



INTERNATIONAL CENTER  
FOR COMPUTATIONAL LOGIC

Technische Universität Dresden faculty of computer science, ICCL

# USING REINFORCEMENT LEARNING TO PLAY ANGRY BIRDS

colloquium

Peter Hirsch

Dresden, 2017/9/26



DRESDEN  
concept  
Erstellen von  
Wissenschaft  
und Kultur

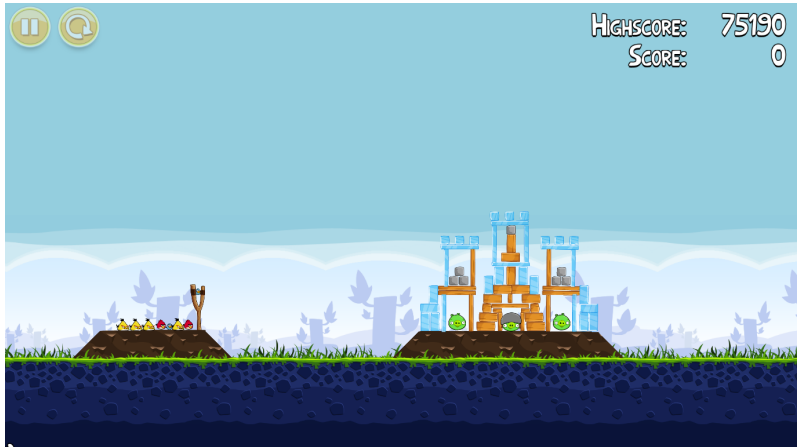


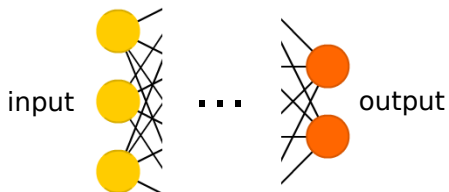
- 1 Angry Birds
- 2 Neural Networks
- 3 Reinforcement Learning
- 4 Deep Deterministic Policy Gradient
- 5 RL Agent: Dr. L. Bird
- 6 Results

# Angry Birds in a Nutshell

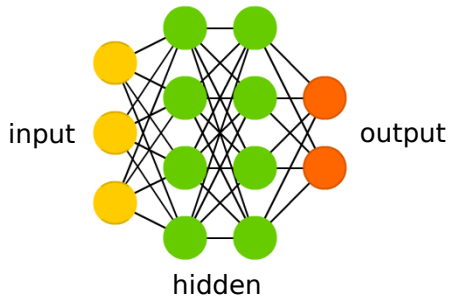


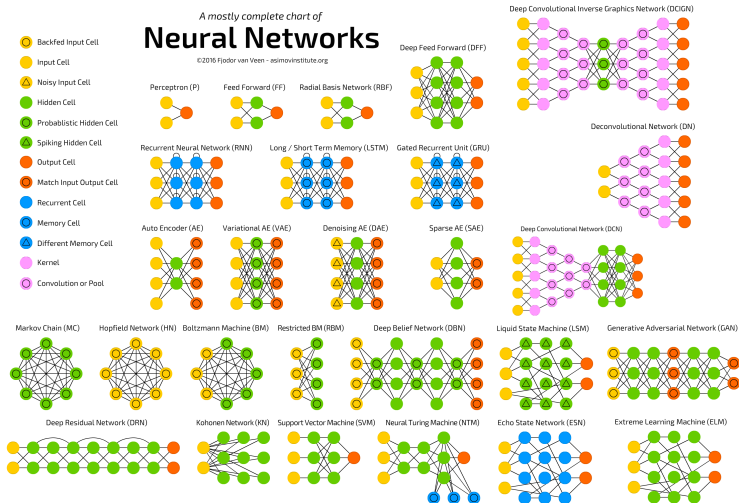
INTERNATIONAL CENTER  
FOR COMPUTATIONAL LOGIC

