



Formation

Introduction au Deep Learning

18

Deep Reinforcement Learning
Tactiques et stratégies



You can also subscribe to :



<http://fidle.cnrs.fr/listeinfo>
Fidle information list



GROUPE **CALCUL**

<https://listes.services.cnrs.fr/wws/info/devlog>
List of ESR* « Software developers » group

<https://listes.math.cnrs.fr/wws/info/calcul>
List of ESR* « Calcul » group

<https://fidle.cnrs.fr>

Powered by CNRS CRIC, and UGA DGDSI
of Grenoble, Thanks !



Course materials (pdf)



Practical work environment^{*}



Corrected notebooks



Videos (YouTube)

(*) Procedure via Docket or pip
Remember to get the latest version !

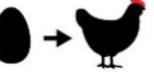
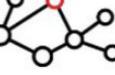
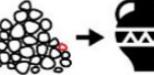


Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)
<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Program

FIDLE



1	 History, Fundamental Concepts		2	3	 High Dimensional Data CNN		4	 Demystify mathematics for neural networks.		5	 Training strategies Evaluation	 Sparse data (text) Embedding		6	 Sequences data RNN	
	 Basic Regression DNN	 Basic Classification DNN	7	8	 PyTorch A small detour with PyTorch.		9	 «Attention is All You Need» Transformers		10	 Graph Neural Network GNN		11	 Variational Autoencoder VAE		
12	► JDLS 2023	 Project session «My project in 180 s»	13	14	 Generative Adversarial Networks GAN		15	 Diffusion Model Text to image		16	 Model and training optimization Resource efficiency		17	 Jean-Zay GPU acceleration		
18	 Deep Reinforcement Learning RL				19	 Physics-Informed Neural Networks PINNs		20	 JDLS 2023 Deep Learning for Science !							

20 Séquences
du 17 novembre
au 14 mai 2023



SAISON
22/23

Roadmap

18



Deep Reinforcement
Learning
RL



18.1

About RL

- RL Introduction
- Tabular RL : Sarsa & Q-learning
- Example : Cliff Walking



18.2

From RL to Deep RL

- Advanced topics in RL
- DQNs



18.3

Policy Gradient

- On-policy Policy Gradient
- Off-policy Policy Gradient
- Example : Car pole



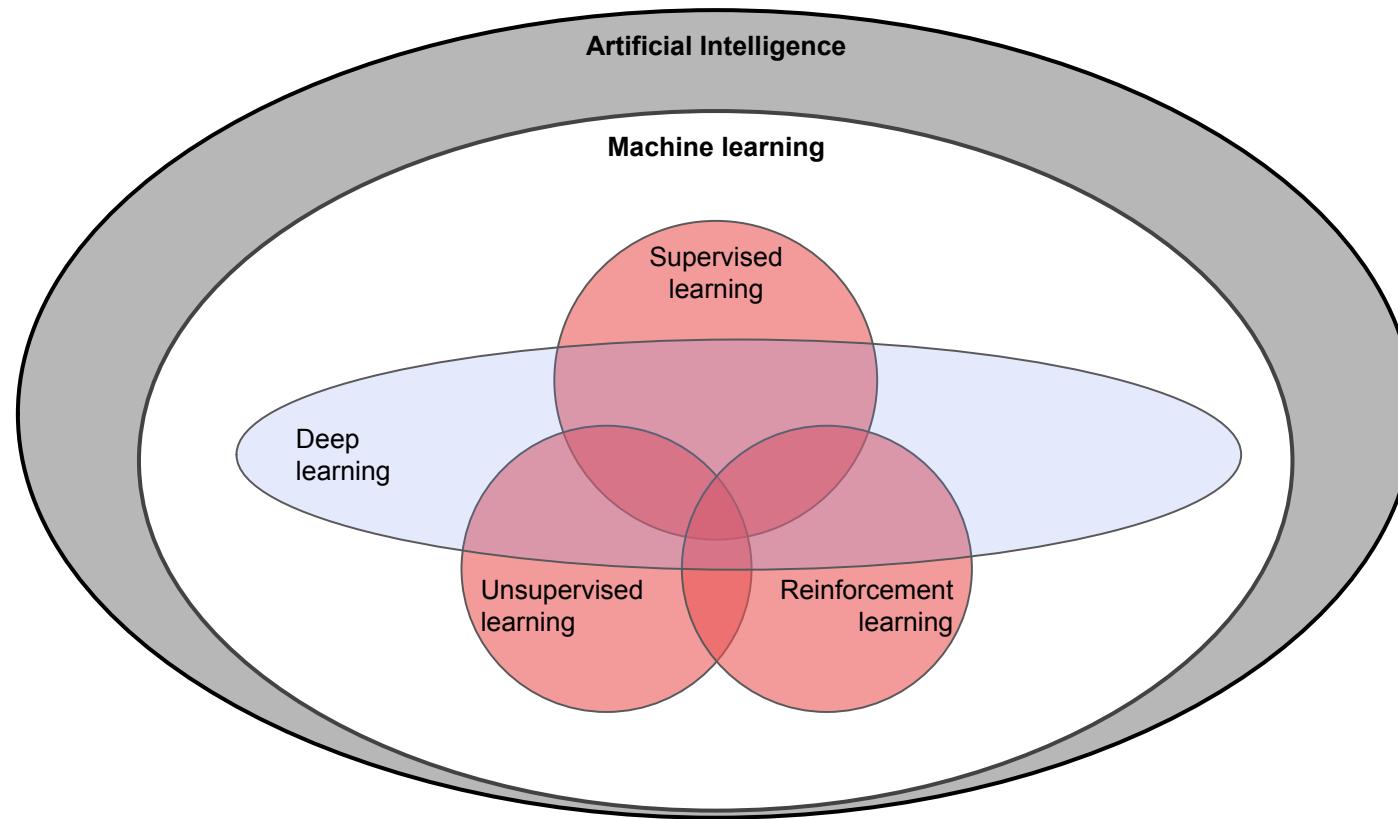
18.4

Success & Perspectives



About RL





1930's

Trial and Error learning

Psychology, animal learning

1950's

Optimal control

Value function, dynamic programming

1950's

Temporal Difference

Animal learning, Game theory



Richard Bellman
(1920-1984)

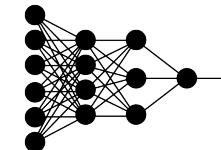
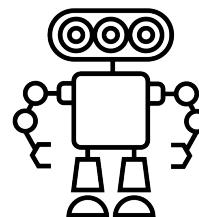


Marvin Minsky
(1927-2016)

18.1

1980's late
Modern Reinforcement Learning (Q-learning, Sarsa)

Robotic



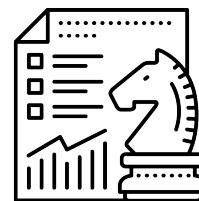
Neural Network

2013

18.2

18.3

DQN Playing Atari with Deep Reinforcement Learning



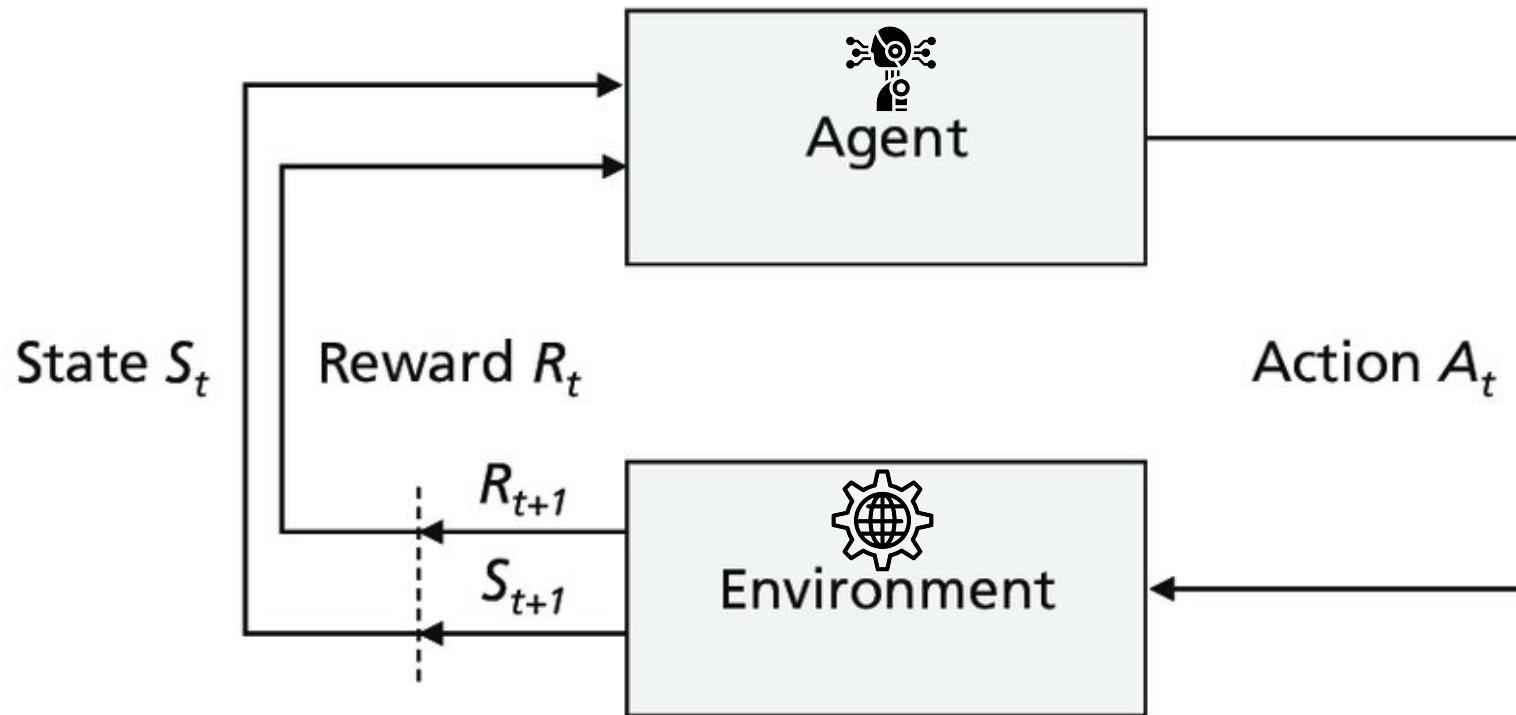
Game Theory



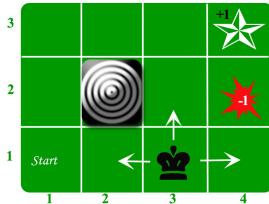
DEEPMIND



RL - Actions/Policies learning



RL - Environments



Grid world



Video Game

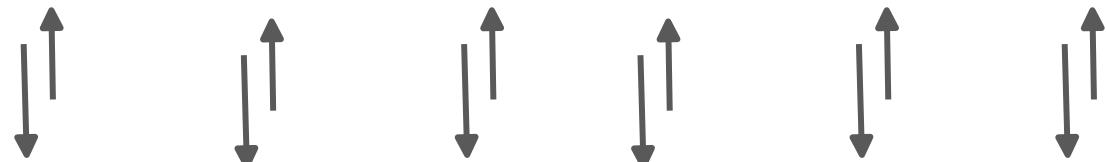


Simulator



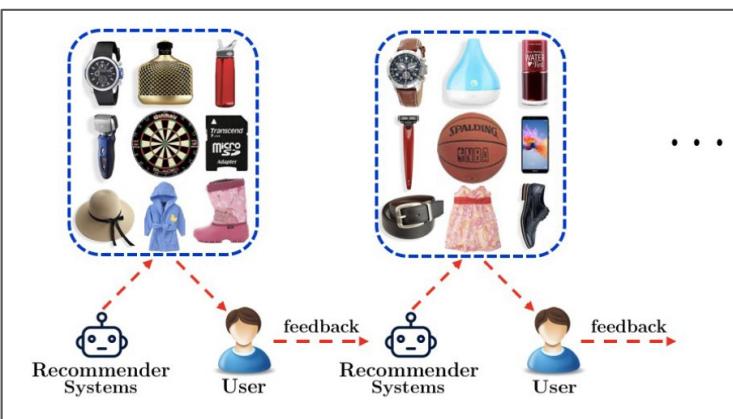
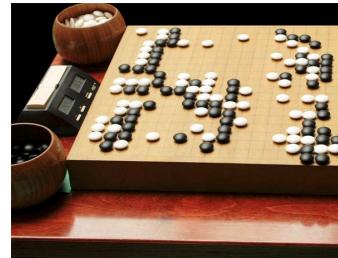
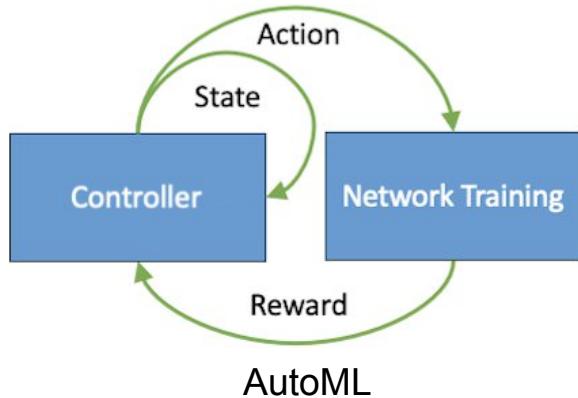
Reality

Deterministic	Fully Observable	Episodic	Static	Discrete Control	Mono-Agent
---------------	------------------	----------	--------	------------------	------------



Stochastic	Partially Observable	Sequential	Dynamic	Continuous Control	Multi-Agents Adversarial Collaborative
------------	----------------------	------------	---------	--------------------	--

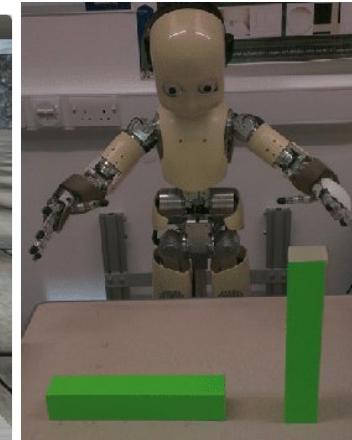
RL - Applications



RL - Limits of Simulation

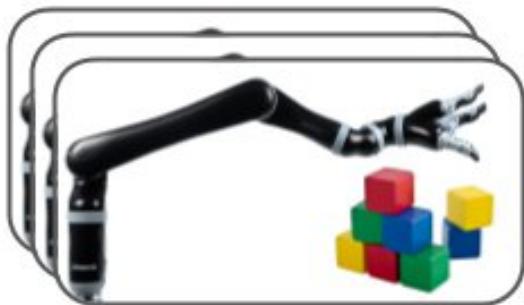
Need simulators to train model

“From robotic simulator to real world” biases



(Self-)Supervised learning

Dataset of experience

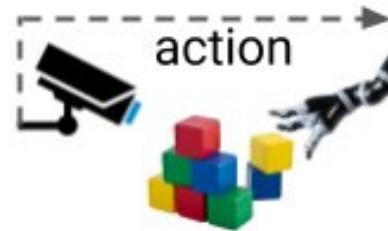


vs

Reinforcement learning

Lifelong learning

Policy



- Hard to collect a representative labeled dataset
- Hard to deal with dynamic environment

- Slow live training
- Reward function hard to design

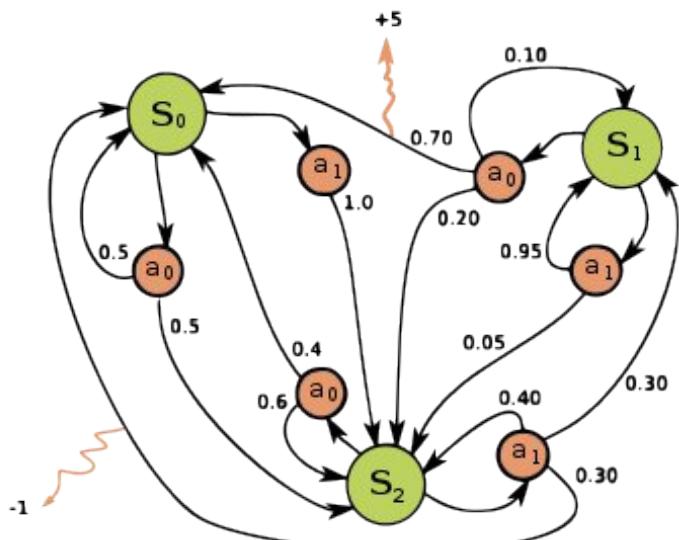
Questions Break



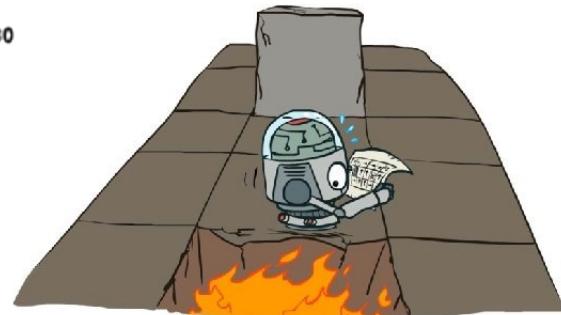
Optimal Control - MDP & Grid World



MDP: Markov Decision Process



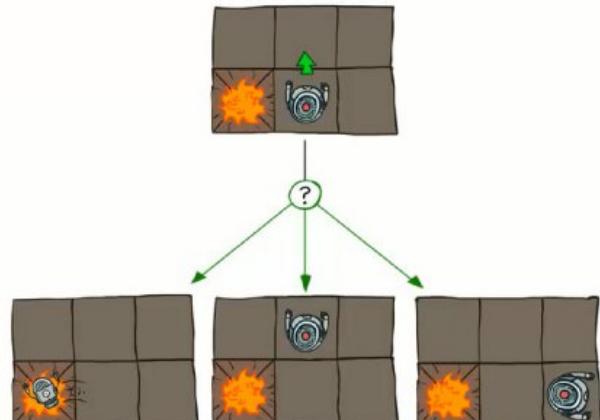
Markov chain



Deterministic Grid World



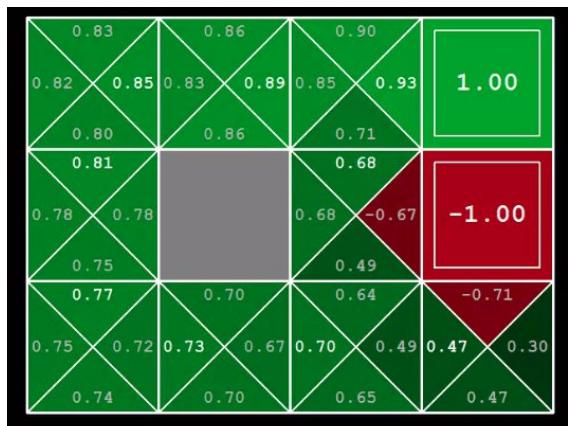
Stochastic Grid World



Optimal Control - Bellman Equations



$Q(s,a)$: Action-State Value function



Q table

$$Q_*(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a \max_{a' \in \mathcal{A}} Q_*(s', a')$$

