

# Deep Reinforcement Learning for Bayes-Adaptive Impulse Control of PDMPs

Orlane Rossini <sup>1</sup>, Alice Cleynen <sup>1,2</sup>, Benoîte de Saporta <sup>1</sup>,  
Régis Sabbadin <sup>3</sup> and 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

November 2025



UNIVERSITÉ DE  
MONTPELLIER

INRAE

IMAG  
INSTITUT MONTPELLIERAIN  
ALEXANDER GROTHENDIECK



anr<sup>®</sup>

# Medical context

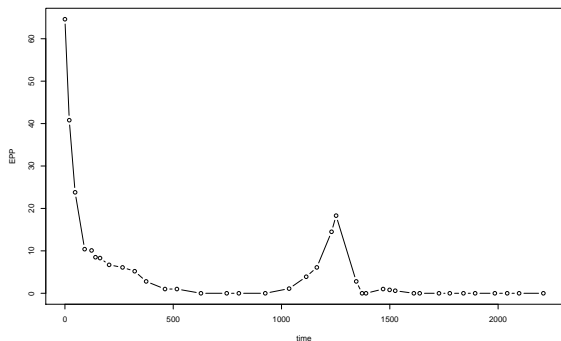


FIGURE: Example of patient data<sup>a</sup>

- Patients who have had **cancer** benefit from **regular follow-up**;
- The concentration of clonal immunoglobulin is measured over time;
- The doctor has to make new **decisions** at each visit.

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

# Medical context

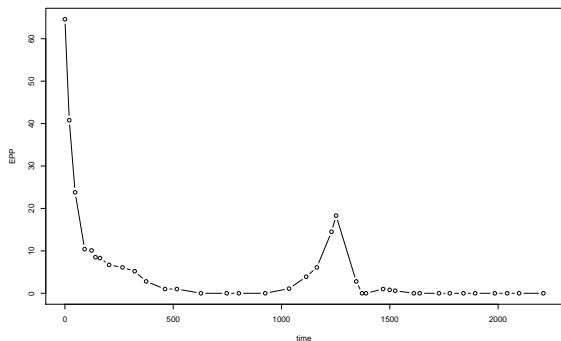


FIGURE: Example of patient data<sup>a</sup>

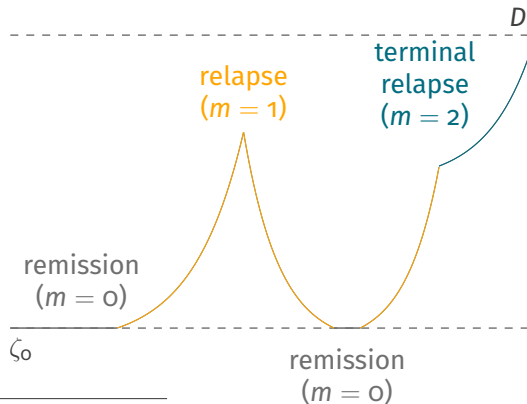
- Patients who have had **cancer** benefit from **regular follow-up**;
- The concentration of clonal immunoglobulin is measured over time;
- The doctor has to make new **decisions** at each visit.

⇒ **Optimising decision-making to ensure the patient's quality of life**

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

# Controlled PDMP<sup>1</sup>

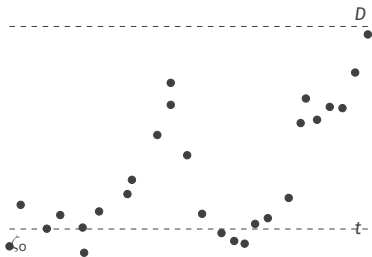
We switch **randomly** from one **deterministic** regime to another.



