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

Orlane Rossini [1], Alice Cleynen [1,2], Benoîte de Saporta [1],
Régis Sabbadin [3] and Meritxell Vinyals [3]

[1]IMAG, Univ Montpellier, CNRS, Montpellier, France
[2]John Curtin School of Medical Research, The Australian National University,
Canberra, ACT, Australia
[3]Univ Toulouse, INRAE-MIAT, Toulouse, France

October 2025

# Medical context



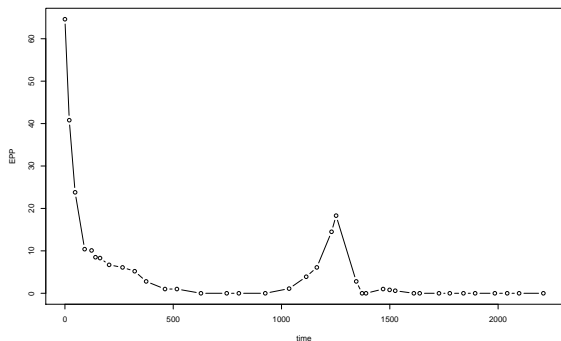FIGURE: Example of patient data[a]

- 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.

---

[a]IUCT Oncopole and CRCT, Toulouse, France

# Medical context



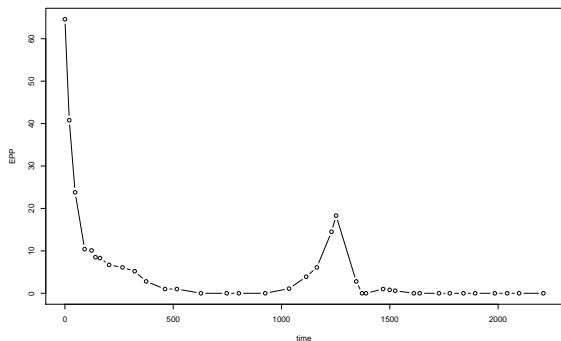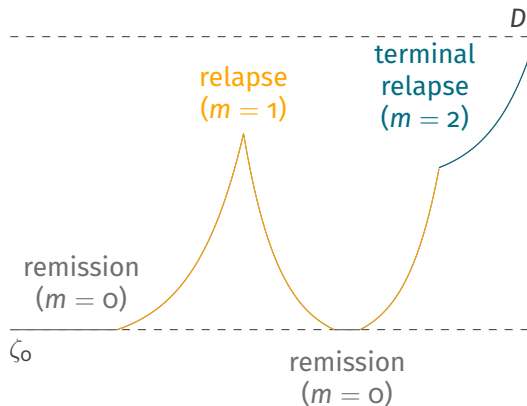FIGURE: Example of patient data[a]

- 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.

$\implies$ **Optimising decision-making to ensure the patient's quality of life**

---

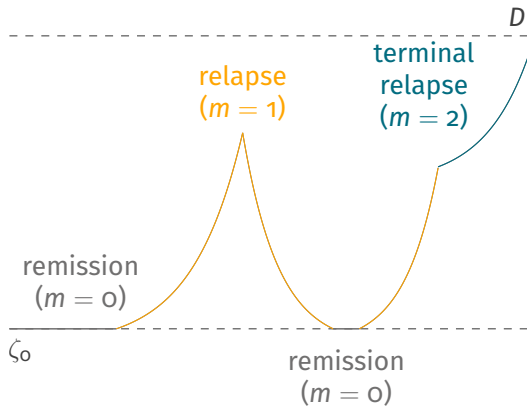[a]IUCT Oncopole and CRCT, Toulouse, France

We switch randomly from one deterministic regime to another.

# Controlled PDMP[1]

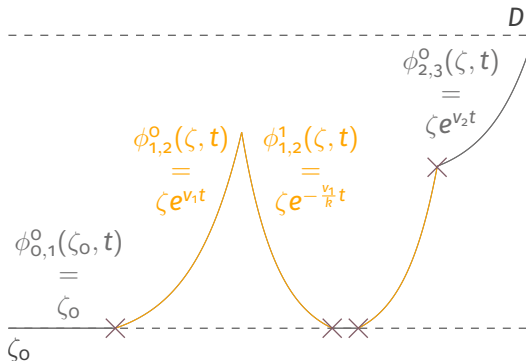We switch randomly from one deterministic regime to another.



