



Développeur d'Intelligence Artificielle Appliquée

Cours #7

www.impactia.org

Structure de la formation

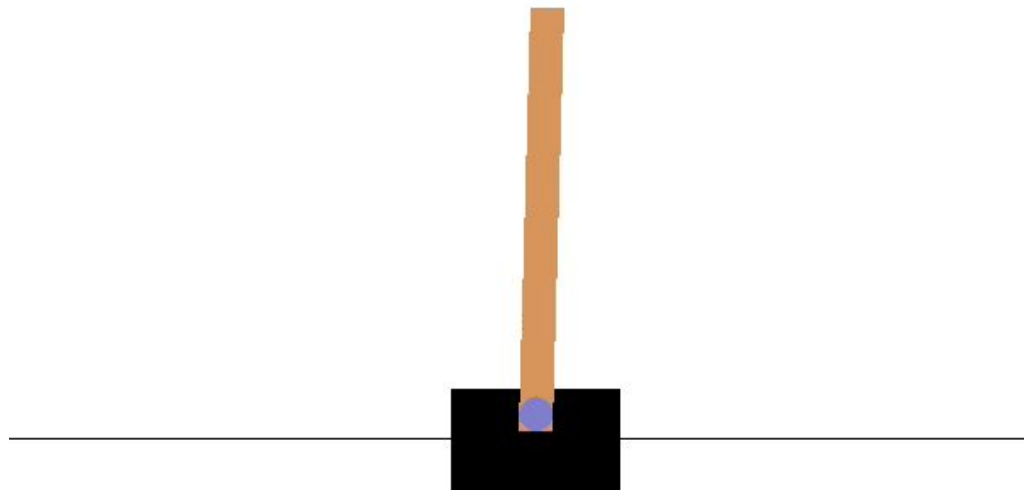
- **#1: Introduction**
- **#2: Vision**
- **#3: Vision**
- **#4: Vision**
- **#5: Renforcement**
- **#6: Renforcement**
- **#7: Renforcement**
- **#8: Langage**
- **#9: Langage**
- **#10: Langage**
- **# 11: Génération d'images**
- **#12: Génération d'images**
- **#13: Génération d'images**
- **#14: Projet**
- **#15: Projet**

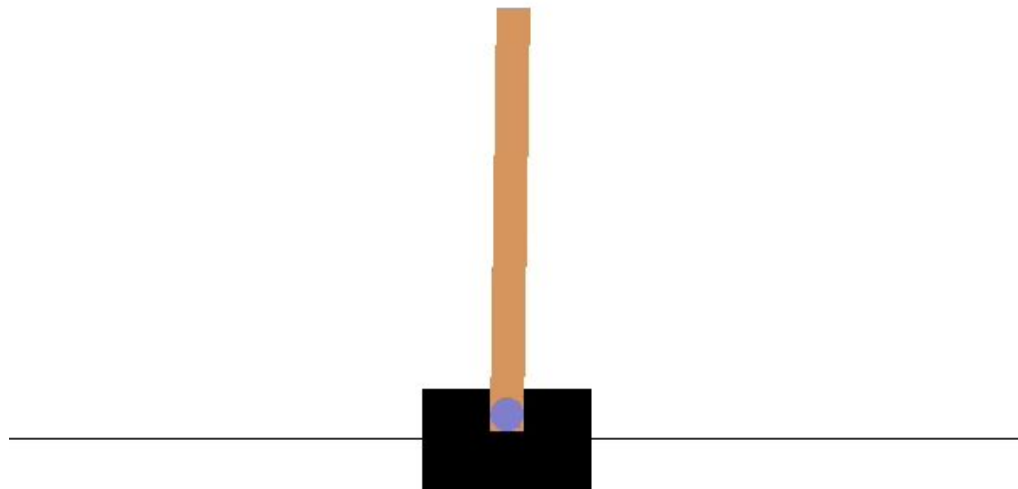
Semaine dernière : Cours #6

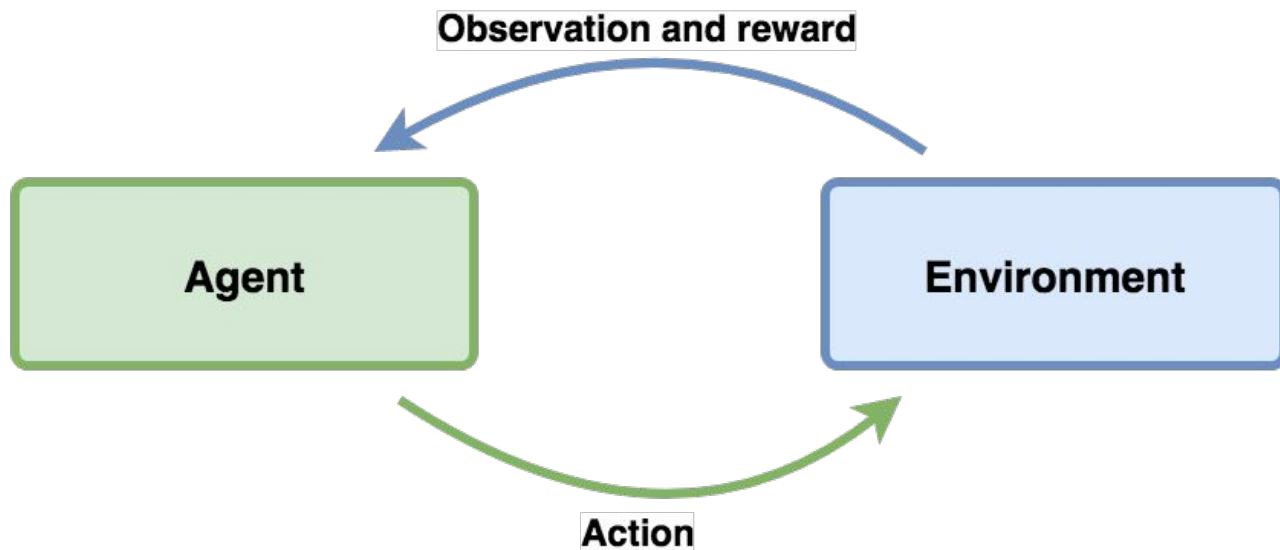
- Théorie
 - L'environnement droneRL
 - Q-learning
- Pratique
 - Entrainement des premiers agents droneRL

Aujourd'hui : Cours #5

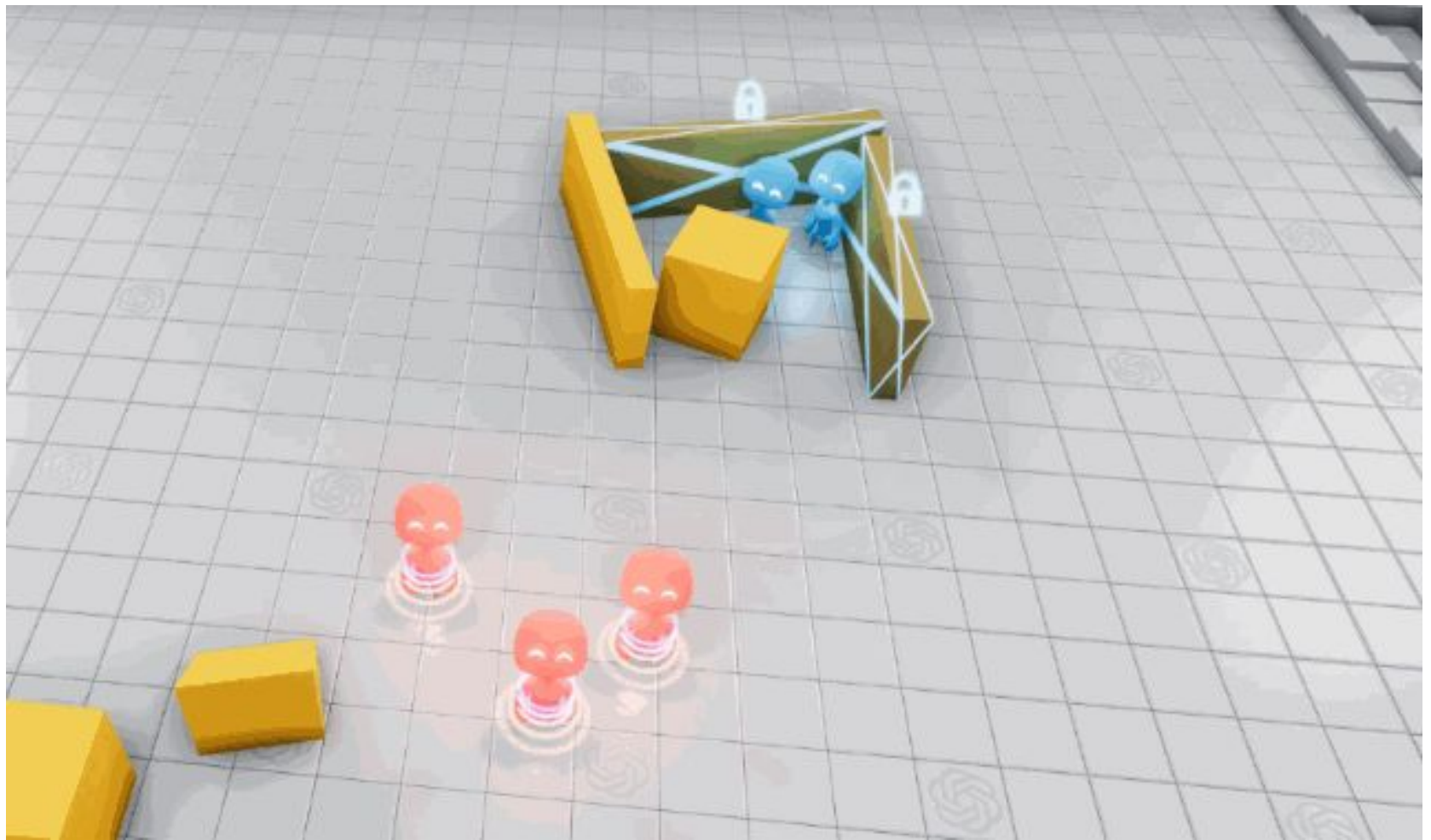
- Théorie
 - Deep Q-Network (DQN)
 - Apprentissage prioritisé
- Pratique
 - Entrainement d'agents droneRL
 - Challenge droneRL ! 💪





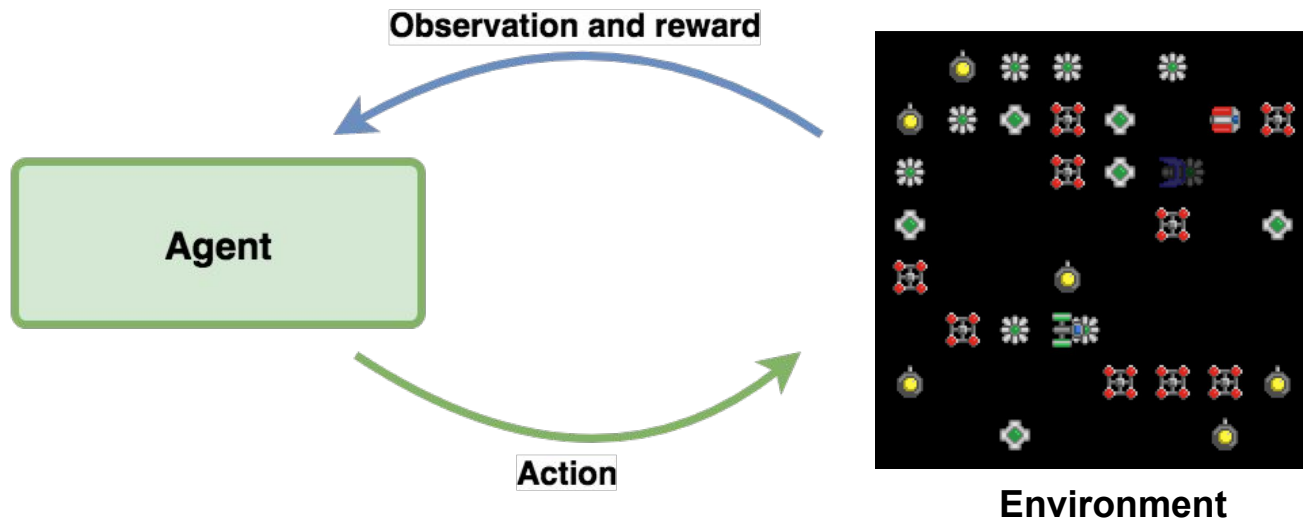


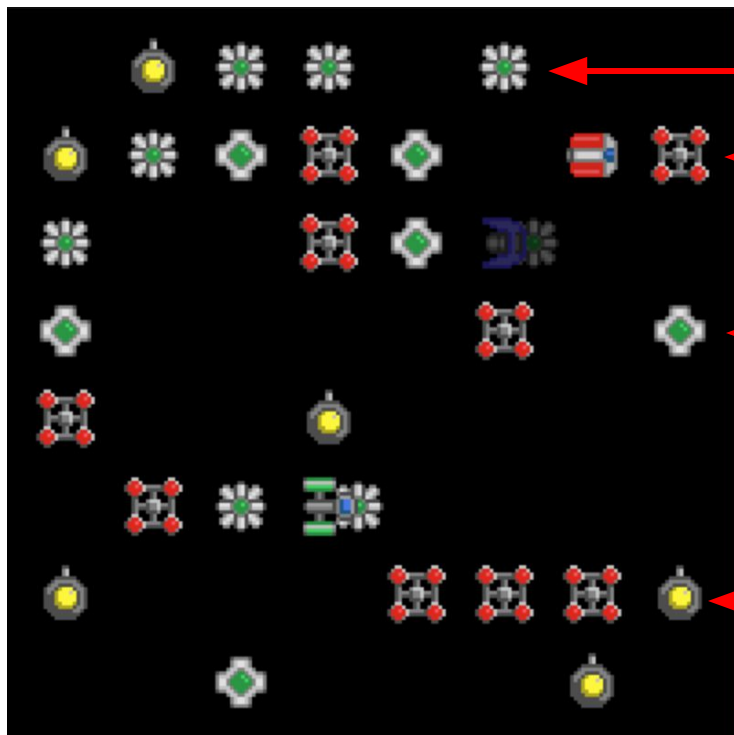
But :
Maximiser la somme des *récompenses* ("rewards")
jusqu'à la fin de l'épisode