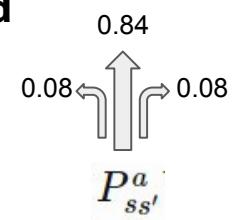
Grid world



$R(s, a)$: Reward

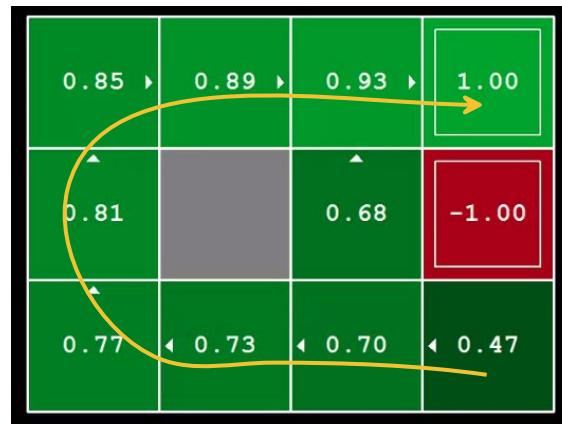
s' : next state

$\gamma = 0.96$
discount rate



$\max_a Q(s,a) = V(s)$

$V(s)$: State Value function



V table

$$V_*(s) = \max_{a \in \mathcal{A}} (R(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_*(s'))$$

Optimal Control - Discount rate



Which discount rate
you apply for
yourself?

day by day:

- 0.1
- 0.5
- 0.9

Reward trajectory

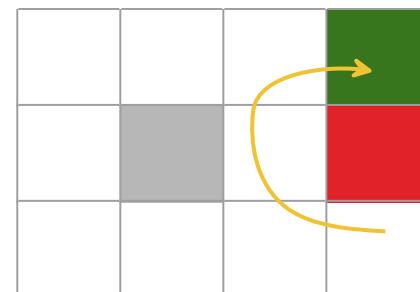
$$R(\tau) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots$$

Return: cumulative reward Gamma: discount rate

Trajectory (read Tau)
Sequence of states and actions

$$R(\tau) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$\gamma = 0.90$



1

Worth Now



γ

Worth Next Step



γ^2

Worth In Two Steps

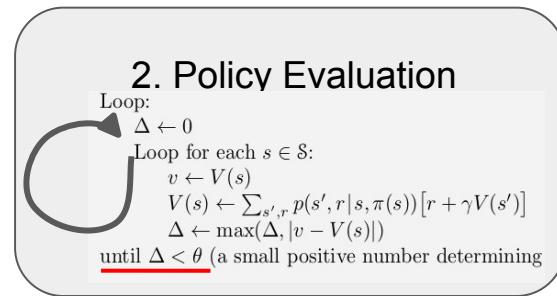


0.	0.	0.	1.00
0.		0.	-1.00
0.	0.	0.	0.

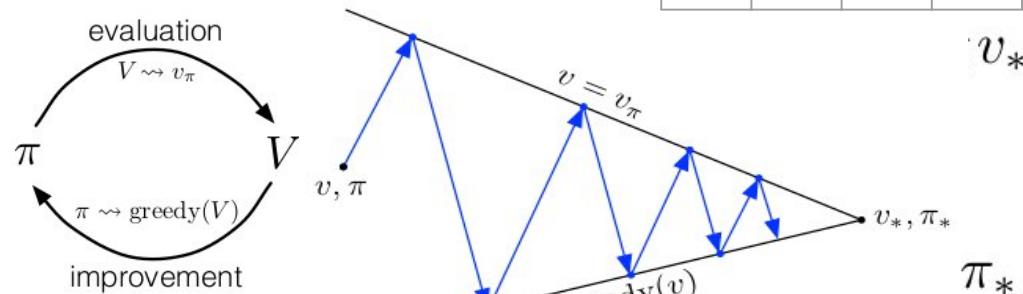
1. Initialisation

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

			1.00
			-1.00



0.85	0.89	0.93	1.00
0.81		0.68	-1.00
0.77	0.73	0.70	0.47



3. Policy Improvement

If policy stable, then stop and return

$V \approx v^*$ and $\pi \approx \pi^*$; else go to 2

			1.00
			-1.00

Dynamic Programming is **not** Machine learning and have very limited applications because it need :

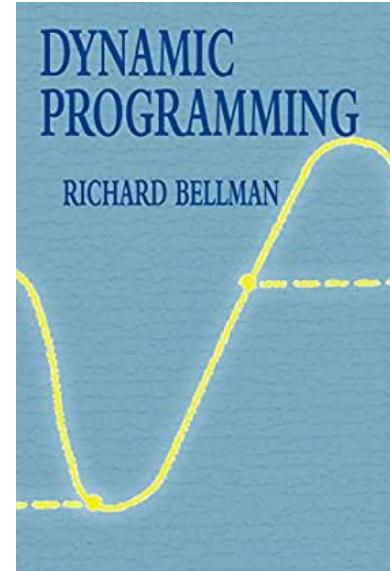
- Finite MDP : Static, Episodic
- Perfect known Model : **Fully Observable**



Is computational expensive if environment is large



Then, Let's go back to the Reinforce Learning world : like an embodied agent !!



"... breaks problem into simpler steps at different points in time."



- **Monte Carlo** G_t : reward trajectory until the **game over** $= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

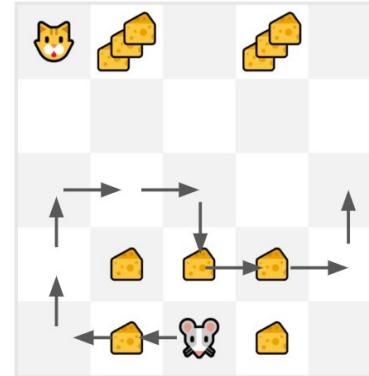
New value of state t

Former estimation of value of state t (= Expected return starting at that state)

Learning Rate

Return at timestep t

Former estimation of value of state t (= Expected return starting at that state)



- **Temporal Difference**

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

New value of state t

Former estimation of value of state t

Learning Rate

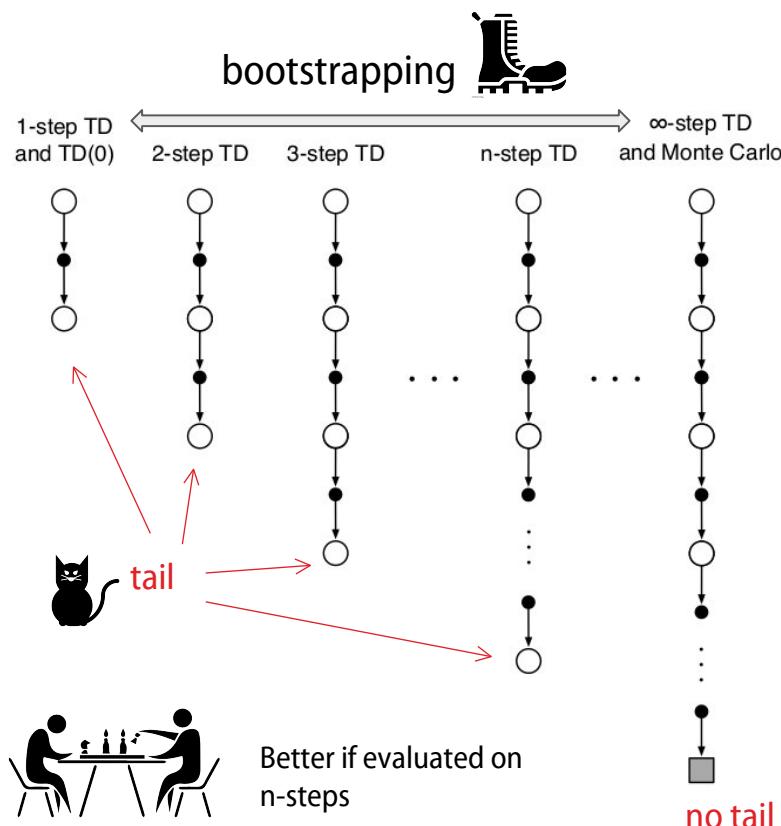
Reward
Discounted value of next state

TD Target





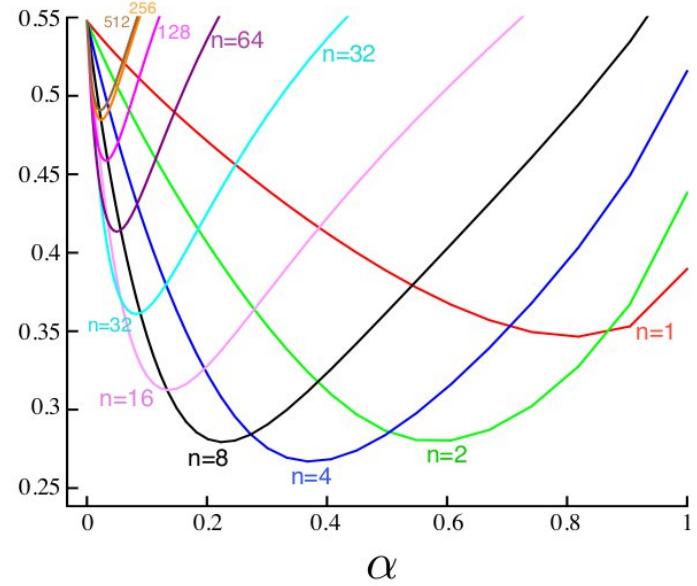
Temporal Difference(n)



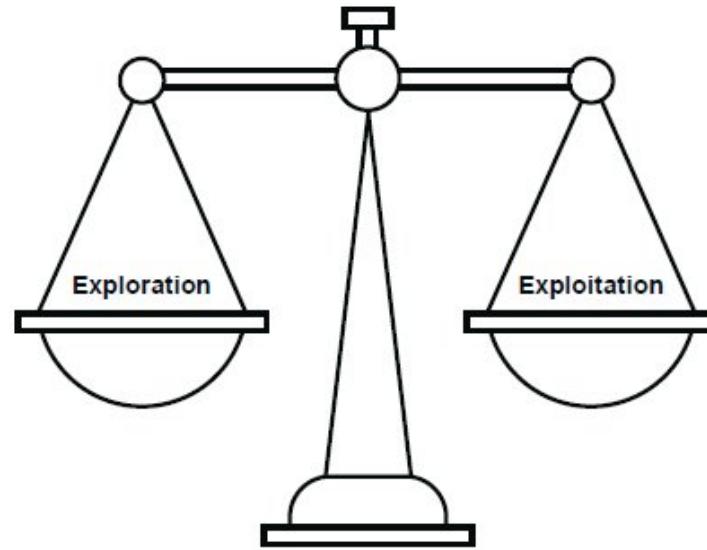
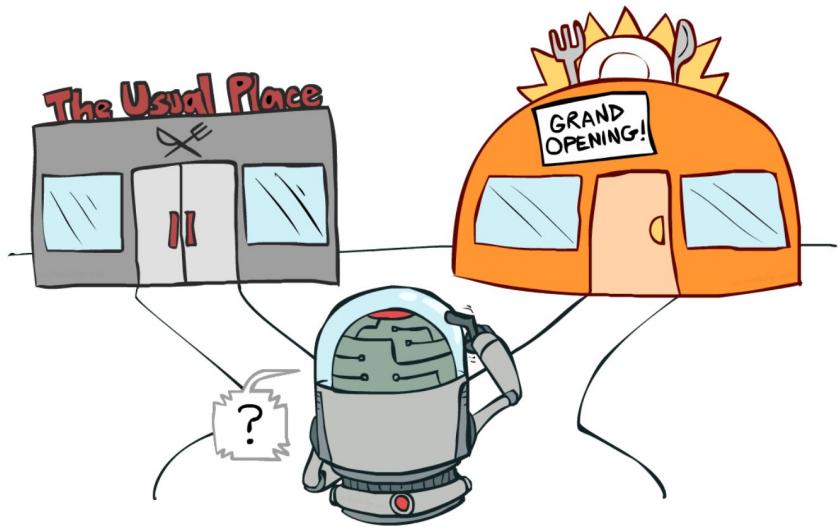
$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \boxed{\gamma^n V_{t+n-1}(S_{t+n})}$$

n-step TD Target

Average RMS error over 19 states and first 10 episodes

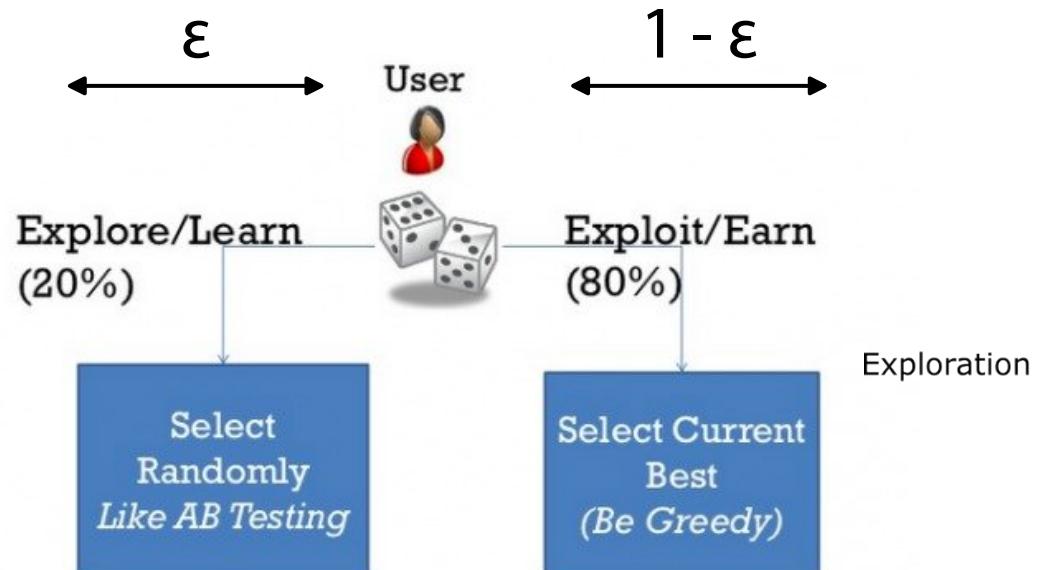


Trial and Error Learning

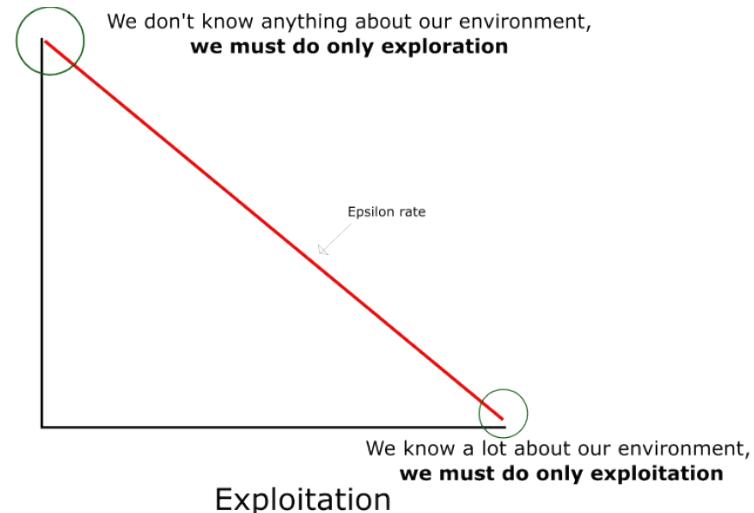


Exploration / Exploitation Trade-off

Trial and Error Learning ϵ -greedy

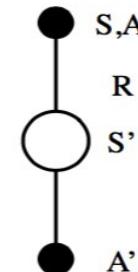


Decrease epsilon in time !!





$2 \times \varepsilon\text{-greedy}!$



Initialize Q table



Choose an action a



Perform action



Measure reward



Update Q

Choose A from S using policy derived from Q (e.g., ε -greedy)

At the end of the training

Good Q*table

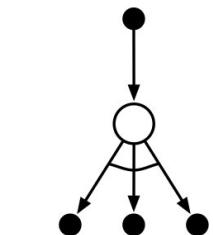
Choose A' from S' using policy derived from Q (e.g., ε -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$$S \leftarrow S'; A \leftarrow A';$$



1 x ϵ -greedy !



