

# Représentations prédictives pour la prise de décision dans l'incertain

Abdeslam Boularias

Mercredi, le 28 juillet 2010



## Contrôle d'un système dynamique

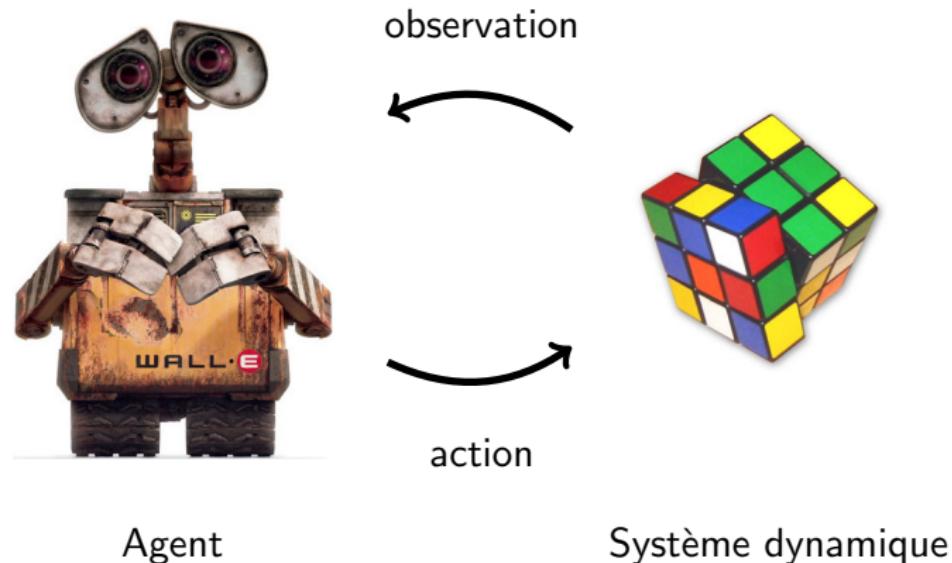


FIG.: L'interaction entre un agent et un système dynamique.

## Prise de décision séquentielle

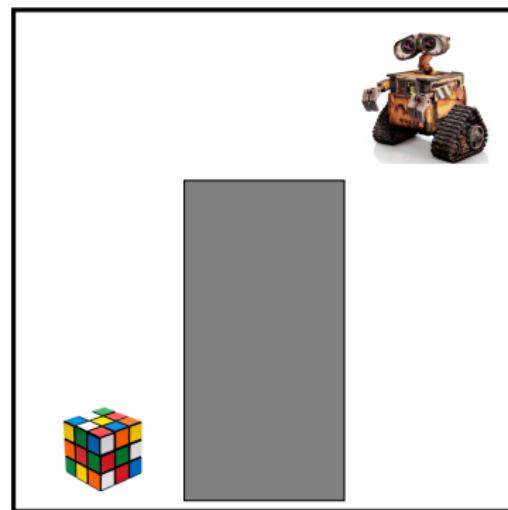


FIG.: Un simple problème de prise de décision séquentielle.

## Prise de décision séquentielle

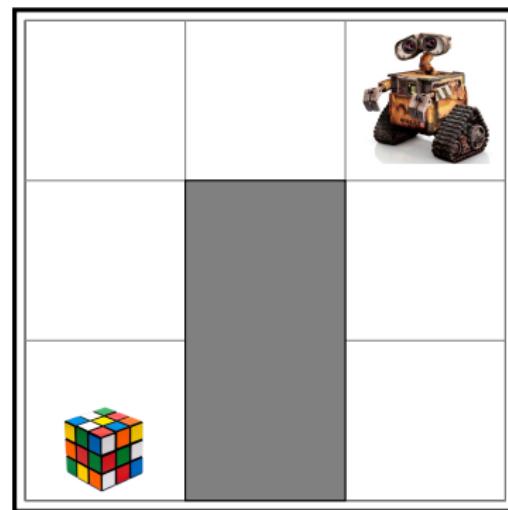


FIG.: Un simple problème de prise de décision séquentielle.

## Prise de décision séquentielle

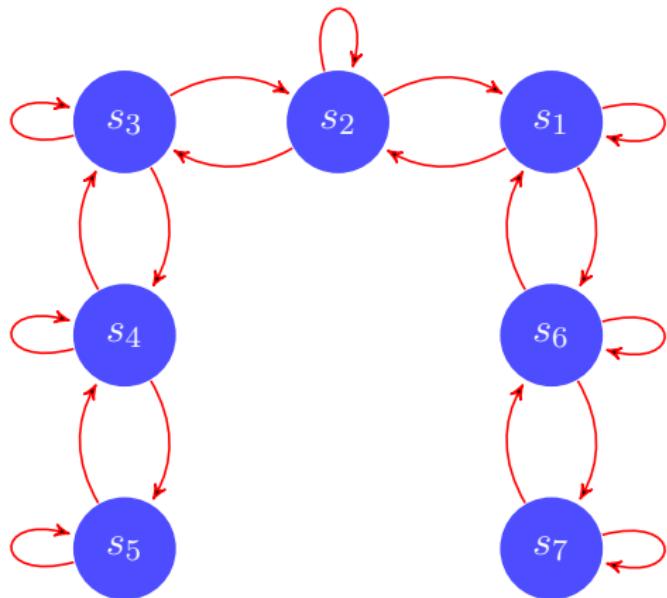
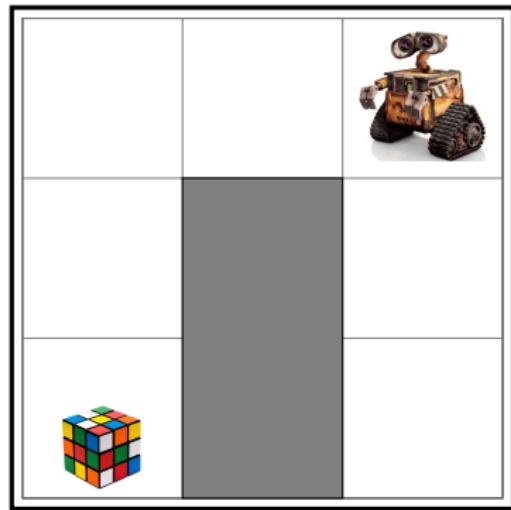


FIG.: Un processus décisionnel de Markov (MDP).

## Processus décisionnels de Markov

Un processus décisionnel de Markov (MDP) est défini par :

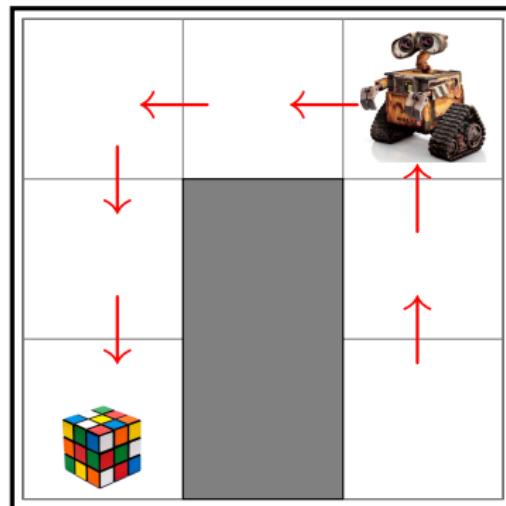
- $\mathcal{S}$  : un ensemble des états.
- $\mathcal{A}$  : un ensemble des actions.
- $T^a$  : une fonction de transition stochastique pour chaque action  $a \in \mathcal{A}$ .

$$T^a(s, s') = \Pr(s_{t+1} = s' | s_t = s, a_t = a)$$

- $R$  : une fonction de récompense (ou de coût).

## Politiques

- Une politique est une fonction qui associe une action à chaque état.



## Politiques

- La valeur (désirabilité) d'une politique  $\pi$  est donnée par :

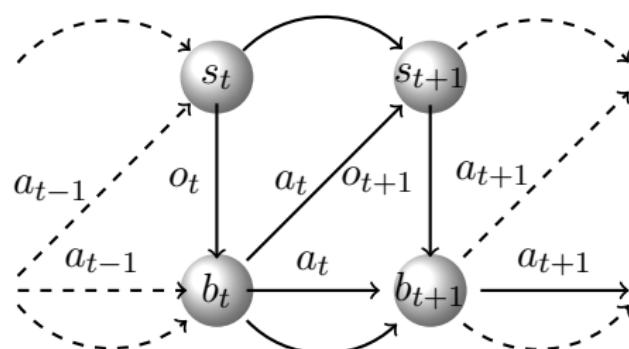
$$V(\pi) = \mathbb{E}_{s_t} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) | \pi, T \right]$$

- $\gamma \in [0, 1)$  est un facteur d'escompte (ou d'amortissement).
- Résoudre un MDP  $\Rightarrow$  Trouver une politique optimale.

## Processus décisionnels de Markov partiellement observables (POMDPs)

- Les observations peuvent être **partielles** et **incorrectes**.
- Un MDP partiellement observable (POMDP) est un MDP avec un ensemble d'observations  $\mathcal{O}$ , et une **fonction d'observation**  $Z$ .
- $Z(s, o)$  est la probabilité d'observer  $o$  dans l'état  $s$ .

États réels (cachés) :



États de croyance :

## Exemple d'un POMDP

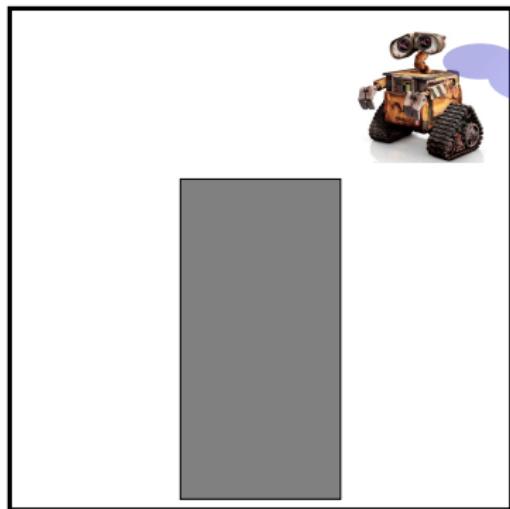


FIG.: État de croyance initial.

## Exemple d'un POMDP

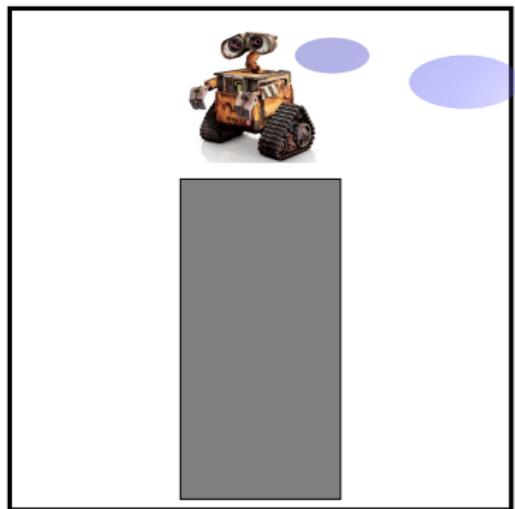


FIG.: État de croyance après avoir exécuté une action.

## Exemple d'un POMDP

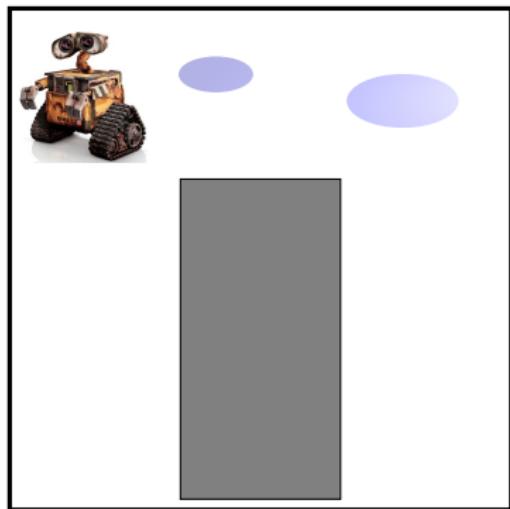
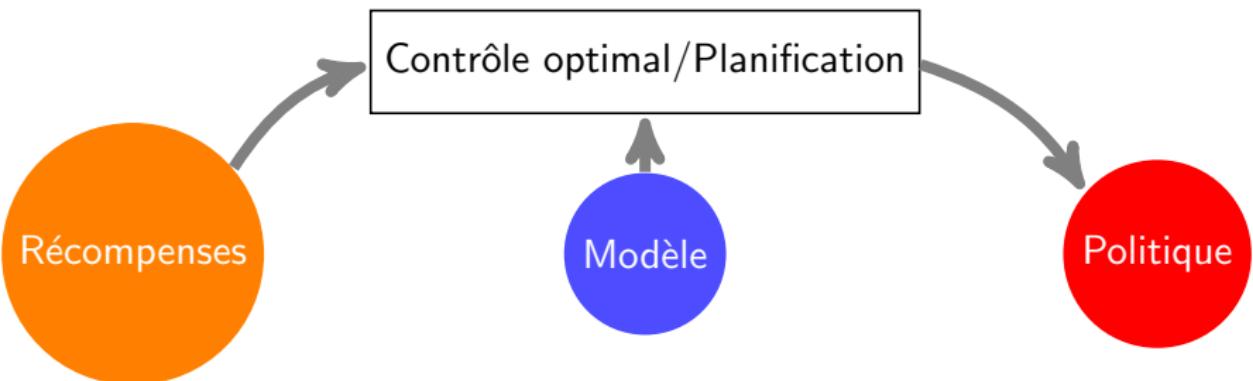


