

# Apprentissage par renforcement pour le contrôle de processus de Markov déterministe par morceaux

Application à l'optimisation d'un traitement médical

Orlane Rossini <sup>1</sup>, Alice Cleyen <sup>1,2</sup>, Benoîte de Saporta <sup>1</sup>,  
Régis Sabbadin <sup>3</sup> et Meritxell Vinyals <sup>3</sup>

<sup>1</sup>IMAG, Univ Montpellier, CNRS, Montpellier, France

<sup>2</sup>John Curtin School of Medical Research, The Australian National University,  
Canberra, ACT, Australia

<sup>3</sup>Univ Toulouse, INRAE-MIAT, Toulouse, France

Octobre 2024



UNIVERSITÉ DE  
MONTPELLIER

INRAE

IMAG  
INSTITUT MONTPELLIERAIN  
ALEXANDER GROTHENDIECK



anr<sup>®</sup>

# Le contexte médical

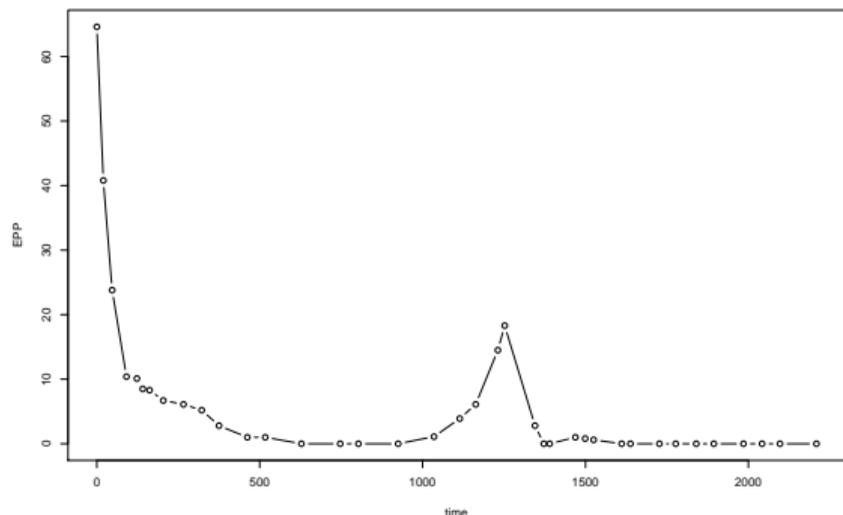


FIGURE: Exemple de donnée d'un patient<sup>a</sup>

- Des patients ayant eu un **cancer** bénéficient d'un **suivi régulier**;
- La concentration d'**immunoglobuline clonale** est mesurée **dans le temps**;
- Le médecin doit prendre de nouvelles **décisions** à chaque visite.

<sup>a</sup>IUCT Oncopole et CRCT, Toulouse, France

# Le contexte médical

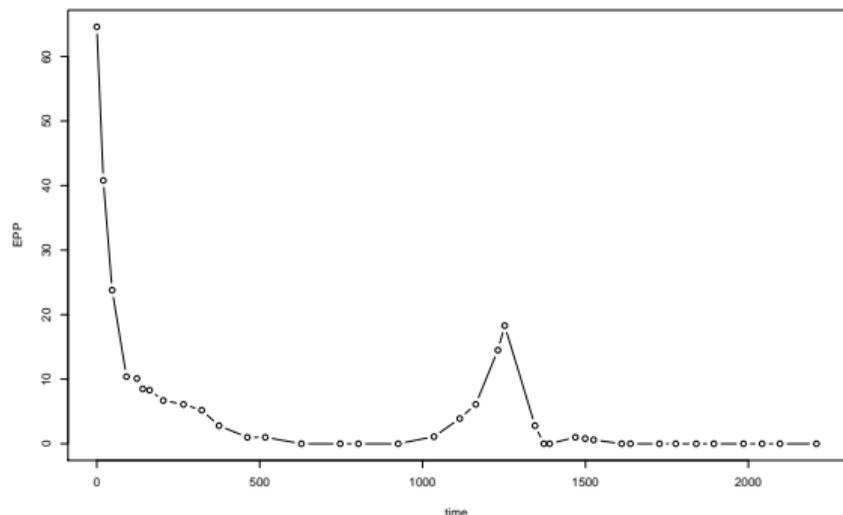
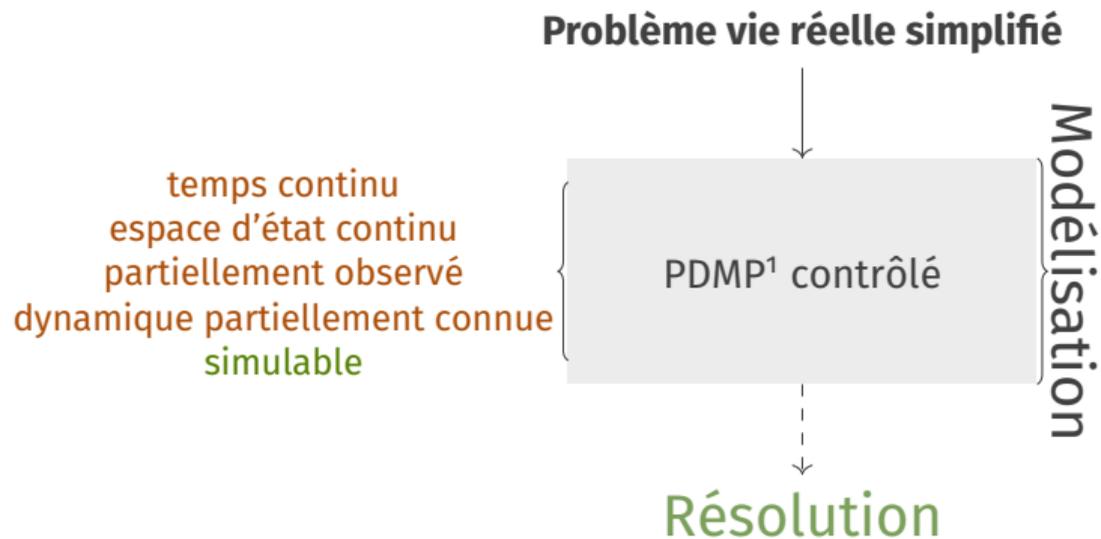


