# Expérimentations et Réflexions autour de l'Apprentissage par Renforcement Développemental

Sém. Apprentissage Développemental, Lyon.
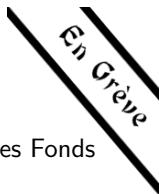
Alain Dutech

Equipe MAIA - LORIA - INRIA
Nancy, France
Web : http://maia.loria.fr
Mail : Alain.Dutech@loria.fr

16 déc 2014

En Grève

# En grève !? Mais pourquoi ??

*En Grève*

- ▶ Le Ministère ne finance pas assez les Universités
- ⤳ LRU + RCE ⇒ A la charge de l'Université de trouver des Fonds

- ▶ UL (comme ailleurs) : Le Président n'assume pas :o(
- ⤳ Obei au Ministère et Management Autocratique
- ⤳ "Fait redescendre" sur pôles et collégium
- ⤳ Gel des postes, Refus solution alternative
- ⤳ "Fusion" ou Mutualisation "forcée"
- ⤳ (à terme ... droits d'inscriptions)

Difficilement mais sûrement, lutte et solidarité se mettent en place.
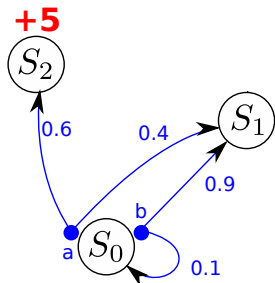
## Outline

En Grève

### Experiment

(3)

- ▶ **Context : RL, Developmental approach**
- ▶ Learning Architectures
- ▶ Robotic experiments

### Discussion

- ▶ What is "Developmental" ?
- ▶ Problems and Questions

# Reinforcement Learning

[Puterman, 1994], [Sutton and Barto, 1998], [Groupe PDMIA, 2008], ...

En Grève

4



▶ States $\mathcal{S}$, Actions $\mathcal{A}$, probabilistic transitions.

▶ $E_{s,a\sim\pi}\left[\sum_{t=1}^{\infty}\gamma^t r_t | s_0 = s, a_0 = a\right]$

▶ Find the optimal policy $\pi$.

⇝ Action a or b in $S_0$ ?

# Reinforcement Learning

[Puterman, 1994], [Sutton and Barto, 1998], [Groupe PDMIA, 2008], ...

En Grève

④



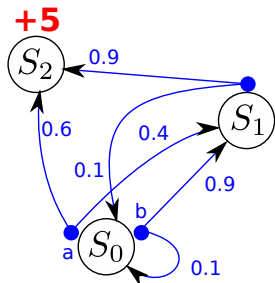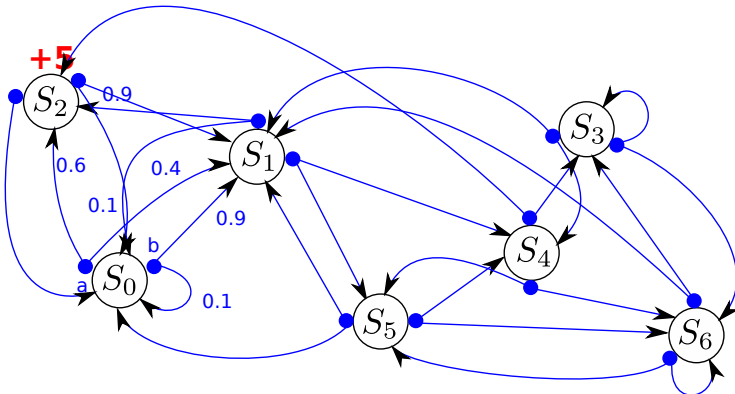- States $\mathcal{S}$, Actions $\mathcal{A}$, probabilistic transitions.

- $E_{s,a \sim \pi} \left[ \sum_{t=1}^{\infty} \gamma^t r_t | s_0 = s, a_0 = a \right]$

- Find the optimal policy $\pi$.
- $\rightsquigarrow$ Action a or b in $S_0$ ?

# Reinforcement Learning

[Puterman, 1994], [Sutton and Barto, 1998], [Groupe PDMIA, 2008], ...



Compute directly the **optimal** value function (as a solution to):

$$Q^*(s, a) = \mathbf{r(s, a)} + \gamma \sum_{s' \in \mathcal{S}} \mathbf{p(s'|a, s)} \max_{a' \in \mathcal{A}}[Q^*(s', a')]$$

Context
○●○○○

DevRL
○○○

XP
○○○○○○○

Conclusion
○○○○○○○

Références
○○

## Q-Learning [Watkins, 1989]

En Grève

States $\mathcal{S}$, Actions $\mathcal{A}$, probabilistic transitions.