Q-learning

Initialize Q table



Choose an action a

Choose A from S using policy derived from Q (e.g., ϵ -greedy)



Perform action



Measure reward



Update Q

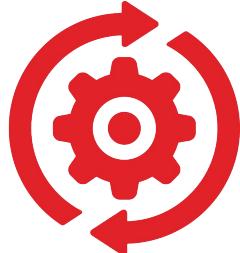
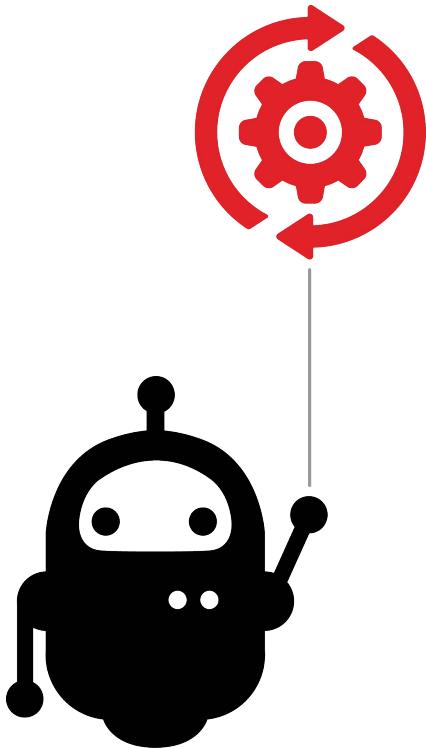
$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

$$S \leftarrow S'$$

At the end of the training

Good Q*table

Example : Cliff Walking



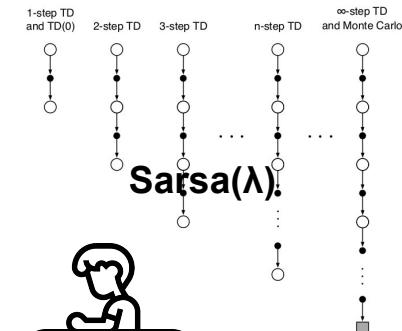
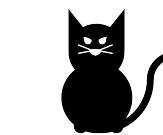
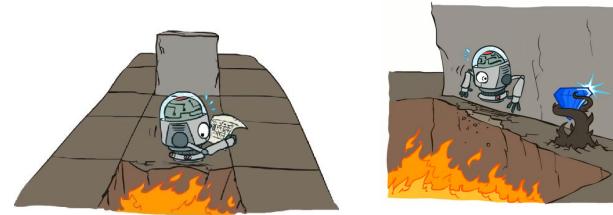
[RL1] Cliff Walking : Sarsa vs Q-learning



- Q-learning: Off-Policy ($1 \times \varepsilon$ -greedy)
- SARSA: On-Policy ($2 \times \varepsilon$ -greedy)

To have in Mind

- Bellman's Value function, Q (value-action) function and discount rate
- Exploration strategy : ϵ -greedy
- Temporal Difference, TD(n) with tail, and Monte Carlo without tail
- On-policy, Off-policy concepts



Tabular RL resolved environments



Deterministic	Fully Observable	Episodic	Static	Discrete Control	Mono-Agent
---------------	------------------	----------	--------	------------------	------------

OK
Bellman's equations **OK**
Trajectory Exploration **OK**
TD Learning **OK**
Lifelong Exploration **NOK** **NOK**

Stochastic	Partially Observable	Sequential	Dynamic	Continuous Control	Multi-Agents Adversarial Collaborative
------------	----------------------	------------	---------	--------------------	--

Questions Break

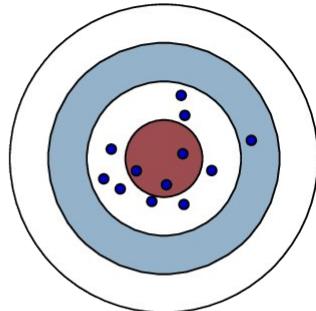


From RL to Deep RL

Variance vs Bias trade-off



High Variance



More Explorative
Slower Learning
Hard convergence



On-policy

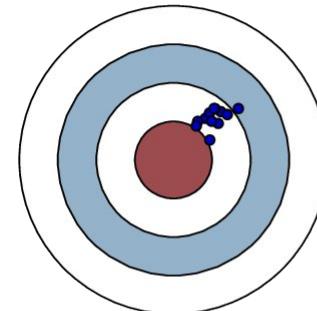
Monte Carlo

Actor

Reinforcement Learning Trade-off



High Bias

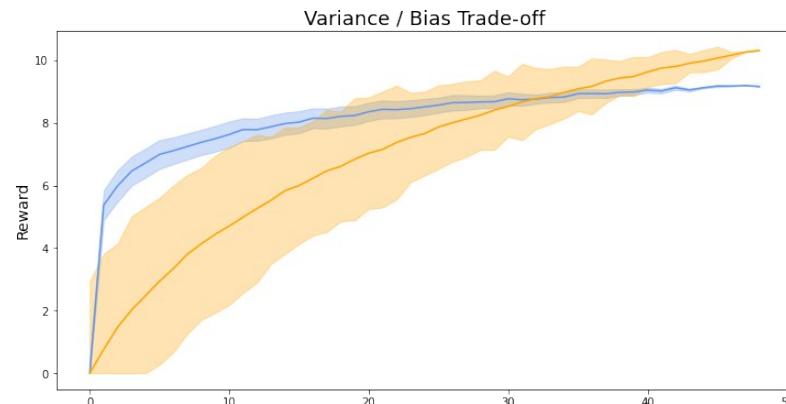


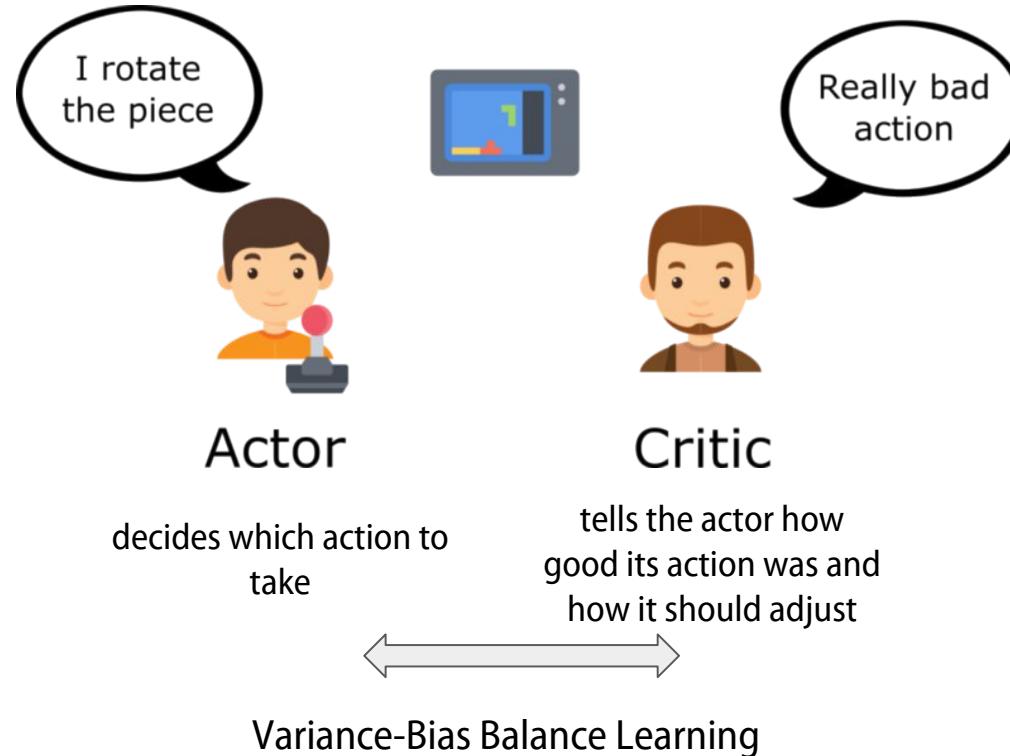
More Exploitative
Faster Learning
Easy convergence



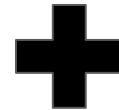
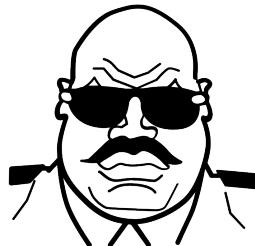
TD Learning

Critic





Advantage function



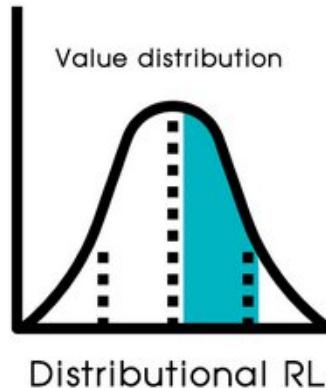
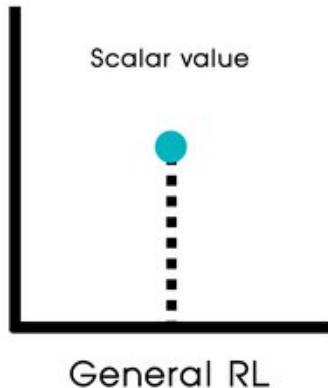
$Q(s, a)$
Value-Action function
Q guevara

$V(s)$
Value function
Strategist
Visionary
Critic
How reach the best state at
the end ?

$A(s, a)$
Advantage function
Tactician
Actor
From every state, what
is the best action ?



Distributional RL

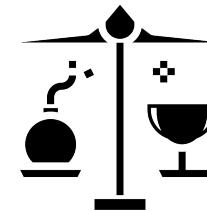


Bellman Equation: $Q(s, a) = R(s, a) + \gamma Q(s', a')$

Distributional Bellman Equation: $Z(s, a) = R(s, a) + \gamma Z(s', a')$

Distributional Q-learning evaluate:

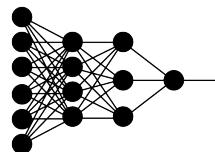
- the expected reward
- **and the risk to reach it**



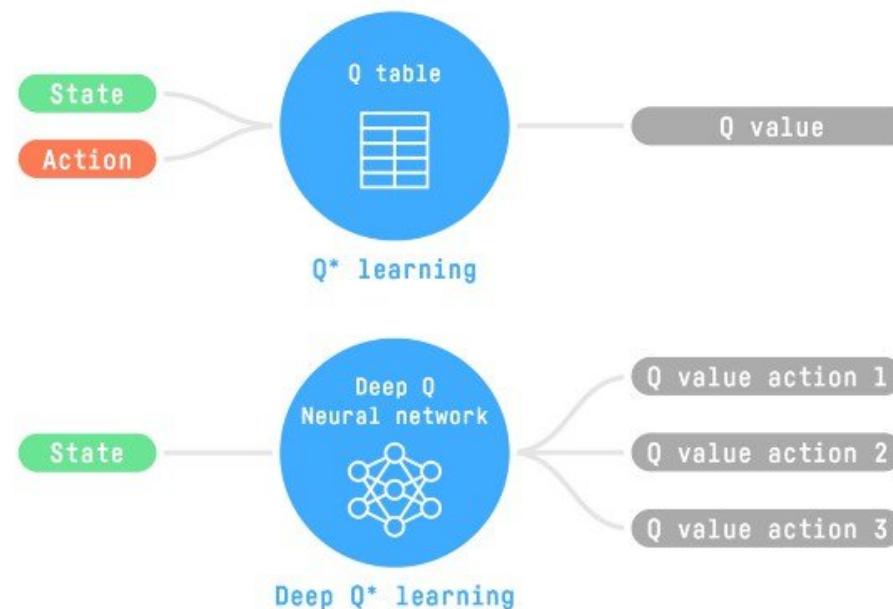
Q Tabular / Approximation Q function

Q table is a nonlinear
function...

!?



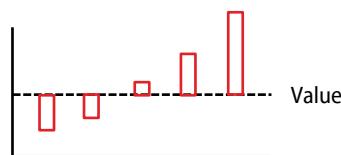
Why don't use a neural
network ?



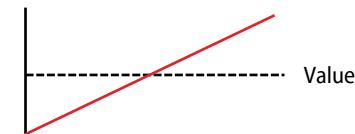
Approximation function benefit

Imagine this state-value function $V=[-2.5, -1.1, 0.7, 3.2, 7.6]$

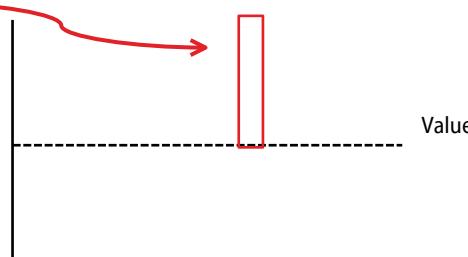
With Q-table each value is independent.



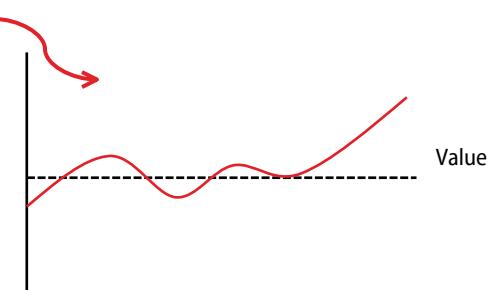
With function approximation, the underlying relationship of the states can be learned and exploited



With Q-table, the update only changes one step.

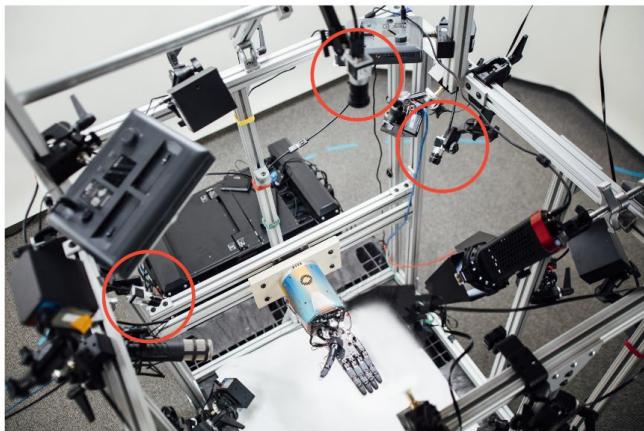


With function approximation, the update changes multiple steps.



Deep learning perspective

Tabular representation is limited. But Deep learning enable to use **high dimensional data** : Image, text, ...



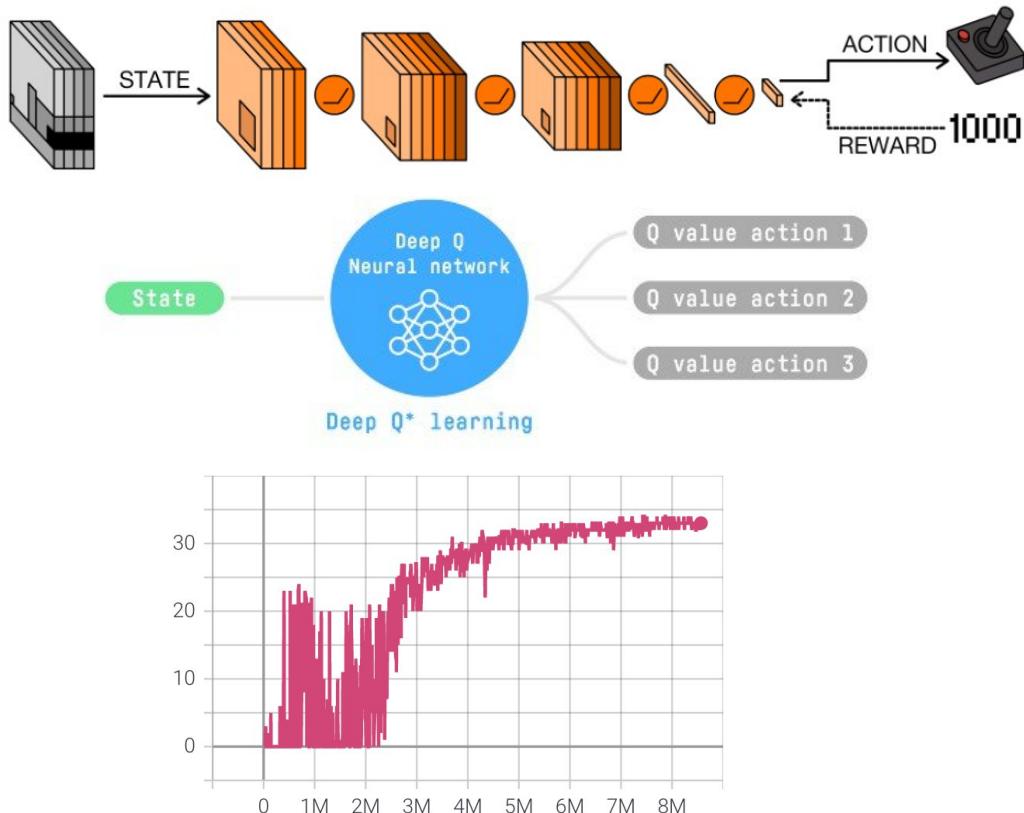
DQN : Deep Q-learning Network

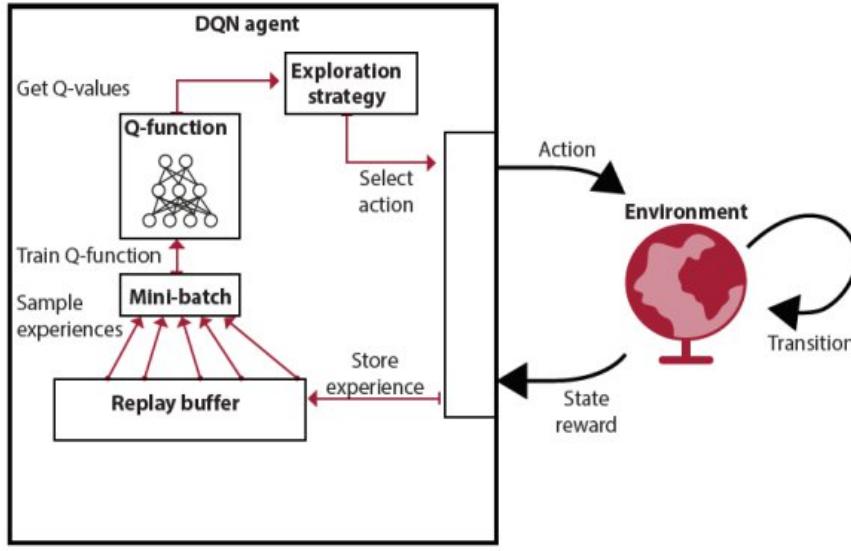


Playing Atari with Deep Reinforcement Learning (2013)