Let $x = (m, \ell, k, \zeta, u)$ the patient's condition:
- $m$ the patient's condition;
- $\ell$ the current treatment;
- $k$ the number of treatments;
- $\zeta$ the biomarker;
- $u$ the time since the last jump.

[1]Piecewise Deterministic Markov Processes

A PDMP is defined by three local characteristics.



$$\Phi^\ell(x, t) = (m, k, \ell, \phi^\ell_{m,k}(\zeta, t), u + t)$$

FLOW

Description of the deterministic part of the process.

²Piecewise Deterministic Markov Processes

# Local Characteristics of a PDMP[2]

A PDMP is defined by three local characteristics.



$$\phi_{2,3}^0(\zeta, t) = \zeta e^{v_2 t}$$

$$\phi_{1,2}^0(\zeta, t) = \zeta e^{v_1 t}$$

$$\phi_{1,2}^1(\zeta, t) = \zeta e^{-\frac{v_1}{k} t}$$

$$\phi_{0,1}^0(\zeta_0, t) = \zeta_0$$

$D$

$\zeta_0$

### JUMP INTENSITY

Description of the process jump mechanisms.

- Boundary jump (deterministic)

$$t^\star(x) = t_{m,k}^{\ell\star}(\zeta) = \inf\{t > 0 : \phi_{m,k}^\ell(\zeta, t) \in \{\zeta_0, D\}\}$$

- Random jump

$$\mathbb{P}(T > t) = e^{-\int_0^t \lambda_{m,k}^\ell(\Phi(x,s))\,\mathrm{d}s}$$

---

[2]Piecewise Deterministic Markov Processes

# Local Characteristics of a PDMP[2]
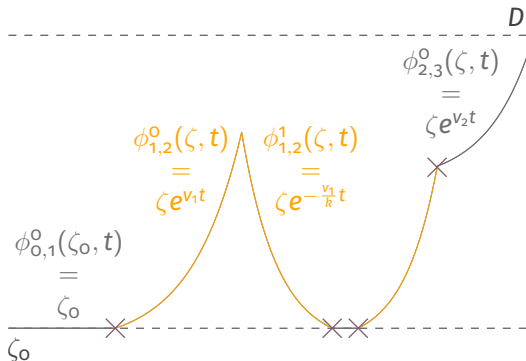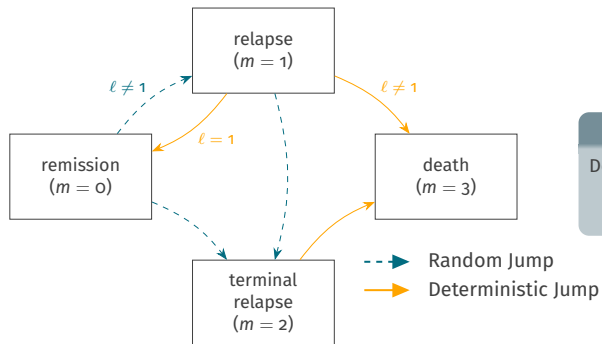
A PDMP is defined by three local characteristics.



```
          relapse
          (m = 1)

ℓ ≠ 1                        ℓ ≠ 1

                ℓ = 1
remission                        death
(m = 0)                          (m = 3)


          terminal
          relapse
          (m = 2)
```
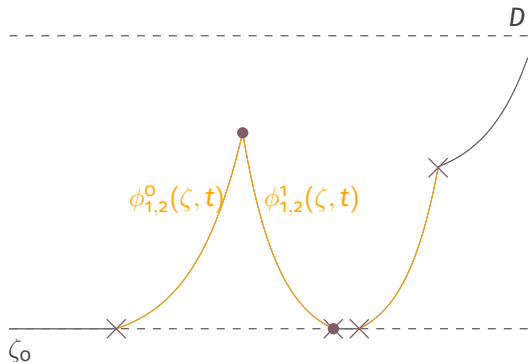
- - - → Random Jump
──→ Deterministic Jump

### MARKOV KERNEL

Description of the state of the process after each jump.

$$\mathbb{P}(X' \in A | X = x) = \int_A Q^d_{m,k}(\Phi^\ell(x, T), \mathrm{d}x')$$
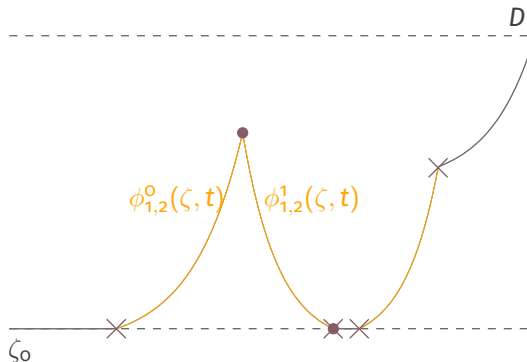
---

[2] Piecewise Deterministic Markov Processes

Choosing a new starting point :
- date for the next impulse;
- point from which to restart the process.