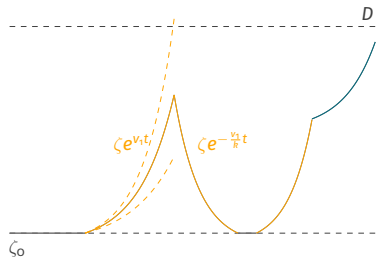
<sup>1</sup>Piecewise Deterministic Markov Processes

# Difficulties

## Partial observation

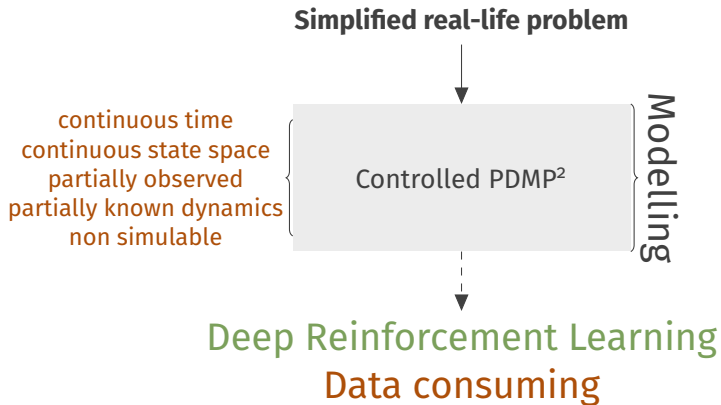


## Partially known dynamics



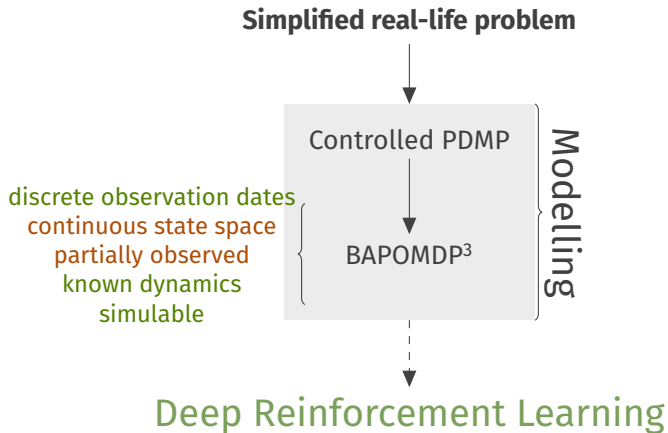
Hypothesis:  $v_1 \sim \text{Log-Normal}(\mu, \sigma^{-2})$ , with  $\mu$  and  $\sigma$  unknown.

# Methods



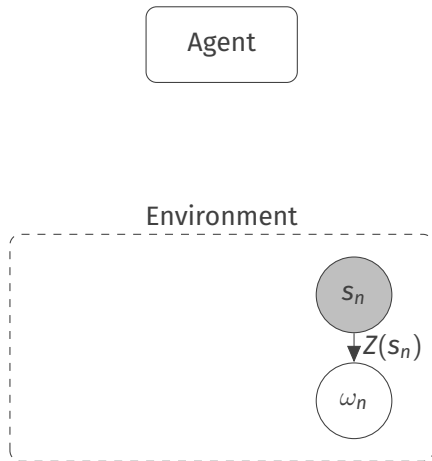
<sup>2</sup>Piecewise Deterministic Markov Processes

# Methods



<sup>3</sup>Bayes-Adaptive Partially Observed Markov Decision Process

# Characteristics of a POMDP<sup>4</sup>



## POMDP DEFINITION

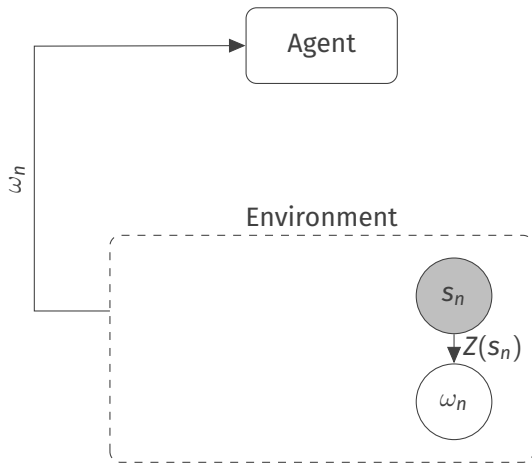
A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{A}, P, \Omega, Z, c)$ .

