



REINFORCEMENT LEARNING INTRODUCTION

Odalric-Ambrym Maillard

DEC. 2021

Inria Scool

TABLE OF CONTENTS

WHAT IS R.L.?

WHAT IS R.L. **Media**

Applications in SCOOL
First concepts
State, Action, Rewards



Robots?



Set of tools?

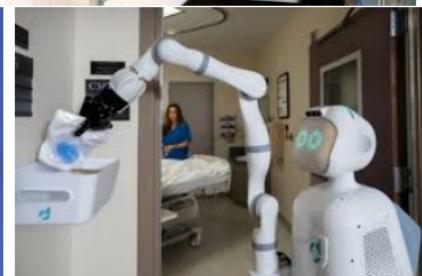
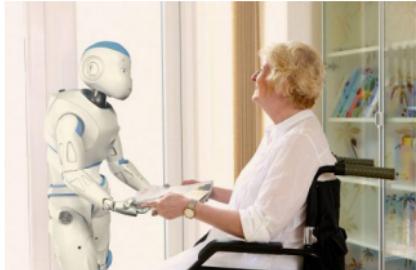


Content Recommendation?

+ Fears: HAL, Terminator, Big-Brother?

ROBOTS EVERYWHERE?

Human bias!



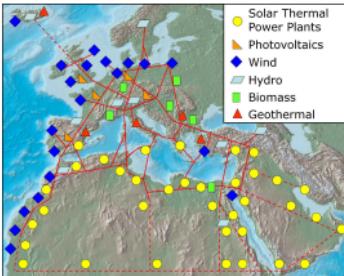
- ▷ +: Repetitive tasks, accurate.
- ▷ -: Expensive, **local** perception, energy, maintenance costs, ethics?

Mimic/replace human expertise vs assisting their limited abilities?

RL: ROBOTS AND GAMES

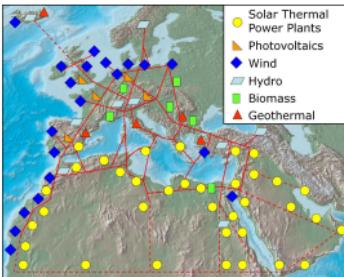


GLOBAL SENSING AND ACTING



- ▷ Making sense of millions of data.
- ▷ Identifying similar contexts.
- ▷ Global testing, (safe) trial-and-error.

GLOBAL SENSING AND ACTING



- ▷ Making sense of millions of data.
- ▷ Identifying similar contexts.
- ▷ Global testing, (safe) trial-and-error.



+10 million users
smartphone app.
Help identify pathogens.
Suggest preventive/curative actions.

Plantix

Go



10^{100} board configurations

Atari



Complex visual input
Very different games

Starcraft

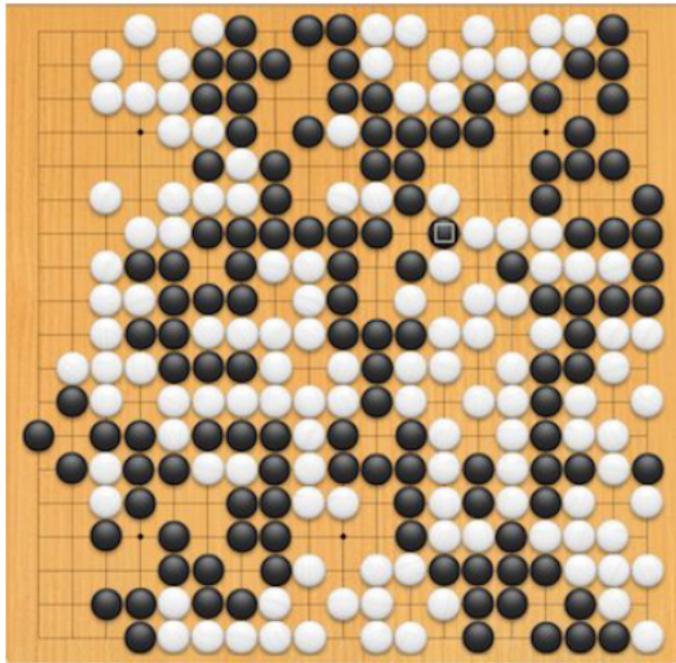
AI becomes grandmaster in 'fiendishly complex' StarCraft II

DeepMind's AlphaStar masters game dubbed 'next grand challenge for AI' in just 44 days