L'environnement *DeliveryDrones*





colis

gratte-ciel

point de dépôt

chargeur

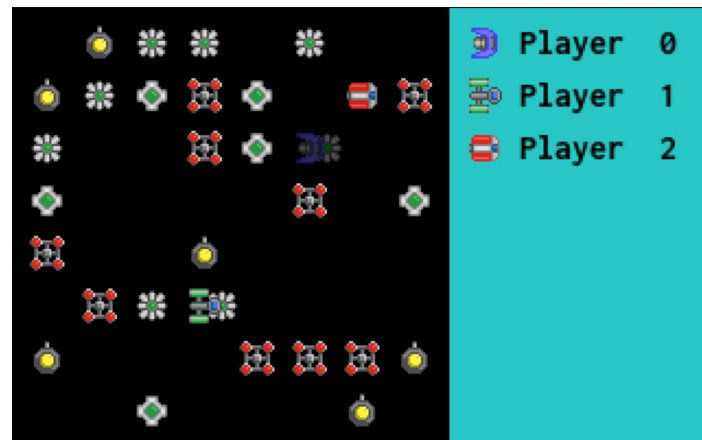
paramètres de l'environnement

```
self.env_params = {  
    'pickup_reward': 0,  
    'delivery_reward': 1,  
    'crash_reward': -1,  
    'charge_reward': -0.1,  
    'discharge': 10,  
    'charge': 20,  
    'drone_density': 0.05,  
    'packages_factor': 3,  
    'dropzones_factor': 2,  
    'stations_factor': 2,  
    'skyscrapers_factor': 3,  
    'rgb_render_rescale': 1.0  
}
```

wrapper

```
env = WindowedGridView(env, radius=3)
```

Paramètres utilisés pour l'évaluation sur Alcrowd !





Comment “apprendre un comportement” ?

Apprendre un comportement



Une information à disposition



5 actions possibles - que faire ?

LEFT

DOWN

RIGHT

UP

STAY

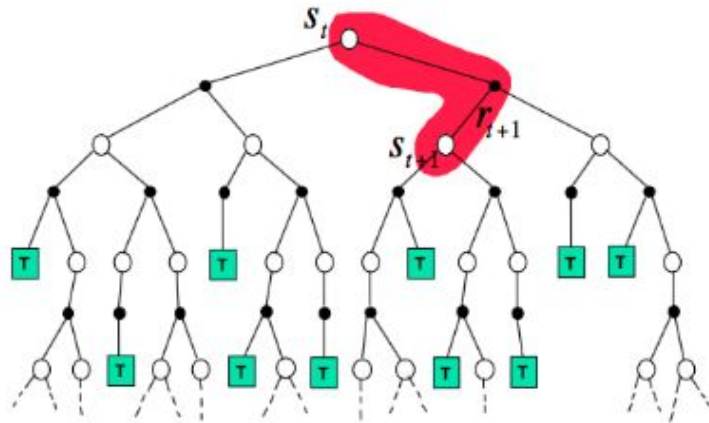
La “solution” d’un environnement

Si on connaît les valeurs de ce tableau,
il est facile de suivre un comportement optimal !

		Actions				
		LEFT	DOWN	RIGHT	UP	STAY
Observation	→	1.7	1.4	2.6	0.86	1.5
	↘	1.3	1.6	2.4	0.78	1.5
	↑	1.9	1.6	1.4	2.5	1.7
	↗	1.4	0.58	1.8	1.6	1.6
	↖	2.2	1.3	0.97	1.6	1.5
	↙	2.5	2	1.3	0.9	1.5
	↓	1.8	2.9	1.7	1.2	2.1
	←	2.6	1.5	2.1	1.2	2.2

← “Q-table”

Comment apprendre la Q-table ?



Q-learning

Actions

LEFT DOWN RIGHT UP STAY

→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2



“Q-table”

Somme attendue des récompenses



Actions

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

Observation



“Q-table”

Somme attendue des récompenses



Actions

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

Observation

On sélectionne l'action qui maximise la somme attendue des récompenses

$$a = \operatorname{argmax}_a Q(s, a)$$



“Q-table”

Somme attendue des récompenses

Actions

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

Observation

“Q-table”

Somme attendue des récompenses

On sélectionne l'action qui maximise la somme attendue des récompenses

$$a = \operatorname{argmax}_a Q(s, a)$$



Actions

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

Observation

“Q-table”

Somme attendue des récompenses

On sélectionne l'action qui maximise la somme attendue des récompenses

$$a = \operatorname{argmax}_a Q(s, a)$$

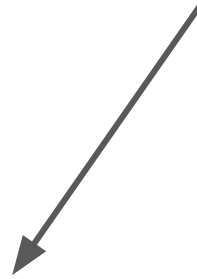


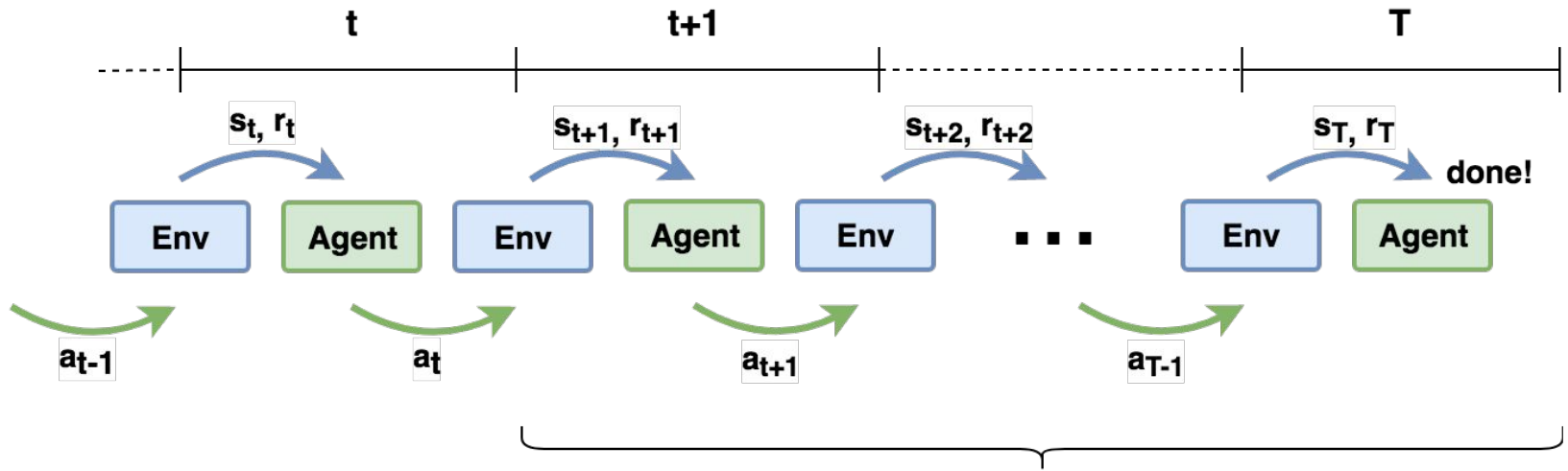
Q-learning

“Equation de Bellman”

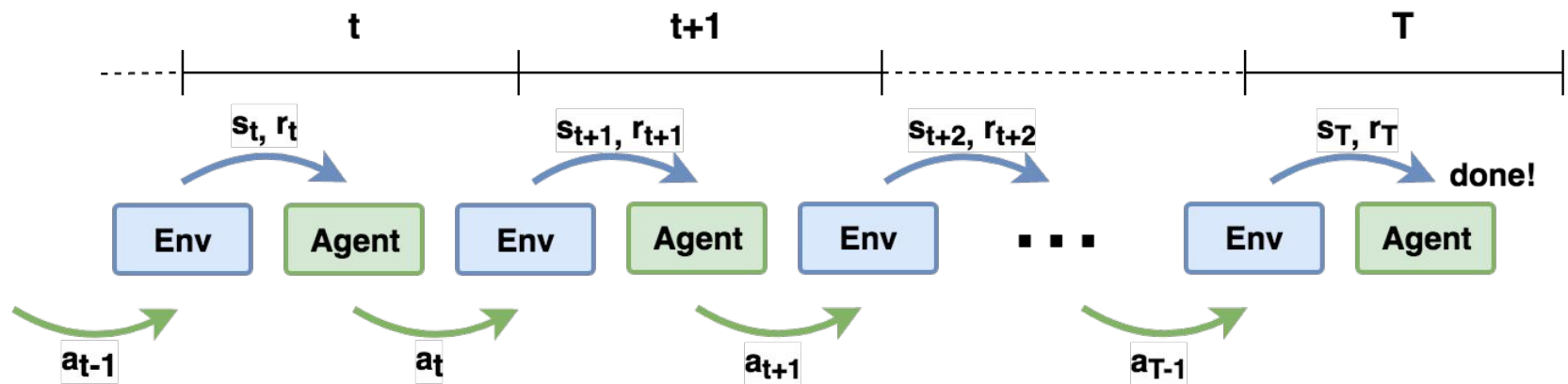
$$Q_{n+1}(s_t, a_t) = r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a)$$

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2





$$Q(s_t, a_t) = \sum r_{t+1} + r_{t+2} + \dots + r_T$$



$$Q_{n+1}(s_t, a_t) = r_{t+1} + \max_a Q_n(s_{t+1}, a)$$

Q-learning

Actions

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

States

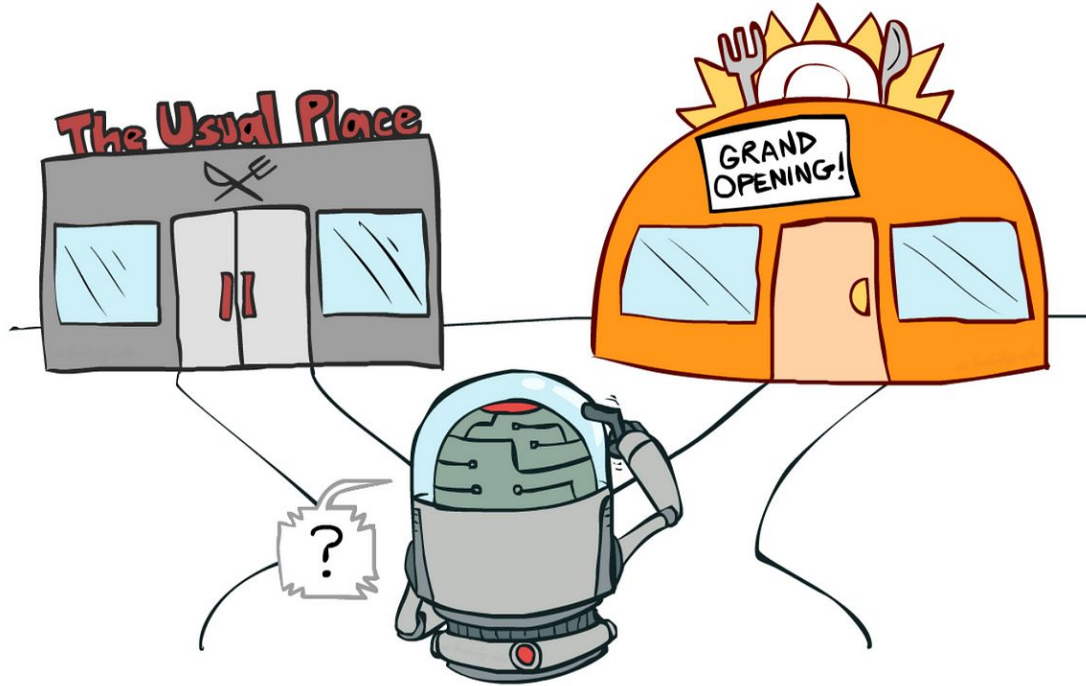
$$a = \operatorname{argmax}_a Q(s_t, a)$$



Equation de Bellman

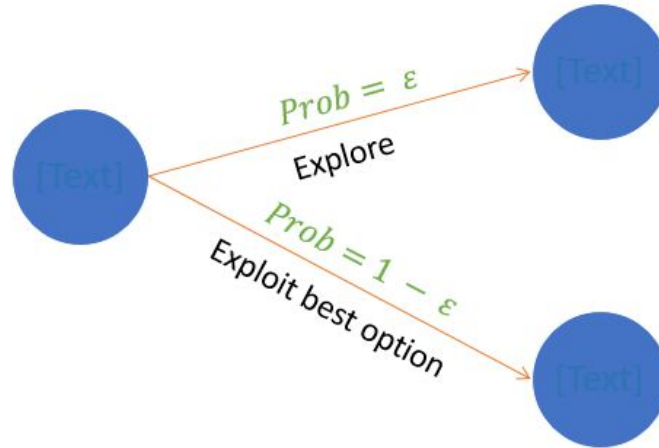


“Explorer” ou “exploiter” ?



Explorer... parfois !

