

2017年臺大電機系大學部專精獎

強化式學習與卷積神經網絡於訓練 視覺化電腦代理人之應用

電機四
劉廷緯 溫明浩



GAMES & AI

AlphaGo - Google DeepMind



GAMES & AI

AlphaGo - Google DeepMind

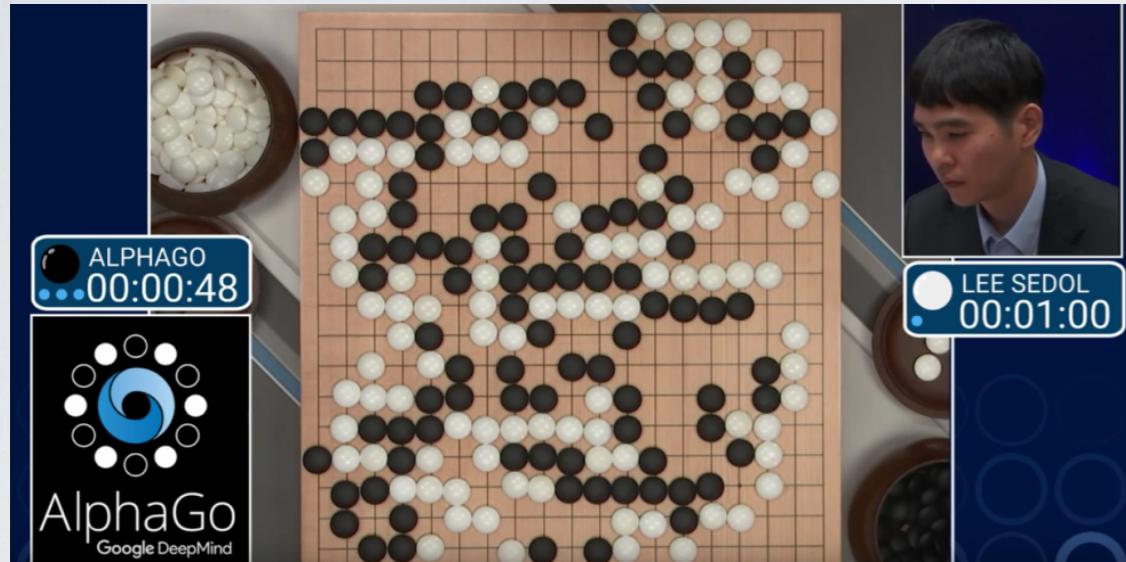


StarCraft II - Google DeepMind



GAMES & AI

AlphaGo - Google DeepMind



StarCraft II - Google DeepMind



GTA V - Harrison Kinsley

GAMES & AI

AlphaGo - Google DeepMind



StarCraft II - Google DeepMind



GTA V - Harrison Kinsley



Doom - Facebook Ai Research

Slither.io

an online massively multiplayer browser game

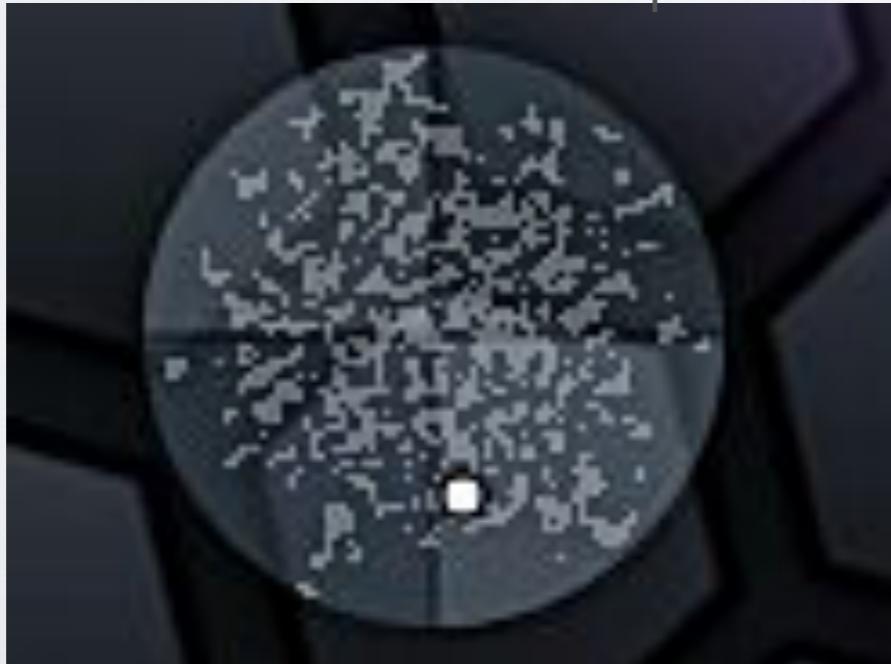


Slither.io

an online massively multiplayer browser game



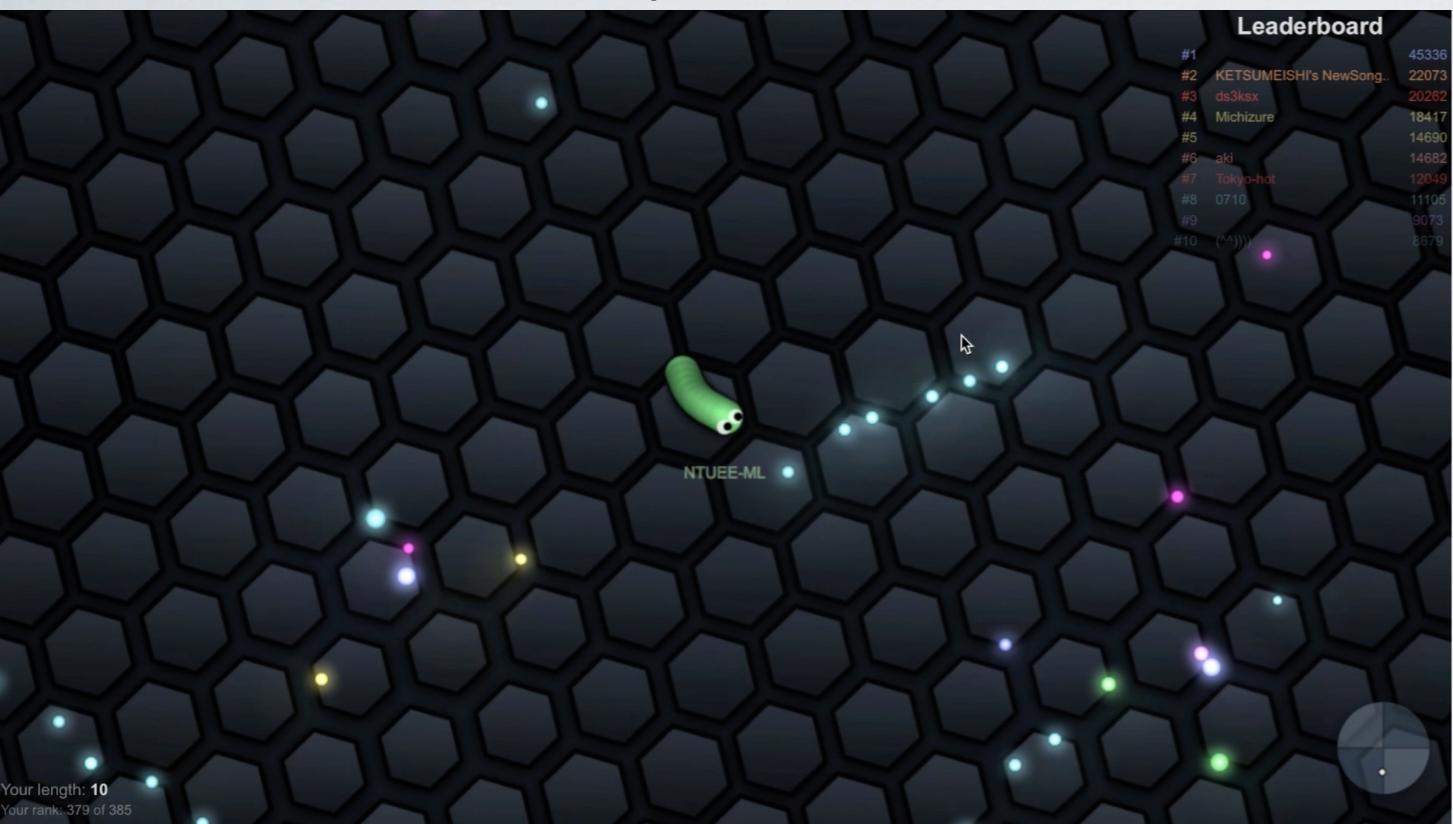
Circular map



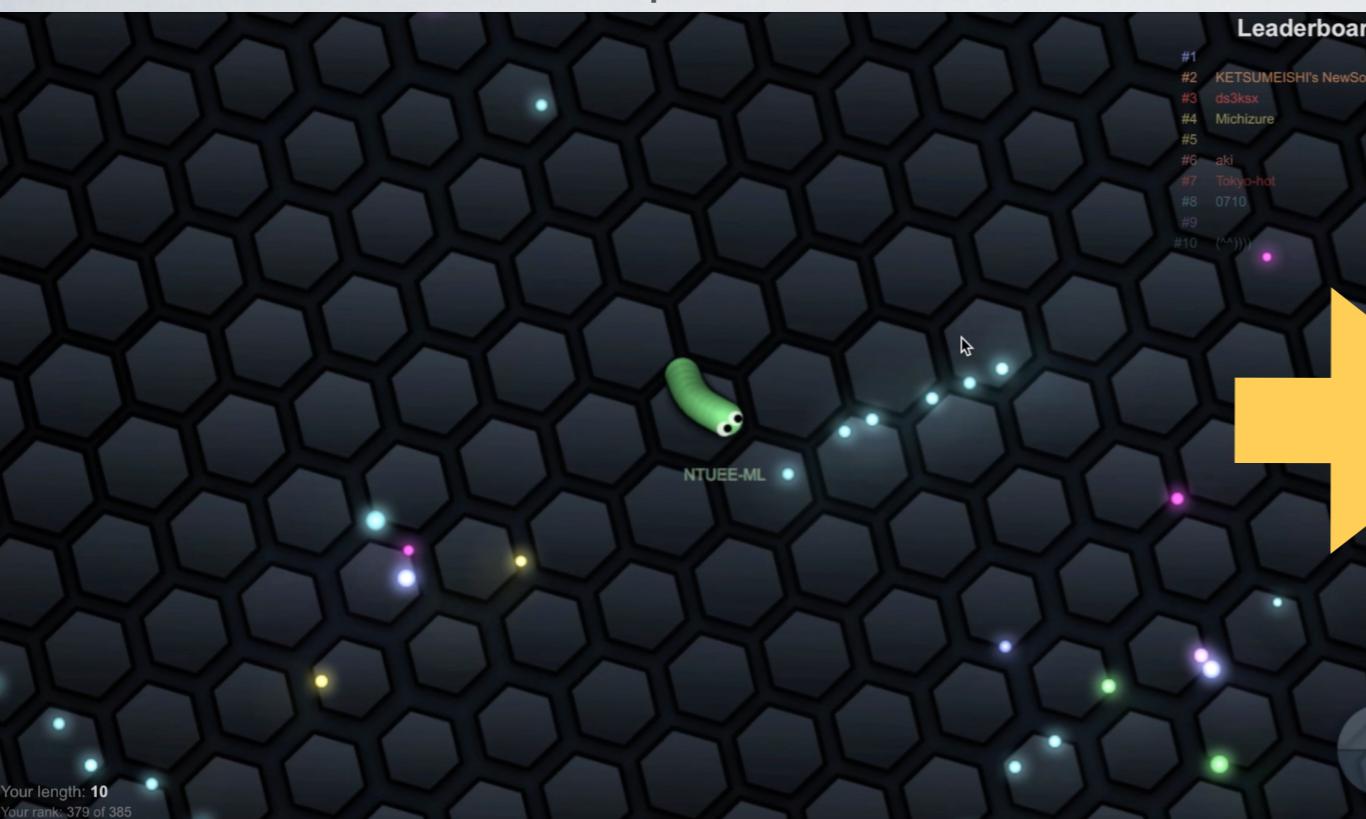
Deadly Outer rim



Spawn



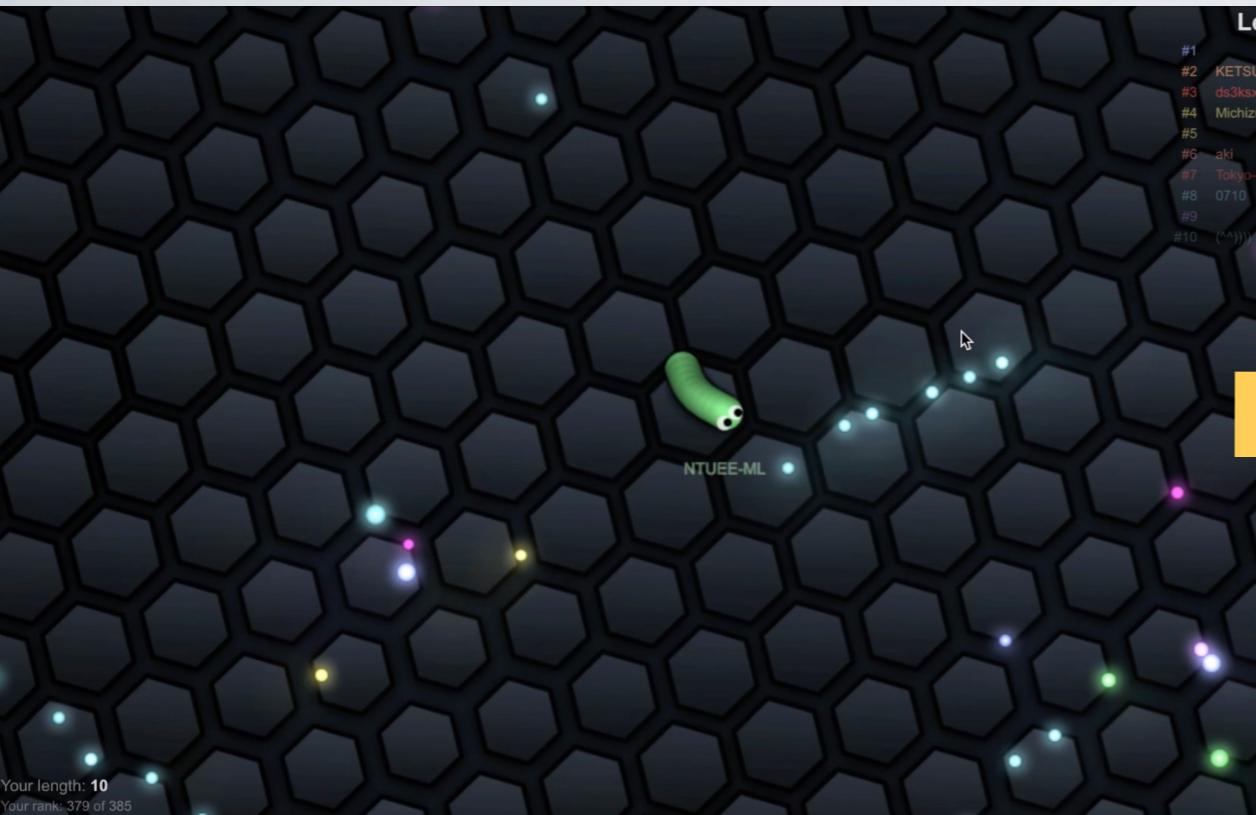
Spawn