- Patient condition  $s = (m, k, \zeta, u) \in \mathcal{S}$ ;
- Actions  $a = (\ell, r) \in \mathcal{A}$ ;
- Transition function  $P(s'|s, a)$ ;
- **Observation**  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- **Observation function**  $Z(\omega|s)$ ;
- Cost function  $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>4</sup>Partially Observed Markov Decision Process



# Characteristics of a POMDP<sup>4</sup>



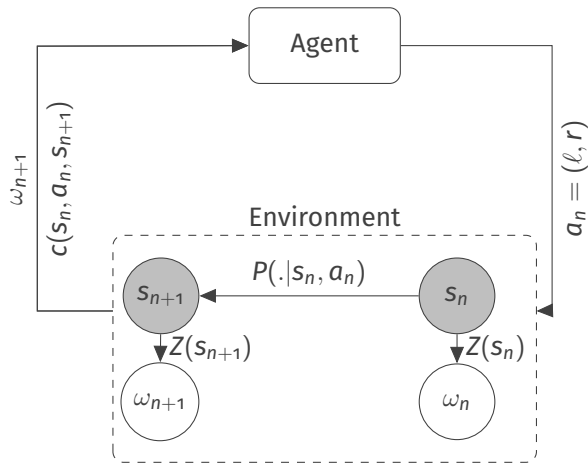
## POMDP DEFINITION

A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{A}, P, \Omega, Z, c)$ .

- Patient condition  $s = (m, k, \zeta, u) \in \mathcal{S}$ ;
- Actions  $a = (\ell, r) \in \mathcal{A}$ ;
- Transition function  $P(s'|s, a)$ ;
- **Observation**  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- **Observation function**  $Z(\omega|s)$ ;
- Cost function  $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>4</sup>Partially Observed Markov Decision Process

# Characteristics of a POMDP<sup>4</sup>



## POMDP DEFINITION

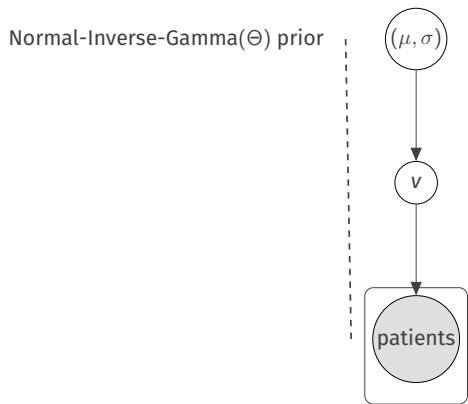
A POMDP is defined by a tuple  $(\mathcal{S}, \mathcal{A}, P, \Omega, Z, c)$ .

- Patient condition  $s = (m, k, \zeta, u) \in \mathcal{S}$ ;
- Actions  $a = (\ell, r) \in \mathcal{A}$ ;
- Transition function  $P(s' | s, a)$ ;
- **Observation**  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- **Observation function**  $Z(\omega | s)$ ;
- Cost function  $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

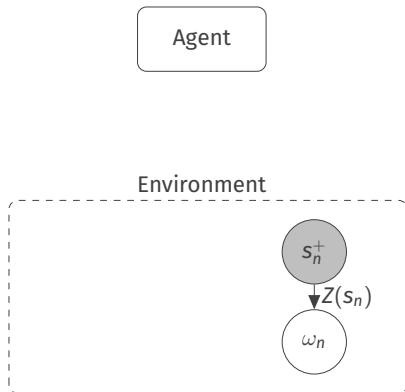
The transition function  $P(s' | s, a)$  is a combination of PDMP local characteristics.

<sup>4</sup>Partially Observed Markov Decision Process

# Handle uncertainty with Bayesian framework



# Characteristics of a BAPOMDP<sup>5</sup>



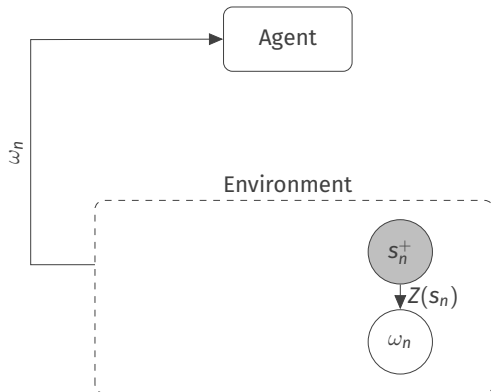
## BAPOMDP DEFINITION

Un BAPOMDP se définit par un tuple  $(\mathcal{S}^+, \mathcal{A}, P^+, \Omega, Z, c)$ .

