# Deep Reinforcement Learning
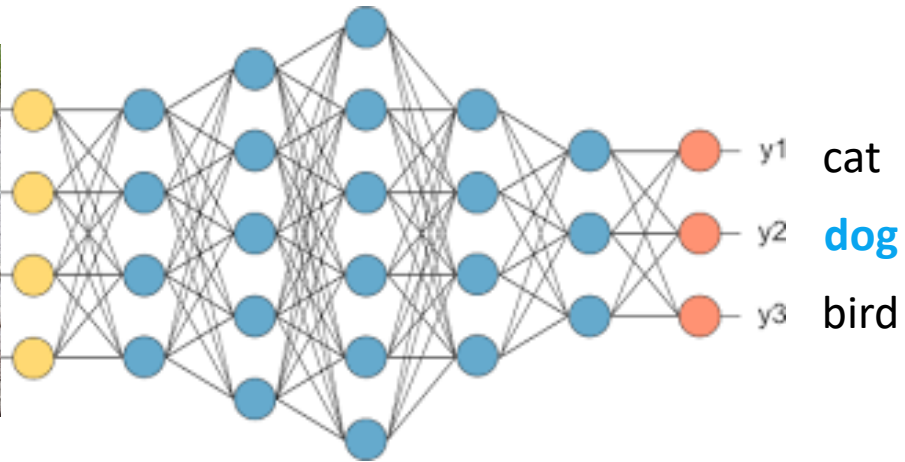


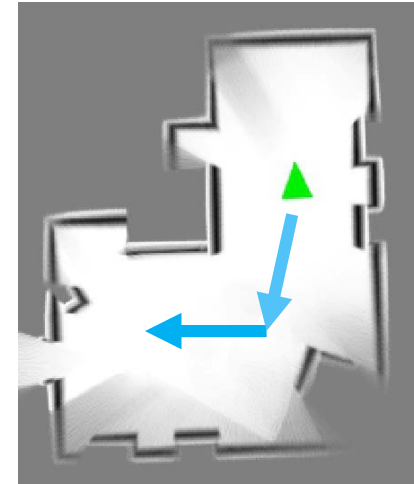school of ai

Beril Sirmacek

7th March 2019, Enschede

# Outline

- Reinforcement Learning (RL)
- Deep reinforcement learning (DRL)
- Practical applications of RL

# What is reinforcement learning?

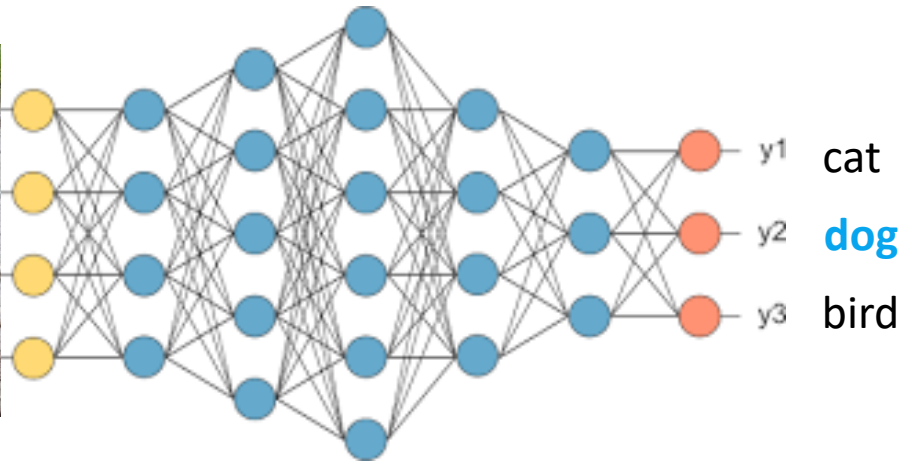- Similarity with other deep learning nets and RL



Supervised learning



Reinforcement learning
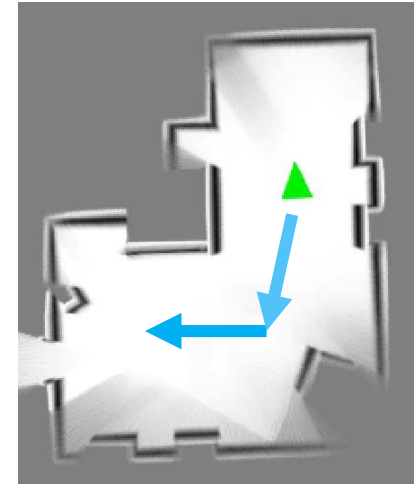
decision making!