Choosing a new starting point :
- date for the next impulse;
- point from which to restart the process.

Restrictions:
- the delay between consecutive impulses is in a finite set;
- only change the current treatment $\ell$.

# Solving impulse control for PDMP

**Identify an $\epsilon$-optimal strategy** $\mathcal{S} = (\tau_n, \chi_n)_{n \geq 1}$

$$\underbrace{\mathcal{V}(\mathcal{S}, x)}_{\text{Expected cost of strategy } \mathcal{S}} = \mathbb{E}_x^{\mathcal{S}} \left[ \int_0^{+\infty} e^{-\gamma t} \underbrace{c_R(X_t)}_{\text{running cost}} dt + \sum_{n=1}^{\infty} \underbrace{c_I}_{\text{impulse cost}} \left( X_{\tau_n}, X_{\tau_n^+} \right) \right],$$
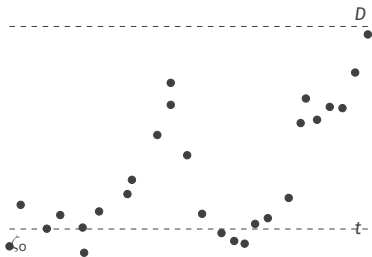
# Solving impulse control for PDMP

**Identify an $\epsilon$-optimal strategy** $\mathcal{S} = (\tau_n, \chi_n)_{n \geq 1}$

$$\underbrace{\mathcal{V}(\mathcal{S}, x)}_{\text{Expected cost of strategy } \mathcal{S}} = \mathbb{E}_x^{\mathcal{S}} \left[ \int_0^{+\infty} e^{-\gamma t} \underbrace{c_R(X_t)}_{\text{running cost}} dt + \sum_{n=1}^{\infty} \underbrace{c_I}_{\text{impulse cost}} \left( X_{\tau_n}, X_{\tau_n^+} \right) \right],$$
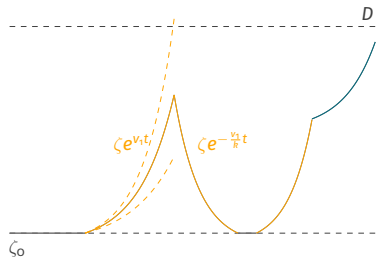
$$\mathcal{V}^{\star}(x) = \inf_{\mathcal{S} \in S} \mathcal{V}(\mathcal{S}, x)$$

**Partial observation**



**Partially known dynamics**



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

**Simplified real-life problem**

continuous time
continuous state space
partially observed
partially known dynamics
simulable

Controlled PDMP[3]

Modelling

Deep Reinforcement Learning
Data consuming

---
[3]Piecewise Deterministic Markov Processes

**Simplified real-life problem**

Controlled PDMP

Modelling

discrete observation dates
continuous state space
partially observed
known dynamics
simulable

POMDP[4]

Reinforcement Learning

---

[4]Partially Observed Markov Decision Process
[5]de Saporta B, Thierry d'Argenlieu A, Sabbadin R, Cleynen A (2024) A Monte-Carlo planning strategy for medical follow-up optimization: Illustration on multiple myeloma data. PLOS ONE 19(12): e0315661

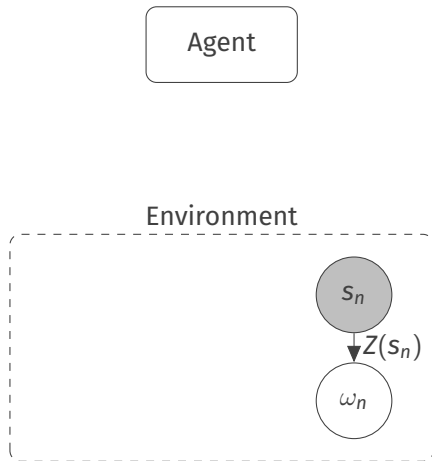**Simplified real-life problem**

Controlled PDMP

discrete observation dates
continuous state space
partially observed
known dynamics
simulable

