

# *DQNViz: A Visual Analytics Approach to Understand Deep Q-Networks*

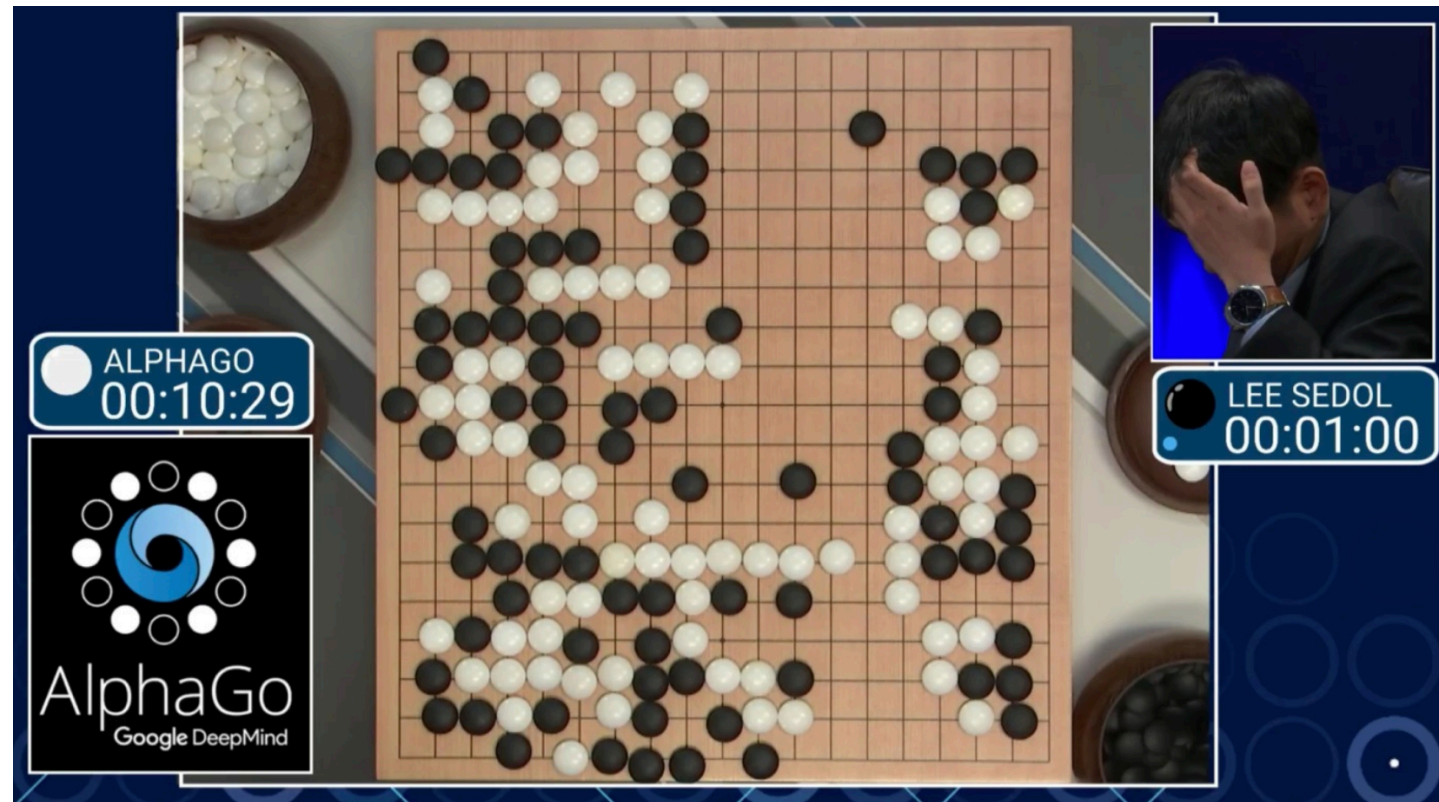
Junpeng Wang<sup>1</sup>, **Liang Gou**<sup>2</sup>, Han-Wei Shen<sup>1</sup>, Hao Yang<sup>2</sup>

