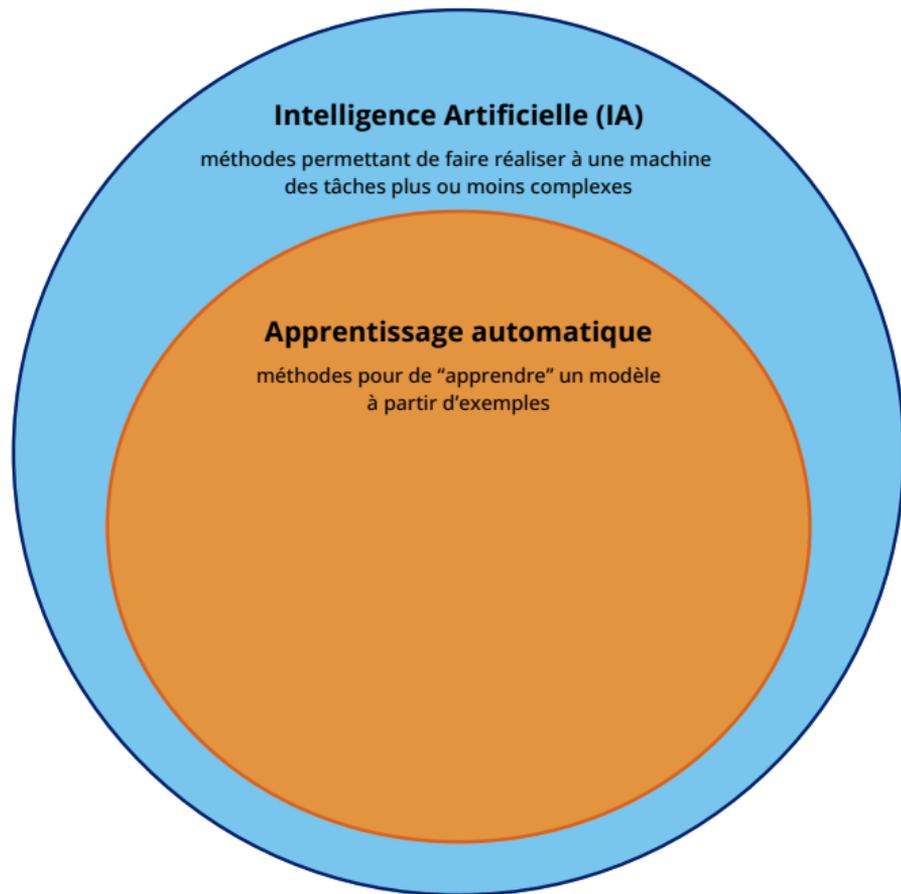
- **Confiance**
- **Explicabilité**

# Ma trajectoire de recherche



## ENJEUX SOCIÉTAUX EN IA

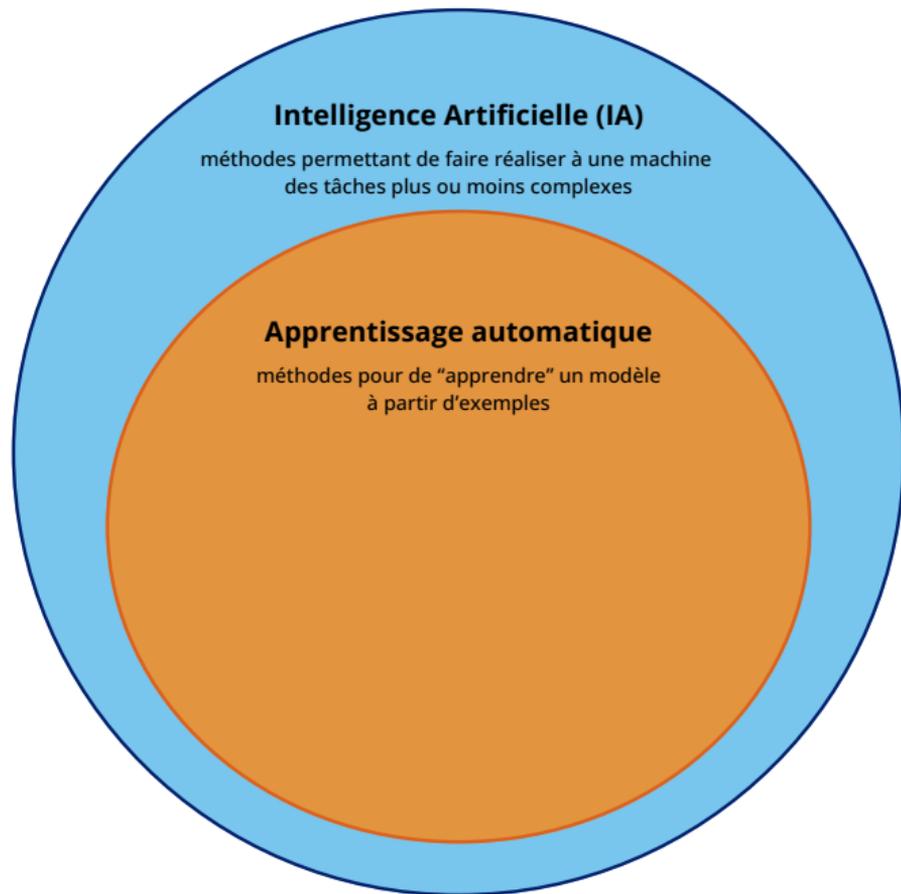
- **Confiance**
- **Explicabilité**
- **Adaptabilité**



## ENJEUX SOCIÉTAUX EN IA

- **Confiance**
- **Explicabilité**
- **Adaptabilité**
- **Robustesse**

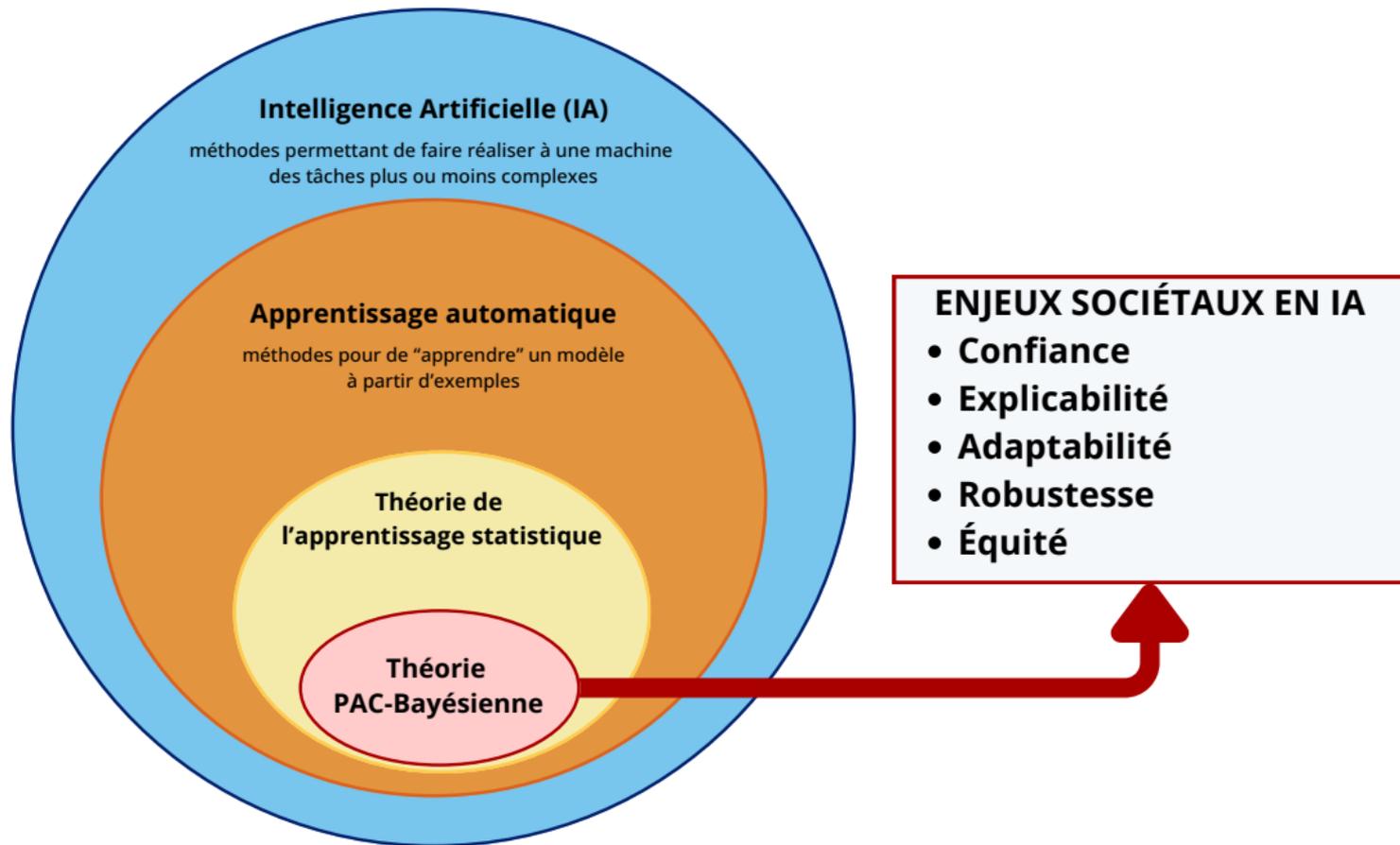
# Ma trajectoire de recherche



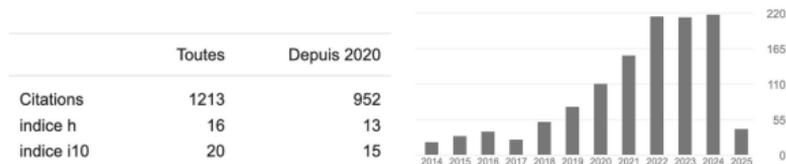
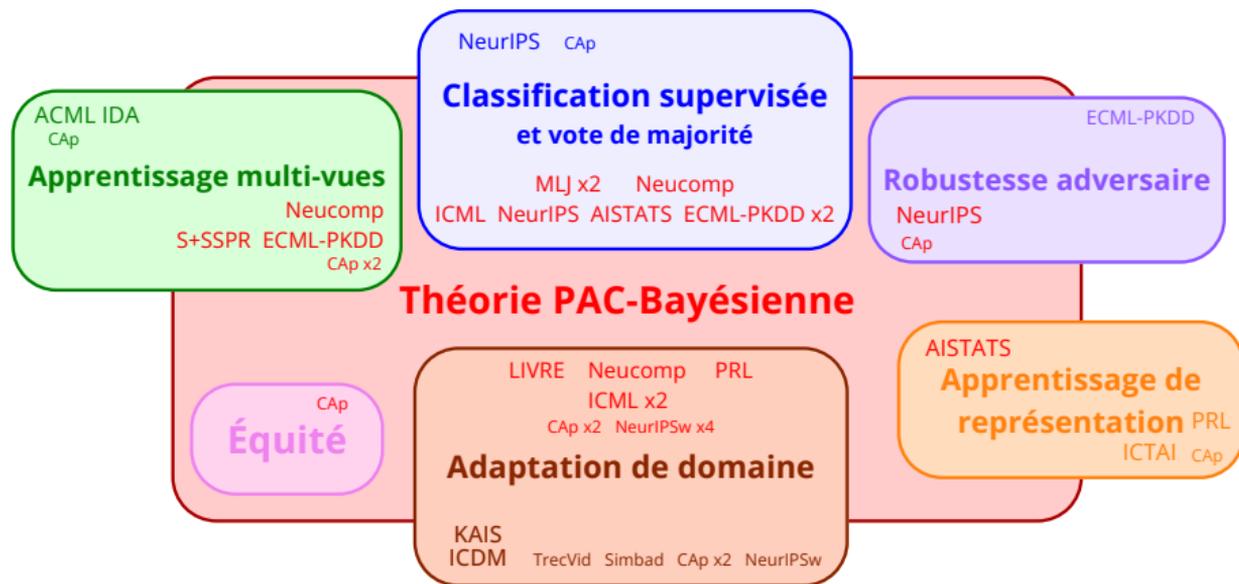
## ENJEUX SOCIÉTAUX EN IA

- **Confiance**
- **Explicabilité**
- **Adaptabilité**
- **Robustesse**
- **Équité**

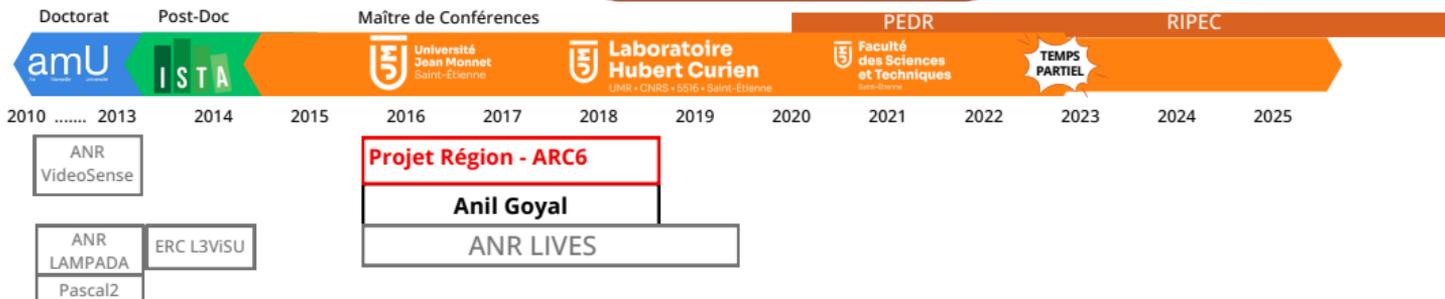
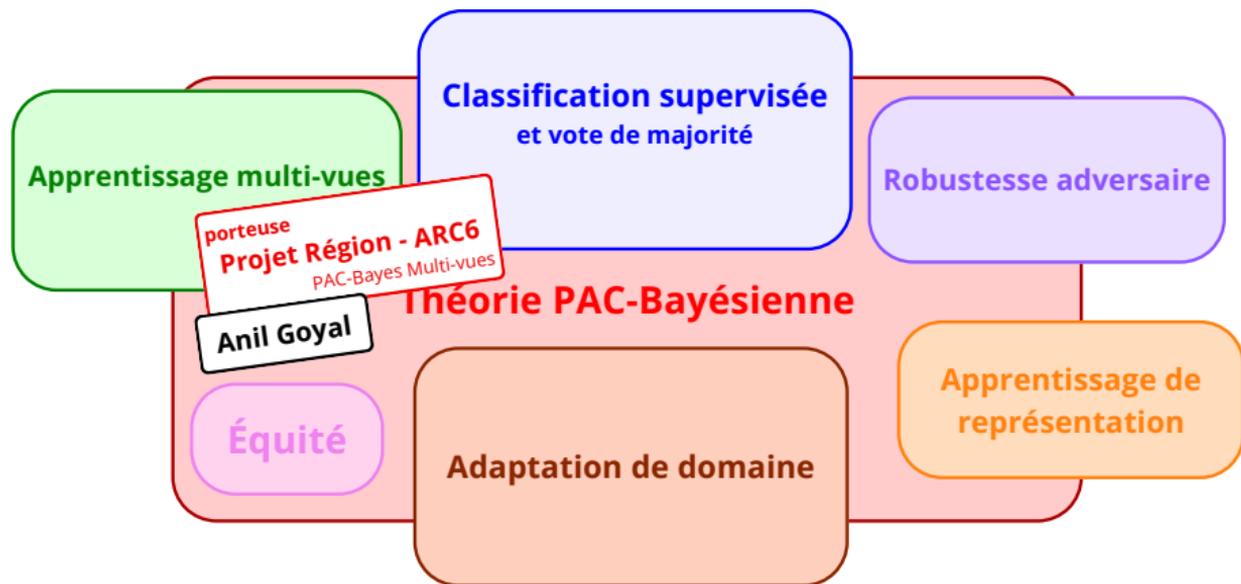
# Ma trajectoire de recherche



