

# Deep Reinforcement Learning: An Overview

Author: Yuxi Li

Presenter: Jiaru Zhang

Report Date: July 19th

# Contents

1

**Introduction**

---

2

**Core Elements**

---

3

**Important Mechanisms**

---

4

**Applications**

---

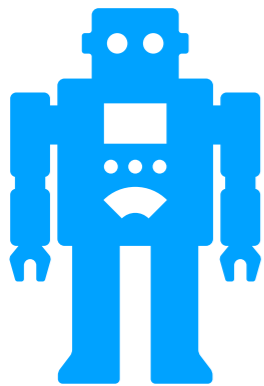
# Introduction

## Machine Learning:

Learning from data and making predictions and/or decisions

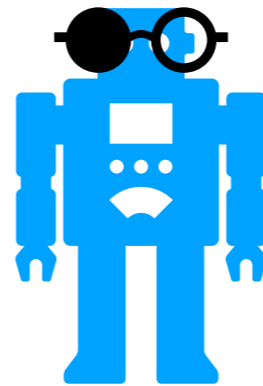
**Supervised**

I need **labels**.



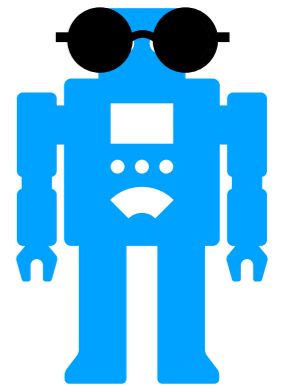
**Reinforcement**

Tell me only the **feedback**.



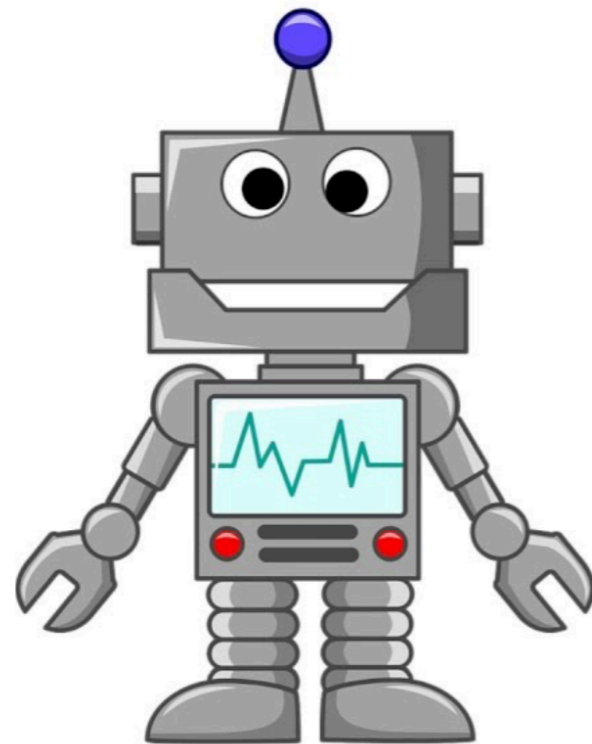
**Unsupervised**

I need **nothing**.

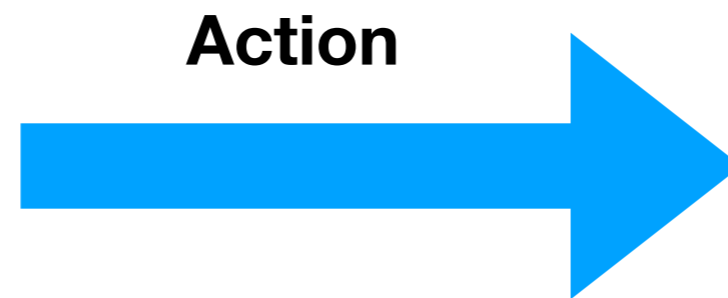


# Introduction

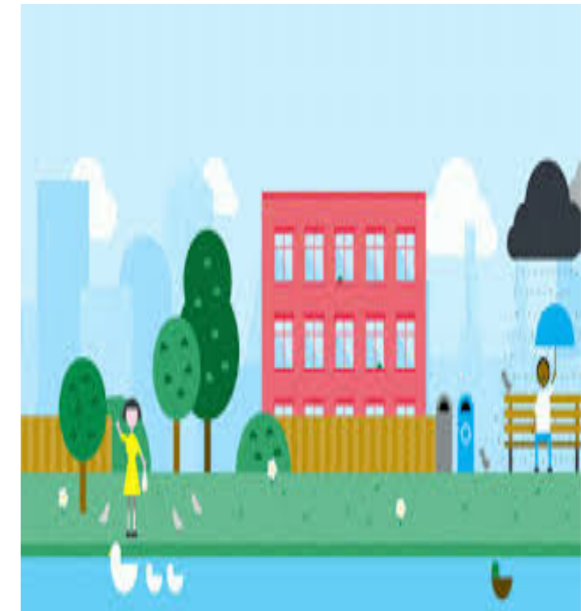
## What is RL?