Méthode “epsilon-greedy”



Q-learning

Actions

LEFT DOWN RIGHT UP STAY

→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

States

$$a = \underset{+}{\operatorname{argmax}}_a Q(s_t, a)$$

epsilon-greedy



Bellman Equation



Q-learning



**Fonctionne dans
tous les cas !***

Actions

LEFT DOWN RIGHT UP STAY

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

States

$$a = \underset{+}{\operatorname{argmax}}_a Q(s_t, a)$$

epsilon-greedy



Bellman Equation



Q-learning



Fonctionne dans
tous les cas !*

States

		Actions				
		LEFT	DOWN	RIGHT	UP	STAY
States	→	1.7	1.4	2.6	0.86	1.5
	↘	1.3	1.6	2.4	0.78	1.5
	↑	1.9	1.6	1.4	2.5	1.7
	↗	1.4	0.58	1.8	1.6	1.6
	↖	2.2	1.3	0.97	1.6	1.5
	↙	2.5	2	1.3	0.9	1.5
	↓	1.8	2.9	1.7	1.2	2.1
	←	2.6	1.5	2.1	1.2	2.2

$$a = \underset{+}{\operatorname{argmax}}_a Q(s_t, a)$$

epsilon-greedy



Bellman Equation



* peut prendre un temps infini 😊

Bellman Equation

$$\underbrace{Q_{n+1}(s_t, a_t)}_{\text{new Q-value}} = r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a)$$

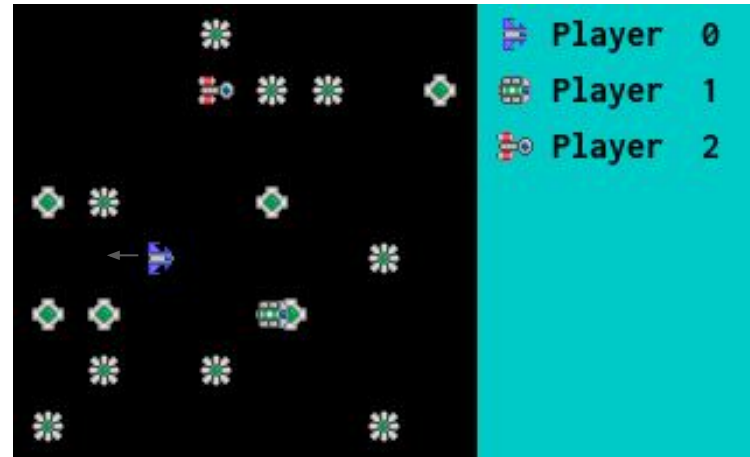
	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2



Bellman Equation

$$\underbrace{Q_{n+1}(s_t, a_t)}_{\text{new Q-value}} = r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a)$$

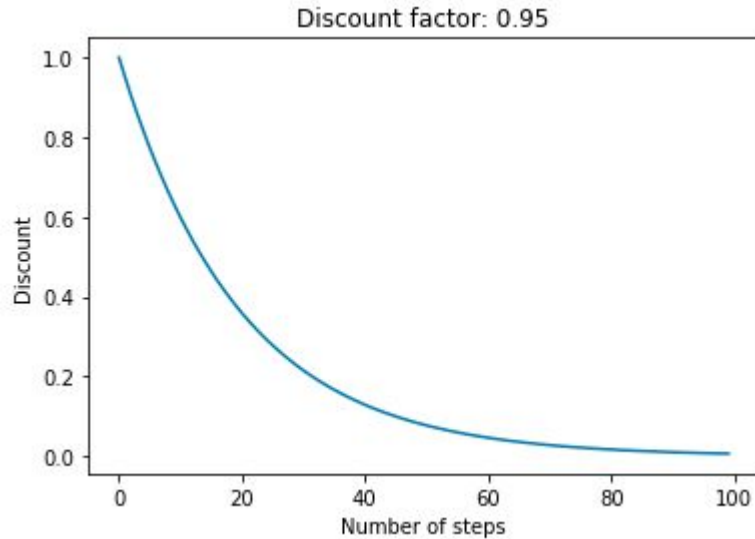
	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2



Equation de Bellman

$$Q_{n+1}(s_t, a_t) = r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a)$$

Discount factor
("facteur d'actualisation")



Learning rate



$$\underbrace{Q_{n+1}(s_t, a_t)}_{\text{new Q-value}} = (1 - \alpha)Q_n(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a))$$

alpha = 0 # No update

alpha = 1 # Full update

alpha = 0.5 # Mean between old/new

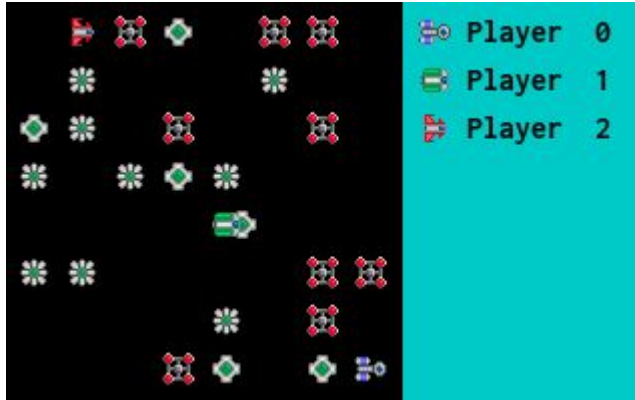
Mise en pratique

Introduction à Q-learning

<https://colab.research.google.com/drive/1KSei3ZdjNyyUYsPMqTyli4G2rm9v-P2E>



Q-table size



```
{'target_dir': 2, 'lidar', [0, 0, 0, 0, 1, 1, 1, 1]}
```

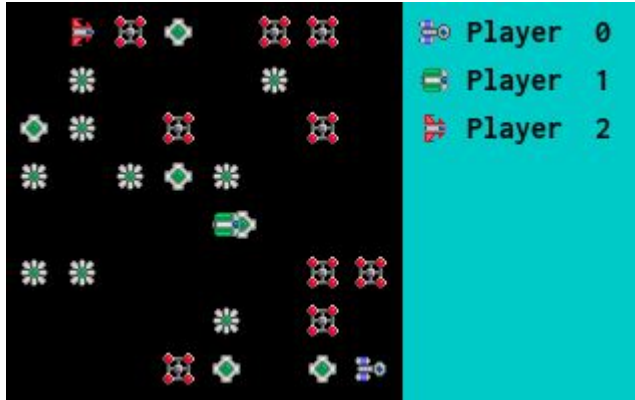


```
Dict(target_dir:Discrete(8), lidar:MultiBinary(8))
```



```
n_states = 8 * (2**8) = 2048
```

Q-table size



```
{'target_dir': 2, 'lidar', [0, 0, 0, 0, 1, 1, 1, 1]}
```



```
Dict(target_dir:Discrete(8), lidar:MultiBinary(8))
```



```
n_states = 8 * (2**8) = 2048
```

DQN

Deep Q-Network

Rappel : Q-learning

Actions

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

Observations

On sélectionne l'action qui maximise la somme attendue des récompenses

$$a = \operatorname{argmax}_a Q(s, a)$$



Rappel : Q-learning

	LEFT	DOWN	RIGHT	UP	STAY
→	1.7	1.4	2.6	0.86	1.5
↘	1.3	1.6	2.4	0.78	1.5
↑	1.9	1.6	1.4	2.5	1.7
↗	1.4	0.58	1.8	1.6	1.6
↖	2.2	1.3	0.97	1.6	1.5
↙	2.5	2	1.3	0.9	1.5
↓	1.8	2.9	1.7	1.2	2.1
←	2.6	1.5	2.1	1.2	2.2

“Equation de Bellman”

$$Q_{n+1}(s_t, a_t) = r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a)$$



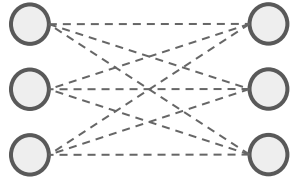
Deep Q-network (DQN)

charge: 16.18%

lidar_left: 0

lidar_down: 1

...

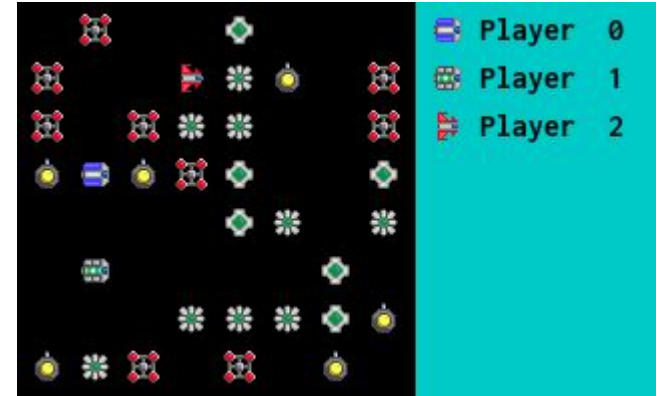


$Q(\text{Up}) = 0.33$

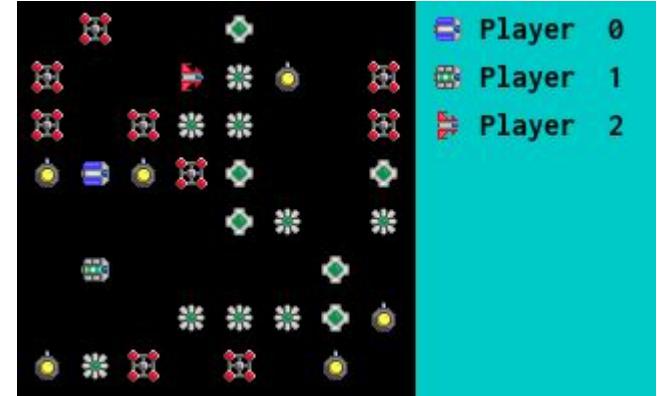
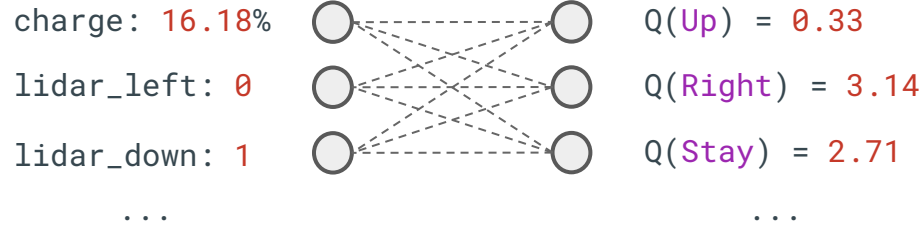
$Q(\text{Right}) = 3.14$

$Q(\text{Stay}) = 2.71$

...



Deep Q-network (DQN)

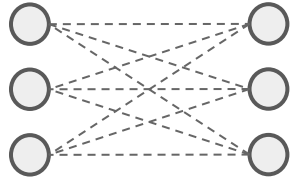


“Q-network”

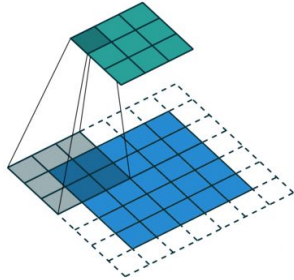
$$\delta = Q_n(s_t, a_t) - (r_{t+1} + \gamma \max_a Q_n(s_{t+1}, a))$$

δ est l'erreur à minimiser

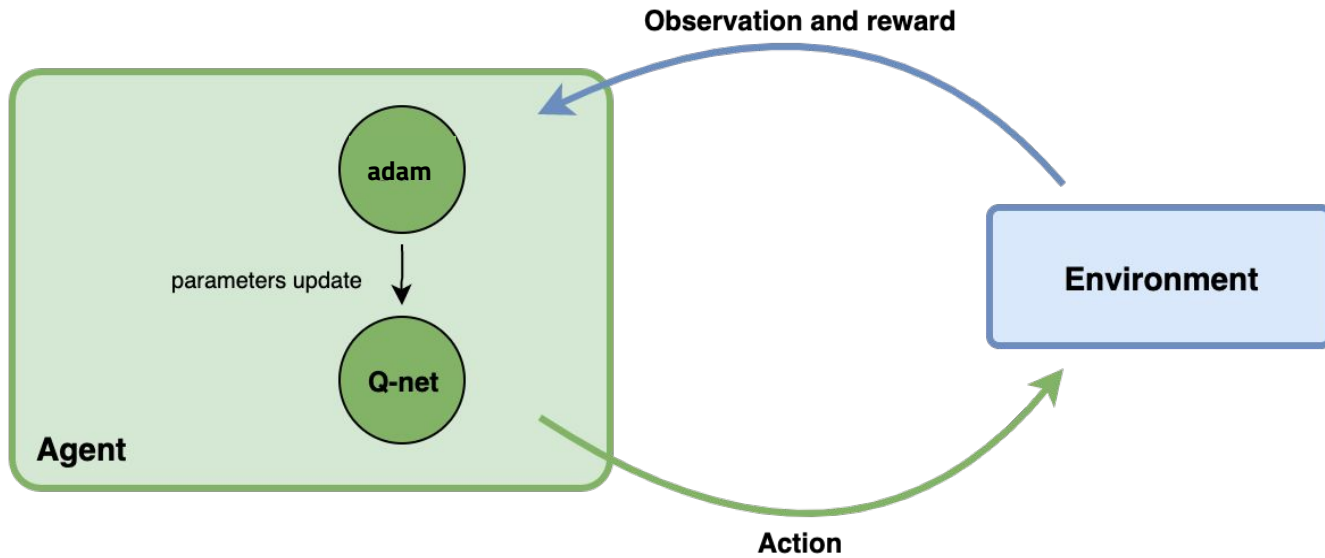
Deep Q-network (DQN)