- **Space of hyperstate**  $\mathcal{S}^+ = \mathcal{S} \times \Theta$ ;
- Actions  $a = (\ell, r) \in \mathcal{A}$ ;
- **Transition function**  $P^+(s', \theta' | s, a, \theta)$ ;
- Observation  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Observation function  $Z(\omega | s)$ ;
- Cost function  $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>5</sup>Bayes Adaptive Partially observed Markov decision process

# Characteristics of a BAPOMDP<sup>5</sup>



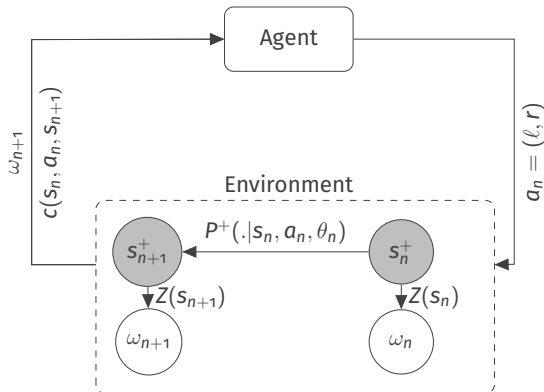
## BAPOMDP DEFINITION

Un BAPOMDP se définit par un tuple  $(\mathcal{S}^+, \mathcal{A}, P^+, \Omega, Z, c)$ .

- **Space of hyperstate**  $\mathcal{S}^+ = \mathcal{S} \times \Theta$ ;
- Actions  $a = (\ell, r) \in \mathcal{A}$ ;
- **Transition function**  $P^+(s', \theta' | s, a, \theta)$ ;
- Observation  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Observation function  $Z(\omega | s)$ ;
- Cost function  $c : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

<sup>5</sup>Bayes Adaptive Partially observed Markov decision process

# Characteristics of a BAPOMDP<sup>5</sup>



## BAPOMDP DEFINITION

Un BAPOMDP se définit par un tuple  $(\mathcal{S}^+, \mathbb{A}, P^+, \Omega, Z, c)$ .

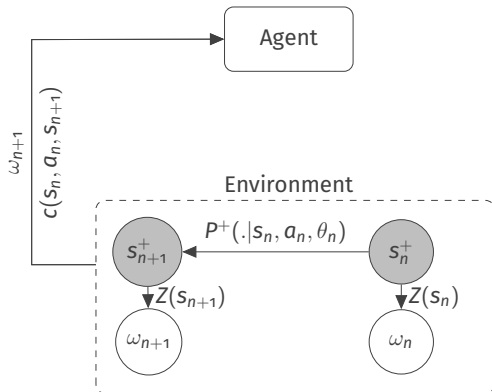
- **Space of hyperstate**  $\mathcal{S}^+ = \mathcal{S} \times \Theta$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- **Transition function**  $P^+(s', \theta' | s, a, \theta)$ ;
- Observation  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Observation function  $Z(\omega | s)$ ;
- Cost function  $c : \mathcal{S} \times \mathbb{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

$$P^+((s', \theta') \in B_E \times B_\Theta | (s, \theta), a)$$

$$= \int_{B_E} \mathbf{1}_{B_\Theta} \mathcal{U}(\theta, s, a, s') \times P(ds' | s, a, \theta).$$

<sup>5</sup>Bayes Adaptive Partially observed Markov decision process

# Characteristics of a BAPOMDP<sup>5</sup>



## BAPOMDP DEFINITION

Un BAPOMDP se définit par un tuple  $(\mathcal{S}^+, \mathbb{A}, P^+, \Omega, Z, c)$ .

- **Space of hyperstate**  $\mathcal{S}^+ = \mathcal{S} \times \Theta$ ;
- Actions  $a = (\ell, r) \in \mathbb{A}$ ;
- **Transition function**  $P^+(s', \theta' | s, a, \theta)$ ;
- Observation  $\omega = (k, F(\zeta, \epsilon), \mathbb{1}_{m=3}) \in \Omega$ ;
- Observation function  $Z(\omega | s)$ ;
- Cost function  $c : \mathcal{S} \times \mathbb{A} \times \mathcal{S} \rightarrow \mathbb{R}$ .