⑤

Stochastic approximation of the Q-values of an optimal policy.

1. For a given state $s$
2. Choose and apply exploratory action $a$
3. Environment gives back new state $s'$ and reward $r$
4. Update

$$Q(s, a) \quad \longleftarrow \quad Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') - Q(s, a)]$$

5. Goto (1) with $s \longleftarrow s'$

# RL problems with huge state $\times$ action space

Typical with robots for example.

*En Grève*

$$Q(s, a) \quad \longleftarrow \quad Q(s, a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') - Q(s, a)] \qquad (6)$$

1. Continuous environment. RL easier with discrete states *and actions*.
   $\rightsquigarrow$ approximation
2. Costly experiences. Time, energy to try $(s, a)$ vs algorithms with huge iteration needs.
   $\rightsquigarrow$ Re-use.
3. Harmfull experiences. Robot destruction ?
   $\rightsquigarrow$ Special low-level behavior.
4. Sparse Reward. Non-zero reward is difficult to get.
   $\rightsquigarrow$ Eligibility traces.
5. Rich environment. Too many "states" to consider.
   $\rightsquigarrow$ Factorization, aggregation
6. Partial observability. No guarantee for learning.
   $\rightsquigarrow$ (State-extension to an MDP), incremental.

# RL problems with huge state × action space

Typical with robots for example.

$$Q(s,a) \quad \longleftarrow \quad Q(s,a) + \alpha[r + \gamma \max_{a' \in \mathcal{A}} Q(s',a') - Q(s,a)] \qquad \text{\textcircled{6}}$$

1. **Continuous environment**. RL easier with discrete states *and actions*.
   ↝ approximation
2. Costly experiences. Time, energy to try $(s, a)$ vs algorithms with huge iteration needs.
   ↝ Re-use.
3. Harmfull experiences. Robot destruction ?
   ↝ Special low-level behavior.
4. **Sparse Reward**. Non-zero reward is difficult to get.
   ↝ Eligibility traces.
5. **Rich environment**. Too many "states" to consider.
   ↝ Factorization, aggregation
6. Partial observability. No guarantee for learning.
   ↝ (State-extension to an MDP), incremental.

## A developmental point of view

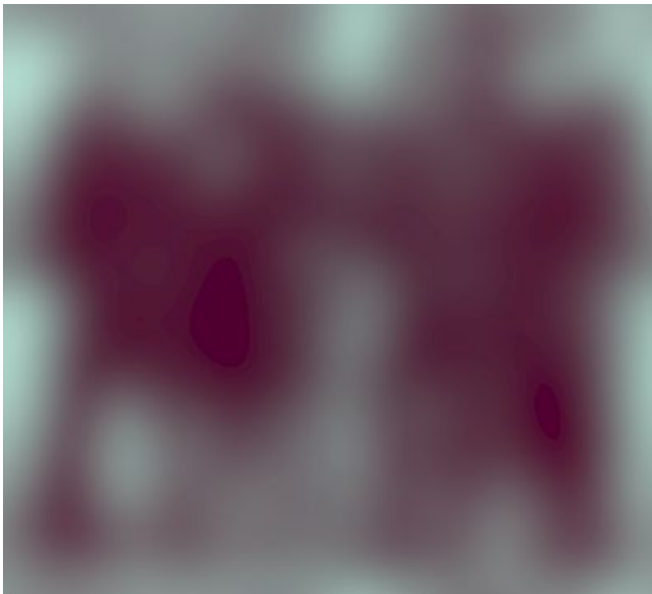# A developmental point of view

Context
00000

DevRL
000

XP
0000000

Conclusion
0000000

Références
00

## A developmental point of view

# Developmental Reinforcement Learning

⤳ **Ease the exploration of state** × **action space**

*En Grève*

⑧

## Concept

Robot's task, perceptual and motor skills **increase** when learned behavior becomes **more efficient**.

(Only one aspect of Developmental Robotic, see [Lungarella et al., 2003]).

# Developmental Reinforcement Learning

⇝ **Ease the exploration of state × action space**

*En Grève*

⑧

**Concept**

Robot's task, perceptual and motor skills **increase** when learned behavior becomes **more efficient**.
(Only one aspect of Developmental Robotic, see [Lungarella et al., 2003]).

New problems arise:

- (Learn approximation of $Q(s, a)$).
- How to deal with increased number of actions ?
- How to deal with increased number of states ?
- How to changes goals, sometimes drastically ?
- How to "transfer" learned behavior to more complex tasks ?