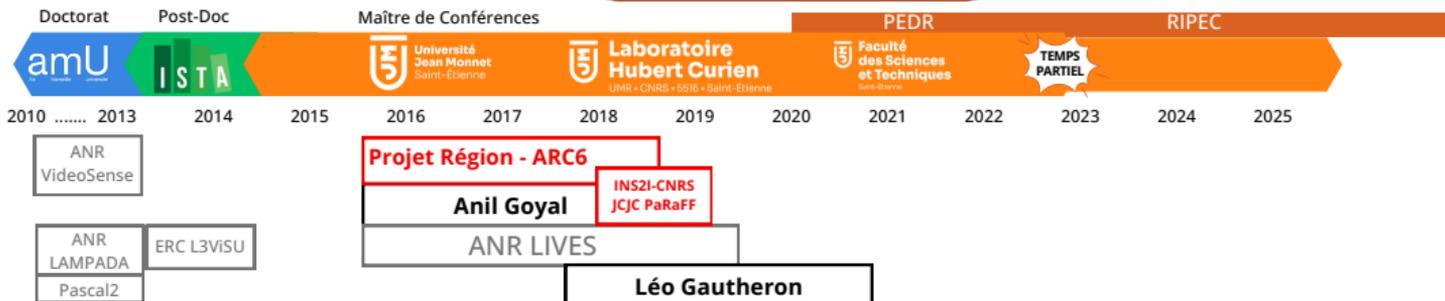
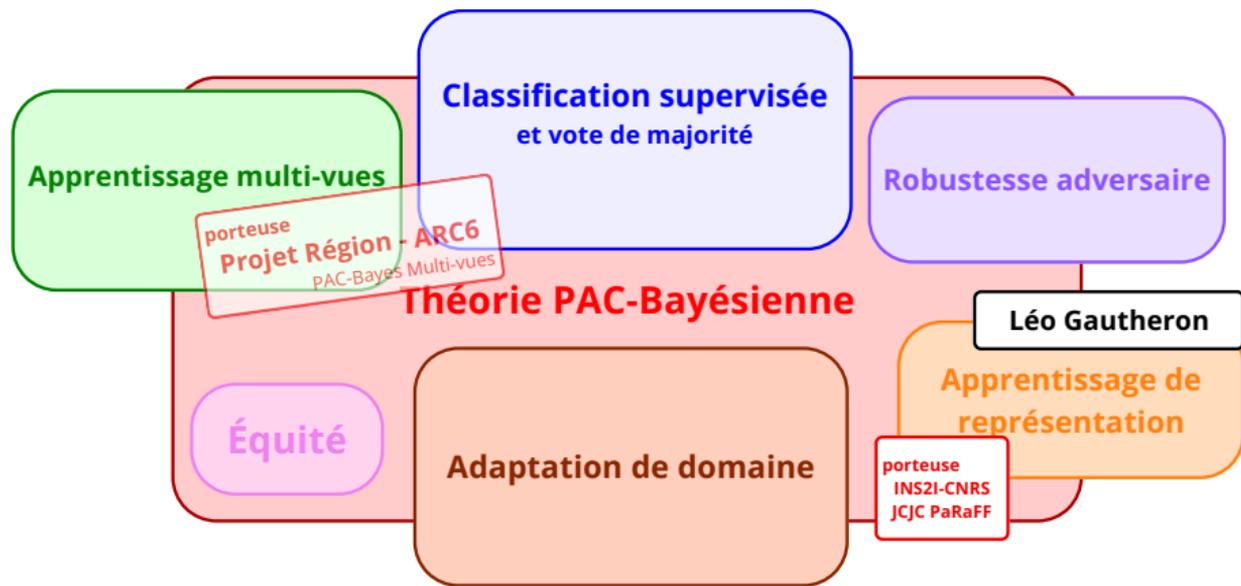
# Thèmes principaux et publications



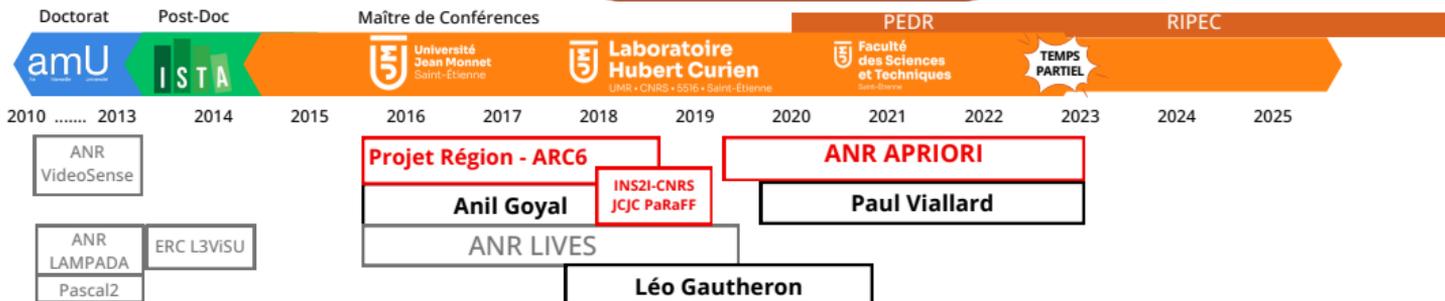
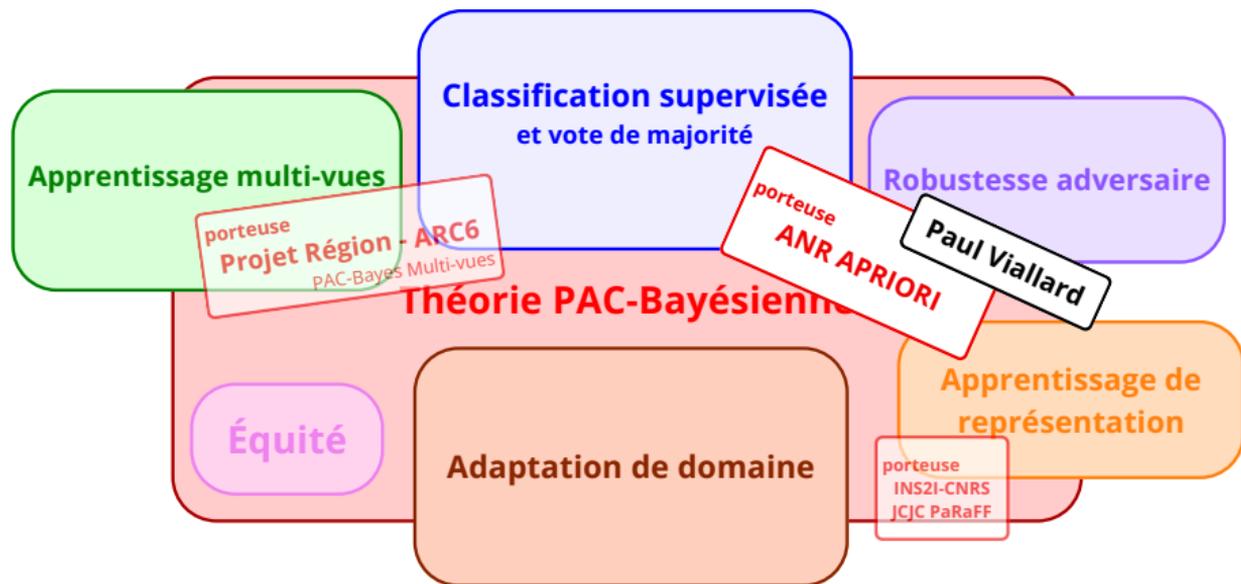
# Thèmes principaux et projets



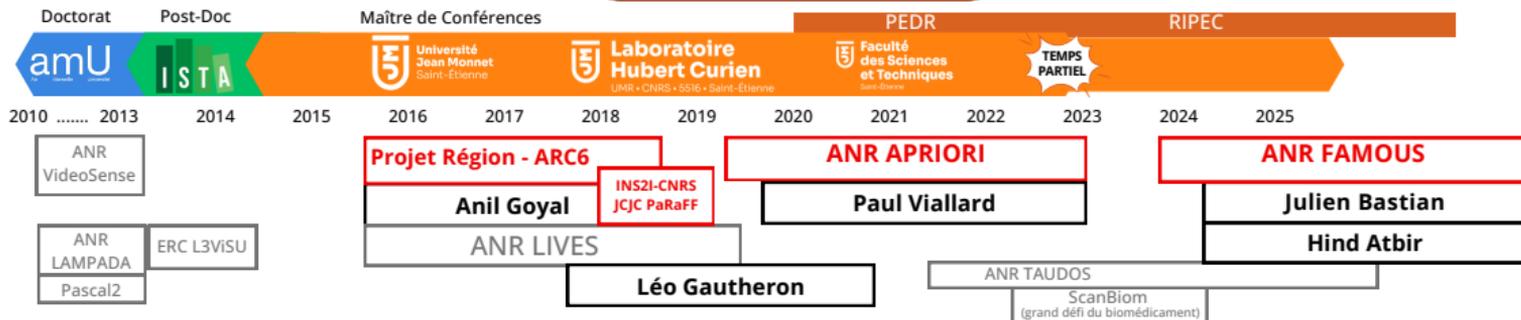
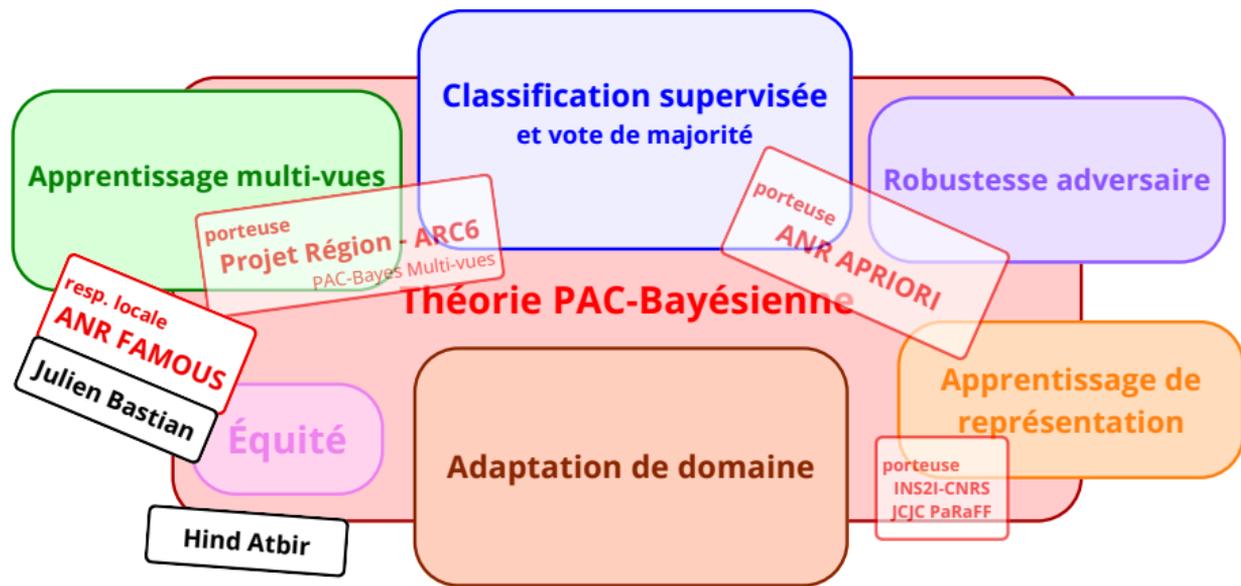
# Thèmes principaux et projets



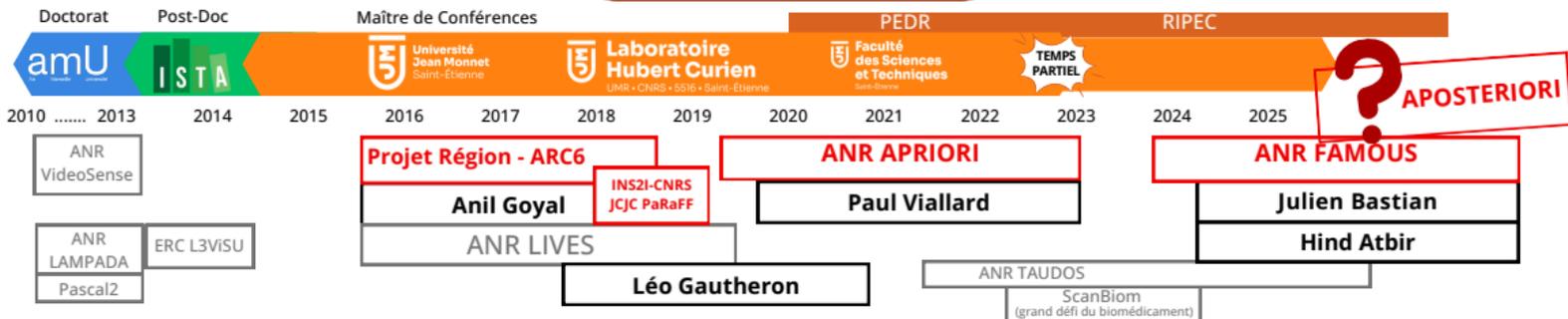
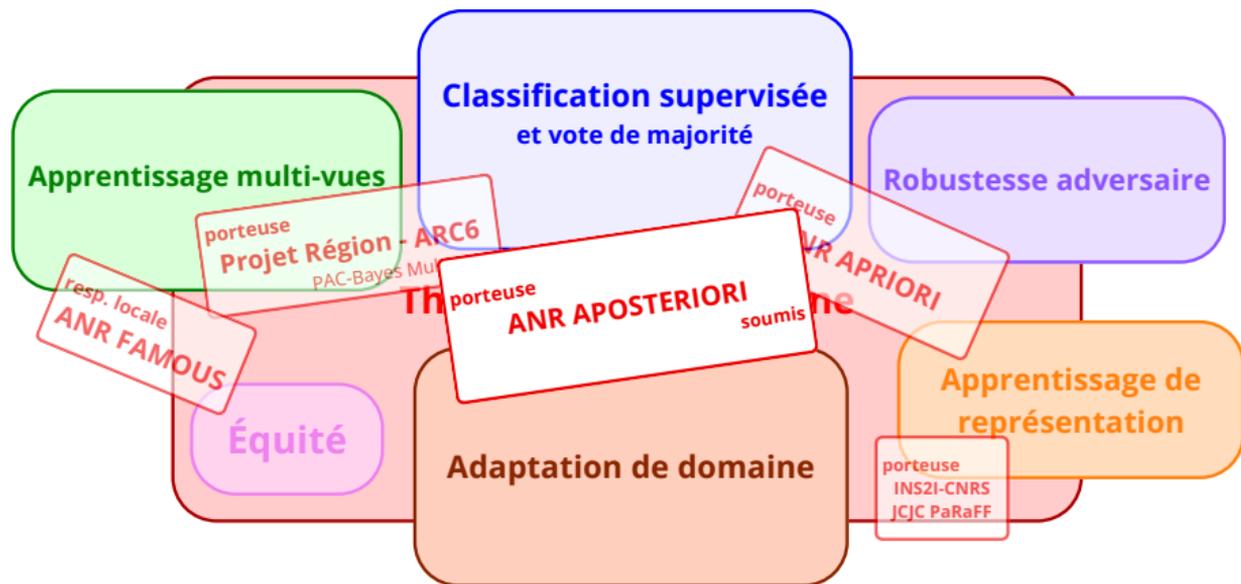
# Thèmes principaux et projets