Consume Pallets to grow



Spawn



Consume Pallets to grow



silvergames.com

Grow the largest snake

Spawn



Consume Pallets to grow



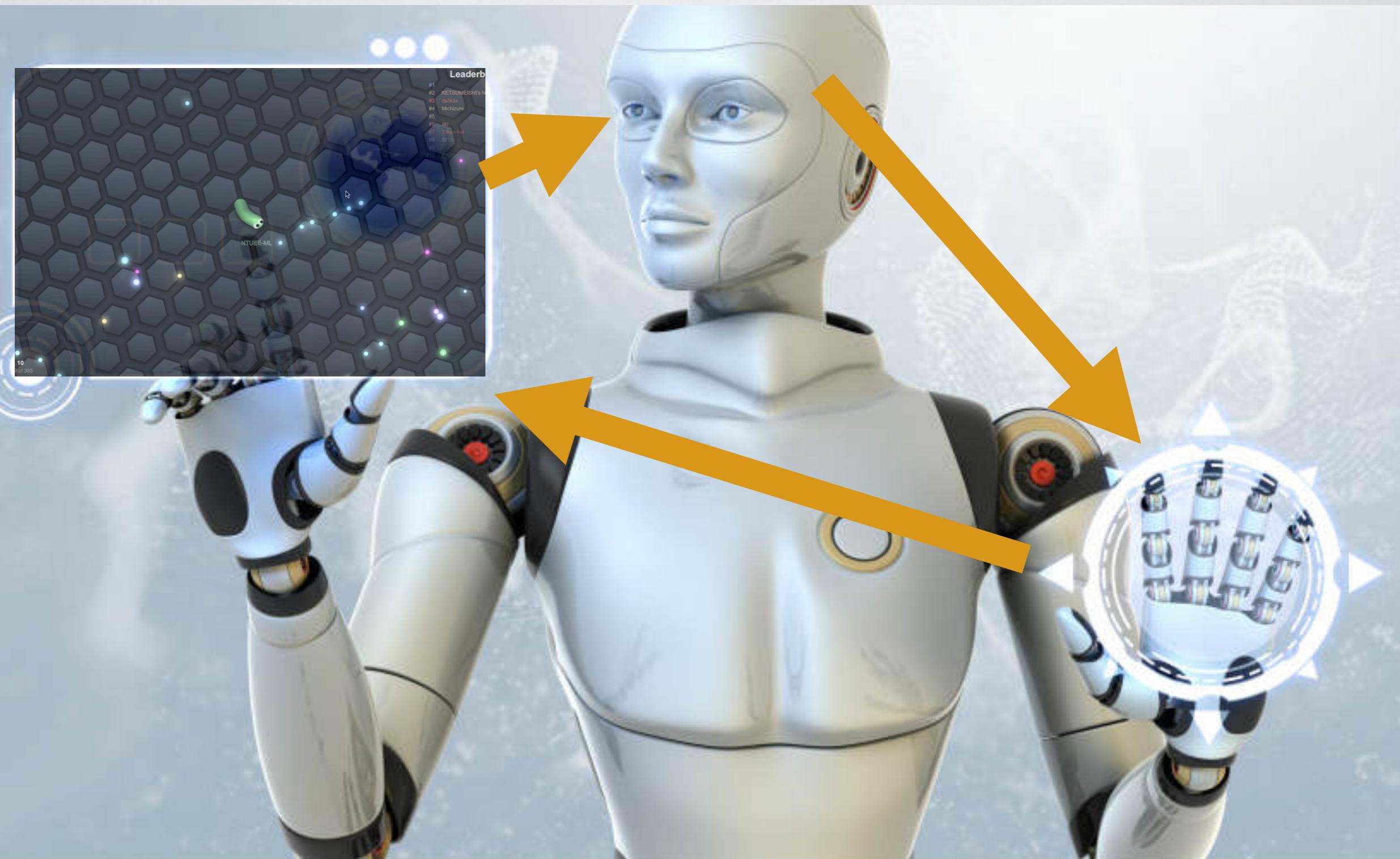
Dies during head collision

Grow the largest snake

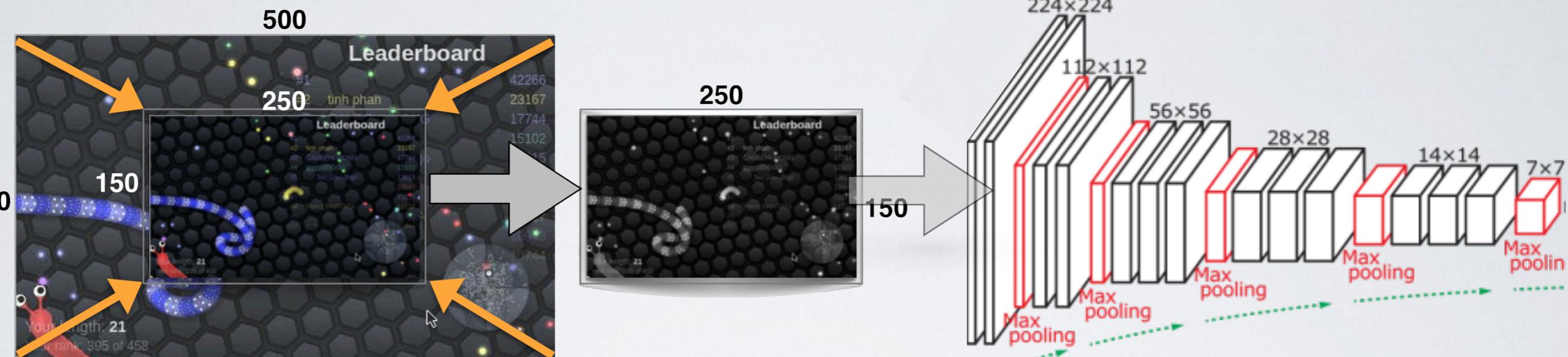
VISION-BASED AGENT



VISION-BASED AGENT



CNN



Screen Processing

Convolutional Neural
Network

ACTOR CRITIC MODEL

Goal:

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} r_t \right]$$

ACTOR CRITIC MODEL

Goal:

$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} r_t \right]$$

Actor

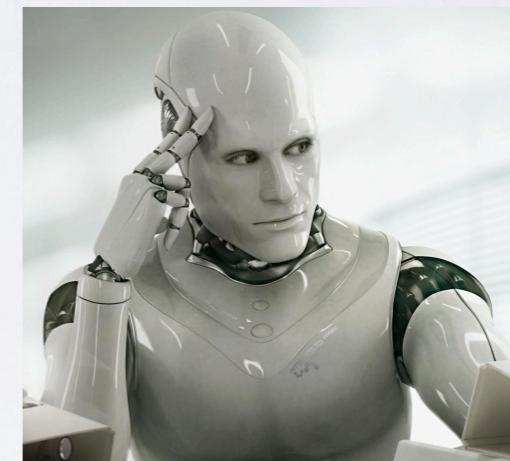


Policy function

$$\pi(a|s)$$

Learns to choose the most
rewardable action

Critic



Value function

$$V(s)$$

Learns the value of each
state



ACTOR CRITIC MODEL

Goal:

