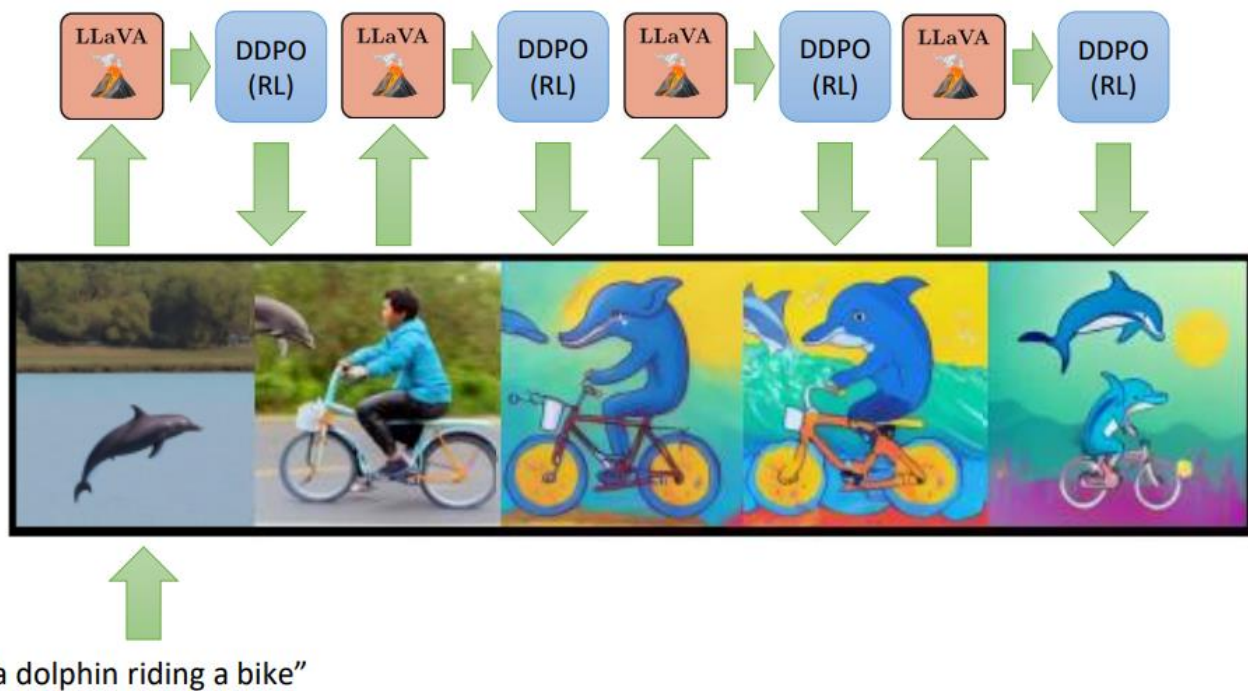


# Introduction To Reinforcement Learning



"a dolphin riding a bike"

Lancaster AI Reading Group – 22<sup>nd</sup> January 2025

*Based on CS285 Berkeley Course*

# Machine Learning

A computer program is said to learn from **experience**  $E$  with respect to some class of **tasks**  $T$  and **performance measure**  $P$  if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .

By Tom Mitchell (Prof. Carnegie Mellon University):



Experience



Experience

$\{(x_1, y_1), \dots, (x_N, y_N)\}$



# Experience

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

$\{x_1, \dots, x_N\}$



# Experience

$$\{(x_1, y_1), \dots, (x_N, y_N)\}$$

$$\{x_1, \dots, x_N\}$$

That is not always the case!

# The 3 ML Settings

Supervised Learning   Unsupervised Learning   Reinforcement Learning

# The 3 ML Settings

Supervised Learning

Unsupervised Learning

Reinforcement Learning

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

Learn  $p(y|x)$

$\{x_1, \dots, x_N\}$

Learn  $p(x)$

?