**Agent**



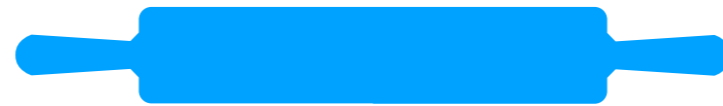
**State Reward**



**Environment**

# Introduction

Reinforcement  
Learning



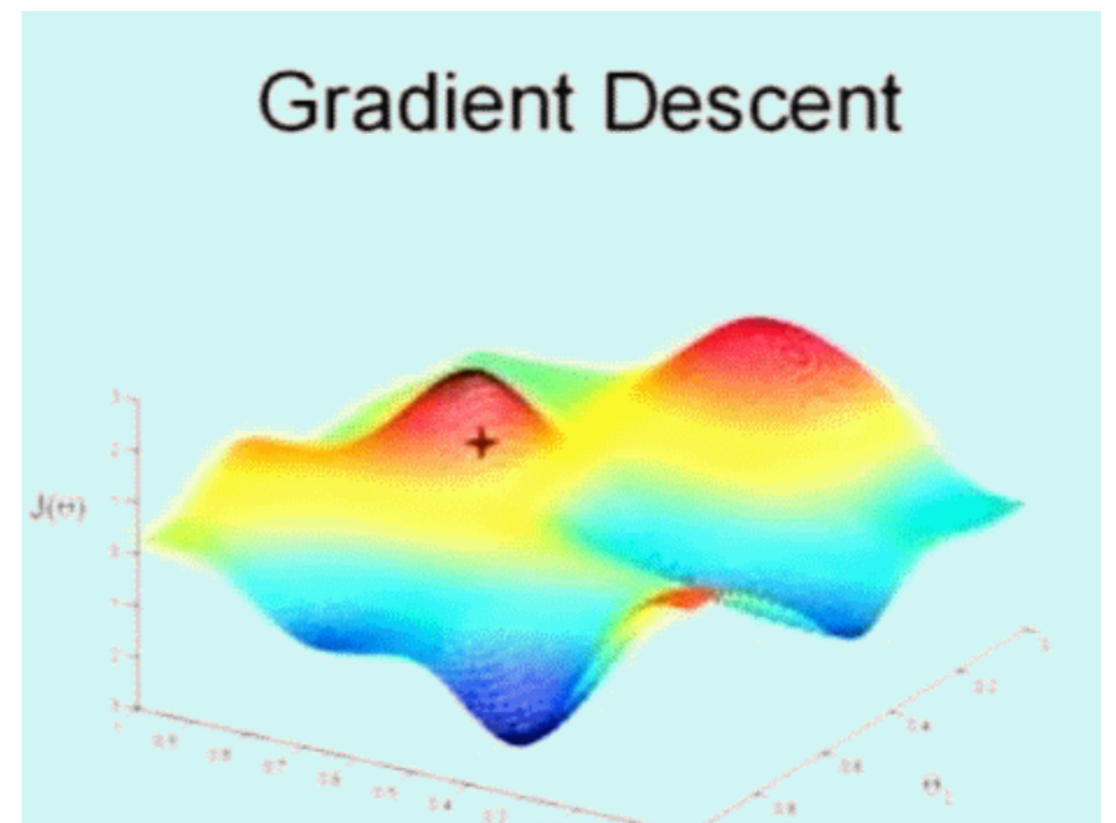
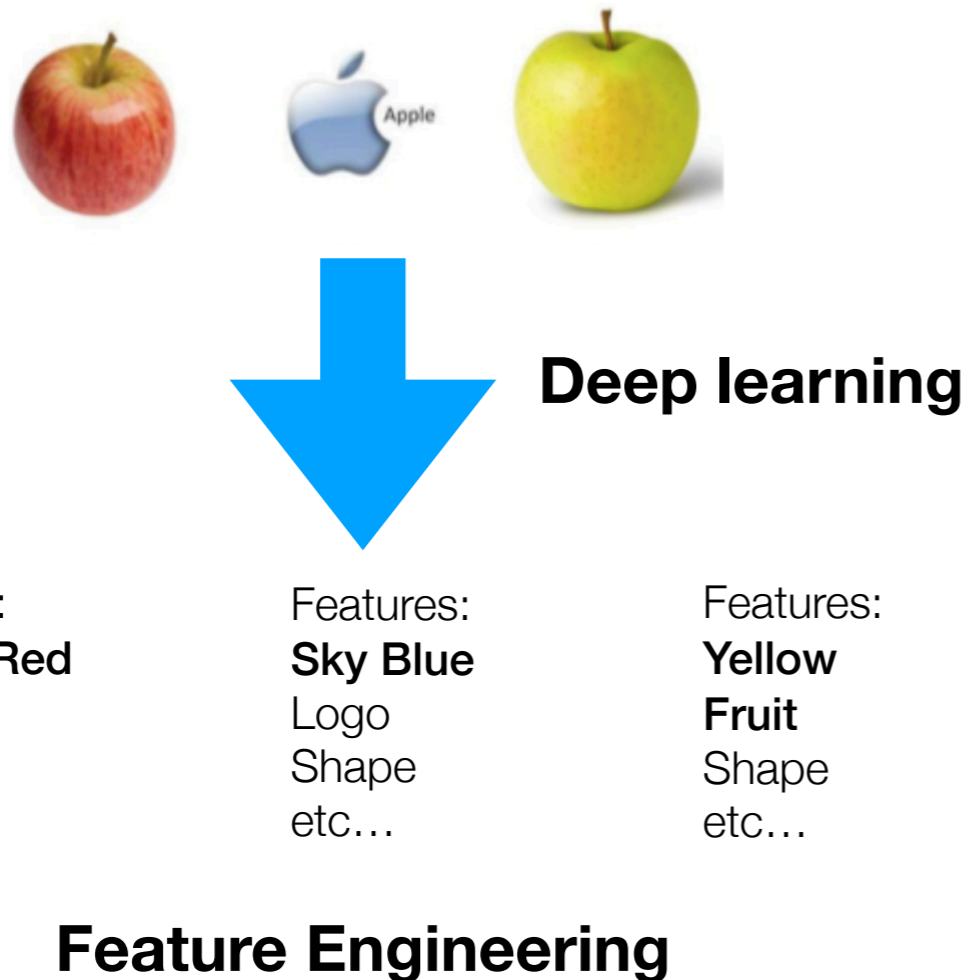
Deep Learning



- Big data
- Powerful computation
- New algorithmic techniques
- Mature software packages and architectures
- .....

# Introduction

## Why is deep learning so significant?



**End-to-end learning  
through gradient descent**

# Introduction

## Basic Reinforcement Algorithms

Value based

Policy based

Input: state  $S$  and action

Output: Value  $V(S, a; \theta)$

# Introduction

## Basic Reinforcement Algorithms

Value based

Policy based

### Temporal Difference(TD) learning

$$\begin{aligned} V(s) &\leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)] \\ &= (1 - \alpha)V(s) + \alpha(r + \gamma V(s')) \end{aligned}$$



# Introduction

## Basic Reinforcement Algorithms

Value based

Policy based

### SARSA

$$\begin{aligned} Q(s, a) &\leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)] \\ &= (1 - \alpha)Q(s, a) + \alpha(r + \gamma Q(s', a')) \end{aligned}$$

# Introduction

## Basic Reinforcement Algorithms

Value based

Policy based

### Q learning

$$\begin{aligned} Q(s, a) &\leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)] \\ &= (1 - \alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a')) \end{aligned}$$

# Introduction

## Basic Reinforcement Algorithms

Value based

Policy based

Q learning

Offline!

$$\begin{aligned} Q(s, a) &\leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)] \\ &= (1 - \alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a')) \end{aligned}$$

# Introduction

## Basic Reinforcement Algorithms

Value based

Policy based

### TD learning with Function approximation

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [r + \gamma \hat{v}(s', \mathbf{w}) - \hat{v}(s, \mathbf{w})] \nabla \hat{v}(s, \mathbf{w})$$