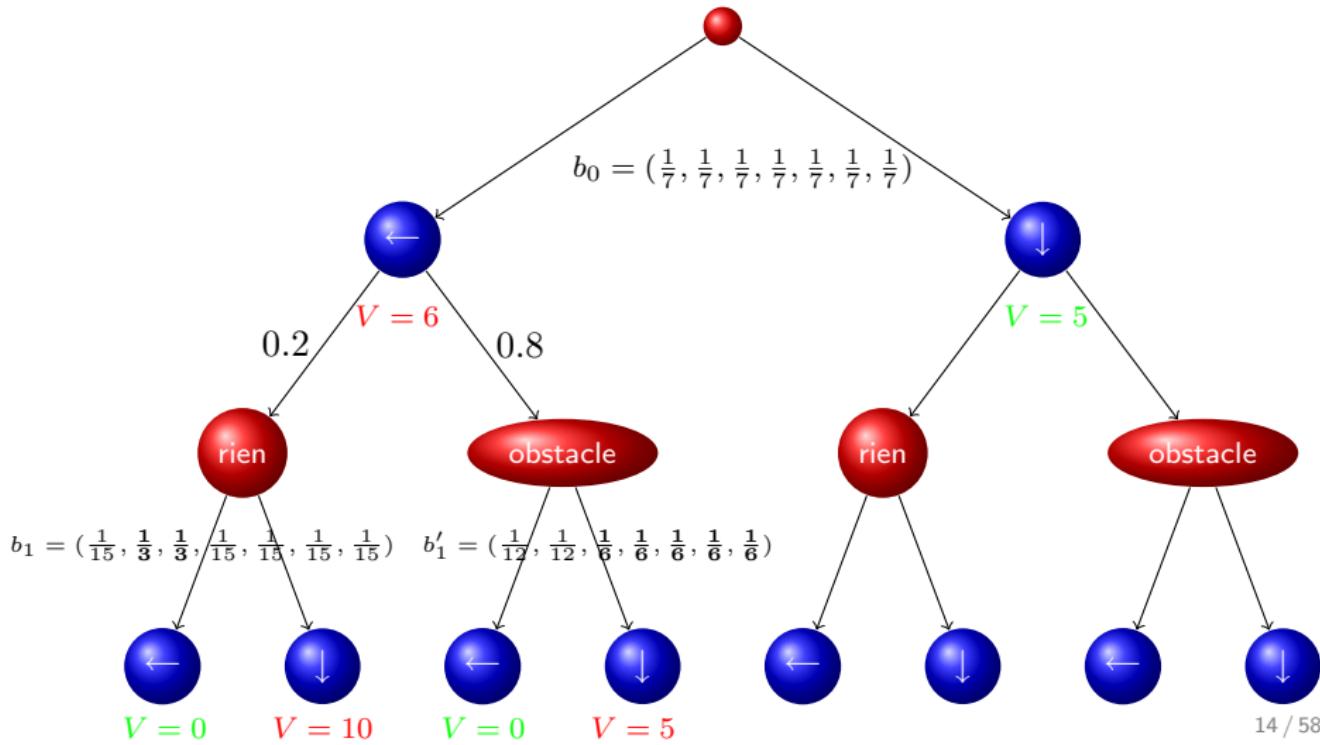
FIG.: État de croyance après avoir exécuté deux actions.

## Résoudre un (PO)MDP



## Planification avec les POMDPs

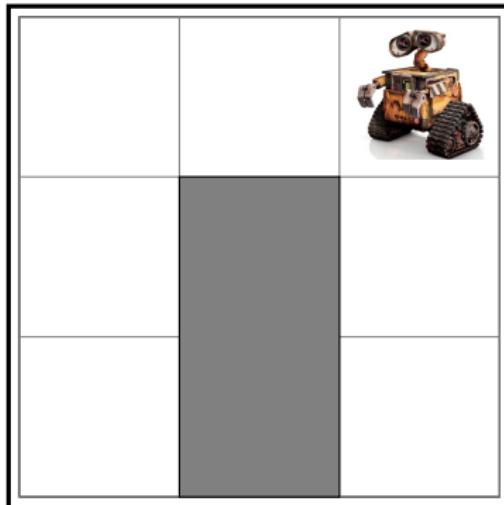
$$V^*(b) = \max_{a \in \mathcal{A}} \left[ R(s, a) + \gamma \sum_{o \in \mathcal{O}} Pr(o|b, a) V^*(b_{a,o}) \right]$$



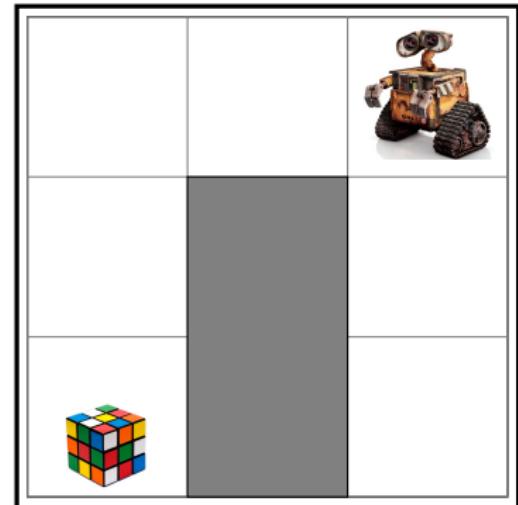


## La malédiction de la dimension

La dimension de l'espace des états est exponentielle en nombre de variables qui le définissent.



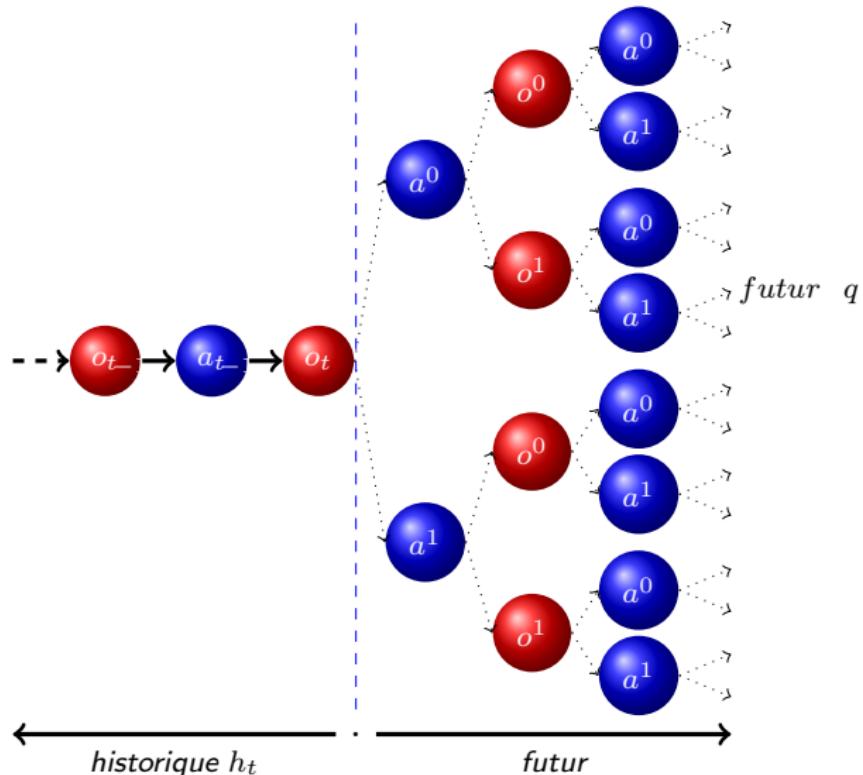
(a) Nombre d'états = 7



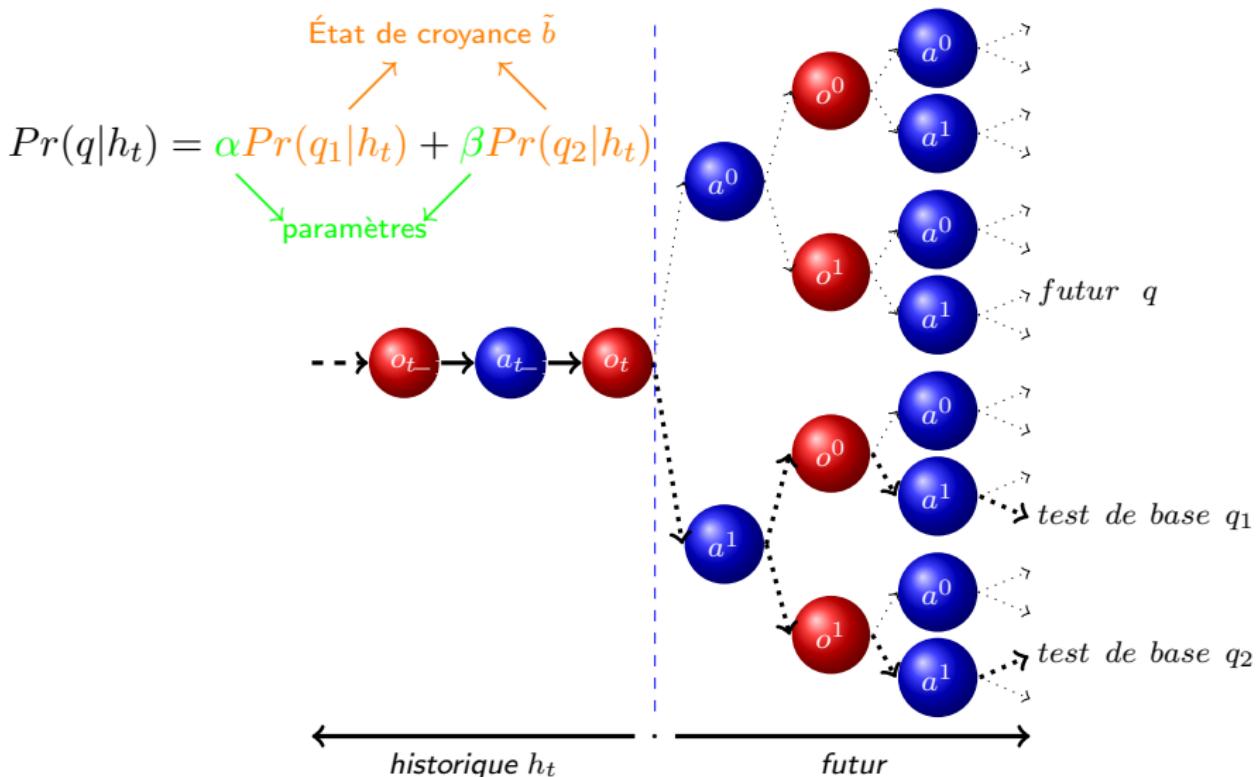
(b) Nombre d'états =  $7 \times 7$

## Représentation prédictive des états (PSR) [Littman et al., 2002]

$$Pr(q|h_t) = ?$$



## Représentation prédictive des états (PSR) [Littman et al., 2002]



## Représentation prédictive des états (PSR) [Littman et al., 2002]

$$U(q|s) = \Pr(q|s)$$

$$U = \begin{matrix} & q_1 & q_2 & q_3 & q_4 & q_5 & q_6 & q_7 \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \\ s_6 \\ s_7 \end{matrix} & \left( \begin{matrix} 0.2 & 0.5 & 1 & 0 & 0.2 & 0.5 & 1 & \dots \\ 0.5 & 0.8 & 0.3 & 0.7 & 0.5 & 0.8 & 0.3 & \dots \\ 0.5 & 0.9 & 0.1 & 0.3 & 0.5 & 0.9 & 0.1 & \dots \\ 1 & 0 & 1 & 0.2 & 1 & 0 & 1 & \dots \\ 0.5 & 0.3 & 0.8 & 0.9 & 0.5 & 0.3 & 0.8 & \dots \\ 0.9 & 0.6 & 0.7 & 0.3 & 0.9 & 0.6 & 0.7 & \dots \\ 0.9 & 0.9 & 0.5 & 0 & 0.9 & 0.9 & 0.5 & \dots \end{matrix} \right) \end{matrix}$$

## Représentation prédictive des états (PSR) [Littman et al., 2002]

$$U(q|s) = Pr(q|s)$$

$$U = \begin{matrix} & \begin{matrix} q_1 & q_2 & q_3 & q_4 & q_5 & q_6 & q_7 \end{matrix} \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \\ s_6 \\ s_7 \end{matrix} & \left( \begin{matrix} 0.2 & 0.5 & 1 & 0 & 0.2 & 0.5 & 1 & \dots \\ 0.5 & 0.8 & 0.3 & 0.7 & 0.5 & 0.8 & 0.3 & \dots \\ 0.5 & 0.9 & 0.1 & 0.3 & 0.5 & 0.9 & 0.1 & \dots \\ 1 & 0 & 1 & 0.2 & 1 & 0 & 1 & \dots \\ 0.5 & 0.3 & 0.8 & 0.9 & 0.5 & 0.3 & 0.8 & \dots \\ 0.9 & 0.6 & 0.7 & 0.3 & 0.9 & 0.6 & 0.7 & \dots \\ 0.9 & 0.9 & 0.5 & 0 & 0.9 & 0.9 & 0.5 & \dots \end{matrix} \right) \end{matrix}$$