FIGURE: Exemple de donnée d'un patient<sup>a</sup>

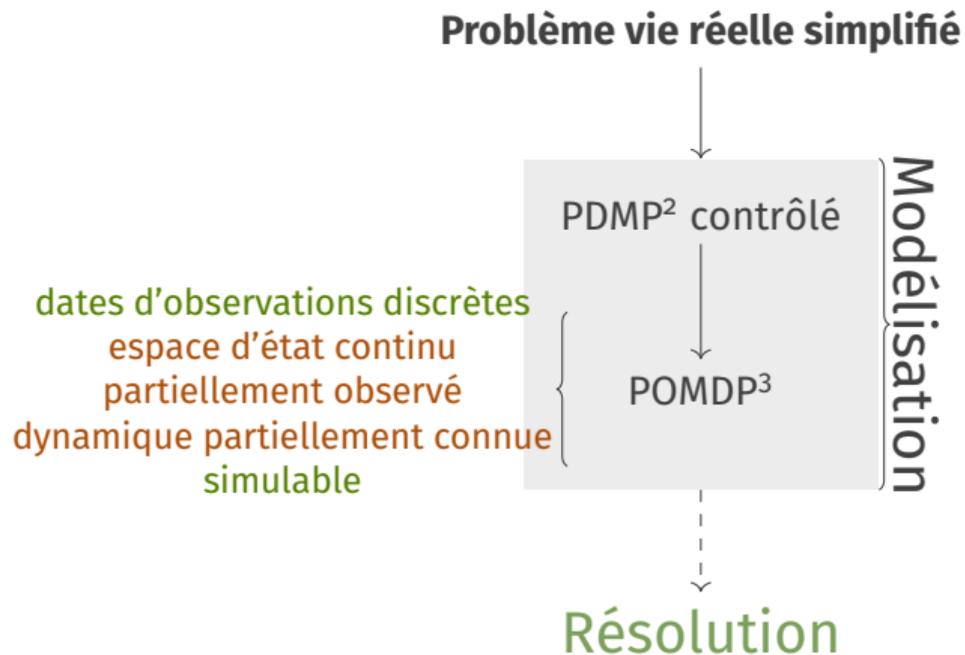
- Des patients ayant eu un **cancer** bénéficient d'un **suivi régulier**;
- La concentration d'**immunoglobuline clonale** est mesurée **dans le temps**;
- Le médecin doit prendre de nouvelles **décisions** à chaque visite.