# What is reinforcement learning?

- Difference from other deep learning nets and RL
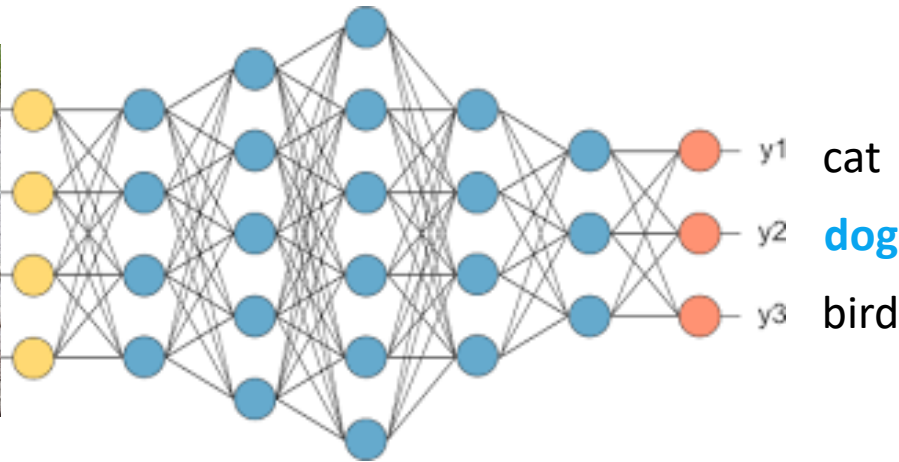


Supervised learning



Reinforcement learning

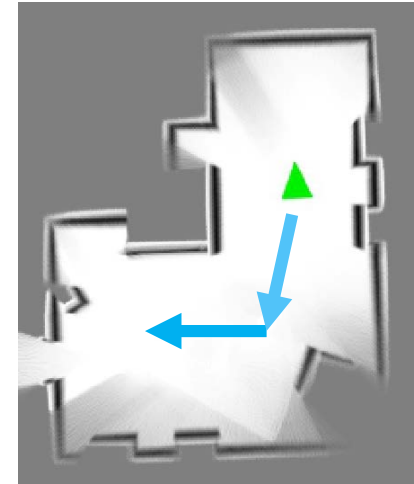A chain of decisions (actions) & dependency on the previous decisions

# What is reinforcement learning?

- Difference from other deep learning nets and RL



Where is UNSUPERVISED LEARNING?
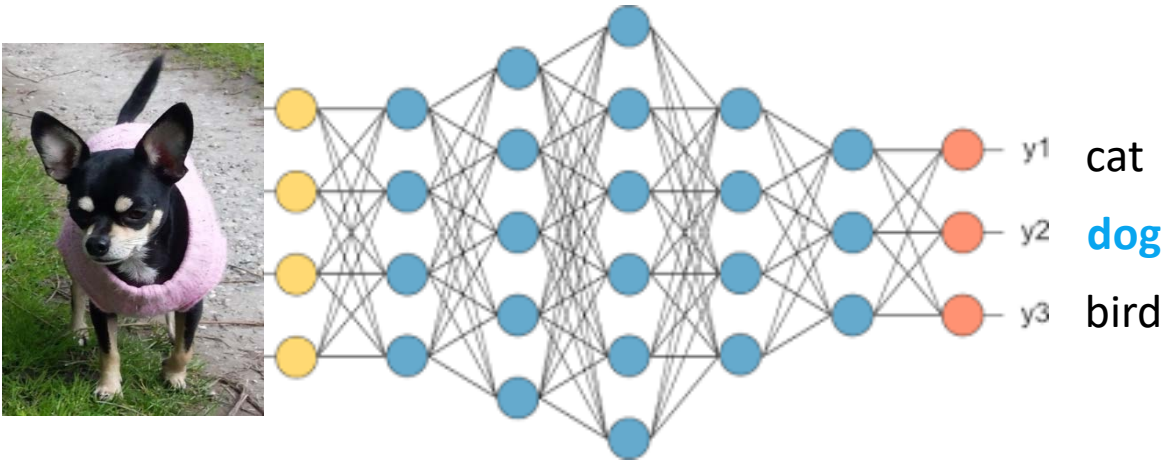
Supervised learning