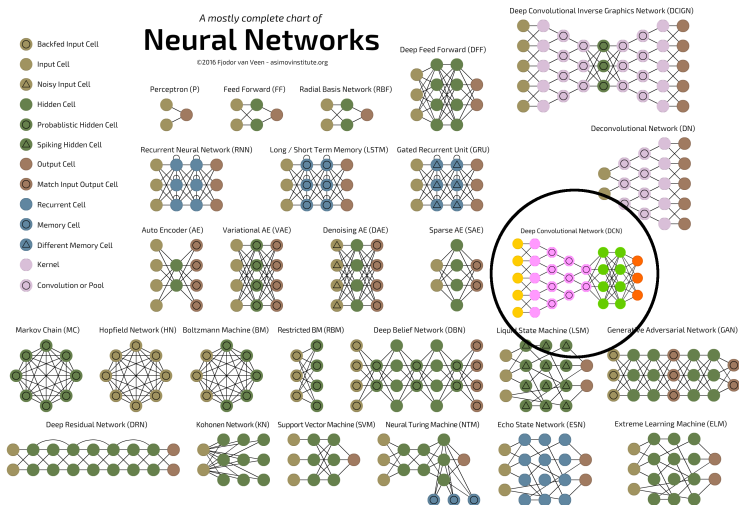


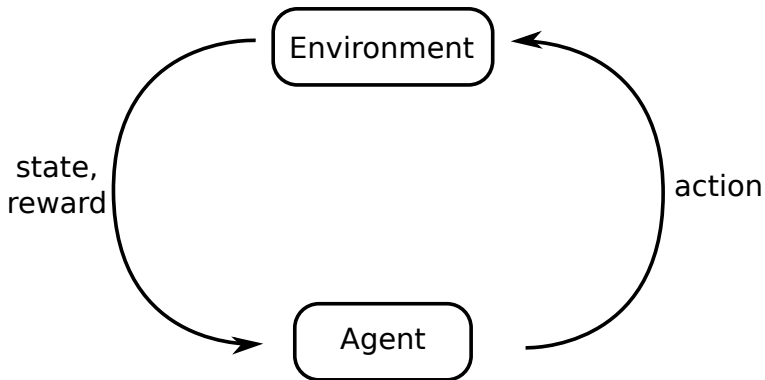


$$\text{output} = \sum \text{weights} \cdot \text{inputs}$$

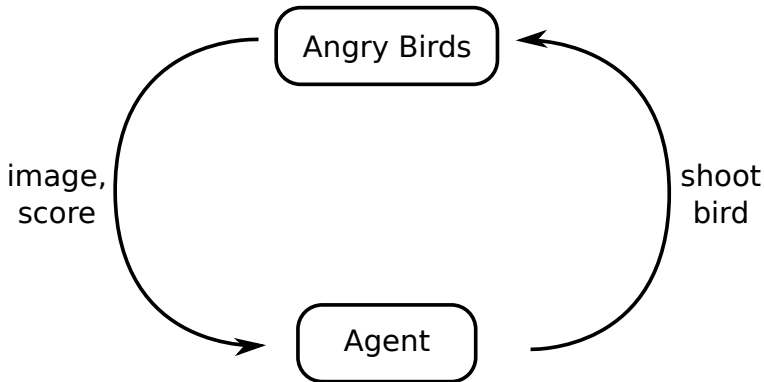


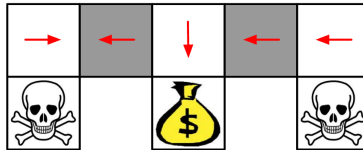




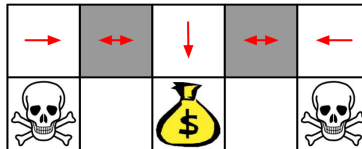






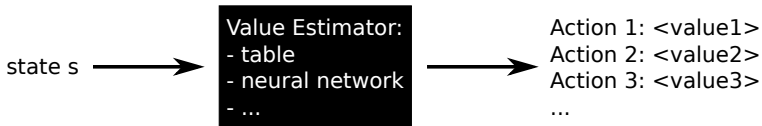


deterministic policy

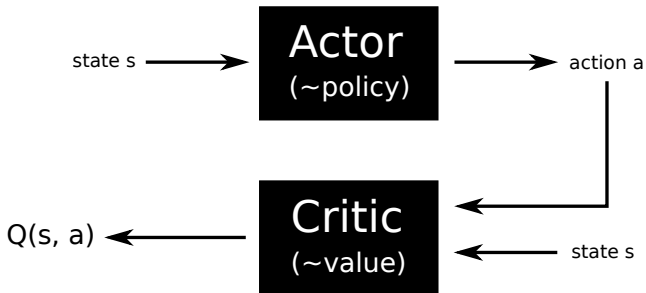


stochastic policy

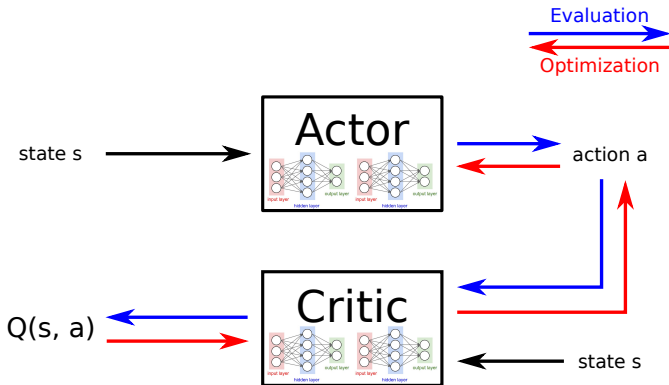
(with state approximation, gray states not distinguishable)