⇒ **Optimiser la prise de décision pour assurer la qualité de vie du patient**

<sup>a</sup>IUCT Oncopole et CRCT, Toulouse, France



<sup>1</sup>Piecewise Deterministic Markov Process



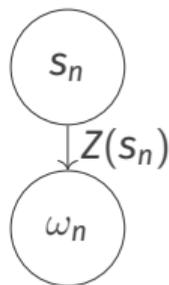
<sup>2</sup>Piecewise-deterministic Markov process

<sup>3</sup>Partially Observable Markov Decision Process

# Caractéristiques d'un POMDP

Agent

Environnement

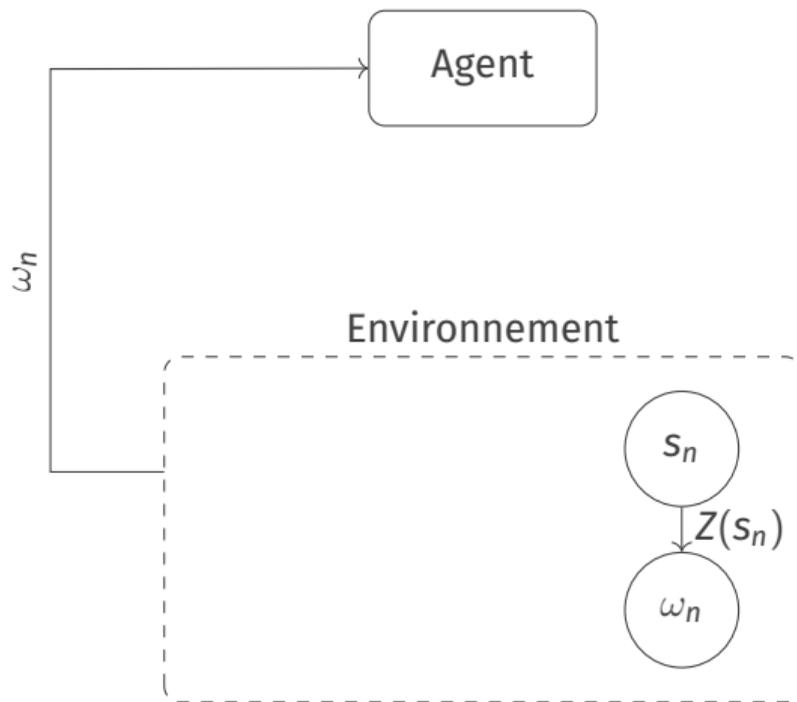


## DÉFINITION D'UN POMDP

Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

# Caractéristiques d'un POMDP

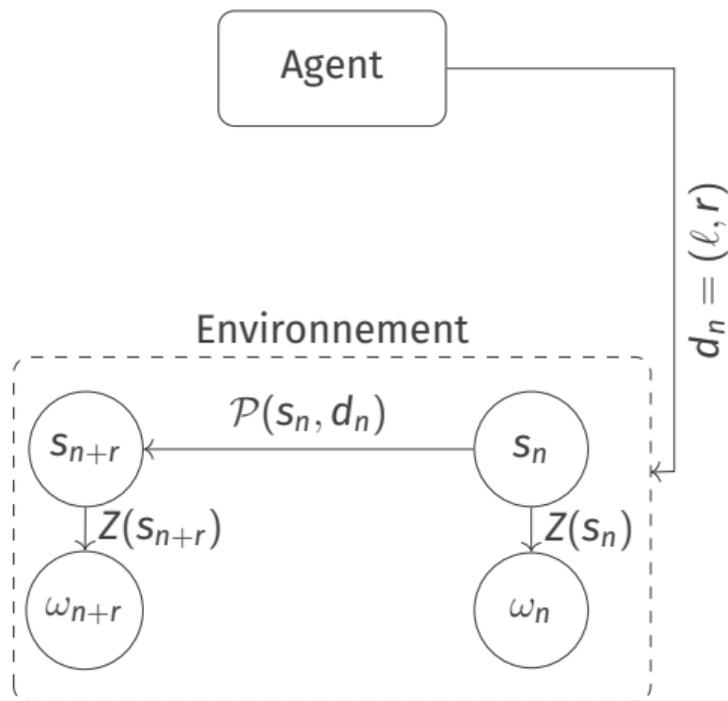


## DÉFINITION D'UN POMDP

Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

# Caractéristiques d'un POMDP

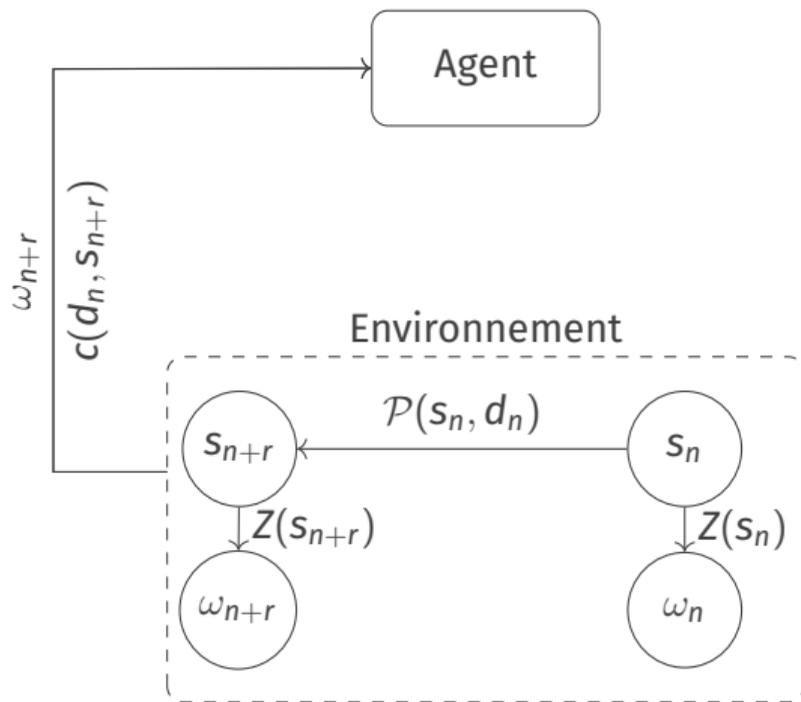


## DÉFINITION D'UN POMDP

Un POMDP se définit par un tuple  $(\mathcal{S}, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in \mathcal{S}$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

# Caractéristiques d'un POMDP



## DÉFINITION D'UN POMDP

Un POMDP se définit par un tuple  $(S, \mathcal{D}, \mathcal{K}, \mathcal{P}, \Omega, \mathcal{Z}, C)$ .

- Etat du patient  $s = (m, k, \zeta, u, t, \tau) \in S$ ;
- Décisions  $d = (\ell, r) \in \mathcal{D}$ ;
- $\mathcal{K}(s) \subseteq \mathcal{D}$  l'espace des décisions admissibles dans l'état  $s$ ;