*F*

$$Pr(q|s) = aPr(q_1|s) + bPr(q_2|s) + cPr(q_3|s) + dPr(q_4|s)$$

## Représentation prédictive des états (PSR) [Littman et al., 2002]

$$U(q|s) = Pr(q|s)$$

$$U = \begin{pmatrix} & q_1 & q_2 & q_3 & q_4 & q_5 & q_6 & q_7 \\ q_1 & 0.2 & 0.5 & 1 & 0 & 0.2 & 0.5 & 1 & \dots \\ q_2 & 0.5 & 0.8 & 0.3 & 0.7 & 0.5 & 0.8 & 0.3 & \dots \\ q_3 & 0.5 & 0.9 & 0.1 & 0.3 & 0.5 & 0.9 & 0.1 & \dots \\ q_4 & 1 & 0 & 1 & 0.2 & 1 & 0 & 1 & \dots \\ q_5 & 0.5 & 0.3 & 0.8 & 0.9 & 0.5 & 0.3 & 0.8 & \dots \\ q_6 & 0.9 & 0.6 & 0.7 & 0.3 & 0.9 & 0.6 & 0.7 & \dots \\ q_7 & 0.9 & 0.9 & 0.5 & 0 & 0.9 & 0.9 & 0.5 & \dots \end{pmatrix}$$

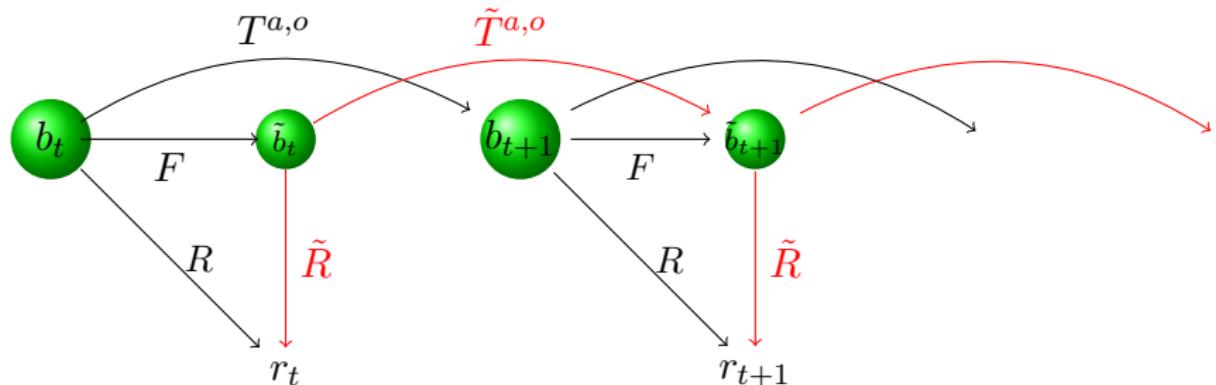
*F*

$$Pr(q|s) \approx a' Pr(q_1|s) + b' Pr(q_2|s) + c' Pr(q_3|s) \quad \cancel{+ d' Pr(q_4|s)}$$

## Compression de l'espace des états



Minimiser  $\|F\tilde{T}^{a,o} - T^{a,o}F\|_\infty$ ,  $\|F\tilde{R} - R\|_\infty$ , et  $\|F\tilde{Z} - Z\|_\infty$ .



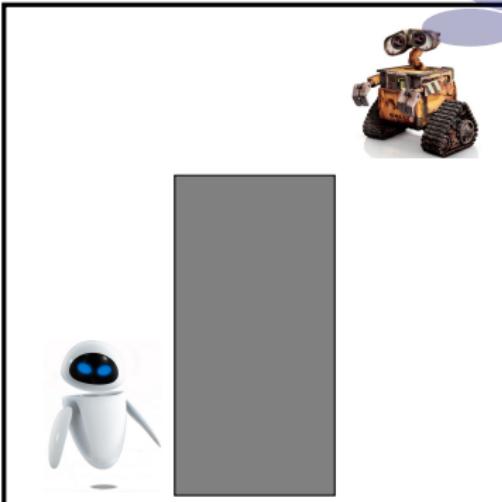
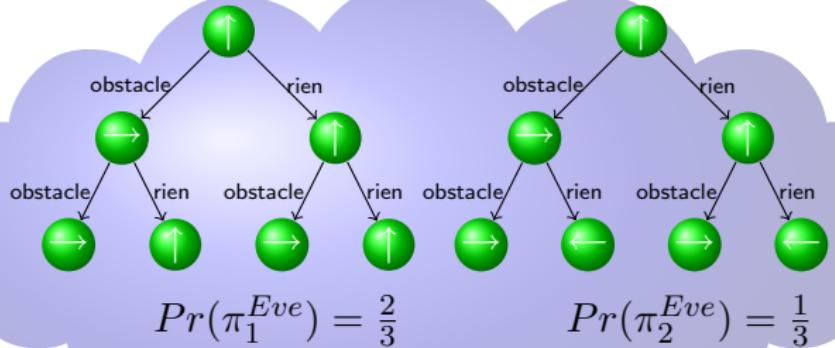
Ce problème peut être formalisé comme un programme linéaire.

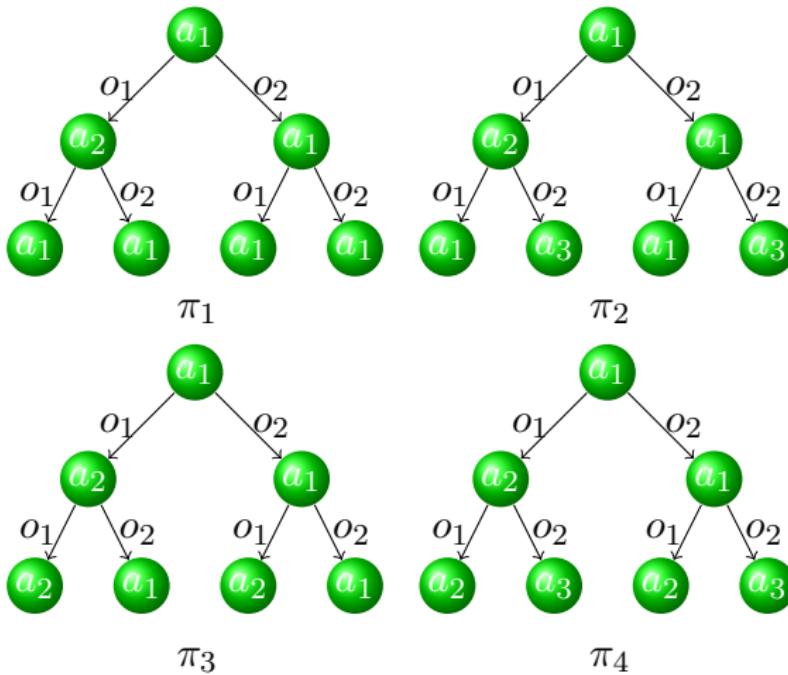
## Résultats empiriques [Boularias et al., ICMLA 2008]

Domaine	PSR approximé			POMDP	
	Compression	Temps	Récompense	Temps	Récompense
<b>Hallway</b>	33%	15.89	0.00	27.58	0.06
	66%	20.70	0.06		
<b>Hallway2</b>	22%	1.56	0.00	80.20	0.02
	50%	5.54	0.01		

TAB.: Temps moyen, en millisecondes, et récompense par étape.

# POMDPs décentralisés





La malédiction de la dimension : La dimension de l'espace des politiques est exponentielle en nombre d'observations et doublement exponentielle en horizon de planification.



## Représentation prédictive des politiques (PPRs)

$$U = \begin{pmatrix} \pi_1 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ \pi_2 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 \\ \pi_3 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 \\ \pi_4 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix} \text{ dimension } = 4$$

The matrix U represents a 4x9 matrix where each row corresponds to a policy vector  $\pi_i$  and each column corresponds to a state-action pair  $a_{10^i a_1 a_2 \dots a_{i-1}}$ . The columns are labeled from left to right as follows:  $a_{10^1 a_1 a_2}, a_{10^1 a_1 a_2 a_3}, a_{10^1 a_1 a_2 a_3 a_4}, a_{10^2 a_1 a_2 a_3}, a_{10^2 a_1 a_2 a_3 a_4}, a_{10^2 a_1 a_2 a_3 a_4 a_5}, a_{10^3 a_1 a_2 a_3 a_4 a_5 a_6}, a_{10^3 a_1 a_2 a_3 a_4 a_5 a_6 a_7}$ .



## Représentation prédictive des politiques (PPRs)

$$U = \begin{pmatrix} \pi_1 & q_1 & q_2 & q_3 \\ \pi_2 & a_{101}a_{201}a_1 & a_{101}a_{201}a_2 & a_{101}a_{201}a_3 \\ \pi_3 & a_{101}a_{202}a_1 & a_{101}a_{202}a_2 & a_{101}a_{202}a_3 \\ \pi_4 & a_{102}a_{102}a_1 & a_{102}a_{102}a_2 & a_{102}a_{102}a_3 \end{pmatrix} \text{dimension} = 4$$



## Représentation prédictive des politiques (PPRs)

$$U = \begin{pmatrix} \pi_1 & \begin{matrix} q_1 & q_2 & q_3 \\ a_{101} & a_{101} & a_{101} \end{matrix} \\ \pi_2 & \begin{matrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{matrix} \\ \pi_3 & \begin{matrix} 0 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 1 \end{matrix} \\ \pi_4 & \begin{matrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \end{matrix} \end{pmatrix} \text{dimension} = 4$$

$$= \begin{pmatrix} 1 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{pmatrix} \times \begin{pmatrix} 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 \end{pmatrix} \text{dimension} = 3$$

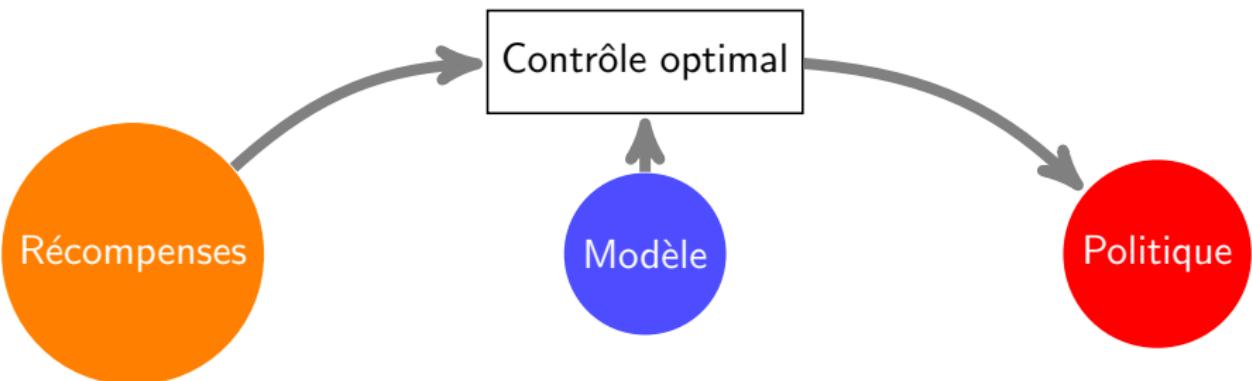
*F*

## Résultats empiriques [Boulaïas & Chaib-draa, ICAPS 2008]

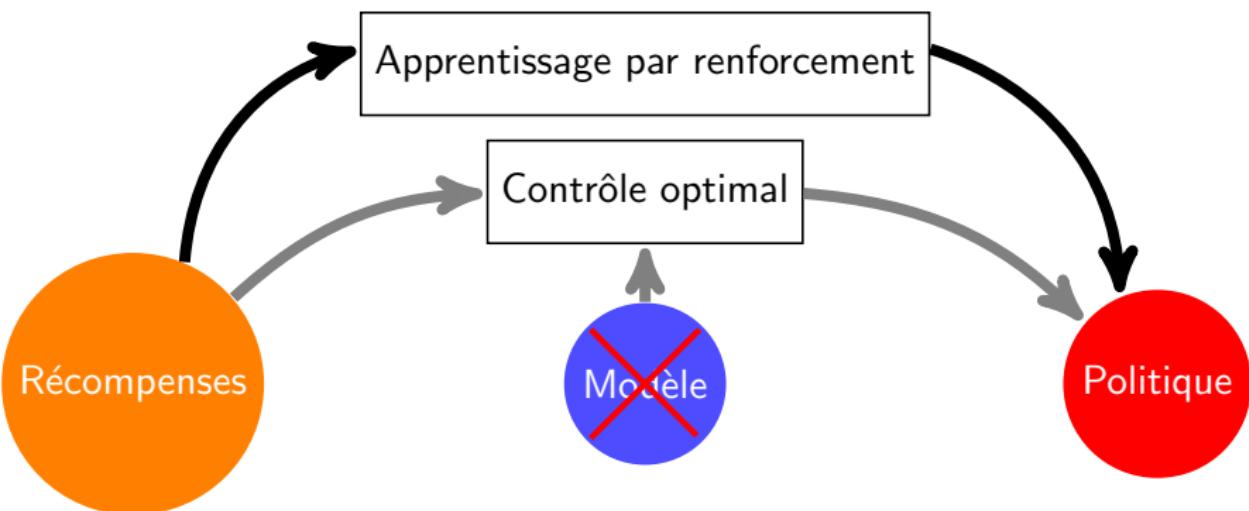
		Arbres de décisions		Représentations prédictives	
Problème	T	temps	arbres	temps	séquences
<b>Multiagent Tiger</b>	2	0.20	(27,27)	0.17	(18,18)
	3	2.29	(675,675)	1.79	(90,90)
	4	-	-	534.90	(540,540)
<b>Multiagent Broadcasting Channel</b>	2	0.12	(8,8)	0.14	(8,8)
	3	0.46	(72,72)	0.36	(24,24)
	4	17.59	(1800,1458)	4.59	(80,80)

**TAB.:** Temps d'exécution (en secondes) et nombre de politiques et de séquences retenues par l'algorithme de programmation dynamique.

