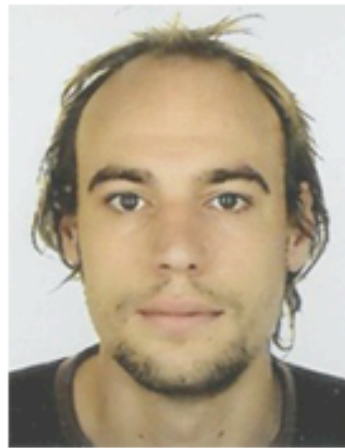


# POLICY SEARCH IN A SPACE OF SIMPLE CLOSED-FORM FORMULAS: TOWARDS INTERPRETABILITY OF REINFORCEMENT LEARNING

FRANCIS MAES RAPHAEL FONTENEAU



Université  
de Liège



LOUIS WEHENKEL DAMIEN ERNST



**OPTIMAL SEQUENTIAL DECISION MAKING** IS A CENTRAL PROBLEM OF COMPUTER SCIENCE AND HAS A HUGE NUMBER OF APPLICATIONS



## MEDICAL THERAPY OPTIMIZATION



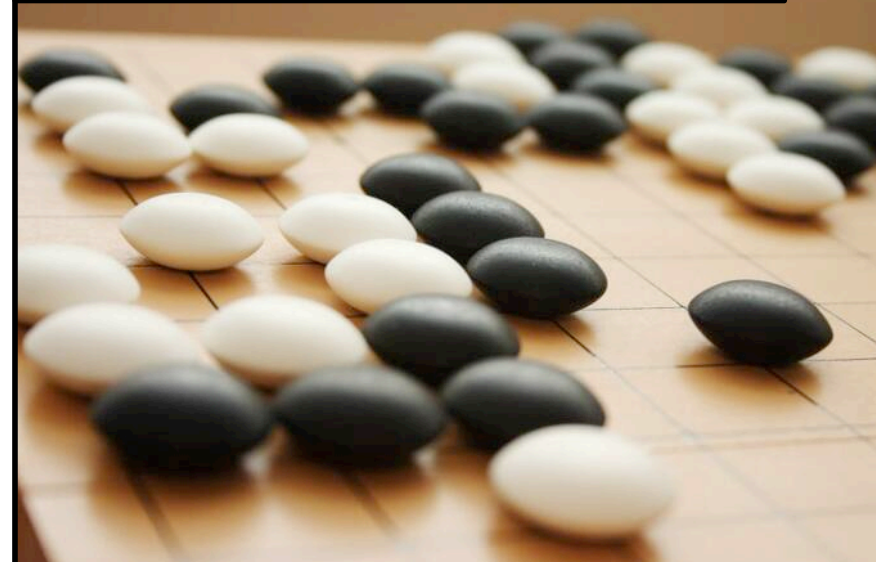
## FINANCIAL TRADING



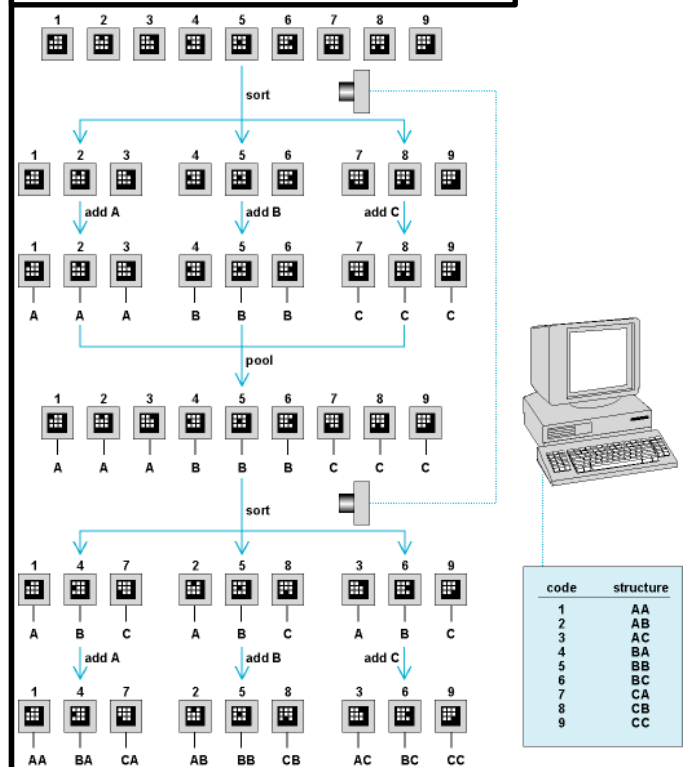
## ROBOT OPTIMAL CONTROL



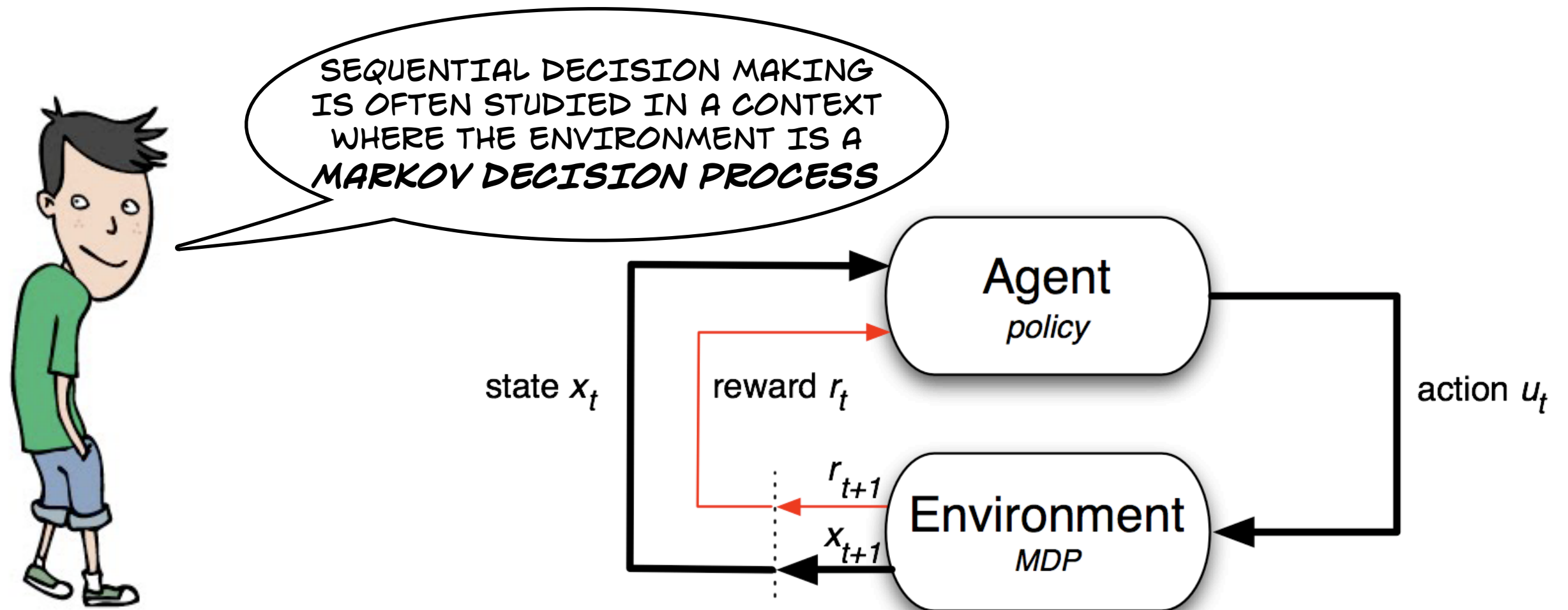
## AUTOMATIC GAME PLAYING



## COMBINATORIAL SEARCH







AT EACH TIME STEP  $t$ , THE AGENT IS IN A STATE  $x_t$  AND SELECTS AN ACTION  $u_t$ .