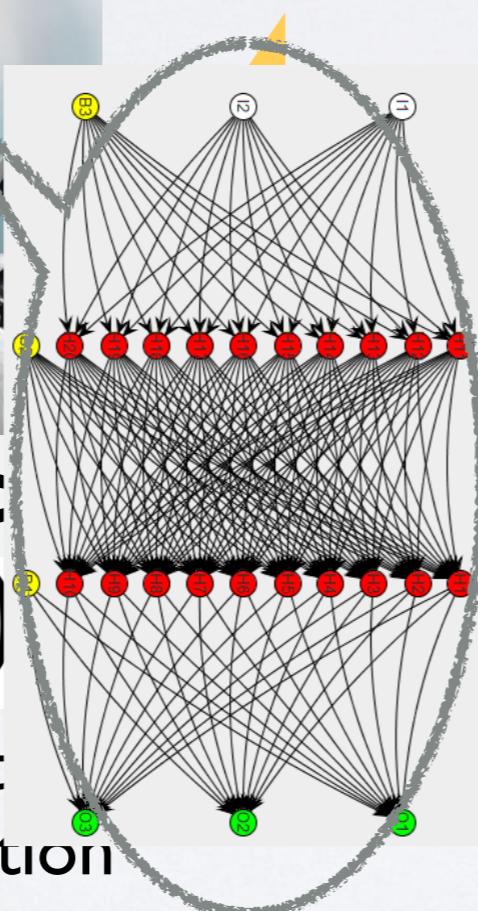
$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} r_t \right]$$