'The challenge was to play like a human': AI takes on the gamers



Complex observations
and Planning



19×19 possible actions, 10^{100} board configurations



Complex visual inputs, very different games

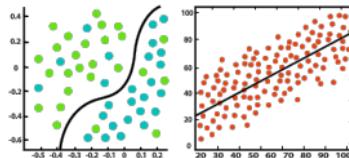


Complex observations and planning

Supervised ("Learn to predict")

- ▷ Classify, Regress

Perception



Unsupervised ("Learn to represent")

- ▷ Cluster, Transform

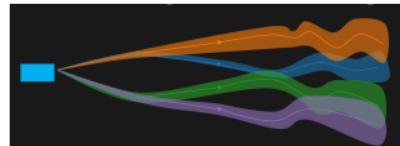
Representations



Reinforced ("Learn to interact")

- ▷ Act, Plan

Decisions



Reinforcement Learning is about automating **Decision** making.

RL builds the **brain**.



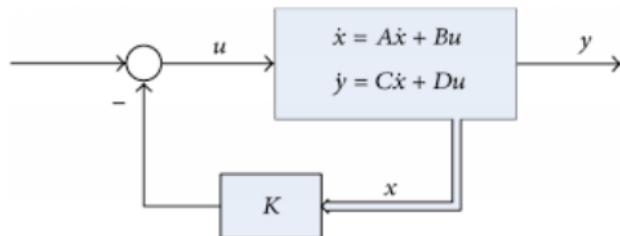
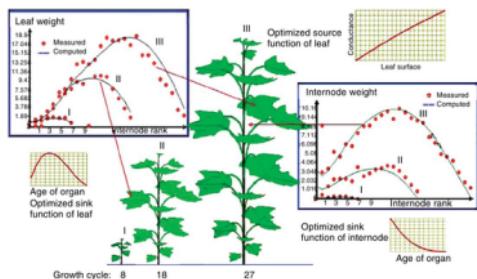
When? **Uncertain** world (black-box, unknown, noisy)
How? Statistics, optimization, control theory.

LEARNING BY EXPERIMENTATION

Reinforcement Learning: The agent does *not* know the dynamics of the system beforehand (Black-box). Learned by **trial and error**.



Control theory: The system's dynamics is perfectly known beforehand (white-box).



WHAT IS R.L.?

Media

Applications in SCool

First concepts

State, Action, Rewards



SUSTAINABLE DECISION MAKING

Sample-efficient, energy-efficient, for **real-world** challenges



SOCIETAL DECISION MAKING

Decision companions, ethics (no human farming!), **collaborative** decisions.

- ▷ Assist vs Replace?
- ▷ Recommend vs Influence ?
- ▷ Crowd-source vs Enslave?
- ▷ Fairness, gender/etc. bias in data.
- ▷ Energetic, Environmental costs (total/per user)?

- ▶ B4H (Bandits For Health): Suggest medical consultation or treatment based on smart meters and many patients.

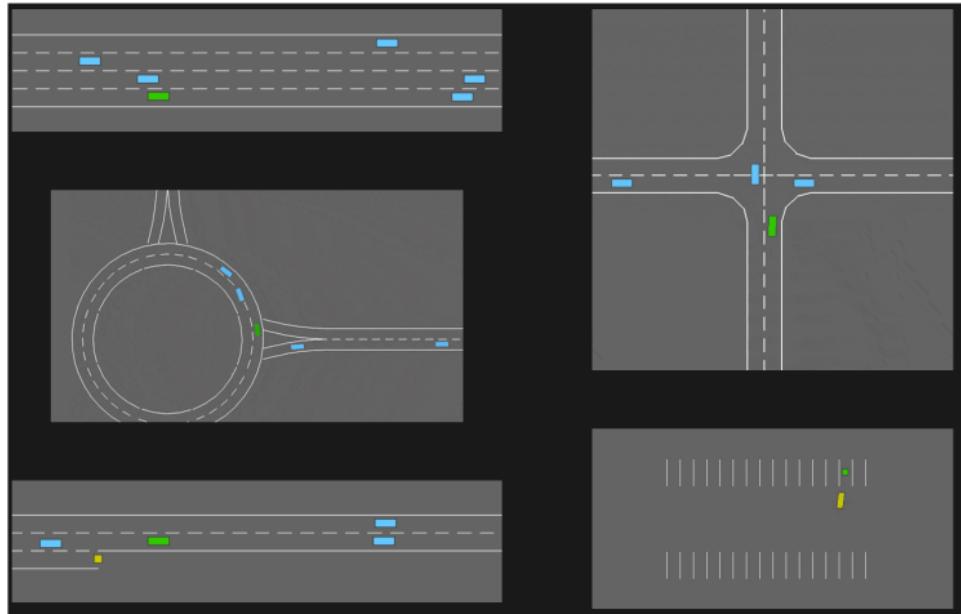


- ▶ What sequence of exercises will maximize the learning progress of this student?