# Introduction

## Basic Reinforcement Algorithms

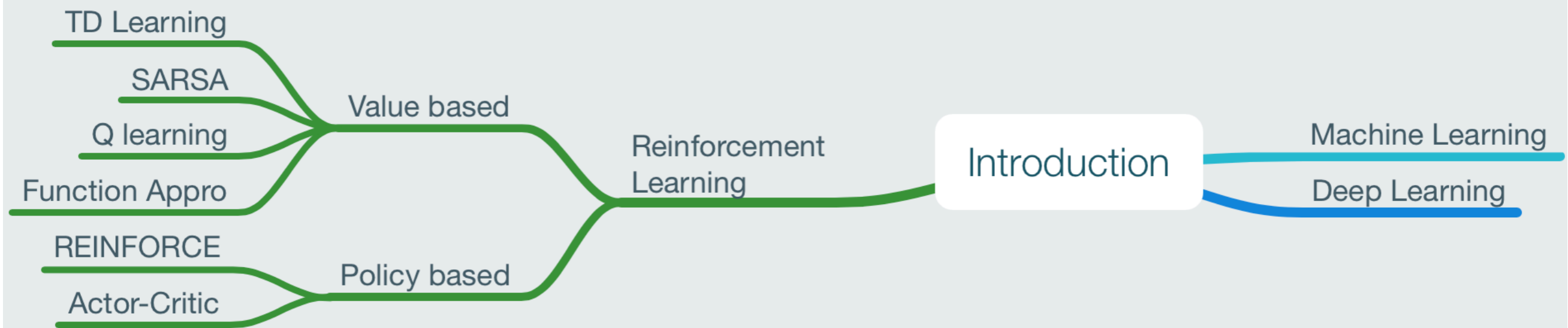
Value based

Policy based

Input: state  $S$

Output: Policy  $a(S)$

- REINFORCE
- Actor-Critic



# Contents

1

**Introduction**

---

2

**Core Elements**

---

3

**Important Mechanisms**

---

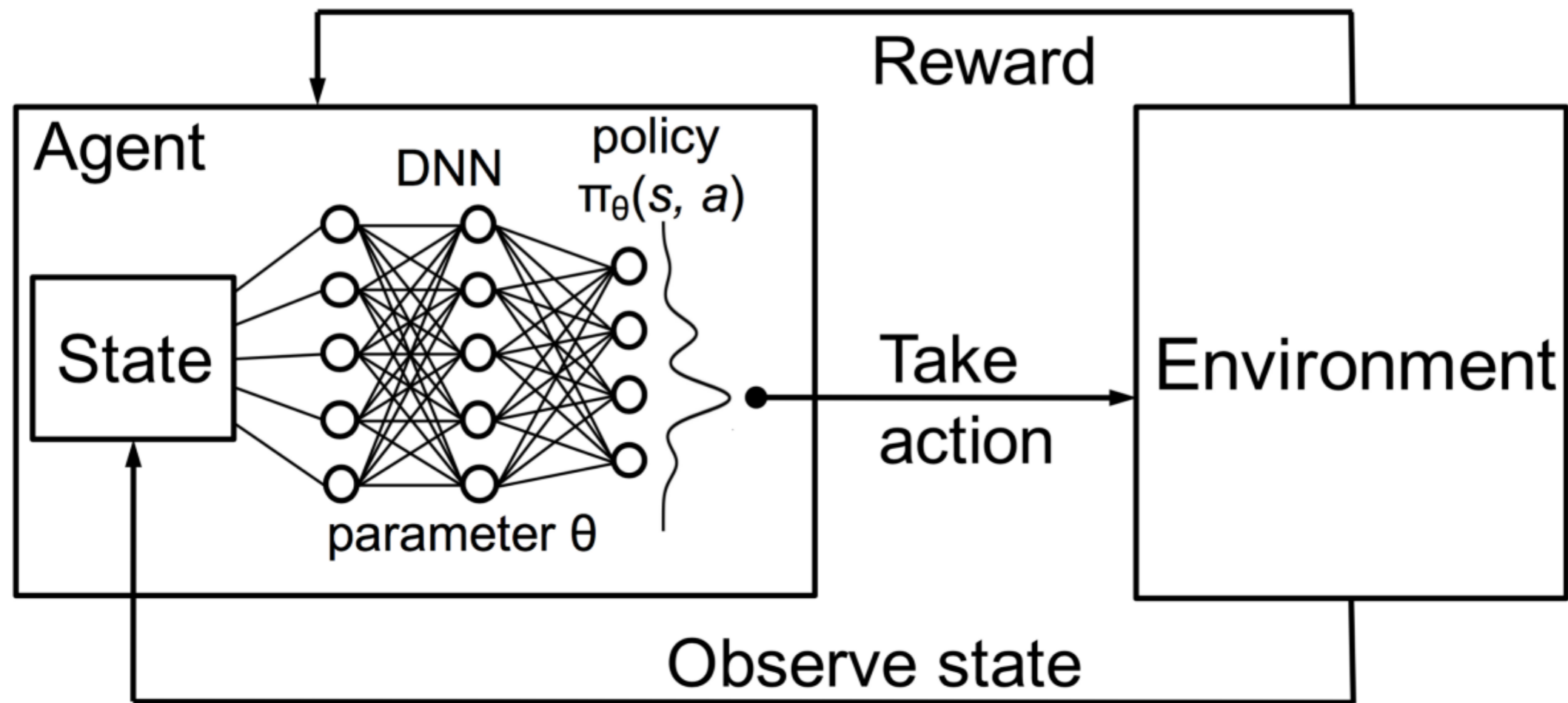
4

**Applications**

---

# Core Elements

## Deep Q-Learning(DQN)





# Core Elements

## Deep Q-Learning(DQN)

### Advantages:

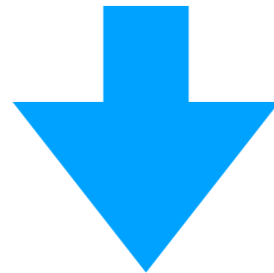
- Experiment Replay
- End-to-end RL approach
- General network framework

# Core Elements

## Double DQN(D-DQN)

Original DQN:

$$y_t^Q = r_{t+1} + \gamma Q(s_{t+1}, \underset{a}{\operatorname{arg\,max}} Q(s_{t+1}, a; \theta_t); \theta_t)$$



**Over Estimation!**

# Core Elements

## Double DQN(D-DQN)

Solution:

Use one network to select,  
use another network to evaluate.

$$y_t^{D-DQN} = r_{t+1} + \gamma Q(s_{t+1}, \underset{a}{\operatorname{arg\,max}} Q(s_{t+1}, a; \theta_t); \theta_t^-)$$

# Core Elements

## Double DQN(D-DQN)

Two methods:

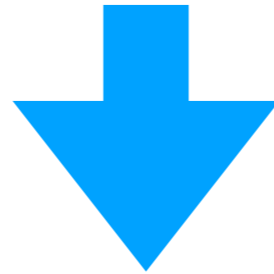
Four networks: Double the original architecture.

Two networks: Just use the target network to evaluate.