Reinforcement learning

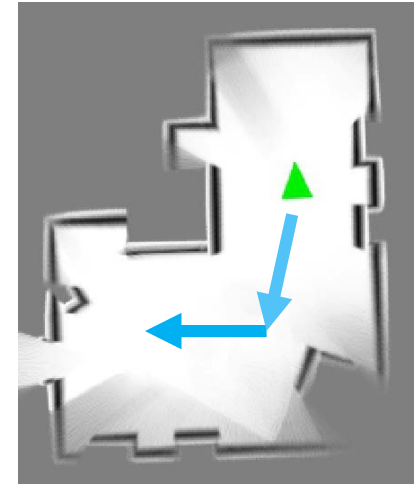A chain of decisions (actions) & dependency on the previous decisions

# What is reinforcement learning?

- Here are the examples / Loss function
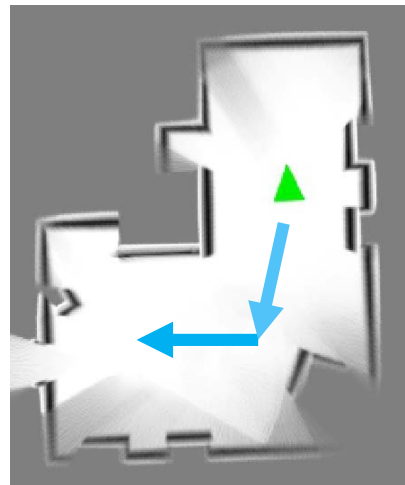- Learn it (learning by example)

- Here is the environment
- Learn it (learning by experience)



Supervised learning



Reinforcement learning

# Problem description

- What series of actions will give me the highest reward?
  (or helps me to reach to my goal faster/safer/better?)

Reinforcement learning

# Problem formulation



Agent



Environment



Actions

# Problem formulation



Agent



Environment



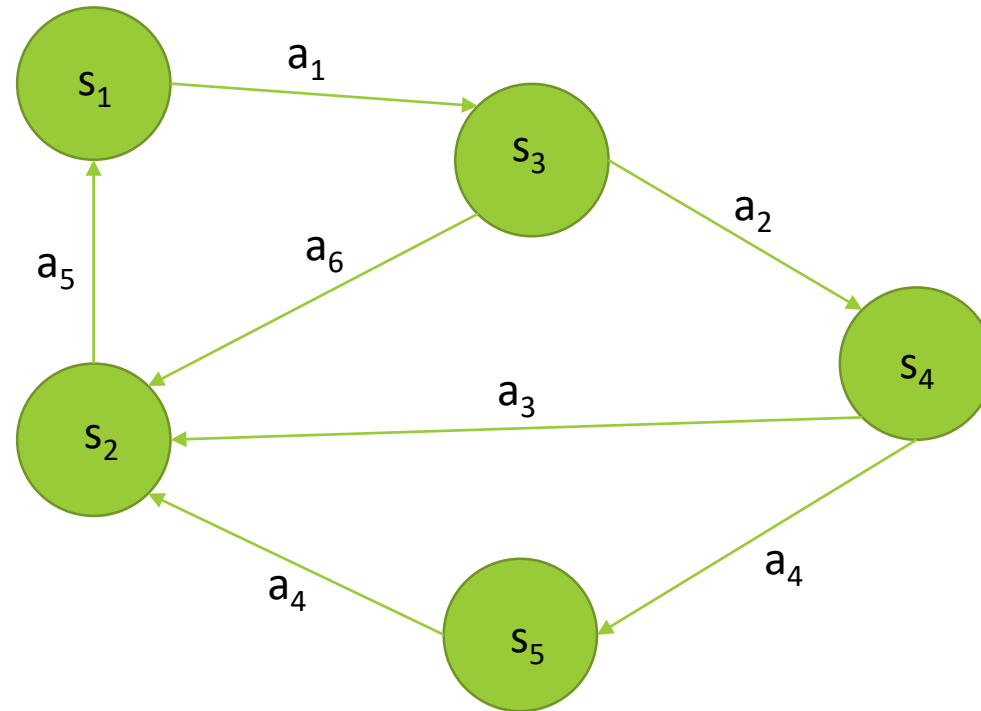Actions

How to decide which actions to take?