THE ENVIRONMENT SENDS IN RETURN AN INSTANTANEOUS REWARD  $r_t$ .

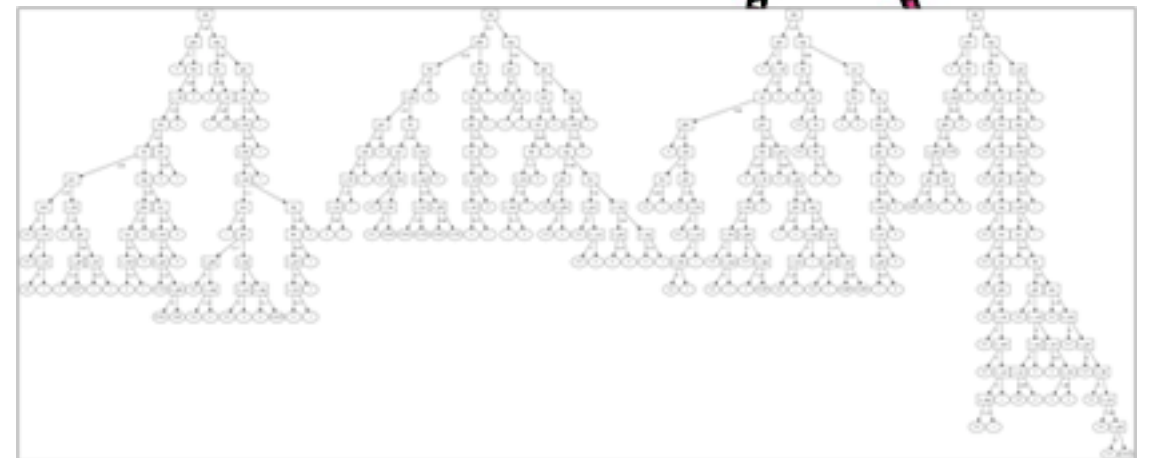
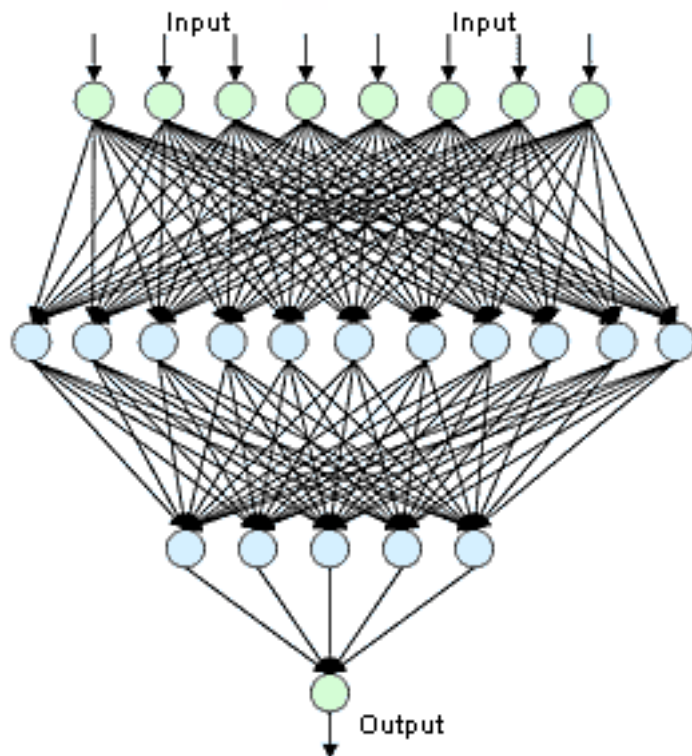
THE AIM OF THE AGENT IS TO SELECT ACTIONS SO AS TO MAXIMIZE THE **LONG-TERM** SUM OF REWARDS



REINFORCEMENT LEARNING IS A CLASS OF **SAMPLING-BASED** APPROACHES TO LEARN DECISION MAKING POLICIES

WHEN THE STATE SPACE IS LARGE, **FUNCTION APPROXIMATION** TECHNIQUES ARE USED TO COMPACTLY STORE THE DECISION POLICY

BEST PERFORMING RL TECHNIQUES RELY ON **BLACK-BOX** FUNCTION APPROXIMATORS, MOSTLY COMING FROM THE FIELD OF SUPERVISED LEARNING: NEURAL NETWORKS, RANDOM FORESTS, RADIAL BASIS FUNCTIONS, ...