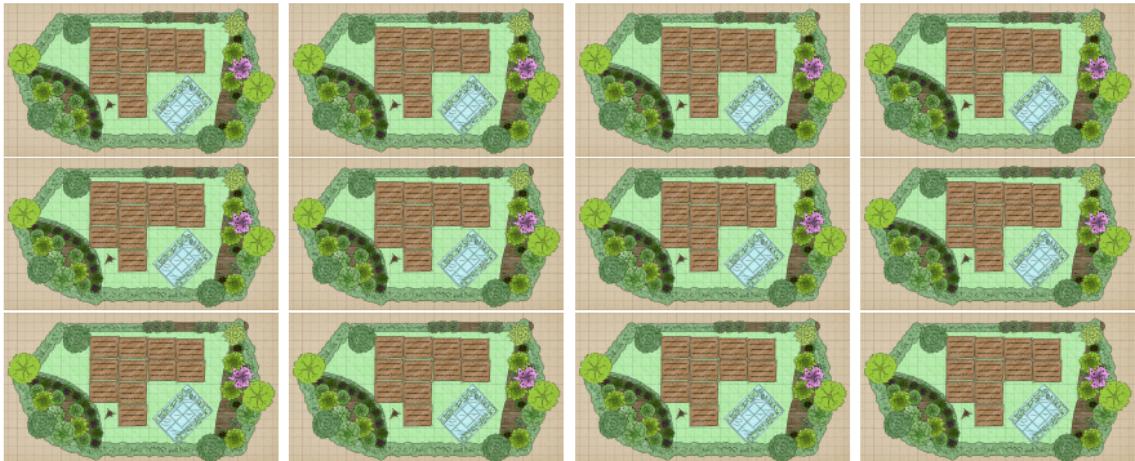
- ▶ pix.fr, lelivrescolaire.fr

- ▶ **Robust** driving in the presence of other vehicles (Renault).



- ▶ Propagation of uncertainty, obstacle avoidance, efficiency vs safety.

- ▶ SR4SG: Sequential Recommendation for Sustainable Gardening.



- ▶ Recommend **good practice** between farms, share/build knowledge.
- ▶ Towards Massively Collaborative Agriculture.
 - ▶ (Re-)discover local practice and expertise.
 - ▶ Experiment/Propagate novel practices (invasive species, diseases, etc.)
 - ▶ Transfer practice between regions (climate evolution).

WHAT IS R.L.
Media
Applications in SCOOT
First concepts
State, Action, Rewards

THE LEARNING GAME



Continuous interaction.

Agent does not know environment: *learning*.

System, called the *Environment*, with which we interact:

- ▶ we get **observations**

System, called the *Environment*, with which we interact:

- ▶ we get **observations**
- ▶ we output **decisions**

System, called the *Environment*, with which we interact:

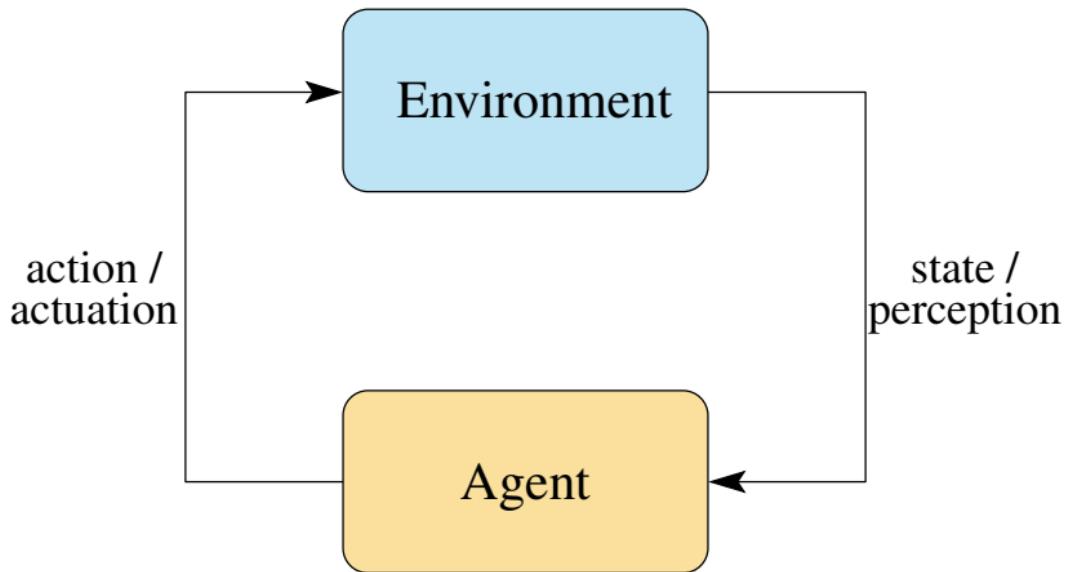
- ▶ we get **observations**,
- ▶ we output **decisions**,
- ▶ we get **new** observations,

System, called the *Environment*, with which we interact:

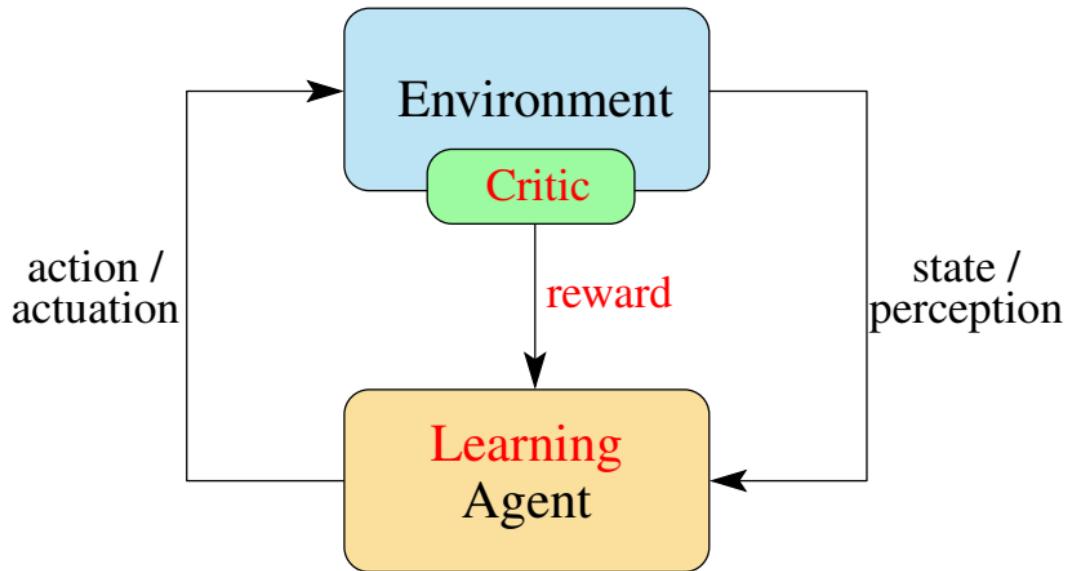
- ▶ we get **observations**,
- ▶ we output **decisions**,
- ▶ we get **new** observations,
- ▶ and so on and so forth.

We receive a signal telling how good/bad are decisions: **reward**

THE REINFORCEMENT LEARNING MODEL



THE REINFORCEMENT LEARNING MODEL



for $t = 1, \dots, n$ **do**

The agent perceives state s_t

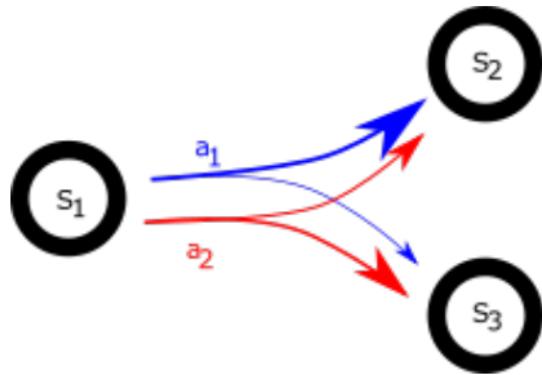
The agent performs action a_t

The environment evolves to s_{t+1}

The agent receives reward r_t

end for

- ▷ *States* s_1, s_2, \dots : "How you describe your system from observations"
- ▷ *Actions* a_1, a_2, \dots : "What you are allowed to do"
- ▷ *Transitions* (unknown): "how the system changes when you choose some action"