# The 3 ML Settings

Supervised Learning

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

Learn  $p(y|x)$

Unsupervised Learning

$\{x_1, \dots, x_N\}$

Learn  $p(x)$

Reinforcement Learning

- Environment

# The 3 ML Settings

Supervised Learning

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

Learn  $p(y|x)$

Unsupervised Learning

$\{x_1, \dots, x_N\}$

Learn  $p(x)$

Reinforcement Learning

- Environment
- Agent

# The 3 ML Settings

Supervised Learning

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

Learn  $p(y|x)$

Unsupervised Learning

$\{x_1, \dots, x_N\}$

Learn  $p(x)$

Reinforcement Learning

- Environment
- Agent
- Actions

# The 3 ML Settings

Supervised Learning

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

Learn  $p(y|x)$

Unsupervised Learning

$\{x_1, \dots, x_N\}$

Learn  $p(x)$

Reinforcement Learning

- Environment
- Agent
- Actions
- Rewards

# The 3 ML Settings

Supervised Learning

$\{(x_1, y_1), \dots, (x_N, y_N)\}$

Learn  $p(y|x)$

Unsupervised Learning

$\{x_1, \dots, x_N\}$

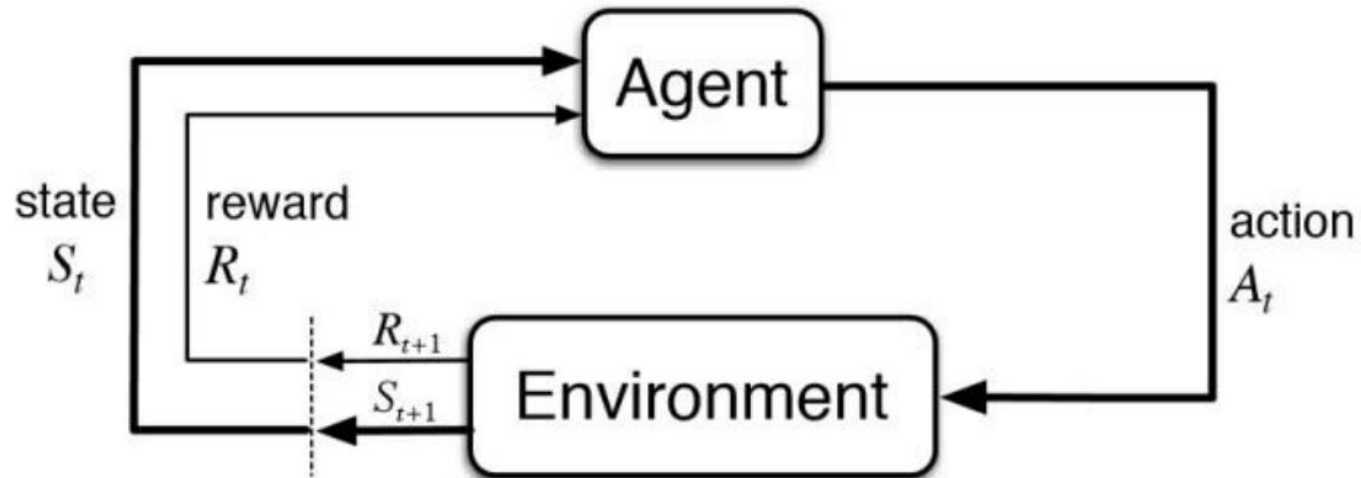
Learn  $p(x)$

Reinforcement Learning

- Environment
- Agent
- Actions
- Rewards

How to learn a policy that maximise the cumulative reward?

# The RL Framework

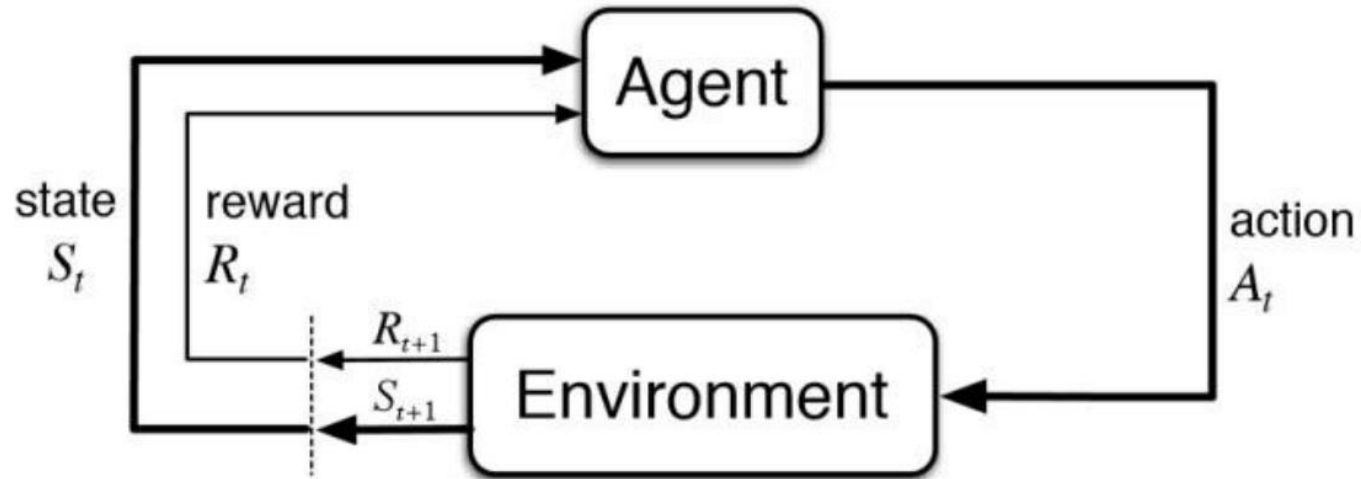


## Reinforcement Learning

- Environment
- Agent
- Actions
- Rewards

How to learn a policy that maximise the cumulative reward?

# The RL Framework



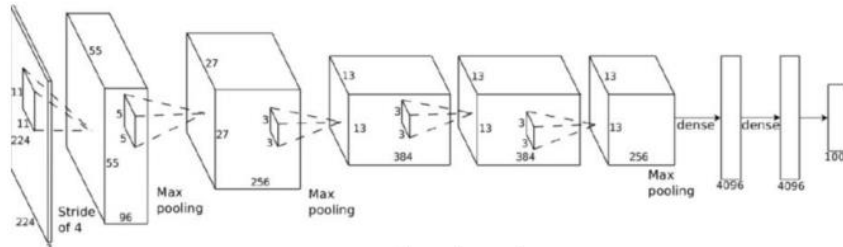
## Reinforcement Learning

- Environment
- Agent
- Actions
- Rewards

How to learn a policy that maximise the cumulative reward?



$\mathbf{o}_t$



$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$



$\mathbf{a}_t$

# RL Examples





# RL Examples



Actions: muscle contractions  
Observations: sight, smell  
Rewards: food



Actions: motor current or torque  
Observations: camera images  
Rewards: task success measure (e.g., running speed)



Actions: what to purchase  
Observations: inventory levels  
Rewards: profit

What about the environments?

What about the environments?