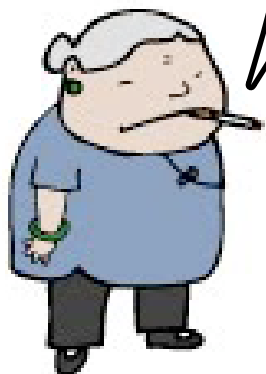
HUMANS TEND TO DISTRUST  
COMPUTER GENERATED  
DECISIONS, ESPECIALLY WHEN  
THE POLICY IS A BLACK-BOX



NO WAY I'M GONNA USE  
A BLACK-BOX DECISION  
POLICY FOR MY FINANCIAL  
PROBLEM



NEITHER ME FOR  
CLINICAL DECISIONS  
THAT IMPACT THE  
PATIENTS' HEALTH



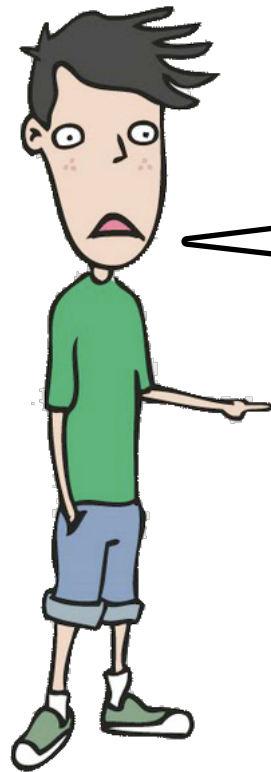
HOSPITAL

EMERGENCY ROOM

FURTHERMORE,  
APPLYING A DECISION  
POLICY TO THE REAL WORLD  
OFTEN INVOLVES ISSUES  
BEYOND THE SCOPE OF  
COMPUTER SCIENCE (E.G  
ETHICAL, POLITICAL,  
IDEOLOGICAL)







THERE ARE MANY EFFICIENT RL ALGORITHMS, SOME OF THEM WITH STRONG THEORETICAL GUARANTEES, BUT THESE ALGORITHMS DO NOT LEAVE THE LABORATORIES

IF WE WANT TO CHANGE THIS, WE MUST PROVIDE DECISION POLICIES THAT HUMANS CAN UNDERSTAND AND EVENTUALLY TRUST



WE THUS NEED RL ALGORITHMS PRODUCING **INTERPRETABLE POLICIES**



INTERPRETABILITY IS AN OLD TOPIC IN SUPERVISED LEARNING, BUT SURPRISINGLY NEARLY ABSENT FROM THE FIELD OF RL

# PROPOSED APPROACH

WE PROPOSE A  
DIRECT POLICY  
SEARCH SCHEME IN A  
SPACE OF  
INTERPRETABLE  
POLICIES

WE CONSIDER  
INDEX-BASED  
POLICIES DEFINED  
BY SIMPLE CLOSED-  
FORMED FORMULAS

AND WE SOLVE THE  
LEARNING PROBLEM  
USING MULTI-ARMED  
BANDITS





WE FOCUS ON  
PROBLEMS WITH  
FINITE NUMBER OF  
ACTIONS AND  
CONTINUOUS STATE  
SPACES

## MDP DYNAMICS

$$x_{t+1} \sim p_f(\cdot | x_t, u_t) \quad t = 0, 1, \dots$$

$$r_t \sim p_\rho(\cdot | x_t, u_t)$$

## OBJECTIVE: MAXIMIZE RETURN

$$J^\pi = \mathbb{E}_{p_0(\cdot), p_f(\cdot), p_\rho(\cdot)} [\mathcal{R}^\pi(x_0)]$$

$$\mathcal{R}^\pi(x_0) = \sum_{t=0}^{\infty} \gamma^t r_t$$

## POLICY

$$\pi(x_t) \sim p_\pi(\cdot | x_t)$$





AN INDEX FUNCTION  $I(\cdot, \cdot)$  IS A MAPPING:

$$I : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$$

GIVEN AN INDEX FUNCTION, WE CAN DEFINE THE INDEX-BASED POLICY:

$$\forall x \in \mathcal{X}, \pi_I(x) \in \arg \max_{u \in \mathcal{U}} I(x, u)$$

WE FOCUS ON THE CLASS OF INDEX-BASED POLICIES WHOSE INDEX FUNCTIONS ARE DEFINED BY ***SIMPLE CLOSED-FORM FORMULAS***



A FORMULA  $F \in \mathbb{F}$  IS:

EITHER A BINARY EXPRESSION  $F = B(F', F'')$ ,  
OR AN UNARY EXPRESSION  $F = U(F')$ ,  
OR A VARIABLE  $F = V$ ,  
OR A CONSTANT  $F = C$ .

WE USE THE FOLLOWING  
OPERATORS AND CONSTANTS:

$$\mathbb{B} = \{+, -, \times, \div, \min, \max\}$$

$$\mathbb{U} = \{\sqrt{\cdot}, \ln(\cdot), |\cdot|, -\cdot, \frac{1}{\cdot}\}$$

$$\mathbb{C} = \{1, 2, 3, 5, 7\}$$

WE CONSIDER  
TWO DIFFERENT  
SETTINGS FOR  
THE VARIABLES:

LOOKAHEAD FREE

$$\mathbb{V} = \mathbb{V}_{LF} = \left\{ x_t^{(1)}, \dots, x_t^{(d_{\mathcal{X}})}, u_t^{(1)}, \dots, u_t^{(d_{\mathcal{U}})} \right\}$$

ONE-STEP LOOKAHEAD (MODEL ACCESSIBLE)

$$\mathbb{V} = \mathbb{V}_{OL} = \left\{ x_t^{(1)}, \dots, x_t^{(d_{\mathcal{X}})}, u_t^{(1)}, \dots, u_t^{(d_{\mathcal{U}})}, r_t, x_{t+1}^{(1)}, \dots, x_{t+1}^{(d_{\mathcal{X}})} \right\}$$