# Thèmes principaux et projets



# Thèmes principaux et projets

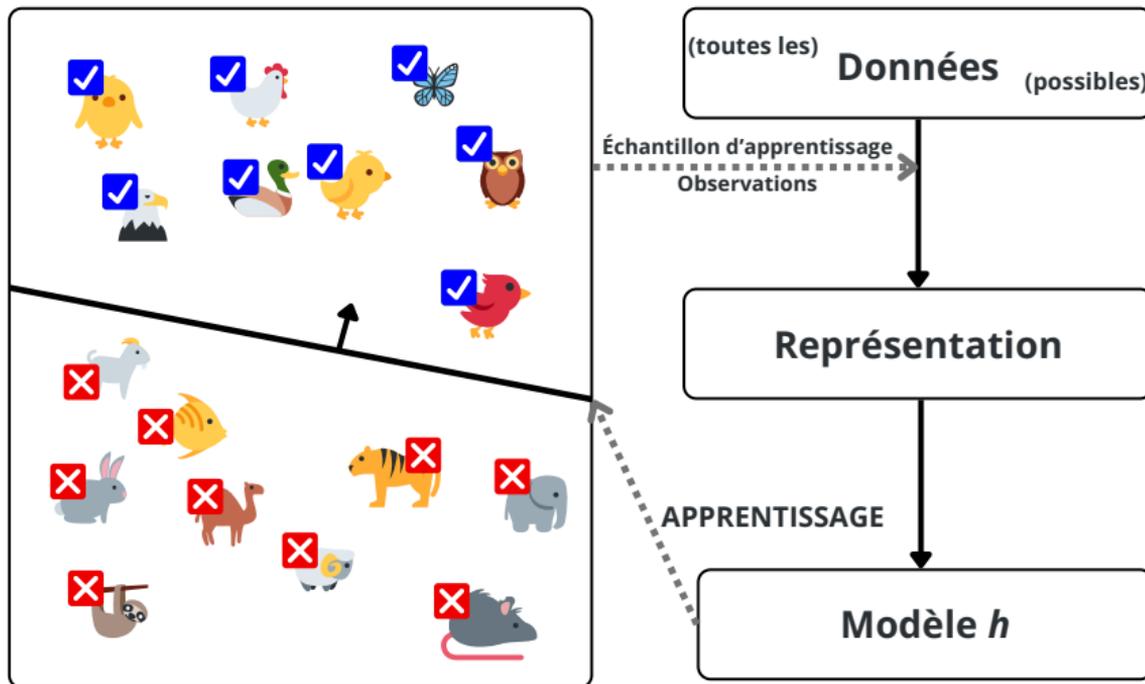


# CONTEXTE SCIENTIFIQUE GÉNÉRAL

# Classification supervisée

Tâche de **classification supervisée**

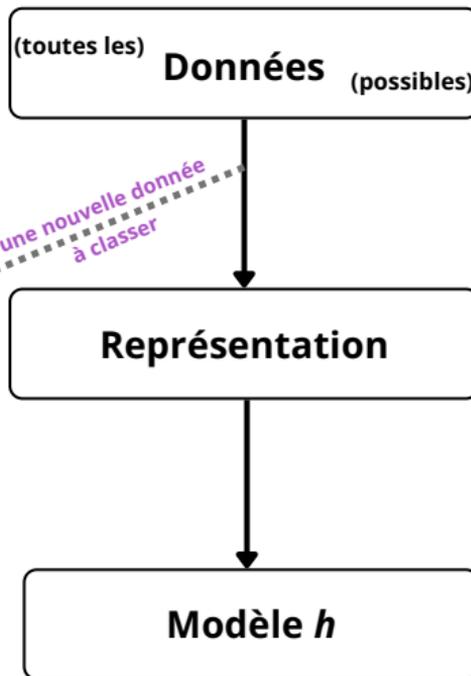
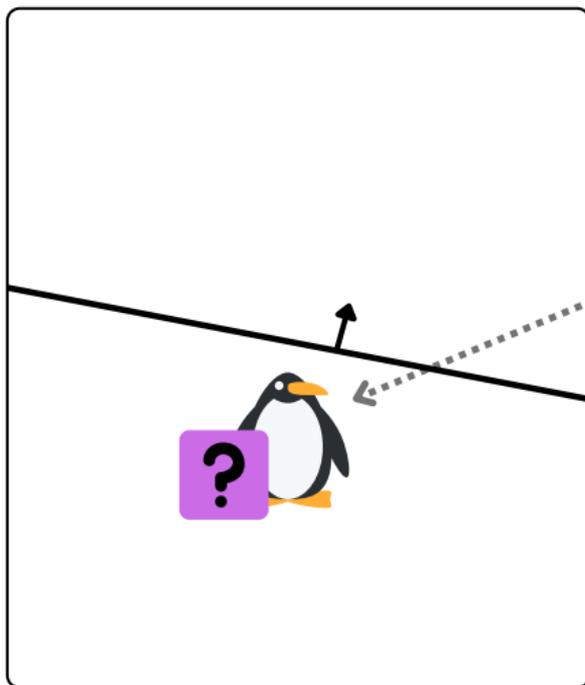
Classer les animaux en 2 catégories :  avec des ailes ou  sans ailes



# Classification supervisée

Tâche de **classification supervisée**

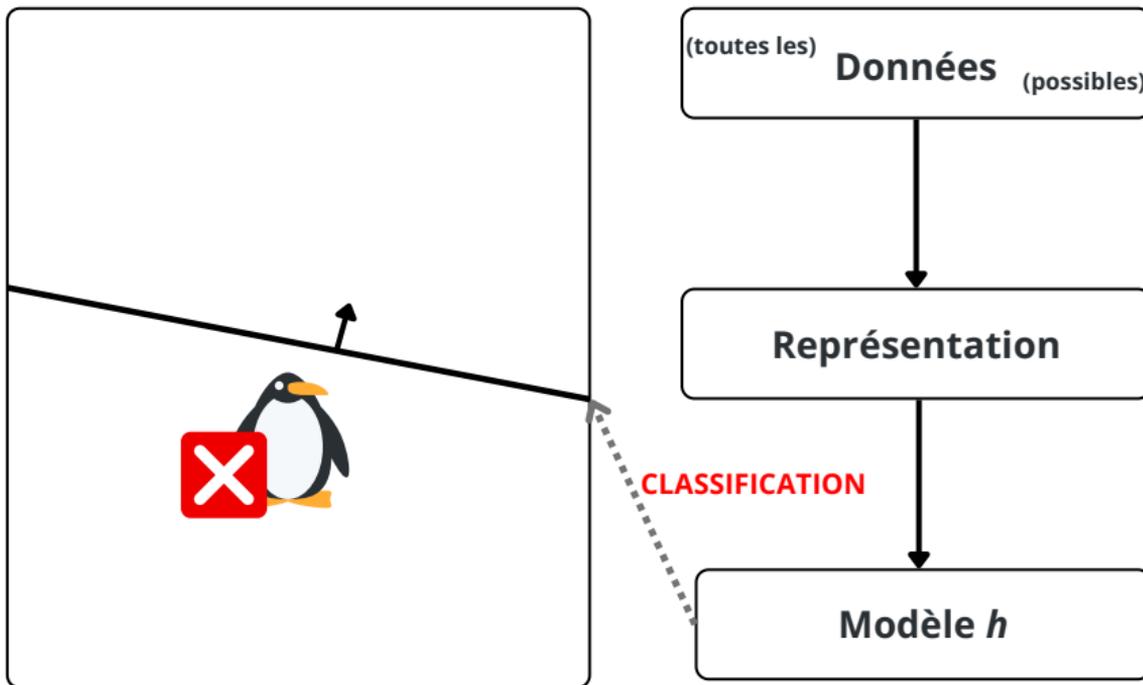
Classer les animaux en 2 catégories :  avec des ailes ou  sans ailes



# Classification supervisée

Tâche de **classification supervisée**

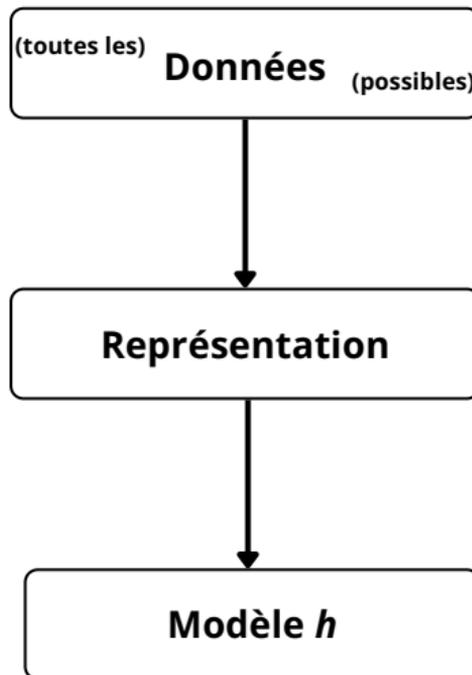
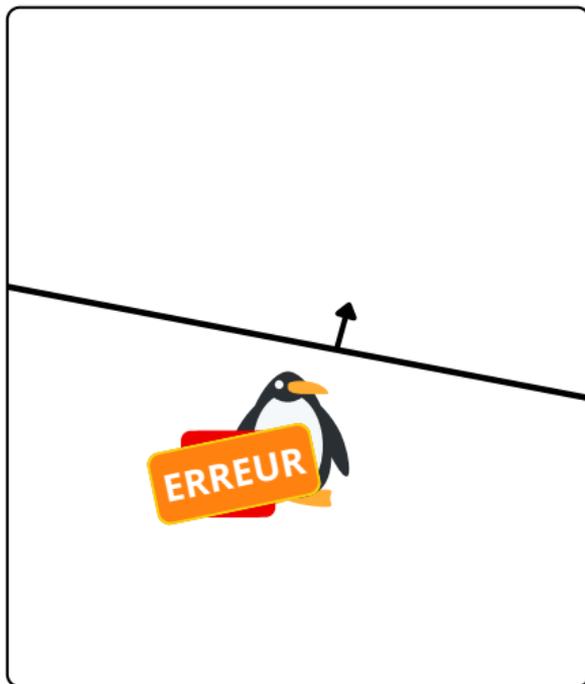
Classer les animaux en 2 catégories :  avec des ailes ou  sans ailes



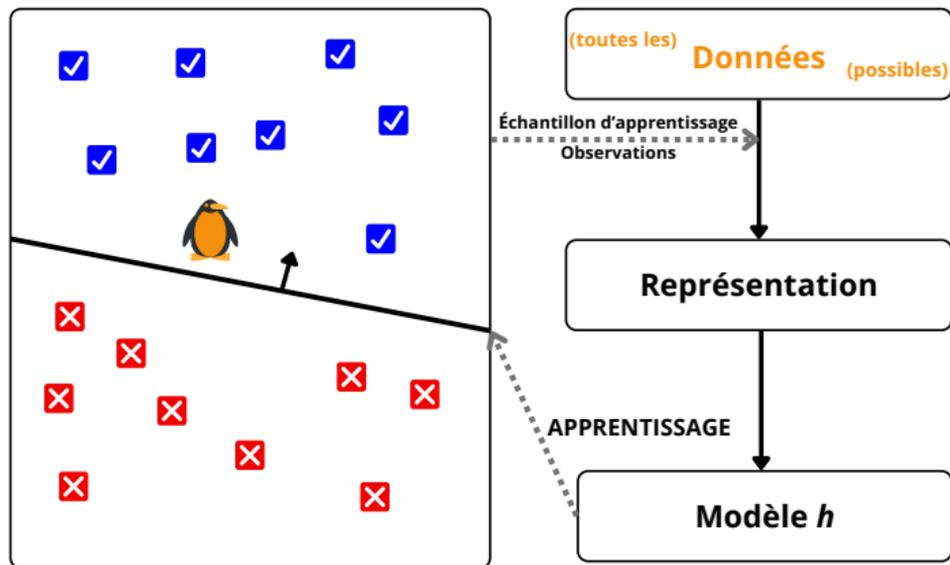
# Classification supervisée

Tâche de **classification supervisée**

Classer les animaux en 2 catégories :  avec des ailes ou  sans ailes

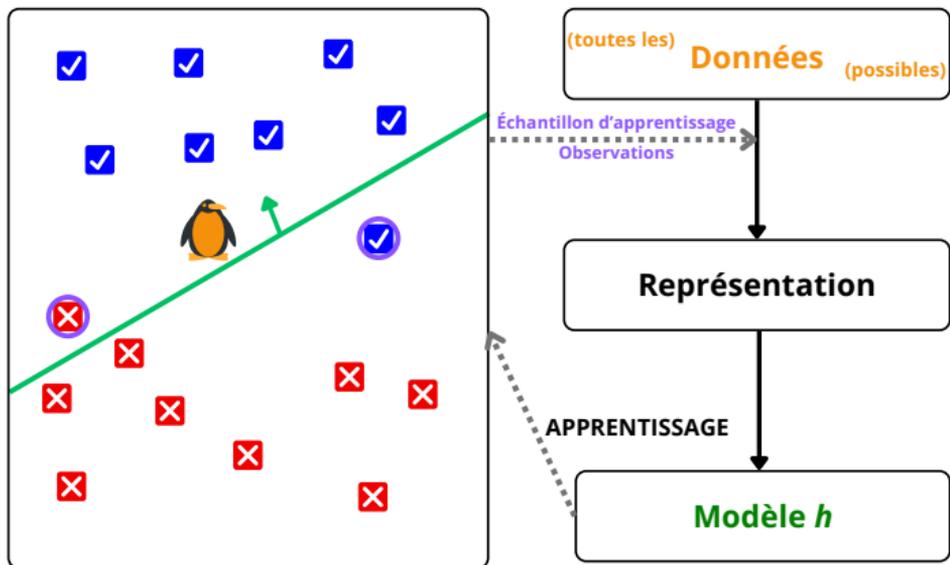


# Bornes en généralisation



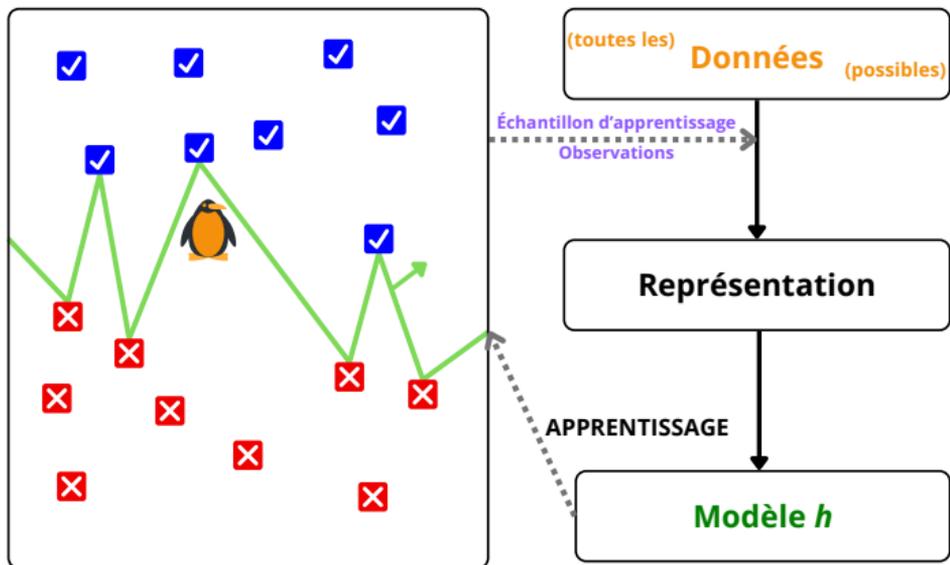
erreur sur  
toutes les données  $\leq$   
NON CALCULABLE