# What about the environments?

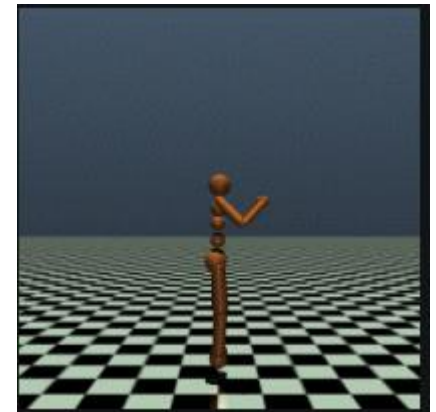
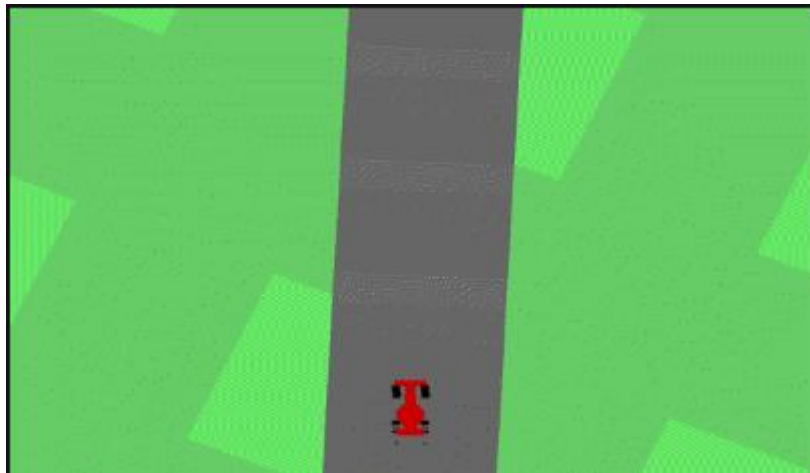
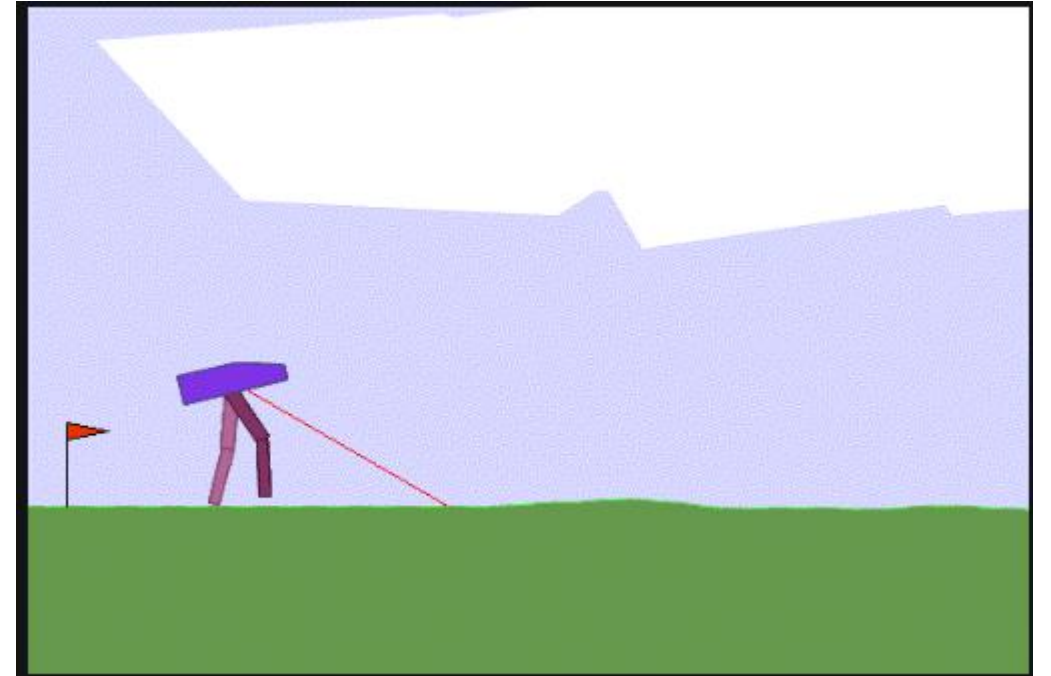
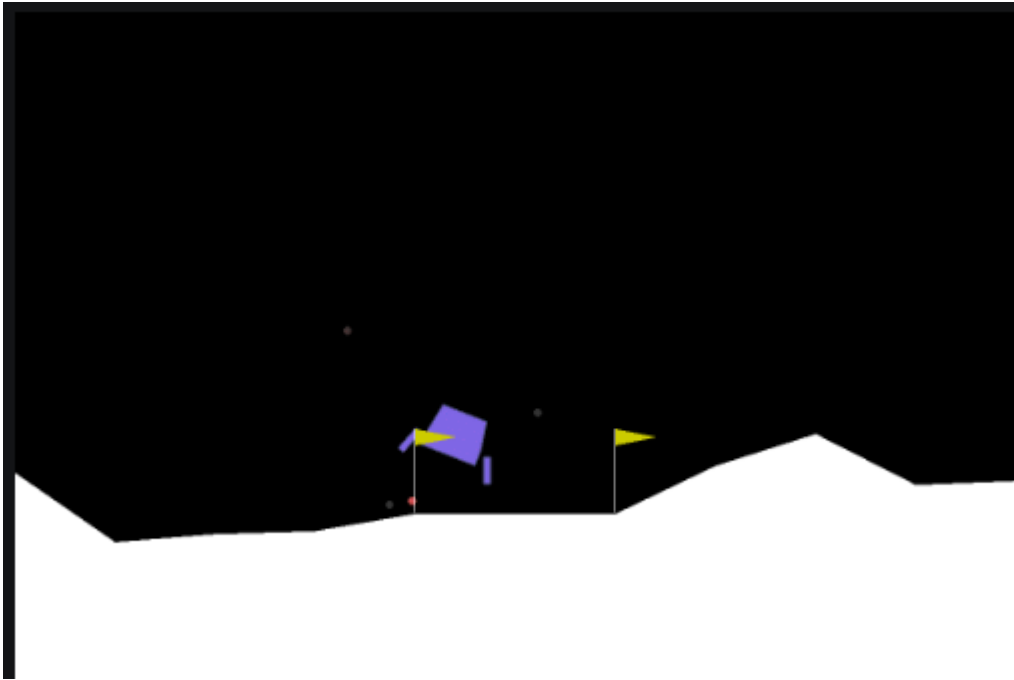


```
import gym

# Create an environment
env = gym.make("Ant-v4", new_step_api=True)
# Reset the environment to an initial state
obs = env.reset()
# Take a random action
action = env.action_space.sample()
# Update the environment according to the action
obs, reward, done, info = env.step(action)

env.close()
```