Actor

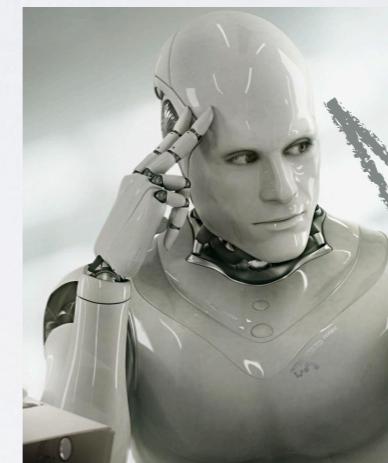


Policy func
 $\pi(a|s)$

Learns to choose the most rewardable action

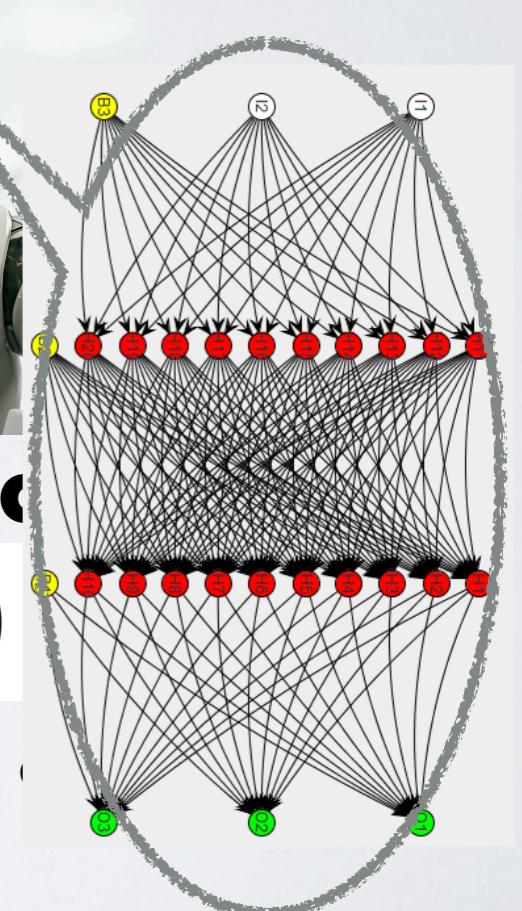


Critic



Value func
 $V(s)$

Learns the value of each state

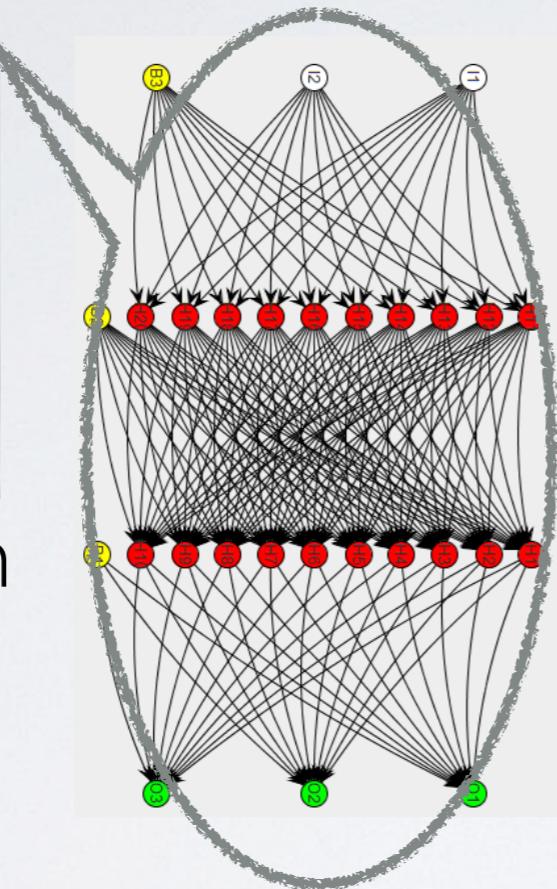


Actor

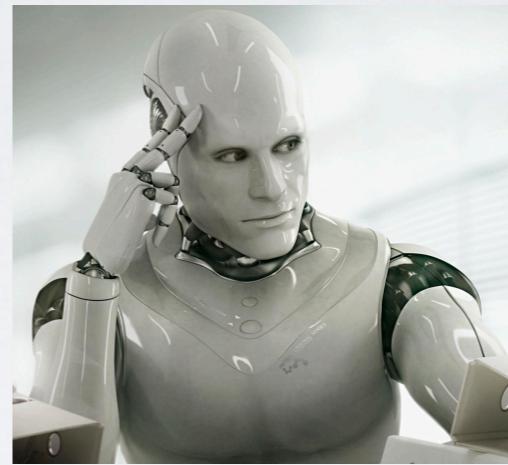


Policy function

$$\pi(a|s)$$

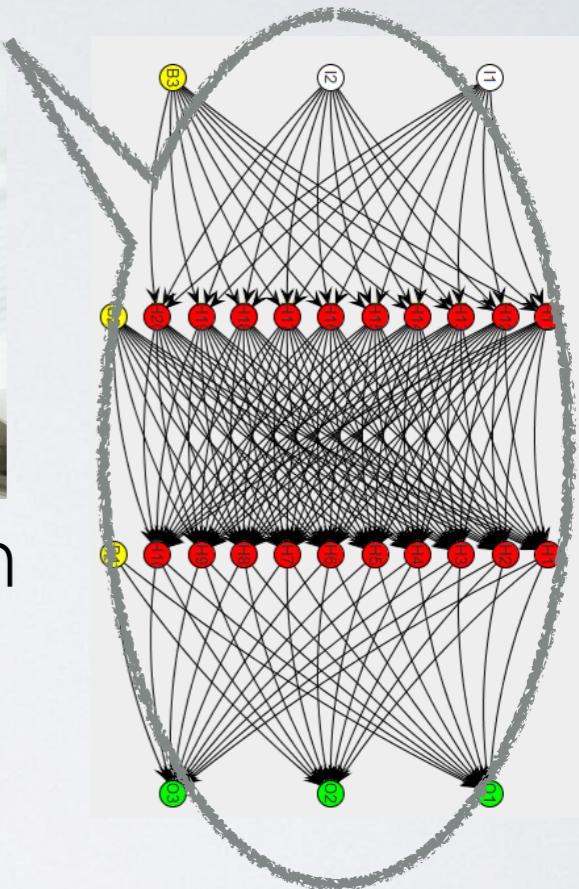


Critic



Value function

$$V(s)$$



Game Frames as states

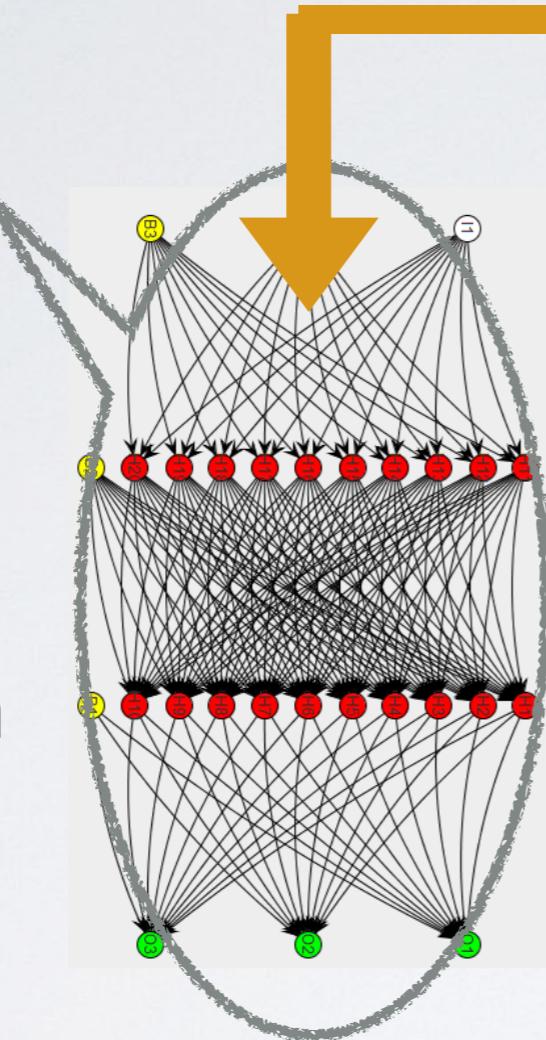


Actor

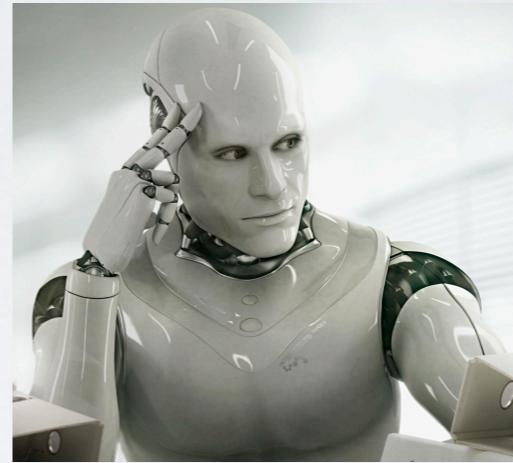


Policy function

$$\pi(a|s)$$

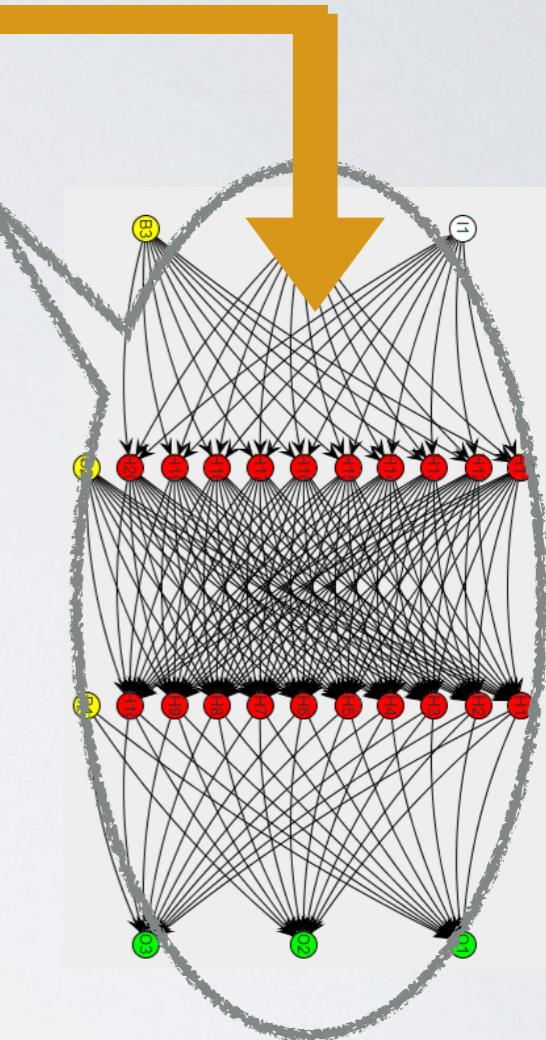


Critic



Value function

$$V(s)$$



Game Frames as states

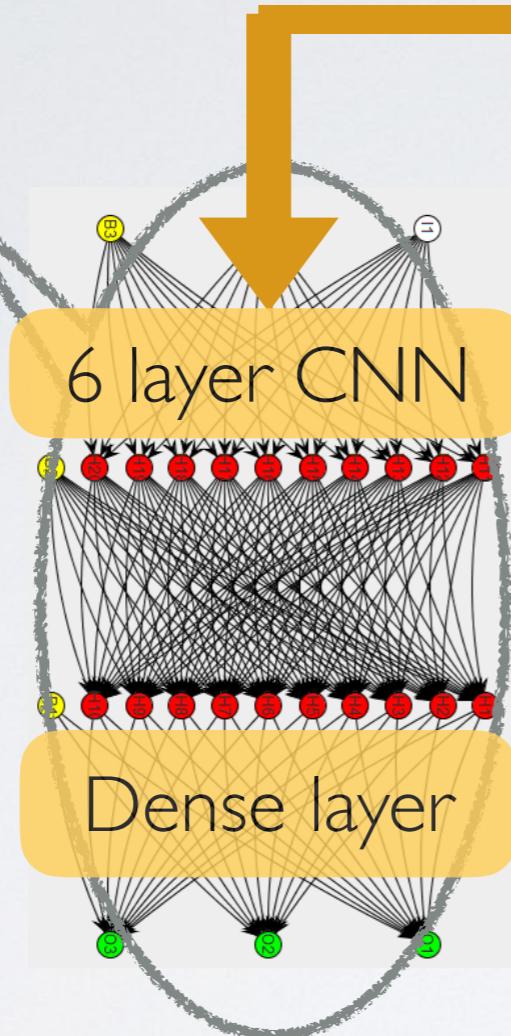


Actor

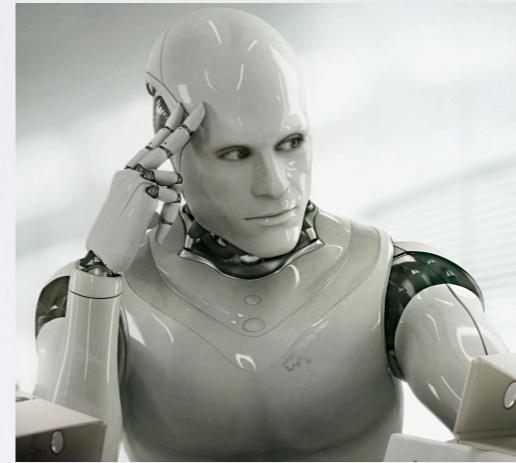


Policy function

$$\pi(a|s)$$

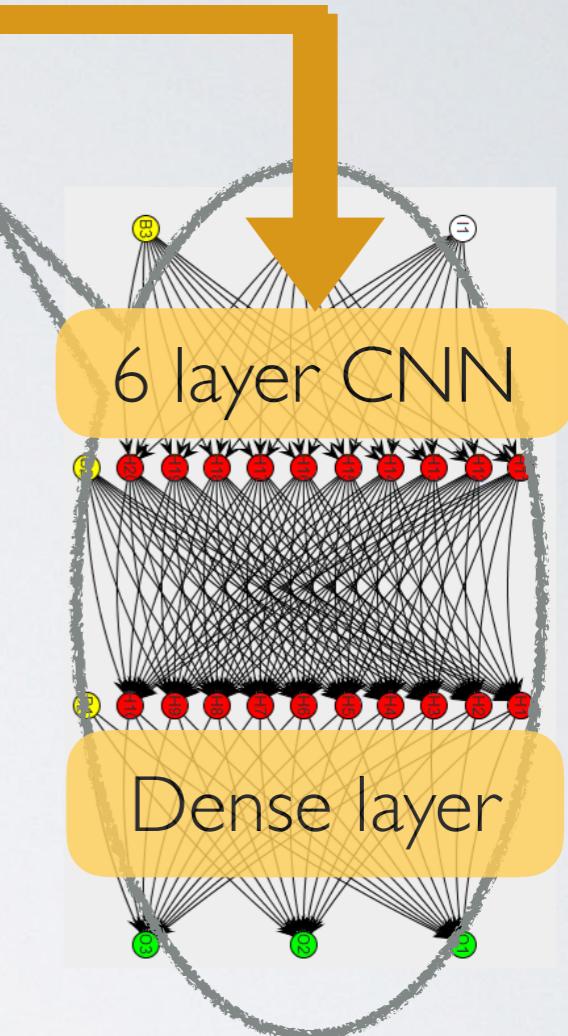


Critic



Value function

$$V(s)$$



Game Frames as states

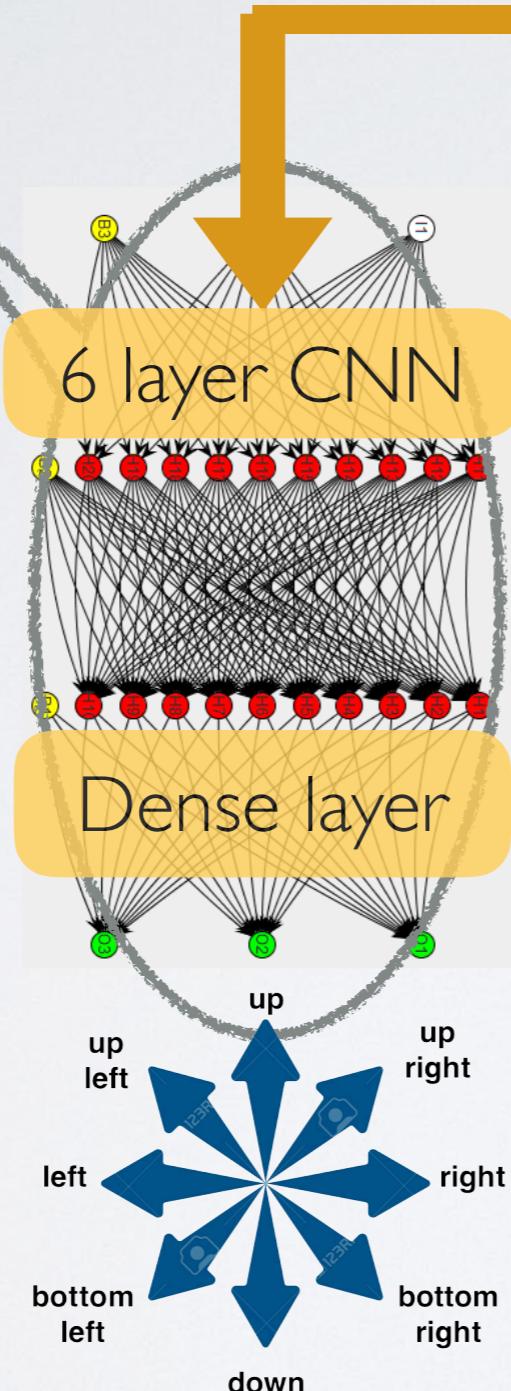


Actor



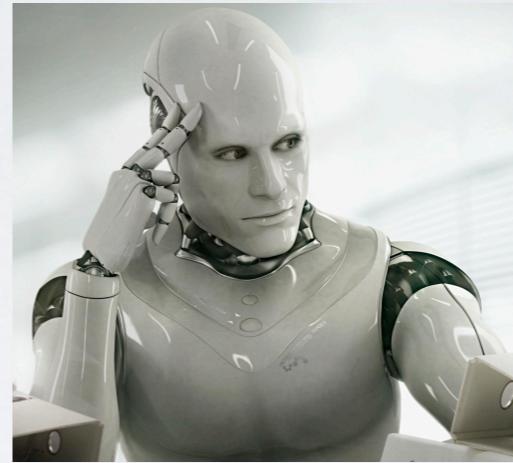
Policy function

$$\pi(a|s)$$



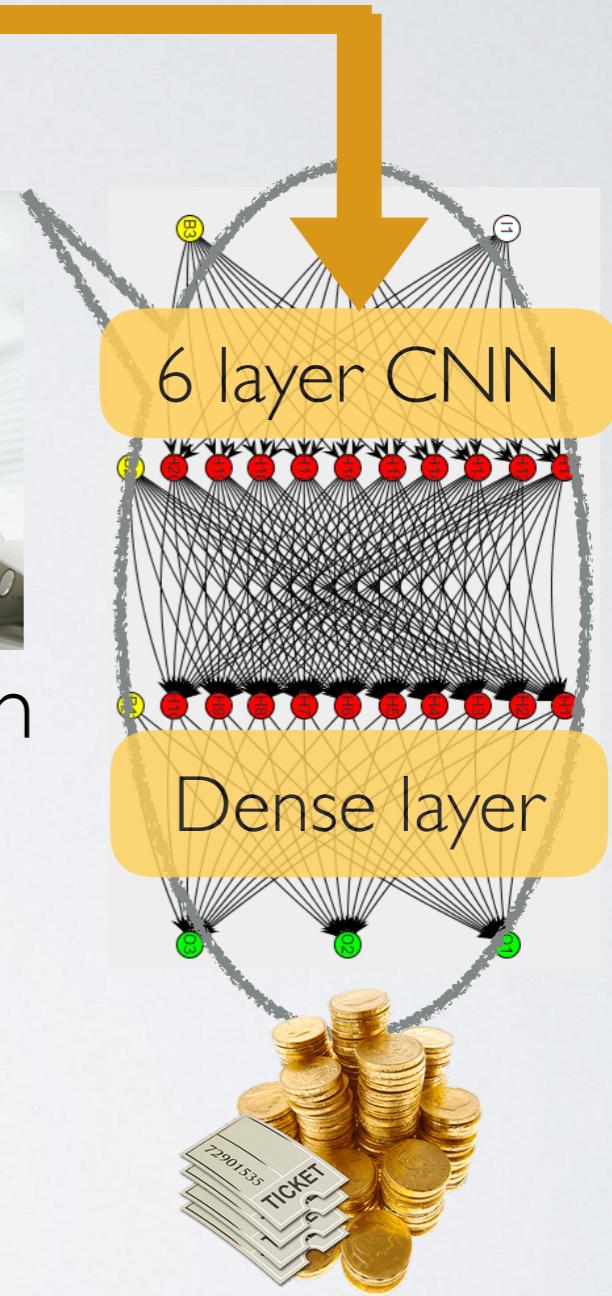
Softmax Output

Critic



Value function

$$V(s)$$



Linear Regression