# Bornes en généralisation



$$\underbrace{\text{erreur sur toutes les données}}_{\text{NON CALCULABLE}} \leq \underbrace{\text{erreur sur les données observées} + f \left( \text{nombre de données observées} \right)}_{\text{CALCULABLE (ou estimable)}}$$

# Bornes en généralisation



erreur sur  
toutes les données

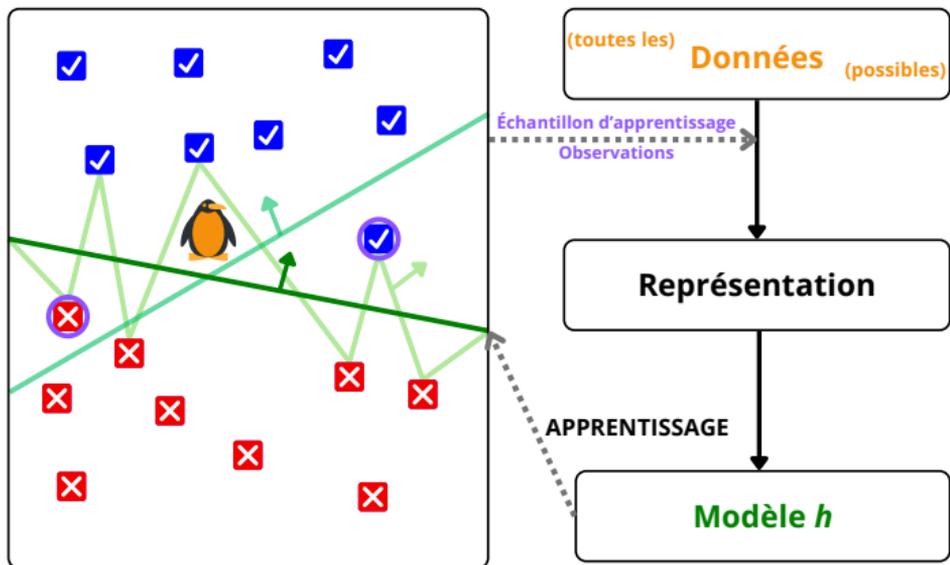
NON CALCULABLE

$\leq$

$+ f \left( \begin{array}{c} \text{complexité} \\ \text{de } h \end{array}, \right)$

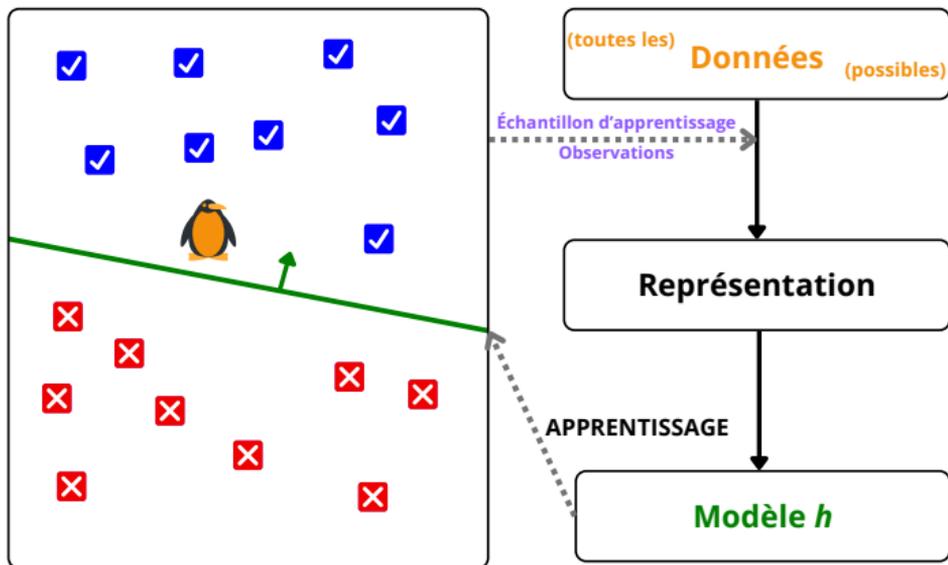
CALCULABLE (ou estimable)

# Bornes en généralisation



$$\underbrace{\text{erreur sur toutes les données}}_{\text{NON CALCULABLE}} \leq \underbrace{\text{erreur sur les données observées} + f \left( \begin{array}{l} \text{complexité} \\ \text{de } h \end{array}, \begin{array}{l} \text{nombre de} \\ \text{données observées} \end{array} \right)}_{\text{CALCULABLE (ou estimable)}}$$

# Bornes en généralisation



$$\text{erreur sur toutes les données} \leq \text{erreur sur les données observées} + f \left( \begin{array}{l} \text{complexité} \\ \text{de } h \end{array}, \begin{array}{l} \text{nombre de} \\ \text{données observées} \end{array} \right)$$

Pour apprendre un modèle avec de **bonnes garanties en généralisation**, on peut minimiser le compromis entre **erreur sur les données observées** et **complexité du modèle**

# Formalisation pour la classification supervisée

$\mathcal{X} \in \mathbb{R}^d$  espace d'entrée

$\mathcal{Y} = \{-1, +1\}$  espace de sortie

$\mathbb{H}$  ensemble de modèles tel que  $\forall h \in \mathbb{H}, h : \mathcal{X} \rightarrow \mathcal{Y}$

$\mathcal{D}$  distribution sur  $\mathcal{X} \times \mathcal{Y}$

$\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \stackrel{iid}{\sim} (\mathcal{D})^m$  ensemble d'apprentissage

Risque empirique associé :  $\hat{R}_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(\mathbf{x}_i) \neq y_i]$

## Classification supervisée

**Objectif :** Trouver le **modèle**  $h$  de  $\mathbb{H}$  qui minimise l'**erreur sur toutes les données**  $R_{\mathcal{D}}(h)$

$$R_{\mathcal{D}}(h) = \underbrace{\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[h(\mathbf{x}) \neq y]}_{\text{Risque réel}}$$

# Formes des bornes en généralisation PAC (Valiant, 1984)

Borne en généralisation dites **Probably Approximately Correct**

Étant donné une distribution  $\mathcal{D}$  sur  $\mathcal{X} \times \mathcal{Y}$ , pour tout  $\delta \in (0, 1]$ , on a

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \underbrace{R_{\mathcal{D}}(h)}_{\text{Risque réel}} \leq \underbrace{\hat{R}_{\mathcal{S}}(h)}_{\text{Risque empirique}} + \underbrace{\Phi(h, \mathcal{S}, \delta)}_{\text{Complexité}} \right] \geq 1 - \delta$$

Avec une grande probabilité d'au moins  $1 - \delta$ , le risque de  $h$  est plus petit que  $\hat{R}_{\mathcal{S}}(h) + \Phi(h, \mathcal{S}, \delta)$

On veut que les bornes soient

- **Informatives**  $\iff \hat{R}_{\mathcal{S}}(h) + \Phi(h, \mathcal{S}, \delta) \leq 1$
- **Précises**  $\iff |R_{\mathcal{D}}(h) - \hat{R}_{\mathcal{S}}(h)| \simeq \Phi(h, \mathcal{S}, \delta)$
- **Fiables**  $\iff$  un petit  $\delta$  augmente la fiabilité de la borne, mais diminue sa précision
- **Calculables (facilement)**  $\iff \hat{R}_{\mathcal{S}}(h) + \Phi(h, \mathcal{S}, \delta)$  est calculable/majorable avec  $\mathcal{S}$

# Formes des bornes en généralisation PAC (Valiant, 1984)

Borne en généralisation dites **Probably Approximately Correct**

Étant donné une distribution  $\mathcal{D}$  sur  $\mathcal{X} \times \mathcal{Y}$ , pour tout  $\delta \in (0, 1]$ , on a

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \underbrace{R_{\mathcal{D}}(h)}_{\text{Risque réel}} \leq \underbrace{\hat{R}_{\mathcal{S}}(h)}_{\text{Risque empirique}} + \underbrace{\Phi(h, \mathcal{S}, \delta)}_{\text{Complexité}} \right] \geq 1 - \delta$$

Avec une grande probabilité d'au moins  $1 - \delta$ , le risque de  $h$  est plus petit que  $\hat{R}_{\mathcal{S}}(h) + \Phi(h, \mathcal{S}, \delta)$

Trois grandes familles de bornes en généralisation PAC :

- Bornes en convergence uniforme (VAPNIK et al., 1971 ; BARTLETT et al., 2002)
- Bornes dépendantes d'un algorithme (BOUSQUET et al., 2002 ; XU et al., 2012)
- Bornes PAC-Bayésiennes (SHAWE-TAYLOR et al., 1997 ; MCALLESTER, 1998)

# Les trois familles de bornes

Borne en convergence uniforme pour un ensemble de modèles  $\mathbb{H}$

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall h \in \mathbb{H}, R_{\mathcal{D}}(h) \leq \hat{R}_{\mathcal{S}}(h) + \Phi(\mathbb{H}, \mathcal{S}, \delta) \right] \geq 1 - \delta$$

Borne en pire cas

Fixe pour tout  $h \in \mathbb{H}$

- Valable pour tous les modèles de  $\mathbb{H}$ 
  - ↪ Bornes en pire cas  $\rightsquigarrow$  souvent trop pessimistes
- $\Phi(h, \mathcal{S}, \delta)$  fixe pour tous les modèles de  $\mathbb{H}$ , par exemple
  - ▶ Dimension de Vapnik-Chervonenkis (VC-dim) de  $\mathbb{H}$ 
    - ✓ Peut être simple à interpréter (capacité d'apprentissage de  $\mathbb{H}$ )
    - ✗ Peut être très grande voir infinie  $\rightsquigarrow$  Borne non informative
  - ▶ Complexité Rademacher de  $\mathbb{H}$ 
    - ✓ Plus précise que la VC-dim (prise en compte de  $\mathcal{D}$ )
    - ✗ Peut être difficile à calculer en pratique

# Les trois familles de bornes

Borne dépendante d'un algorithme  $\mathcal{A} : \mathcal{S} \mapsto h_{\mathcal{S}}$

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ R_{\mathcal{D}}(h_{\mathcal{S}}) \leq \hat{R}_{\mathcal{S}}(h_{\mathcal{S}}) + \Phi(\mathcal{A}, \mathcal{S}, \delta) \right] \geq 1 - \delta$$

Valide uniquement pour  
 $h_{\mathcal{S}}$  appris avec  $\mathcal{A}$  et  $\mathcal{S}$

dépend de caractéristiques  
de l'algorithme  $\mathcal{A}$

- Valable pour un seul modèle  $h_{\mathcal{S}} \rightarrow$  celui appris  $\mathcal{A}$ 
  - $\hookrightarrow$  Plus réaliste mais moins générale
- $\Phi(\mathcal{A}, \mathcal{S}, \delta)$  dépend de propriétés de  $\mathcal{A}$ , par exemple
  - ▶ Stabilité algorithmique
    - ✓ Capture la résistance à de faibles variations dans les données
    - ✗ Peut-être difficile à estimer
  - ▶ Robustesse algorithmique
    - ✓ Capture la résistance aux variations dans une même zone de l'espace
    - ✗ Peut être difficile à calculer en pratique

# Les trois familles de bornes

## Borne PAC-Bayésiennes

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall \rho, \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h) \leq \mathbb{E}_{h \sim \rho} \hat{R}_{\mathcal{S}}(h) + \Phi(\Delta(\rho, \pi), \mathcal{S}, \delta) \right] \geq 1 - \delta$$

Borne en espérance sur  $\mathbb{H}$

Dépend d'une distance  $\Delta$   
entre  $\rho$  et un a priori  $\pi$

- Borne en espérance sur  $\mathbb{H}$ 
  - ↪ souvent plus précises que les bornes en pire cas
  - ↪ borne sur un risque stochastique sur  $\mathbb{H}$
- Mais l'espérance sur  $\mathbb{H}$  est étroitement liée au vote de majorité sur  $\mathbb{H}$
- Peut être facilement calculable (ou majorable)
  - ↪ et minimisable  $\Rightarrow$  algorithme auto-certifié
- $\Phi(\Delta(\rho, \pi), \mathcal{S}, \delta)$  dépend
  - ▶ du choix d'un *a priori*  $\pi$  sur l'espérance
  - ▶ d'une distance  $\Delta(\rho, \pi)$  entre  $\rho$  et  $\pi$  (ex. KL-divergence)

# Les trois familles de bornes

Borne en convergence uniforme pour un ensemble de modèles  $\mathbb{H}$

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall h \in \mathbb{H}, R_{\mathcal{D}}(h) \leq \hat{R}_{\mathcal{S}}(h) + \Phi(\mathbb{H}, \mathcal{S}, \delta) \right] \geq 1 - \delta$$

Borne en pire cas

Fixe pour tout  $h \in \mathbb{H}$

Borne dépendante d'un algorithme  $\mathcal{A} : \mathcal{S} \mapsto h_{\mathcal{S}}$

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ R_{\mathcal{D}}(h_{\mathcal{S}}) \leq \hat{R}_{\mathcal{S}}(h_{\mathcal{S}}) + \Phi(\mathcal{A}, \mathcal{S}, \delta) \right] \geq 1 - \delta$$

Valide uniquement pour  $h_{\mathcal{S}}$  appris avec  $\mathcal{A}$  et  $\mathcal{S}$

dépend de caractéristiques de l'algorithme  $\mathcal{A}$

Borne PAC-Bayésienne

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall \rho, \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h) \leq \mathbb{E}_{h \sim \rho} \hat{R}_{\mathcal{S}}(h) + \Phi(\Delta(\rho, \pi), \mathcal{S}, \delta) \right] \geq 1 - \delta$$