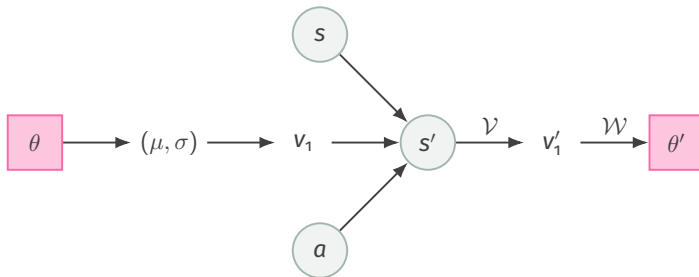
$$P^+((s', \theta') \in B_E \times B_\Theta | (s, \theta), a)$$

$$= \int_{B_E} \mathbf{1}_{B_\Theta} \mathcal{U}(\theta, s, a, s') \times P(ds' | s, a, \theta).$$

<sup>5</sup>Bayes Adaptive Partially observed Markov decision process

# Generate transition from prior

$$\mathcal{U}(\theta, s, a, s') = \mathcal{W}(\theta, \mathcal{V}(s, a, s')),$$





# Solving a BAPOMDP<sup>6</sup>

**Identify an optimal policy  $\pi^*$**

$$\underbrace{c(s, a, s')}_{\text{Cost function}} = \underbrace{C_V}_{\text{visit cost}} + \underbrace{C_D(H - t') \times \mathbb{1}_{m'=3}}_{\text{death cost}} + \underbrace{\kappa_C \times r \times \mathbb{1}_{\ell=a}}_{\text{treatment cost}}$$

---

<sup>6</sup>Bayes Adaptative Partially Observable Markov Decision Process

# Solving a BAPOMDP<sup>6</sup>

**Identify an optimal policy  $\pi^*$**

$$\underbrace{V(\pi, s)}_{\text{Optimization criterion}} = \underbrace{\mathbb{E}_s^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, A_n, S_n) \right]}_{\text{Expected total cost as a result of the policy } \pi}$$

---

<sup>6</sup>Bayes Adaptive Partially Observable Markov Decision Process

# Solving a BAPOMDP<sup>6</sup>

**Identify an optimal policy  $\pi^*$**

$$\underbrace{V(\pi, s)}_{\text{Optimization criterion}} = \underbrace{\mathbb{E}_s^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, A_n, S_n) \right]}_{\text{Expected total cost as a result of the policy } \pi}$$

$$\underbrace{V^*(s)}_{\text{Value function}} = \underbrace{\min_{\pi \in \Pi} V(\pi, s)}_{\text{Minimisation across policy space}}$$

---

<sup>6</sup>Bayes Adaptative Partially Observable Markov Decision Process

# Solving a BAPOMDP<sup>6</sup>

**Identify an optimal policy  $\pi^*$**

In reality, we do not observe state space!

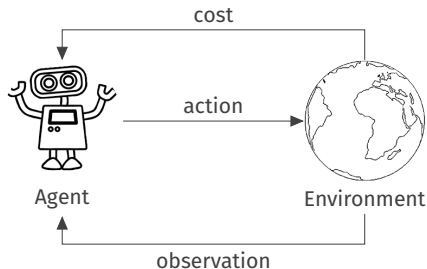
Let  $h_n = (\omega_0, a_0, \omega_1, a_1, \dots, \omega_n)$  be the history

$$\underbrace{V^*(h)}_{\text{Value function}} = \underbrace{\min_{\pi \in \Pi} V(\pi, h)}_{\text{Minimisation across policy space.}}$$

---

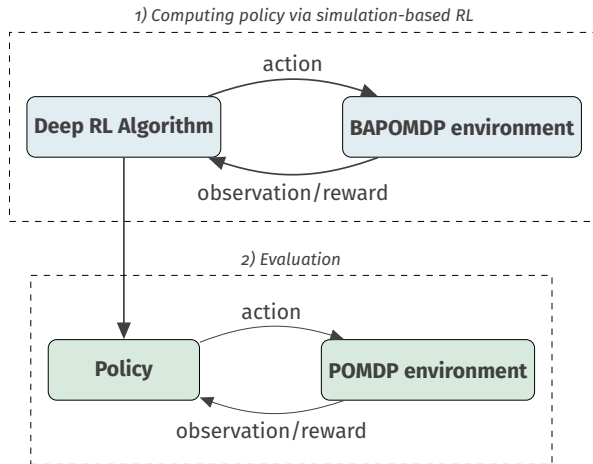
<sup>6</sup>Bayes Adaptative Partially Observable Markov Decision Process