Output



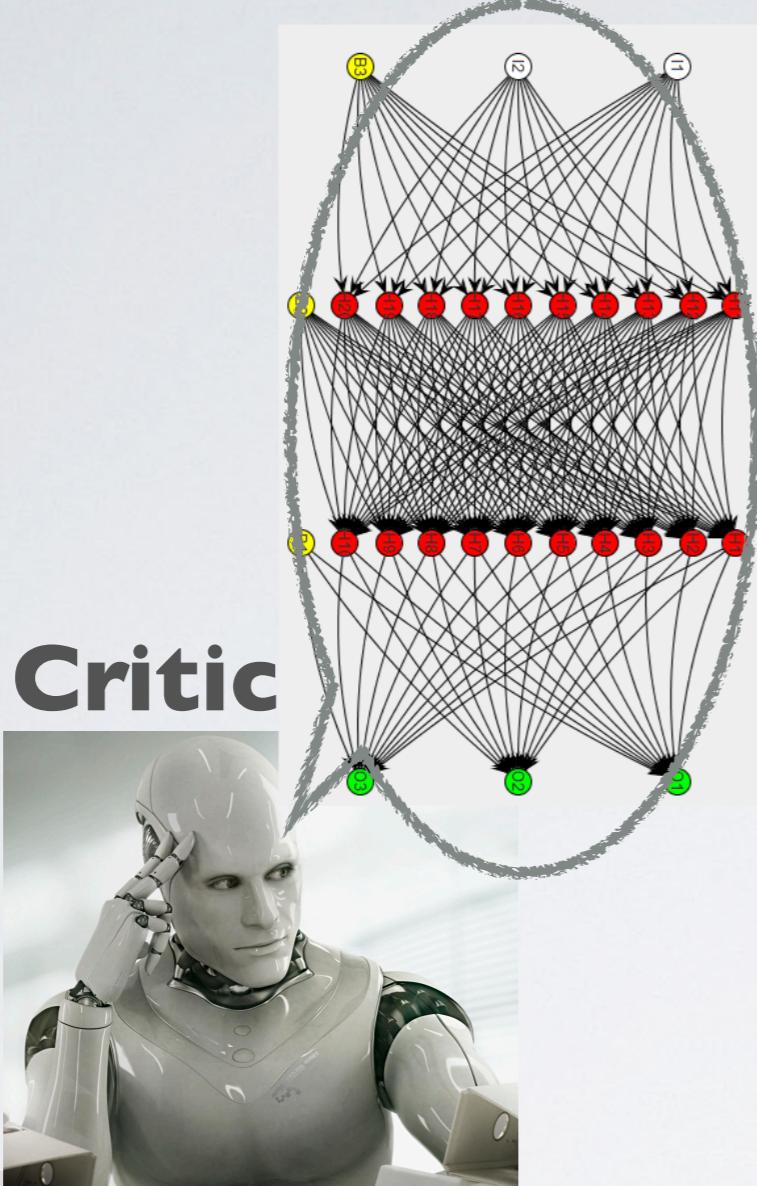
Actor

Value function

$$V(s)$$

Policy function

$$\pi(a|s)$$



Critic



Setup Game Environment for Agent to interact



Implement CNN network and Actor Critic Model



Implement Training Algorithm

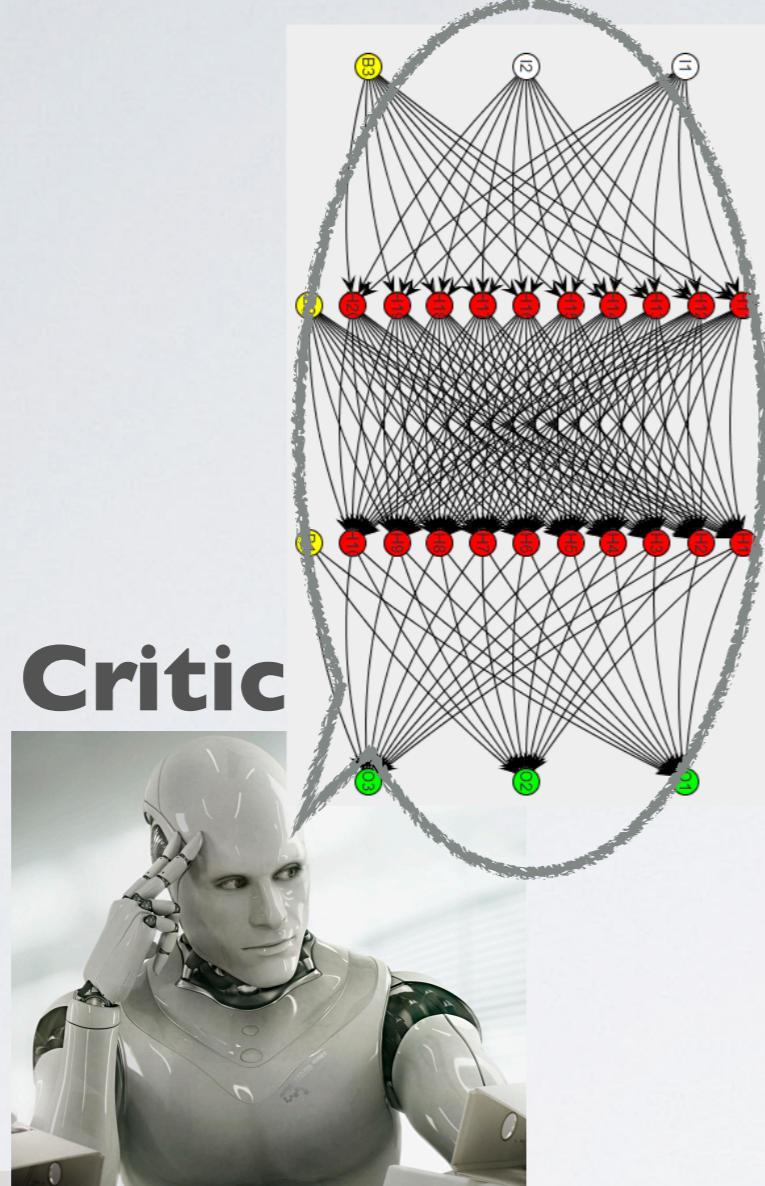


Actor

Value function

$$V(s)$$

Policy function
 $\pi(a|s)$



Critic



2D Convolution: filters=32, kernel=7*7 strides=2

2D Convolution: filters=64, kernel=7*7 strides=2

Max Pooling: pool=3*3 strides=2

2D Convolution: filters=128, kernel=3*3 strides=1

Max Pooling: pool=3*3 strides=2

2D Convolution: filters=192, kernel=3*3 strides=1

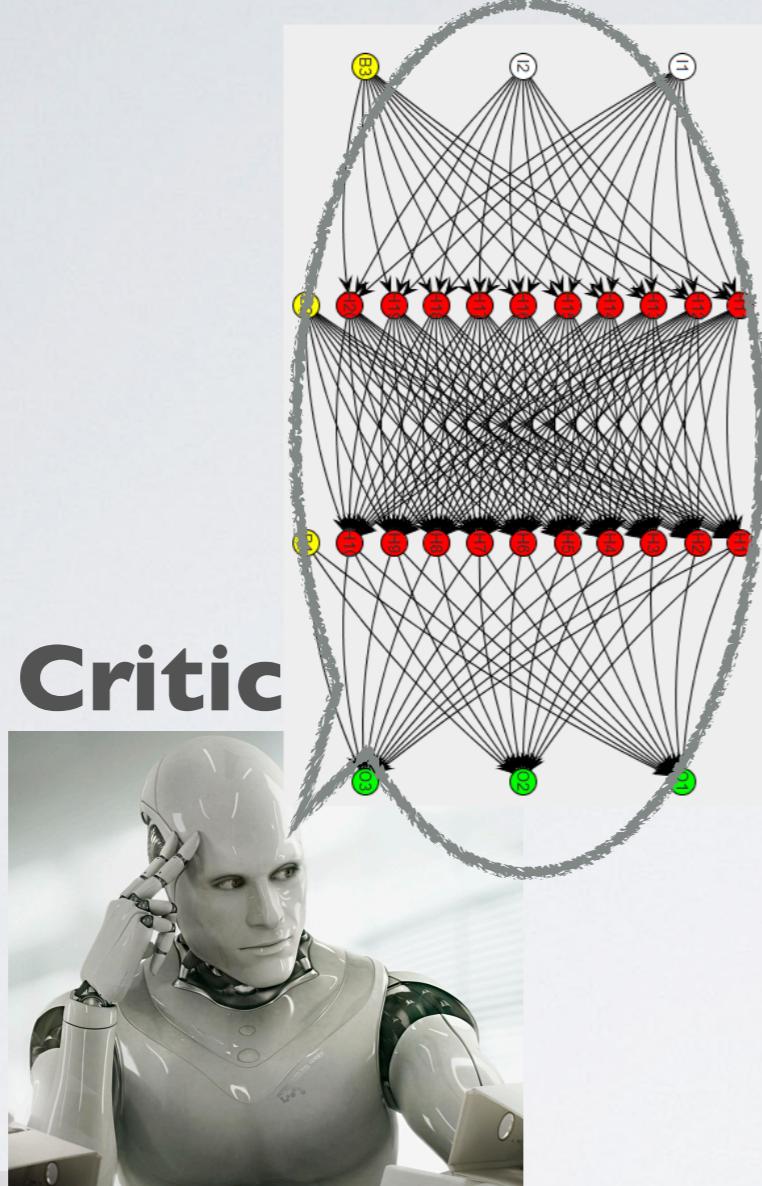
Actor



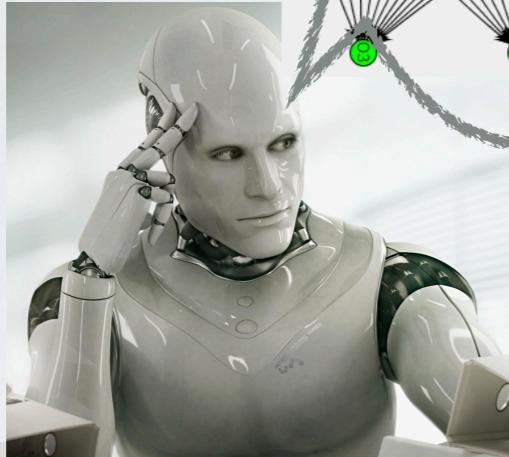
Value function

$$V(s)$$

Policy function
 $\pi(a|s)$



Critic



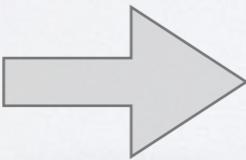
- 2D Convolution: filters=32, kernel=7*7 strides=2
- 2D Convolution: filters=64, kernel=7*7 strides=2
- Max Pooling: pool=3*3 strides=2
- 2D Convolution: filters=128, kernel=3*3 strides=1
- Max Pooling: pool=3*3 strides=2
- 2D Convolution: filters=192, kernel=3*3 strides=1
- Dense Layer: units=1024