BAPOMDP[6]

Modelling

Deep Reinforcement Learning

---

[6]Bayes-Adaptive Partially Observed Markov Decision Process

```
┌─────────────────────────────────────────────┐
│              BAPOMDP                          │
│   ┌─────────────────────────────────────┐    │
│   │             POMDP                    │    │
│   │   ┌─────────────────────────────┐   │    │
│   │   │  Markov Decision Process    │   │    │
│   │   └─────────────────────────────┘   │    │
│   └─────────────────────────────────────┘    │
└─────────────────────────────────────────────┘
```
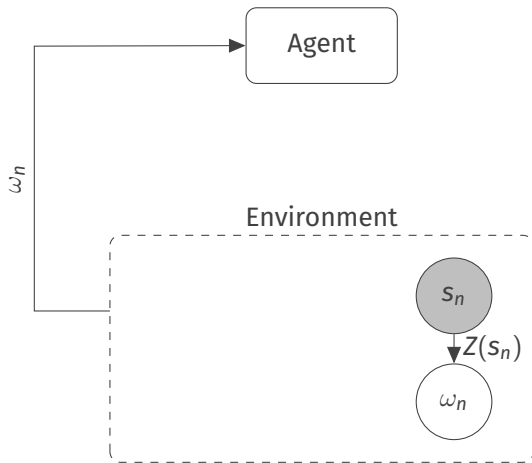
---

[7]Markov Decision Process

Agent

Environment



### POMDP Definition

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

- Patient condition $s = (m, k, \zeta, u) \in \mathbb{S}$;
- Actions $a = (\ell, r) \in \mathbb{A}$;
- 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 $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$.

[8]Partially Observed Markov Decision Process

Agent
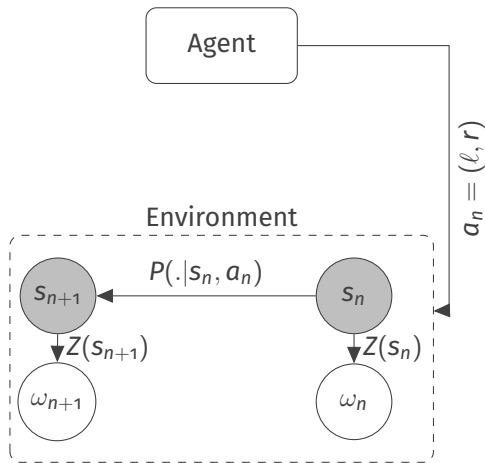
$\omega_n$

Environment

$s_n$

$Z(s_n)$

$\omega_n$

### POMDP Definition

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

- Patient condition $s = (m, k, \zeta, u) \in \mathbb{S}$;
- Actions $a = (\ell, r) \in \mathbb{A}$;
- 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 $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$.

[8] Partially Observed Markov Decision Process
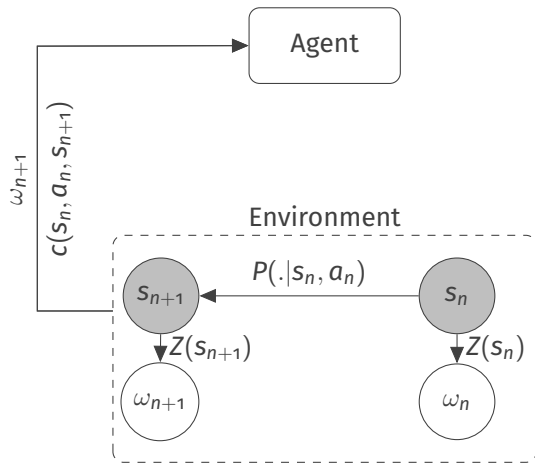
# Characteristics of a POMDP[8]



## POMDP DEFINITION

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

- Patient condition $s = (m, k, \zeta, u) \in \mathbb{S}$;
- Actions $a = (\ell, r) \in \mathbb{A}$;
- 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 $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$.