# Example Environments



How to learn a “good” policy?

How to learn a "good" policy?

A simple RL algorithm: **REINFORCE**

# How to learn a “good” policy?

A simple RL algorithm: **REINFORCE**

- Sample a trajectory from the current policy.



# How to learn a “good” policy?

A simple RL algorithm: **REINFORCE**

$$\nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(s_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | s_{i,t}) \right) \left( \sum_{i=1}^T r(s_{i,t}, \mathbf{a}_{i,t}) \right)$$

- Sample a trajectory from the current policy.
- Estimate the gradient of the objective (use backpropagation in the case of Neural Networks).

# How to learn a “good” policy?

A simple RL algorithm: **REINFORCE**

$$\nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(s_t, a_t) \right] \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{i=1}^T r(s_{i,t}, a_{i,t}) \right)$$

- Sample a trajectory from the current policy.
- Estimate the gradient of the objective (use backpropagation in the case of Neural Networks).
- Update the policy using a gradient based optimization algorithm.

# The necessity of RL

- Some tasks are hard to get supervised pairs from (e.g. movements of a robot).

# The necessity of RL

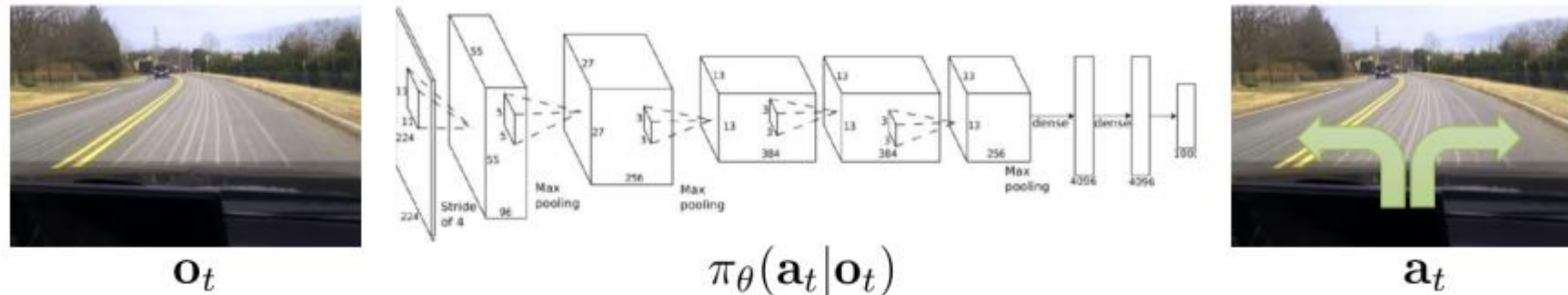
- Some tasks are hard to get supervised pairs from (e.g. movements of a robot).
- Supervised learning is not ideal for sequential decision making.

# The necessity of RL

- Some tasks are hard to get supervised pairs from (e.g. movements of a robot).
- Supervised learning is not ideal for sequential decision making.
- RL allows for novel solutions.