Actor

Value function

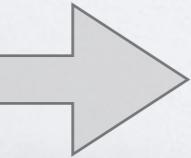
$$V(s)$$



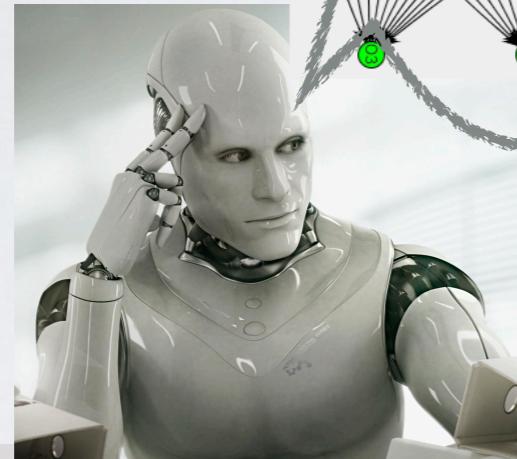
$$\text{loss}_V = (r + \gamma V(s_{t+1}) - V(s_t))^2$$

Policy function

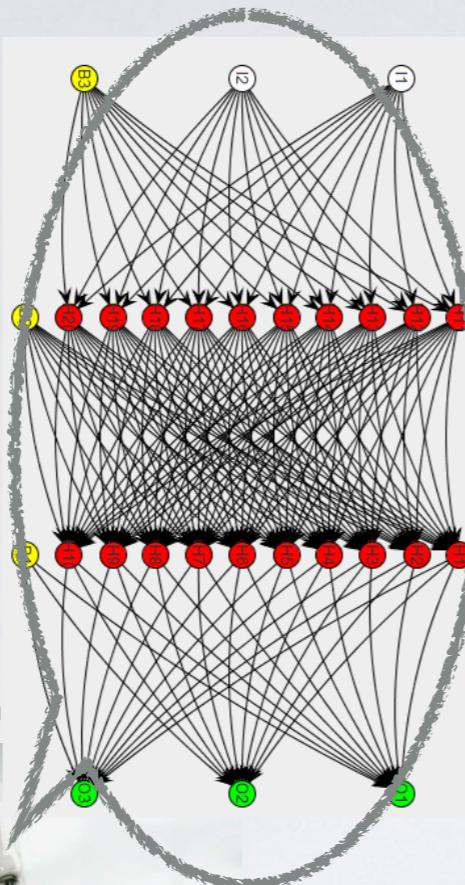
$$\pi(a|s)$$



$$\text{loss}_\pi = -(r + \gamma V(s_{t+1}) - V(s_t)) \log \pi(a|s)$$



Critic



2D Convolution: filters=32, kernel=7*7 strides=2

2D Convolution: filters=64, kernel=7*7 strides=2

Max Pooling: pool=3*3 strides=2

2D Convolution: filters=128, kernel=3*3 strides=1

Max Pooling: pool=3*3 strides=2

2D Convolution: filters=192, kernel=3*3 strides=1

Dense Layer: units=1024

Critic

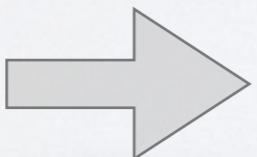


Actor



Value function

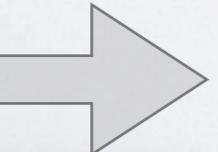
$$V(s)$$



$$\text{loss}_V = (r + \gamma V(s_{t+1}) - V(s))^2$$

Policy function

$$\pi(a|s)$$



$$\text{loss}_\pi = -(r + \gamma V(s_{t+1}) - V(s)) \ln \pi(a|s)$$

2D Convolution: filters=32, kernel=7*7 strides=2

2D Convolution: filters=64, kernel=7*7 strides=2

Max Pooling: pool=3*3 strides=2

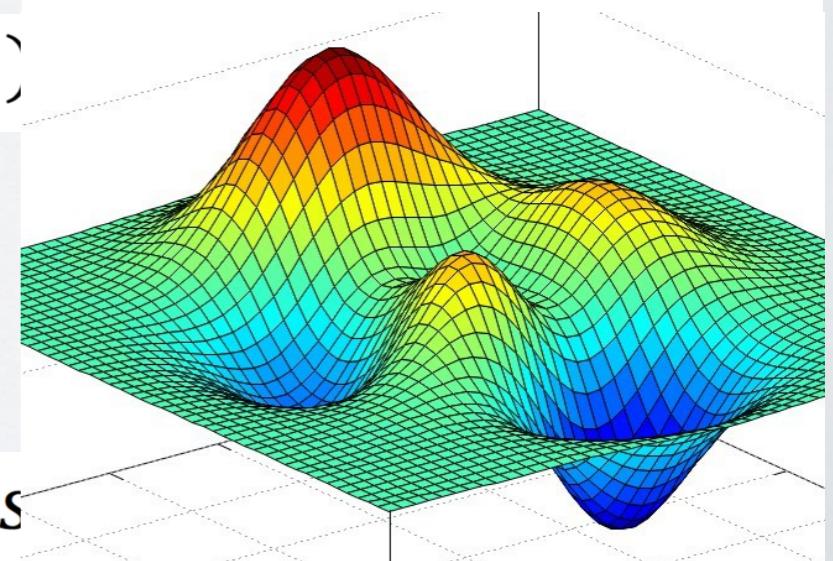
2D Convolution: filters=128, kernel=3*3 strides=1

Max Pooling: pool=3*3 strides=2

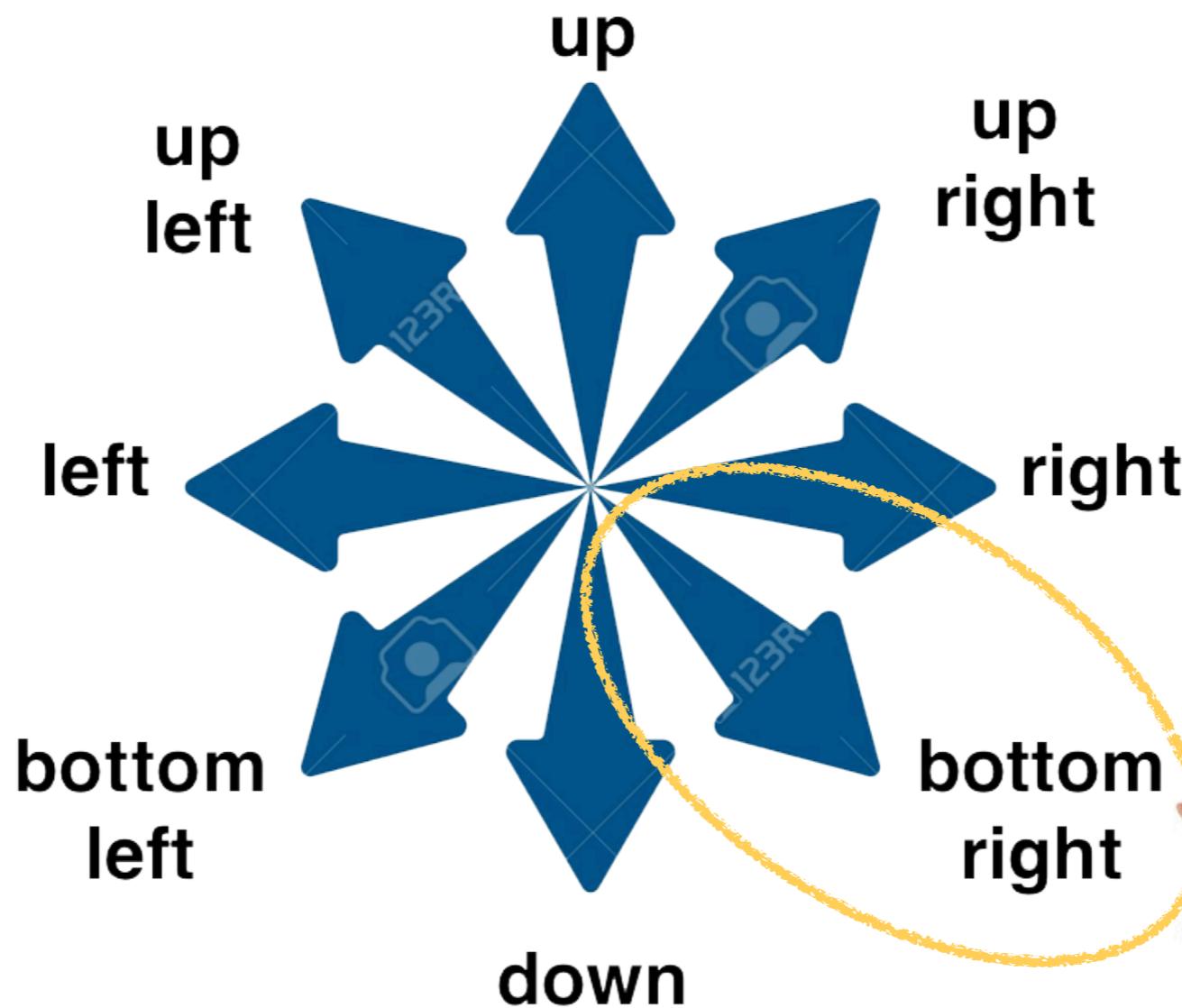
2D Convolution: filters=192, kernel=3*3 strides=1

Dense Layer: units=1024

Optimize with
Gradient Descent



HIGHLY PEAKED POLICY FUNCTION DILEMMA



The same action was repeatedly chosen, our Agent move only walk in a straight line.

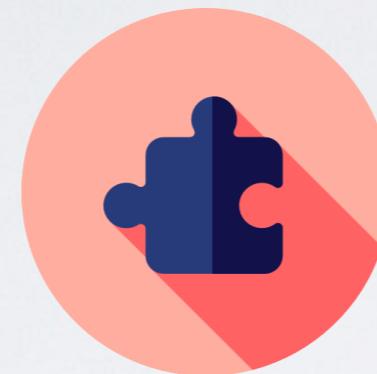
THE DILEMMA



Early convergence discourages exploration, if an action is over reinforced too early, we may be missing something that made another action better

EXPLORATION TUNING

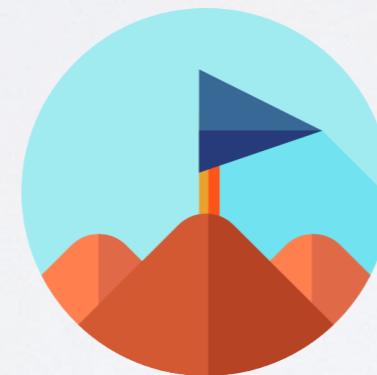
Add an Entropy Term



Adjust Action Policy



Reward Shaping

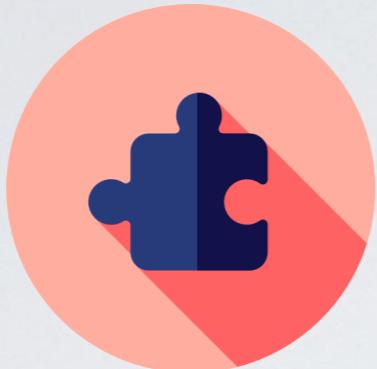


Add Exploration Factor

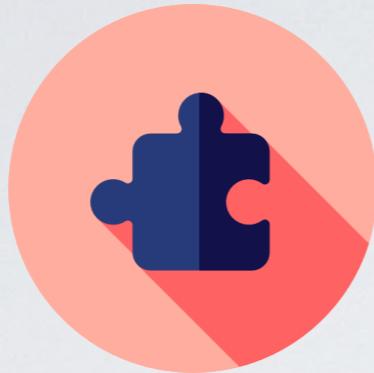


Encourage exploration! Keep the search alive!
Prevent the agent from converging to a single action