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

---

[8]Partially Observed Markov Decision Process

Agent

$c(s_n, a_n, s_{n+1})$

$\omega_{n+1}$

Environment

$P(.|s_n, a_n)$

$s_{n+1}$     $s_n$

$Z(s_{n+1})$    $Z(s_n)$
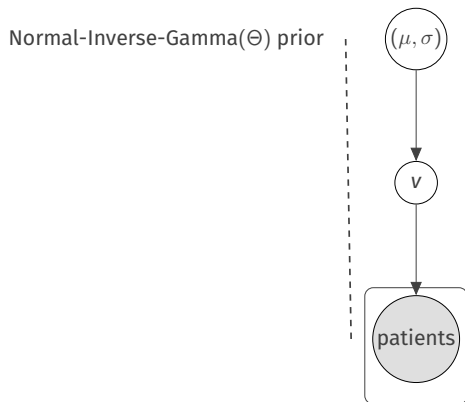
$\omega_{n+1}$     $\omega_n$

---

### POMDP DEFINITION

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

- Patient condition $s = (m, k, \zeta, u) \in \mathbb{S}$;
- Actions $a = (\ell, r) \in \mathbb{A}$;
- 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 $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$.

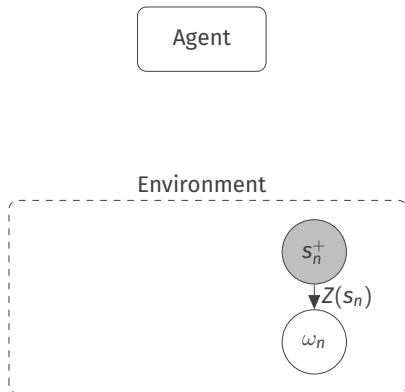The transition function $P(s'|s, a)$ is a combination of PDMP local characteristics.

---

[8] Partially Observed Markov Decision Process

Normal-Inverse-Gamma($\Theta$) prior

$(\mu, \sigma)$

$v$

patients

Agent

Environment



### BAPOMDP Definition

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

- **Space of hyperstate** $\mathbb{S}^+ = \mathbb{S} \times \Theta$;
- 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 $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \rightarrow \mathbb{R}$.

---

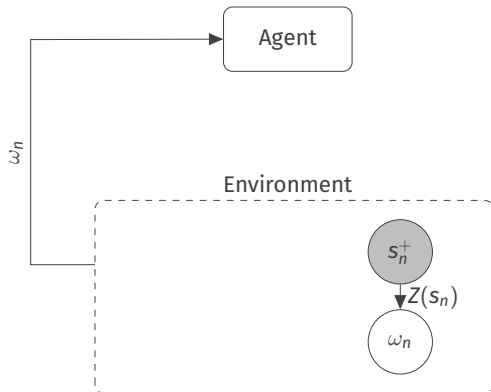[9]Bayes Adaptive Partially observed Markov decision process

### BAPOMDP Definition

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

- **Space of hyperstate** $\mathbb{S}^+ = \mathbb{S} \times \Theta$;
- 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 $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$.

---

[9]Bayes Adaptive Partially observed Markov decision process

# Characteristics of a BAPOMDP[9]



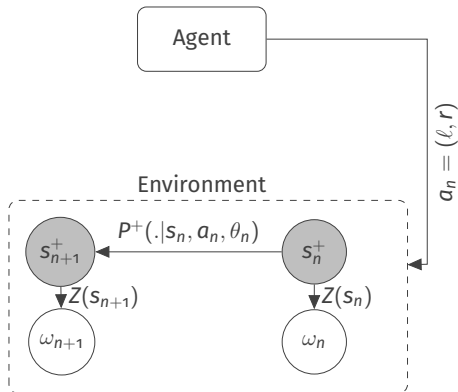### BAPOMDP Definition

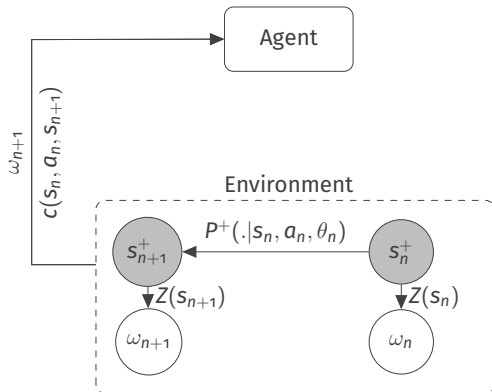Un BAPOMDP se définit par un tuple $(\mathbb{S}^+, \mathbb{A}, P^+, \Omega, Z, c)$.
- **Space of hyperstate** $\mathbb{S}^+ = \mathbb{S} \times \Theta$;
- Actions $a = (\ell, r) \in \mathbb{A}$;
- **Transition function** $P^+(s', \theta'|s, a, \theta)$;
- Observation $\omega = (k, F(\zeta, \epsilon),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