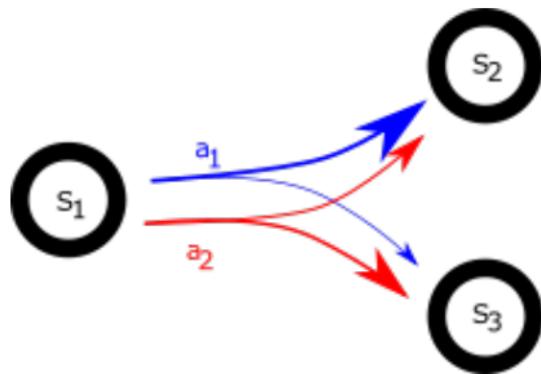
- ▷ *Initial* states: "Where you start from". Final state (optional): "Where you stop"
- ▷ *Reward* (unknown): "How good is an action". Linked to desirable states.

SIMPLIFIED EXAMPLE 1: CLINICAL TRIALS

- Environment is a human
- States*: $s_1 = \text{sick}$, $s_2 = \text{cured}$, $s_3 = \text{uncured}$



- Actions*: $a_1 =$ [image of capsules], $a_2 =$ [image of blister pack], etc.
- Transitions*:



- Initial* state: sick. Final states: cured and uncured.
- Reward*: +1 if cured, 0 else.

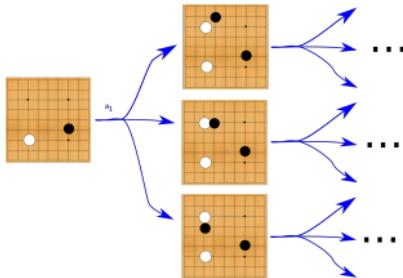
Actions have different success probabilities. (States do not really matter here.)

SIMPLIFIED EXAMPLE 2: GAME OF GO

- Environment is board.



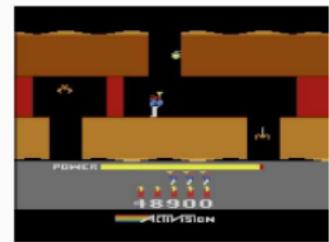
- States*: board (goban) configurations.
- Actions*: all possible moves allowed by game rules.
- Transitions*: e.g. play white at 2×2 , then opponent plays, then ...



- Initial* state: Empty board. *Final* states: end game configurations.
 - Reward*: +1 if win, 0 else.
- Many possible states ($3^{19 \times 19}$), trajectories of (states, actions, next states).

SIMPLIFIED EXAMPLE 3: ATARI ENVIRONMENTS

- Atari games



- Observations:* Screen pixels *Actions:* *Rewards:* Score update
- Transitions:* internal rules of each game (very different).
Complex visual inputs ($256^{640 \times 480}$ possible screens).

GAMES WITH NATURE?



Diverse situations: *Diverse games*
Q: Actions? Observations? Rewards?

- The system is **partially** known (Uncertainty)

DIFFICULTY

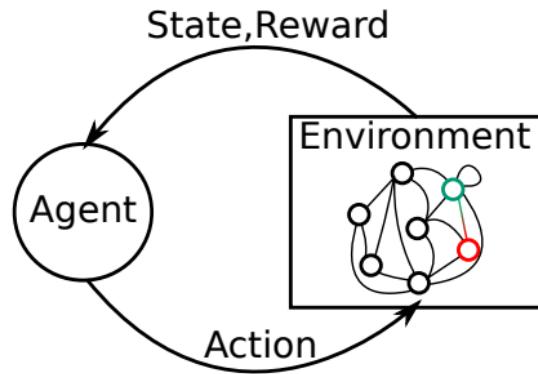
- ▶ The system is **partially** known (Uncertainty)
- ▶ We get to know it thanks to **observations** only. (Statistics)

DIFFICULTY

- ▶ The system is **partially** known (Uncertainty)
- ▶ We get to know it thanks to **observations** only. (Statistics)
- ▶ Observations are collected **actively**. (Algorithm)