Freeway Game





Initialize Network

Initialize network Q
 Initialize target network \hat{Q}
 Initialize experience replay memory D
 Initialize the *Agent* to interact with the Environment

Choose an action a

using policy ϵ -greedy(Q)

Perform action

Measure reward

Update Q

Store transition $(s, a, r, s', done)$ in the experience replay memory D

Calculate the loss $\mathcal{L} = 1/N \sum_{i=0}^{N-1} (Q(s_i, a_i) - y_i)^2$
 Update Q using the SGD algorithm by minimizing the loss \mathcal{L}

Q-label

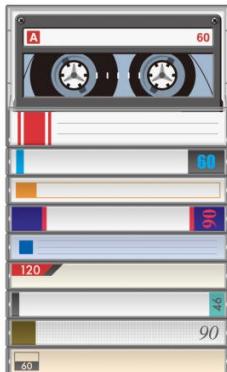
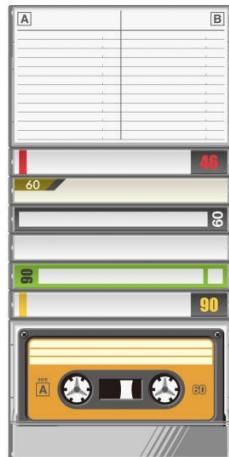
$$y_i = r + \gamma \max_a Q(s', a)$$

At the end of the training





Experience Replay Buffer



Experience replay Buffer

Sample



Batch of experiences

→

DQN
Agent

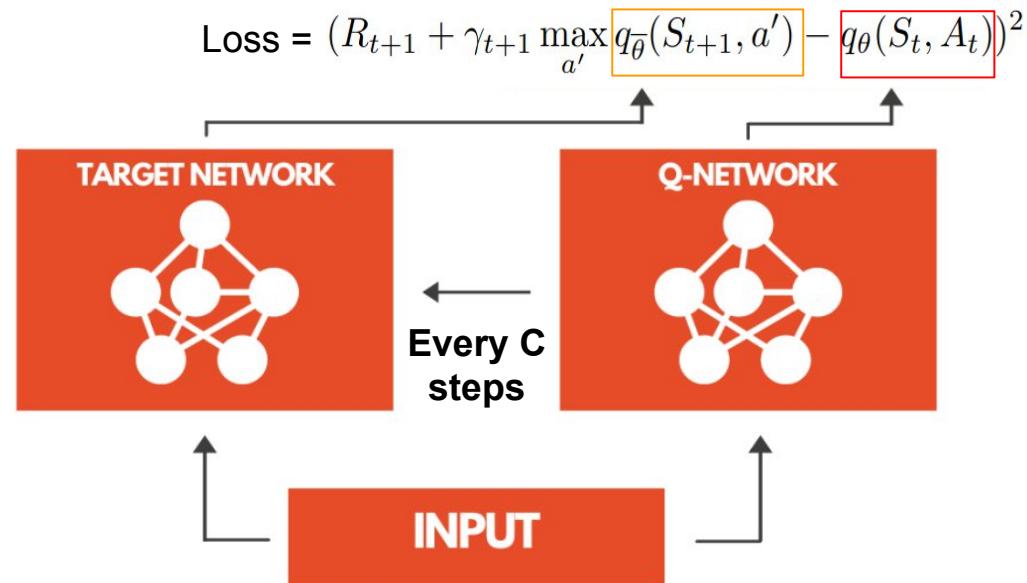
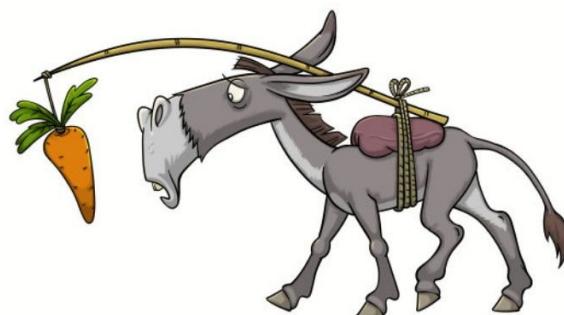
Random sampling





Off-policy

Target network / online network





Prioritizing Experience Replay Buffer



Prioritize the examples with larger TD error, and so with more information

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha} \quad \text{Probability sampling}$$

$$p_i = \frac{1}{\text{rank}(i)} \quad \text{rank according to how big is the TD error}$$

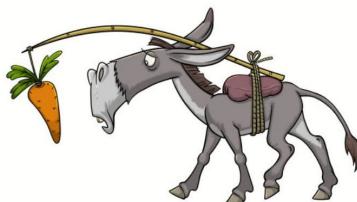
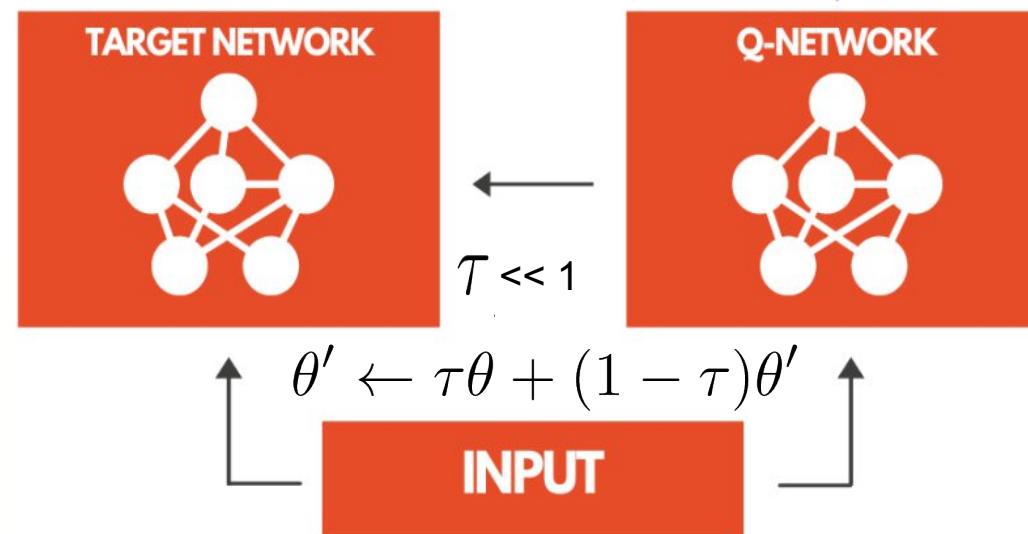
Double DQN



To fix Q value overestimation issue
and to smooth the learning

$$\text{Loss} = (R_{t+1} + \gamma_{t+1} \max_{a'} q_{\theta}(S_{t+1}, a') - q_{\theta}(S_t, A_t))^2$$

$$\text{Loss} = (R_{t+1} + \gamma_{t+1} q_{\bar{\theta}}(S_{t+1}, \operatorname{argmax}_{a'} q_{\theta}(S_{t+1}, a')) - q_{\theta}(S_t, A_t))^2$$

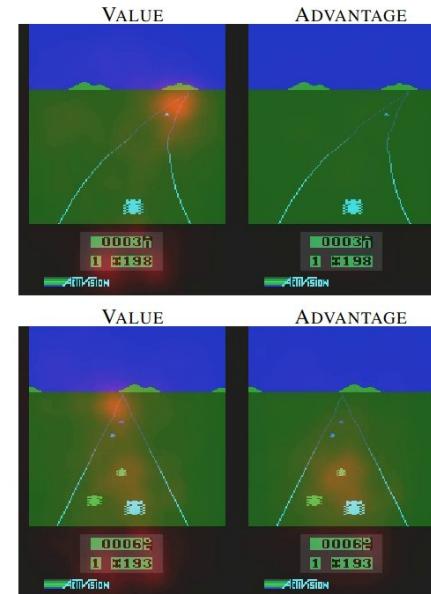
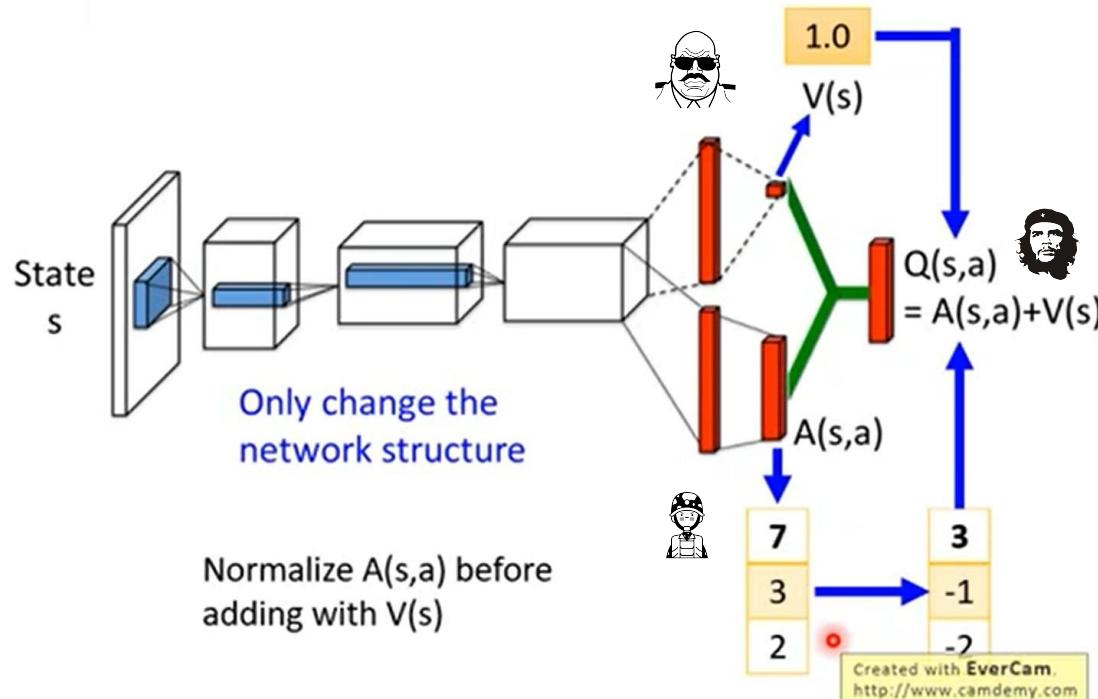


Dueling-DDQN



Off-policy

$$Q(s, a) = V(s) + A(s, a)$$



The value stream learns to pay **attention to the road**.

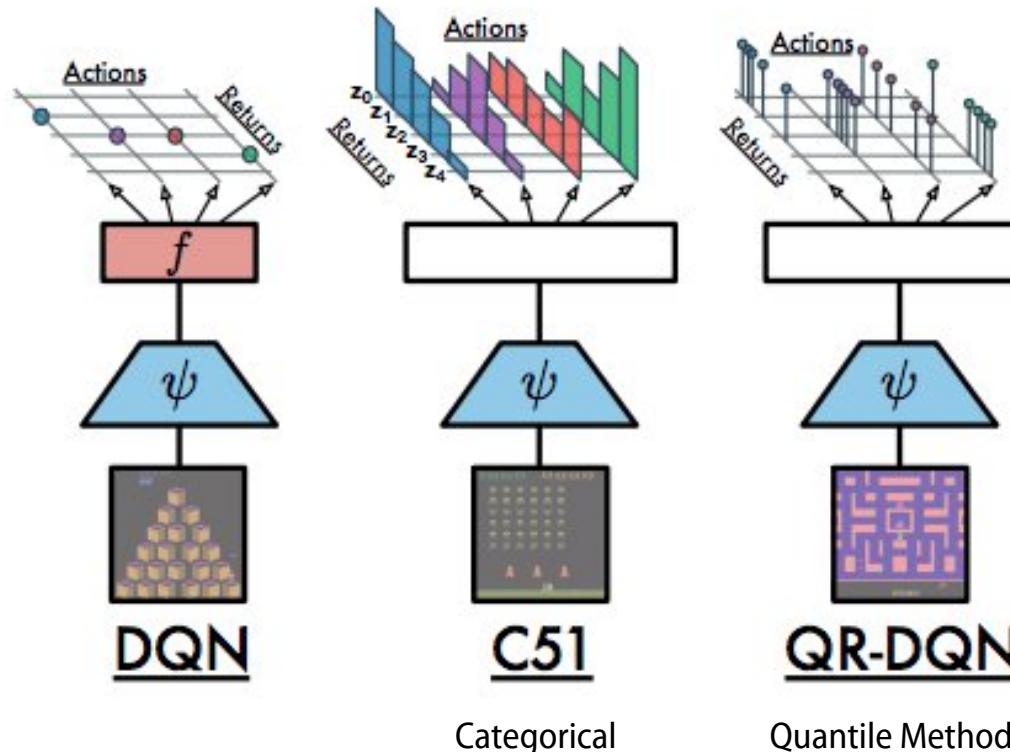
The advantage stream learns to pay attention only when there are cars immediately in front, so as to **avoid collisions**.

Distributional DQNs



Distributional RL aims to model the distribution over returns...

Distributional Bellman Equation: $Z(s, a) = R(s, a) + \gamma Z(s', a')$



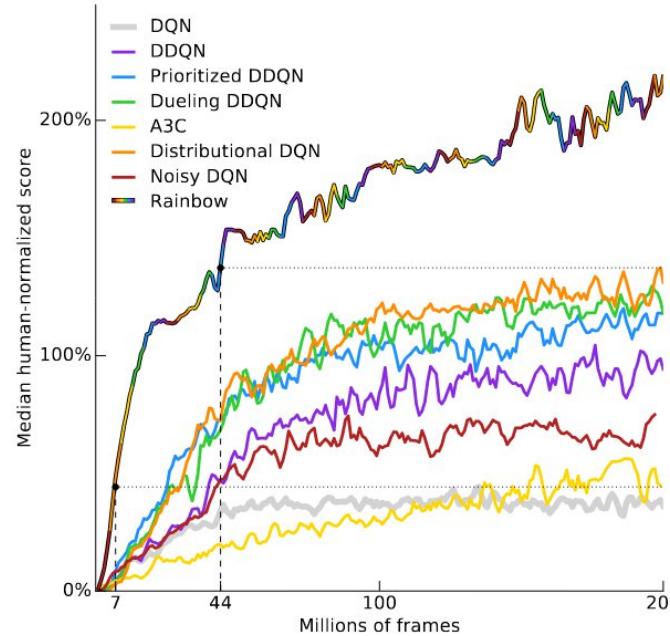


<https://arxiv.org/pdf/1710.02298.pdf>

- DQN
- DDQN
- Prioritized DDQN
- Dueling DDQN
- Distributional DQN
- **Noisy DQN** (Add noise for exploration)