## Résoudre un (PO)MDP



## Résoudre un (PO)MDP



## Gradient de la politique

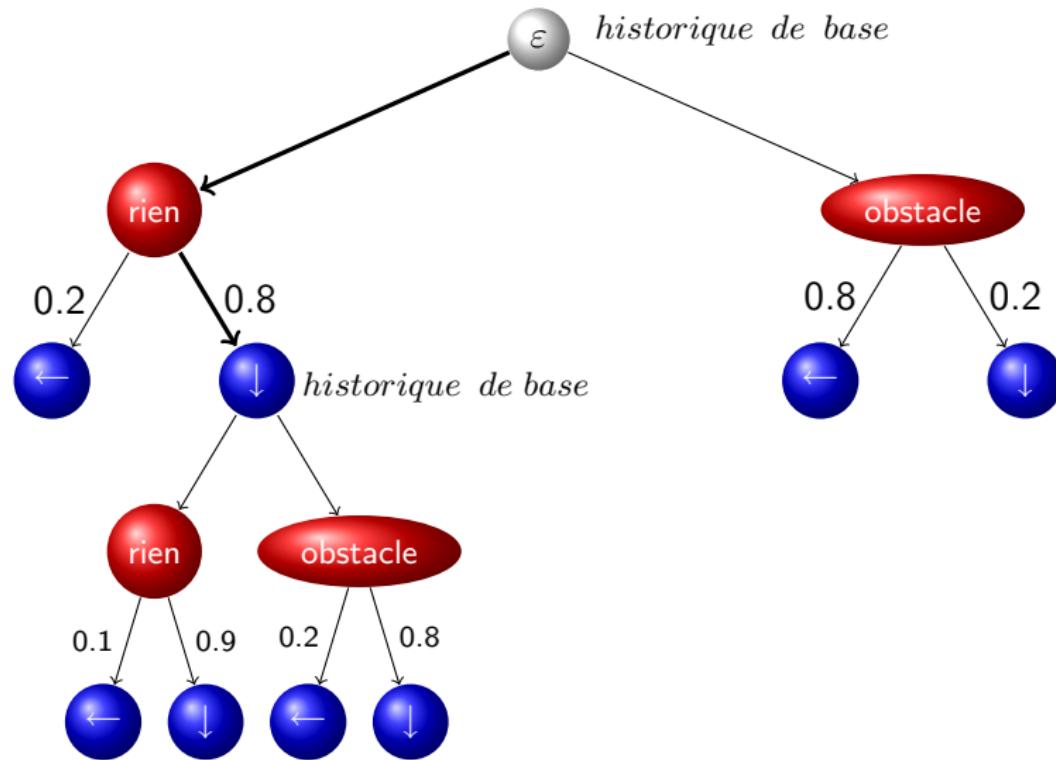
Le gradient de la fonction de valeur d'une politique définie par le vecteur de paramètres  $\theta$  est donné par :

$$\frac{\partial V(\theta)}{\partial \theta_i} = \sum_{t=0}^{\infty} \gamma^t \sum_{a_1 o_1 \dots a_t o_t a_{t+1}} \underbrace{\prod_{i=0}^t Pr(o_i | a_1 o_1 \dots a_i, \theta) \mathbb{E}(r_t)}_{\text{environnement}} \\ \times \underbrace{\frac{\partial \prod_{i=0}^t Pr(a_{i+1} | a_1 o_1 \dots a_i o_i, \theta)}{\partial \theta_i}}_{\text{politique}}$$

Le terme  $\prod_{i=0}^t Pr(o_i | a_1 o_1 \dots a_i, \theta) \mathbb{E}(r_t)$  peut être estimé à partir des trajectoires en utilisant l'échantillonnage préférentiel (Importance Sampling).