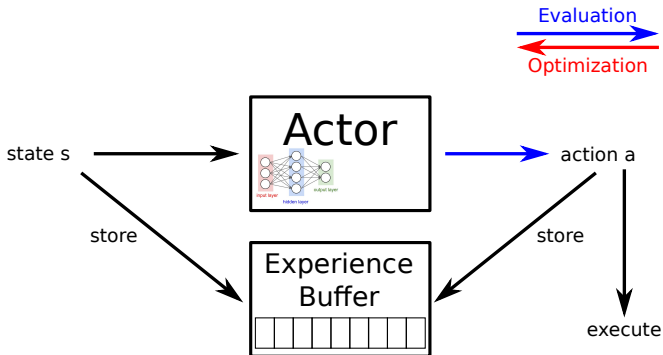


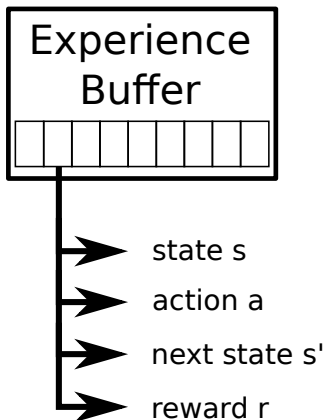
a **greedy** or  $\epsilon$ -**greedy** policy is used to act  
(a.k.a. go to neighboring state with highest (Q-)value)



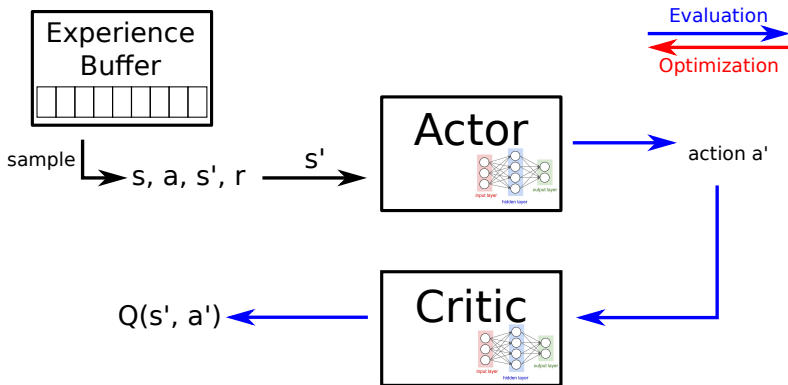
Combination of policy-based and value-based