Add an
Entropy Term



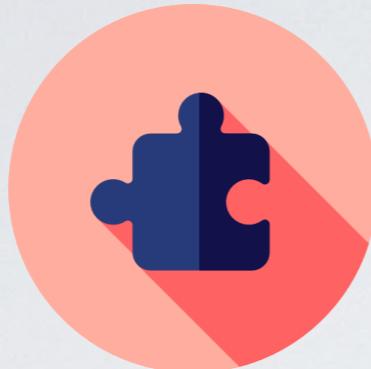
Add an
Entropy Term



Original Policy Network Loss Function:

$$loss_{\pi} = -[(r + \gamma V(s_{t+1}) - V(s_t)) \log \pi(a|s)]$$

Add an
Entropy Term



Uniform



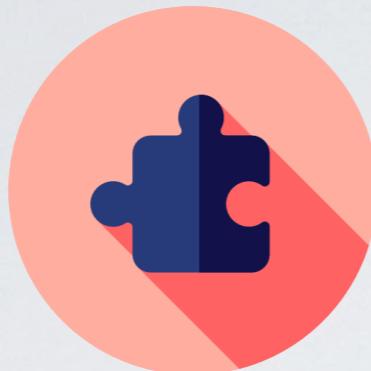
Highly peaked

$$H(\pi(\cdot|s)) = -\text{Sum}[\pi(\cdot|s) \log(\pi(\cdot|s) + \rho)]$$

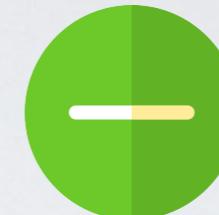
Original Policy Network Loss Function:

$$loss_{\pi} = -[(r + \gamma V(s_{t+1}) - V(s_t)) \log \pi(a|s)]$$

Add an
Entropy Term



Uniform



Highly peaked

$$H(\pi(\cdot | s)) = -\text{Sum}[\pi(\cdot | s) \log(\pi(\cdot | s) + \rho)]$$



$$\text{loss}_\pi = -[(r + \gamma V(s_{t+1}) - V(s_t)) \log \pi(a|s) + \beta H(\pi(\cdot | s))]$$

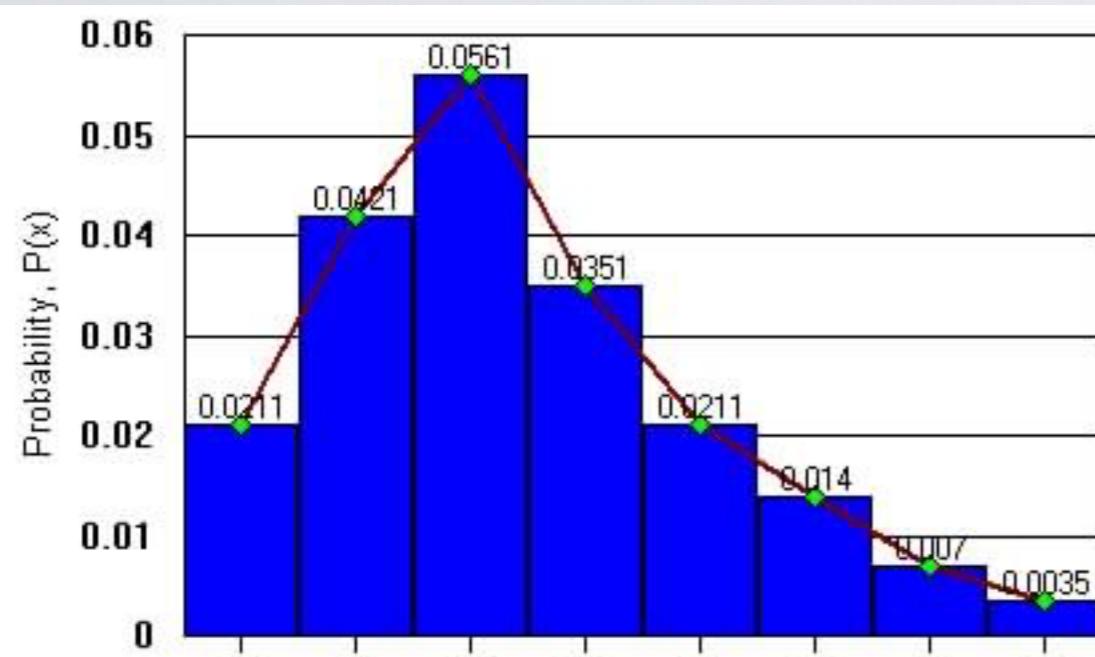
maximizing the entropy term



encouraging uniform action probability

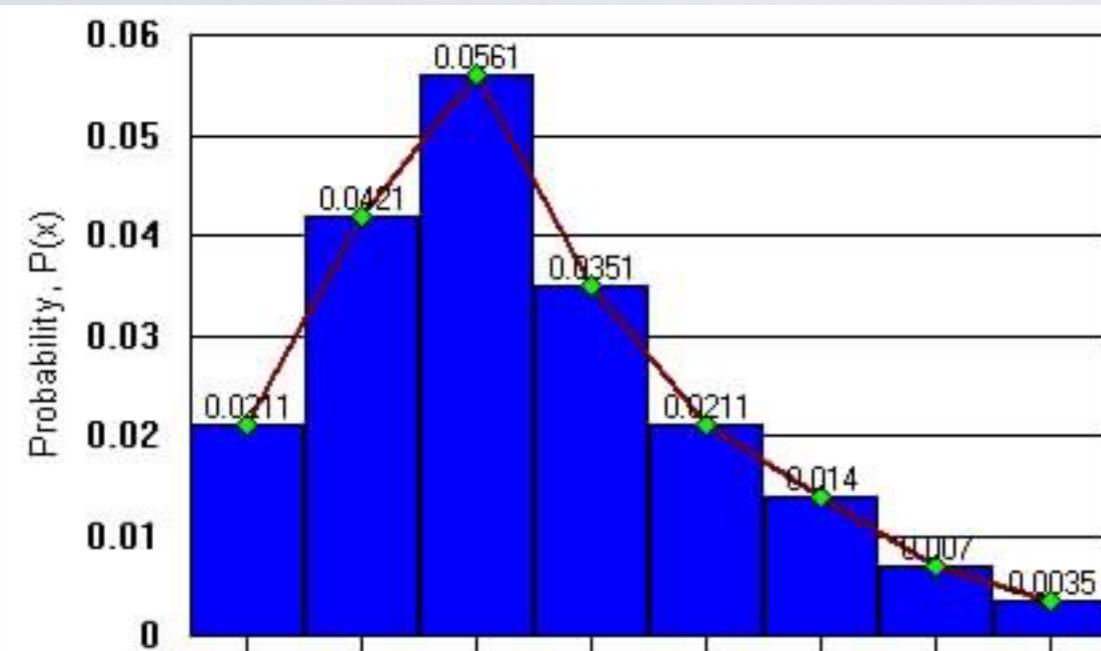


Adjust
Action Policy



Adjust
Action Policy

chosen randomly with
respect to probability
distribution



Adjust Action Policy

chosen randomly with respect to probability distribution

Activate **Epsilon Greedy Policy During** Training

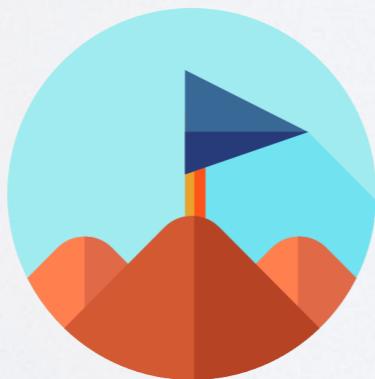
%



20% Random Explore

80% Greedy Exploit

Reward
Shaping



in addition to basic rewards, we apply immediate rewards:

Condition	Reward, r
$r > 0$	$r \times 2$
repeated action > 3	$r - = \text{penalty}$, $\text{penalty} += 0.25/\text{repeated action}$
action \neq action'	$r + = 1$
die	$r = -10$

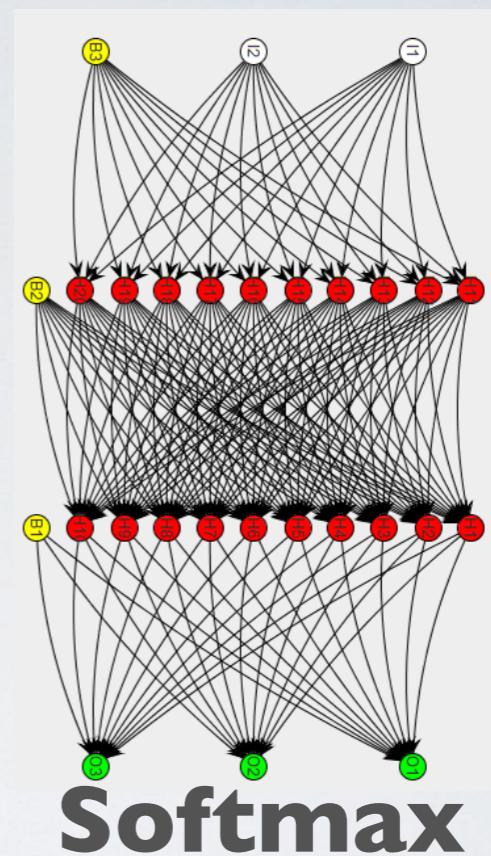
Effective when rewards
are sparse and delayed

Reward
Shaping





Add Exploration
Factor

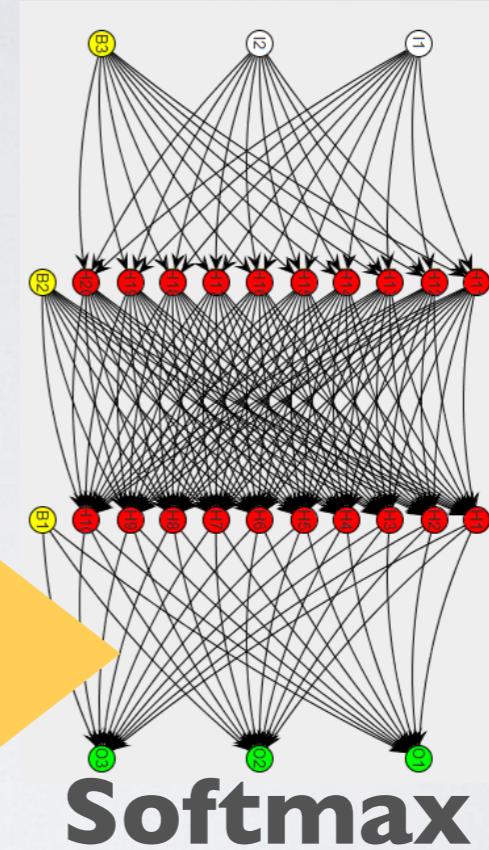


Add Exploration
Factor

0.2



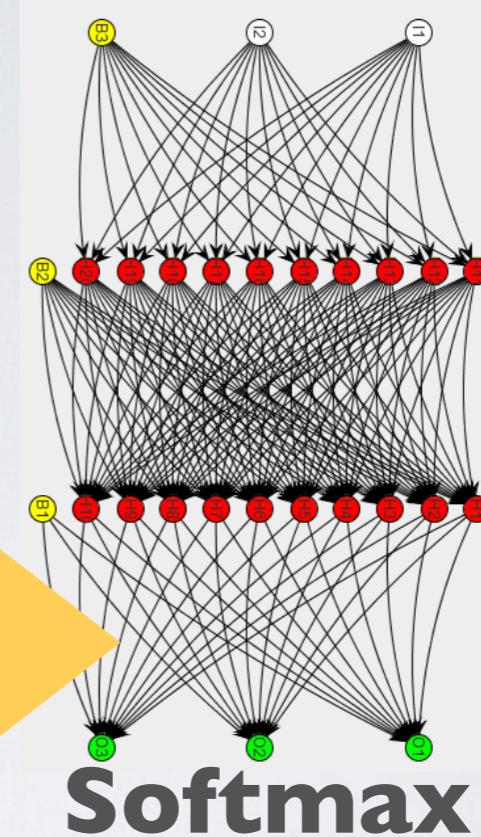
Exploration factor



Add Exploration
Factor

0.2 ✗

Exploration factor



Softmax

- scales down the affect of the highly-peaked policy function
- makes the output of the softmax more uniform