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

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

---

Agent

Environment

$\omega_{n+1}$

$c(s_n, a_n, s_{n+1})$

$P^+(.|s_n, a_n, \theta_n)$

$s_{n+1}^+$     $s_n^+$

$Z(s_{n+1})$     $Z(s_n)$

$\omega_{n+1}$     $\omega_n$

---

### BAPOMDP DEFINITION

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

- **Space of hyperstate** $\mathbb{S}^+ = \mathbb{S} \times \Theta$;
- 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 $Z(\omega|s)$;
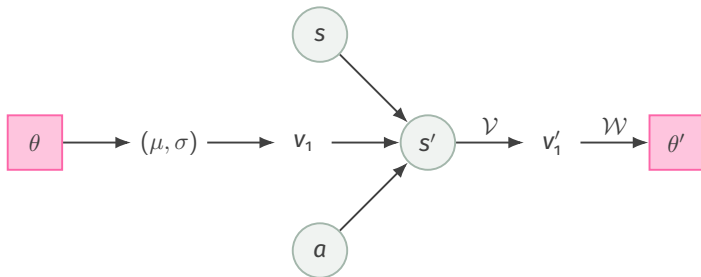- Cost function $c : \mathbb{S} \times \mathbb{A} \times \mathbb{S} \to \mathbb{R}$.

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

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

---

[9]Bayes Adaptive Partially observed Markov decision process

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

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

$$\underbrace{c(s, a, s')}_{\text{Cost function}} = \underbrace{C_V}_{\text{visit cost}}$$

$$+ \underbrace{C_D(H - t') \times \mathbb{1}_{m'=3}}_{\text{death cost}}$$

$$+ \underbrace{\kappa_C \times r \times \mathbb{1}_{\ell=a}}_{\text{treatment cost}}$$

---

[10] Bayes Adaptative Partially Observable Markov Decision Process

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

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

---

[10] Bayes Adaptative Partially Observable Markov Decision Process

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

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

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

---

[10] Bayes Adaptative Partially Observable Markov Decision Process
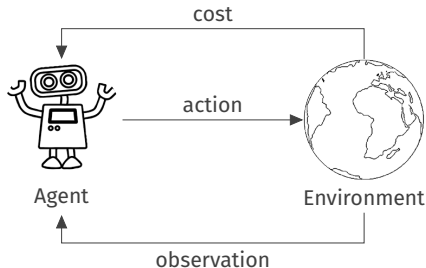
**Identify an optimal policy** $\pi^\star$

In reality, we do not observe state space!

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

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

---

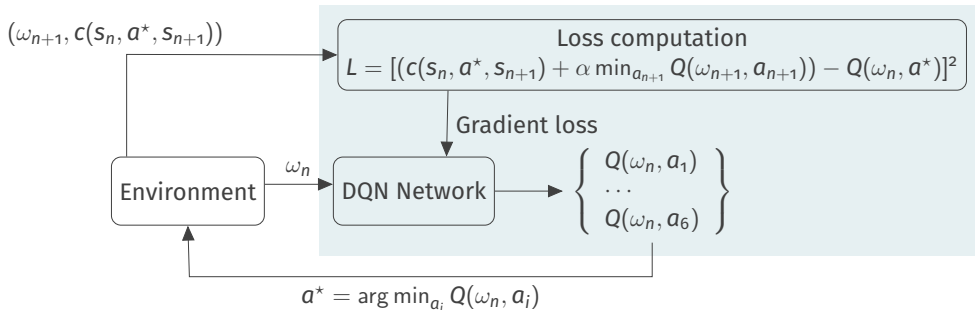[10] Bayes Adaptative Partially Observable Markov Decision Process