# Core Elements

## Prioritized Experience Replay

$$TDerror = \Delta Q$$



prioritization  $p$

Choose samples according to  $p$  by SumTree

# Core Elements

## Dueling Architecture

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + (A(s, a; \theta, \alpha) - \frac{a}{|\mathcal{A}|} A(s, a'; \theta, \alpha))$$

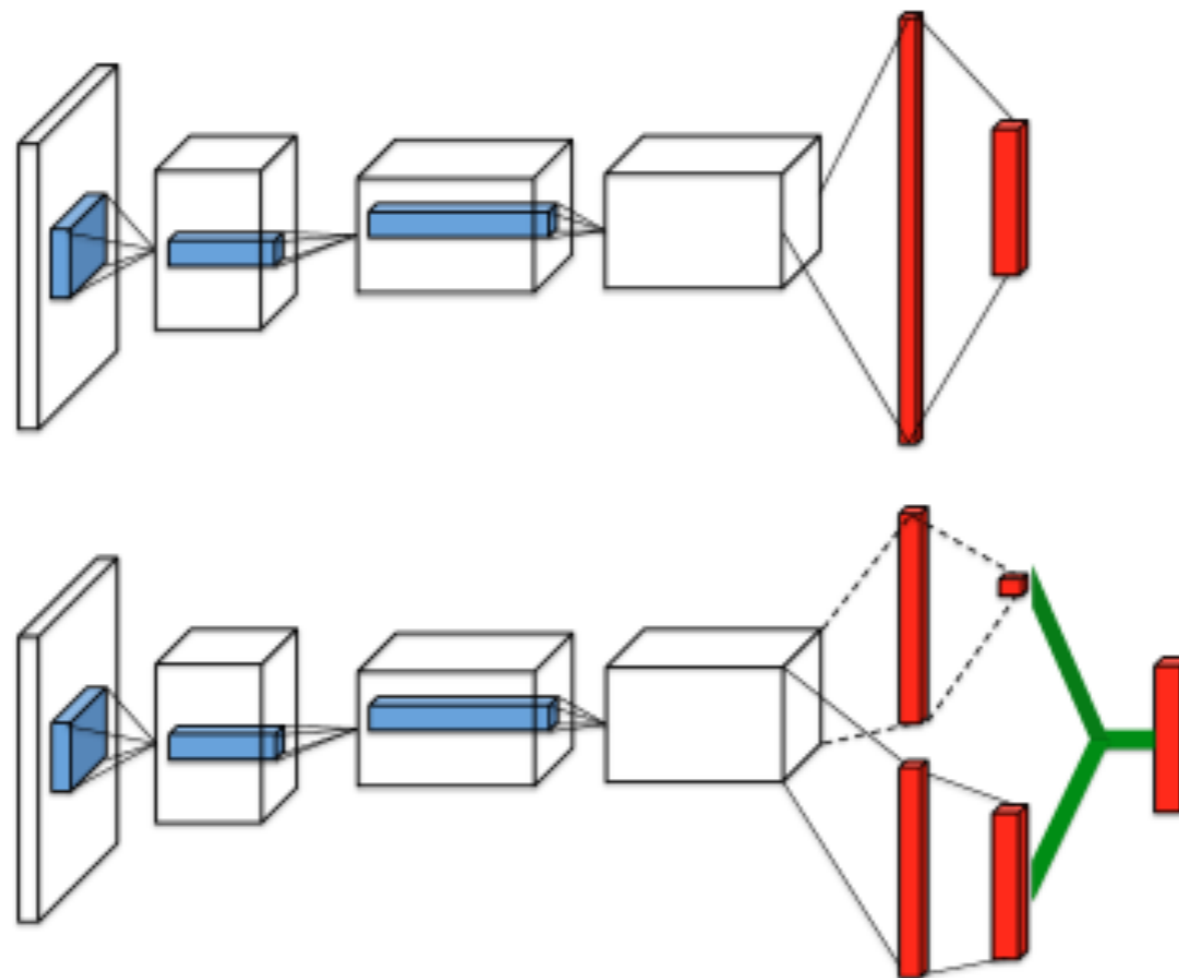


Figure Source:

Wang, Z., Schaul, T., Hessel, M., van Hasselt, H., Lanctot, M., & de Freitas, N. (2015). Dueling Network Architectures for Deep Reinforcement Learning, (9). <https://doi.org/10.1109/MCOM.2016.7378425>

# Core Elements

## Dueling Architecture

Value: Only consider state

Advantage: Consider both state and action

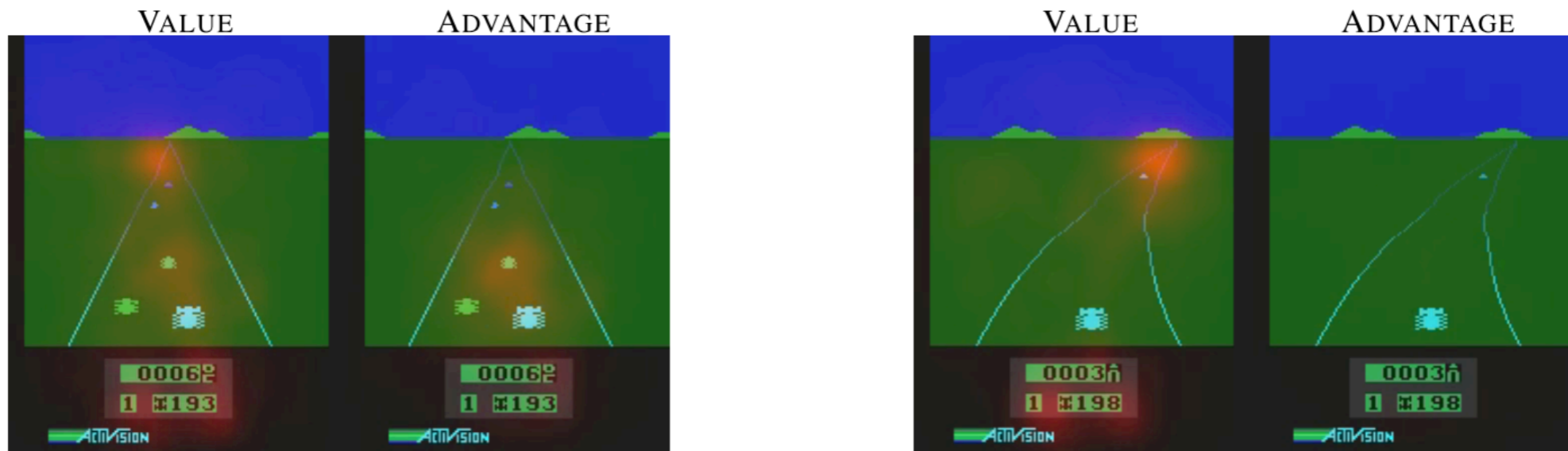


Figure Source:

Wang, Z., Schaul, T., Hessel, M., van Hasselt, H., Lanctot, M., & de Freitas, N. (2015). Dueling Network Architectures for Deep Reinforcement Learning, (9). <https://doi.org/10.1109/MCOM.2016.7378425>

# Core Elements

## Policy based Deep algorithms

- Policy Gradient
- Asynchronous Advantage Actor-critic(A3C)