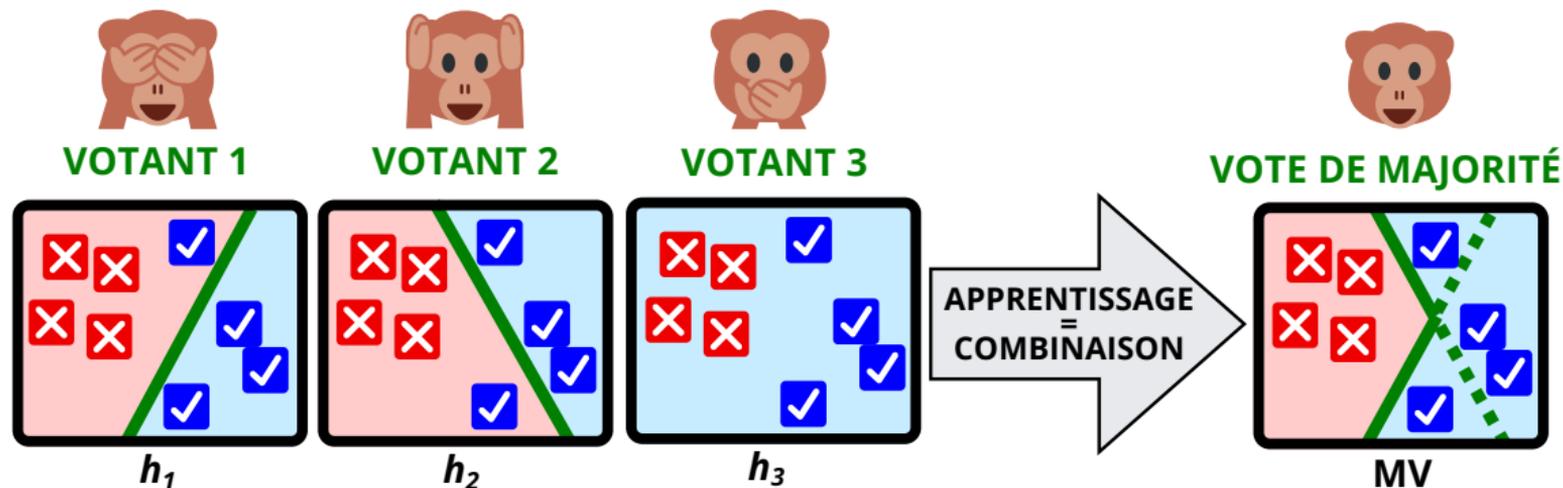
Borne en espérance sur  $\mathbb{H}$

Dépend d'une distance  $\Delta$  entre  $\rho$  et un apriori  $\pi$

# LA THÉORIE PAC-BAYÉSIENNE

Un peu plus en détails

# La théorie PAC-Bayésienne — Le vote de majorité



$\mathbb{H}$  : Ensemble de modèles simples (*i.e.*, de **votants**), ici 3 votants  $\{h_1, h_2, h_3\}$

$\implies$  Apprendre une combinaison **pondérée** des votants  $\iff$  Un vote de majorité **MV**

$$MV(\cdot) = \text{sign} \left[ \sum_{h \in \mathbb{H}} \overbrace{\rho(h)}^{\text{poids de } h} h(\cdot) \right]$$

# Classification supervisée PAC-Bayésienne

$\mathcal{X} \in \mathbb{R}^d$  espace d'entrée

$\mathcal{Y} = \{-1, +1\}$  espace de sortie

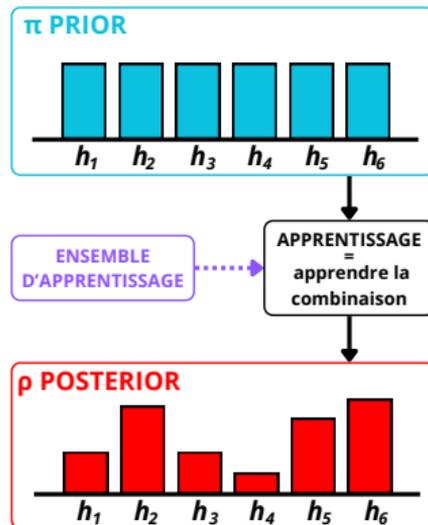
$\mathbb{H}$  ensemble de votants tel que  $\forall h \in \mathbb{H}, h : \mathcal{X} \rightarrow \mathcal{Y}$

$\mathcal{D}$  distribution sur  $\mathcal{X} \times \mathcal{Y}$

$\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \stackrel{iid}{\sim} \mathcal{D}^m$  ensemble d'apprentissage

$\pi$  **distribution prior** sur  $\mathbb{H}$

$\rho$  **distribution posterior** sur  $\mathbb{H}$



## Classification supervisée PAC-Bayésienne

**Objectif :** Trouver le vote de majorité pondéré MV sur  $\mathbb{H}$  qui minimise  $R_{\mathcal{D}}(\text{MV})$

$$R_{\mathcal{D}}(\text{MV}) = \underbrace{\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[\text{MV}(\mathbf{x}) \neq y]}_{\text{Risque réel}} \quad \text{où } \text{MV}(\mathbf{x}) = \text{sign} \left[ \mathbb{E}_{h \sim \rho} h(\mathbf{x}) \right]$$

# Un point clé en PAC-Bayes : Le risque de Gibbs

Risque de Gibbs = espérance des risques individuels

$$R_{\mathcal{D}}(MV) \leq 2 \times \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h)$$

## Avantages :

- Les garanties en **garanties en espérances** sur l'ensemble des votants  $\mathbb{H}$   
 $\implies$  permettent d'avoir des garanties sur un modèle déterministe
- Permet de faire intervenir la **diversité** des votants

## Important :

- D'autres bornes supérieures de  $R_{\mathcal{D}}(MV)$  existent

# Les ingrédients d'une borne PAC-Bayésienne

Rappel : Forme d'une borne en généralisation PAC

$$\text{risque réel} \leq \text{risque empirique} + f \left( \text{complexité}(h), \frac{1}{m} \right)$$

- Risque de Gibbs réel sur la distribution  $\mathcal{D}$  :  $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h)$
- Risque empirique sur  $m$  exemples  $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \sim (\mathcal{D})^m$  :

$$\mathbb{E}_{h \sim \rho} \widehat{R}_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim \rho} \mathbb{I}[h(\mathbf{x}_i) \neq y_i] \quad (\text{risque de Gibbs empirique})$$

- Complexité (divergence de Kullback-Leiber) :  $\text{KL}(\rho \parallel \pi) = \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)}$
- Une mesure de déviation convexe entre les risques réel et empirique  $D : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$

# Borne PAC-Bayésienne générale

avec  $D(a, b) = (a - b)^2$  (borne de la forme de McAllester)

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[ \begin{array}{l} \text{Pour toute distribution posterior } \rho \text{ sur } \mathbb{H}, \\ \underbrace{\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h)}_{\text{risque réel}} \leq \underbrace{\mathbb{E}_{h \sim \rho} \hat{R}_S(h)}_{\text{risque empirique}} + \sqrt{\frac{1}{2m} \left( \underbrace{\text{KL}(\rho \parallel \pi)}_{\text{complexité}} + \ln \frac{2\sqrt{m}}{\delta} \right)} \end{array} \right] \geq 1 - \delta$$

# Borne PAC-Bayésienne générale

avec  $D(a, b) = (a - b)^2$  (borne de la forme de McAllester)

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[ \begin{array}{l} \text{Pour toute distribution posterior } \rho \text{ sur } \mathbb{H}, \\ \underbrace{\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h)}_{\text{risque réel}} \leq \underbrace{\mathbb{E}_{h \sim \rho} \hat{R}_S(h)}_{\text{risque empirique}} + \sqrt{\frac{1}{2m} \left( \underbrace{\text{KL}(\rho \parallel \pi)}_{\text{complexité}} + \ln \frac{2\sqrt{m}}{\delta} \right)} \end{array} \right] \geq 1 - \delta$$

## Théorème PAC-Bayésien général

Pour toute distribution  $\mathcal{D}$  on  $\mathcal{X} \times \mathcal{Y}$ , pour tout ensemble de votants  $\mathbb{H}$ , pour toute distribution prior  $\pi$  sur  $\mathbb{H}$ , pour tout  $\delta \in (0, 1]$ , pour toute fonction convexe  $D : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ , on a

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[ \begin{array}{l} \text{Pour toute distribution posterior } \rho \text{ sur } \mathbb{H}, \\ D \left( \underbrace{\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h)}_{\text{risque réel}}, \underbrace{\mathbb{E}_{h \sim \rho} \hat{R}_S(h)}_{\text{risque empirique}} \right) \leq \frac{1}{m} \left[ \text{KL}(\rho \parallel \pi) + \ln \left( \frac{1}{\delta} \mathbb{E}_{S \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} e^{mD(R_{\mathcal{D}}(h), \hat{R}_S(h))} \right) \right] \end{array} \right] \geq 1 - \delta$$

SUR LA FLEXIBILITÉ DU PAC-BAYES

# Principe simplifié de mes contributions

- Dérivation de **nouvelles bornes** en généralisation PAC-Bayésiennes pour différents cadres
- Dérivation d'algorithmes d'**apprentissage supervisé**
  - ▶ Soit en s'inspirant des bornes
  - ▶ Soit en minimisant directement les bornes  $\rightsquigarrow$  **Algorithmes auto-certifiés**

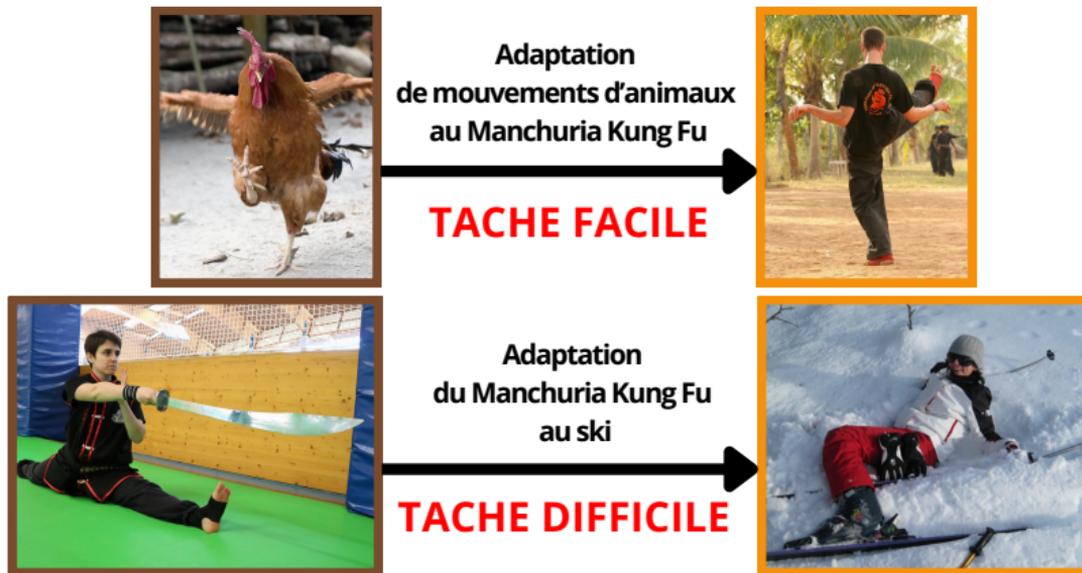
Fil rouge : une borne PAC-Bayésienne

Pour toute distribution  $\mathcal{D}$  sur  $\mathbb{X} \times \mathbb{Y}$ , pour tout ensemble de votants  $\mathbb{H}$ ,  
pour toute **distribution prior**  $\pi$  sur  $\mathbb{H}$ , pour tout  $\delta \in (0, 1]$ , on a

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall \text{ distribution posterior } \rho \text{ sur } h, \right. \\ \left. \underbrace{\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h)}_{\text{risque réel}} \leq \underbrace{\mathbb{E}_{h \sim \rho} \hat{R}_{\mathcal{S}}(h)}_{\text{risque empirique}} + \sqrt{\frac{1}{2m} \left( \underbrace{\text{KL}(\rho \parallel \pi)}_{\text{complexité}} + \ln \frac{2\sqrt{m}}{\delta} \right)} \right] \geq 1 - \delta$$

# Adaptation de domaine d'une tâche source vers une tâche cible

L'être humain est capable de s'adapter à une nouvelle tâche à partir de connaissances acquises



En apprentissage automatique

Se traduit par le besoin d'une mesure de capacité d'adaptation entre les tâches

# Adaptation de domaine d'une tâche source vers une tâche cible

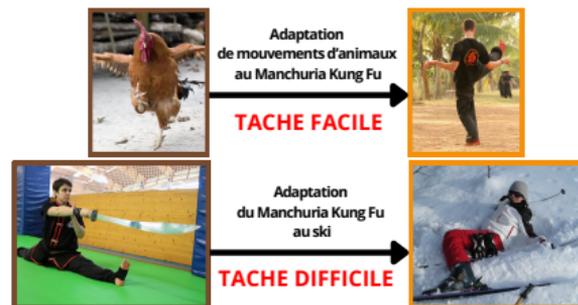
Deux distributions de données sur  $\mathcal{X} \times \mathcal{Y}$

$\mathcal{S}$  domaine source     $\mathcal{T}$  domaine cible

Deux échantillons d'observations

$\mathcal{S} \sim \mathcal{S}$  étiqueté     $\mathcal{T} \sim \mathcal{T}_{\mathcal{X}}$  non étiqueté

$\text{dist}(\mathcal{S}, \mathcal{T})$  : Distance entre domaines



La distance entre les domaines permet de relier les deux domaines

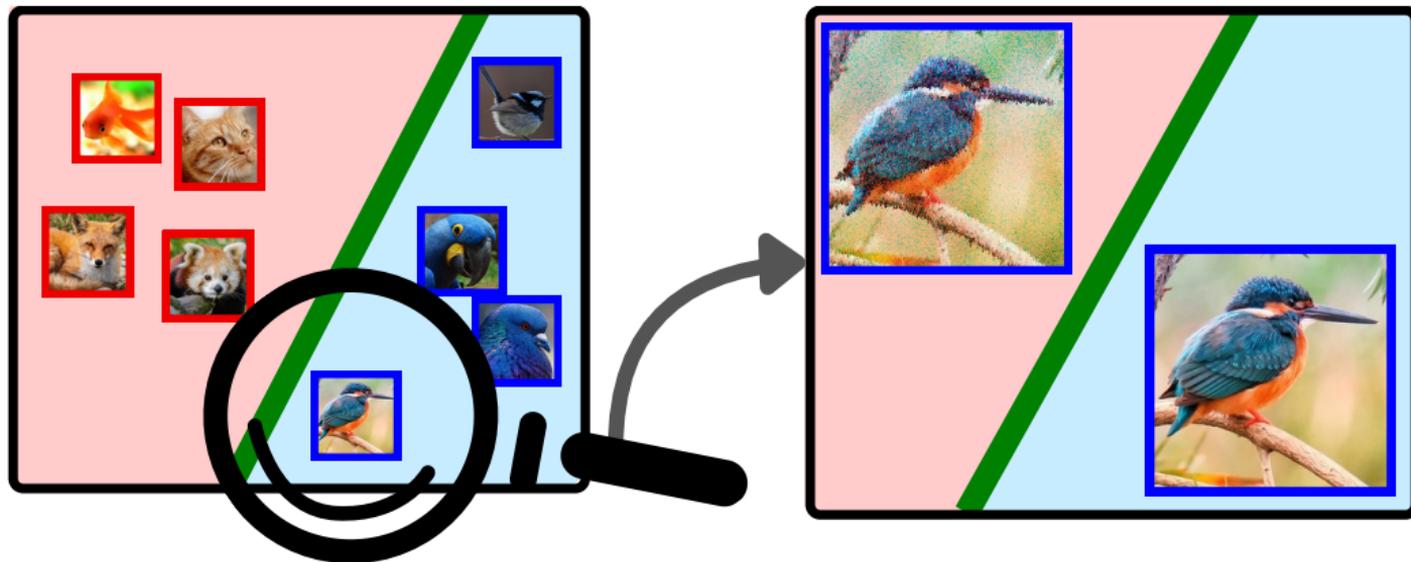
$\Leftrightarrow$  À quel point peut-on utiliser les étiquettes sources pour minimiser le risque cible ?

$$\mathbb{P}_{\substack{\mathcal{S} \sim \mathcal{S} \\ \mathcal{T} \sim \mathcal{T}_{\mathcal{X}}}} \left[ \begin{array}{l} \forall \text{ posterior } \rho \text{ sur } \mathbb{H}, \\ \mathbb{E}_{h \sim \rho} R_{\mathcal{T}}(h) \leq \underbrace{\mathbb{E}_{h, h' \sim \rho} \hat{R}_{\mathcal{T}}(h, h')}_{\text{désaccord (sans étiquette)}} + \frac{1}{\text{dist}(\mathcal{S}, \mathcal{T})} \times \underbrace{\mathbb{E}_{h, h' \sim \rho} \hat{R}_{\mathcal{S}}(h, h')}_{\text{erreur jointe (avec étiquettes)}} + \Phi(\text{KL}(\rho \| \pi), (\mathcal{S} \cup \mathcal{T}), \delta) \end{array} \right] \geq 1 - \delta$$

$$= \sum_{\mathbf{x} \in \mathcal{T}} \frac{1}{|\mathcal{T}|} \mathbb{I}[h(\mathbf{x}) \neq h'(\mathbf{x})] \qquad = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{S}} \frac{1}{|\mathcal{S}|} \mathbb{I}[h(\mathbf{x}) \neq \mathbf{y}] \mathbb{I}[h'(\mathbf{x}) \neq \mathbf{y}]$$

# PAC-Bayes pour la robustesse adverse

Comment garantir que le modèle sera résistant à des attaques malveillantes imperceptibles ?