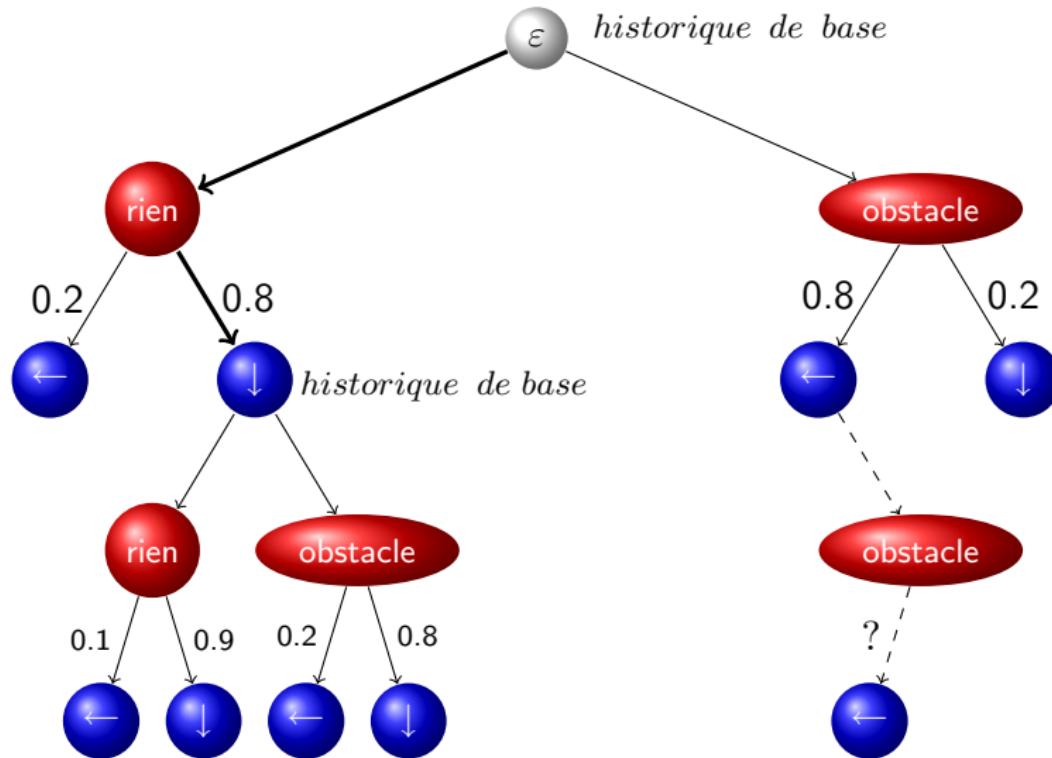


## Représentations prédictives des politiques (PPRs)



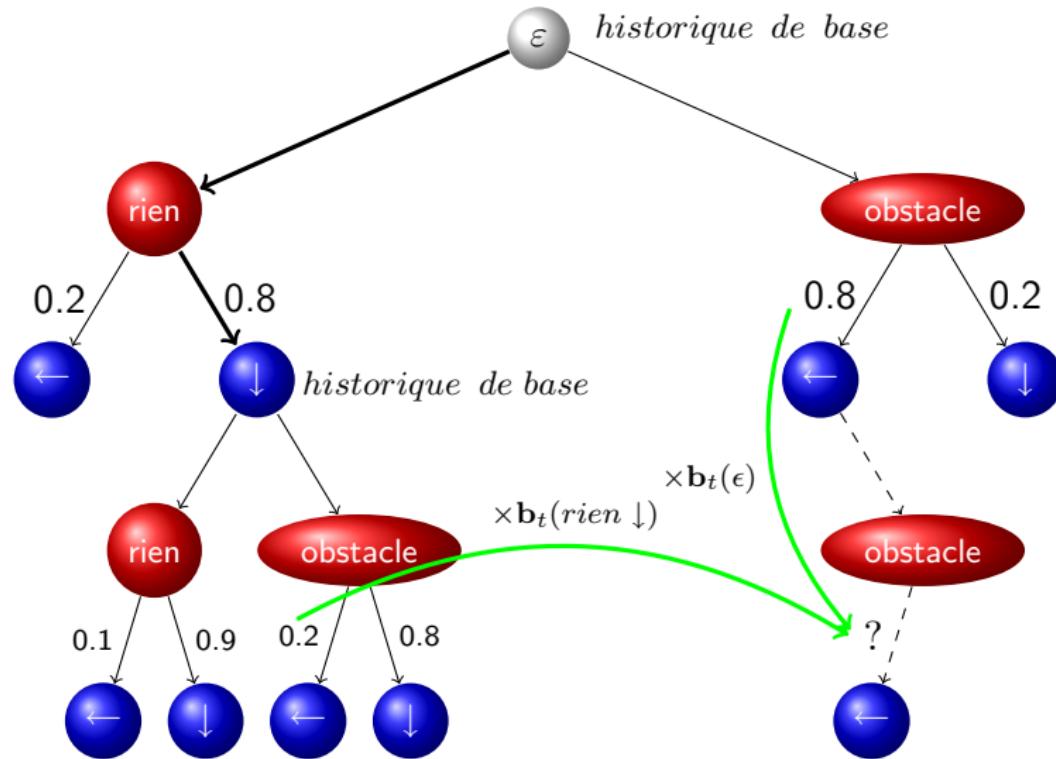


## Représentations prédictives des politiques (PPRs)



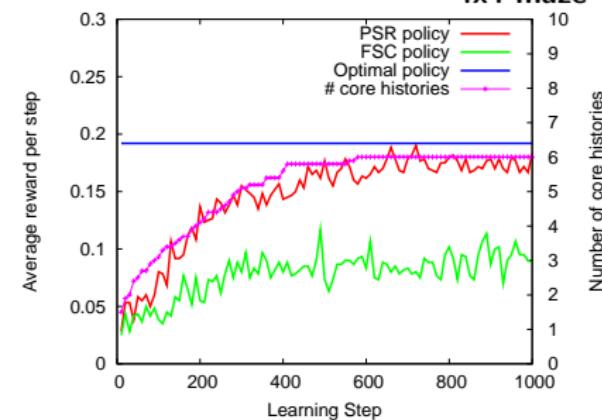


## Représentations prédictives des politiques (PPRs)

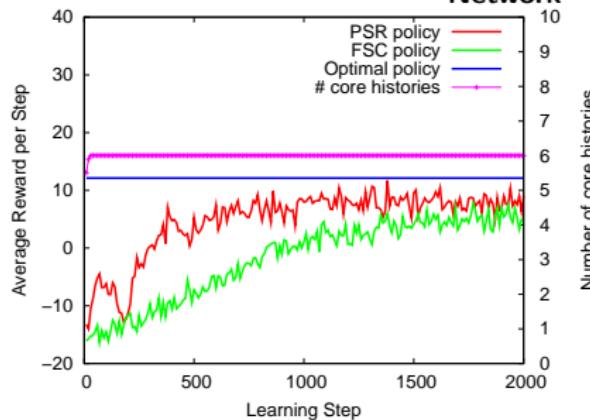


## Résultats empiriques [Boularias & Chaib-draa, ICML 2009]

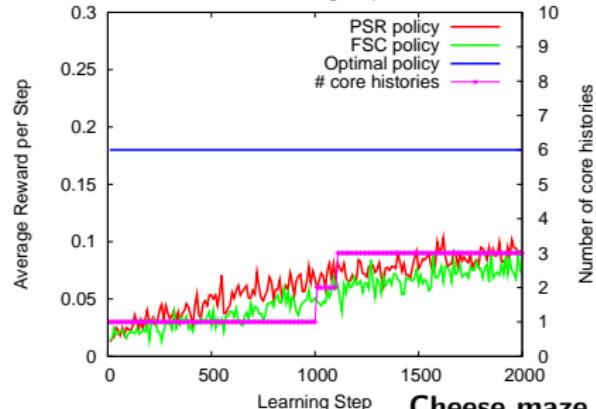
4x4 maze



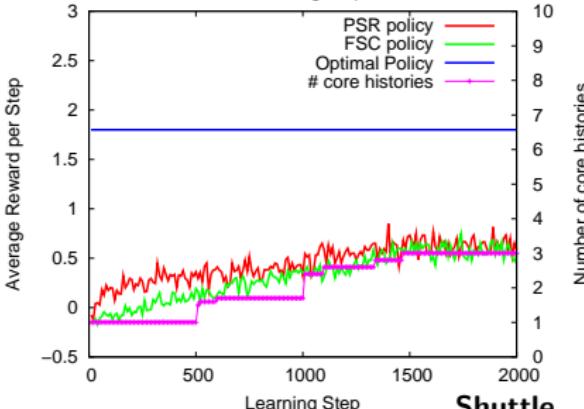
Network



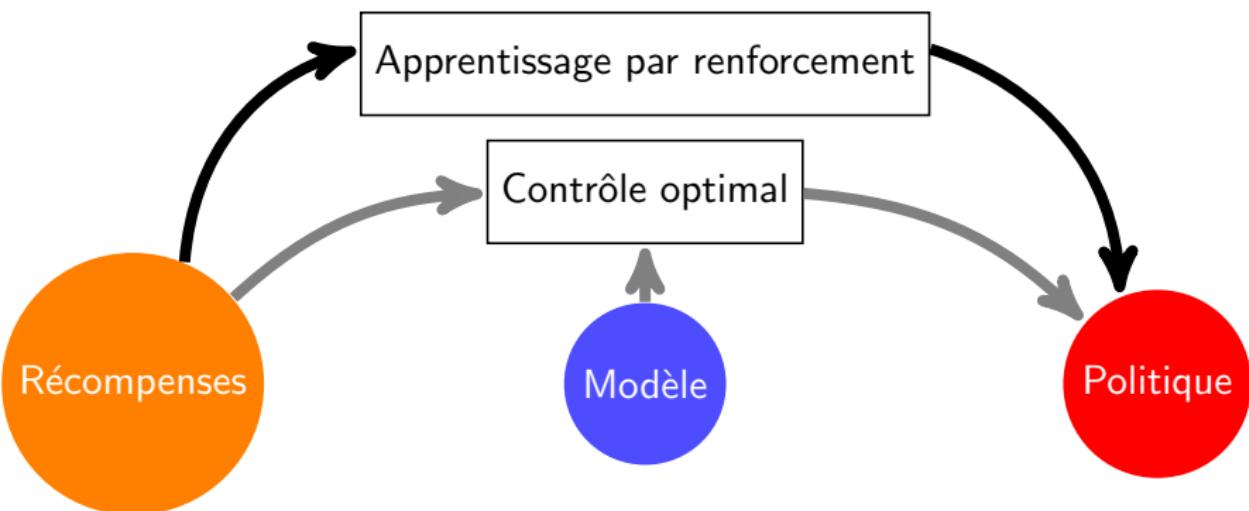
Cheese maze



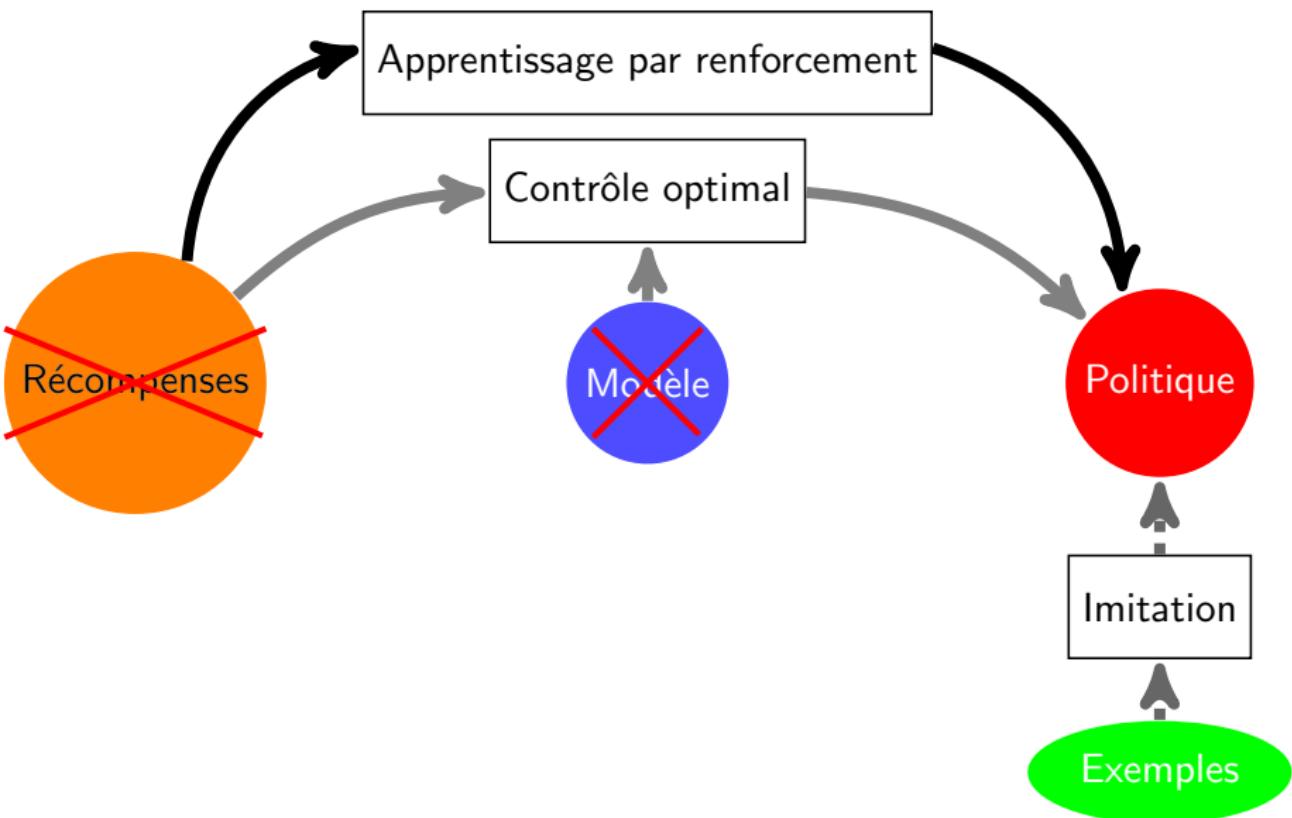
Shuttle



## Résoudre un (PO)MDP



## Résoudre un (PO)MDP



## Apprentissage par imitation direct (Behavioral cloning)

- ✗ Ce problème a été peu étudié dans le contexte d'un environnement **partiellement observable**.
- ✓ L'expert peut être modélisé comme un système dynamique  
⇒ **Représentation prédictive de la politique de l'expert.**

## Apprentissage par imitation avec les représentations prédictives

$$U(q_i|h_j) = \Pr(q_i|h_j)$$

$$U = \begin{pmatrix} & q_1 & q_2 & q_3 & q_4 & q_5 & q_6 & q_7 \\ h_1 & 0.2 & 0.5 & 1 & 0 & 0.2 & 0.5 & 1 & \dots \\ h_2 & 0.5 & 0.8 & 0.3 & 0.7 & 0.5 & 0.8 & 0.3 & \dots \\ h_3 & 0.5 & 0.9 & 0.1 & 0.3 & 0.5 & 0.9 & 0.1 & \dots \\ h_4 & 1 & 0 & 1 & 0.2 & 1 & 0 & 1 & \dots \\ h_5 & 0.5 & 0.3 & 0.8 & 0.9 & 0.5 & 0.3 & 0.8 & \dots \\ h_6 & 0.9 & 0.6 & 0.7 & 0.3 & 0.9 & 0.6 & 0.7 & \dots \\ h_7 & 0.9 & 0.9 & 0.5 & 0 & 0.9 & 0.9 & 0.5 & \dots \\ & \vdots & \end{pmatrix}$$

## Apprentissage par imitation avec les représentations prédictives

$$U(q_i|h_j) = \Pr(q_i|h_j)$$

$$U = \begin{pmatrix} & q_1 & q_2 & q_3 & q_4 & q_5 & q_6 & q_7 \\ h_1 & 0.2 & 0.5 & 1 & 0 & 0.2 & 0.5 & 1 & \dots \\ h_2 & 0.5 & 0.8 & 0.3 & 0.7 & 0.5 & 0.8 & 0.3 & \dots \\ h_3 & 0.5 & 0.9 & 0.1 & 0.3 & 0.5 & 0.9 & 0.1 & \dots \\ h_4 & 1 & 0 & 1 & 0.2 & 1 & 0 & 1 & \dots \\ h_5 & 0.5 & 0.3 & 0.8 & 0.9 & 0.5 & 0.3 & 0.8 & \dots \\ h_6 & 0.9 & 0.6 & 0.7 & 0.3 & 0.9 & 0.6 & 0.7 & \dots \\ h_7 & 0.9 & 0.9 & 0.5 & 0 & 0.9 & 0.9 & 0.5 & \dots \\ \vdots & \end{pmatrix}$$

$F$

## Résultats empiriques [Boularias, ICMLA 2008]

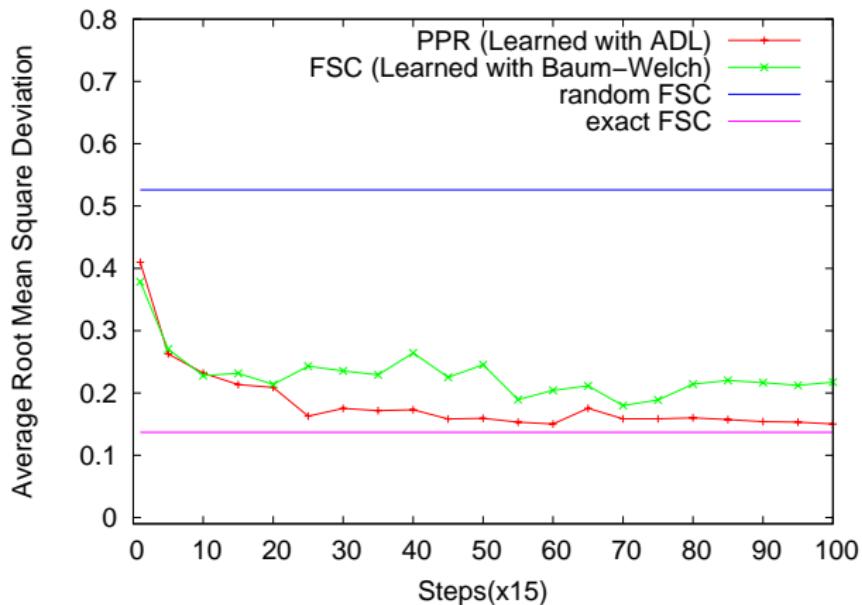


FIG.: Erreur quadratique moyenne de la prédiction des actions de l'expert.

## Apprentissage par imitation direct (Behavioral cloning)

- L'optimalité d'une action dans un état donné dépend de :
  - ① les caractéristiques de l'état.
  - ② les caractéristiques des états voisins.
  - ③ la dynamique du système.
- L'apprentissage par imitation est un problème de prédiction de sortie structurée (**structured output prediction**).

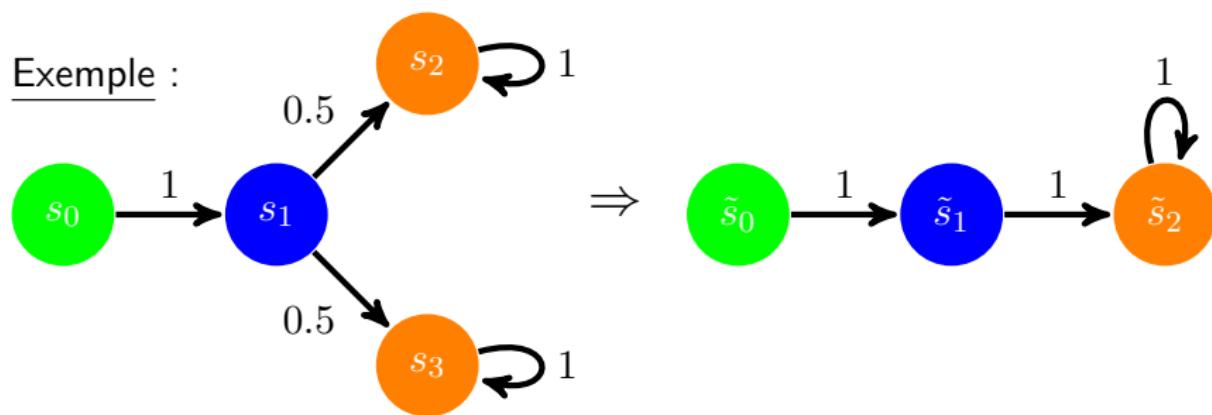
## Morphisme des processus décisionnels de Markov [Ravindran, 2004]

Un morphisme entre deux MDPs  $M(\mathcal{S}, \mathcal{A}, \{T^a\})$  et  $\tilde{M}(\tilde{\mathcal{S}}, \tilde{\mathcal{A}}, \{\tilde{T}^a\})$  est une fonction surjective  $f$  définie par :

$$f : \mathcal{S} \rightarrow \tilde{\mathcal{S}}$$

$$\tilde{T}^a(f(s), \tilde{s}) = \sum_{s' \in \mathcal{S}, f(s') = \tilde{s}} T^a(s, s')$$

Exemple :



## Morphismes souples [Sorg & Singh, 2009]

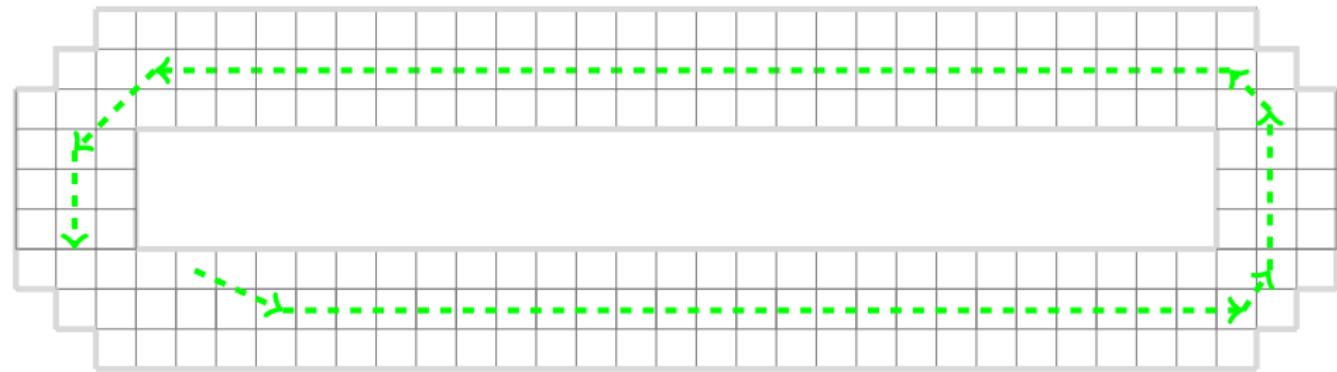
Un morphisme souple entre deux MDPs  $M$  et  $\tilde{M}$  est une fonction  $f$  définie par :

$$f : \mathcal{S} \times \tilde{\mathcal{S}} \rightarrow [0, 1]$$
$$\left\| \sum_{\tilde{s} \in \tilde{\mathcal{S}}} f(s, \tilde{s}) \tilde{T}^a(\tilde{s}, \tilde{s}') - \sum_{s' \in \mathcal{S}} T^a(s, s') f(s, s') \right\| \leq \epsilon$$

- ✓ Le problème de trouver un morphisme souple peut être résolu à l'aide d'un programme linéaire.