Couche dense (“fully connected”)



Couche à convolutions



Timestep: 1, action: left

Reward: 1.0



Timestep: 2, action: right

Reward: 1.0



Timestep: 3, action: left

Reward: 1.0



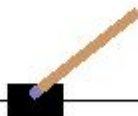
Timestep: 4, action: left

Reward: 1.0



Timestep: 5, action: left

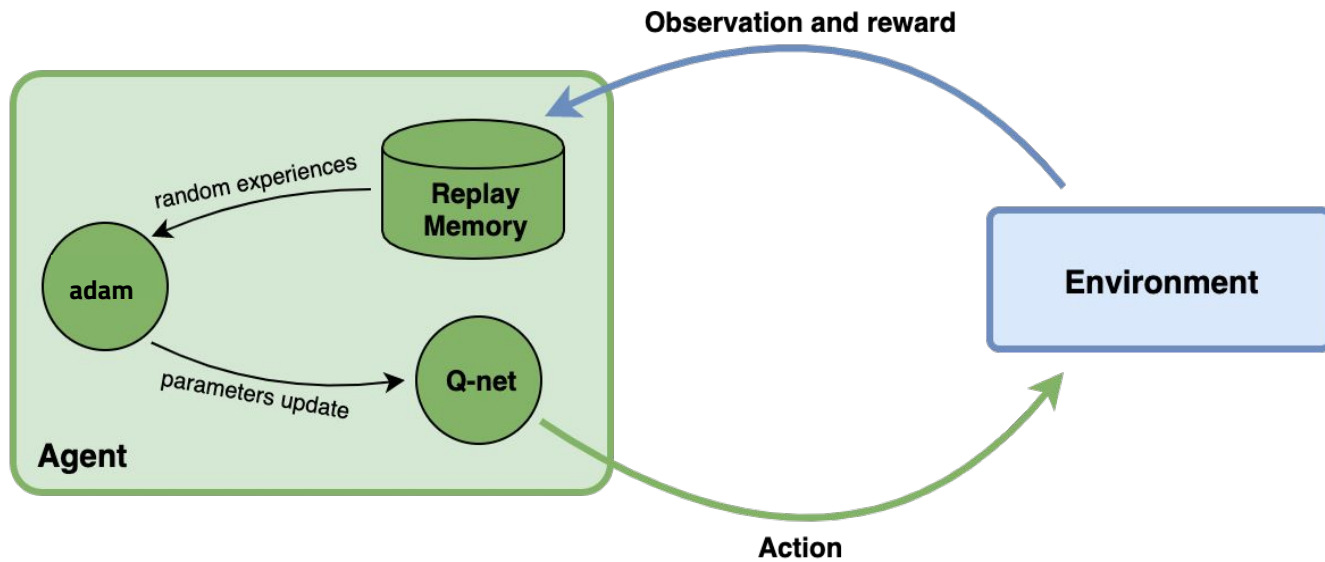
Reward: 1.0

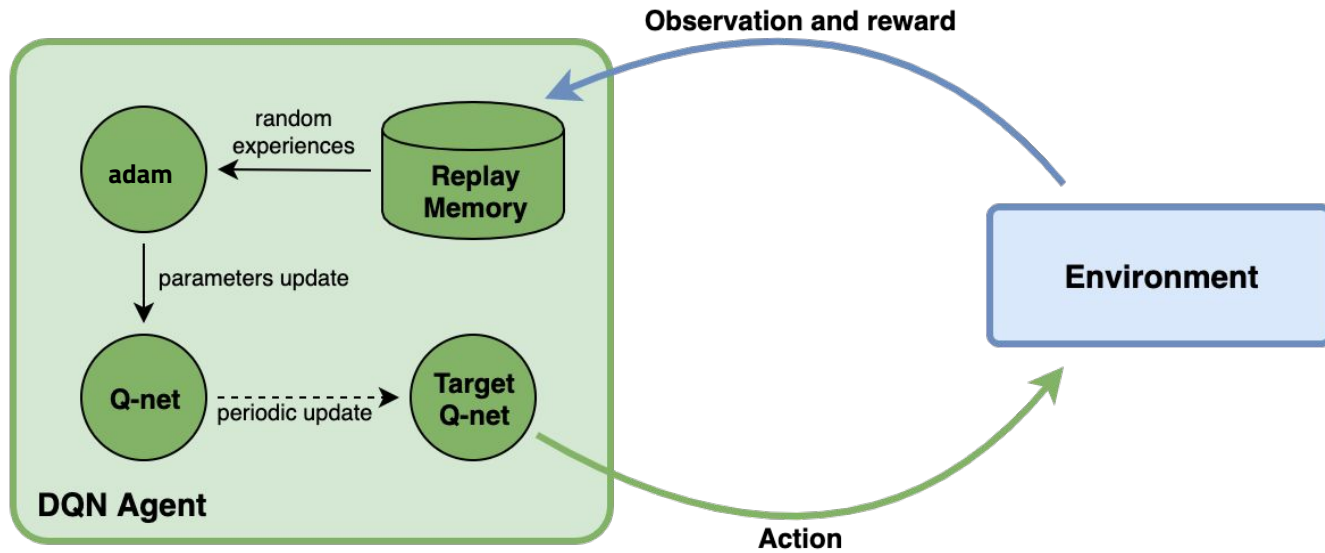


Timestep: 6, action: right

Reward: 0.0







A vous de jouer !



Mise en pratique

Deep Q-Network

<https://colab.research.google.com/drive/1R9X7vTaoTzDI1QGoArBMnMwgPB-BXK2C>



Peut-on faire mieux ?



Les difficultés de l'AR

- Entraînement long
- Sensibilité des paramètres
- Manque de généralisation
- Dilemme exploration/exploitation

Trucs et Astuces

- “Reward engineering”
 - Ajustement à la main des récompenses
 - Mais, risque d'introduire des biais !
- Entraîner dans différents environnements
 - Appeler plusieurs fois `trainer.train()`
- Faire varier le paramètre epsilon
 - Le faire descendre puis remonter de façon cyclique ?
- Entraîner plus n'améliore pas toujours les choses !
 - Évaluer régulièrement, puis garder le meilleur agent ?

Les difficultés de l'AR

- Entraînement long
- Sensibilité des paramètres
- Manque de généralisation
- Dilemme exploration/exploitation

Apprentissage Priorité !



Prioritized Experience Replay

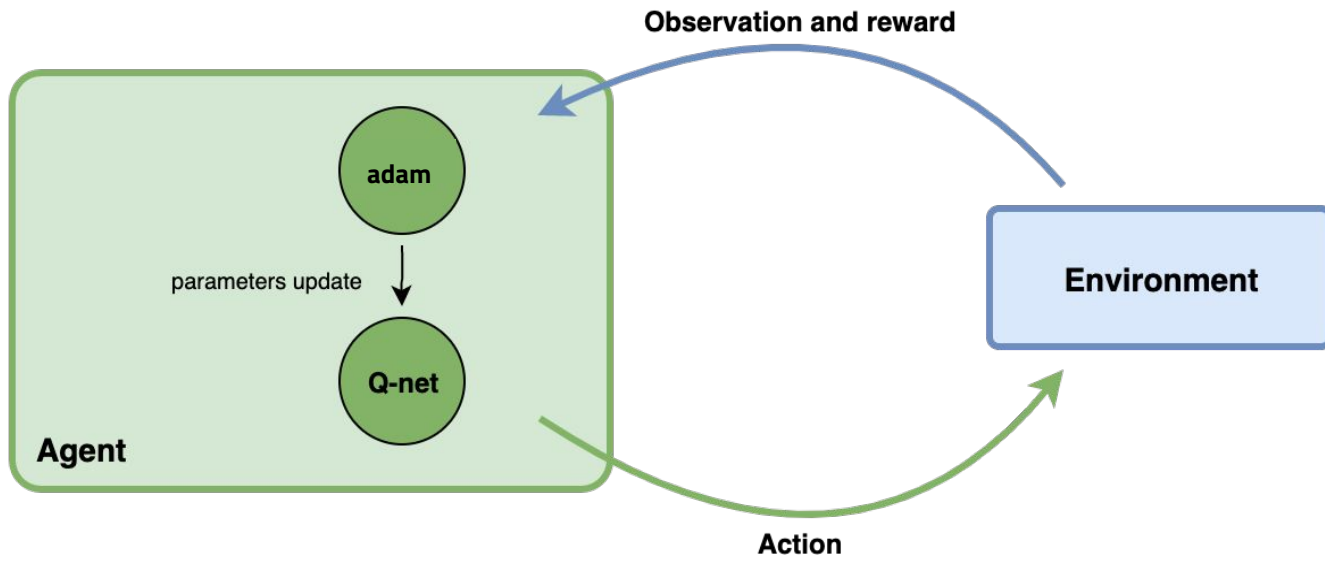
(“Rejeu d'Expériences Priorisées”)