Policy

# Mathematical formulation



Agent

States: $s_1, s_2, s_3, s_4, ...$
Actions: $a_1, a_2, a_3, ...$
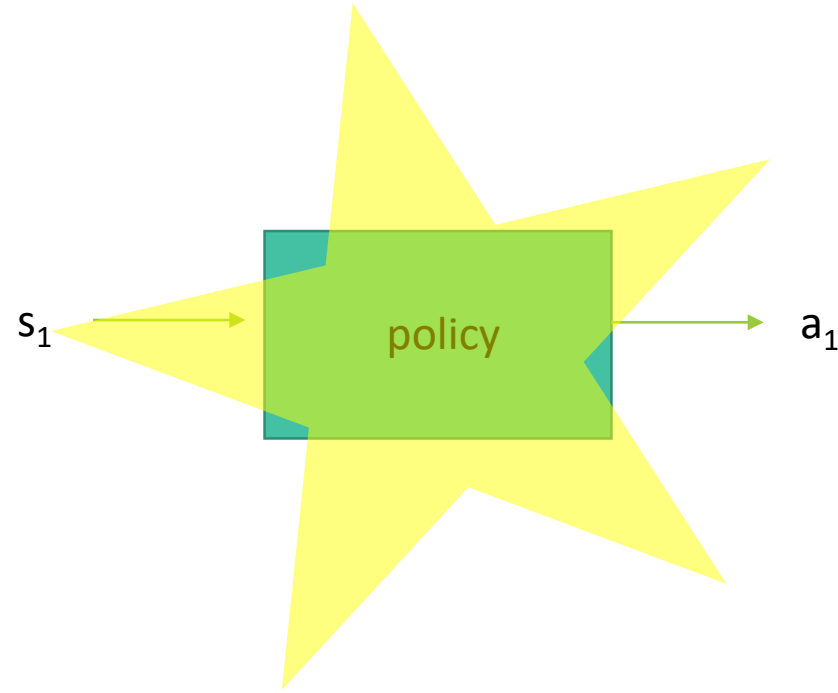Rewards: $r(s_n, a_m)$

# Mathematical formulation



Agent

$s_1$ → policy → $a_1$

**Policy:** Action that determines the next state.

**Reward:** Measure of the outcome of taking an action.

# Mathematical formulation



$s_1$ → policy → $a_1$

Agent

Choosing the **optimal policy**?

# Choosing the optimal policy

Agent

Greedy strategy:

# Choosing the optimal policy



Agent

Random strategy:
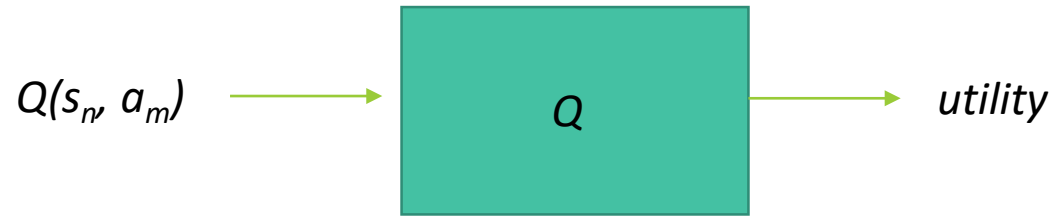


Multiple random scenarios should be tested before choosing the next action.

# Mathematical formulation



$Q(s_n, a_m)$ ⟶ Q ⟶ *utility*

Agent

**Utility:** The long term reward is called a utility.
The utility of performing an action a at *a* state *s* is written as a function *Q(s, a)*, called the utility function.

# Mathematical formulation

Agent

If you were given the utility function Q(s,a), how would you use it to derive a policy function?

# Mathematical formulation



Agent

If you were given the utility function Q(s,a), how would you use it to derive a policy function?