# Core Elements

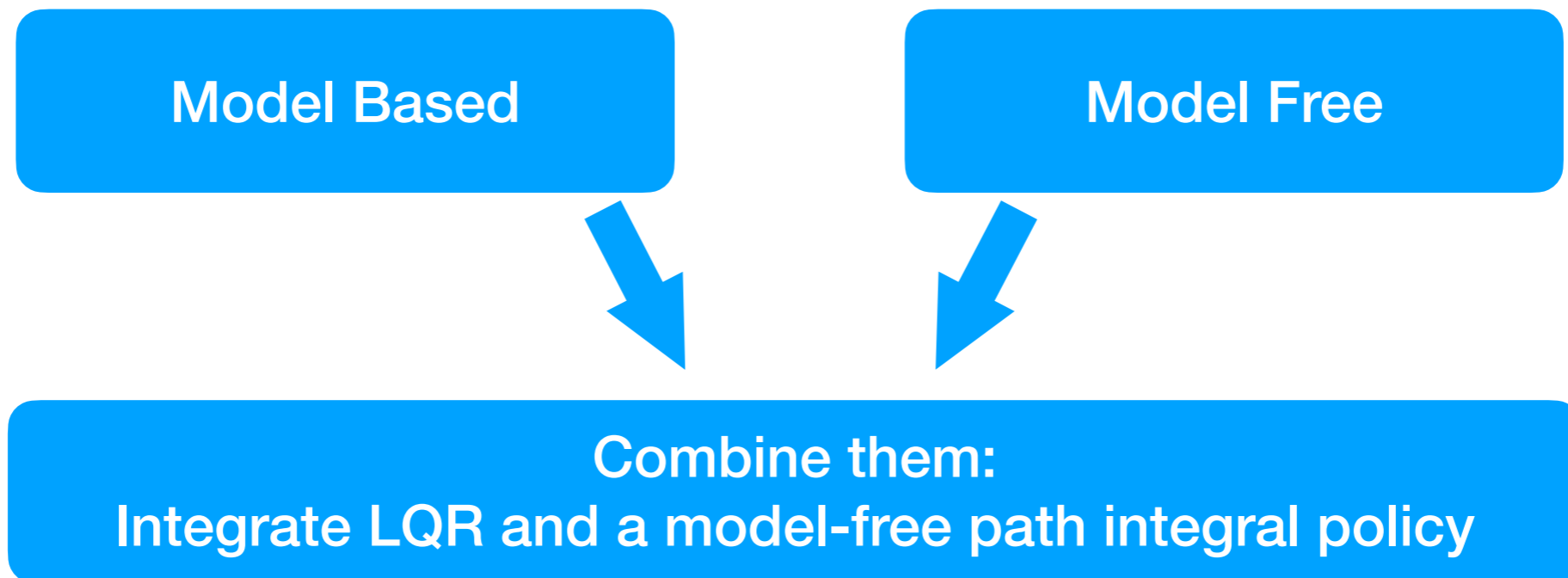
## Reward

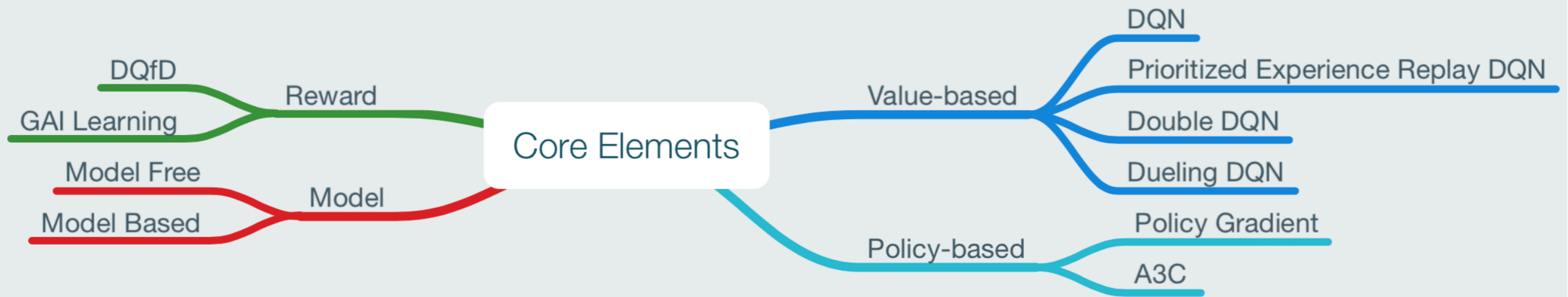
- Learning from demonstration:  
Deep Q-learning from Demonstrations(DQfD)
- Generative Adversarial Imitation Learning

# Core Elements

## Model

Model: An agent's representation of the environment.  
e.g. transition model and the reward model.