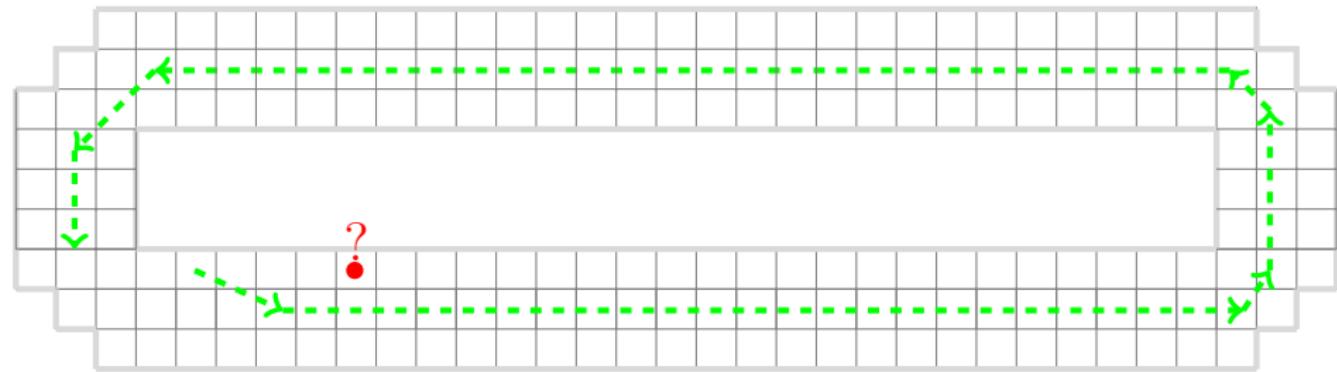
## L'exemple du circuit

 *trajectoire d'un expert*



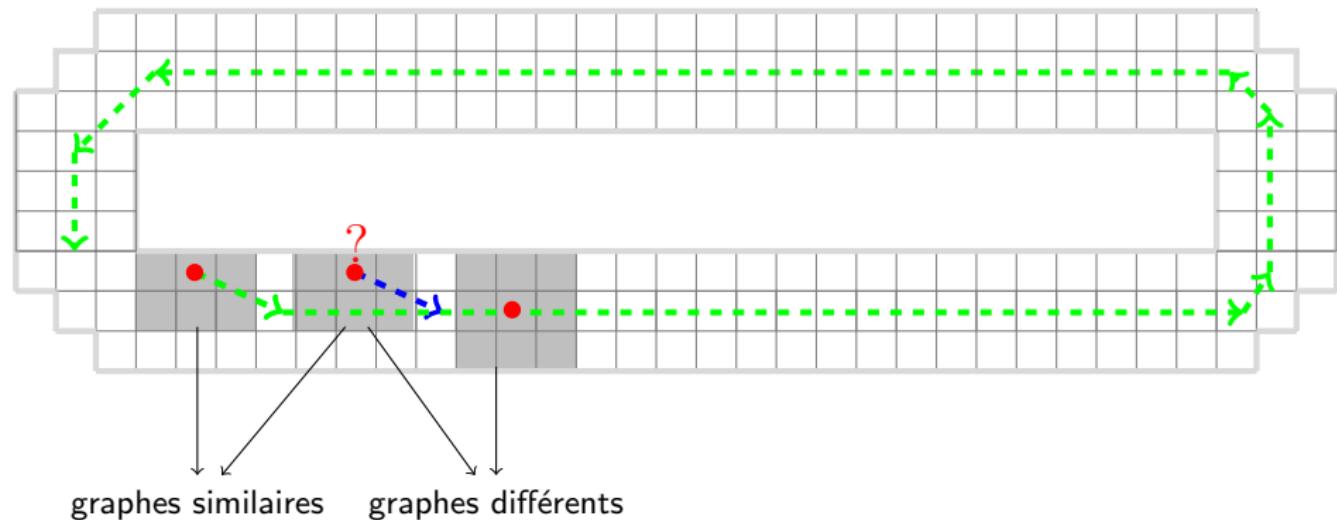
## L'exemple du circuit

 *trajectoire d'un expert*

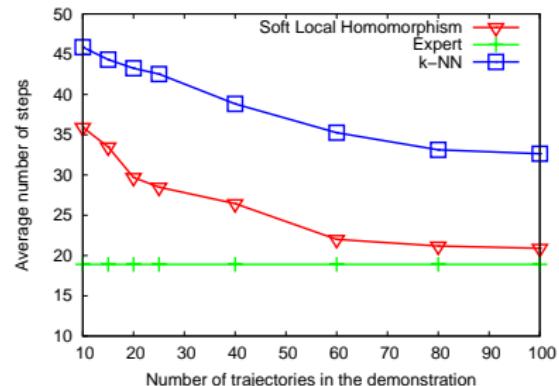


## L'exemple du circuit

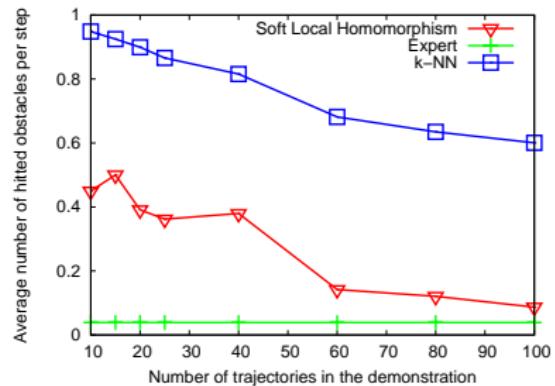
*----- trajectoire d'un expert*



## Résultats empiriques [Boularias & Chaib-draa, ICRA 2010]

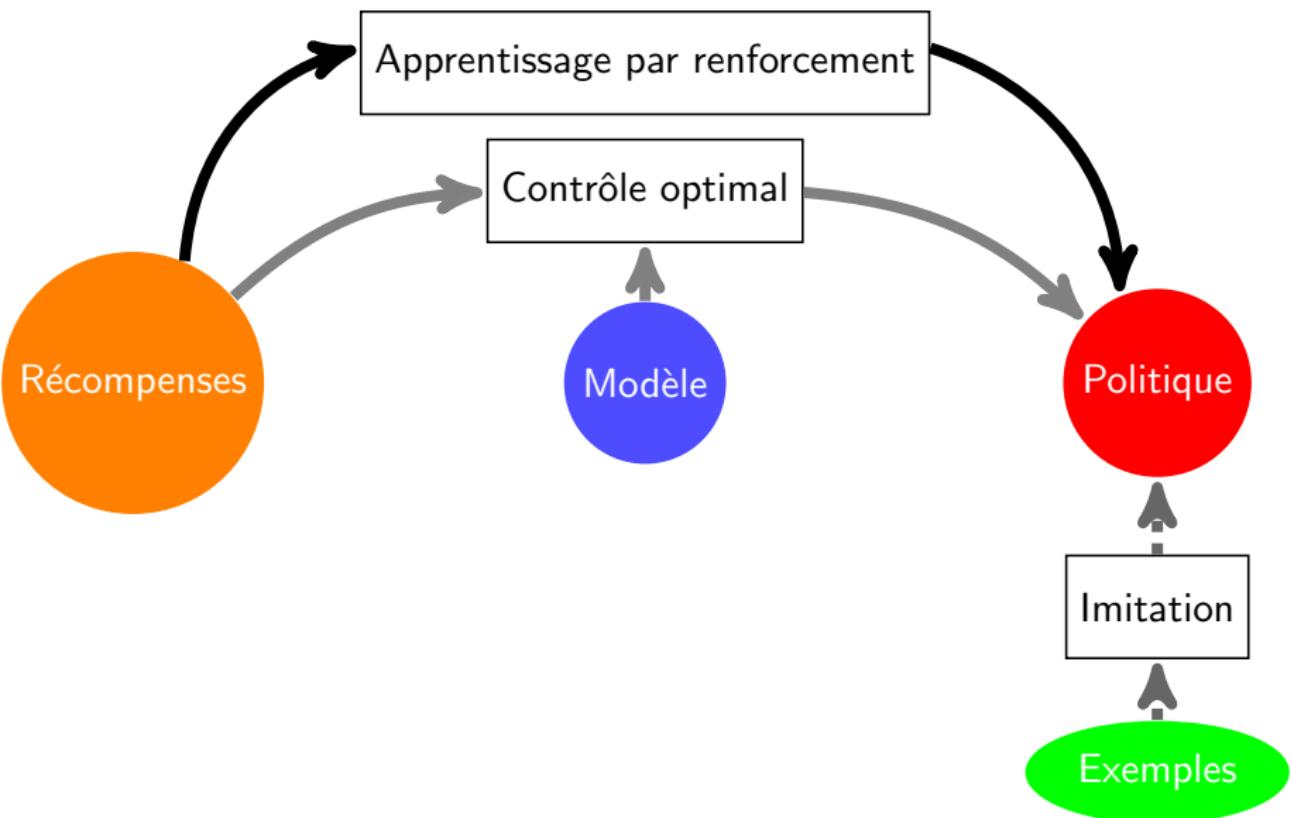


(a) Durée moyenne d'une course

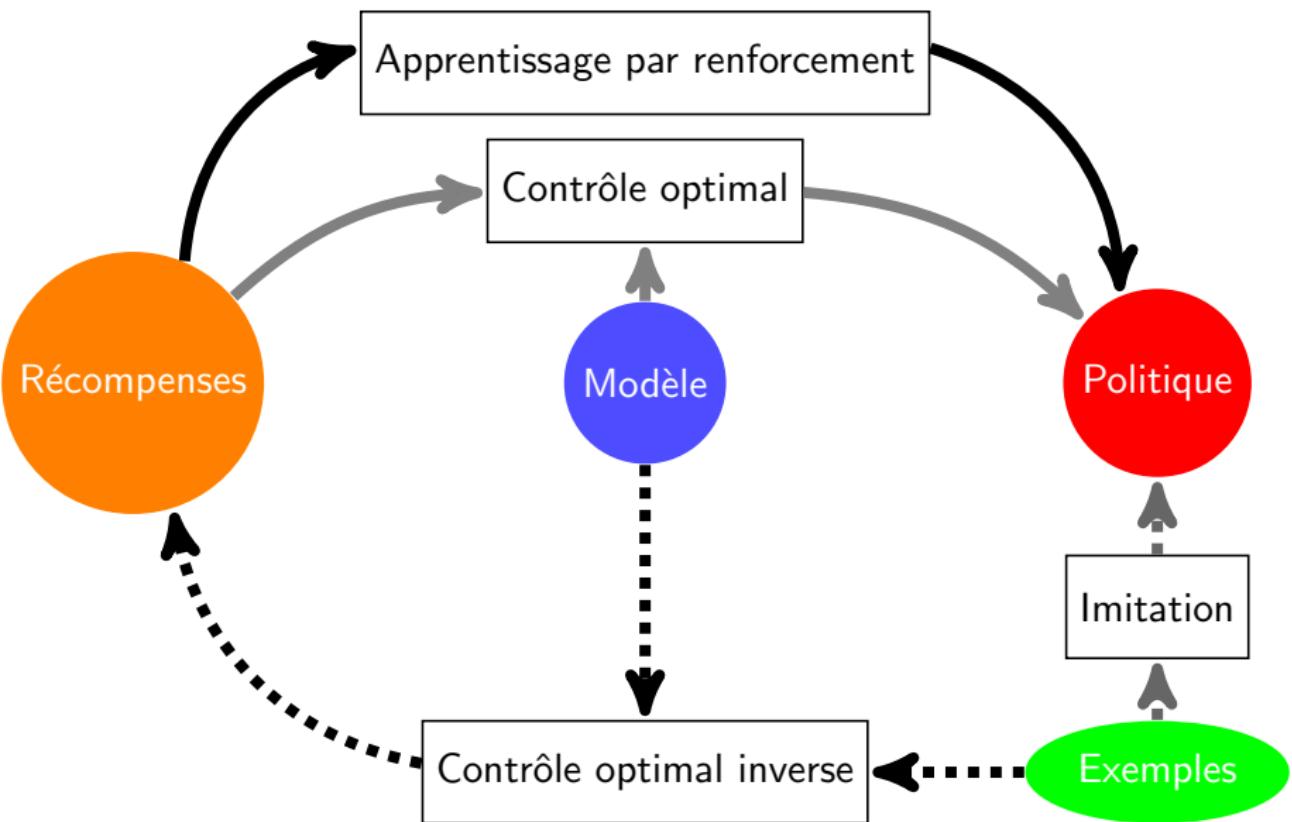


(b) Fréquence des sorties de piste

## Résoudre un (PO)MDP



## Résoudre un (PO)MDP



## Apprentissage par renforcement inverse

Étant donnés des exemples d'une politique  $\pi^*$  de l'expert, trouver un vecteur de récompenses  $R$  tel que :