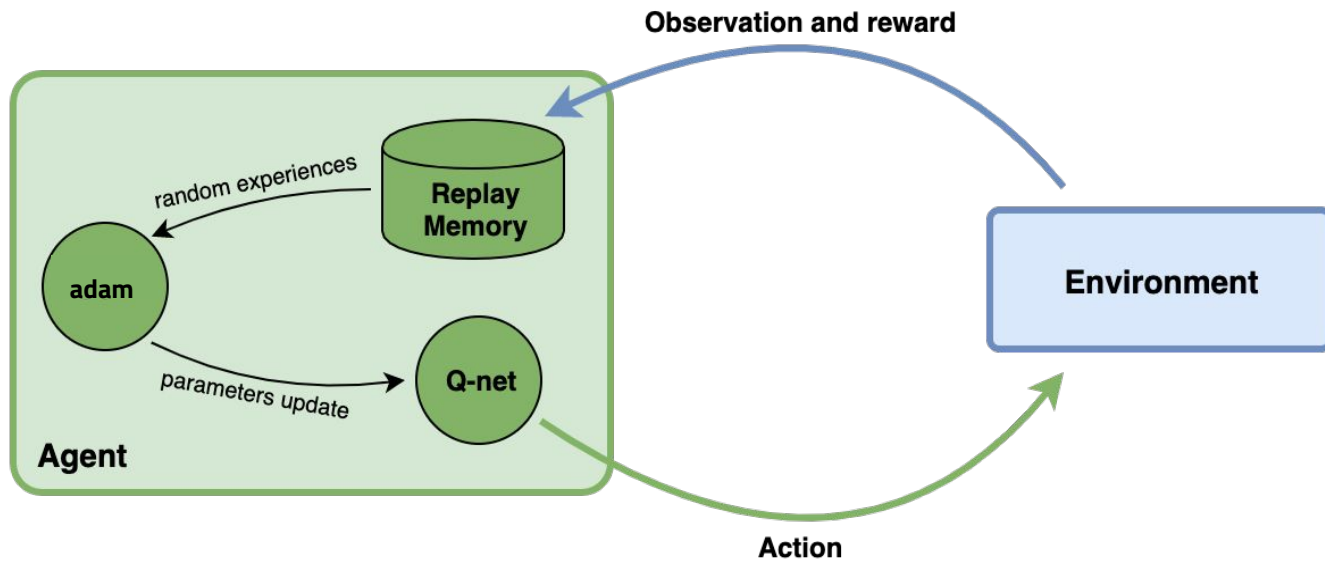
Schaul et al. 2015

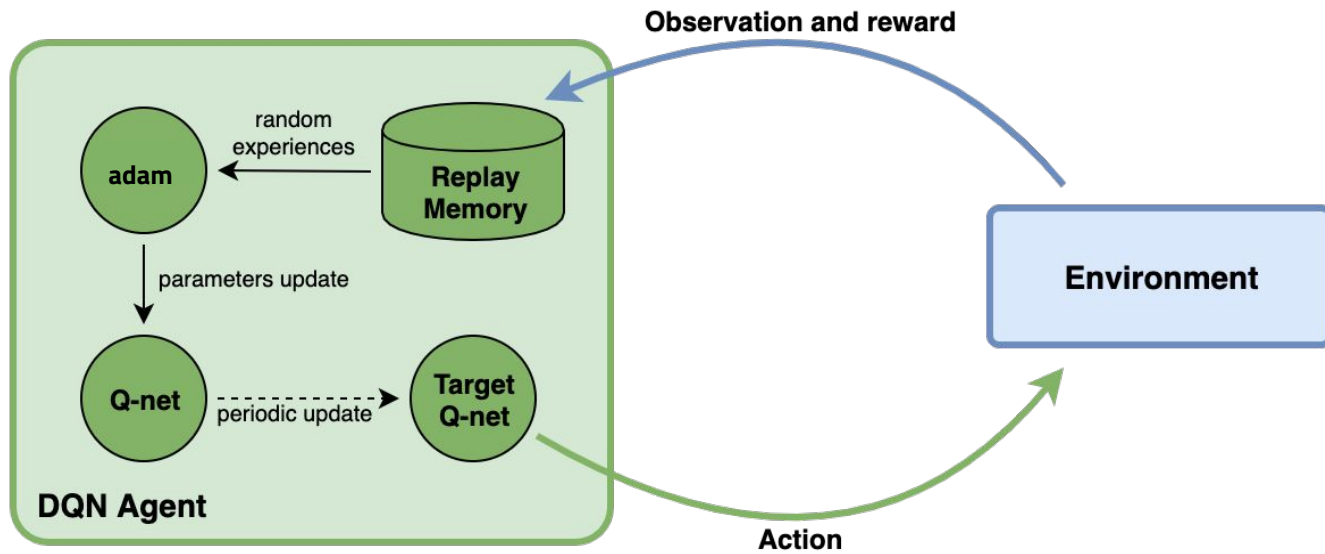
<https://arxiv.org/abs/1511.05952>

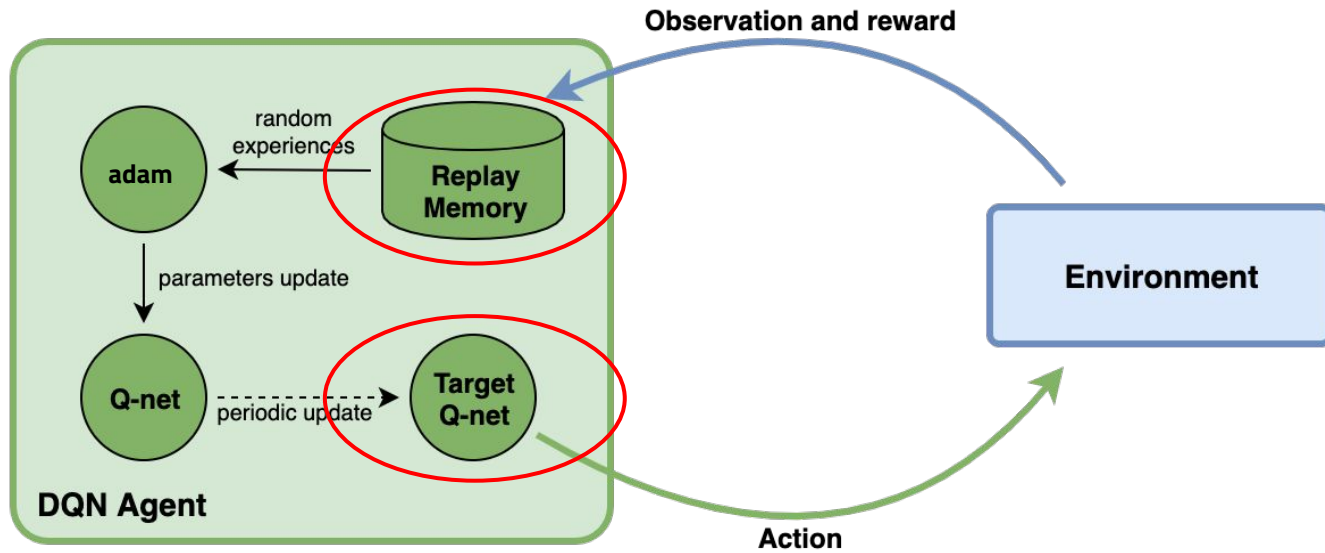
But :
apprendre plus efficacement

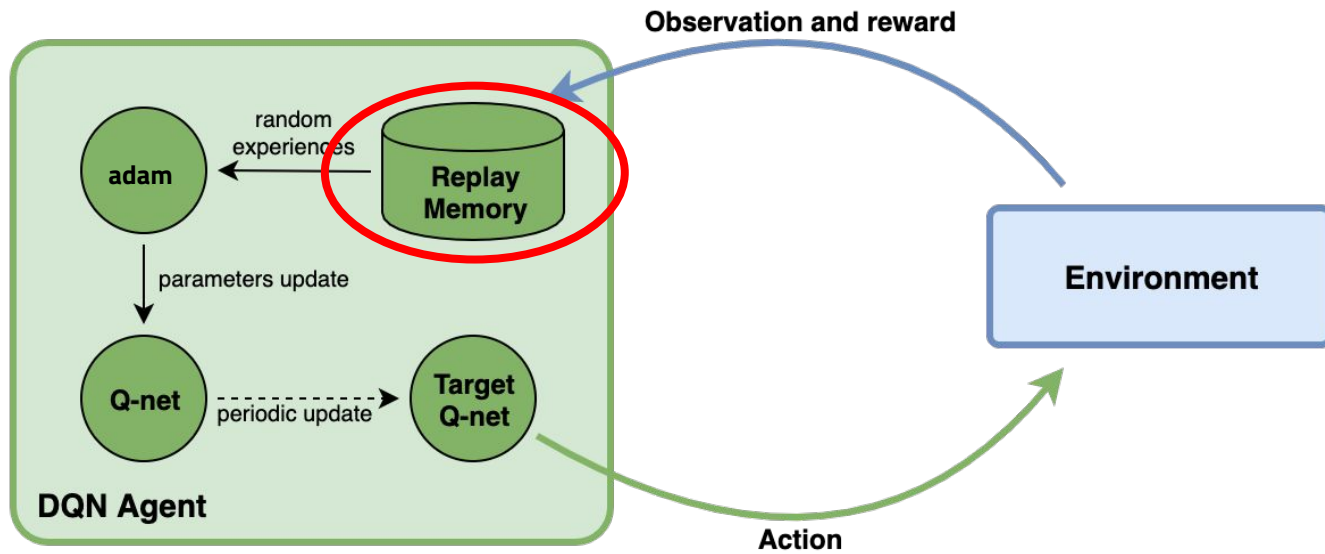






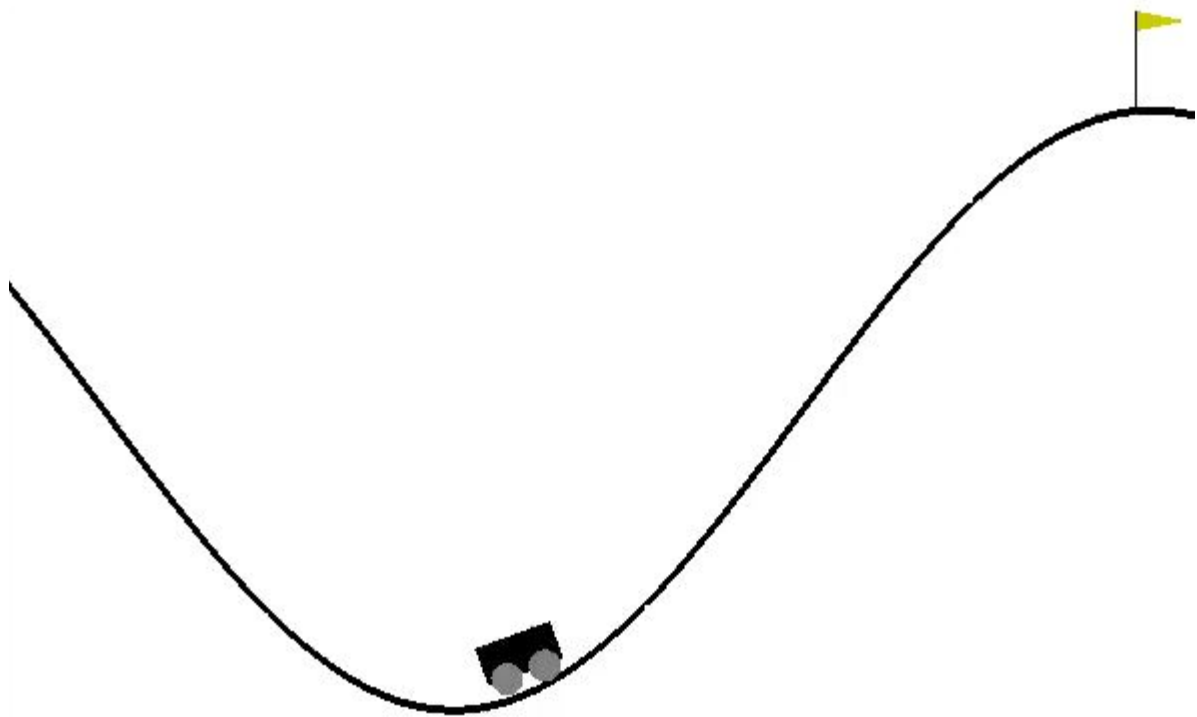


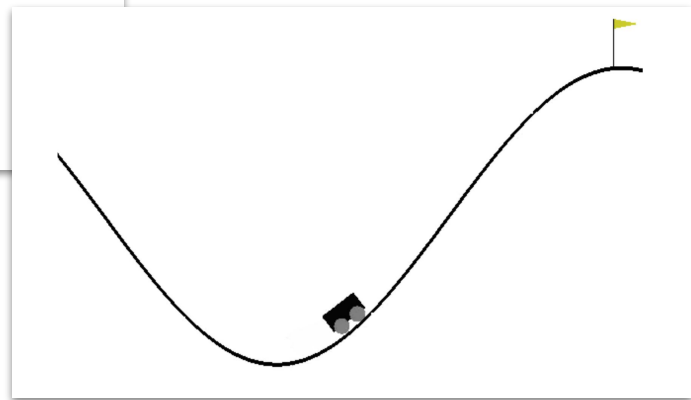
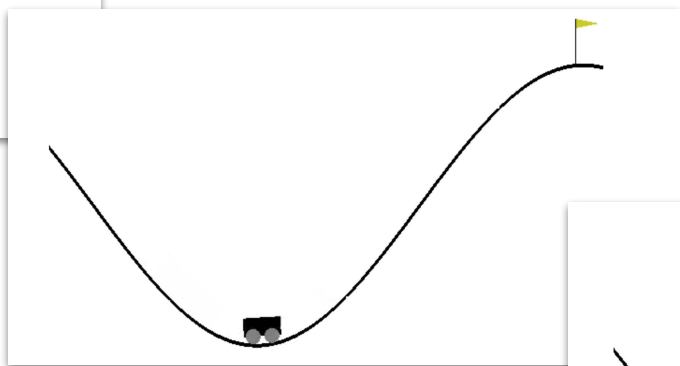
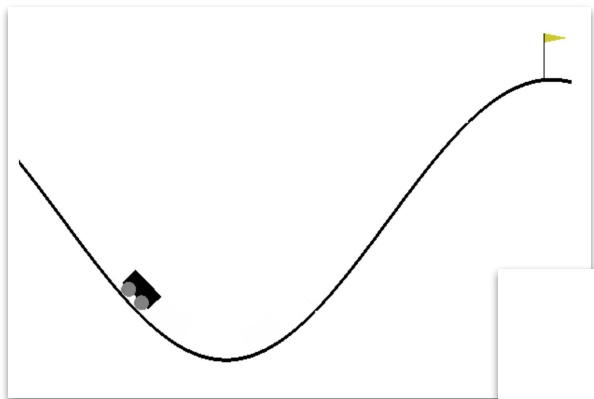




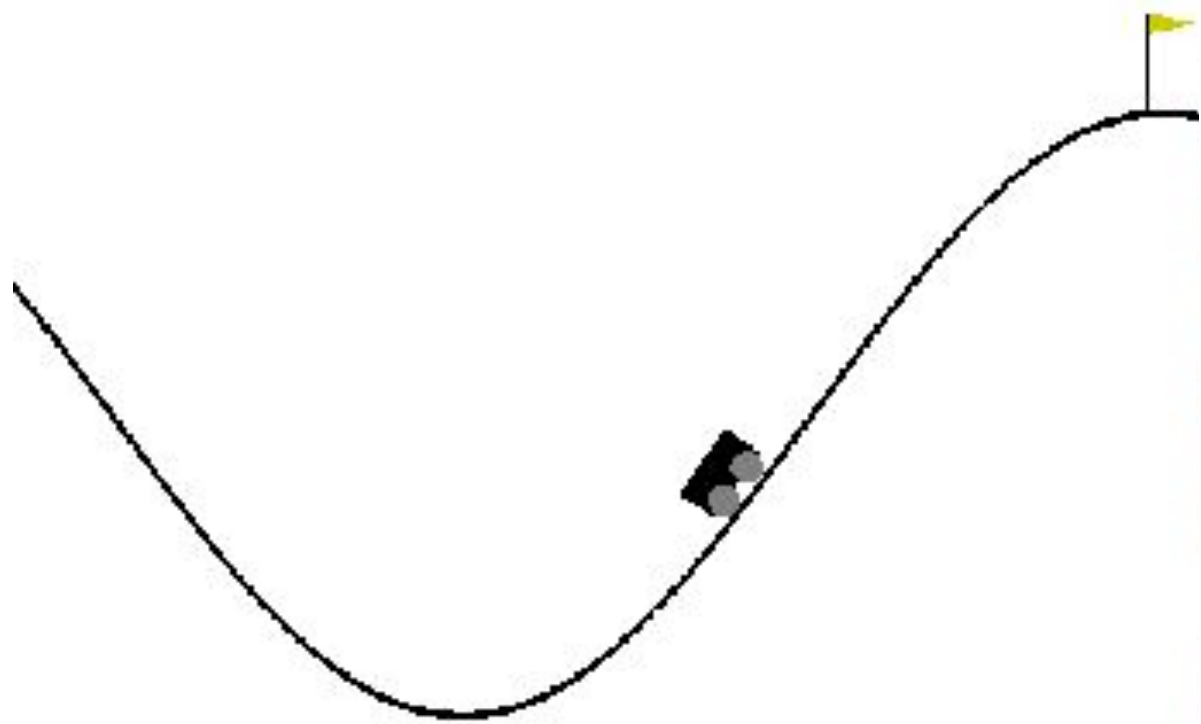
“Toutes les expériences ne sont pas égales”