# Contents

1

**Introduction**

---

2

**Core Elements**

---

3

**Important Mechanisms**

---

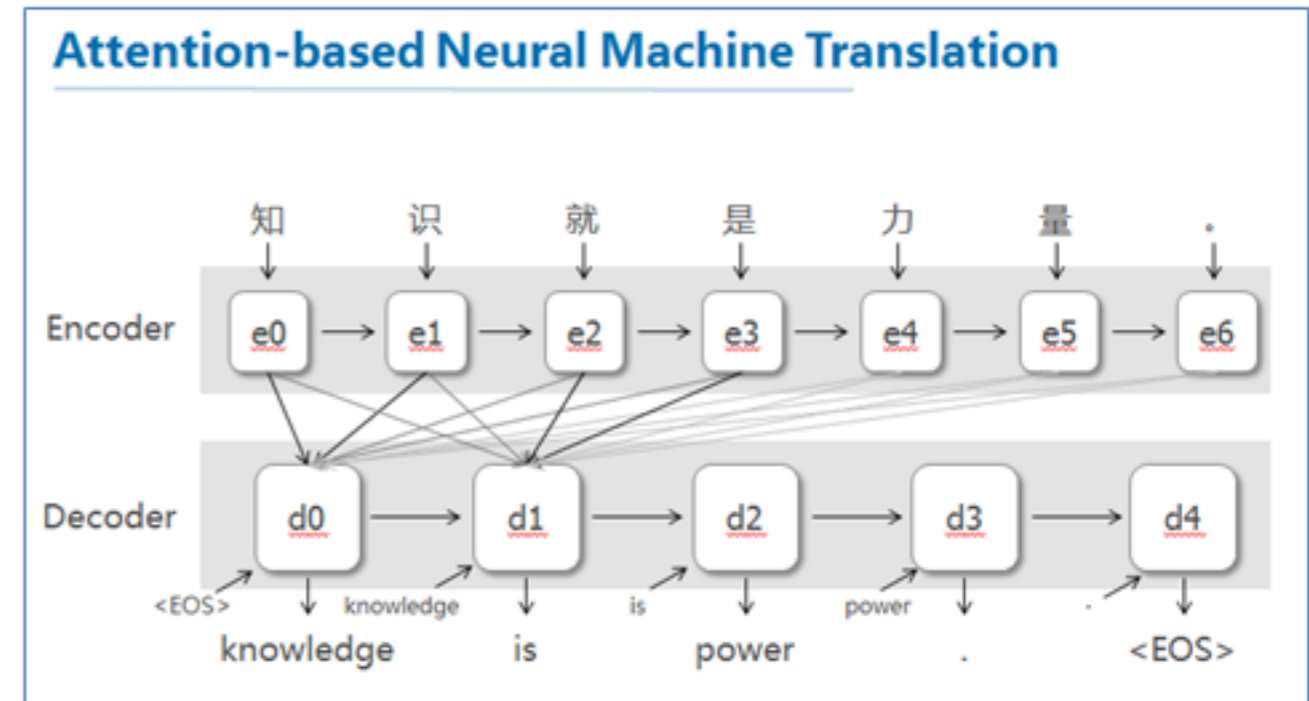
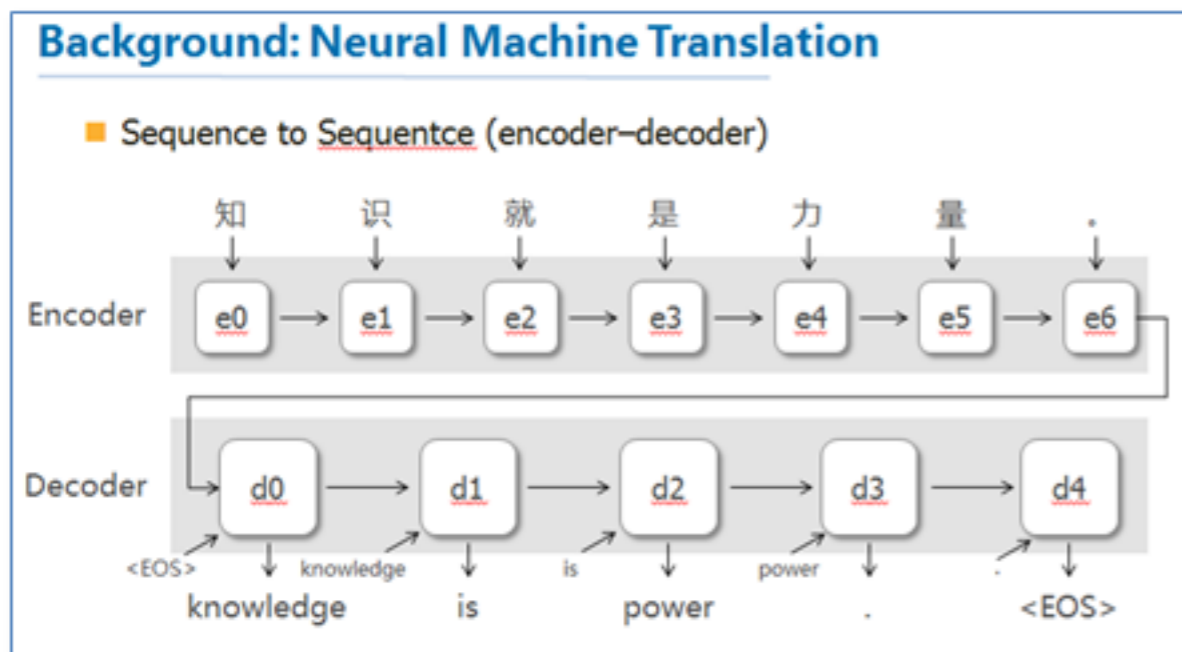
4

**Applications**

---

# Important Mechanisms

## Attention



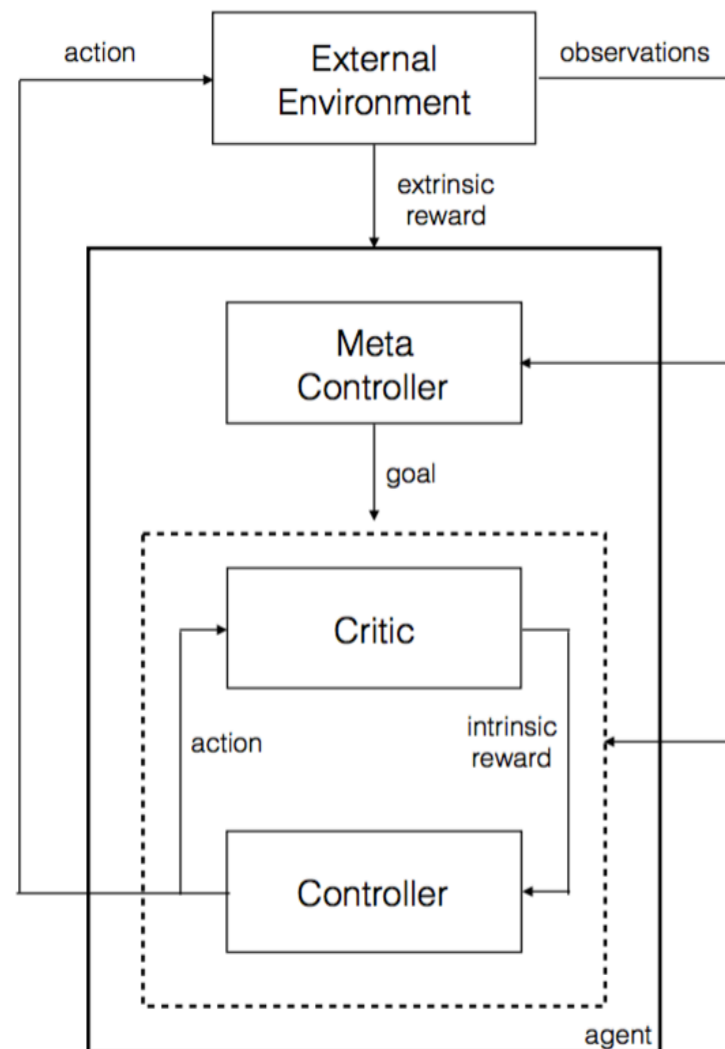
# Important Mechanisms

## Unsupervised Learning

- Unsupervised REinforcement and Auxiliary Learning(UNREAL)
- Generative Adversarial Networks(GAN)

# Important Mechanisms

## Hierarchical Reinforcement Learning



Use two networks:  
One to determine destination,  
another to determine concrete action.