On veut garantir que le modèle **se trompe le moins possible sur des données bruitées**

# PAC-Bayes pour la robustesse adverse

Espace des données  $\times$  Espace des bruits

$$\underbrace{(\mathcal{X} \times \mathcal{Y})}_{\mathcal{E}} \times \mathcal{B}$$

$\mathcal{E}$  distribution sur  $(\mathcal{X} \times \mathcal{Y}) \times \mathcal{B}$

Ensemble d'apprentissage perturbé/bruité

$$\hat{\mathcal{S}} = \left\{ \left( (\mathbf{x}_i, y_i), \{\epsilon_j^i\}_{j=1}^n \right) \right\}_{i=1}^m \sim (\mathcal{E}^n)^m$$

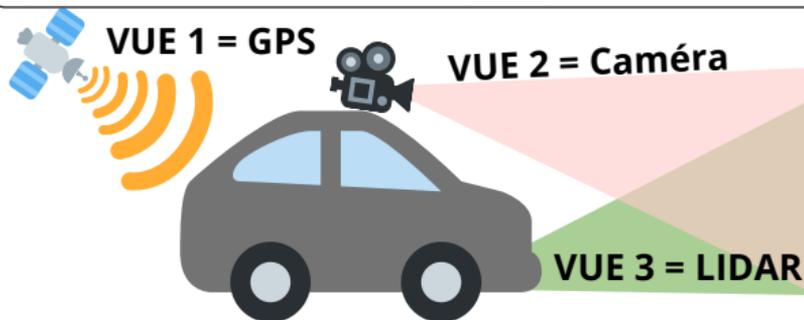
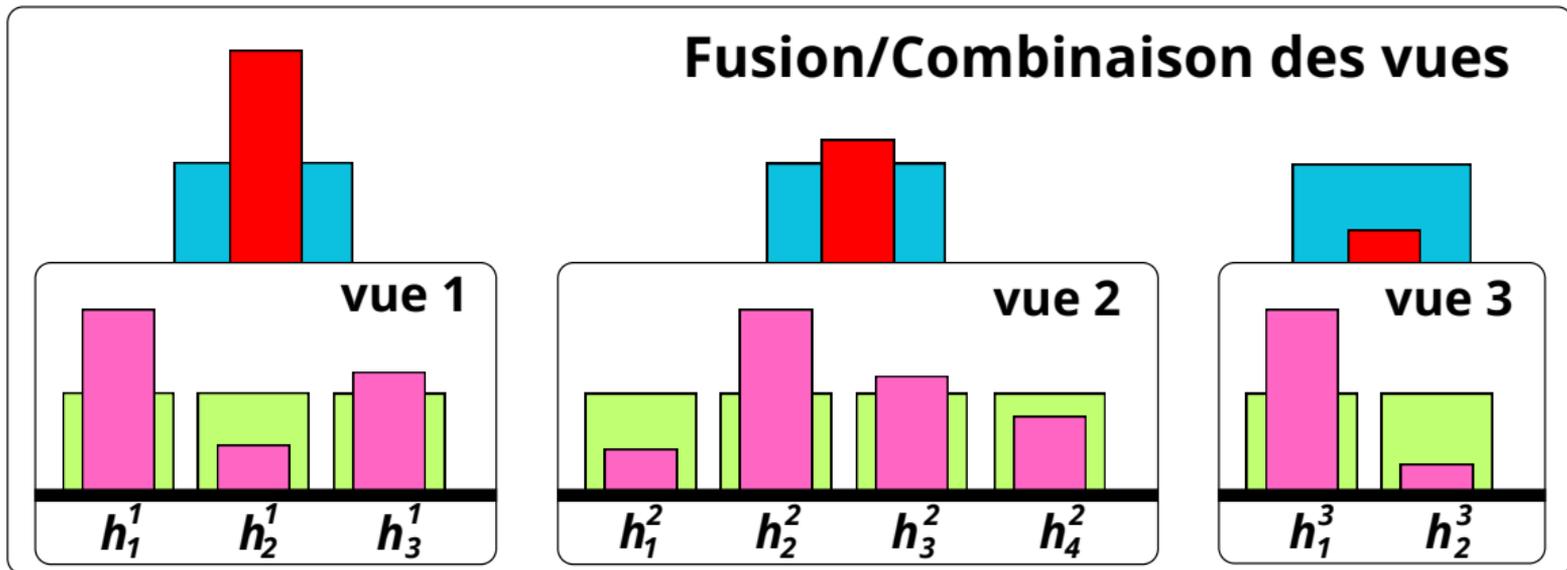
$\Leftrightarrow$  Chaque exemple  $(\mathbf{x}_i, y_i)$  est perturbé avec  $n$  bruits



$$\mathbb{P}_{\hat{\mathcal{S}} \sim (\mathcal{E}^n)^m} \left[ \forall \text{ posterior } \rho \text{ sur } \mathbb{H}, \mathbb{E}_{h \sim \rho} R_{\mathcal{E}}(h) \leq \mathbb{E}_{h \sim \rho} \hat{R}_{\hat{\mathcal{S}}}(h) + \sqrt{\frac{1}{m} \left[ \text{KL}(\rho \| \pi) + \ln \frac{m+1}{\delta} \right]} \right] \geq 1 - \delta$$

$$\text{où } R_{\mathcal{E}}(h) = \mathbb{E}_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} \mathbb{I}[h(\mathbf{x} + \epsilon) \neq y] \quad \text{et} \quad \hat{R}_{\hat{\mathcal{S}}}(h) = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}[h(\mathbf{x}_i + \epsilon_j^i) \neq y_i]$$

## Fusion/Combinaison des vues



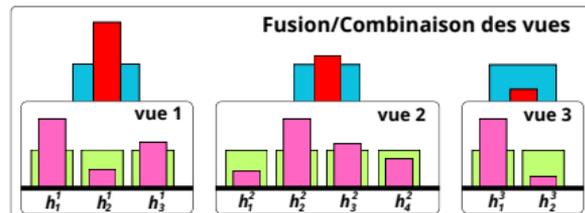
# Apprentissage multi-vues PAC-Bayésien

$\mathbf{V}$  ensemble de  $V$  vues

$\mathcal{X} = \mathcal{X}^1 \times \dots \times \mathcal{X}^V$  espace d'entrée multi-vues

$\forall v \in \mathbf{V}, \mathbb{H}^v$  un ensemble de votants par vue

$\forall h \in \mathbb{H}^v, h : \mathcal{X}^v \rightarrow \mathbb{R}$



## Hiéarchie de distributions sur les votants

$\pi$  distribution hyper-prior sur  $\mathbf{V}$

$\forall v \in \mathbf{V}, P^v$  distribution prior sur  $\mathbb{H}^v$

$\rho$  distribution hyper-posterior sur  $\mathbf{V}$

$\forall v \in \mathbf{V}, Q^v$  distribution posterior sur  $\mathbb{H}^v$

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[ \forall v \in \mathbf{V}, \forall \text{posterior } Q^v \text{ sur } \mathbb{H}^v, \forall \text{hyper-posterior } \rho \text{ sur } \mathbf{V}, \right. \\ \left. \mathbb{E}_{\rho} \mathbb{E}_{h \sim Q^v} R_{\mathcal{D}}(h) \leq \mathbb{E}_{\rho} \mathbb{E}_{h \sim Q^v} \hat{R}_S(h) + \sqrt{\frac{1}{2m} \left( \mathbb{E}_{\rho} \text{KL}(Q^v \| P^v) + \text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right)} \right] \geq 1 - \delta$$

# Le vote de majorité stochastique PAC-Bayésien

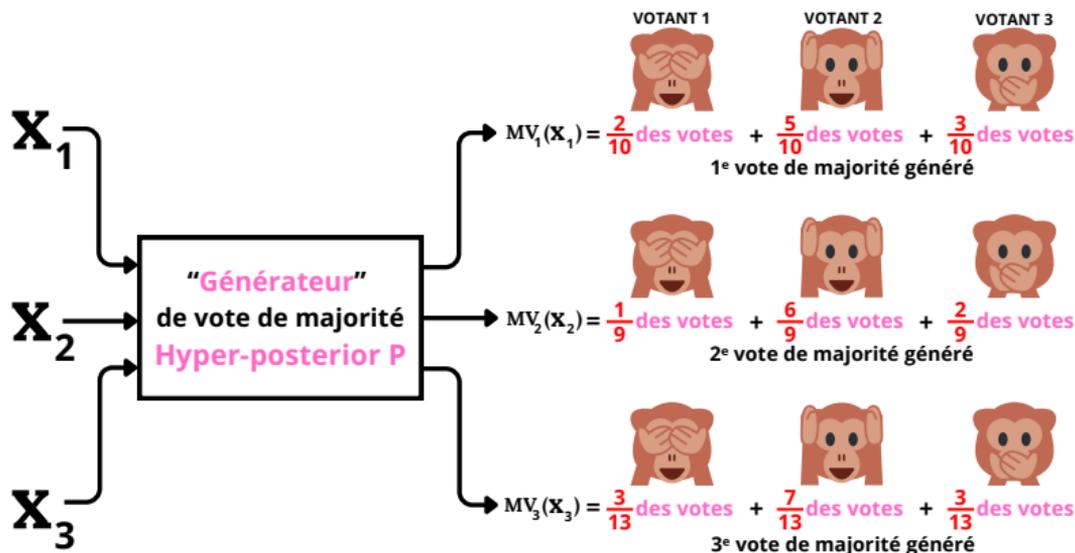
## Rappel

$\mathbb{E}_{h \sim \rho} R(h) \rightsquigarrow$  Les bornes PAC-Bayésiennes sont probabilistes sur  $\mathbb{H}$

$\rightsquigarrow$  Elles portent sur un modèle stochastique, plutôt qu'un modèle déterministe

**MAIS** elles permettent de borner le risque du vote de majorité  $MV(\cdot) = \text{sign} \left[ \mathbb{E}_{h \sim \rho} h(\cdot) \right]$

$\rightarrow$  Via des relaxations *plus ou moins* précises et pas nécessairement facile à minimiser



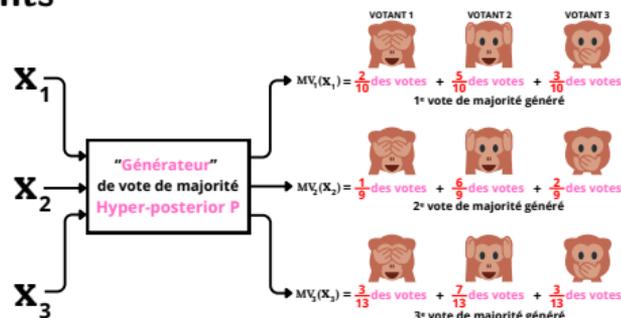
# Le vote de majorité stochastique PAC-Bayésien

## Une sorte de hiérarchie de distributions sur les votants

$\Pi$  hyper-prior sur  $\mathbb{H}$  (distribution de dirichlet)

$\mathbb{P}$  hyper-posterior  $\mathbb{H}$   $\rho$  posterior tirée selon  $\mathbb{P}$

$$\hookrightarrow MV(\cdot) = \text{sign} \left[ \mathbb{E}_{h \sim \rho} h(\cdot) \right]$$



$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall \text{ hyper-posterior } \mathbb{P} \text{ sur } \mathbb{H}, \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{D}}(MV) \leq \mathbb{E}_{\rho \sim \mathbb{P}} \hat{R}_{\mathcal{S}}(MV) + \sqrt{\frac{1}{2m} \left( \text{KL}(\mathbb{P} \parallel \Pi) + \ln \frac{2\sqrt{m}}{\delta} \right)} \right] \geq 1 - \delta$$

# Désintégration des bornes PAC-Bayésiennes

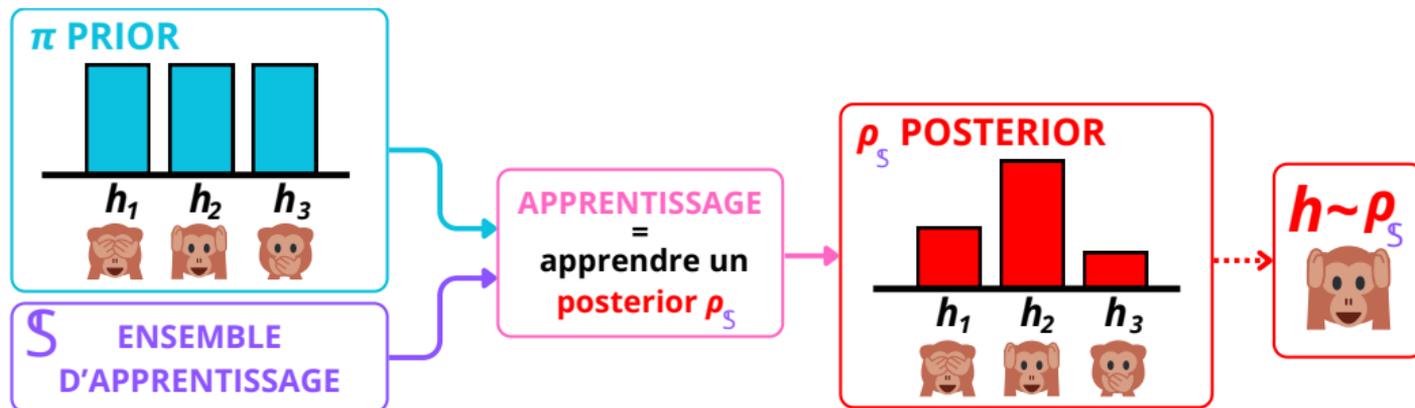
Rappel

$\mathbb{E}_{h \sim \rho} R(h) \rightsquigarrow$  Les bornes PAC-Bayésiennes sont probabilistes sur  $\mathbb{H}$