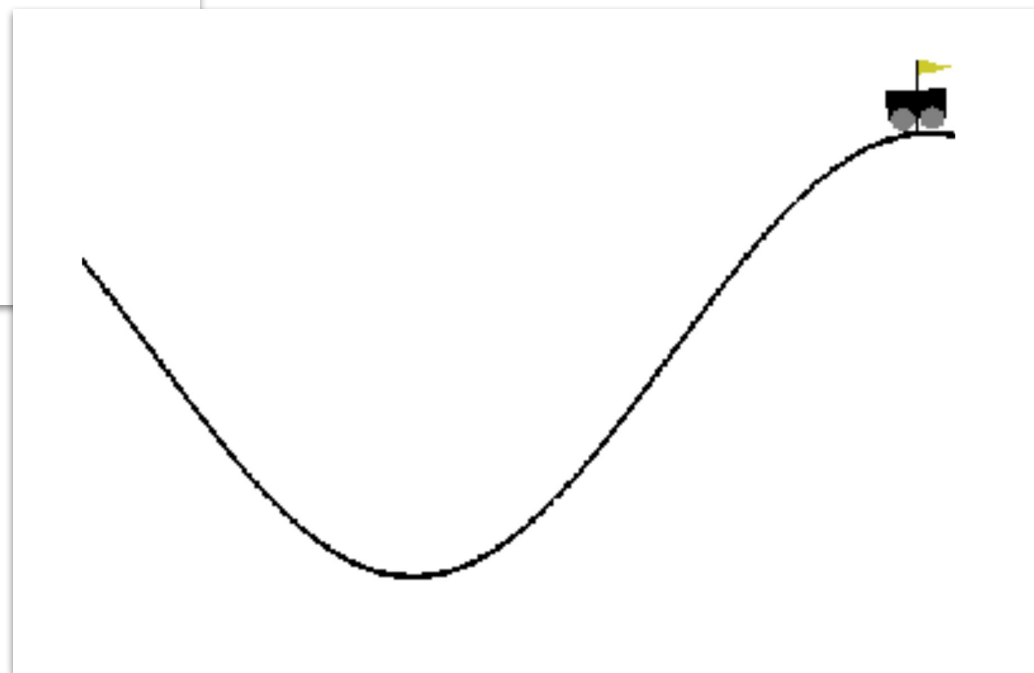
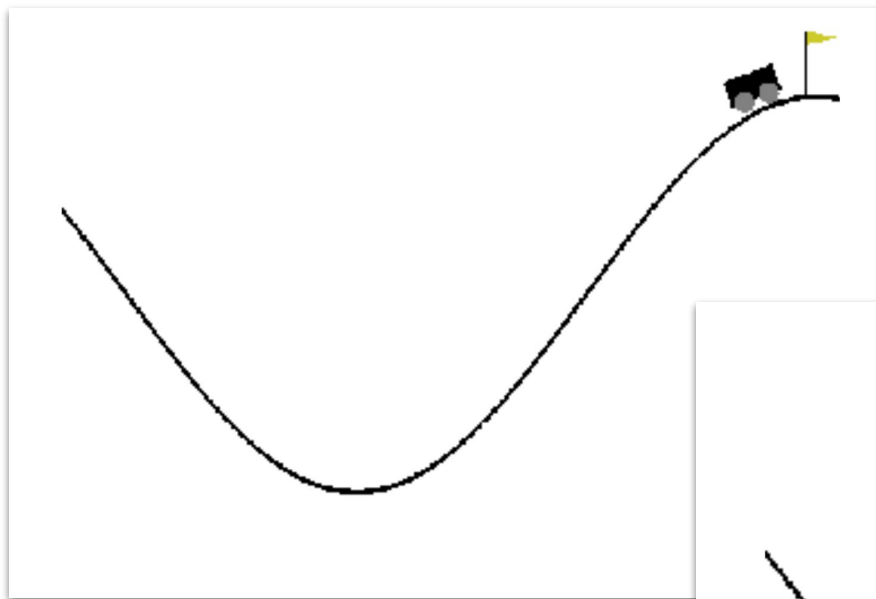
A Katharopoulos





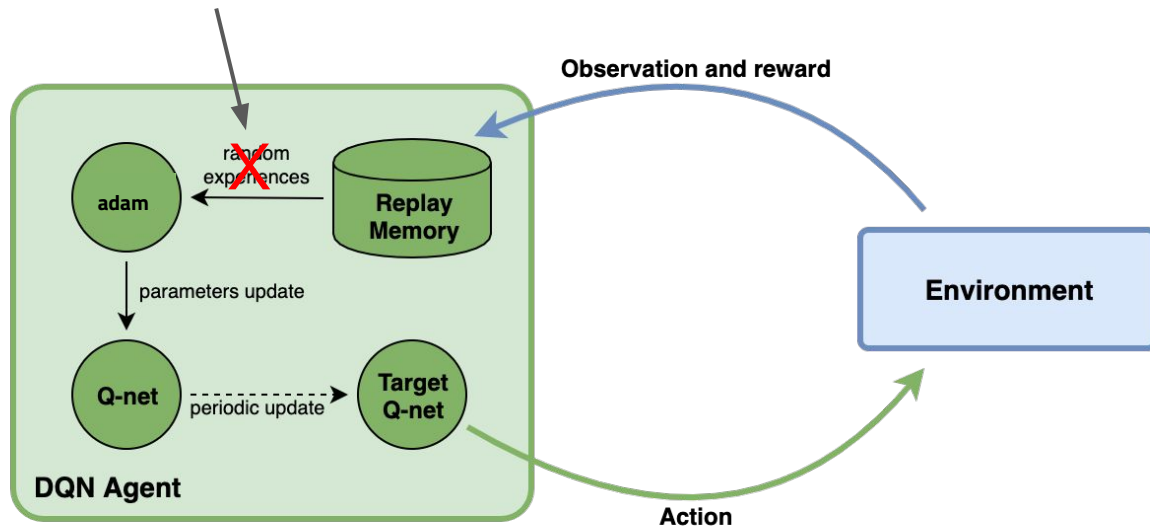
99% du replay buffer

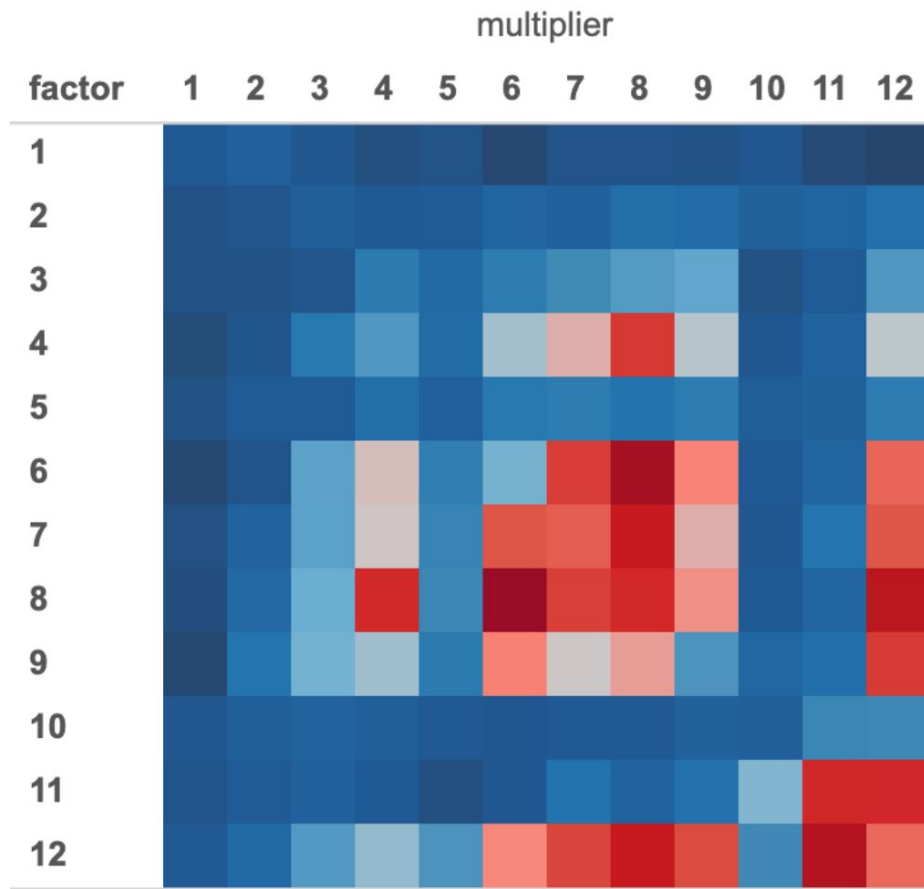


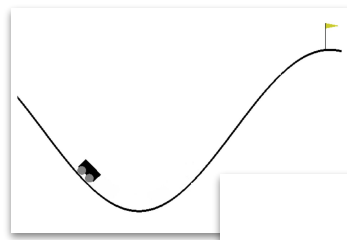


Les 1% !

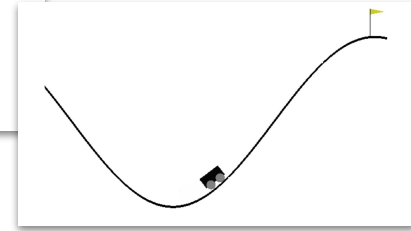
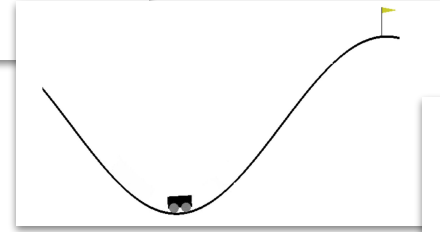
Sample in smarter way!



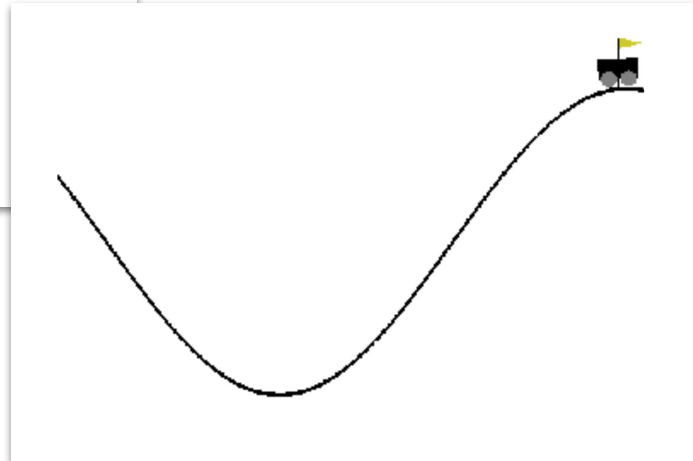
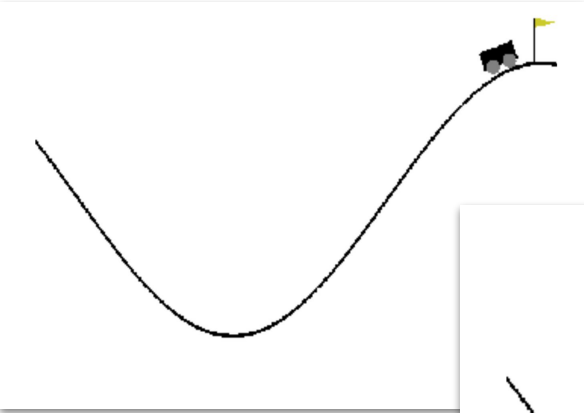


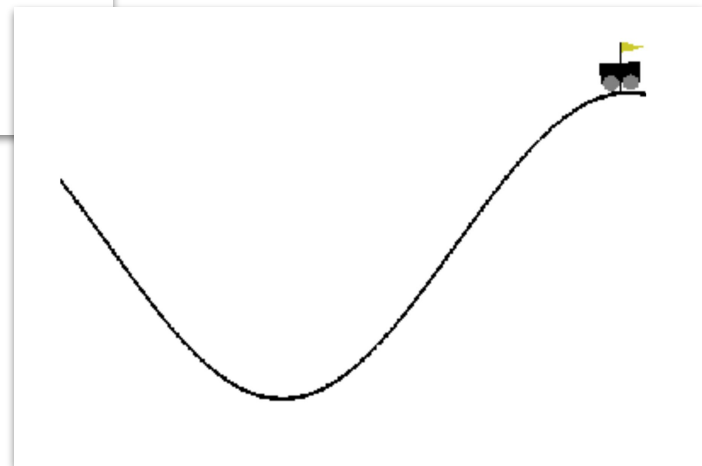
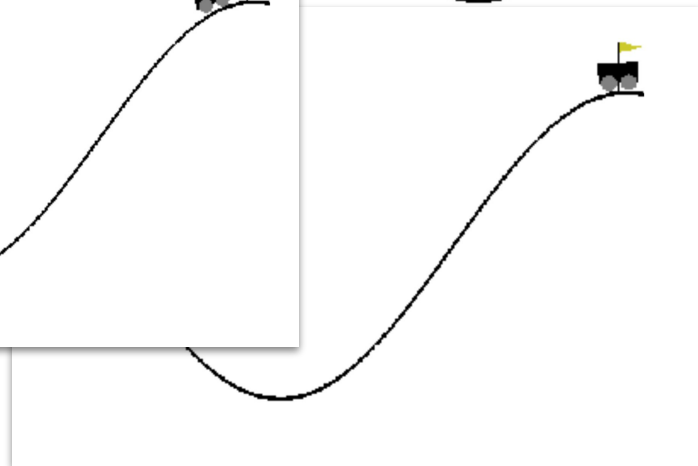
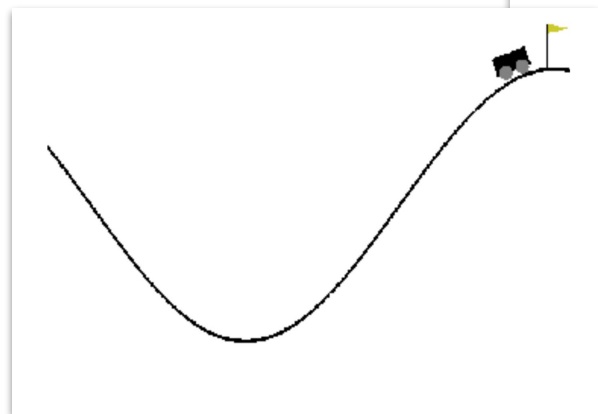
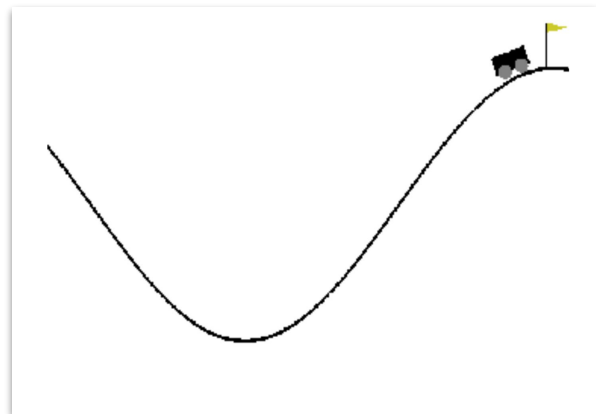
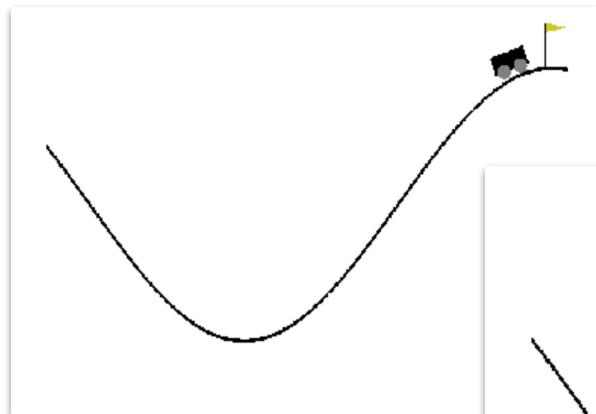