GIVEN A POLICY  $\pi$ , WE DEFINE:

$$D_F(\pi) = \{F \in \mathbb{F} \mid \pi_F = \pi\}$$

THE KOLMOGOROV  
COMPLEXITY OF  $\pi$  IS:

$$\kappa(\pi) = \min_{F \in D_F(\pi)} |F|$$

GIVEN  $K$ , OUR SPACE OF INTERPRETABLE POLICIES  
IS DEFINED BY:

$$\Pi_{int}^K = \{\pi \mid D_F(\pi) \neq \emptyset \text{ and } \kappa(\pi) \leq K\}$$

I REALLY LIKE  
THIS SLIDE







CONSTRUCTING  $\Pi_{int}^K$  IS NON TRIVIAL (EXCEPT FOR FINITE STATE SPACES)

WE INSTEAD APPROXIMATE THIS SPACE BY COMPARING POLICIES ON A FINITE SET OF SAMPLES

1. WE ENUMERATE ALL FORMULAS  $|F| \leq K$

2. GIVEN A FINITE SET OF STATE POINTS  $S = \{s_i\}_{i=1}^S$ ,

WE CLUSTERIZE FORMULAS. TWO FORMULAS  $F$  AND  $F'$  ARE EQUIVALENT IFF:

$$\forall s \in \{s_1, \dots, s_S\},$$

$$\arg \max_{u \in \mathcal{U}} F(s, u, r, y) = \arg \max_{u \in \mathcal{U}} F'(s, u, r, y)$$

3. AMONG EACH CLUSTER, WE SELECT A FORMULA OF MINIMAL LENGTH

4. WE GATHER ALL SELECTED FORMULAS OF MINIMAL LENGTH AND WE DENOTE:

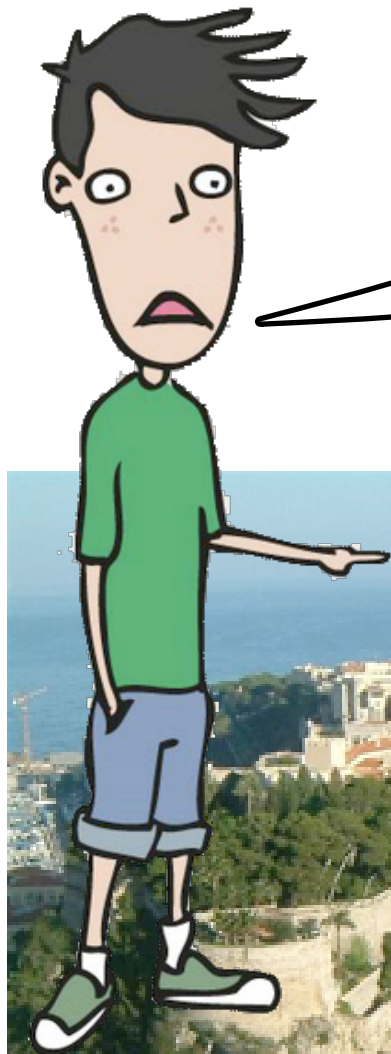
$$\tilde{\Pi}_{int}^K = \{\pi_{F_1}, \dots, \pi_{F_N}\}$$

WE NOW HAVE N CANDIDATE  
POLICIES. TO IDENTIFY THE BEST  
ONE, WE HAVE TO OPTIMIZE:

$$J^\pi = \mathbb{E}_{p_0(\cdot), p_f(\cdot), p_r(\cdot)} [\mathcal{R}^\pi(x_0)]$$

THIS MAY BE  
EXCESSIVELY LONG WITH  
A NAIVE MONTE-CARLO  
APPROACH

I MAY HAVE A  
SOLUTION TO THIS  
PROBLEM !





A MULTI-ARMED BANDIT PROBLEM IS A SEQUENTIAL GAME

- N ARMS WITH UNKNOWN REWARD DISTRIBUTIONS
- AT EACH STEP, THE PLAYER SELECTS ONE OF THE ARMS AND RECEIVES A REWARD DRAWN FROM THE ASSOCIATED DISTRIBUTION
- BEST ARM IDENTIFICATION: WITHIN A FINITE NUMBER OF STEPS  $T$ , SELECT WHICH ARMS TO PLAY SO AS TO IDENTIFY THE ARM WHICH PERFORMS BEST ON AVERAGE

=> EXPLORATION / EXPLOITATION DILEMMA


**EXPLORATION:**  
TRYING ARMS THAT MAY  
POTENTIALLY BE GOOD

**EXPLOITATION:**  
FOCUSING ON ARMS THAT WE  
ALREADY KNOW TO BE GOOD

A GOOD MULTI-ARMED BANDIT  
STRATEGY BALANCES EXPLORATION  
AND EXPLOITATION








IN OUR CASE:

- ONE ARM PER CANDIDATE POLICY
- PLAYING AN ARM: SELECTING AN INITIAL STATE, PERFORMING A TRAJECTORY WITH THE POLICY ASSOCIATED TO THE ARM AND OBSERVING THE RETURN AS A REWARD



NOTE THAT A MAJORITY OF OUR FORMULAS TYPICALLY ENCODE POLICIES PERFORMING VERY BAD

SUCH BAD POLICIES CAN BE DISCARDED QUICKLY

HENCE, THE COMPUTATIONAL BUDGET CAN RAPIDLY BE SPENT ON THE FEW GOOD PERFORMING POLICIES



HERE IS OUR POLICY LEARNING ALGORITHM

SINCE IT RELIES ON DIRECT EVALUATION OF POLICY RETURNS, IT IS AN INSTANCE OF "DIRECT POLICY SEARCH"

1. CONSTRUCT THE APPROXIMATE SET OF CANDIDATE POLICIES

$$\tilde{\Pi}_{int}^K = \{\pi_{F_1}, \dots, \pi_{F_N}\}$$

2. PLAY EACH ARM ONCE (= DRAW ONE TRAJECTORY PER CANDIDATE POLICY)

3. WHILE THERE IS TRAINING TIME:

A) SELECT THE ARM WHICH MAXIMIZES  $A_{n,t} = \bar{r}_{n,t} + \frac{\alpha}{\theta_{n,t}}$