## Outline

En Grève

(9)

### Experiment

- ▶ Context : RL, Developmental approach
- ▶ **Learning Architecture**
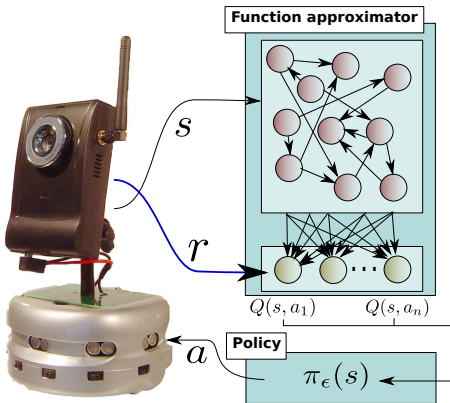- ▶ Robotic experiments

### Discussion

- ▶ What is "Developmental" ?
- ▶ Problems and Questions

## Function approximation
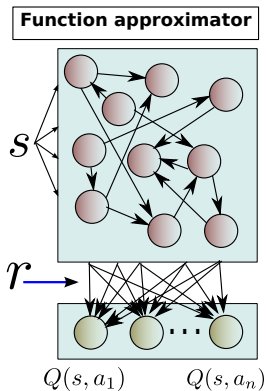
Dynamic Self-Organizing Map and a linear readout



"Supervised learning" using Bellman error:
$$\Delta Q_{\mathrm{NET}}(s, a) = \alpha[r + \gamma \max_{a'} Q_{\mathrm{NET}}(s', a') - Q_{\mathrm{NET}}(s, a)]$$

Context
00000

DevRL
0●0

XP
0000000

Conclusion
0000000

Références
00

# Dynamic Self-Organizing Maps (DSOM)

[Rougier and Boniface, 2011]

En Grève

(11)

**Function approximator**

$s$

$r$

$Q(s, a_1)$ $\quad$ $Q(s, a_n)$

- ▶ Present input $v$
- ▶ Winner neuron $w_w$

$$s \longleftarrow \operatorname*{argmin}_{i \in \mathcal{N}}(||v - w_i||)$$

- ▶ Learn "lateral" connexions towards $v$

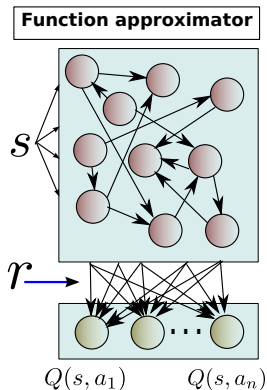$$\delta w_i \longleftarrow \epsilon_D ||v - w_i|| h_\eta(i, w_s, v)(v - w_i)$$

- ▶ Readout

$$\mathrm{out}_i \longleftarrow \sum \omega_i . w_i$$

- ▶ Learning readout weights

$$\delta \omega_i \longleftarrow -\epsilon_L . \mathrm{error} . w_i$$

Context
00000

DevRL
00●

XP
0000000

Conclusion
0000000

Références
00

# Grow sensori-motor space

En Grève



**Function approximator**

$s$

$r$

$Q(s, a_1) \qquad Q(s, a_n)$

### Increase input dimension

(12)

- ▶ Start with *n* input neurons
  (where *n* is the maximum input dimension)
- ▶ At start, some input are cloned.
- ▶ Then, inputs are discriminated

### Increase nb of actions

- ▶ Add an output neuron
- ▶ Init weights
  ⤳ random, copy
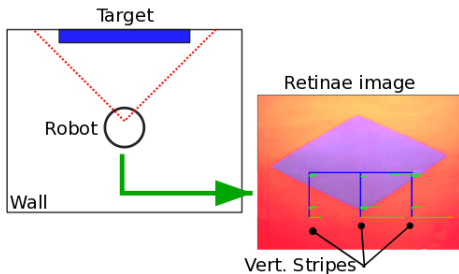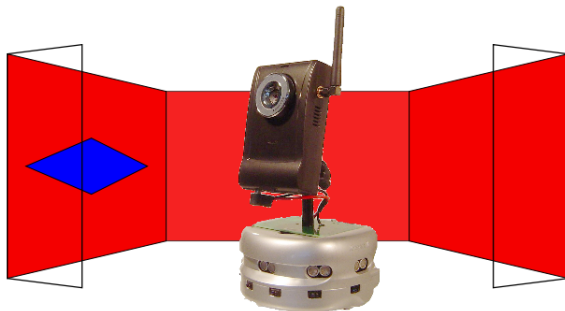
## Outline

En Grève

(13)

### Experiment

- ► Context : RL, Developmental approach
- ► Learning Architectures
- ► **Robotic experiments**