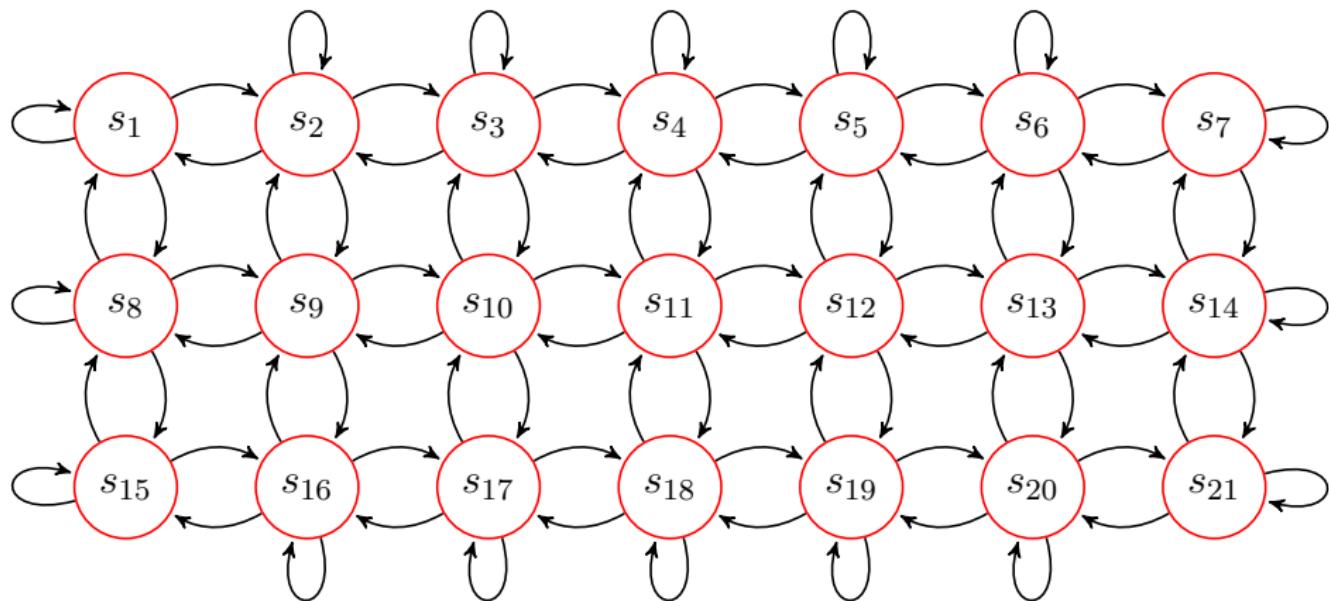
$$R^T \mu_{\pi^*} \geq \max_{\pi \in \Pi(M)} R^T \mu_\pi$$

Deux problèmes principaux :

- ✗ L'espace des solutions est infini (par exemple,  $R = (0, 0, \dots, 0)$  est une solution triviale. Cela est un problème de régularisation.
- ✗ La politique de l'expert  $\pi^*$  n'est pas complètement connue, la distribution stationnaire  $\mu_{\pi^*}$  ne peut pas être calculée analytiquement.



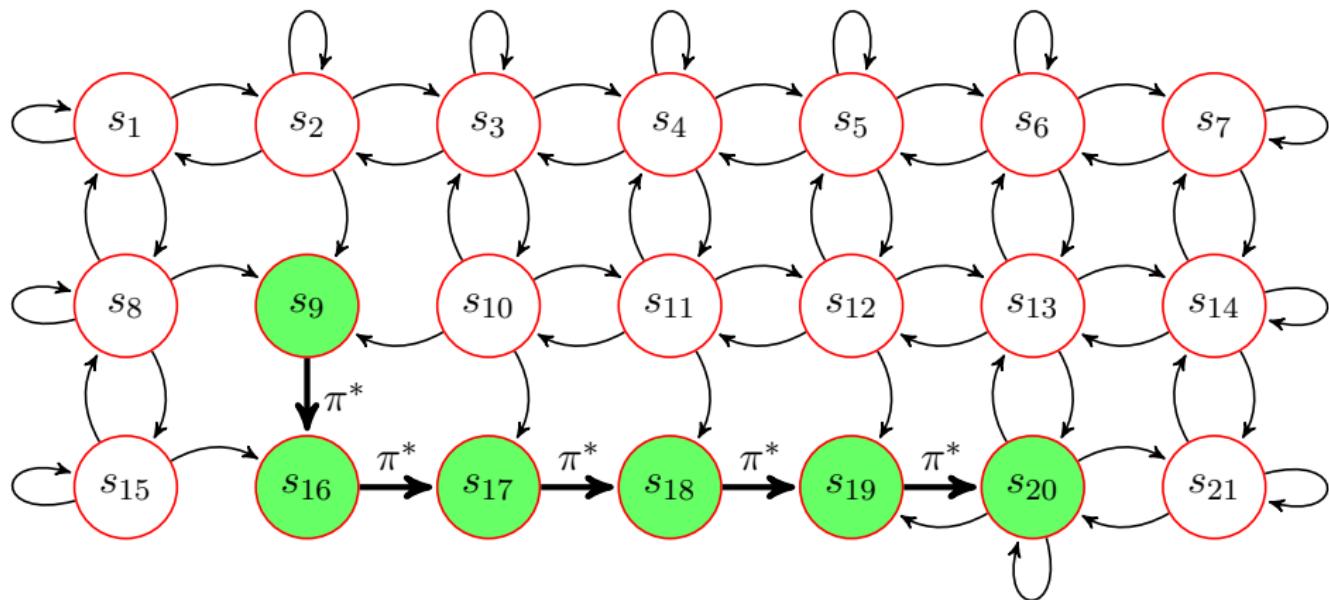
## Bootstrapping des exemples de l'expert



Un processus décisionnel de Markov  $M$ .



## Bootstrapping des exemples de l'expert



Un processus décisionnel de Markov modifié  $M'$ .

$$\max_{\pi \in \Pi(M')} R^T \mu_\pi \geq \max_{\pi \in \Pi(M)} R^T \mu_\pi$$

## Résultats empiriques

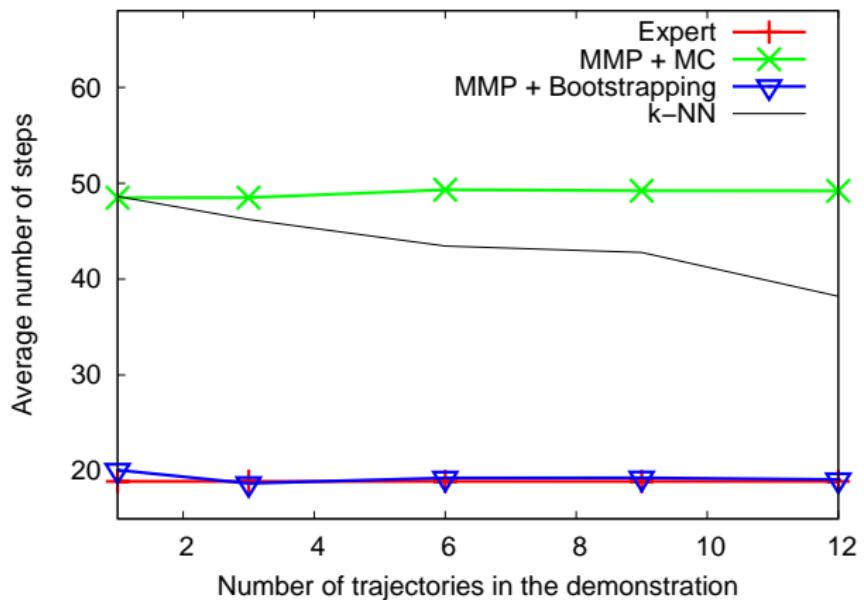


FIG.: Durée moyenne d'une course.

## Conclusion

- ✓ Les représentations prédictives des états et des politiques sont une solution prometteuse aux problèmes de contrôle et d'apprentissage dans les systèmes dynamiques.
- ✓ Les représentations prédictives réduisent la complexité des structures définissant les états et les politiques.
- ✗ Les paramètres des représentations prédictives n'ont pas une interprétation intuitive.
- ✗ Les états de croyance souffrent d'une instabilité numérique.
- ✗ L'apprentissage des tests et des historiques de base reste un problème ouvert.

## Quelques références :

-  Boularias, A. (2008).  
A Predictive Model for Imitation Learning in Partially Observable Environments.  
In *ICMLA'08*, pages 83–90.
-  Boularias, A., Izadi, M., and Chaib-draa, B. (2008).  
Prediction-directed Compression of POMDPs.  
In *ICMLA'08*, pages 99–105.
-  Boularias, A. and Chaib-draa, B. (2008a).  
Exact Dynamic Programming for Decentralized POMDPs with Lossless Policy Compression.  
In *ICAPS'08*, pages 20–27.
-  Boularias, A. and Chaib-draa, B. (2009).  
Predictive Representations for Policy Gradient in POMDPs.  
In *ICML'09*, pages 65–72.
-  Boularias, A. and Chaib-draa, B. (2010).  
Apprenticeship Learning via Soft Local Homomorphisms.  
In *ICRA'10*.

# Questions ?



## Quelques références :

-  Jaeger, H. (2000).  
Observable Operator Models for Discrete Stochastic Time Series.  
*Neural Computation*, 12(6) :1371–1398.
-  Aberdeen, D., & Baxter, J. (2002).  
Scaling Internal-State Policy-Gradient Methods for POMDPs.  
*Proc. 19th Int. Conf. Machine Learning* (pp. 3–10).
-  Littman, M., Sutton, R., & Singh, S. (2002).  
Predictive Representations of State.  
*Advances in Neural Information Processing Systems 14* (pp. 1555–1561).
-  Peters, J. and Schaal, S. (2008).  
Natural Actor-Critic.  
*Neurocomputing*, 71 (pp. 1180–1190).
-  Sutton, R. S., Mcallester, D., Singh, S., & Mansour, Y. (2000).  
Policy Gradient Methods for Reinforcement Learning with Function Approximation.  
*In Advances in Neural Information Processing Systems 12 (NIPS'00)* (pp. 1057–1063).

## Automates finis stochastiques (FSM)

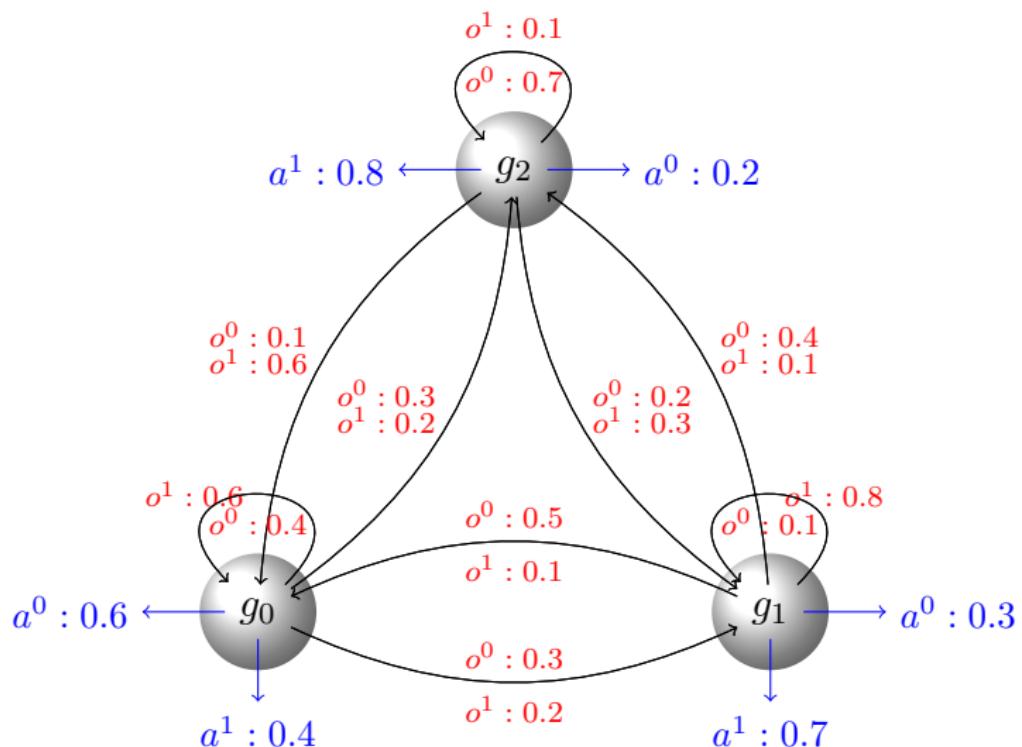


FIG.: Un automate fini stochastique.