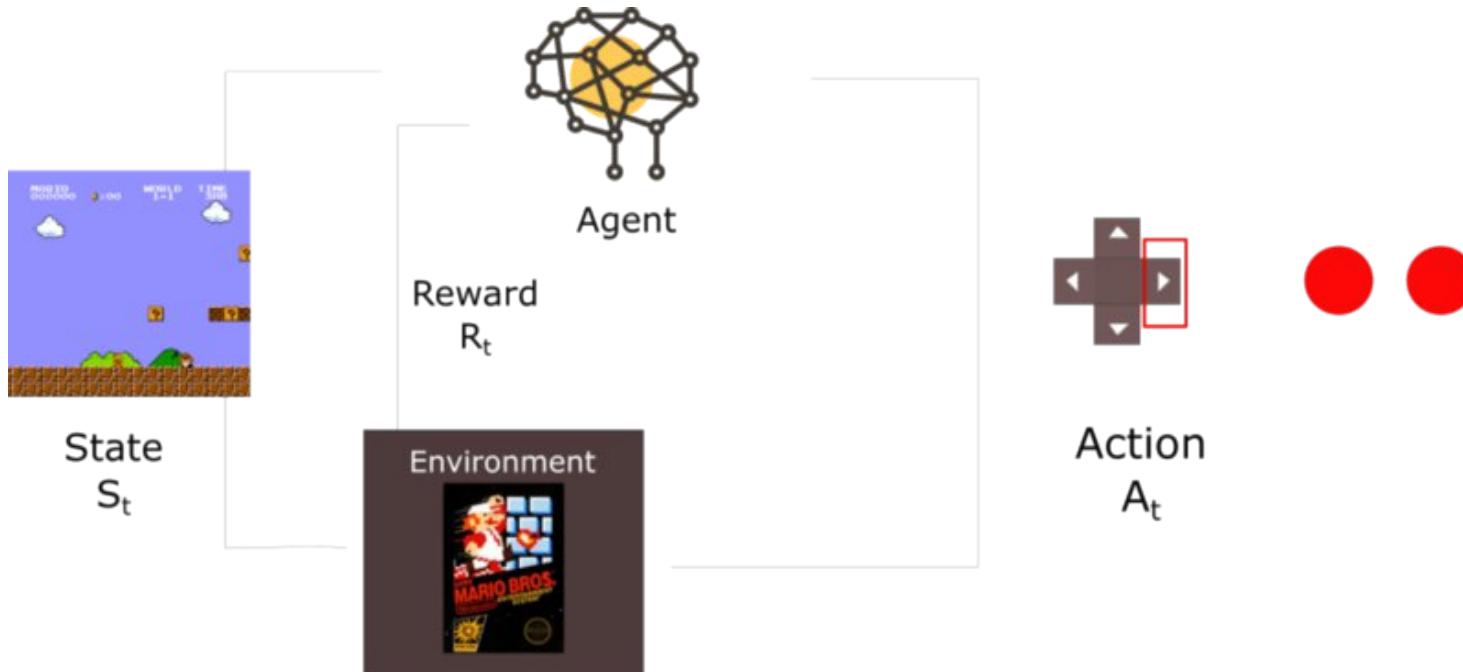


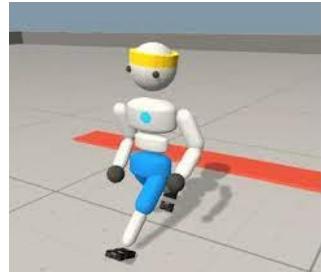
Rainbow DQN



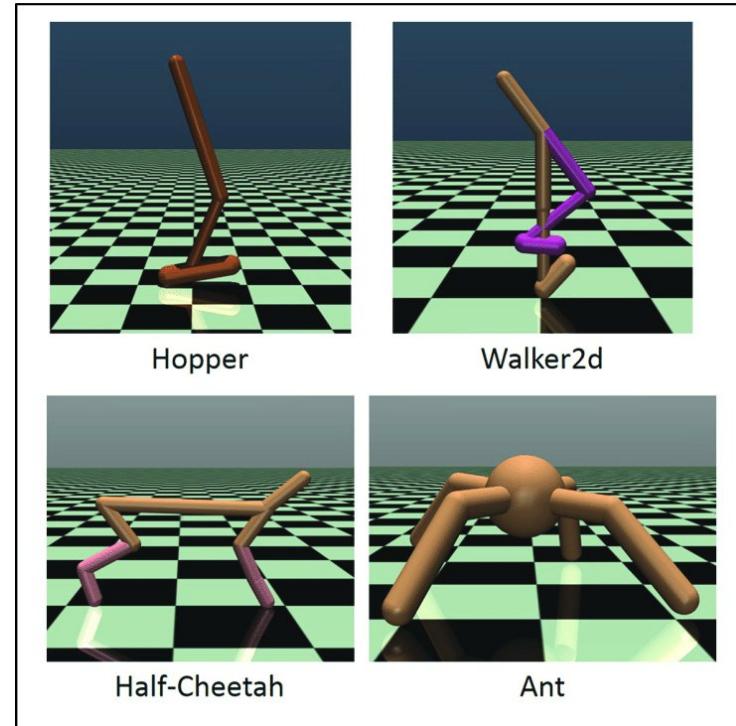
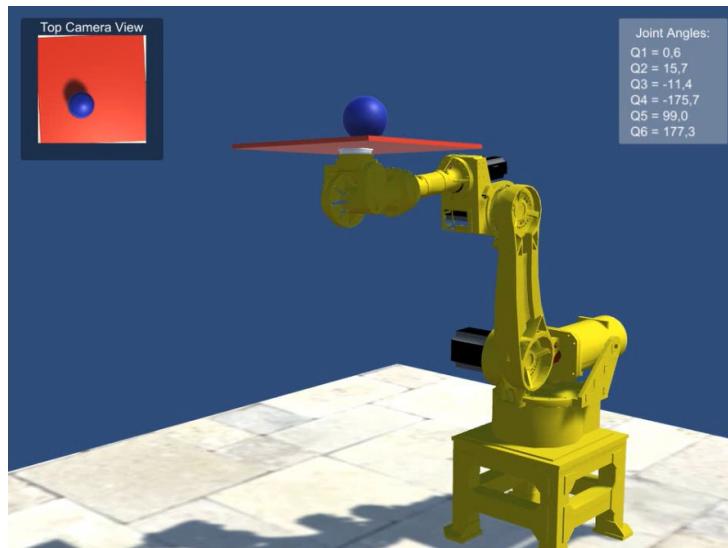
Median human-normalized performance
across 57 Atari games

Environment's actions space must be **discrete** !!!





Environment's actions space **cannot be** continuous !!!



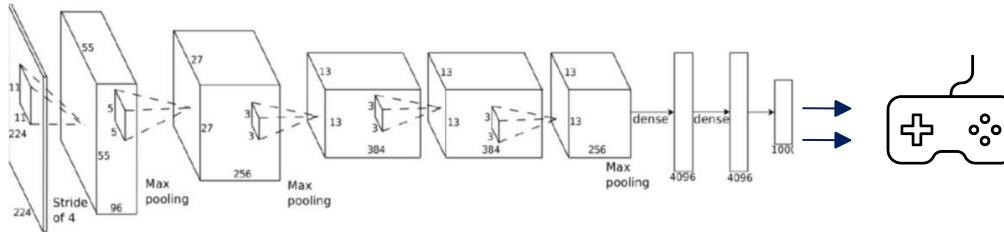
Questions Break



Policy Gradient



=> Action Inference



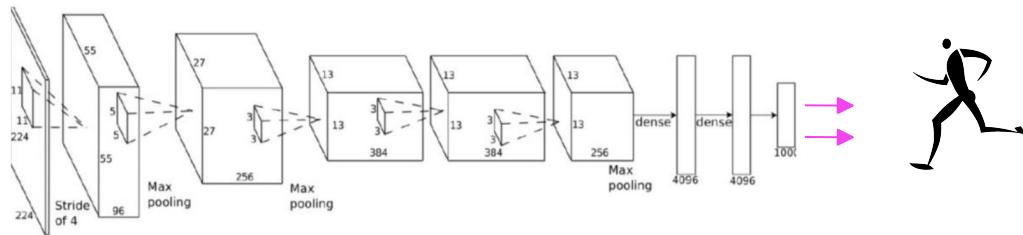
Discrete control

last activation: Softmax, Sigmoid

Stochastic policy

Probability of taking that action

or



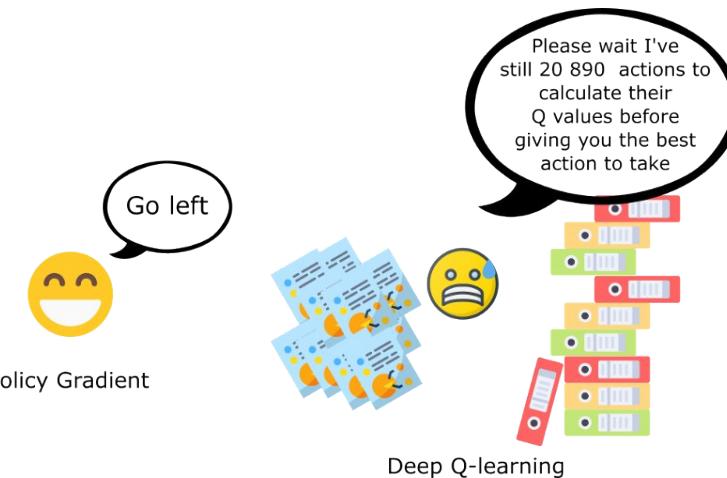
Continuous control

last activation : Sigmoid [0;1] or Tanh [-1;1]

Deterministic policy

directly the accurate action value

Policy Gradient



Policy gradient **simpler** (straightforward) than Q or V based method
Easy to use for **continuous control & high-dimensional action spaces**

⚠ **Exploration** is few efficient especially for deterministic/continuous PG (no ϵ -greedy)

Loss Objective function :

$$J(\theta) = E_{\tau \sim \pi}[R(\tau)]$$

$$R(\tau) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \dots$$

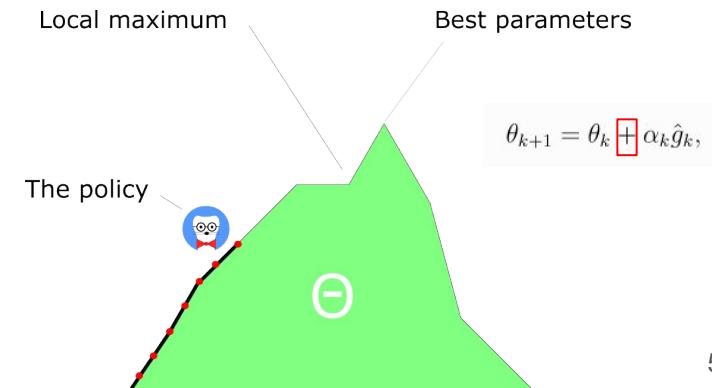
Return: cumulative reward Gamma: discount rate

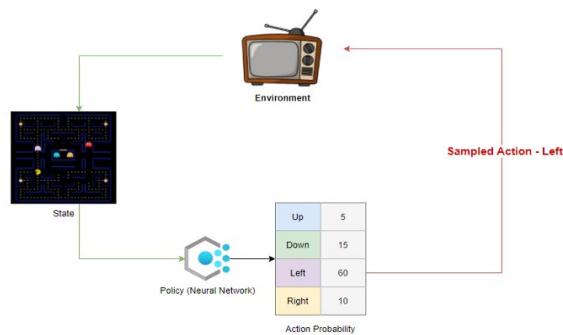
Trajectory (read Tau)
Sequence of states and actions

Gradient Ascent

= negative gradient descent

(**-obj**).backward()
optimizer.step()





Good stuff is made more likely or more intensively
 Bad stuff is made less likely or less intensively

$$\text{Policy gradient : } E_{\pi}[\nabla_{\theta}(\log \pi(s, a, \theta)) R(\tau)]$$

Policy function Score function

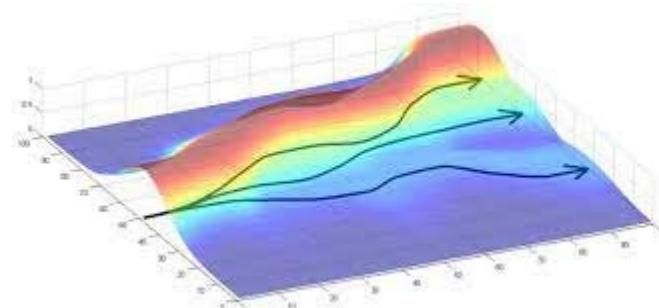
$$\text{Update rule : } \Delta \theta = \alpha * \nabla_{\theta}(\log \pi(s, a, \theta)) R(\tau)$$

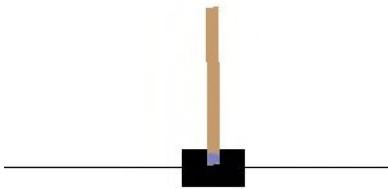
Change in parameters Learning rate

(See Policy Gradient Theorem)

Monte-Carlo policy gradient

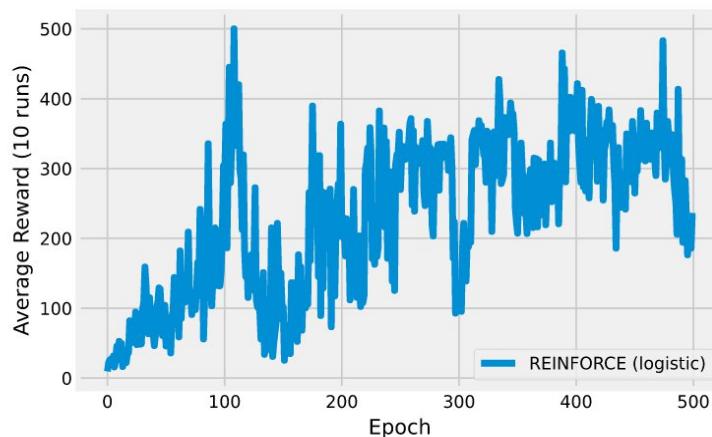
1. Initialize the policy parameter θ at random.
2. Generate one trajectory on policy π_{θ} : $S_1, A_1, R_2, S_2, A_2, \dots, S_T$.
3. For $t=1, 2, \dots, T$:
 1. Estimate the return G_t ;
 2. Update policy parameters: $\theta \leftarrow \theta + \alpha \gamma^t G_t \nabla_{\theta} \ln \pi_{\theta}(A_t | S_t)$





CartPole
Discrete or Continuous
environment

Big issues :
Huge variance & Local maxima !!





reduce the variance of the policy gradient by using an estimate of the **advantage function**, instead of the actual return



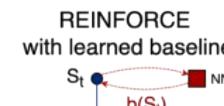
$$\underline{L^{PG}(\theta)} = \underline{E_t}[\log \pi_\theta(a_t|s_t) * \underline{A_t}]$$

Policy Loss Expected log probability of taking that action at that state

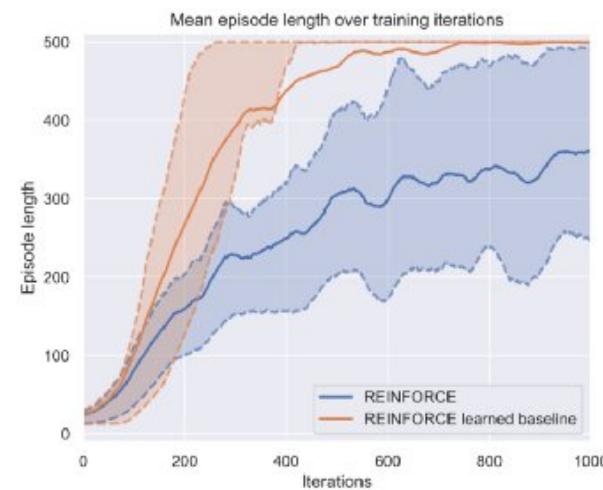
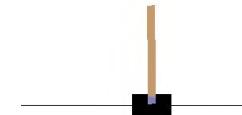
Advantage if $A > 0$, this action is better than the other action possible at that state

Instead of waiting until the end of the episode as we do in Monte Carlo REINFORCE, **we make an update at each step (TD Learning)**.

TD Learning needs a tail !!



= Vanilla Policy Gradient (VPG)



VPG = Actor-Critic



I rotate
the piece

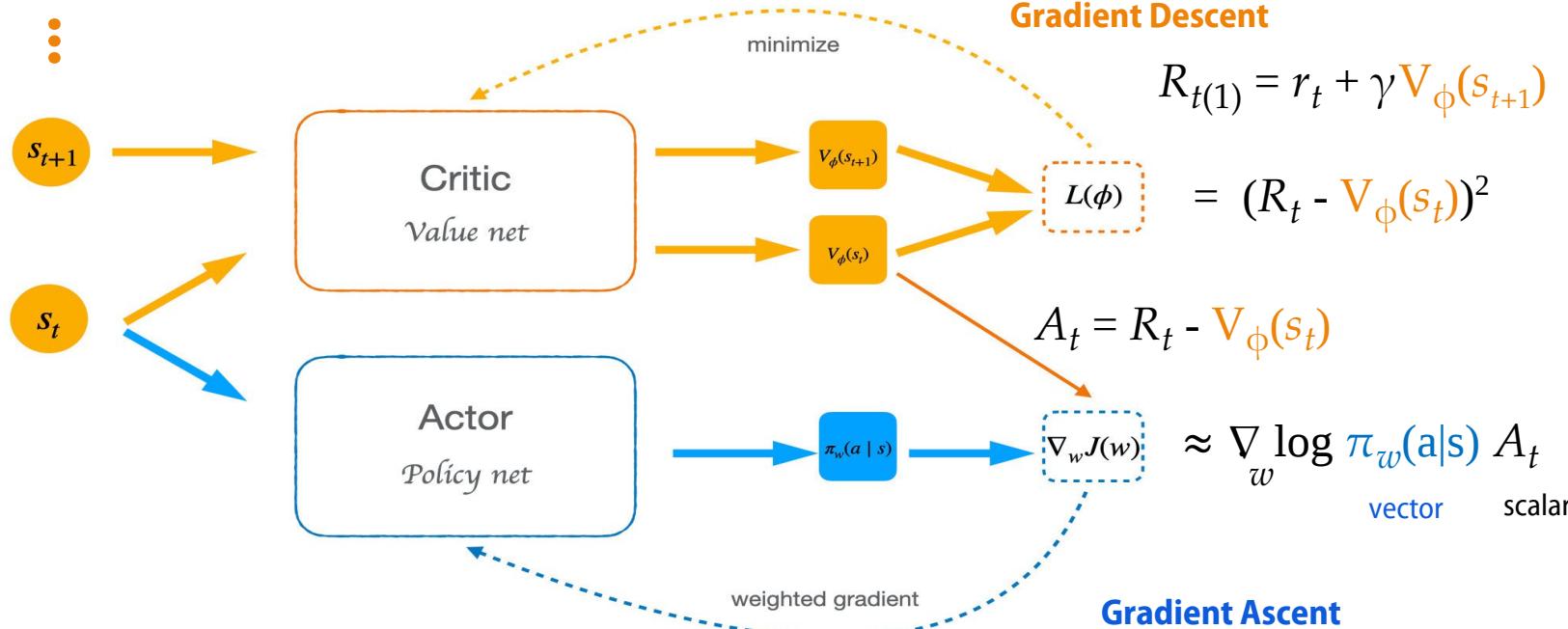
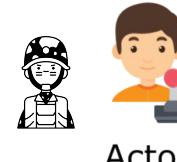