$$Policy(s) = argmax_a \ Q(s_n, a_m)$$

# Mathematical formulation
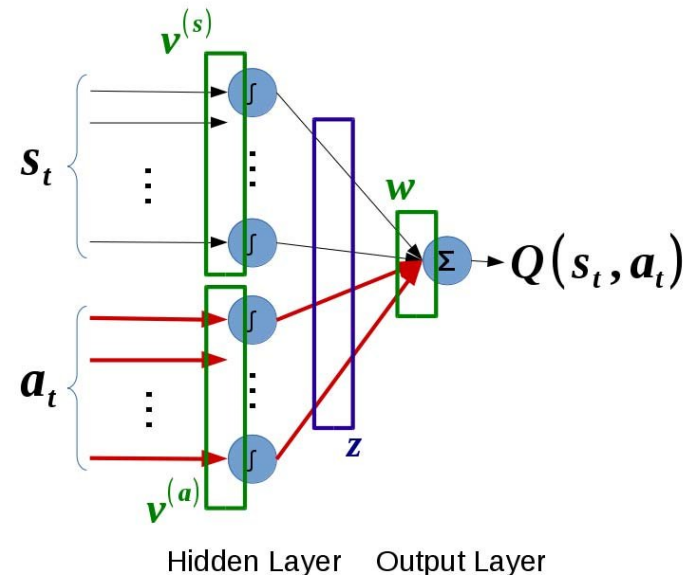
Agent

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

# Where is the deep part?

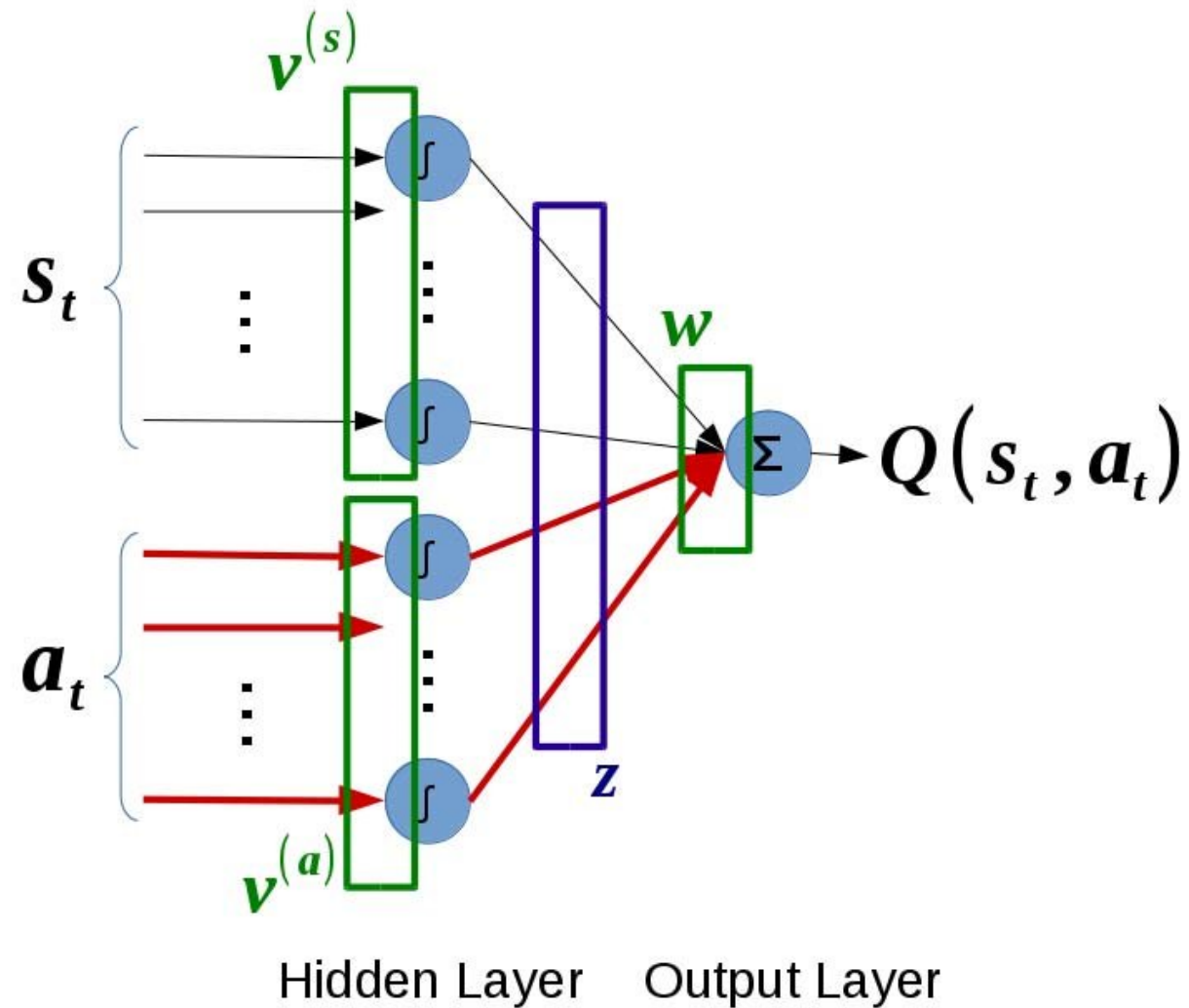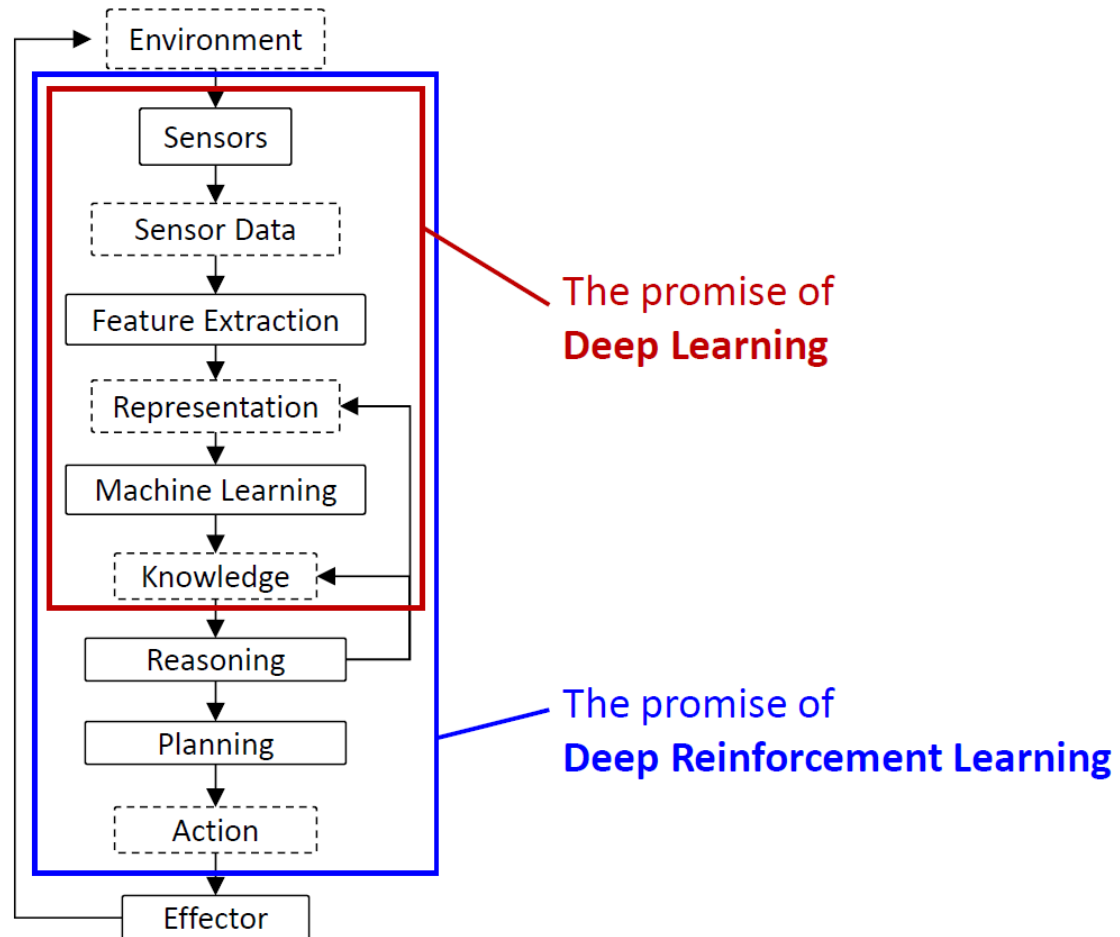If you were given the utility function Q(s,a), how would you use it to derive a policy function?

Agent

$$Policy(s) = argmax_a\ Q(s_n, a_m)$$



Hidden Layer     Output Layer

# Where is the deep part?



Agent

$v^{(s)}$

$s_t$

$v^{(a)}$

$a_t$

$w$

$z$

$\Sigma$

$Q(s_t, a_t)$

Hidden Layer    Output Layer

# Where is the deep part?



The promise of
**Deep Learning**

The promise of
**Deep Reinforcement Learning**

# Do humans use deep reinforcement learning?

# Practice

$ pip install yahoo-finance

# Practice

Last 10 days history of the prices

Budget Shares

# Other applications of reinforcement learning

# Keras RL

# Deep Traffic