- Probabilité de transition  $\mathcal{P}(s, d)(s')$ ;
- Observation  $\omega = (\tau, t, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Fonction d'observation  $\mathcal{Z}(s)(\omega)$ ;
- Fonction coût  $C(s, d, s')$ .

# Résoudre un POMDP

**Identifier une politique  $\pi : h \rightarrow d$  optimale déterministe!**

On n'observe pas l'espace d'état ! Soit l'historique  $h = (\omega_0, d_1, \omega_1, d_2, \dots, \omega_n)$

# Résoudre un POMDP

**Identifier une politique  $\pi : h \rightarrow d$  optimale déterministe!**

On n'observe pas l'espace d'état ! Soit l'historique  $h = (\omega_0, d_1, \omega_1, d_2, \dots, \omega_n)$

$$\underbrace{V(\pi, h)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}_h^\pi \left[ \sum_{n=1}^{H-1} c(D_n, S_n) \right]}_{\text{Coût attendu à long terme suite à la politique menée } \pi}$$

# Résoudre un POMDP

**Identifier une politique  $\pi : h \rightarrow d$  optimale déterministe!**

On n'observe pas l'espace d'état ! Soit l'historique  $h = (\omega_0, d_1, \omega_1, d_2, \dots, \omega_n)$

$$\underbrace{V(\pi, h)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}_h^\pi \left[ \sum_{n=1}^{H-1} c(D_n, S_n) \right]}_{\text{Coût attendu à long terme suite à la politique menée } \pi}$$