WHERE  $\bar{r}_{n,t}$  IS THE EMPIRICAL MEAN OF RETURNS ASSOCIATED TO  $\pi_{F_n}$

$\theta_{n,t}$  IS THE NUMBER OF TIMES  $\pi_{F_n}$  HAS BEEN PLAYED

$\alpha > 0$  IS AN EXPLORATION/EXPLOITATION TRADEOFF CONSTANT

B) DRAW AN INITIAL STATE, PERFORM ONE TRAJECTORY WITH  $\pi_{F_n}$  AND OBSERVE THE RETURN

C) UPDATE  $\bar{r}_{n,t}$  AND  $\theta_{n,t}$

4. RETURN THE POLICY (OR POLICIES) THAT MAXIMIZE(S)  $\bar{r}_{n,t}$

# EXPERIMENTS

WE EXPERIMENT OUR  
APPROACH ON SIX  
CLASSICAL  
BENCHMARKS

WE COMPARE THE  
LOOKAHEAD-FREE  
AND ONE-STEP  
LOOKAHEAD  
VARIANTS...

... AGAINST "NON-  
INTERPRETABLE" RL  
TECHNIQUES





OUR BENCHMARKS: LINEAR POINT, LEFT OR RIGHT, CAR ON THE HILL, ACROBOT SWING UP, BICYCLE BALANCING, HIV THERAPY



BENCHMARK	LP	LoR	CAR	ACR	B	HIV
$d_x$	2	1	2	4	5	6
$d_u$	1	1	1	1	2	2
$m$	2	2	2	2	9	4
Stoch.	no	yes	no	no	yes	no
$\#V_{LF}$	3	2	3	5	7	8
$\#V_{OL}$	6	4	6	10	13	15
$\gamma$	.9	.75	.95	.95	.98	.98
$T$	50	20	1000	100	5e4	300

THIS TABLE REPORTS THE DIMENSIONALITY OF THE STATE AND ACTION SPACES, THE NUMBER OF ACTIONS, STOCHASTICITY, NUMBER OF VARIABLES IN OUR TWO SETTINGS, DISCOUNT FACTOR AND THE HORIZON USED FOR LEARNING





WE TESTED OUR  
APPROACH WITH 4  
SETTINGS

TO DISCRIMINATE  
AMONGS THE FORMULAS,  
WE USE SAMPLES OF 100  
STATE POINTS

WE COMPARE AGAINST THE  
FOLLOWING TECHNIQUES:  
- RANDOM POLICY  
- LOOKAHEAD POLICIES  
- FITTED Q-ITERATION

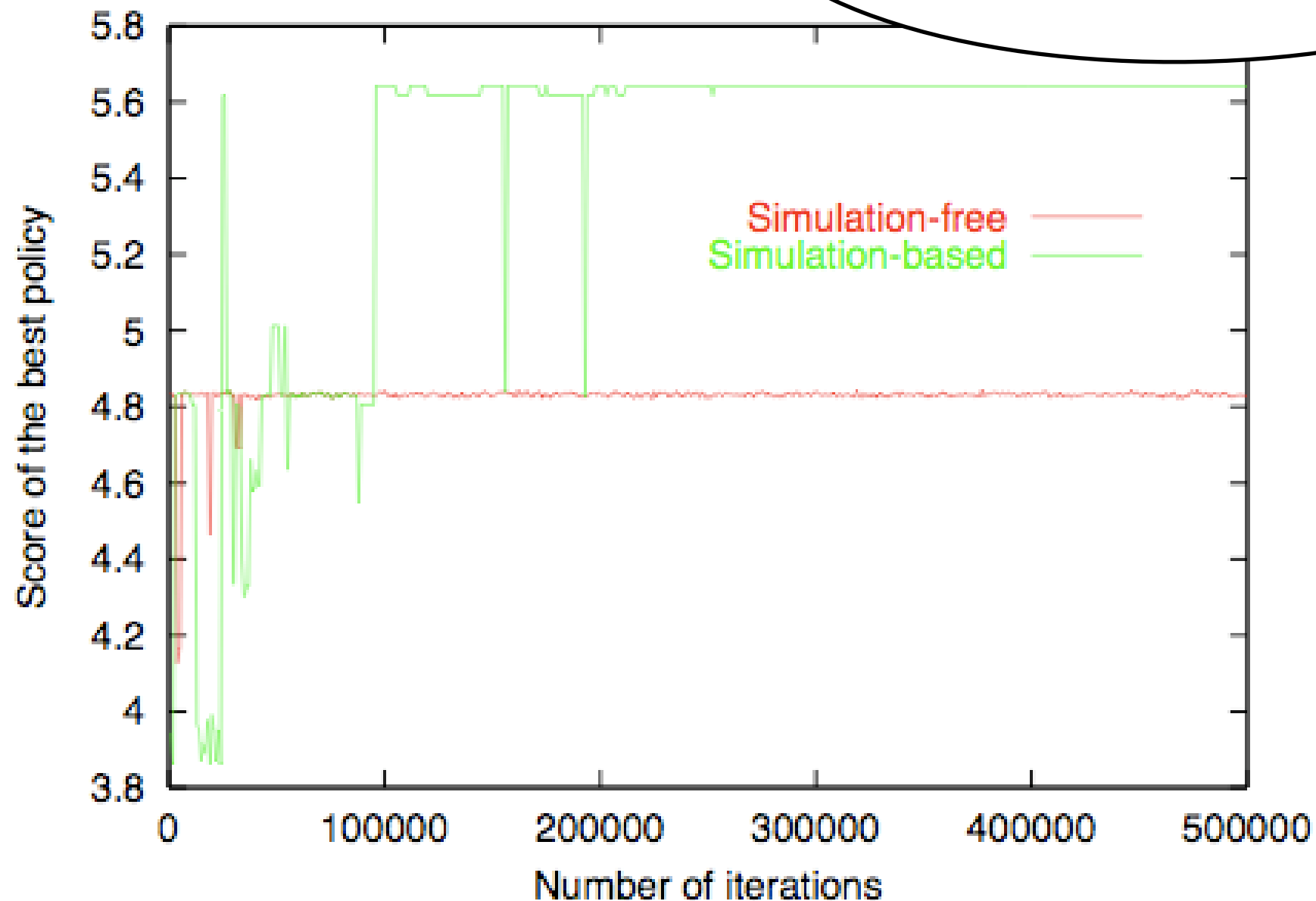
WE PERFORM LEARNING WITH:  
K=5  $\Rightarrow 10^6$  BANDIT STEPS  
K=6  $\Rightarrow 10^7$  BANDIT STEPS  
WE USE A FIXED VALUE  $\alpha = 2$