## Apprentissage par imitation direct (Behavioral cloning)

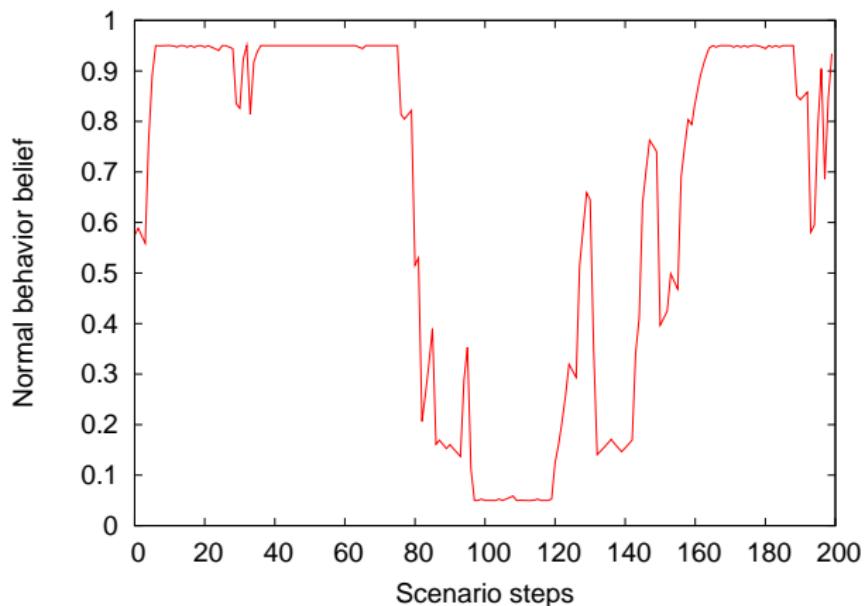
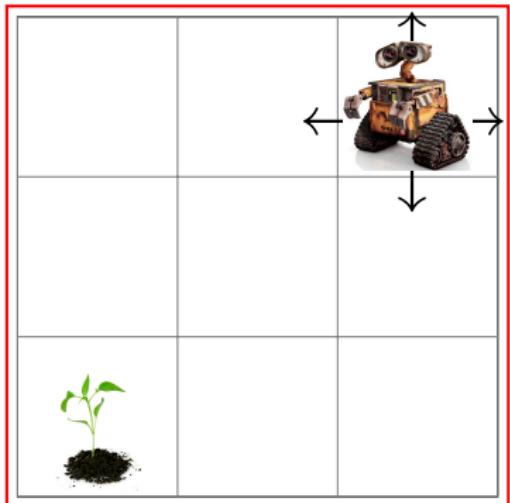


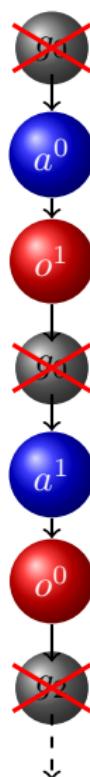
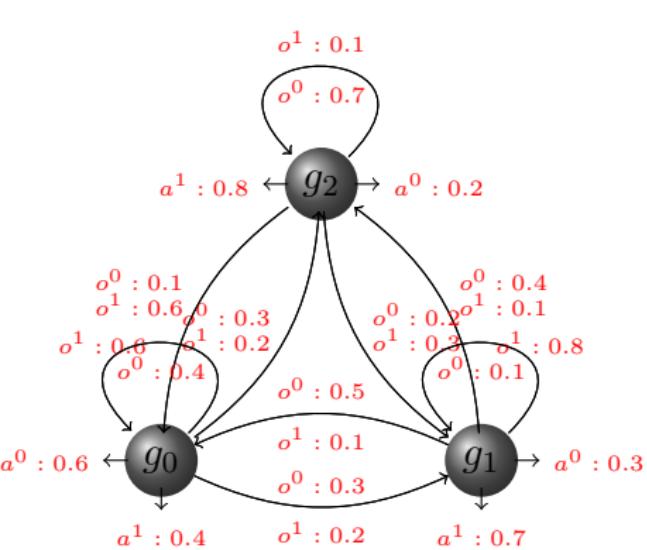
FIG.: Détection d'une déviation du comportement normal.

## Predictive State Representations (PSR)- a simple example



- A POMDP belief state is a probability distribution over the 9 hidden states.
- A PSR belief state is the probabilities of the tests :
  - ① action :↑, observation : wall
  - ② action :↓, observation : wall
  - ③ action :←, observation : wall
  - ④ action :→, observation : wall
- The current belief state of the robot is :
  - ①  $Pr(\text{wall} | \uparrow) = 1$
  - ②  $Pr(\text{wall} | \downarrow) = 0$
  - ③  $Pr(\text{wall} | \leftarrow) = 0$
  - ④  $Pr(\text{wall} | \rightarrow) = 1$

## FSC with Internal Belief States [Aberdeen & Baxter, 2002]



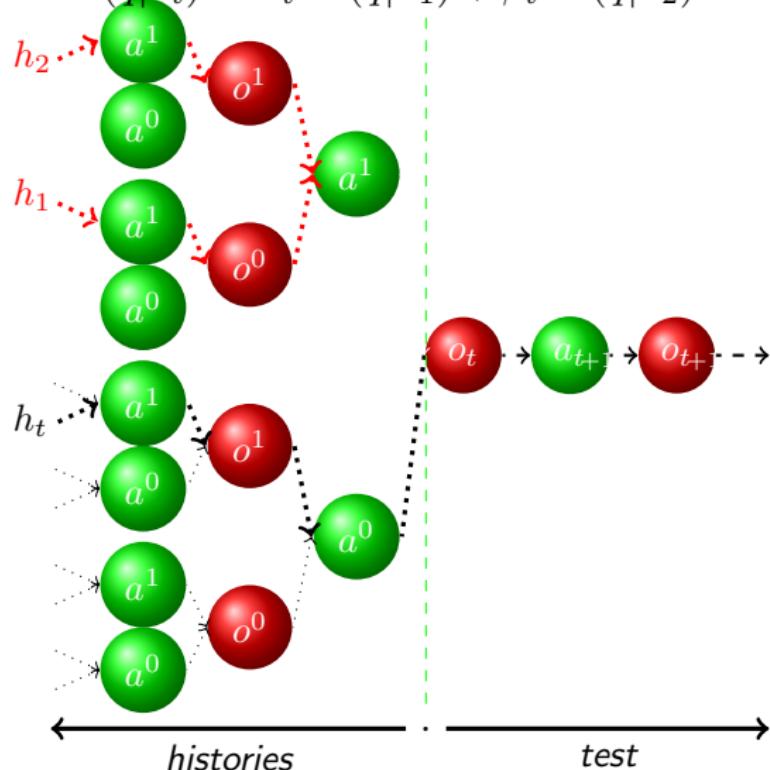
$$Pr(g_0) = 1, Pr(g_1) = 0, Pr(g_2) = 0$$

$$Pr(g_0) = 0.6, Pr(g_1) = 0.2, Pr(g_2) = 0.2$$

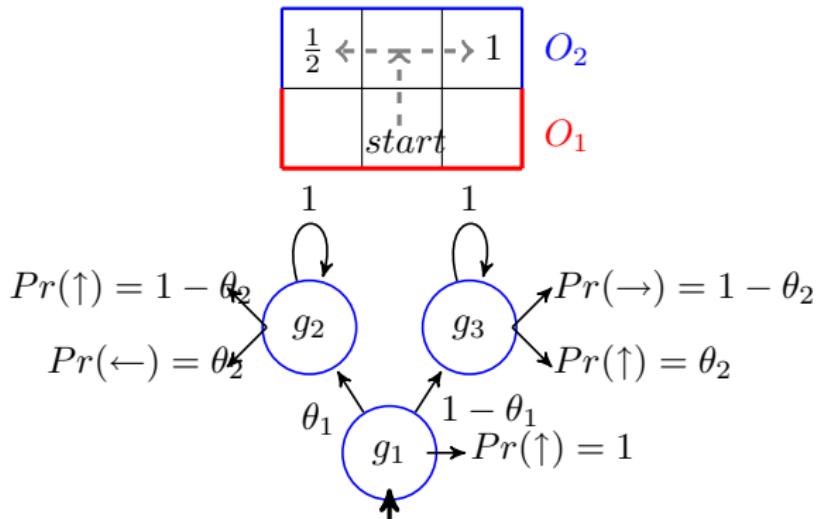
$$Pr(g_0) = 0.36, Pr(g_1) = 0.24, Pr(g_2) = 0.4$$

## Predictive State Representations (PSRs) - core histories

$$\forall q \in \{\mathcal{A} \times \mathcal{O}\}^*: Pr(q|h_t) = \alpha_t Pr(q|h_1) + \beta_t Pr(q|h_2)$$

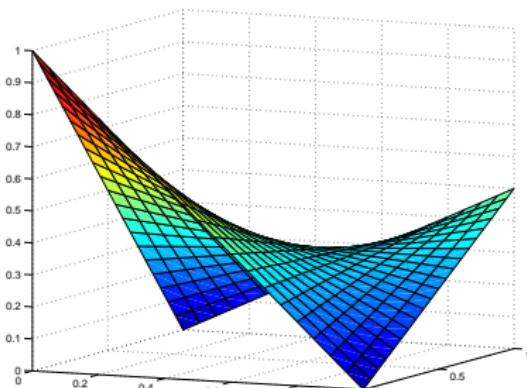


## Small example

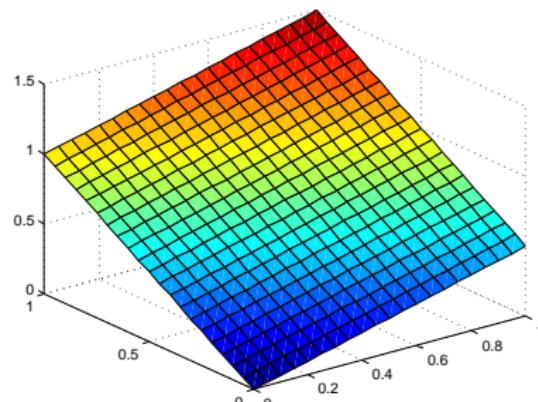


- The equivalent PSR representation has one core history  $h_0 = \emptyset$ , therefore  $b_t(\emptyset) = 1, \forall t \geq 0$ .
- The parameters of the PSR representation are :  $Pr(\leftarrow | O_2, \emptyset) \stackrel{\text{def}}{=} \theta'_1, Pr(\rightarrow | O_2, \emptyset) \stackrel{\text{def}}{=} \theta'_2, Pr(\uparrow | O_2, \emptyset) \stackrel{\text{def}}{=} \theta'_3$

## Example



$$V(\pi_{\theta}^{FSC}) = \frac{1}{2}\theta_1\theta_2 + (1 - \theta_2)(1 - \theta_1)$$



$$V(\pi_{\theta'}^{PSR}) = \frac{1}{2}\theta_1' + \theta_2' \\ \text{s.t. } \theta_1' + \theta_2' \leq 1$$

The value functions of an FSC and its equivalent PSR policy.

## Predictive Reward Representation

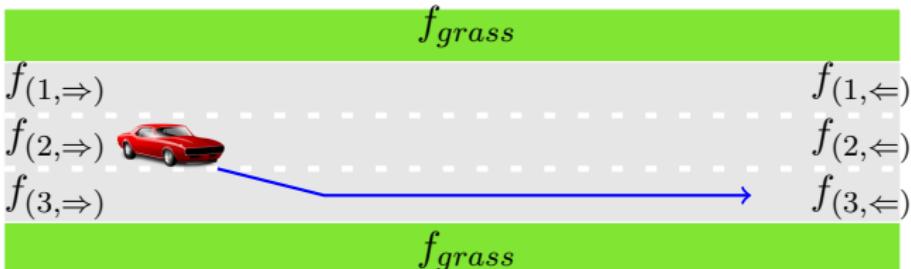


FIG.: A car driving simulation and the corresponding reward features.

- Actions : steer right, steer left, accelerate, decelerate, nothing.
- A high reward is given for driving on the right side of the road.
- However, the feature  $f_{(1,\Leftarrow)}$  does not appear in the demonstration !
- This problem can be solved by adding the feature  $\tilde{f}(steer\ right, f_{grass})$

## Predictive Reward Representation

The set of extended features is denoted by  $\tilde{F}$ .

- ①  $\tilde{F}^0$  is initialized with linearly independent base features  $F$ ;
- ② For  $t > 0$ ,  $\tilde{F}^{t+1} = \tilde{F}^t \cup \{T^a \tilde{f} | a \in \mathcal{A}, \tilde{f} \in \tilde{F}^t\}$ ;
- ③ At the end of each iteration  $t$ , only independent vectors are kept in  $\tilde{F}^t$ ;

This process stops when  $\tilde{F}^{t+1} = \tilde{F}^t$ .

## Results 2

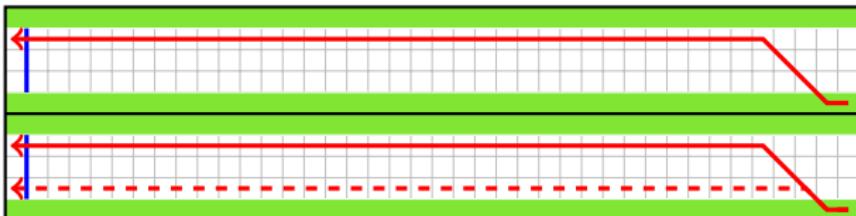


FIG.: The trajectories of the expert cover only the upper half of the road. In the lower half, the policy learned using extended features was able to keep the car on the right side, contrary to the policy learned using base features (dashed line).

- How to determine if a feature is important.
- Computational cost of searching for latent features.
- Other types of latent features (using images as input).
- Nonparametric bayesian method for infinite latent features.