# Important Mechanisms

## Other Topics

- Transfer Learning
- Multi-Agent Reinforcement Learning
- Learning to Learn



Hierarchical Reinforcement Learning

Transfer Learning

Learning to Learn

Important Mechanisms

Multi-Agent Reinforcement Learning

Unsupervised Learning

Attention

UNREAL

GAN



# Contents

1

**Introduction**

---

2

**Core Elements**

---

3

**Important Mechanisms**

---

4

**Applications**

---

# Applications

## Games

Perfect Information Broad Games

Computer Go:  AlphaGo

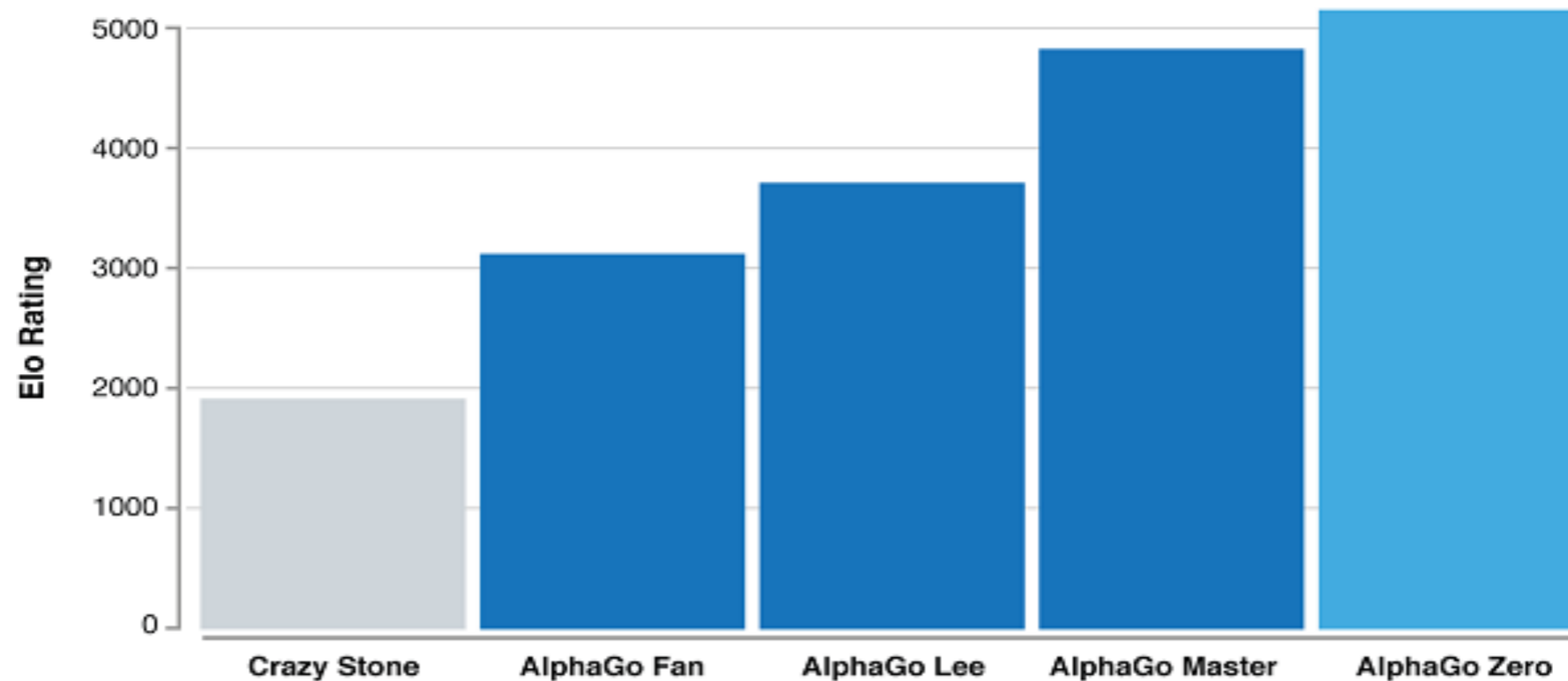


Figure Source:  
<https://deepmind.com/blog/alphago-zero-learning-scratch/>

# Applications

## Games

Computer Go:  AlphaGo

- Supervised Learning(SL) policy network
- Reinforcement Learning(RL) policy network
- Reinforcement Learning(RL) value network
- Monte Carlo tree search(MCTS)

# Applications

## Games

Computer Go:

**AlphaGo Zero**  
Starting from scratch

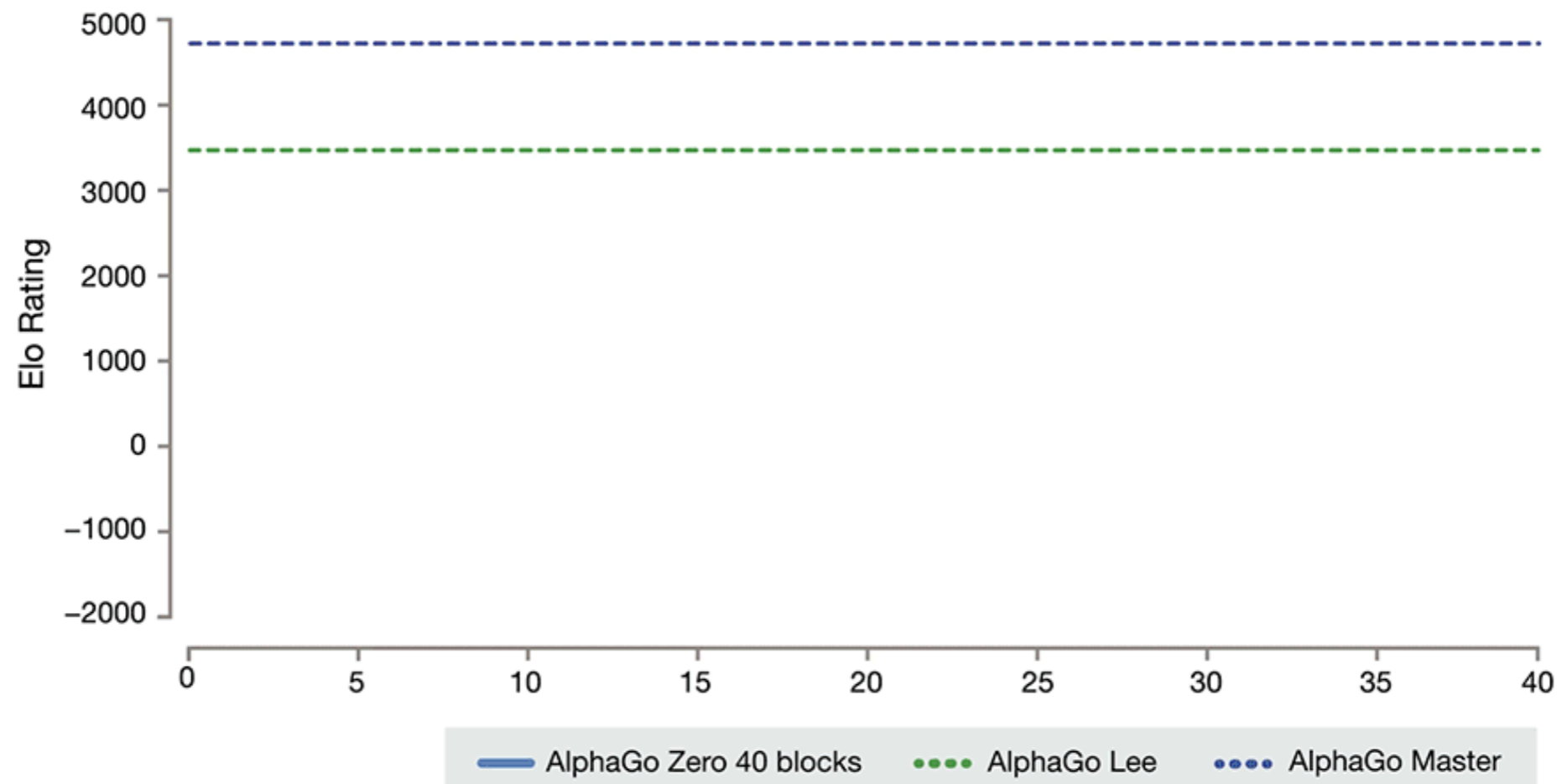


Figure Source:  
<https://deepmind.com/blog/alphago-zero-learning-scratch/>

# Applications

## Games

Computer Go:

The logo for AlphaGo Zero, featuring the text "AlphaGo Zero" in a large, white, sans-serif font, with the subtitle "Starting from scratch" in a smaller, white, sans-serif font below it. The text is set against a dark, rectangular background with a subtle grid pattern.

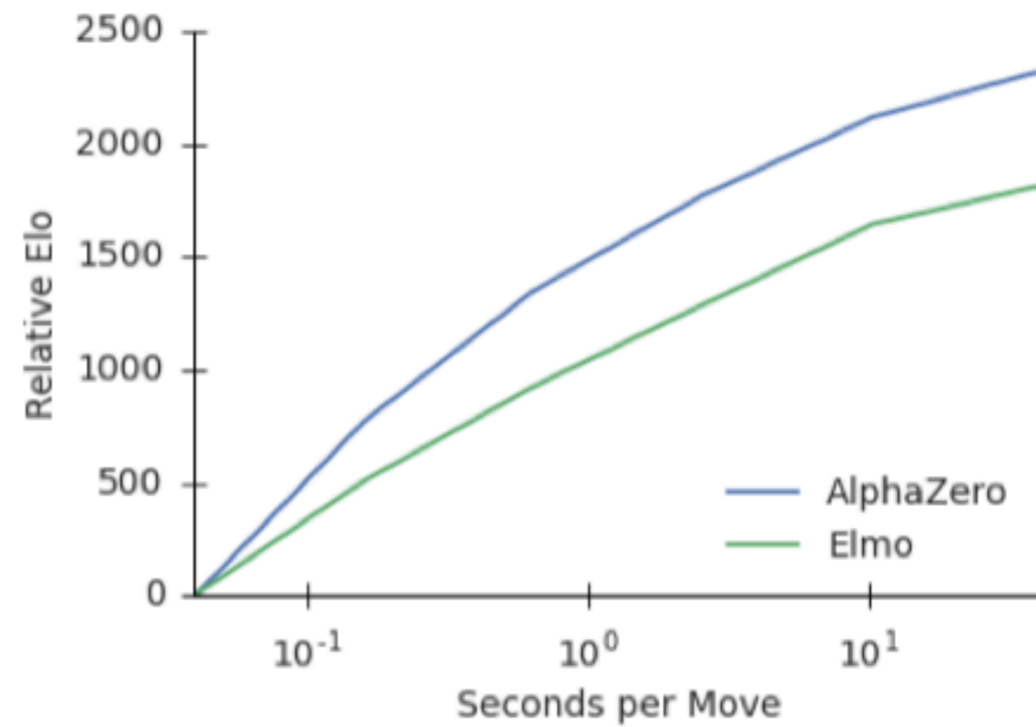
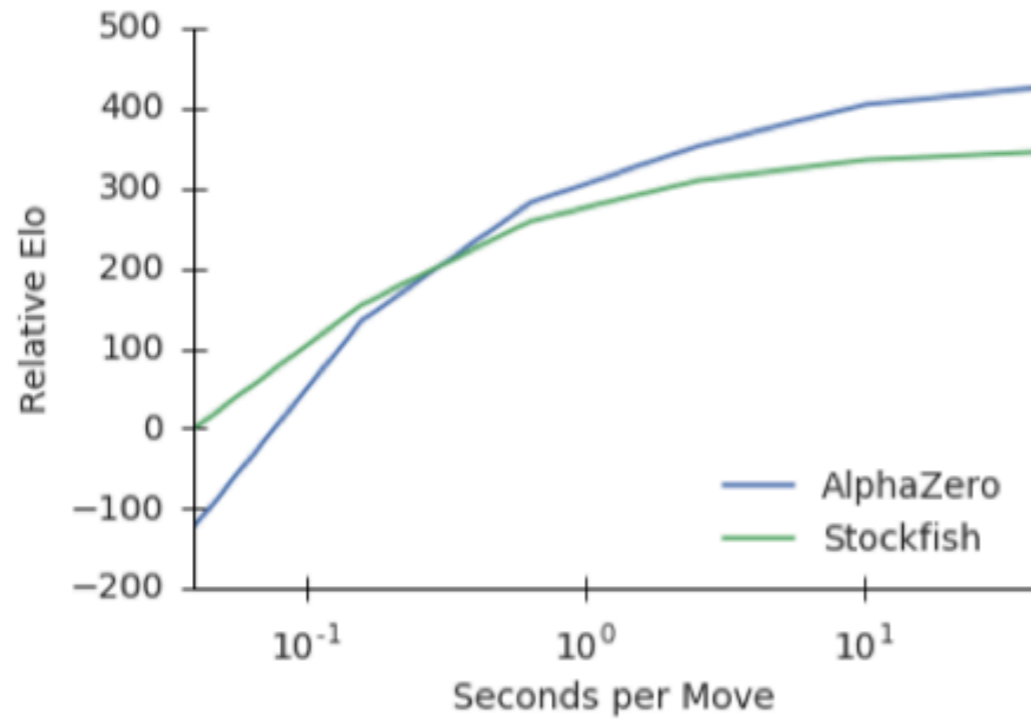
AlphaGo Zero  
Starting from scratch

- No human knowledge
- Only one Reinforcement Learning(RL) network
- Monte Carlo tree search(MCTS)

# Applications

## Games

More Generally: Alpha Zero



# Applications

## Other applications

- Imperfect Information Broad Games
- Video Games
- Robotics
- Natural Language Processing(NLP)
- Computer Vision(CV)
- Neural Architecture Design
- Healthcare
- Computer Systems
- ...