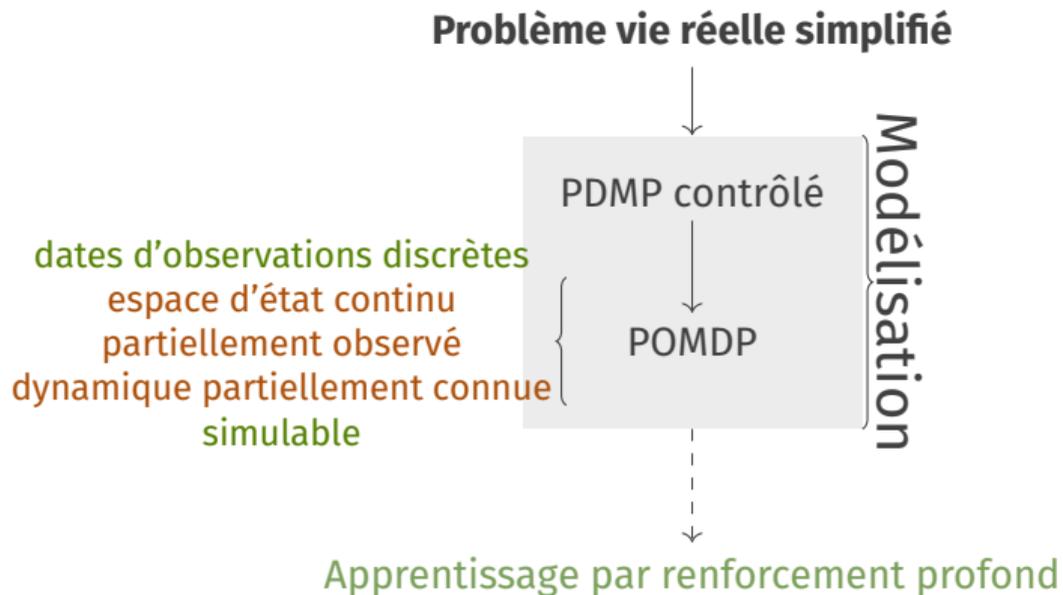
# Résoudre un POMDP

**Identifier une politique  $\pi : h \rightarrow d$  optimale déterministe!**

On n'observe pas l'espace d'état ! Soit l'historique  $h = (\omega_0, d_1, \omega_1, d_2, \dots, \omega_n)$

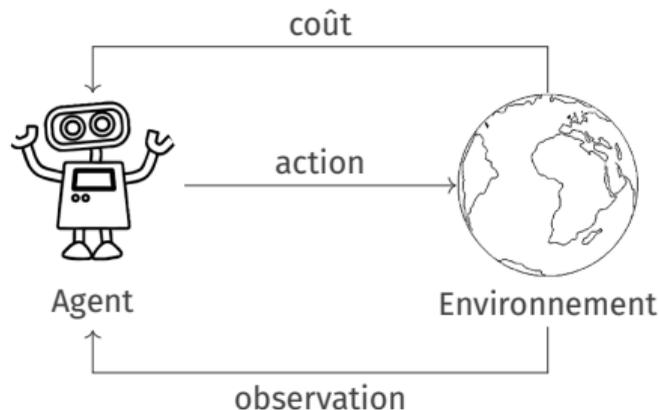
$$\underbrace{V(\pi, h)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}_h^\pi \left[ \sum_{n=1}^{H-1} c(D_n, S_n) \right]}_{\text{Coût attendu à long terme suite à la politique menée } \pi}$$

$$\underbrace{V^*(h)}_{\text{Fonction valeur}} = \underbrace{\min_{\pi \in \Pi} V(\pi, h)}_{\text{Minimisation sur l'ensemble des politiques } \Pi.}$$



<sup>4</sup>Univ Toulouse, INRAE-MIAT, Toulouse, France

# Apprentissage par renforcement



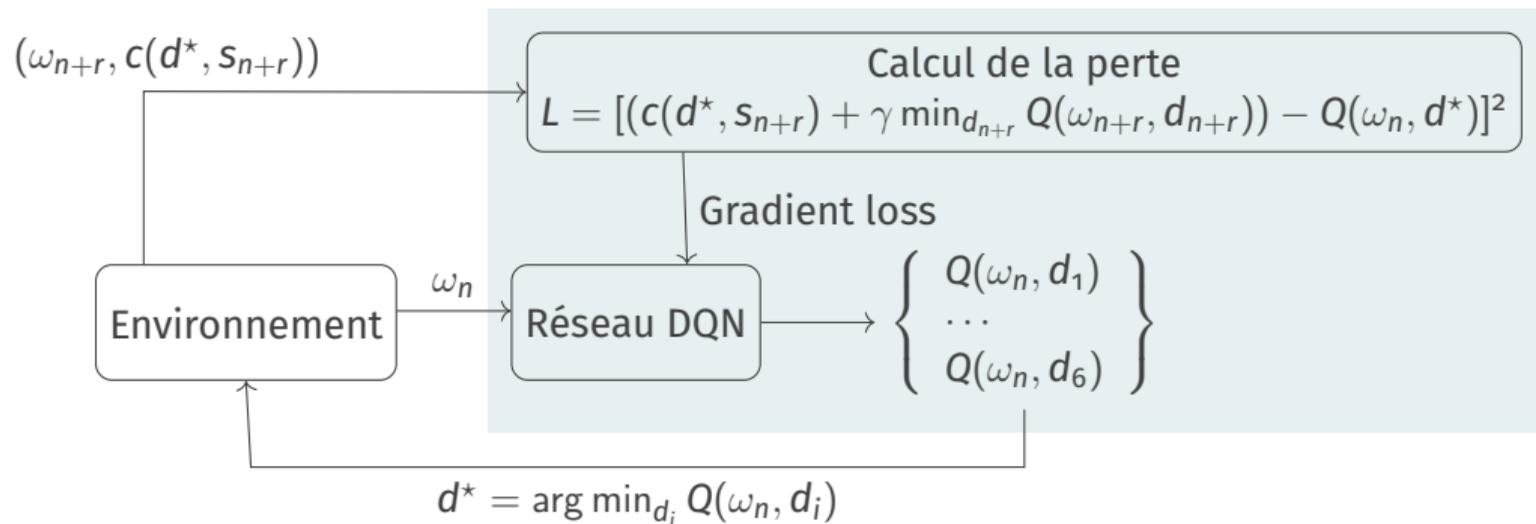
La politique optimale est obtenue à partir des expériences  $\langle \omega, d, \omega', c \rangle$