Erreur faible

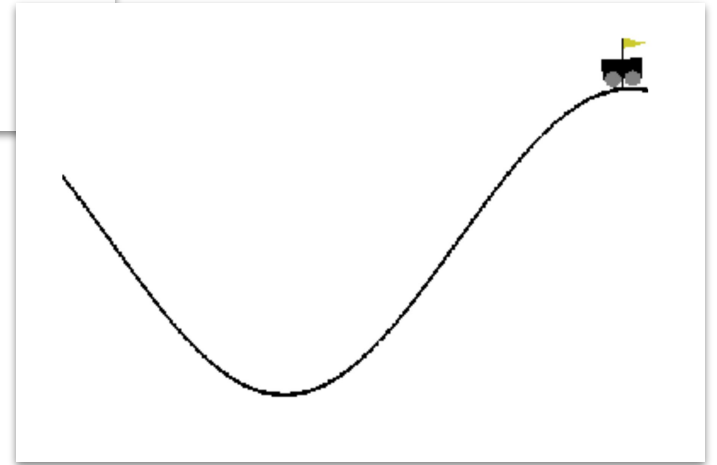
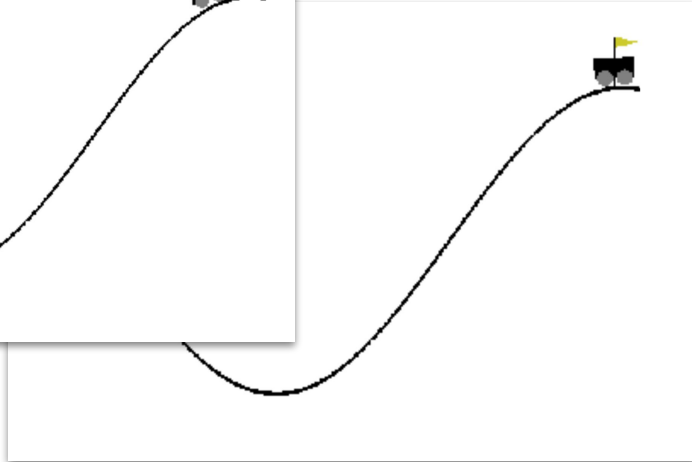
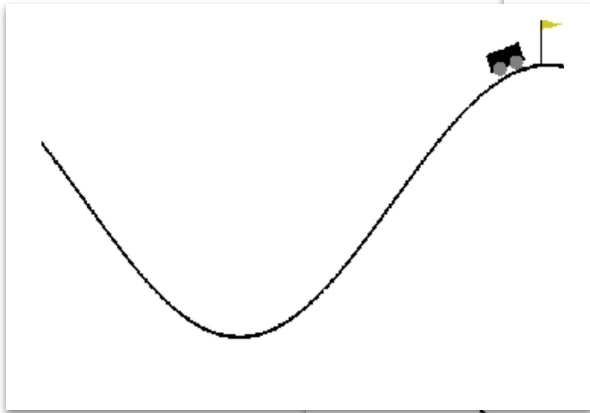
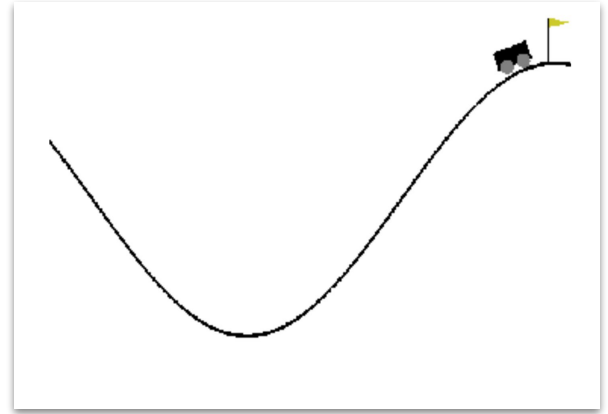
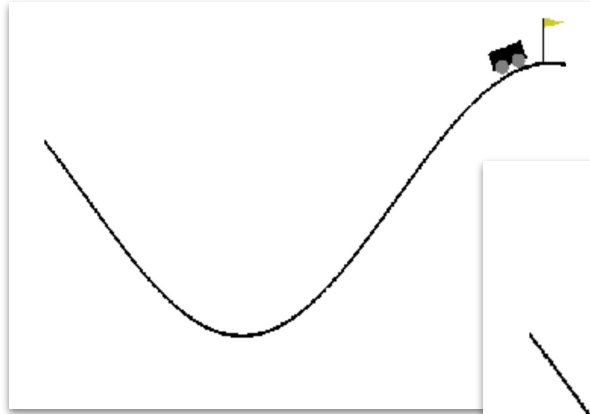


Erreur élevée !





Trop de biais!

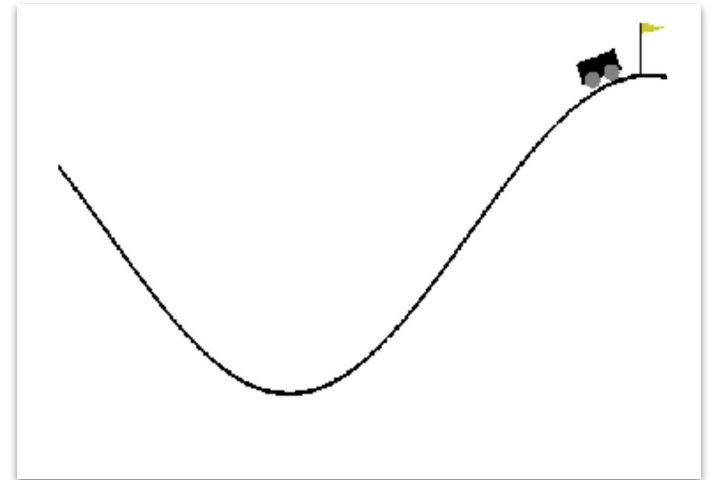
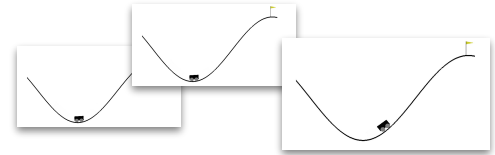


Réduisons le biais...

α

Débiaisage

Avant

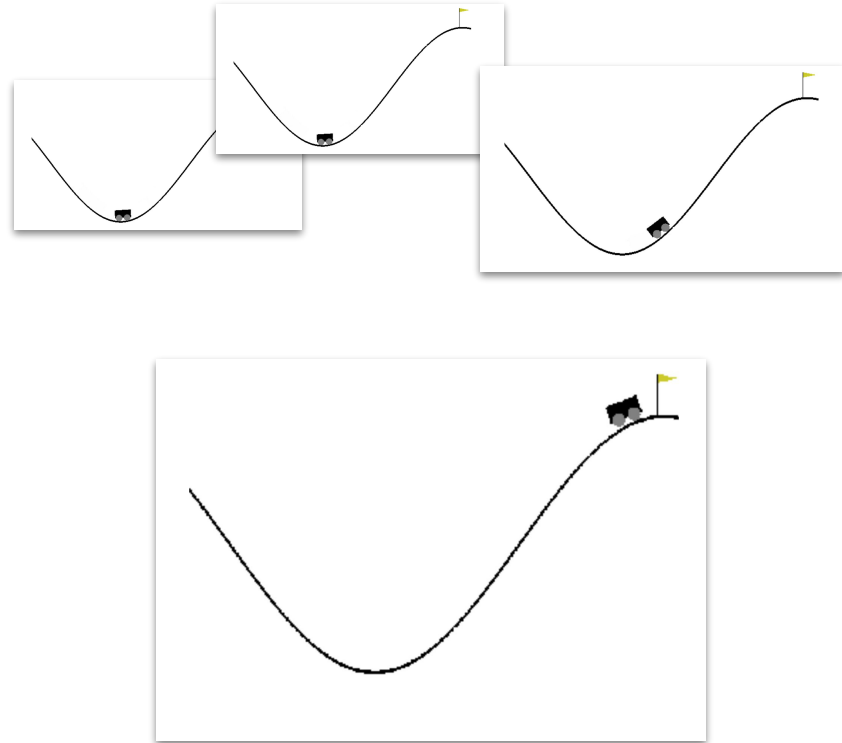


Réduisons le biais...

α

Débiaisage

Après



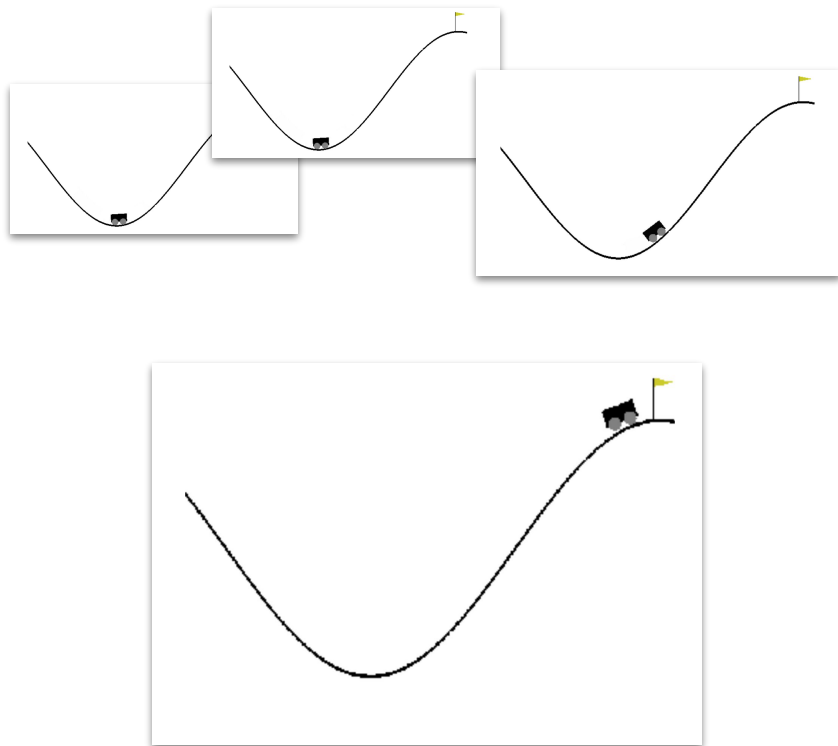
Réduisons le biais...

α

Débiaisage

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

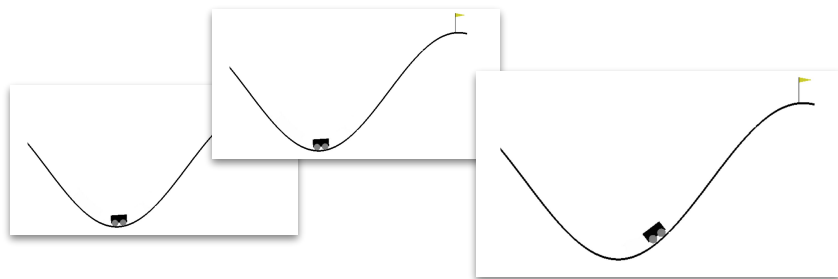
Après



Réduisons le biais...