# Reinforcement Learning



The optimal policy is obtained from the experiments  $\langle \omega, a, \omega', c \rangle$ , generated from  $P^+$  transition function

# Evaluate BAPOMDP framework



# Performance BAPOMDP

Policies	Mean Cost (log)	Survival rate	Treatment number	Visit number
<b>OH</b>	$13.45 \pm 0.01$	99.60%	0.87	120.57
<b>Random</b>	$12.79 \pm 0.01$	92.08%	4.36	67.21
<b>Inactive</b>	$8.21 \pm 0.07$	63.42%	0.00	35.39
<b>Threshold</b>	$7.03 \pm 0.04$	99.98%	5.07	64.67
<b>DQN<sup>7</sup></b>	$5.08 \pm 0.00$	99.70%	19.99	38.99
<b>PPO<sup>8</sup></b>	$5.94 \pm 0.01$	99.80%	19.99	58.99

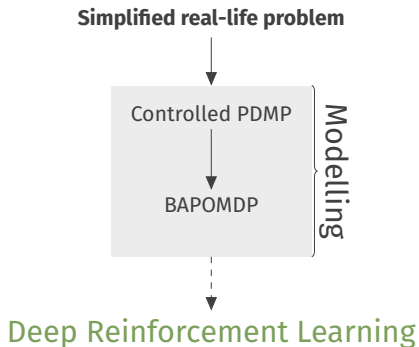
TABLE: Policy evaluation performance on 5000 Monte-Carlo simulations

---

<sup>7</sup>Deep Q-Network

<sup>8</sup>Proximal Policy Optimization

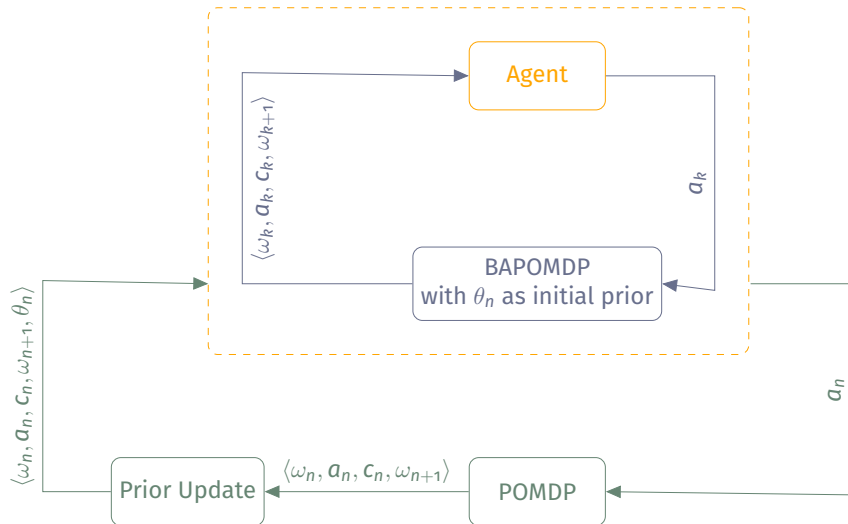
# Conclusion



- Bayes-adaptive method to address the PDMP control problem
- No explicit policies
- No estimates of unknown parameters



# Future work

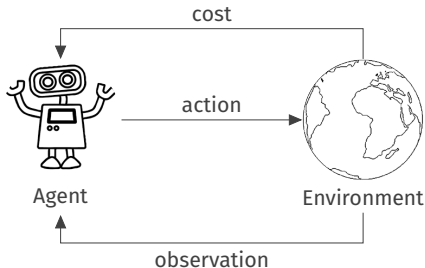


# Policy behavior indicators

TABLE: Summary of policy behavior indicators based on 5 000 Monte-Carlo simulations.