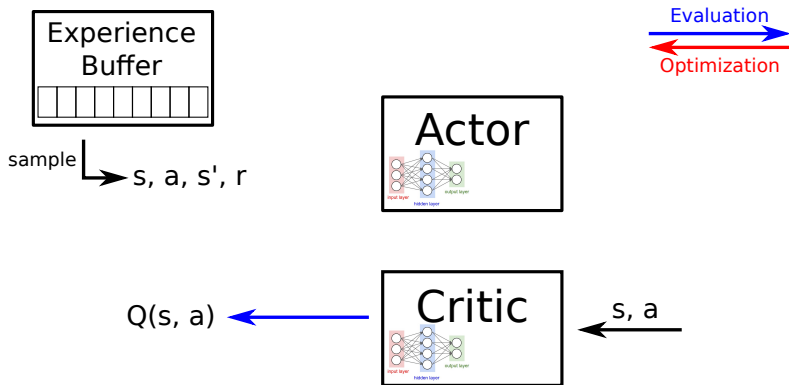


# DDPG: Learning - Part 1

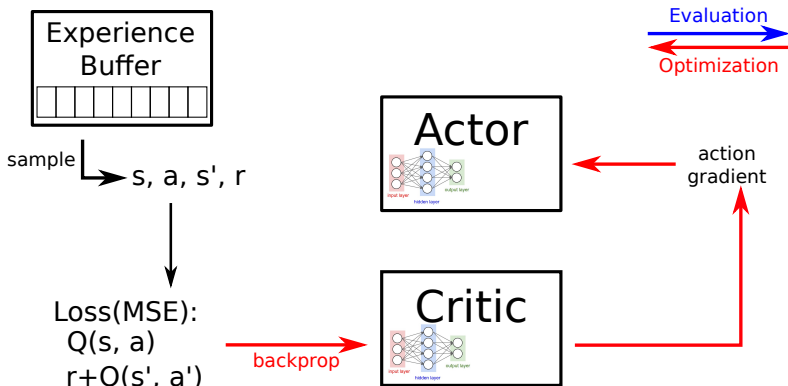




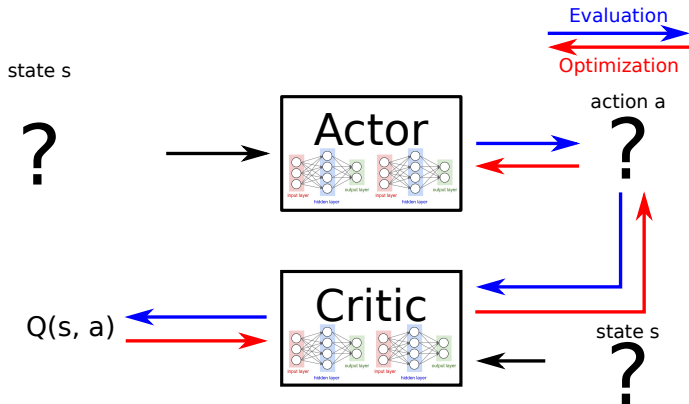
# DDPG: Learning - Part 2

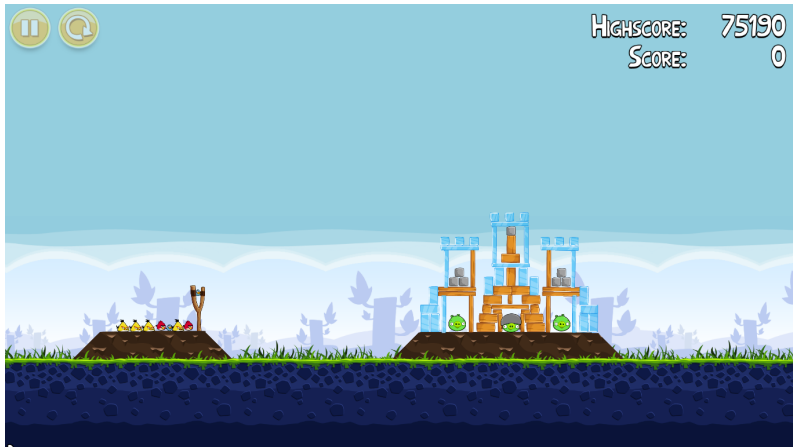


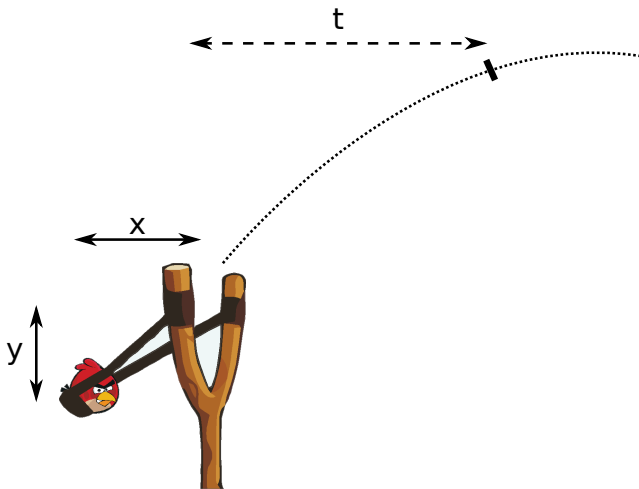
# DDPG: Learning - Part 3

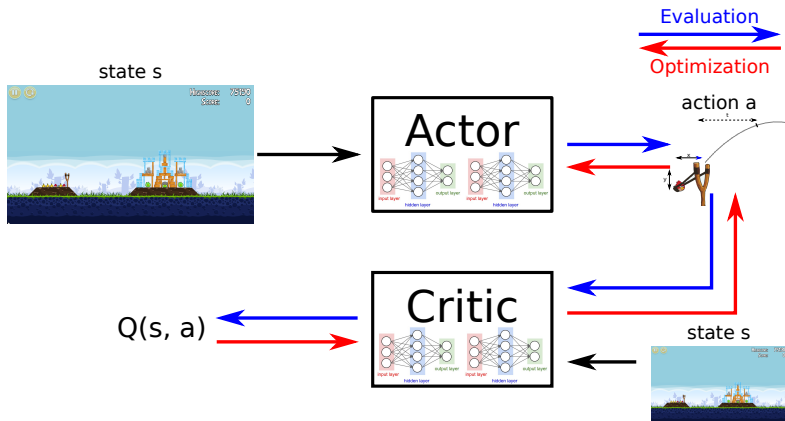


Optimization: Backpropagation using chain rule across the two networks

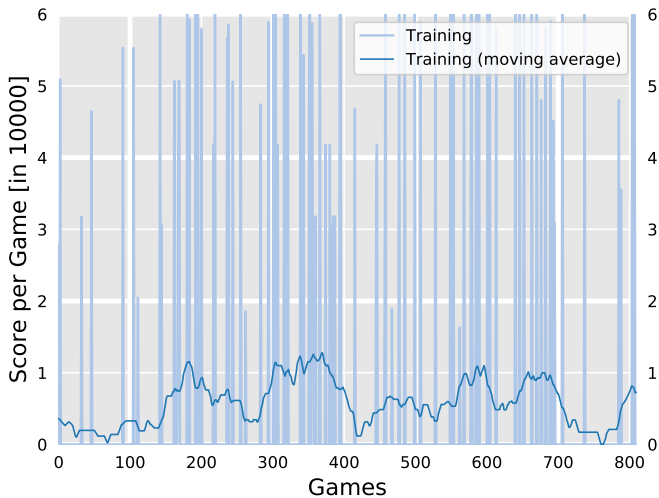




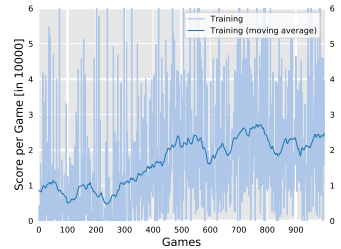
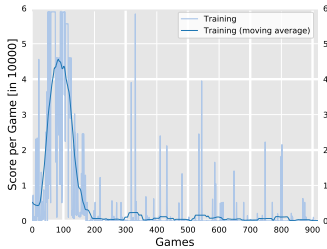
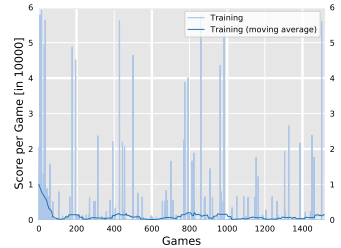
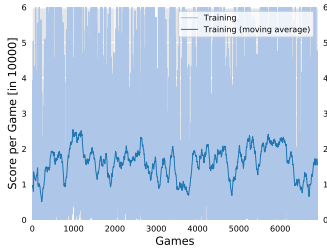




# Results

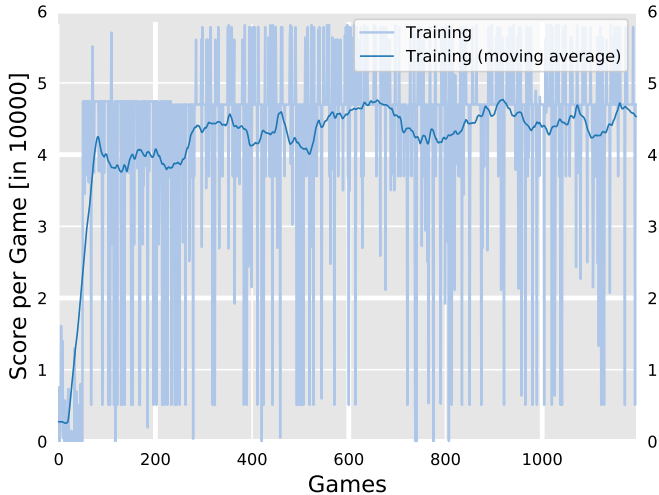


# Results





# One Good Run





- DDPG loss function:
  - Q-Learning-based (off-policy)
  - Sarsa-based (on-policy)
  - TD-based (estimated value)
  - Monte-Carlo-based (cumulative return)
- Stochastic Policy Gradient:
  - policy-based
  - uses statistics of probability distribution
  - output (sampled for action):
    - mean
    - variance
- A3C (asynchronous advantage actor-critic)
  - actor-critic
  - stochastic policy
  - parallel (asynchronous) execution of multiple agents
  - advantage instead of Q-value (relative value of actions)

→ no success so far



Sources neural network schematics:

- <http://www.asimovinstitute.org/neural-network-zoo/>
- <http://cs231n.github.io/neural-networks-1/>

Source policy-based method example:

<http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>



# Discussion