$$\underbrace{Q^\pi(h, d)}_{\text{Critère à optimiser}} = \underbrace{\mathbb{E}\left[\sum_{n=0}^{H-1} c(D_n, S_n) \mid h, d, \pi\right]}_{\text{Valeur d'une action dans un état suivant la politique } \pi}$$

$$\underbrace{Q^*(h, d)}_{\text{Q fonction}} = \min_{\pi \in \Pi} Q^\pi(h, d)$$

$$\underbrace{\pi^*}_{\text{Politique Optimale}} = \arg \min_{d \in \mathcal{D}} Q^*(h, d)$$

# Algorithme DQN<sup>5</sup>



<sup>5</sup>Deep Q-Network

# Résultats

Politique	Coût moyen (log)	IC	Taux de survie	Rechutes
<b>OH</b>	5.05	[3.58, 5.62]	91.88%	2.12
<b>Inactive</b>	6.39	[5.98, 6.67]	0.08%	1.00
<b>Threshold</b>	16.03	[14.83, 16.56]	74.04%	1.01
<b>DQN observé</b>	6.40	[5.99, 6.69]	0.84%	1.01
<b>DQN</b>	5.64	[4.49, 6.16]	76.76%	0.99
<b>R2D2<sup>6</sup></b>	6.49	[6.19, 6.71]	0.08%	1.00

TABLE: Policy evaluation performance on  $10^5$  simulations

---

<sup>6</sup>R2D2  $\approx$  DQN + LSTM

# Conclusion

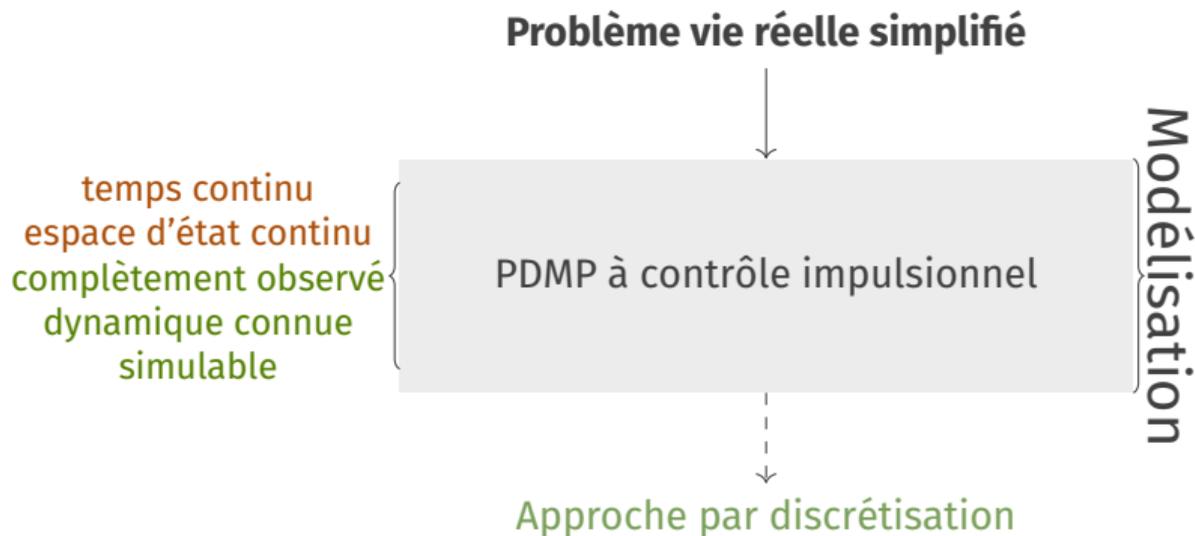
Politique	Coût moyen (log)	IC	Taux de survie	Rechutes
<b>OH</b>	5.05	[3.58, 5.62]	91.88%	2.12
<b>Inactive</b>	6.39	[5.98, 6.67]	0.08%	1.00
<b>Threshold</b>	16.03	[14.83, 16.56]	74.04%	1.01
<b>DQN observé</b>	6.40	[5.99, 6.69]	0.84%	1.01
<b>DQN</b>	5.64	[4.49, 6.16]	76.76%	0.99
<b>R2D2</b>	6.49	[6.19, 6.71]	0.08%	1.00