Really bad
action

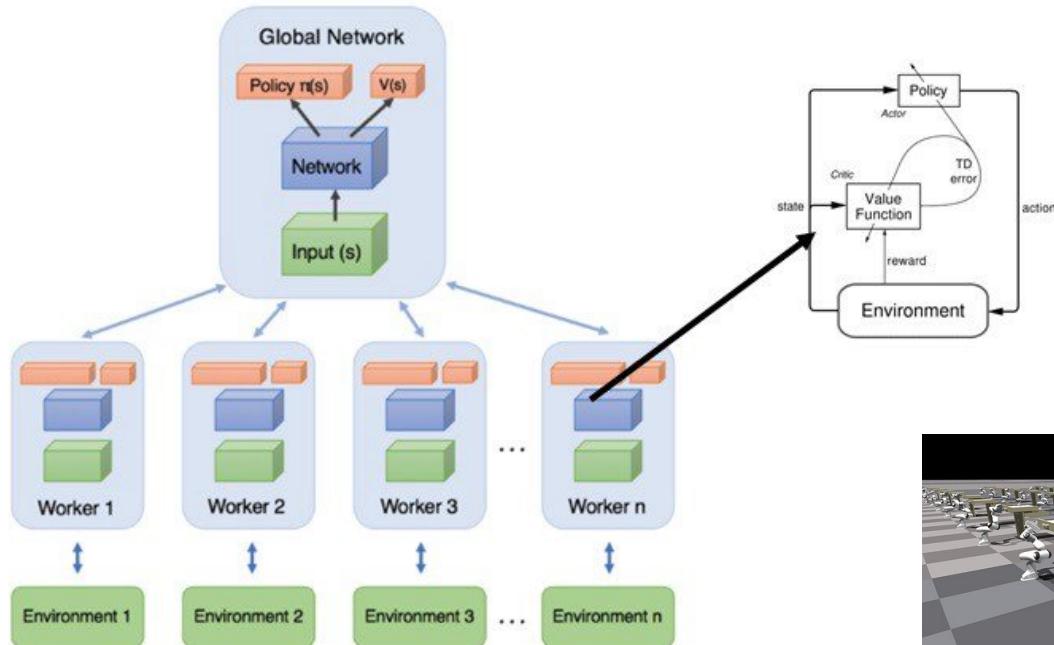


TD(n-step)

$$R_{t(n)} = \sum_{l=0}^{n-1} \gamma^l r_{t+l} + \gamma^n V_\phi(s_{t+n})$$



A3C - Asynchronous Advantage Actor-Critic



like REINFORCE with Baseline (VPG) :



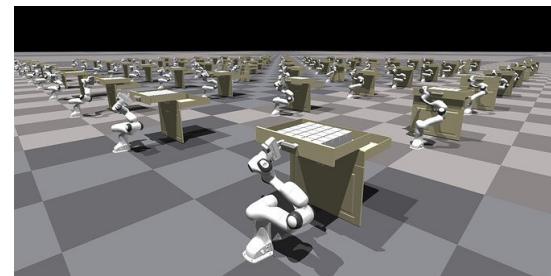
- **TD(n-step) Actor Critic**



- using **Advantage function**

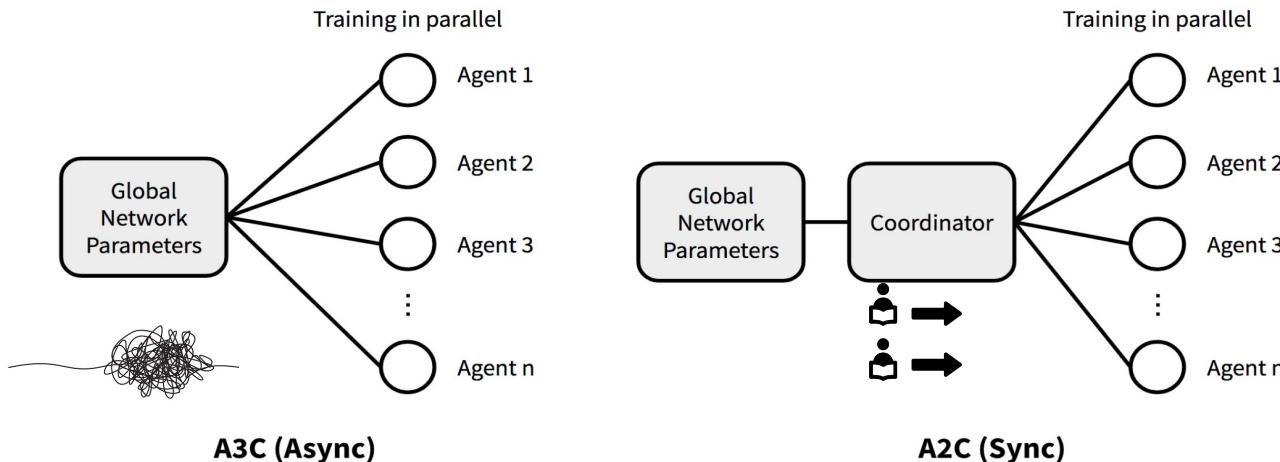


+ Asynchronous **multi-workers** for exploration





- Asynchronous Synchronous multi-workers for exploration



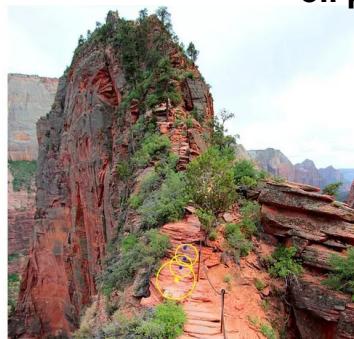
Sometimes the thread-specific agents would be playing with **policies of different versions**.
Update would not be optimal.

The synchronized gradient update keeps **the training more cohesive** and potentially to make **convergence faster**.

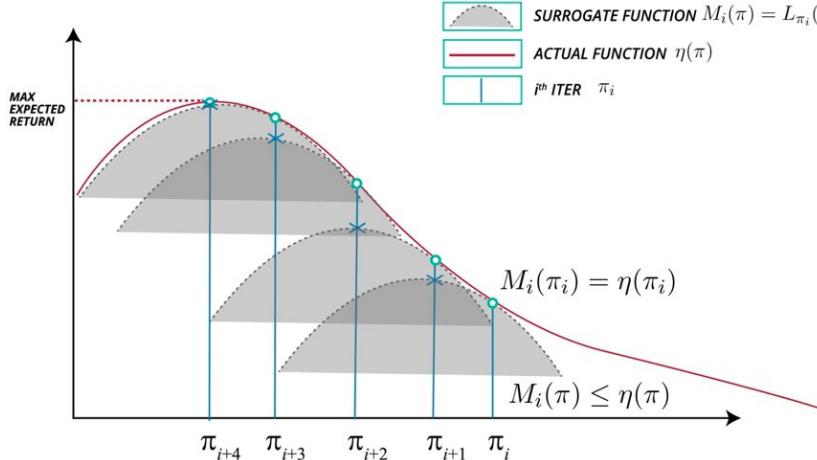
Trust Region Policy Optimization - TRPO



Line search
(like gradient ascent)



Trust region



Maximize

$$\frac{L^{PG}(\theta)}{\text{Policy Loss}} = \hat{\mathbb{E}}_t \left[\log \pi_\theta(a_t | s_t) \hat{A}_t \right].$$

Expected log(probabilities) from the output of the policy network Estimate of the relative value of selected action



“destructively large policy updates.”



The Trust Region Constraint

$$\begin{aligned} & \max_{\pi'} \mathcal{L}_{\pi}(\pi') \\ \text{s.t. } & \mathbb{E}_{s \sim d^{\pi}} [D_{KL}(\pi' || \pi)[s]] \leq \delta \end{aligned}$$

trust region
kullback-leibler divergence

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$

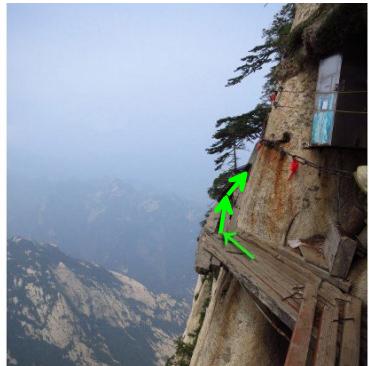
natural policy gradient

Computationally expensive with big model but use conjugate gradient algorithm to approximate the natural policy gradient !!

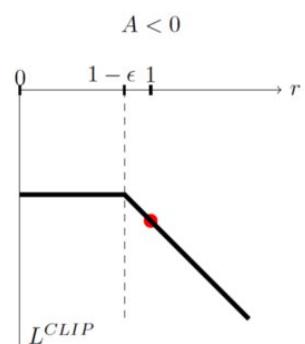
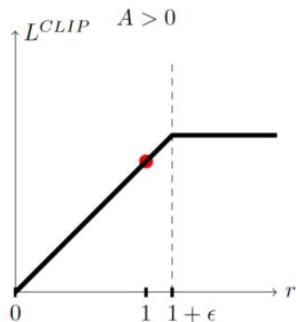
Proximal Policy Optimization - PPO



Line search



Clipped search



Maximize

$$\frac{L^{PG}(\theta)}{\text{Policy Loss}} = \hat{\mathbb{E}}_t \left[\log \pi_\theta(a_t | s_t) \hat{A}_t \right].$$



Expected log(probabilities) from the output of the policy network

Estimate of the relative value of selected action

“destructively large policy updates.”



The Clipped Surrogate Objective Function

$$L^{PG}(\theta) \approx L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

$$L_t^{CLIP+VF+S}(\theta) = \hat{\mathbb{E}}_t \left[L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_\theta](s_t) \right],$$

clipped PG objective MSE of value function entropy term

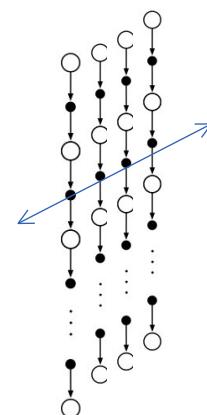
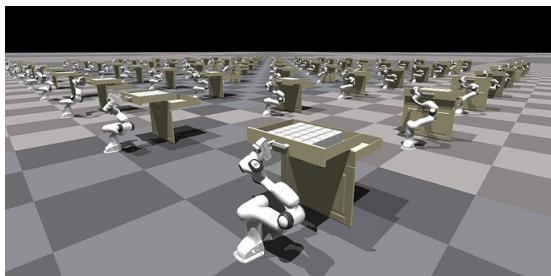
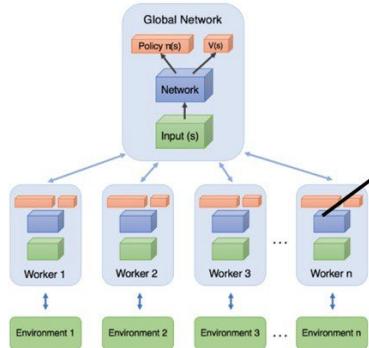
to ensure enough exploration for many agents

How to Minibatch Policy gradient ?



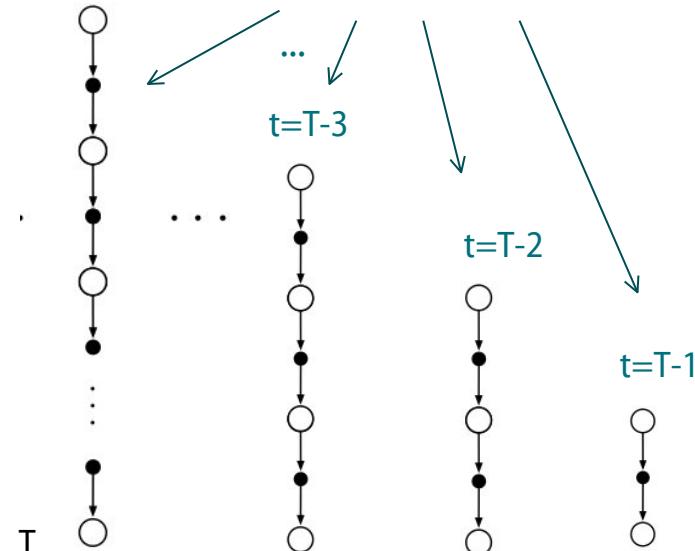
Batches on n epochs

Collect **set of trajectories** by running policy in the different environments.



At each timestep in each trajectory, compute the return of each TD(n-step) trajectory candidates :

$$R_{t(n)} = \sum_{l=0}^{n-1} \gamma^l r_{t+l} + \gamma^n V_\phi(s_{t+n})$$



PPO Applications



"PPO has become the default reinforcement learning algorithm at OpenAI because of its ease of use and good performance. July 20, 2017"



OpenAI Five DOTA2 winners

Questions Break



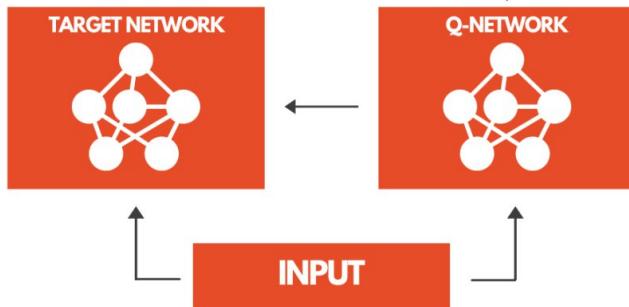
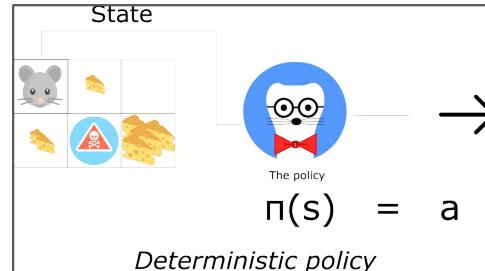
Deep Deterministic Policy Gradient



Off-policy

DDPG is an **off-policy** TD learning algorithm

DDPG can be thought of as being **Deep Q-learning**
Network for continuous control



$S_t, A_t, R_{t+1}, S_{t+1}, p_t$
$S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, p_{t+1}$
$S_{t+2}, A_{t+2}, R_{t+3}, S_{t+3}, p_{t+2}$
$S_{t+3}, A_{t+3}, R_{t+4}, S_{t+4}, p_{t+3}$
...

Experience replay Buffer

Sample



Batch of experiences

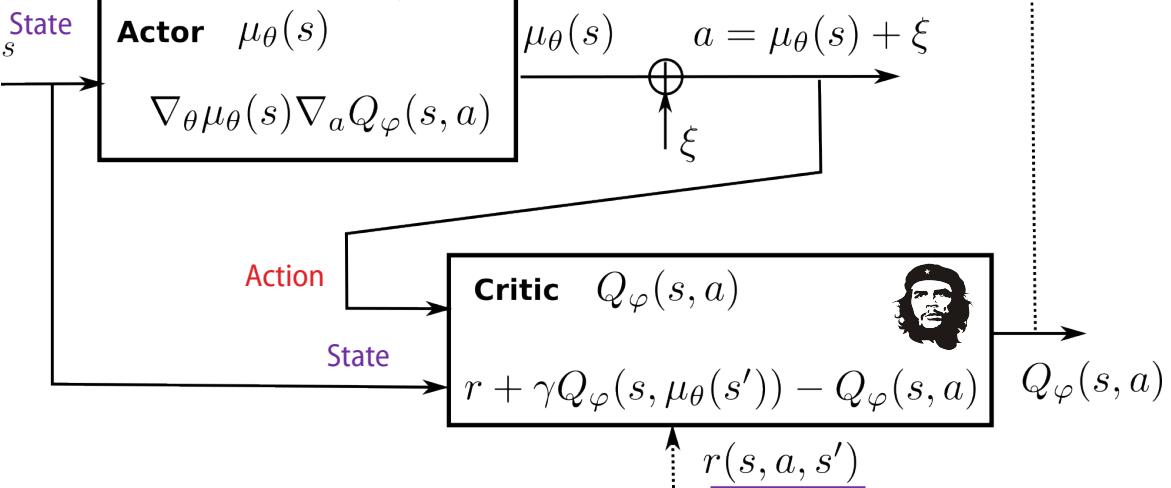
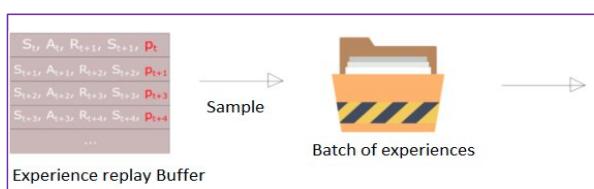


Deep Deterministic Policy Gradient



Off-policy

$$\nabla_a Q_\varphi(s, a)$$



Sample a random minibatch of N transitions (s_i, a_i, r_i, s_{i+1}) from R

Set $y_i = \underline{r_i} + \gamma Q'(s_{i+1}, \mu'(s_{i+1}| \theta^{\mu'})| \theta^{Q'})$

Update critic by minimizing the loss: $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$