### Discussion

- ► What is "Developmental" ?
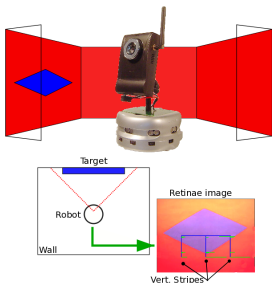- ► Problems and Questions

# Robotic setting



En Grève

(14)

Target

Robot

Wall

Retinae image

Vert. Stripes

Context
○○○○○

DevRL
○○○

XP
○●○○○○○

Conclusion
○○○○○○○

Références
○○

# Task details

En Grève

⑮



- One Khepera3 robot
- **Continuous** perception
  ⤳ camera and blue stripes sensor
- Actions are **discrete**
  - 3 actions: Stop, Left, Right
  - 5 actions: Stop, Left, Right, slowLeft, slowRight
- Reward linked to the blue levels in sensor
  ex: $0.7B1 + B2 + 0.7B3 \geq 1.2$
- Real 2.000 samples (used and re-used) or Simulator

# Learning Scenarii

Learning

- ▶ Generate ($\epsilon_p$) or use transition ($s, a : s', r$)
- ▶ unsupervised DSOM : $\epsilon_D$, elasticity $\eta$
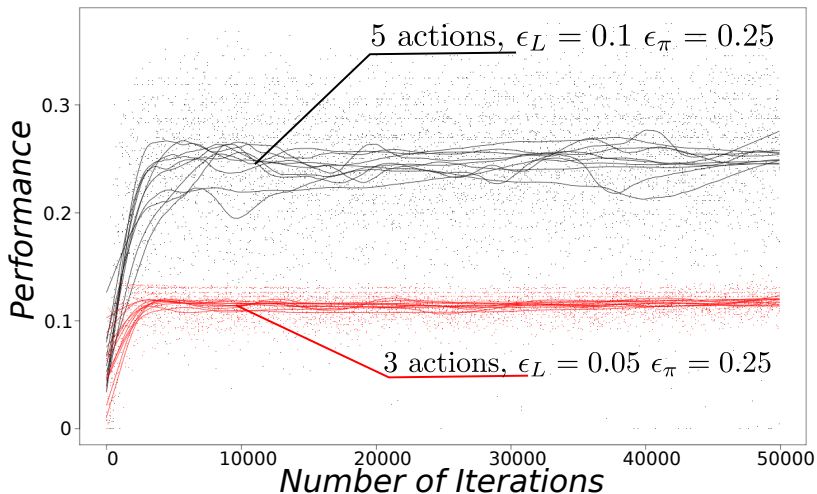- ▶ supervised LIN : $\epsilon_L$, $\alpha$

Evaluation

- ▶ Periodically
- ▶ No learning
- ▶ Mean reward : $\frac{1}{N} \sum_{t=1}^{N} r_t$

Visual evaluation of policy

Context
○○○○○

DevRL
○○○

XP
○○○○●○○○

Conclusion
○○○○○○○

Références
○○

# Better performances with 5 actions

5 actions, $\epsilon_L = 0.1$ $\epsilon_\pi = 0.25$

3 actions, $\epsilon_L = 0.05$ $\epsilon_\pi = 0.25$

*Performance* vs *Number of Iterations*

Context
○○○○○

DevRL
○○○

XP
○○○○●○○

Conclusion
○○○○○○○

Références
○○

# Learned policy, 5 actions



En Grève

18

Context
ooooo
DevRL
ooo
XP
ooooooeo
Conclusion
ooooooo
Références
oo

# Variability vs. Speed

5 actions, $\epsilon_L = 0.1$, $\epsilon_\pi = 0.25$

5 actions, $\epsilon_L = 0.002$, $\epsilon_\pi = 0.25$

Context
○○○○○
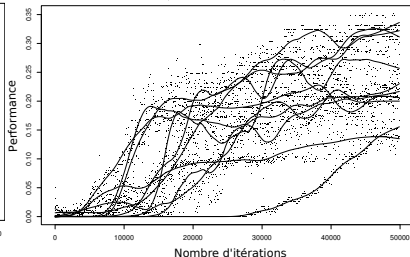
DevRL
○○○

XP
○○○○○○○●

Conclusion
○○○○○○○

Références
○○

# Developmental learning, 3 then 5 actions

En Grève

Direct: 5 actions / DevRL: 2 actions added $N_B = 5000$ ($\epsilon_L$ 0.1 to 0.01)