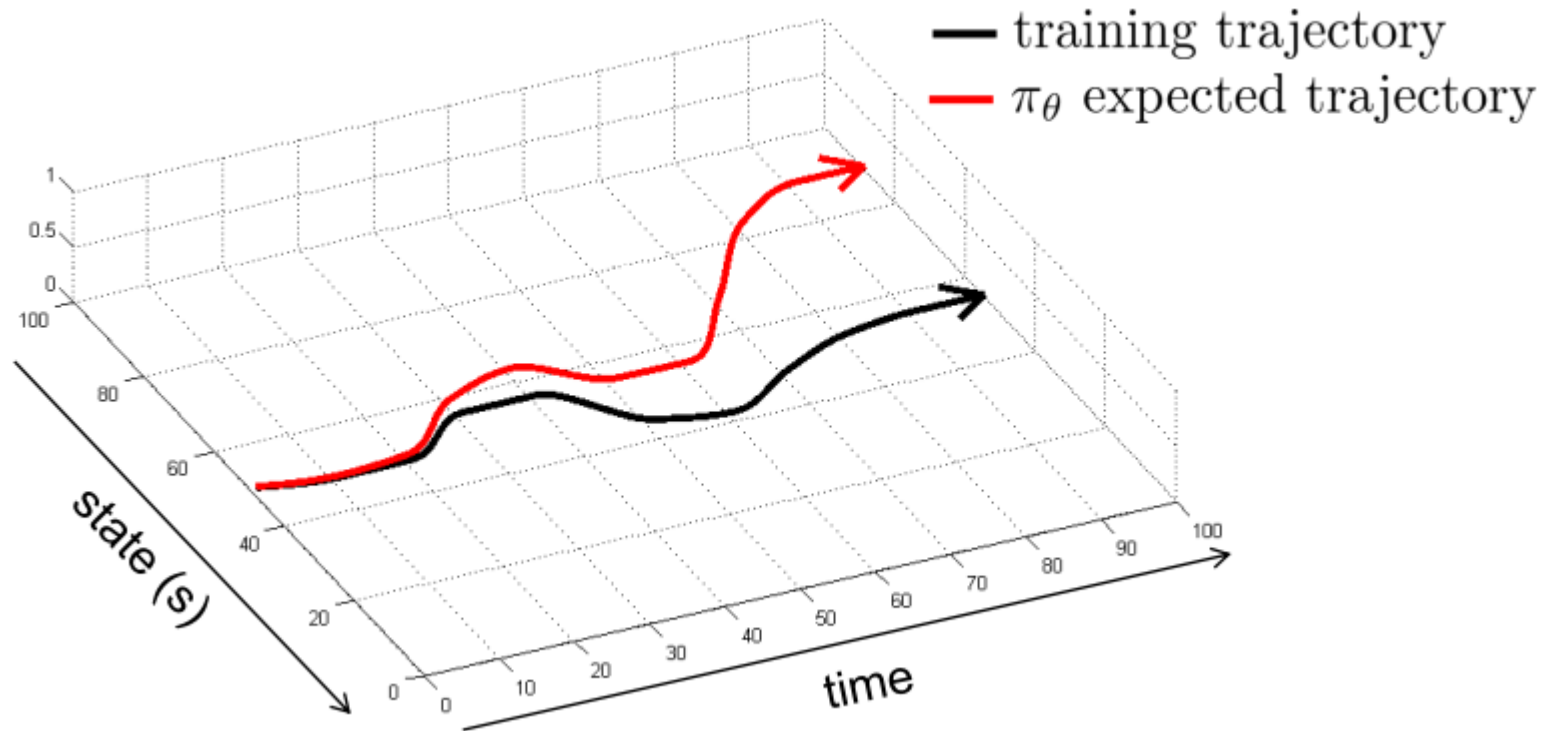
# Imitation Learning

- Supervised learning is not ideal for sequential decision making.



# Imitation Learning

- Supervised learning is not ideal for sequential decision making.



Data is not iid!

# Recent advances in DL vs RL

- RL allows for novel solutions.



vibrant portrait painting of Salvador Dali with a robotic half face

a shiba inu wearing a beret and black turtleneck

a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation

panda mad scientist mixing sparkling chemicals, artstation

a corgi's head depicted as an explosion of a nebula





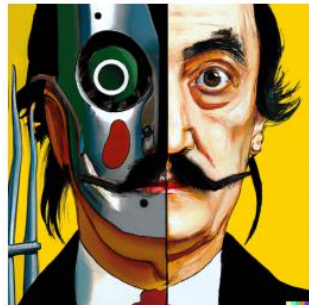
# Recent advances in DL vs RL

- RL allows for novel solutions.

$$p(x)$$

$$p(x|c)$$

$$p(y_n|y_{1:n-1})$$



vibrant portrait painting of Salvador Dali with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



a corgi's head depicted as an explosion of a nebula



# Recent advances in DL vs RL

- RL allows for novel solutions.

Image/Video  $p(x)$

Text  $p(x|c)$

Next token  $p(y_n | y_{1:n-1})$  Prev tokens



# Recent advances in DL vs RL

- RL allows for novel solutions.

Image/Video  
 $p(x)$

Text  
 $p(x|c)$

Next token      Prev tokens  
 $p(v_n | v_{1:n-1})$

Impressive because it looks like something a person might draw!



vibrant portrait painting of Salvador Dali with a robotic half face



a shiba inu wearing a beret and black turtle neck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



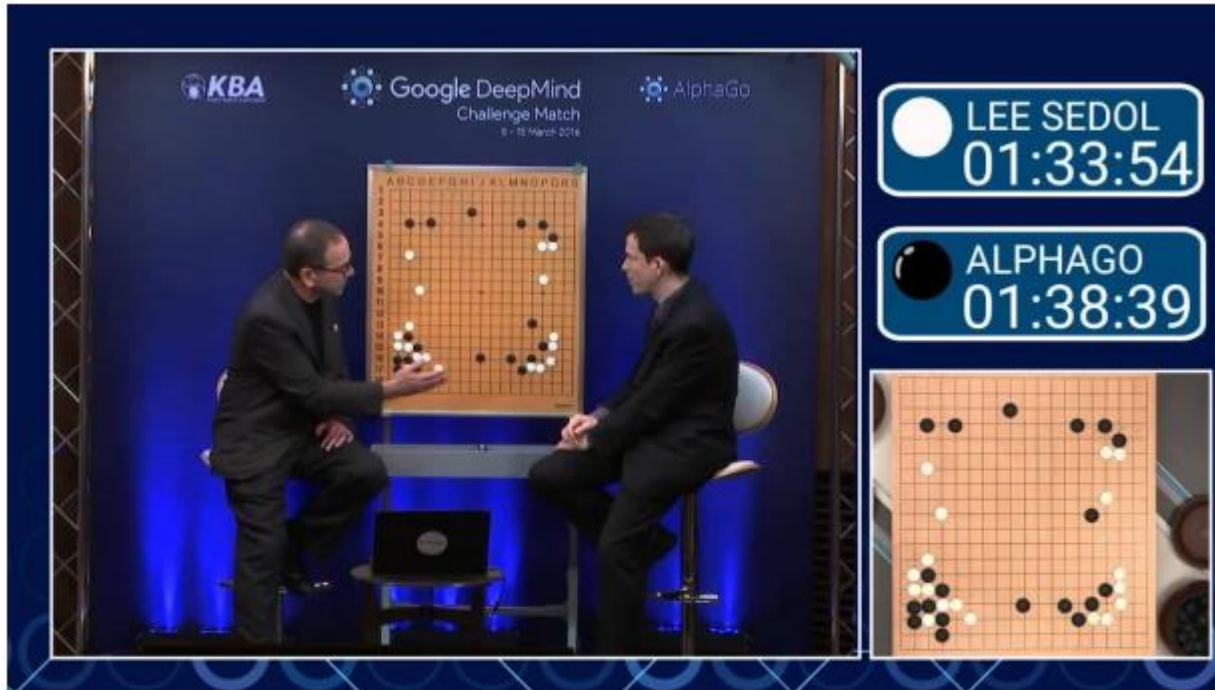
a corgi's head depicted as an explosion of a nebula