- Politiques peu conforme à la réalité
- Paramétrisation de la fonction de coût
- Work in Progress

**Nécessite beaucoup de données pour apprendre la politique optimale !**

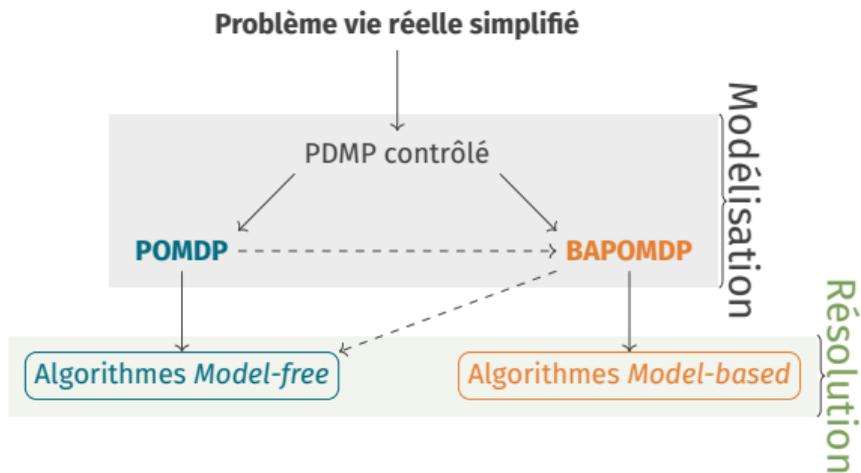
# Perspectives

EVALUER LES POLITIQUES OBTENUES AVEC LES DIFFÉRENTS ALGORITHMES

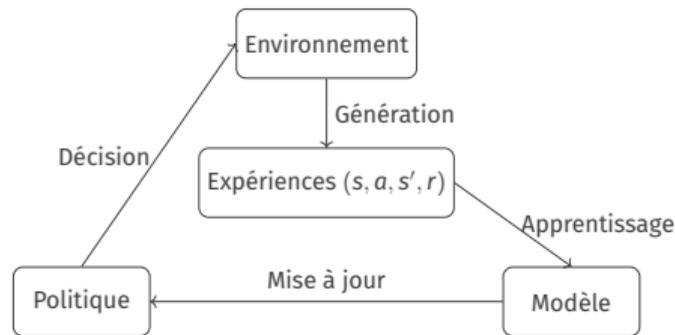


# Perspectives

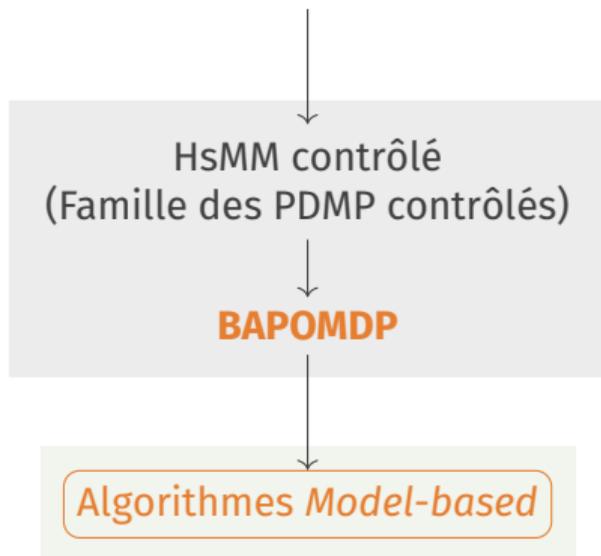
EXPLOITER LA CONNAISSANCE MODÈLE



## Apprentissage par renforcement *model-based*



## Problème vie réelle simplifié

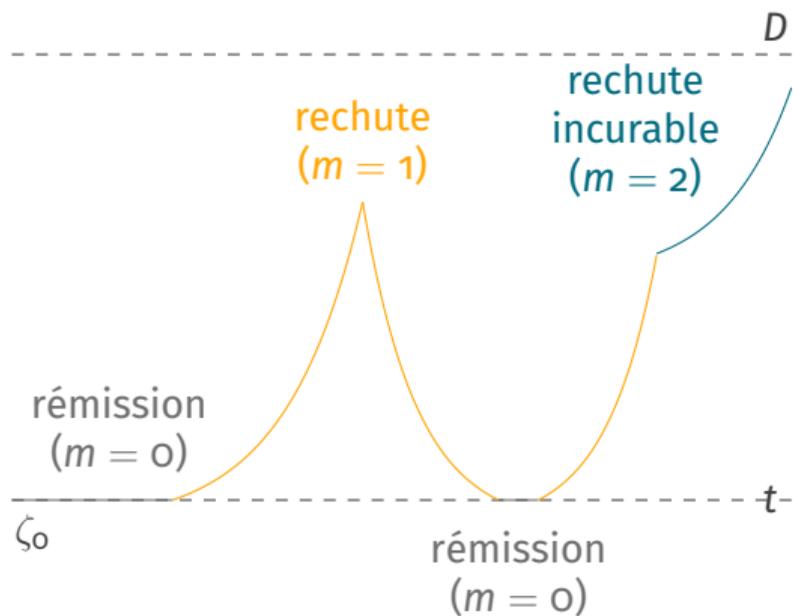


- Politique d'**échantillonnage adaptative** pour **reconstruire un processus stochastique**
- Application à la migration des oiseaux

---

<sup>7</sup>Hidden semi-Markov Model

- Travaux sur les familles de politiques  $\epsilon$ -optimales (publication contrôle)
- Transformation et résolution du PDMP contrôlé en BAPOMDP (acte de conférence ML)
- Extension de la méthode aux HsMM contrôlés - Application sur la migration des oiseaux (publication d'écologie)
- Rédaction et soutenance de thèse

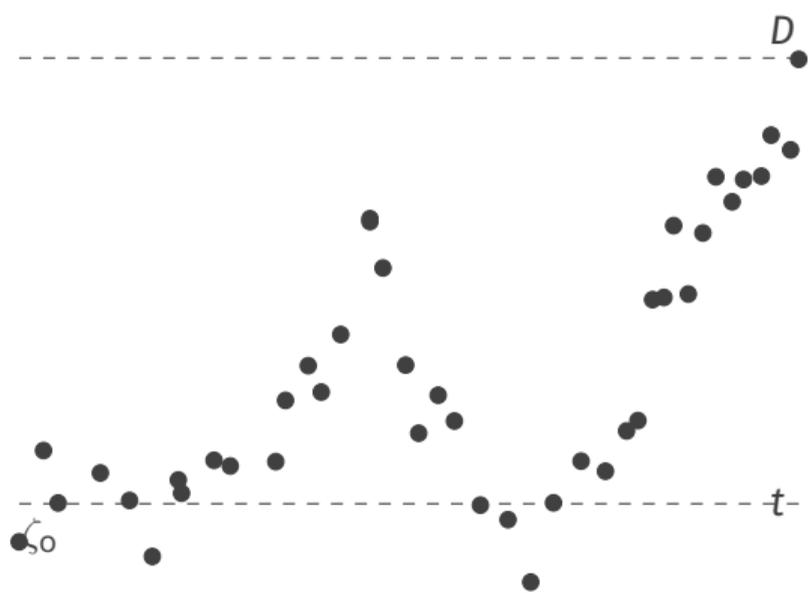


Soit  $s = (m, k, \zeta, u, t, \tau)$  l'état du patient:

- $m$  état général du patient;
- $k$  nombre de rechute;
- $\zeta$  biomarqueur;
- $u$  temps depuis le dernier saut;
- $t$  temps écoulé depuis le début du suivi;
- $\tau$  temps depuis l'application d'un traitement.

Soit  $d$  la **décision** telle que:  $d = (\ell, r)$ :

- $\ell$  traitement (*rien, chimiothérapie*);
- $r$  temps avant la prochaine visite (15, 30, 60 jours).



Soit  $s = (m, k, \zeta, u, t, \tau)$  l'état du patient:

- $m$  état général du patient;
- $k$  nombre de rechute;
- $\zeta$  biomarqueur;
- $u$  temps depuis le dernier saut;
- $t$  temps écoulé depuis le début du suivi;
- $\tau$  temps depuis l'application d'un traitement.

Soit  $d$  la **décision** telle que:  $d = (\ell, r)$ :

- $\ell$  traitement (*rien, chimiothérapie*);
- $r$  temps avant la prochaine visite (15, 30, 60 jours).

$$\underbrace{C(d, s')}_{\text{Fonction de coût}} = \underbrace{C_V}_{\text{coût de la visite}} + \underbrace{C_D(H - t') \times \mathbb{1}_{m'=3}}_{\text{coût de la mort}} + \underbrace{\kappa_C \times r \times \mathbb{1}_{\ell=a}}_{\text{coût de la chimiothérapie}}$$

avec  $C_V = 1$ ,  $C_D = 0.08$  et  $\kappa_C = 0.4$