# Reinforcement Learning



Agent

Environment

cost

action

observation

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

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

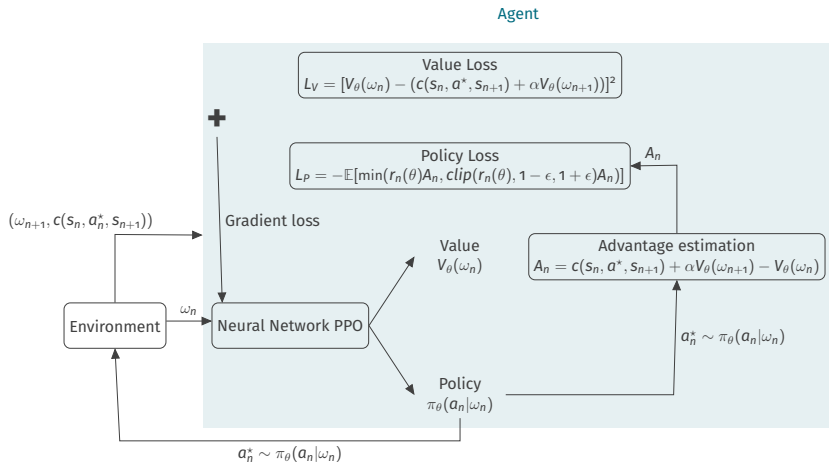$$\underbrace{Q^\star(s, a)}_{\text{Q function}} = \min_{\pi \in \Pi} Q^\pi(s, a)$$

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

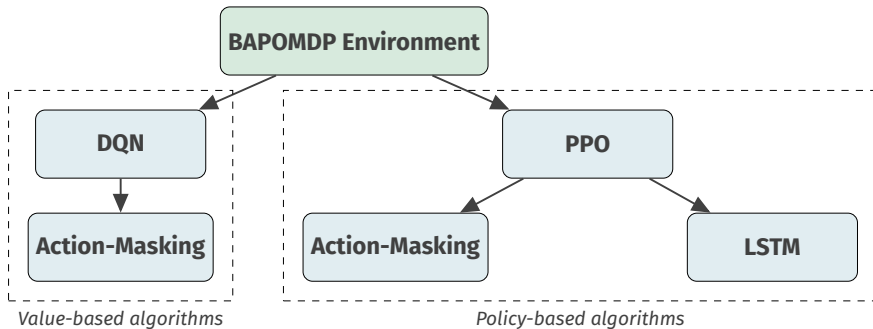The figure shows a block diagram. A box labeled "Environment" outputs $\omega_n$ to the "DQN Network" box. The DQN Network outputs the set:

$$\left\{ \begin{array}{l} Q(\omega_n, a_1) \\ \cdots \\ Q(\omega_n, a_6) \end{array} \right\}$$

Above, the "Loss computation" box contains:

$$L = [(c(s_n, a^\star, s_{n+1}) + \alpha \min_{a_{n+1}} Q(\omega_{n+1}, a_{n+1})) - Q(\omega_n, a^\star)]^2$$

with input $(\omega_{n+1}, c(s_n, a^\star, s_{n+1}))$ and a "Gradient loss" arrow pointing to the DQN Network.

The feedback path is $a^\star = \arg\min_{a_i} Q(\omega_n, a_i)$.

---

[11]Deep Q-Network

# Algorithm example: PPO[12]



Agent

Value Loss
$$L_V = [V_\theta(\omega_n) - (c(s_n, a^*, s_{n+1}) + \alpha V_\theta(\omega_{n+1}))]^2$$

Policy Loss
$$L_P = -\mathbb{E}[\min(r_n(\theta)A_n, clip(r_n(\theta), 1-\epsilon, 1+\epsilon)A_n)]$$

$A_n$

$(\omega_{n+1}, c(s_n, a_n^*, s_{n+1}))$

Gradient loss

Value
$V_\theta(\omega_n)$

Advantage estimation
$$A_n = c(s_n, a^*, s_{n+1}) + \alpha V_\theta(\omega_{n+1}) - V_\theta(\omega_n)$$

$\omega_n$

Environment → Neural Network PPO

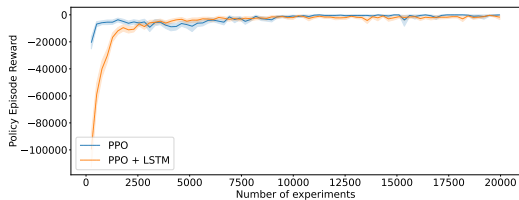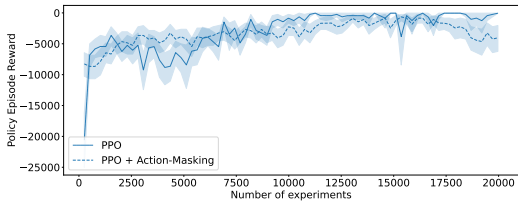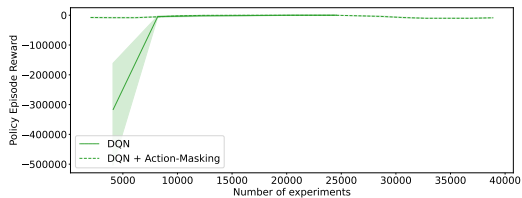$a_n^* \sim \pi_\theta(a_n|\omega_n)$