	CANDIDATE FORMULAS	CANDIDATE POLICIES
K=5, LOOKAHEAD-FREE	80,000 - 340,000	500 - 11,500
K=5, ONE-STEP LOOKAHEAD	140,000 - 990,000	3,800 - 95,000
K=6, LOOKAHEAD-FREE	1,000,000 - 5,500,000	3,600 - 132,000
K=6, ONE-STEP LOOKAHEAD	2,100,000 - 18,500,000	31,000 - 1,200,000



THIS FIGURE SHOWS A TYPICAL  
RUN OF THE ALGORITHM (K=5,  
LINEAR POINT BENCHMARK)

IN THIS SETTING, WE HAVE 907 (RESP.  
12,214) CANDIDATE POLICIES IN THE  
LOOKAHEAD FREE (RESP. ONE-STEP  
LOOKAHEAD) SETTING



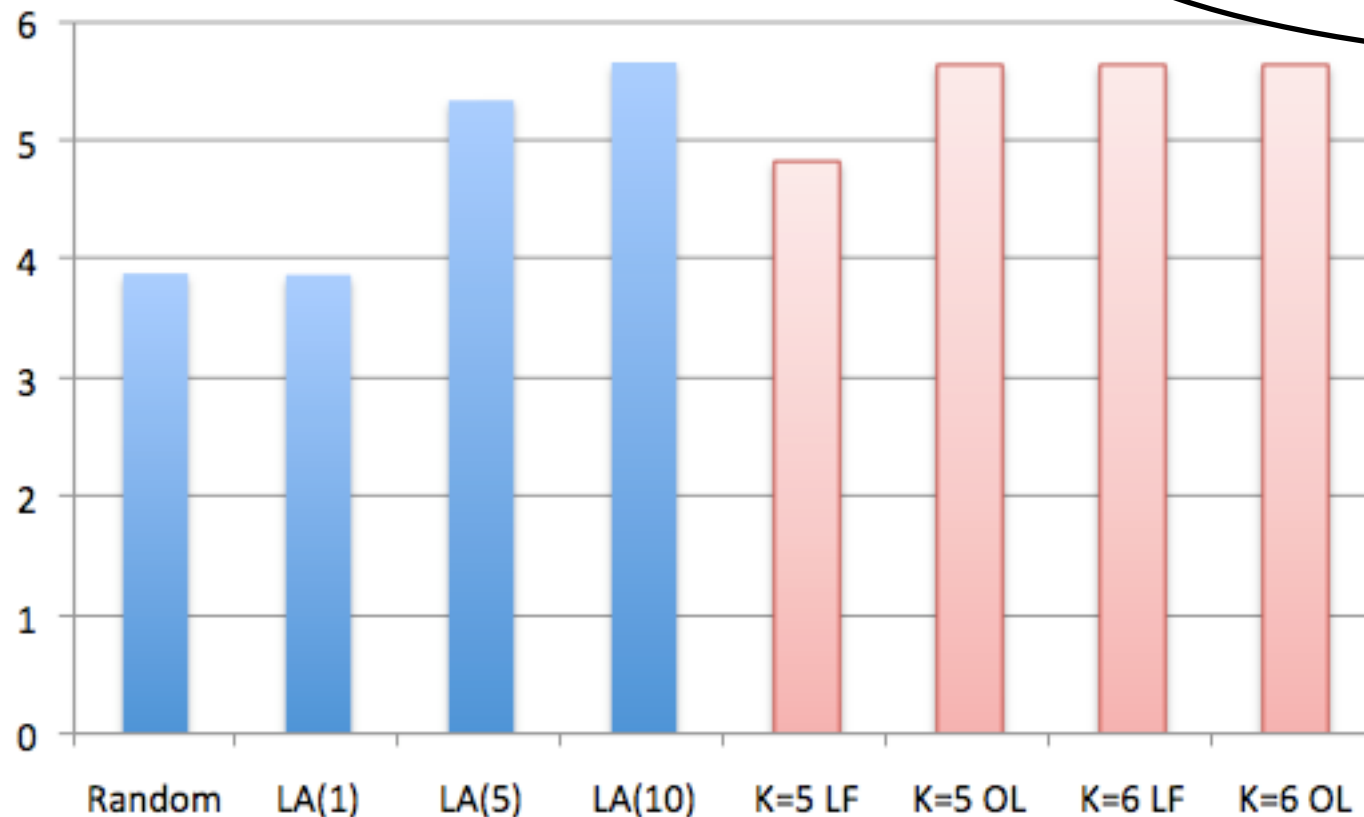


IN THREE CASES, WE  
FIND AS GOOD POLICIES  
AS THE LOOKAHEAD  
POLICY OF DEPTH 10

THE "K=6 LF" DISCOVERED  
FORMULA IS:  $F^* = (-y - v)a$

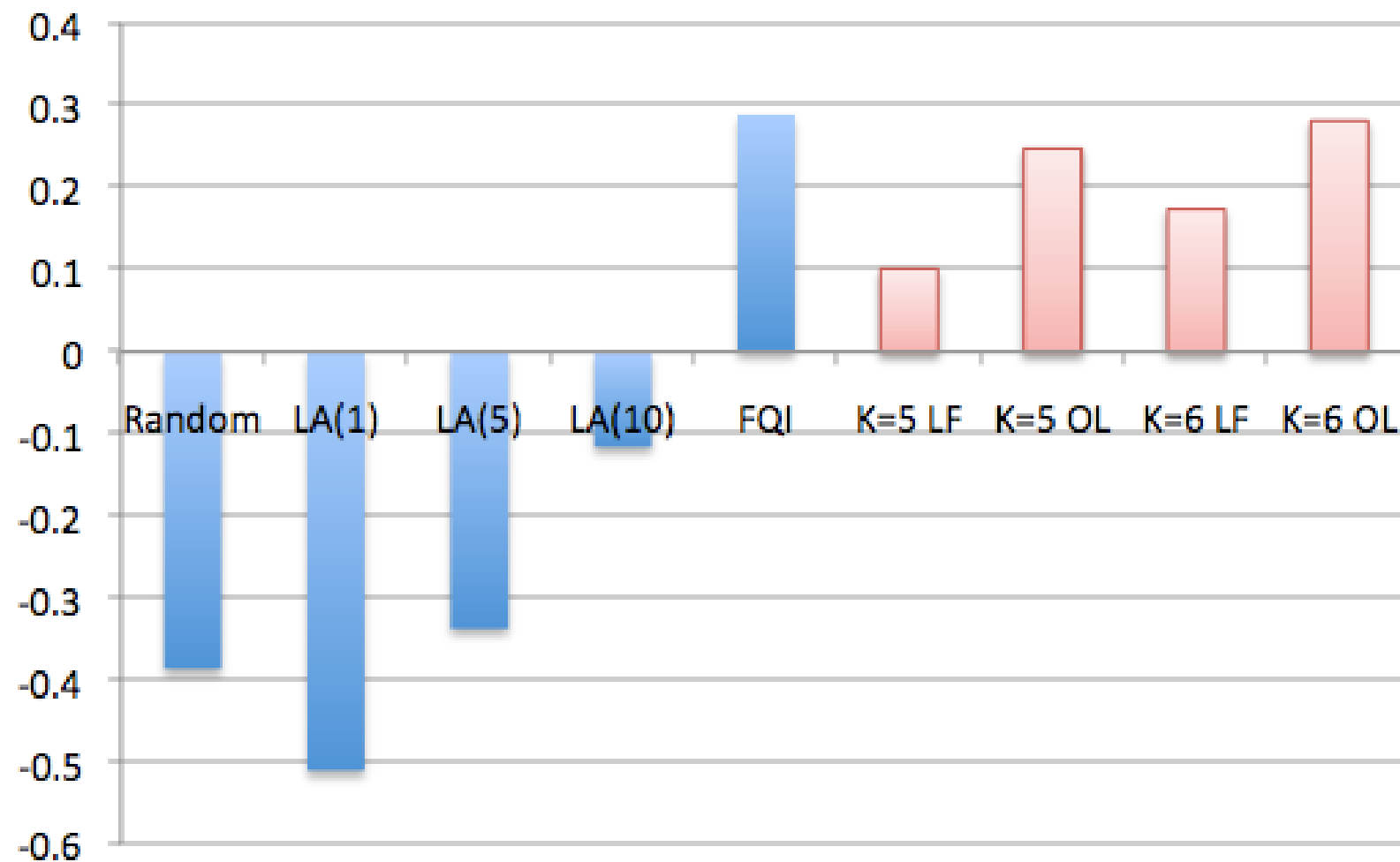
WHERE  $y$  AND  $v$  ARE THE  
STATE VARIABLES AND  $a$  IS THE  
ACTION VARIABLE

SINCE THE ACTION  $a$  IS EITHER -1 OR +1,  
THE POLICY CAN BE TERMED AS:  
"CHOOSE -1 WHEN  $y > -v$  AND +1 OTHERWISE"



OUR BEST INTERPRETABLE  
POLICY PERFORMS NEARLY AS  
WELL AS FQI (0.282 VS 0.29)