ChatGPT

# Recent advances in DL vs RL

- RL allows for novel solutions.



“Move 37” in Lee Sedol AlphaGo match: reinforcement learning “discovers” a move that surprises everyone

# Recent advances in DL vs RL

- RL allows for novel solutions.



Impressive because no person had thought of it!

“Move 37” in Lee Sedol AlphaGo match: reinforcement learning “discovers” a move that surprises everyone

# Q and Value Functions

- State at time  $t$ :  $s_t$
- Action at time  $t$ :  $a_t$
- Reward function:  $r(s_t, a_t)$

# Q and Value Functions

The Q-Function is the **total reward from taking action  $a_t$  in state  $s_t$**  :

- State at time  $t$ :  $s_t$
- Action at time  $t$ :  $a_t$
- Reward function:  $r(s_t, a_t)$

# Q and Value Functions

The Q-Function is the **total reward from taking action  $a_t$  in state  $s_t$**  :

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

- State at time  $t$ :  $s_t$
- Action at time  $t$ :  $a_t$
- Reward function:  $r(s_t, a_t)$



# Q and Value Functions

The Q-Function is the **total reward** from taking action  $a_t$  in state  $s_t$  :

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

The Value Function is the **total reward** from state  $s_t$  :

- State at time  $t$ :  $s_t$
- Action at time  $t$ :  $a_t$
- Reward function:  $r(s_t, a_t)$

# Q and Value Functions

The Q-Function is the **total reward** from taking action  $a_t$  in state  $s_t$  :

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

The Value Function is the **total reward** from state  $s_t$  :

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

- State at time  $t$ :  $s_t$
- Action at time  $t$ :  $a_t$
- Reward function:  $r(s_t, a_t)$

# Q and Value Functions

The Q-Function is the **total reward from taking action  $a_t$**  in state  $s_t$  :

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

The Value Function is the **total reward from state  $s_t$**  :

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

$$V^\pi(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$$

# Q and Value Functions

**Idea 1:**

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

**Idea 2:**

$$V^\pi(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi(\alpha_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$$

# Q and Value Functions

## Idea 1:

If we have a policy  $\pi$  and we know the Q-Function we can improve the policy by setting:

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

$$V^\pi(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$$

# Q and Value Functions

## Idea 1:

If we have a policy  $\pi$  and we know the Q-Function we can improve the policy by setting:

$$\pi'(a|s) = 1 \text{ if } a = \arg \max_a Q^\pi(s, a)$$

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

# Q and Value Functions

## Idea 1:

If we have a policy  $\pi$  and we know the Q-Function we can improve the policy by setting:

$$\pi'(a|s) = 1 \text{ if } a = \arg \max_a Q^\pi(s, a)$$

## Idea 2:

If  $Q^\pi(s, a) > V^\pi(s)$

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

# Q and Value Functions

## Idea 1:

If we have a policy  $\pi$  and we know the Q-Function we can improve the policy by setting:

$$\pi'(a|s) = 1 \text{ if } a = \arg \max_a Q^\pi(s, a)$$

## Idea 2:

If  $Q^\pi(s, a) > V^\pi(s)$  then  $a$  is better than average and we can increase  $\pi(a|s)$ .

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$



# Type of RL Algorithms

- **Policy Gradients:**
- **Value Based:**
- **Actor Critic:**
- **Model-Based:**

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

$$V^\pi(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$$

# Type of RL Algorithms

- **Policy Gradients:** Maximise the expected reward by direct differentiation (REINFORCE).
- **Value Based:**
- **Actor Critic:**
- **Model-Based:**

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

# Type of RL Algorithms

- **Policy Gradients:** Maximise the expected reward by direct differentiation (REINFORCE).
- **Value Based:** Estimate the Value or Q Function of the optimal policy.
- **Actor Critic:**
- **Model-Based:**

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

# Type of RL Algorithms

- **Policy Gradients:** Maximise the expected reward by direct differentiation (REINFORCE).
- **Value Based:** Estimate the Value or Q Function of the optimal policy.
- **Actor Critic:** Estimate the Value or Q Function of the current policy. Use it to improve policy.
- **Model-Based:**

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

# Type of RL Algorithms

- **Policy Gradients:** Maximise the expected reward by direct differentiation (REINFORCE).
- **Value Based:** Estimate the Value or Q Function of the optimal policy.
- **Actor Critic:** Estimate the Value or Q Function of the current policy. Use it to improve policy.
- **Model-Based:** Estimate the transition model and use it for planning/improving policy/other.

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

# Online vs Offline RL

- **Online RL:** The agent interacts with the environment during training.

# Online vs Offline RL

- **Online RL:** The agent interacts with the environment during training.
- **Offline RL:** The agent uses a fixed dataset of previously collected experiences without further interaction with the environment during the training phase.