Policy
$\pi_\theta(a_n|\omega_n)$

$a_n^* \sim \pi_\theta(a_n|\omega_n)$

[12] Proximal policy optimization
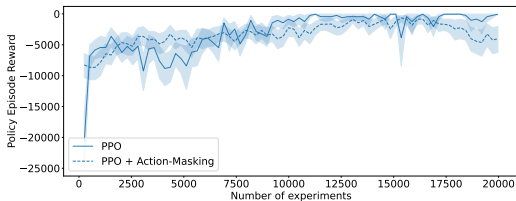
# Results

# Results



$\implies$ DQN with Action-Masking outperformed all baseline algorithms.

# Evaluate BAPOMDP framework



*1) Computing policy via simulation-based RL*

action

**Deep RL Algorithm** → **BAPOMDP environment**

observation/reward

*2) Evaluation*

action

**Policy** → **POMDP environment**

observation/reward

| Prior | $(\mu_0, \kappa_0, \alpha_0, \beta_0)$ |
|---|---|
| $\theta_{weak}$ | $(1, \ 0.001, \ 1.01, \ 1)$ |
| $\theta_{medium}$ | $(-6.785, \ 5.001, \ 3.51, \ 1)$ |
| $\theta_{high}$ | $(-6.23, \ 10, \ 6.01, \ 1)$ |

# Evaluate BAPOMDP framework

*What would happen if we modeled the unknown parameter $v_1$ as a fixed Normal–Gamma random variable in a POMDP?*

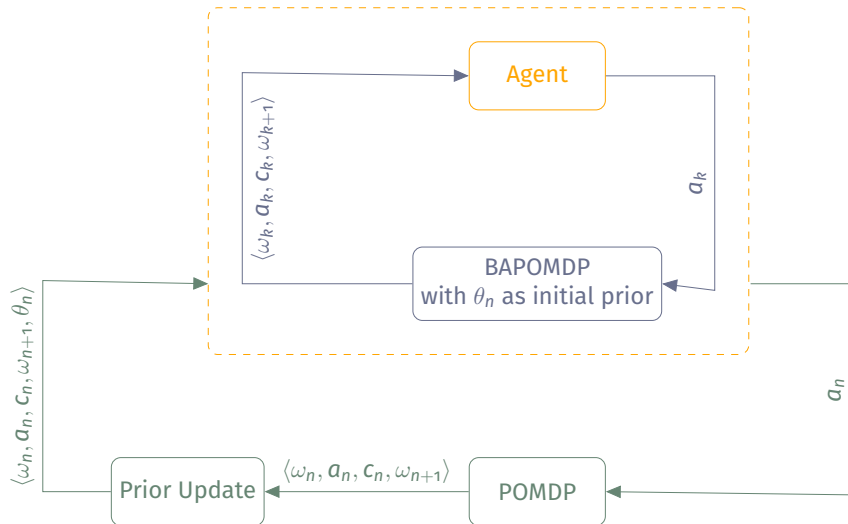| Prior | BAPOMDP | Non-adaptive POMDP |
|---|---|---|
| $\theta_{\mathsf{weak}}$ | $-5.74 \pm 0.00$ | $-5.70 \pm 0.00$ |
| $\theta_{\mathsf{medium}}$ | $-5.66 \pm 0.00$ | $-5.97 \pm 0.00$ |
| $\theta_{\mathsf{high}}$ | $-5.76 \pm 0.00$ | $-5.45 \pm 0.00$ |

# Conclusion

**Simplified real-life problem**

Controlled PDMP

BAPOMDP

Modelling

Deep Reinforcement Learning

- Bayes-adaptive method to address the PDMP control problem
- Comparable test-time performance to non-adaptive models
- No estimates of unknown parameters

# 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 | $33.33 \pm 0.17$ | $100 \pm 0.00$ |