## Outline

En Grève

21

### Experiment

- ▶ Context : RL, Developmental approach
- ▶ Learning Architectures
- ▶ Robotic experiments

### Discussion

- ▶ **What is "Developmental" ?**
- ▶ Problems and Questions

Context
00000

DevRL
000

XP
0000000

Conclusion
●000000

Références
00

# What is "Developmental" ?

En Grève

## Autonomous All-Life-Long Learning Agent

22

⤳ Interactions Body-Brain-Environment

⤳ Coupling Development-Learning

▶ Agent Development

▶ No "exogen" intervention

▶ Build on Previous Acquired "Behaviors" :

# What is "Developmental" ?

En Grève

## Autonomous All-Life-Long Learning Agent

22

⇝ Interactions Body-Brain-Environment

⇝ Coupling Development-Learning

- ▶ Agent Development
  - ▶ "Only" actions, (perceptions)
  - ▶ Body !!
- ▶ No "exogen" intervention
  - ▶ Intrinsic Motivations vs "external" reward signal !
- ▶ Build on Previous Acquired "Behaviors" :
  - ▶ What starting substrate ?
  - ▶ "Knowledge" Transfer

# Facets of Developmental Robotics [Lungarella et al., 2003]

23

1. Development is an incremental process
2. Development as a set of constraints
3. Development as a self-organizing process
4. Degrees of freedom and motor activity
5. Self-exploratory activity
6. Spontaneous activity
7. Anticipatory movements and early abilities
8. Categorization and sensorimotor co-ordination
9. Neuromodulation, values and neural plasticity
10. Social interaction

## Outline

*En Grève*

(24)

### Experiment

- ▶ Context : RL, Developmental approach
- ▶ Learning Architectures
- ▶ Robotic experiments

### Discussion

- ▶ What is "Developmental" ?
- ▶ **Problems and Questions**

Context
ooooo

DevRL
ooo

XP
ooooooo

Conclusion
oooo●ooo

Références
oo

# Problems and Question

En Grève

25

- ▶ Body !!
  - ▶ Simulated vs Real
  - ▶ Richness (Sensors, DOF), but very poor compared to "us"
  - ▶ "Growing"   Maturation
- ▶ Architecture, Learning MechanismS
- ▶ "Sum" of Individually Tested Mechanisms ??
- ▶ Multiple Time Scales : individual and generations
  - ▶ How to care and interact with agents during all that time?
- ▶ Changing Motivations
  - ▶ How to create new motivations ?
  - ▶ Priority of motivations ?

### Main Question

"Just" a question of "right Integration" or Still Lack Brook's "Juice" [**?**, **?**] ?

## Focussing on the Experiment

En Grève

26

- ▶ Limits of Reinforcement Learning in "low level" Robotics
- ▶ Original(s) architecture for Developmental approach
  - ▶ Dynamic Self-Organizing Map : online life-long learning
  - ▶ Linear regression readout
  - ▶ Compatible with growing actions (and perceptions)

- ▶ NEED MORE RESULTS ... (especially for Perceptions)

- ▶ Memorize (replay) useful transitions, sequence of tasks...

## Short-term and long-term work

En Grève

27

- ▶ scenario for Perception growing.
- ▶ parameters influence, especially on stability

- ▶ more complexe setting ⤳ different tasks to learn

- ▶ select and memorize useful transitions

- ▶ recurrent DSOM (???) for sequence learning in non-Markov (??)
⤳ PhD of Matthieu ZIMMER...

Context
○○○○○

DevRL
○○○

XP
○○○○○○○

Conclusion
○○○○○○●

Références
○○

En Grève

28

**A vous :o)**

# Références I

📄 Groupe PDMIA (2008).
*Processus Décisionnels de Markov en Intelligence Artificielle. (Edité par Olivier Buffet et Olivier Sigaud)*, volume 1 & 2.
Lavoisier - Hermes Science Publications.

📄 Lungarella, M., Metta, G., Pfeifer, R., and Sandini, G. (2003).
Developmental robotics: a survey.
*Connection Science*, 15(4):151–190.

📄 Puterman, M. (1994).
*Markov Decision Processes: discrete stochastic dynamic programming*.
John Wiley & Sons, Inc. New York, NY.

📄 Rougier, N. P. and Boniface, Y. (2011).
Dynamic Self-Organising Map.
*Neurocomputing*, 74(11):1840–1847.

📄 Sutton, R. and Barto, A. (1998).
*Reinforcement Learning*.
Bradford Book, MIT Press, Cambridge, MA.

# Références II

📄 Watkins, C. (1989).
*Learning from delayed rewards.*
PhD thesis, King's College of Cambridge, UK.