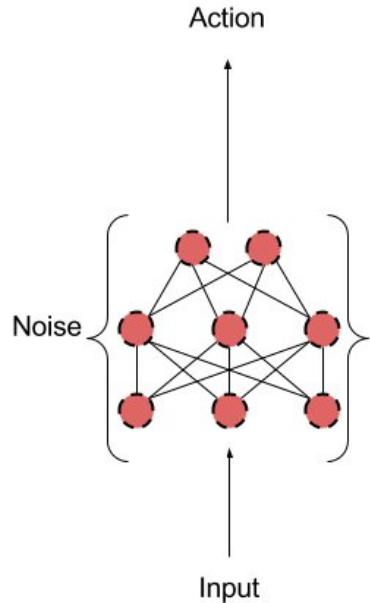
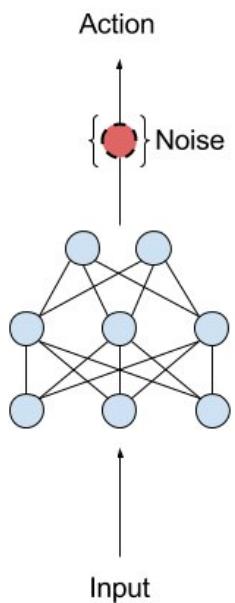
Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \underbrace{\nabla_a Q(s, a | \theta^Q)}_{\text{scalar}}|_{s=s_i, a=\mu(s_i)} \underbrace{\nabla_{\theta^\mu} \mu(s | \theta^\mu)}_{\text{vector}}|_{s_i}$$



Noisy exploration

Select action $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise



TD3 : Twin Delayed Deep Deterministic

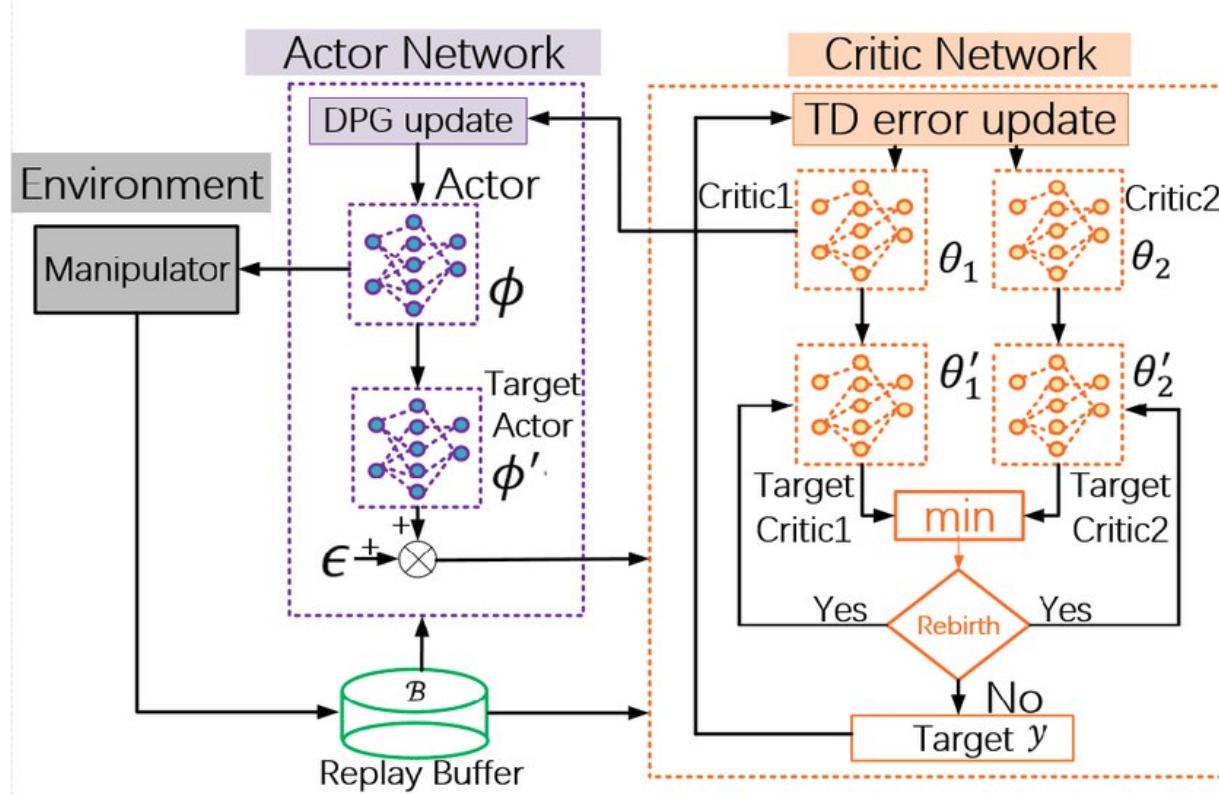


To fix Q value overestimation issue and variance issue : use Twin critic



Double Q-learning

Delayed update of target and policy networks

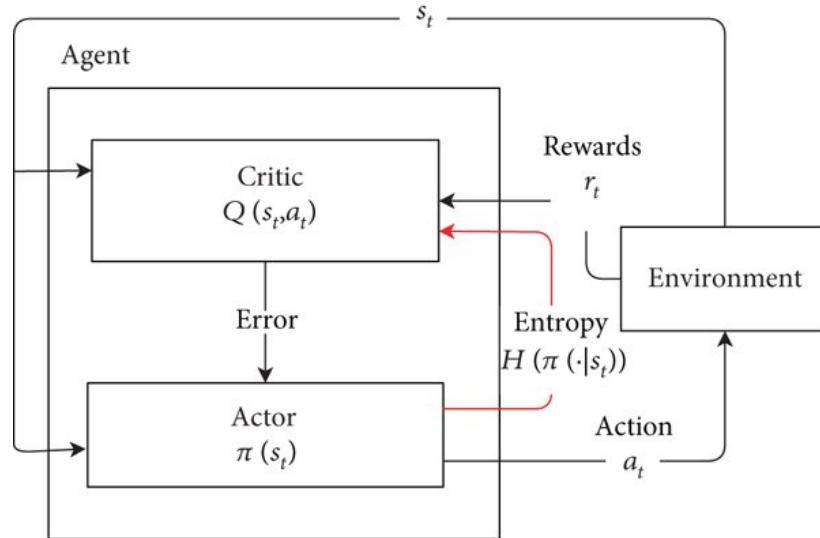


Soft Actor Critic



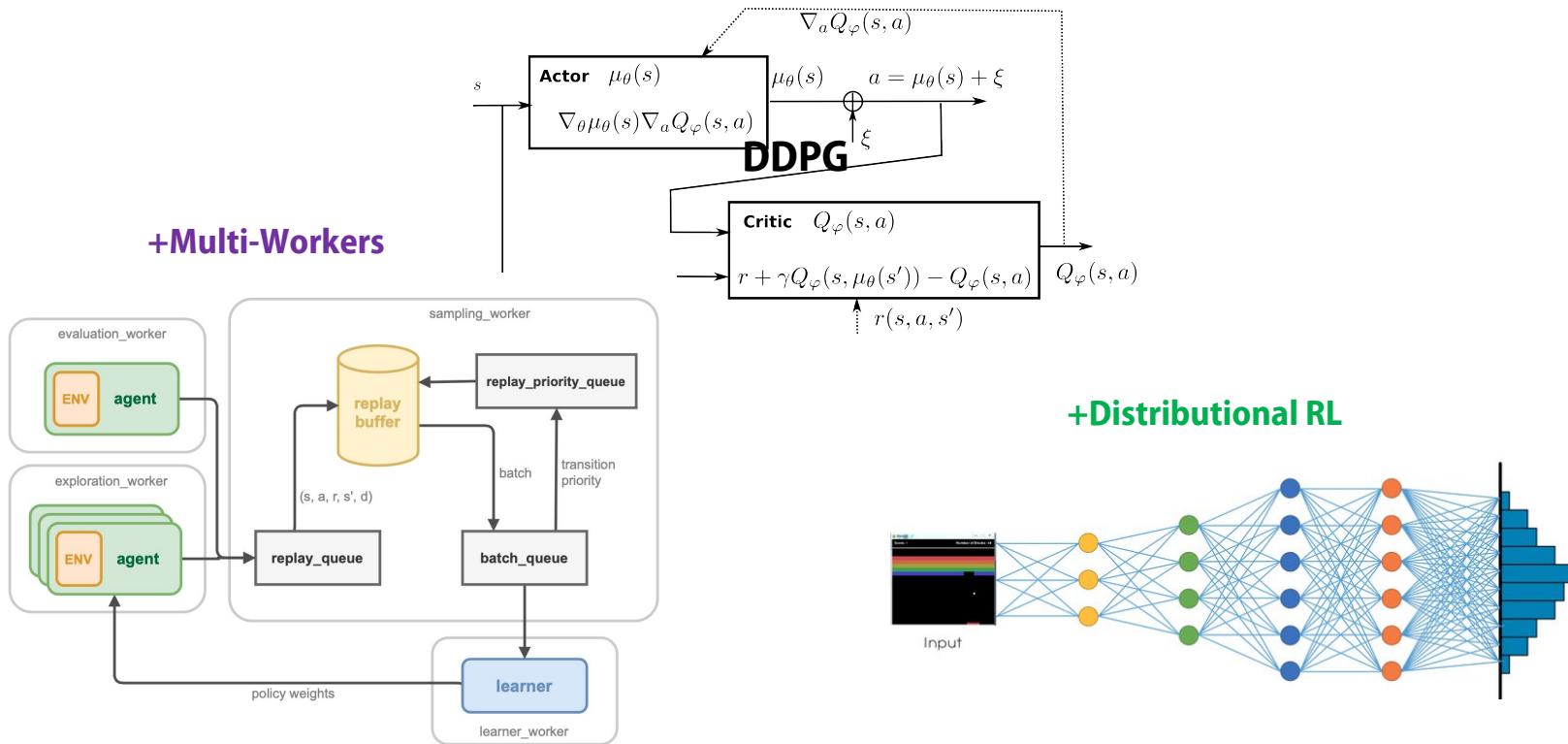
incorporates the **entropy measure of the policy** into the reward to encourage exploration

$$J(\theta) = \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi_\theta}} [r(s_t, a_t) + \underline{\alpha \mathcal{H}(\pi_\theta(\cdot | s_t))}]$$



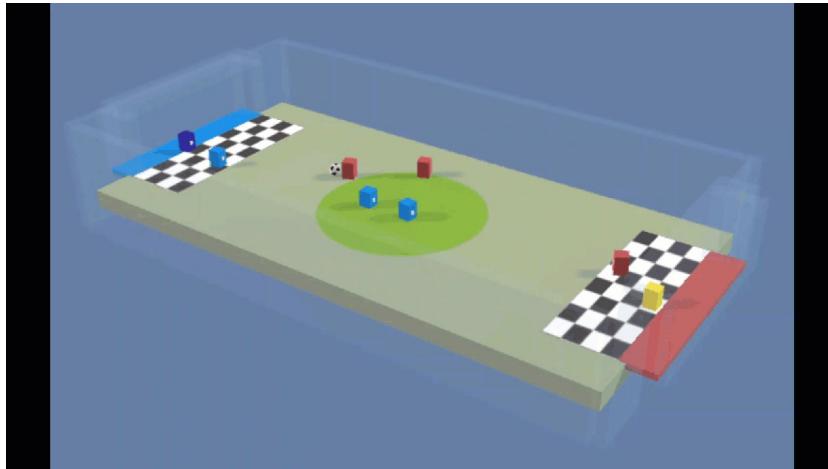


Distributed Distributional Deep Deterministic Policy Gradients

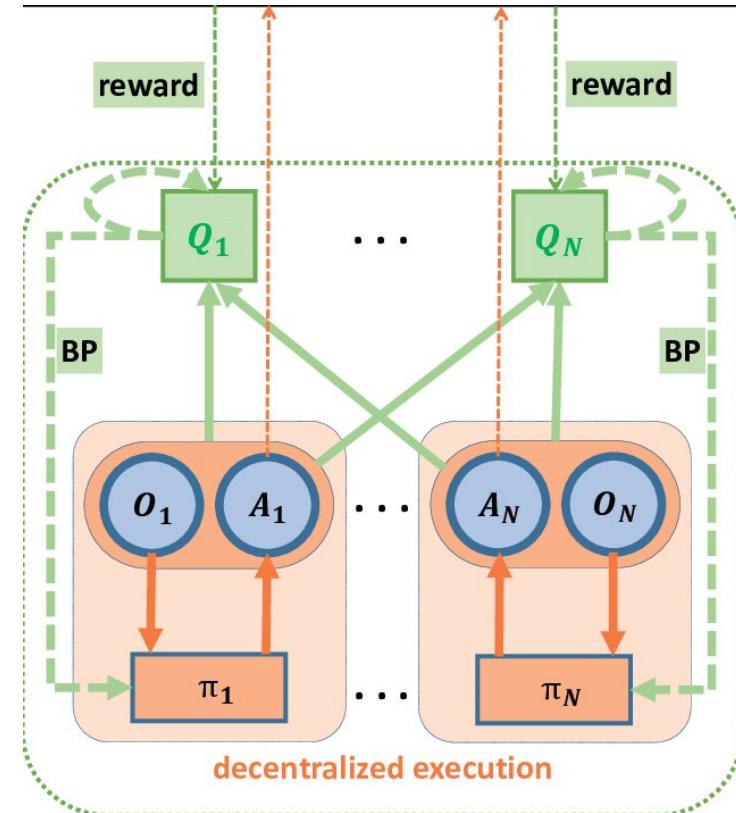




Collaborative and/or adversarial



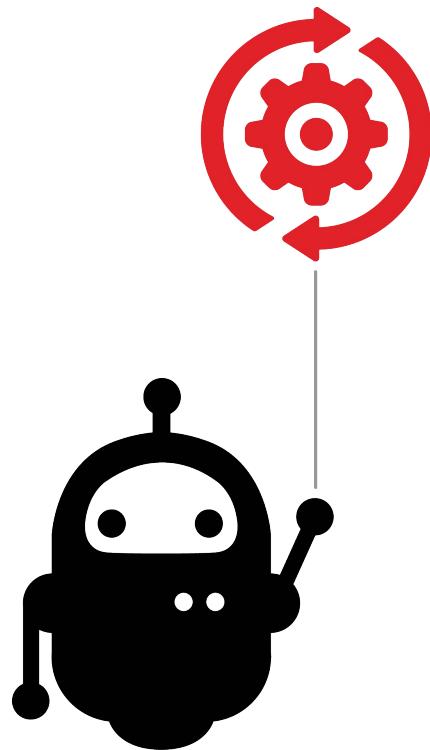
Critics take the actions of every actor



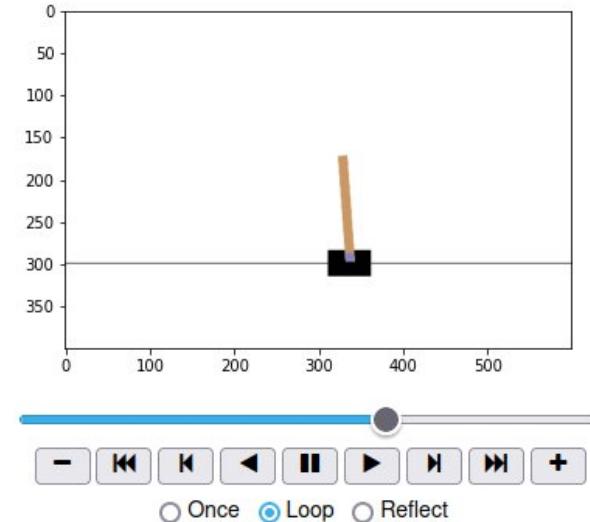
Questions Break



Example : Car Pole



[DRL1] - Solving CartPole with DDQN & Vanilla Policy Gradient



Success & Perspectives



[Dopamine on Tensorflow](#)

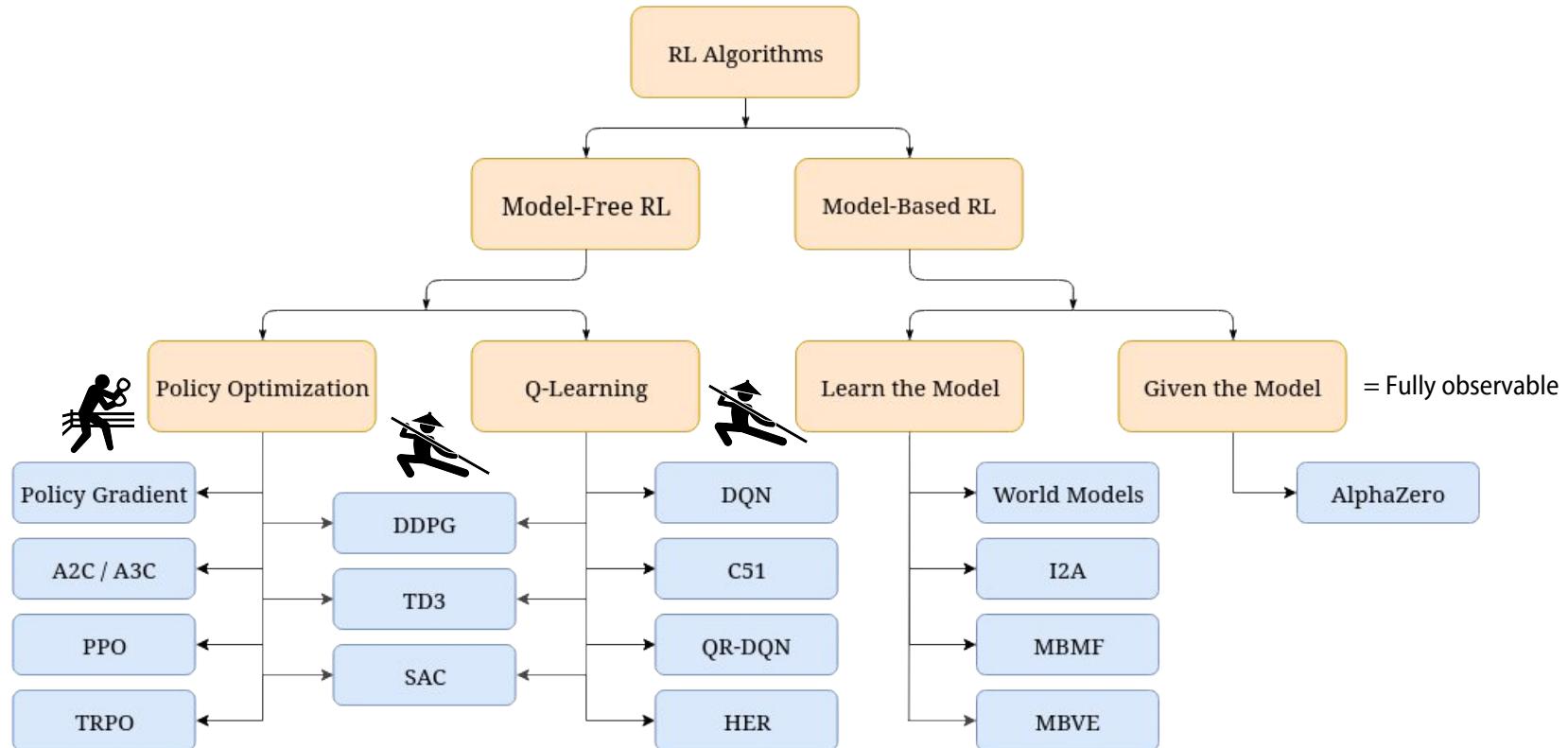


[Stable Baseline3 on Pytorch](#)