THE K=6 OL  
DISCOVERED  
FORMULA IS:  
$$r = \frac{1}{\max(p', s')}$$

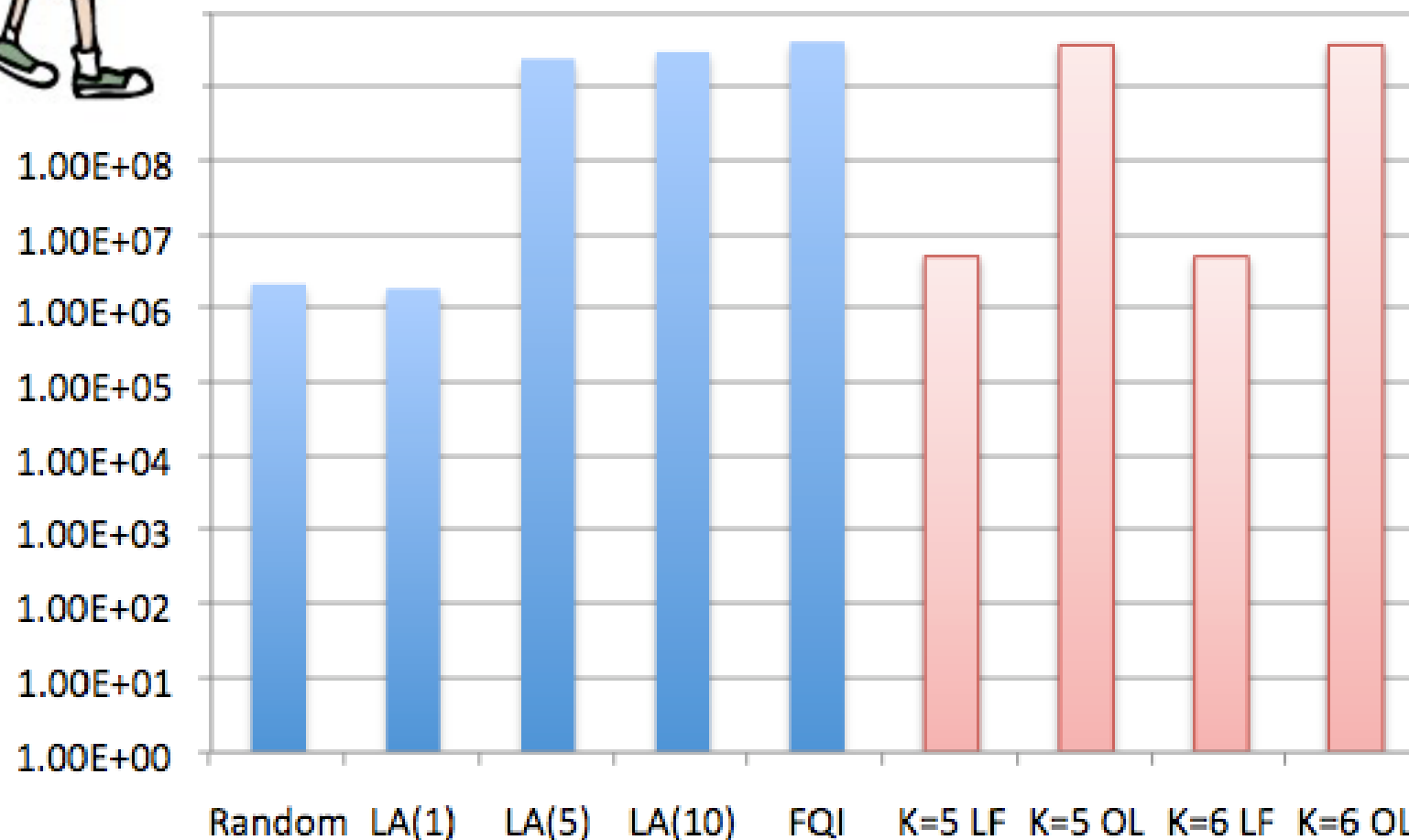




AGAINST ALL OUR EXPECTATIONS, WE DISCOVERED A SIMPLE GOOD PERFORMING POLICY FOR THE HIV BENCHMARK:

$$F^* = \frac{\sqrt{E'}}{\ln(T_1')}$$

THIS POLICY OBTAINS AN AVERAGE RETURN SLIGHTLY BELOW FQI (3.744E9 VS 4.16E9) HOWEVER, THE POLICY FOUND BY FQI REQUIRED 2GB OF MEMORY TO BE STORED



E IS THE CONCENTRATION OF CYTOTOXIC T-LYMPHOCYTES (IN CELLS/ML) AND T1 IS THE CONCENTRATION OF NON-INFECTED CD4 T-LYMPHOCYTES (IN CELLS/ML)



HERE IS AN OVERVIEW  
OF THE FORMULAS WE  
FOUND WITH K=6

### LOOKAHEAD-FREE FORMULAS

$$\begin{aligned} &(-y - v)a \\ &u/(x - \sqrt{5}) \\ &u(\sqrt{7} - s) \\ &\max(\dot{\theta}_2/u, \sqrt{2}) \\ &\psi\dot{\theta} - |d| \\ &1/(\epsilon_1 - \frac{T_2}{T_1^*}) \end{aligned}$$

### ONE-STEP LOOKAHEAD FORMULAS

$$\begin{aligned} &y' - |y + v'| \\ &u/(x - \sqrt{7}) \\ &r - \frac{1}{\max(p', s')} \\ &\dot{\theta}_2|\dot{\theta}'_2| - u \\ &1/(7 - \dot{\theta}'/\dot{\omega}') \\ &\sqrt{E'}/\ln(T'_1) \end{aligned}$$



ON ALL BENCHMARKS, WE  
FOUND AT LEAST ONE SIMPLE  
POLICY PERFORMING BETTER THAN  
A POLICY USING A FULL LOOKAHEAD  
TREE OF DEPTH 10

IT SEEMS LIKE MANY COMPLEX  
SEQUENTIAL DECISION PROBLEMS  
ADMIT SURPRISINGLY SIMPLE  
SOLUTIONS !





## CONCLUSIONS





WE INTRODUCED THE ISSUE OF  
INTERPRETABILITY IN  
REINFORCEMENT LEARNING

ALTHOUGH IT IS AN OLD TOPIC IN  
SUPERVISED LEARNING, THIS ISSUE HAS  
RECEIVED VERY FEW ATTENTION IN THE FIELD  
OF RL

WE FOCUSED ON THE  
CLASS OF INDEX-BASED  
POLICIES DESCRIBED BY  
SMALL INDEX FORMULAS

WE INTRODUCED A  
DIRECT POLIC SEARCH  
SCHEME BASED ON  
MULTI-ARMED BANDITS



WE TESTED OUR APPROACH ON  
SIX BENCHMARK PROBLEMS AND  
DISCOVERED SIMPLE  
- REASONABLY EFFICIENT -  
POLICIES IN EACH CASE





WE EXPERIENCED THE INTERPRETABILITY / EFFICIENCY TRADEOFF: MORE INTERPRETABLE POLICIES DO NOT REACH THE PERFORMANCE OF NON-INTERPRETABLE ONES

HOWEVER, ON ALL THE SIX DOMAINS, THE PERFORMANCE WE OBTAIN IS STILL REASONABLE AND OUR APPROACH ENABLES TO OBTAIN USEFUL INSIGHTS ON THE PROBLEMS

THIS WORK IS ONLY ONE EXAMPLE OF INTERPRETABLE RL. MANY OTHER APPROACHES COULD BE INVESTIGATED (E.G. BASED ON DECISION TREES OR DECISION GRAPHS)





THANK YOU FOR YOUR ATTENTION !



## SEE ALSO ...



[Monte Carlo Search Algorithm Discovery for One Player Games](#). Francis Maes, David Lupien St-Pierre and Damien Ernst.

[Automatic discovery of ranking formulas for playing with multi-armed bandits](#). Francis Maes, Louis Wehenkel and Damien Ernst. In 9th European workshop on reinforcement learning (EWRL'11), Athens, Greece, September 2011.

[Meta-Learning of Exploration/Exploitation Strategies: The Multi-Armed Bandit Case](#). Francis Maes, Louis Wehenkel and Damien Ernst.

[Learning exploration/exploitation strategies for single trajectory reinforcement learning](#). Michael Castronovo, Francis Maes, Raphael Fonteneau and Damien Ernst. In 10th European Workshop on Reinforcement Learning (EWRL'12), Edinburgh, Scotland, June 2012.