$\rightsquigarrow$  Elles portent sur un modèle stochastique, plutôt qu'un modèle déterministe



Désintégrer les bornes permet de ne les faire porter que sur un **seul** modèle



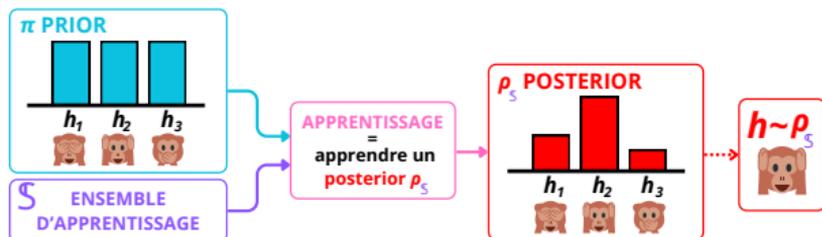
# Désintégration des bornes PAC-Bayésiennes

$\mathcal{S} \sim \mathcal{D}^m$  ensemble d'apprentissage

$\pi$  distribution prior sur  $\mathbb{H}$

$A$  algorithme déterministe

$\hookrightarrow A(\mathcal{S}, \pi) = \rho_{\mathcal{S}}$  distribution posterior sur  $\mathbb{H}$



$$\mathbb{P}_{\substack{\mathcal{S} \sim \mathcal{D}^m \\ h \sim \rho_{\mathcal{S}}}} \left[ R_{\mathcal{D}}(h) \leq \hat{R}_{\mathcal{S}}(h) + \sqrt{\frac{1}{2m} \left[ D_{\alpha}(\rho_{\mathcal{S}} \parallel \pi) + \frac{2\alpha-1}{\alpha-1} \ln \frac{1}{\alpha} + \ln(2\sqrt{m}) \right]} \right] \geq 1 - \delta$$

Divergence de Rényi :

$$D_{\alpha}(\rho_{\mathcal{S}} \parallel \pi) = \frac{1}{\alpha-1} \ln \left[ \mathbb{E}_{h \sim \pi} \left[ \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right]^{\alpha} \right]$$

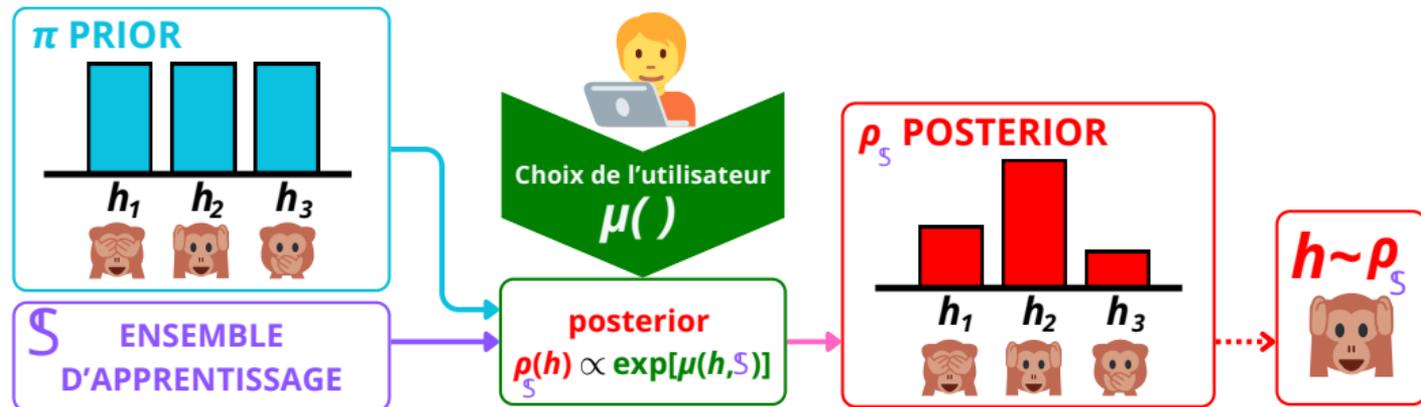
# Complexités arbitraires

Une remarque importante

Le terme de complexité dans les bornes en généralisation dépend du cadre considéré



Comment rendre les bornes plus flexibles ?



# Complexités arbitraires

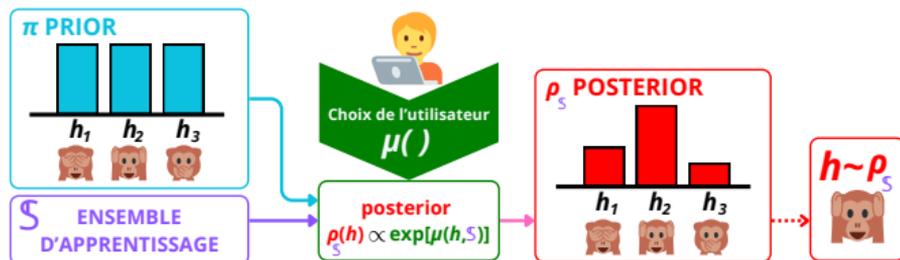
On considère une **fonction personnalisable**  $\mu : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}$  (e.g., une régularisation)

$\mathcal{S} \sim \mathcal{D}^m$  ensemble d'apprentissage

$\pi$  distribution **prior** sur  $\mathbb{H}$

$\rho_{\mathcal{S}}$  distribution **posterior** sur  $\mathbb{H}$

telle que  $\rho_{\mathcal{S}}(h) \propto e^{\mu(h, \mathcal{S})}$



$$\text{Forme générale : } \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[ R_{\mathcal{D}}(h) \leq \hat{R}_{\mathcal{S}}(h) + \Phi_{\mu}(h, \mathcal{S}, \delta) \right] \geq 1 - \delta$$

$$\mathbb{P}_{\substack{h' \sim \pi, \\ \mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}}} \left[ R_{\mathcal{D}}(h) \leq \hat{R}_{\mathcal{S}}(h) + \sqrt{\frac{1}{2m} \left[ \underbrace{\mu(h', \mathcal{S})}_{\text{Prior}} - \underbrace{\mu(h, \mathcal{S})}_{\text{Posterior}} + \frac{8\sqrt{m}}{\delta^2} \right]_+} \right] \geq 1 - \delta$$

# Propriétés clés des bornes PAC-Bayésiennes

- **Flexibles** et adaptables à différents scénarios
  - ▶ Compatibles avec différents types de modèles
    - ▶ Modèle stochastique par nature
    - ▶ Modèle déterministe par nature via les votes de majorité
    - ▶ Modèle déterministe via la désintégration
  - ▶ (Hiérarchie de) distributions sur des ensembles (de modèles ou d'autres objets)
- **Intermédiaires** entre bornes en convergence uniforme et dépendante d'un algorithme
  - ▶ Les complexités arbitraires permettent même de retrouver ces deux types de bornes
- **Précises** et informatives
  - ▶ Particulièrement pour affiner un modèle *a priori*
- **Pratiques** et exploitables  $\rightsquigarrow$  Souvent calculables ou majorables
  - ▶ Minimisables directement via un algorithme auto-certifié
  - ▶ Meilleure explicabilité et confiance des modèles par nature de l'algorithme

# Forme générale des résultats

Forme générale des bornes PAC-Bayésiennes

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \forall \rho, D \left( \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}(h), \mathbb{E}_{h \sim \rho} \hat{R}_{\mathcal{S}}(h) \right) \leq \Phi(\Delta(\rho, \pi), \mathcal{S}, \delta) \right] \geq 1 - \delta$$

**Borne en espérance sur  $\mathbb{H}$**

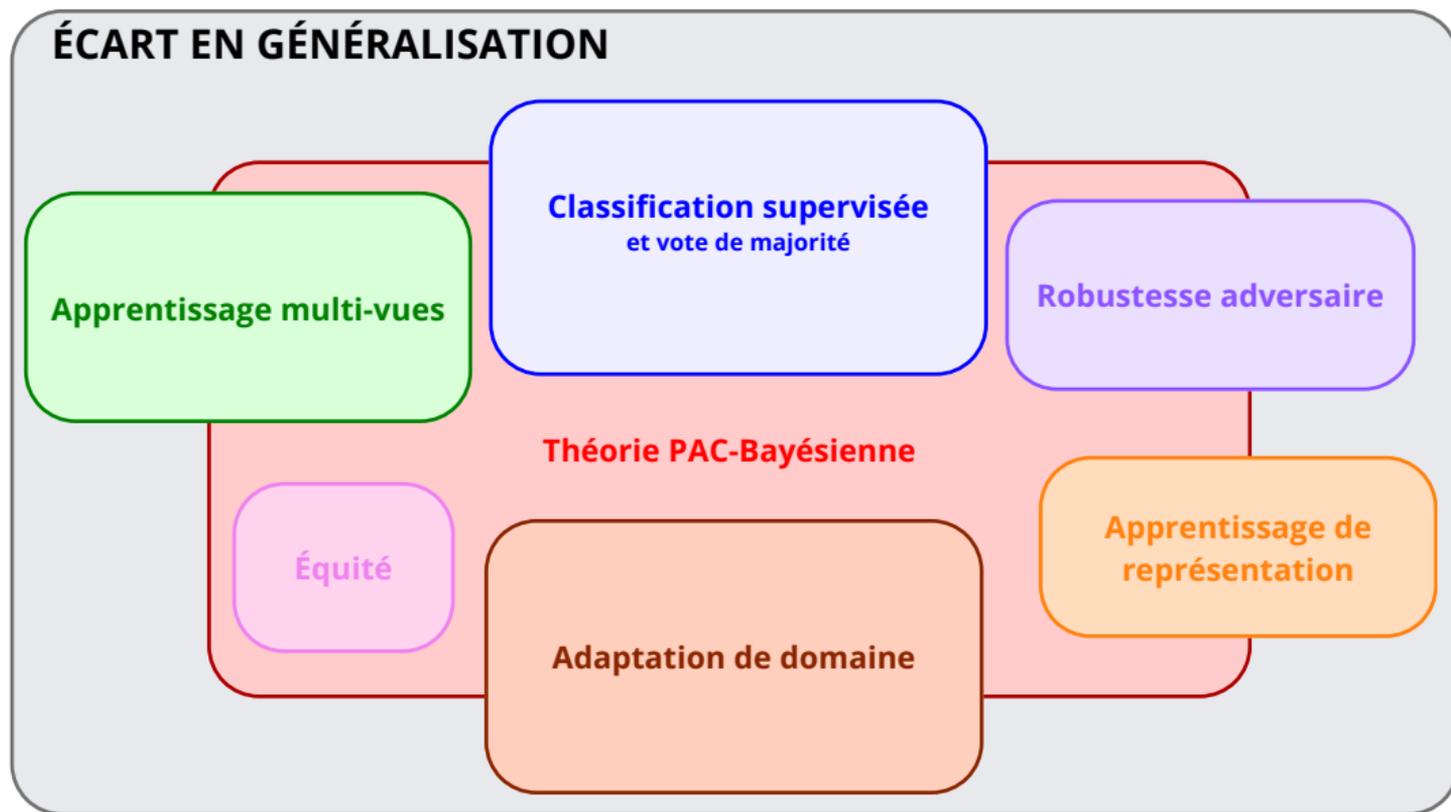
$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[ D \left( R_{\mathcal{D}}(h), \hat{R}_{\mathcal{S}}(h) \right) \leq \Phi(\Delta(\rho_{\mathcal{S}}, \pi), \mathcal{S}, \delta) \right] \geq 1 - \delta$$

**Borne pour un unique  $h \sim \rho_{\mathcal{S}}$**       **OU Peut être personnalisable**



$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[ \underbrace{D \left( R_{\mathcal{D}}(\text{modèle}), \hat{R}_{\mathcal{S}}(\text{modèle}) \right)}_{\text{écart en généralisation}} \leq \text{Borne} \right] \geq 1 - \delta$$

**écart en généralisation**  
distance entre risque réel et empirique



AU-DELÀ DE  
L'ÉCART EN GÉNÉRALISATION

# Projet de recherche : Au-delà de l'écart en généralisation

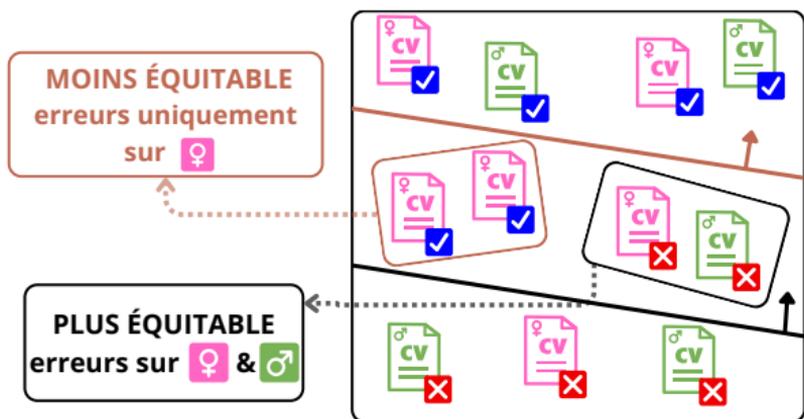
Tâche de classification supervisée équitable

Classer des CVs en 2 catégories :

✓ pertinent ✗ non pertinent

**Objectif** : Trouver un modèle équitable et garantir de bonnes performances

⇒ Ne pas apprendre les biais !



# Projet de recherche : Au-delà de l'écart en généralisation

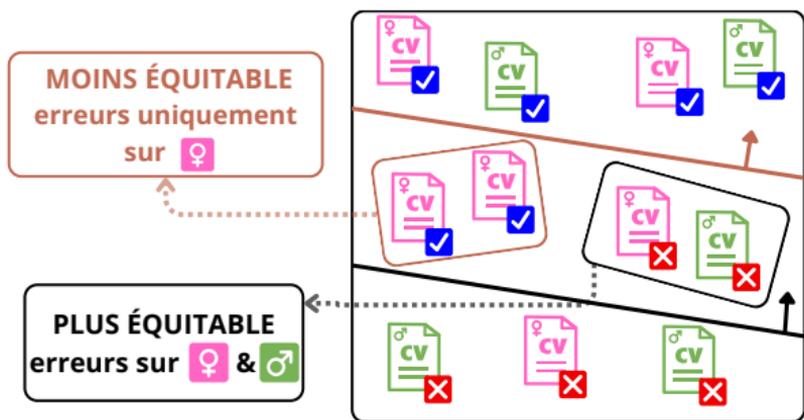
Tâche de classification supervisée équitable

Classer des CVs en 2 catégories :

✓ pertinent ✗ non pertinent

**Objectif** : Trouver un modèle équitable et garantir de bonnes performances

⇒ Ne pas apprendre les biais !