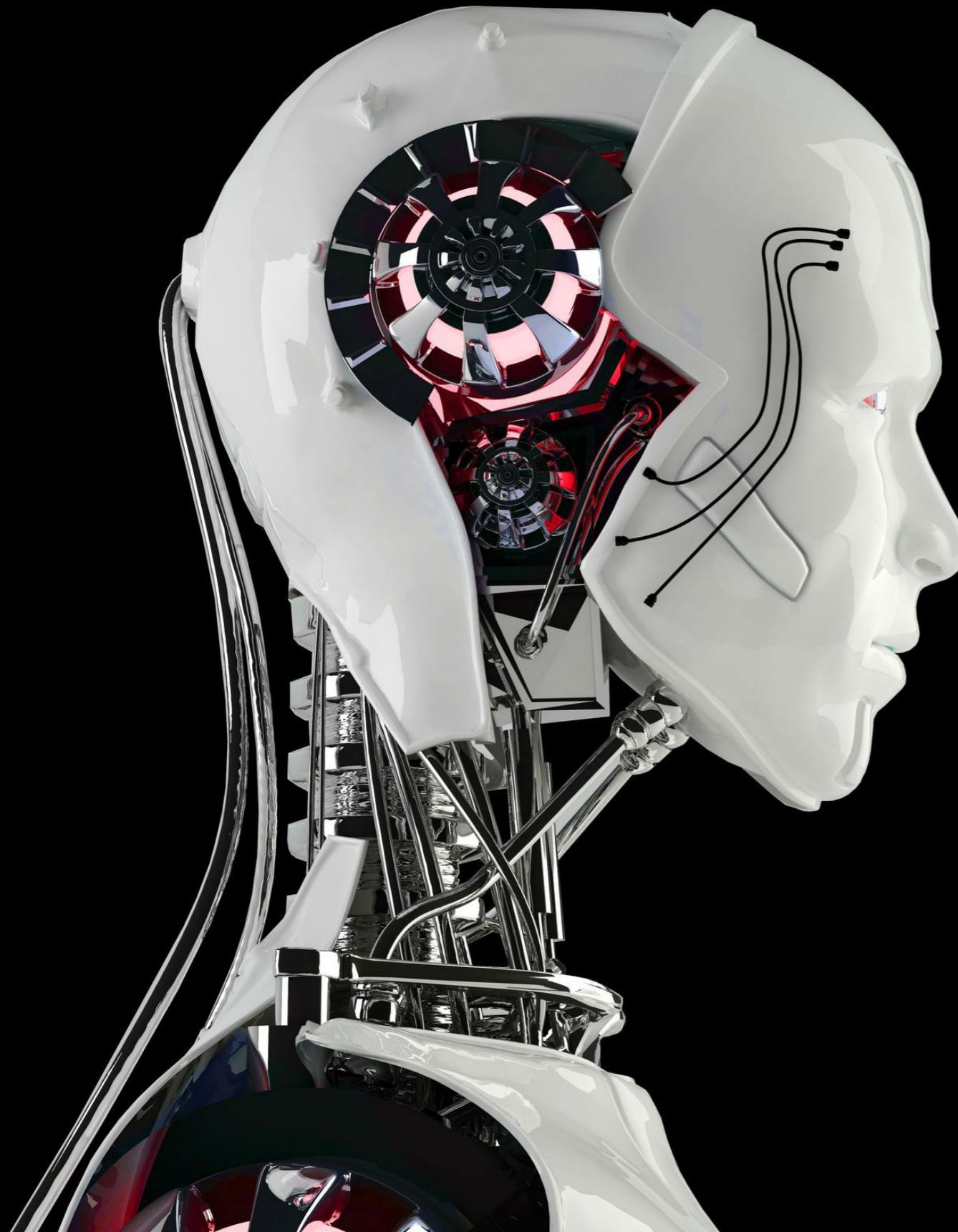
Add Exploration Factor



Result

Well Trained Proficient Agent

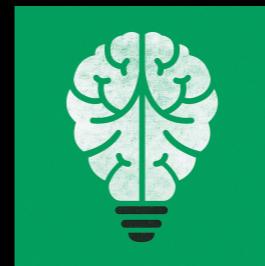
Gif and video of the agents gameplay available



Result

Well Trained Proficient Agent

Gif and video of the agents gameplay available



Applications

- RL used in NLP for training QA models
- One Day, agents will be able to play the largest game in the world - life!

Normal turn



Sharp Turn



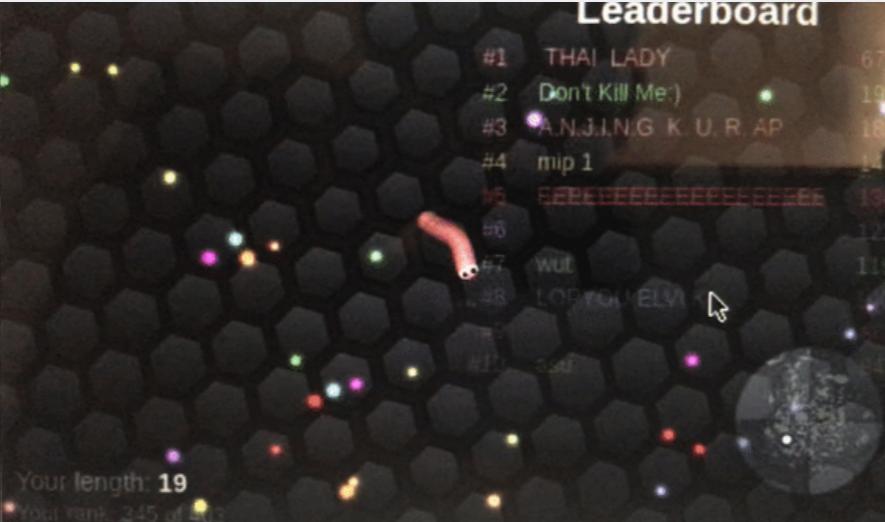
Hard Twist



Snake Follow



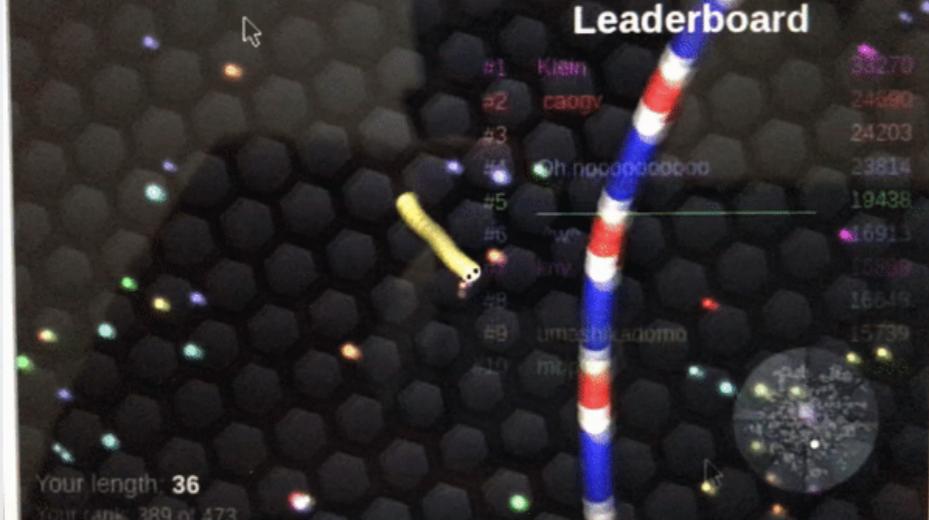
Kill Enemies



Bypass Enemies

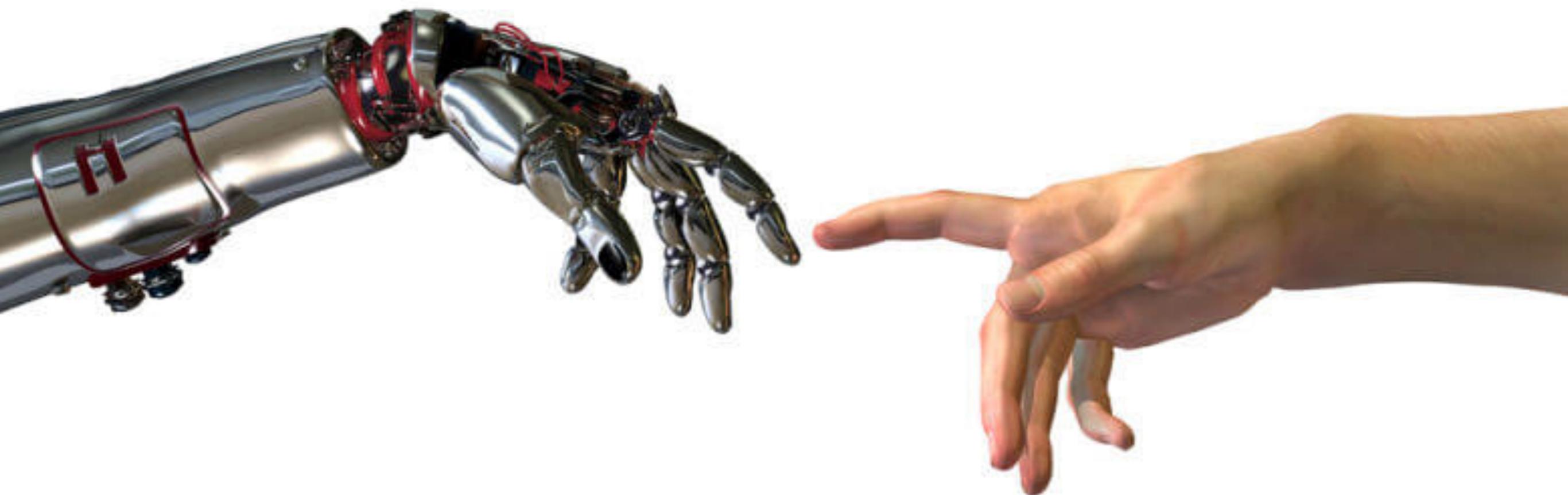


Survive in Chaos



Circular Movement





THANK YOU
for your attention!