# On-Policy vs Off-Policy RL

- **On-Policy RL:** The agent improves the policy currently being used to make decisions.



# On-Policy vs Off-Policy RL

- **On-Policy RL:** The agent improves the policy currently being used to make decisions.
- **Off-Policy RL:** The agent improves a different policy than the one it is using.

# Exploration vs. Exploitation

- **Exploitation:** Go to your favorite restaurant.
- **Exploration:** Try a new restaurant.

# Exploration vs. Exploitation

- **Exploitation:** Go to your favorite restaurant.
- **Exploration:** Try a new restaurant.
  
- **Exploitation:** Doing what you know will yield highest reward.
- **Exploration:** Doing things you haven't done before, in the hopes of getting even higher reward.

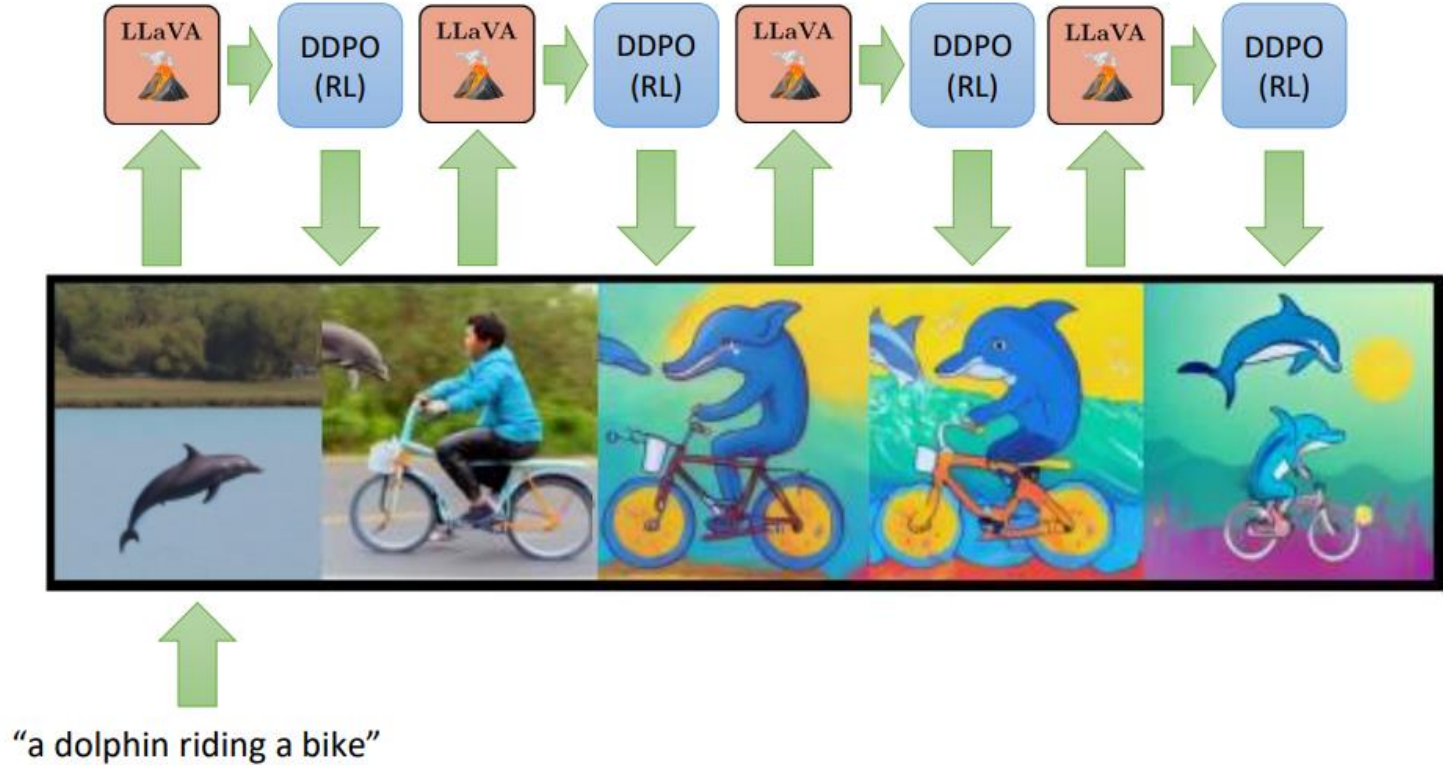
# Exploration vs. Exploitation

- **Exploitation:** Go to your favorite restaurant.
- **Exploration:** Try a new restaurant.

**How much should I explore?**

- **Exploitation:** Doing what you know will yield highest reward.
- **Exploration:** Doing things you haven't done before, in the hopes of getting even higher reward.