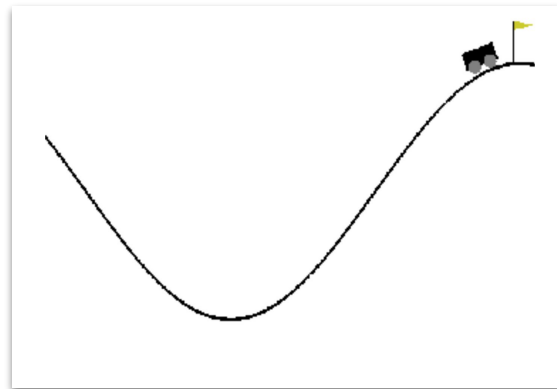
α

Débiaisage



β

Correction de l'échantillonnage

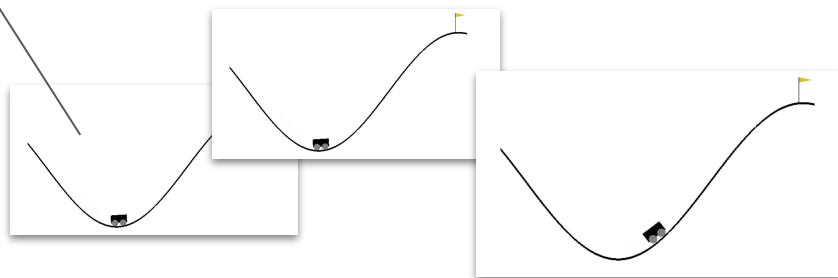


Réduisons le biais...

Plus d'effet

α

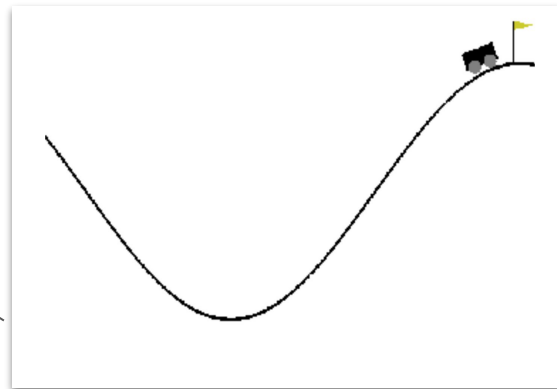
Débiaisage



β

Correction de l'échantillonnage

Moins d'effet

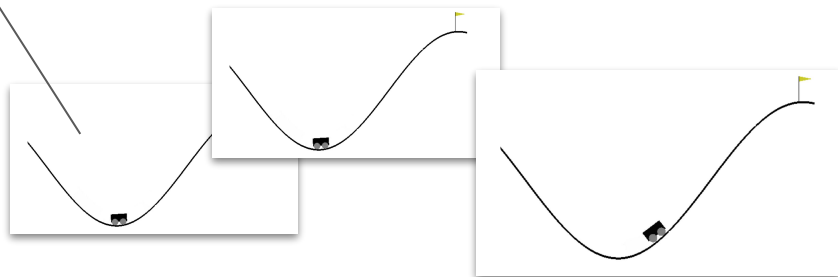


Réduisons le biais...

Plus d'effet

α

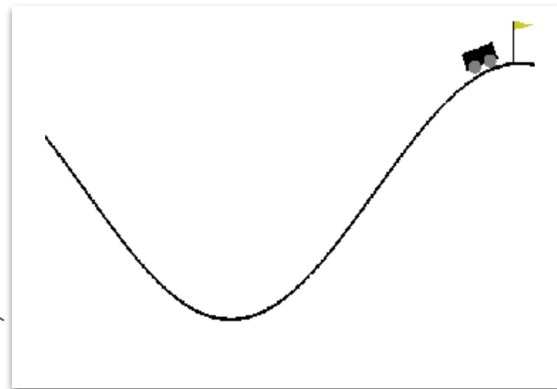
Débiaisage
= 0.6

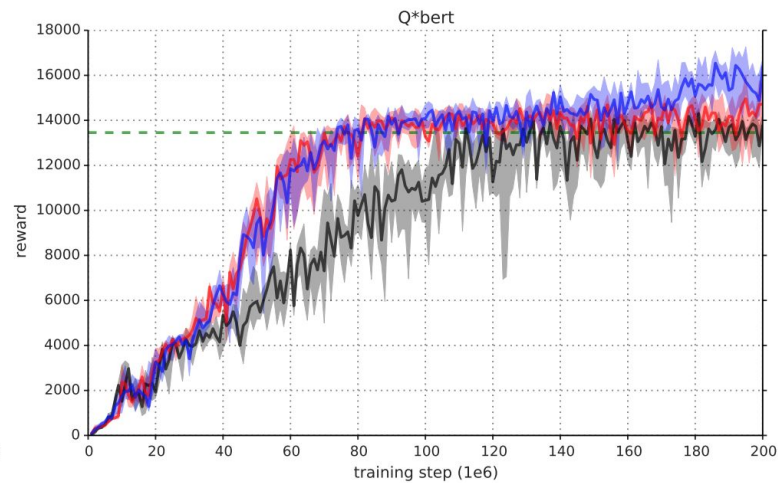
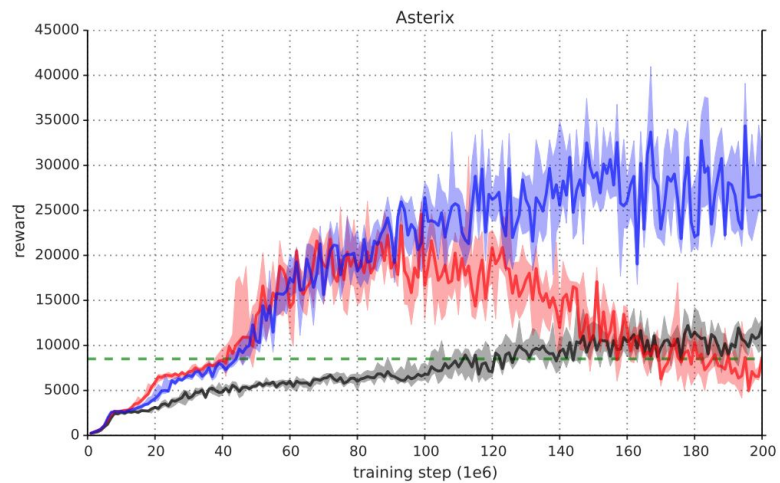
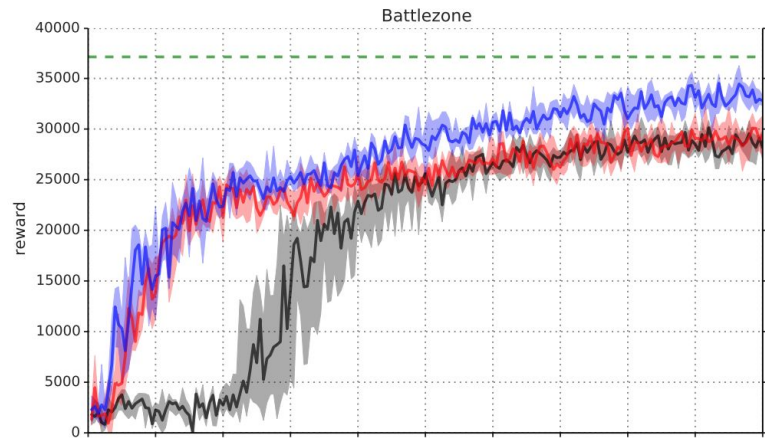
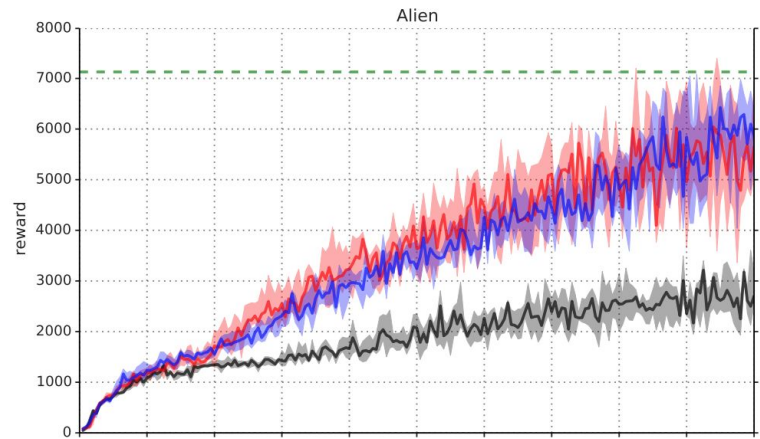


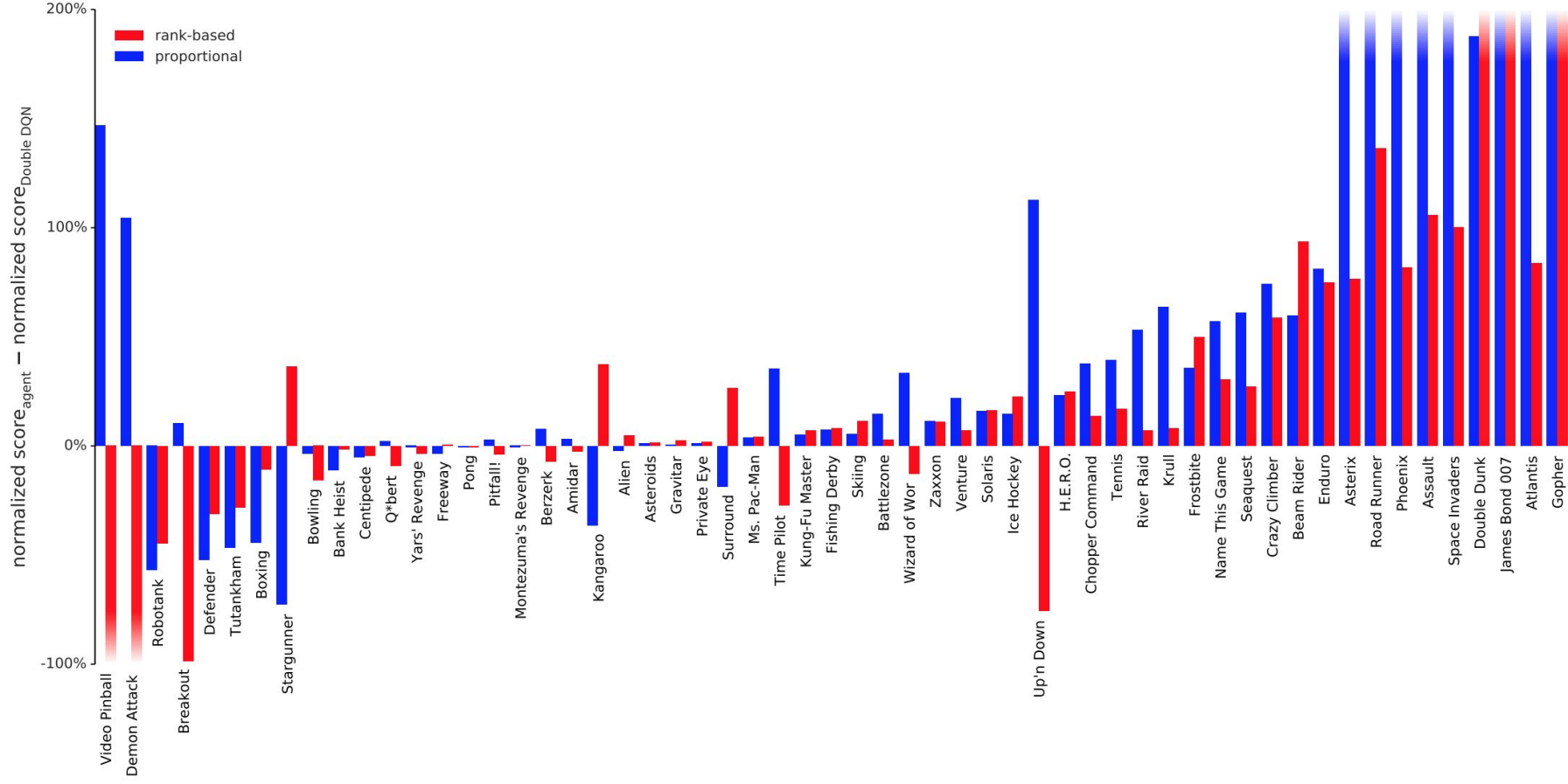
β

Correction de l'échantillonnage
= 0.4

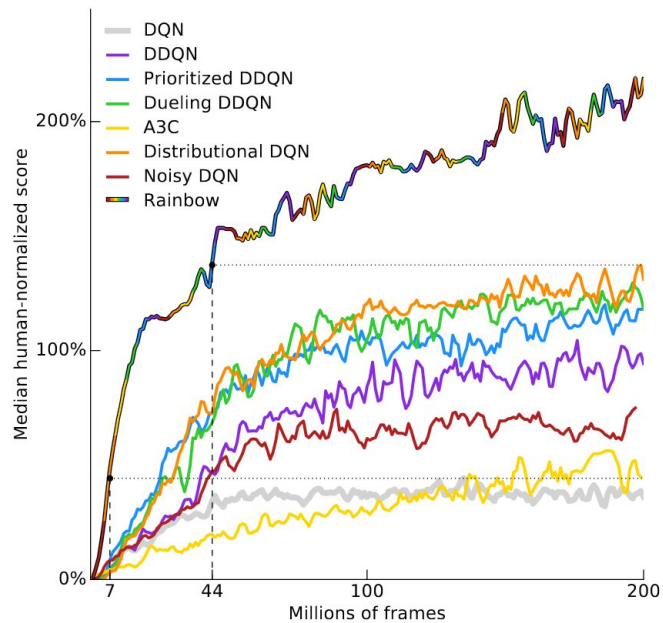
Moins d'effet







Rainbow DQN



A vous de jouer !



Mise en pratique

Apprentissage Prioritis 

[https://colab.research.google.com/drive/
156LKB0aPnKab1loUjeEKwV3VTmL7fguW](https://colab.research.google.com/drive/156LKB0aPnKab1loUjeEKwV3VTmL7fguW)