Indicator	PPO with AM	DQN with AM
Survival rates	99.80% $\pm$ 0.00	99.70% $\pm$ 0.00
Average number of treatment	19.99 $\pm$ 0.00	19.99 $\pm$ 0.01
Average time spend under treatment	1199.63 $\pm$ 00.04	1199.56 $\pm$ 0.05
Average number of visit	58.99 $\pm$ 0.01	38.99 $\pm$ 0.01
Average delay between two visits	40.00 $\pm$ 0.00	60.00 $\pm$ 0.00
Rate of visits occurring within 15 days	0.01 $\pm$ 0.00	0.00 $\pm$ 0.00
Rate of visits occurring within 30 days	66.66 $\pm$ 0.17	0.00 $\pm$ 0.00
Rate of visits occurring within 60 days	33.33 $\pm$ 0.17	100 $\pm$ 0.00

# Reinforcement Learning



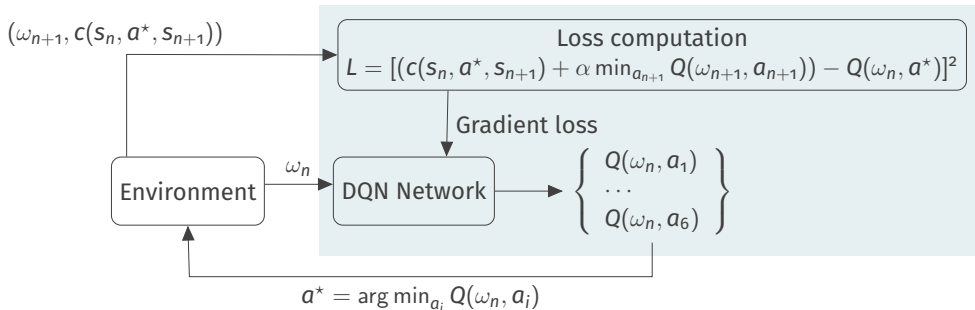
The optimal policy is obtained from the experiments  $\langle \omega, a, \omega', c \rangle$ , generated from  $P^+$  transition function

$$\underbrace{Q^\pi(s, a)}_{\text{Q value}} = \underbrace{\mathbb{E}^\pi \left[ \sum_{n=0}^{H-1} c(S_{n-1}, A_n, S_n) \mid s, a = (\ell, r) \right]}_{\text{Value of an action in a state according to the policy } \pi}$$

$$\underbrace{Q^*(s, a)}_{\text{Q function}} = \min_{\pi \in \Pi} Q^\pi(s, a)$$

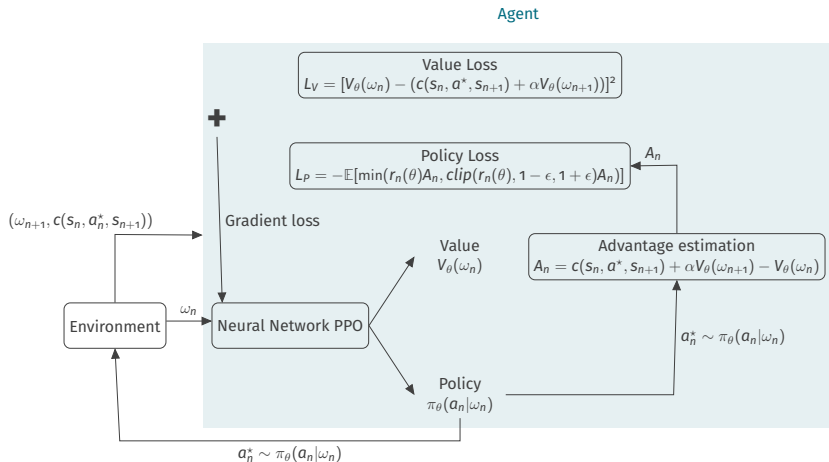
$$\underbrace{A(s, a)}_{\text{Advantage function}} = \underbrace{Q(s, a) - V(s)}_{\text{Extra cost obtained by the agent by taking the action}}$$

# Algorithm example: DQN<sup>9</sup>



<sup>9</sup>Deep Q-Network

# Algorithm example: PPO<sup>10</sup>



<sup>10</sup>Proximal policy optimization