**Une solution** : Considérer une contrainte/mesure d'équité

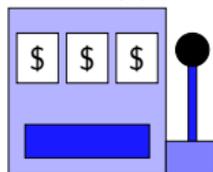
↪ **Équité** = capacité à bien classer malgré les biais dans les données

- Algorithme auto-certifié PAC-Bayésien pour garantir un compromis performance/équité
  - ▶ Contrôle de l'écart entre les **biais observés** et l'**équité espérée**
- Construire des mesures d'équité robustes aux biais

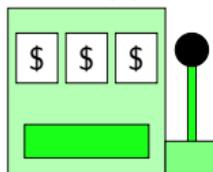
# Projet de recherche : Au-delà de l'écart en généralisation

Les bandits manchots

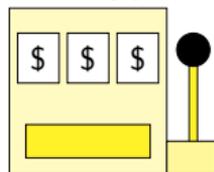
Bras 1 –  $r_t(1) \sim \mathcal{D}_1$



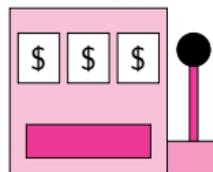
Bras 2 –  $r_t(2) \sim \mathcal{D}_2$



Bras 3 –  $r_t(3) \sim \mathcal{D}_3$



Bras 4 –  $r_t(4) \sim \mathcal{D}_4$



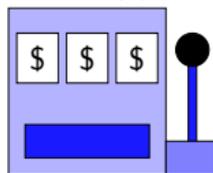
Quel bras/machine choisir à l'étape  $t$  pour maximiser le gain cumulé (*reward*)

**Objectif** : Trouver la meilleure stratégie (*policy*) pour maximiser le gain cumulé

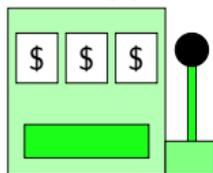
# Projet de recherche : Au-delà de l'écart en généralisation

Les bandits manchots

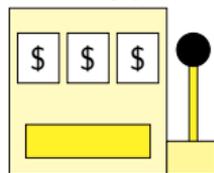
Bras 1 -  $r_t(1) \sim \mathcal{D}_1$



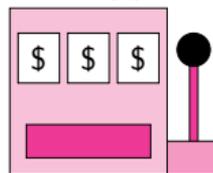
Bras 2 -  $r_t(2) \sim \mathcal{D}_2$



Bras 3 -  $r_t(3) \sim \mathcal{D}_3$



Bras 4 -  $r_t(4) \sim \mathcal{D}_4$



Quel bras/machine choisir à l'étape  $t$  pour maximiser le gain cumulé (*reward*)

**Objectif** : Trouver la meilleure stratégie (*policy*) pour maximiser le gain cumulé

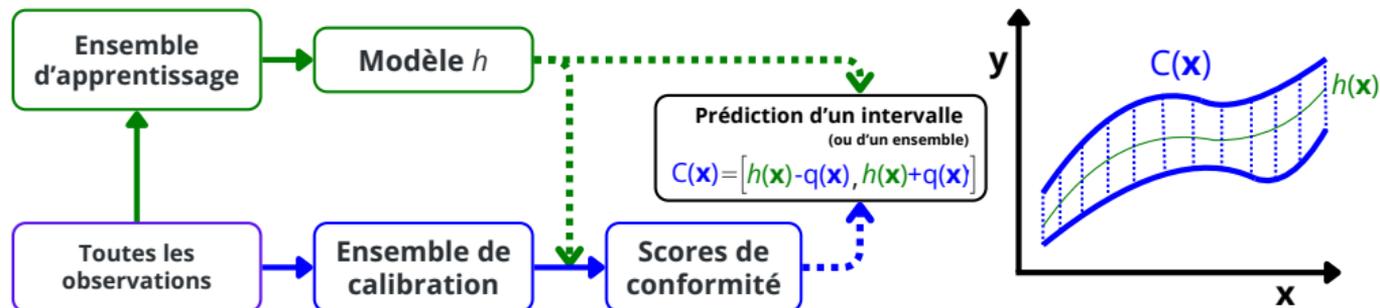
**Une solution** : Minimiser le regret au cours du temps

↪ **Regret** = Écart entre **gain cumulé actuel** et **gain optimal** (si le meilleur bras était tjrs choisi)

- Algorithme auto-certié PAC-Bayésien pour garantir un regret faible
- Analyse dynamique du regret au cours du temps
  - ▶ Étudier comment le PAC-Bayes peut équilibrer exploration et exploitation

# Projet de recherche : Au-delà de l'écart en généralisation

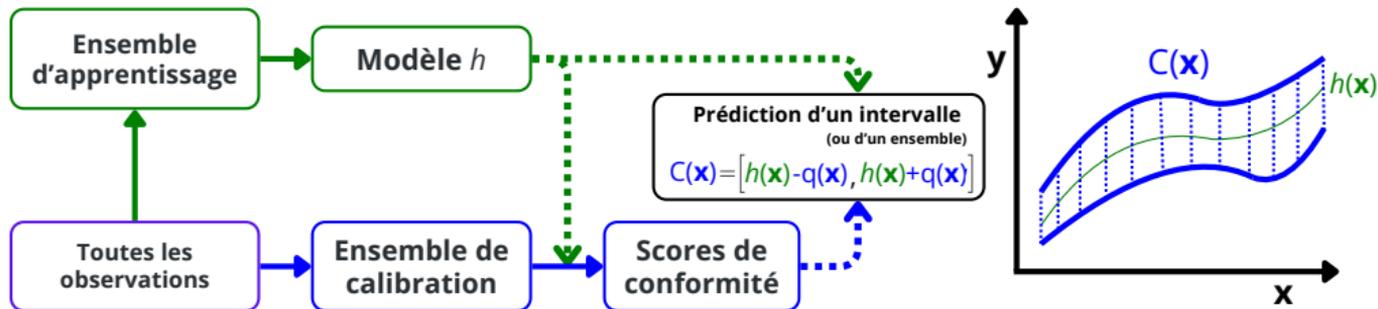
la prédiction conforme : quantification de l'incertitude



**Objectif** : Comment s'assurer qu'un modèle prédit des intervalles fiables, pas trop grands ?

# Projet de recherche : Au-delà de l'écart en généralisation

la prédiction conforme : quantification de l'incertitude



**Objectif** : Comment s'assurer qu'un modèle prédit des intervalles fiables, pas trop grands ?

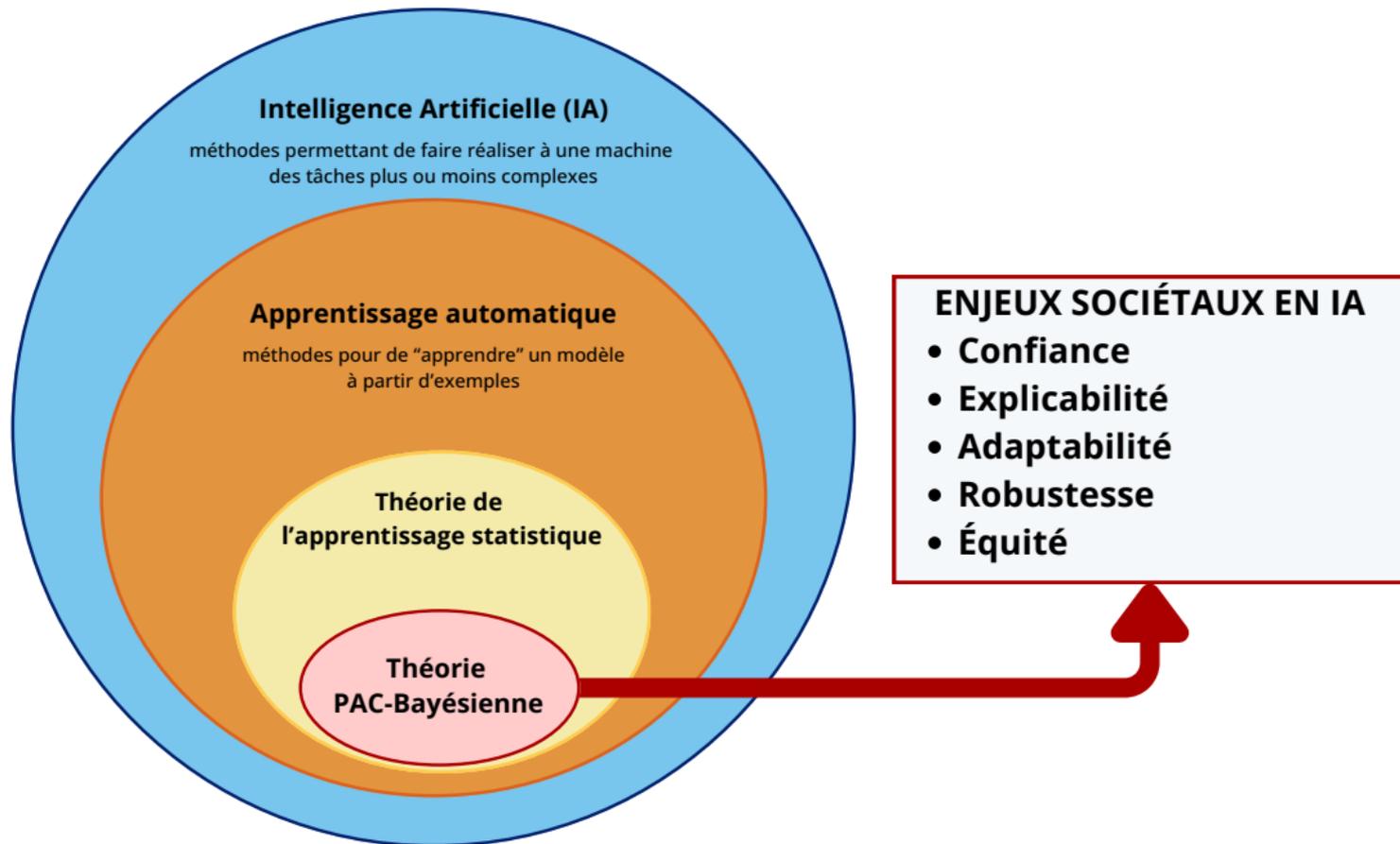
**Une solution** : Contrôler un compromis entre couverture et efficacité

↔ **couverture** = l'intervalle doit contenir la vraie valeur dans  $X\%$  des cas

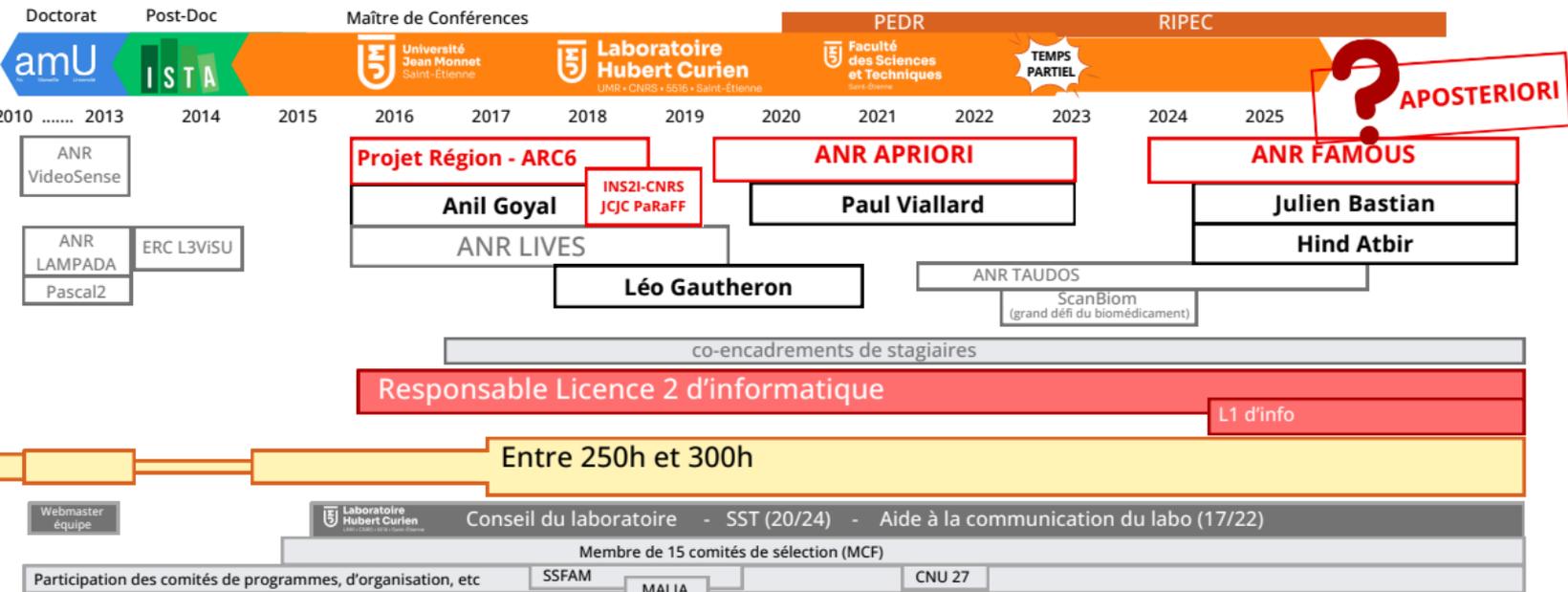
↔ **efficacité** = l'intervalle doit être le plus petit possible

- Algorithme auto-certié pour garantir un compromis couverture/efficacité
- Ajustement de la largeur des intervalles via le PAC-Bayes

# Projet de recherche : Au-delà de l'écart en généralisation



# MERCI POUR VOTRE ATTENTION



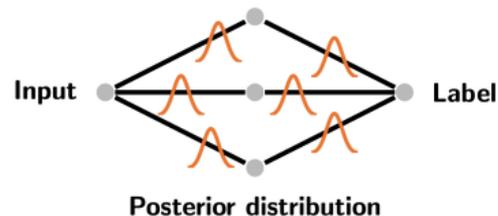
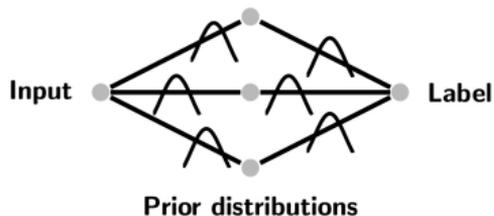
# Références I

-  Vladimir VAPNIK et Alexey CHERVONENKIS. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability and its Applications*. (1971).
-  Leslie VALIANT. A Theory of the Learnable. *Communications of the ACM*. (1984).
-  John SHAWE-TAYLOR et Robert WILLIAMSON. A PAC Analysis of a Bayesian Estimator. *Annual Conference on Learning Theory*. (1997).
-  David MCALLESTER. Some PAC-Bayesian Theorems. *Annual Conference on Learning Theory*. (1998).
-  Peter BARTLETT et Shahar MENDELSON. Rademacher and Gaussian Complexities : Risk Bounds and Structural Results. *Journal of Machine Learning Research*. (2002).
-  Olivier BOUSQUET et André ELISSEEFF. Stability and Generalization. *Journal of Machine Learning Research*. (2002).
-  Huan XU et Shie MANNOR. Robustness and Generalization. *Machine Learning*. (2012).

# Algorithme de minimisation de la borne désintégrée

## Apprentissage d'un réseau de neurone avec une borne PAC-Bayésienne désintégrée

- 1 Prior/posterior Gaussian distributions associated with the weights of the neural network



- 2 Learn  $T$  priors  $\mathbb{P} = \{\pi_t\}_{t=1}^T$  in  $T$  epochs with  $\mathcal{S}'$



- 3 Learn the distribution  $\rho_{\mathcal{S}}$  with  $\mathcal{S}$  from a prior  $\pi_t$  selected with  $\mathcal{S}$

- 4 Sample the neural network  $h \sim \rho_{\mathcal{S}}$