- ▶ The system is **partially** known (Uncertainty)
- ▶ We get to know it thanks to **observations** only. (Statistics)
- ▶ Observations are collected **actively**. (Algorithm)
- ▶ Decisions may have **effect** on next possible observations. (Dynamics)

- ▶ The system is **partially** known (Uncertainty)
- ▶ We get to know it thanks to **observations** only. (Statistics)
- ▶ Observations are collected **actively**. (Algorithm)
- ▶ Decisions may have **effect** on next possible observations. (Dynamics)
- ▶ Possibly noisy, high dimensional, structured, with risk-aversion, not so clear reward signal, dependencies, etc. (Models)



Sequential Decision Making under Uncertainty

WHAT IS R.L.?

Media

Applications in SCOOT

First concepts

State, Action, Rewards

Next observation ?

- ▶ May depend on all past history of interaction: unfathomable?

$$\mathbb{P}\left(O_t \middle| \underbrace{o_1, a_1, r_1, \dots, o_{t-1}}_{\text{exponentially many!!}}, a_{t-1}\right)$$

- ▶ State = **Summary** of history

$$= \mathbb{P}(O_t | s_{t-1}, a_{t-1})$$

"All you need to know to perform action"

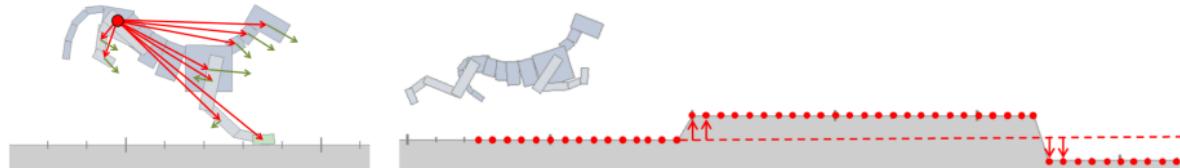
TYPICAL EXAMPLE

State= position

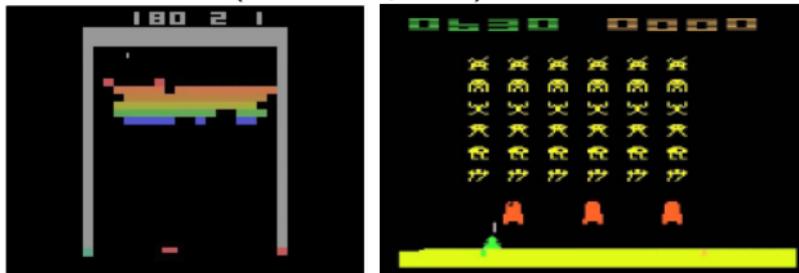
0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
18	19	20	21	22	23
24	25	26	27	28	29
30	31	32	33	34	35

OTHER EXAMPLES

State= relative positions, velocities, and relative heights



State= frame (640x480 pixels) and last few frame differences



ACTIONS

Medical trials



: $\mathcal{A} = \{$



$\}$

Ad-placement industry



: $\mathcal{A} = \{$



$\}$

Plant-health care

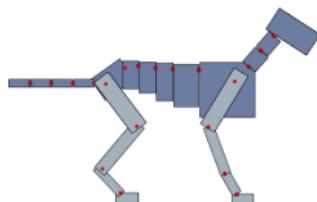


: $\mathcal{A} = \{$



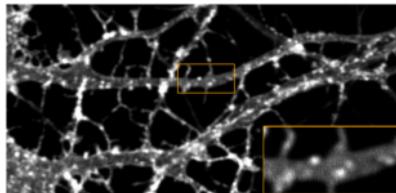
$\}$

Motion-planning:



: $\mathcal{A} = \left\{ \text{accelerations at each joint} \right\}$

High-resolution microscopy:



: $\mathcal{A} = \left\{ \text{exposure time, dosage.} \right\}$

The goal of an agent is to **maximize the rewards** accumulated over time:

$$\sum_{t=1}^T r_t$$

r_t : reward obtained in state s_t when playing action a_t .

Either by:

- ▷ **Finding a policy** that accumulates rewards when played (e.g. Simulator).
- ▷ **Interacting directly** with the real-world (no simulator available).

“The more applied you go, the stronger theory you need”

MERCI

odalricambrym.maillard@inria.fr

odalricambrymmaillard.wordpress.com