# Applications of RL



# Applications of RL

---

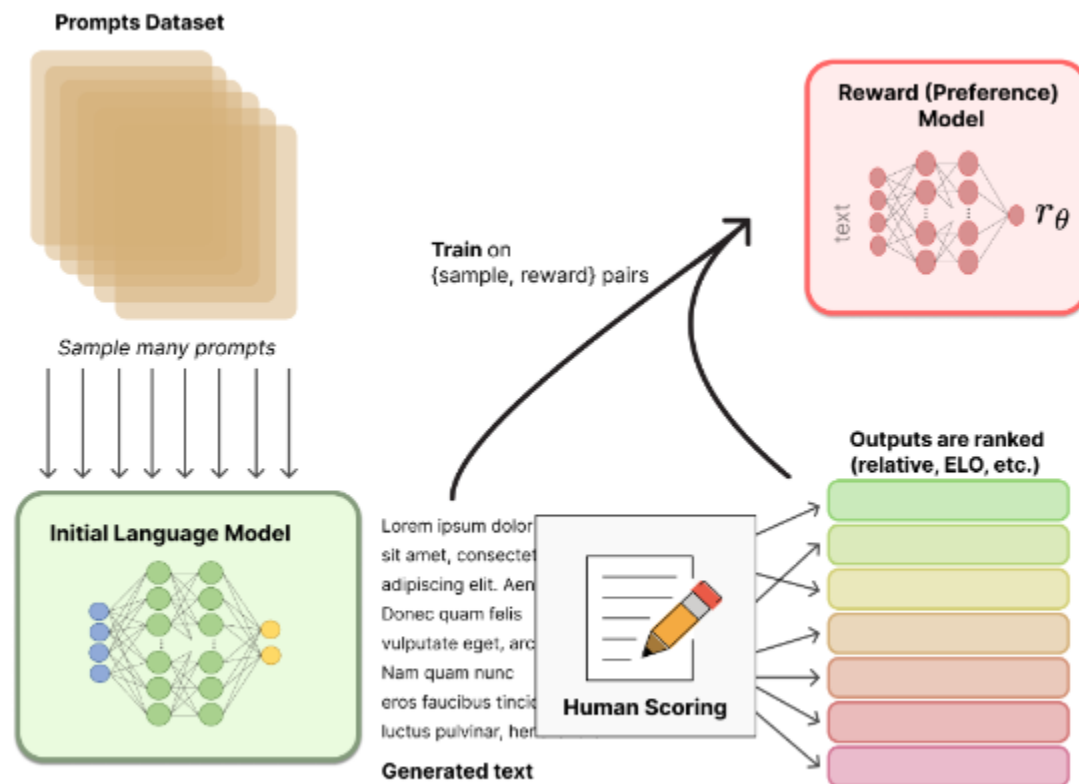
**Imitation-based policies can be sensitive**



**source policy (motion imitation)**

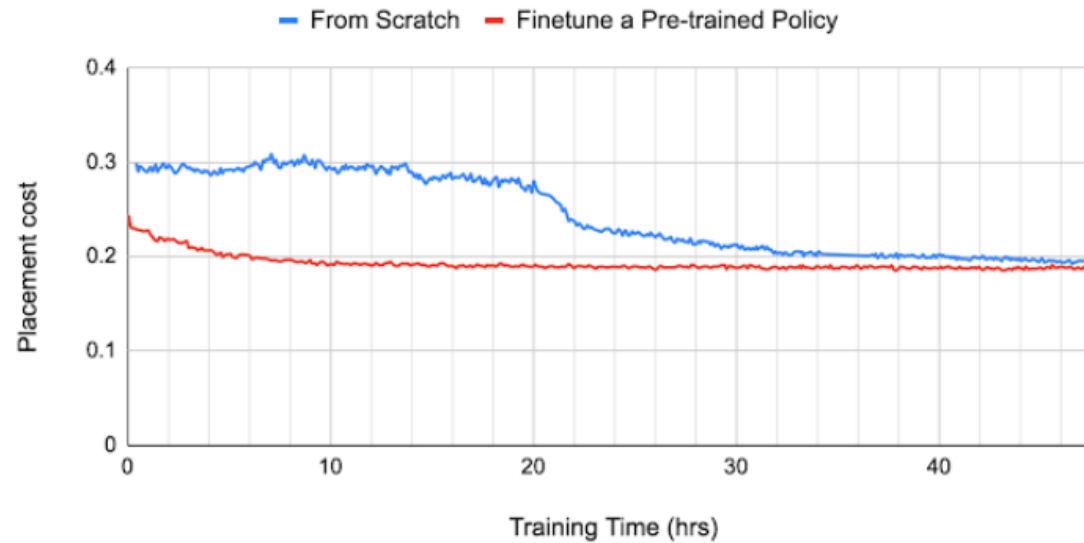
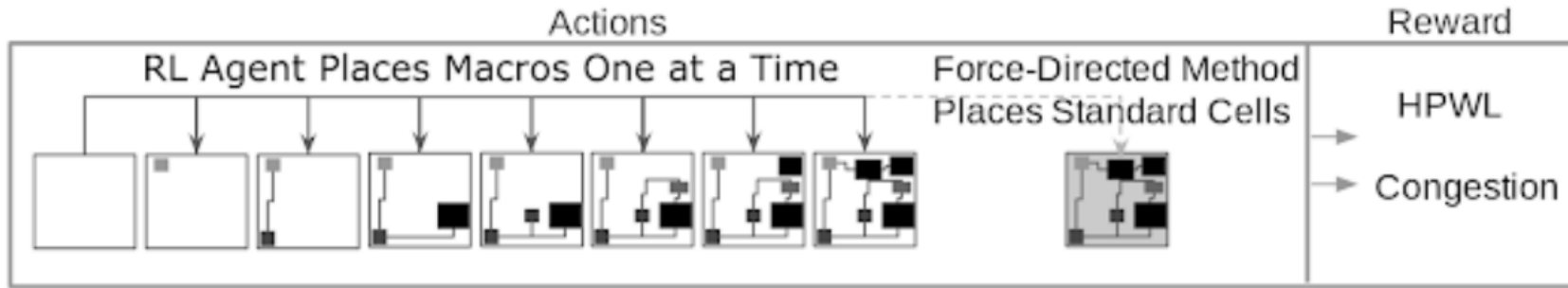
# Applications of RL

RL for fine-tuning LLMs



# Applications of RL

## RL for Chip Design





# Applications of RL

RL for Controlling Traffic