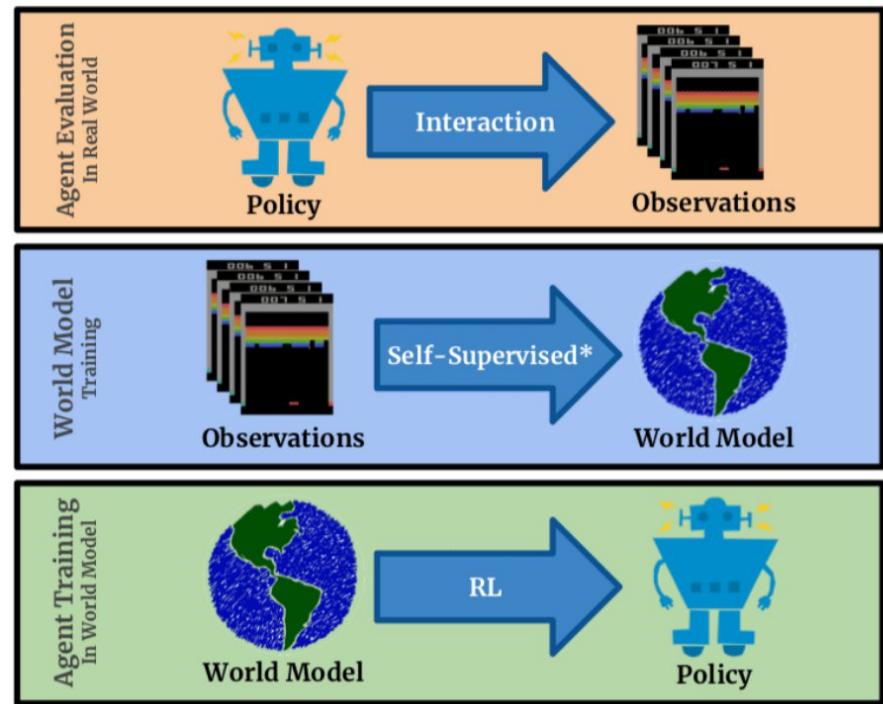
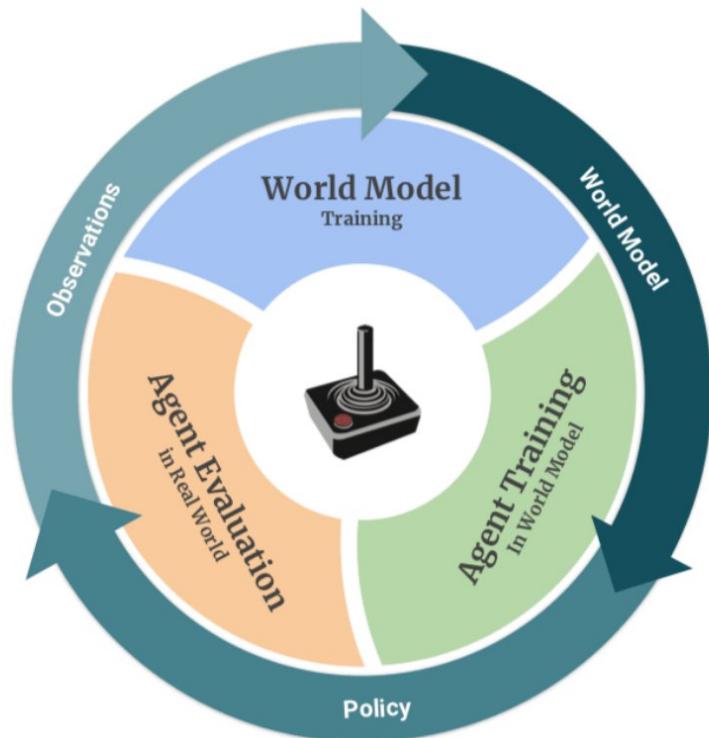
Which algorithm should I use ?

	Discrete control	Continuous control
Single Process	DQN and Distributional DQN. DQN is usually slower to train (regarding wall clock time) but is the most sample efficient (because of its replay buffer)	SAC, TD3 and TQC Please use the hyperparameters in the RL zoo for best results
Multiprocessed	PPO and A2C	PPO, TRPO and A2C Please use the hyperparameters in the RL zoo for best results

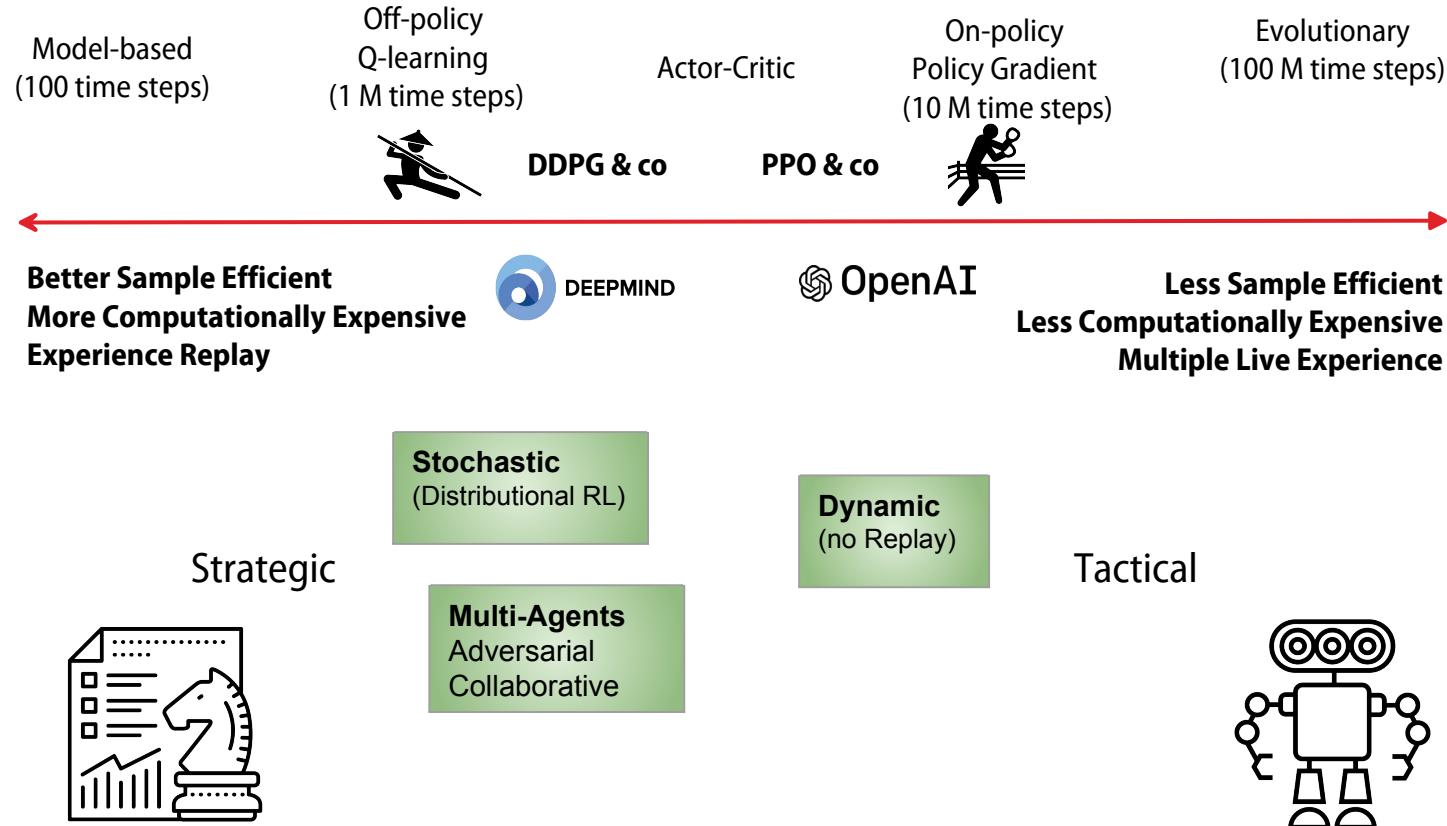
Kinds of DRL Algorithms



Model-Based Reinforcement Learning for Atari



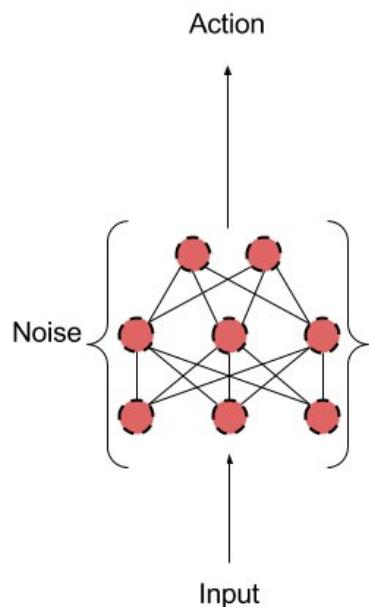
Which algorithm should I use ?



Augmented Random Search (2018)

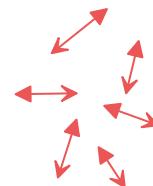
ARS is combination of :

- Evolution Strategies – **derivative free optimization** method, parallelized training
- Natural policy gradients for training linear policies



Exploration in the policy space:

Apply several $+/-$ noises δ



Collect $r[+]$, $r[-]$

Update the weights $\Theta \leftarrow \alpha(r[+] - r[-]). \delta$

Augmented with:

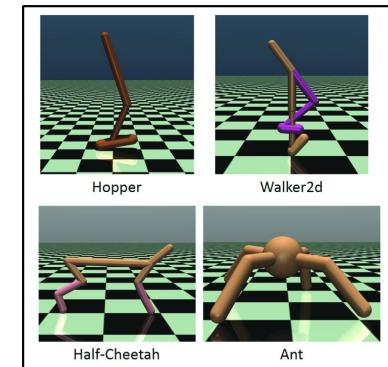
- Dividing by the Standard Deviation σ_r ,
- Normalizing the States,
- Using top performing directions

Sample Efficient

↔

Computationally Cheap

Outperforms largely PPO, DDPG on
MuJoCo locomotion environments !!



- **Reinforcement Learning - Book of R. S. Sutton and A. G. Barto**
- **Grokking Deep Reinforcement Learning - Book of M. Morales**
- [OpenAI Spinning Up](#)
- [Berkeley's Deep Reinforcement Learning course](#)
- [More resources](#)
- [Isaac Gym: High Performance GPU-Based Physics Simulation For Robot Learning](#) 

Policy Gradient Algorithms (2018)

..., PPG, ACER, ACTKR, SVPG, IMPALA

Meta Reinforcement Learning (2019)

generalize to new tasks or new environments with memory

Curriculum for Reinforcement Learning (2020)

step by step complexity guided learning

Exploration Strategies in Deep Reinforcement Learning (2020)

hard environments : rarely provides rewards or have distracting noise

And after 2020 ?

And what ?



After 2020 are we in one **Deep Reinforcement Learning Winter ?**

 **OpenAI**



DEEPMIND



TRANSFORMERS

ex : **GATO**, AlphaFold, GPT, PALM-E

New trend is **RL fine-tuning !!**

For impossible tasks with Reinforcement Learning from scratch !!

Video labeled pre-train then fine-tune with RL !

Video PreTraining (VPT): Learning to Act by Watching Unlabeled Online Videos

Bowen Baker^{*†}
bowen@openai.com

Ilge Akkaya^{*†}
ilge@openai.com

Peter Zhokhov^{*†}
peterz@openai.com

Joost Huizinga^{*†}
joost@openai.com

Jie Tang^{*†}
jietang@openai.com

Adrien Ecoffet^{*†}
adrien@openai.com

Brandon Houghton^{*†}
brandon@openai.com

Raul Sampedro^{*†}
raulsamg@gmail.com

Jeff Clune^{*†‡}
jclune@gmail.com

jun. 2022

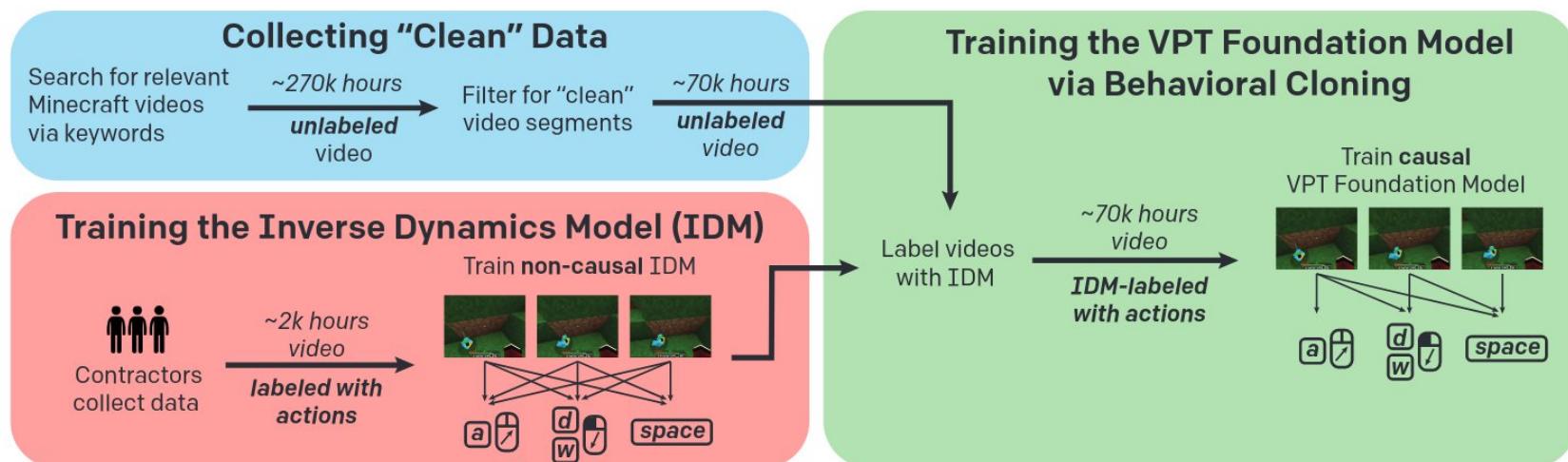


Figure 2: Video Pretraining (VPT) Method Overview.

Step 1

Collect demonstration data and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain reinforcement learning to a 6 year old.



We give treats and punishments to teach...



This data is used to fine-tune GPT-3.5 with supervised learning.

Step 2

Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

In reinforcement learning, the agent is...

C In machine learning...

D We give treats and punishments to teach...

B Explain rewards.



D > C > A > B

RM

D > C > A > B

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.

Write a story about otters.



Once upon a time...



r_k

The PPO model is initialized from the supervised policy.

The policy generates an output.

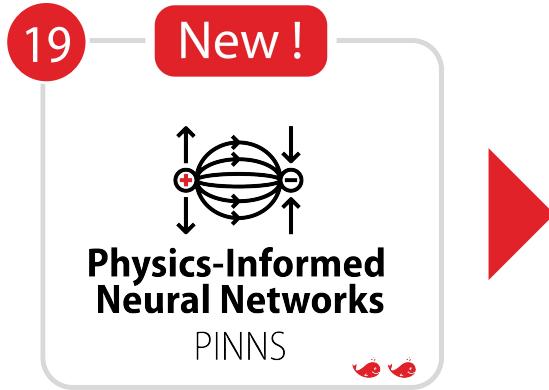
The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

Questions Break



Next, on Fidle :



Jeudi 4 mai,

Episode 19 :

Des neurones pour la physique : les PINNs (New !)

Que sont les Physics-Informed Neural Networks ?
Equilibrage de la Loss pour améliorer l'exactitude
Sampling & autres astuces

Durée : 2h



19

New !



Physics-Informed
Neural Networks
PINNs

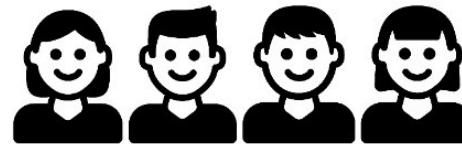


Jeudi 4 mai,

Episode 19 :

Des neurones pour la physique : les PINNs (New !)

Merci !



To be continued...



Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0)
<https://creativecommons.org/licenses/by-nc-nd/4.0/>