1. The Ohio State University

2. Visa Research

# Introduction

- Deep Reinforcement Learning + AlphaGo

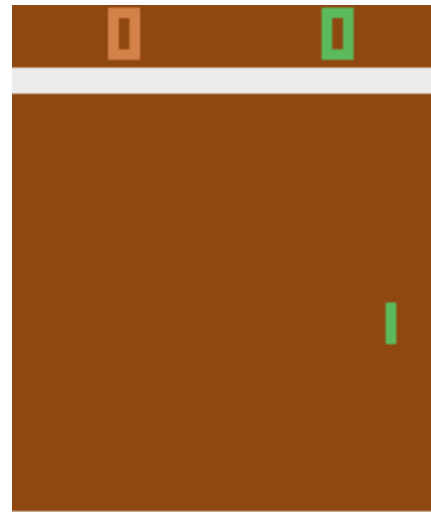


# Introduction

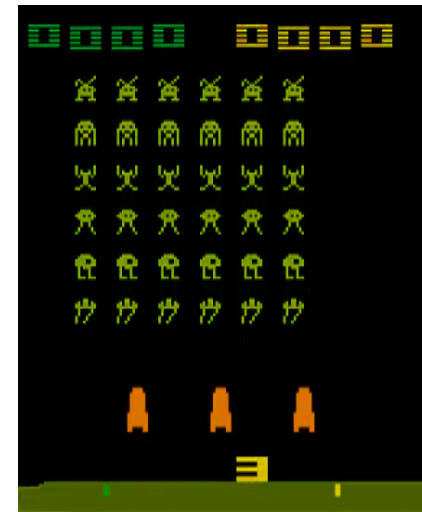
- Deep Reinforcement Learning + AlphaGo
- Deep Q Networks + Atari Games



Breakout



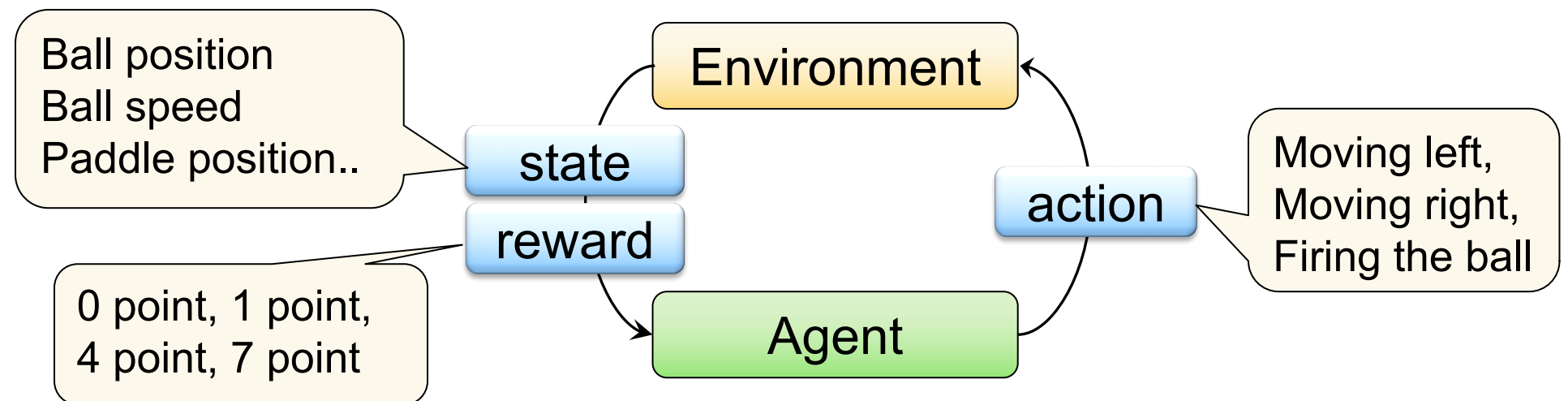
Pong



Space-Invader

# Introduction

- Deep Reinforcement Learning + AlphaGo
- Deep Q Networks + Atari Games



$$a_1, s_1, r_1, a_2, s_2, r_2, \dots, a_n, s_n, r_n$$

# Background: Q-Learning with Bellman Equation

Data:  $s_0, a_0, r_1, s_1, a_1, r_2, \dots, r_n, s_n$

Total reward:  $R = r_1 + r_2 + \dots + r_n$

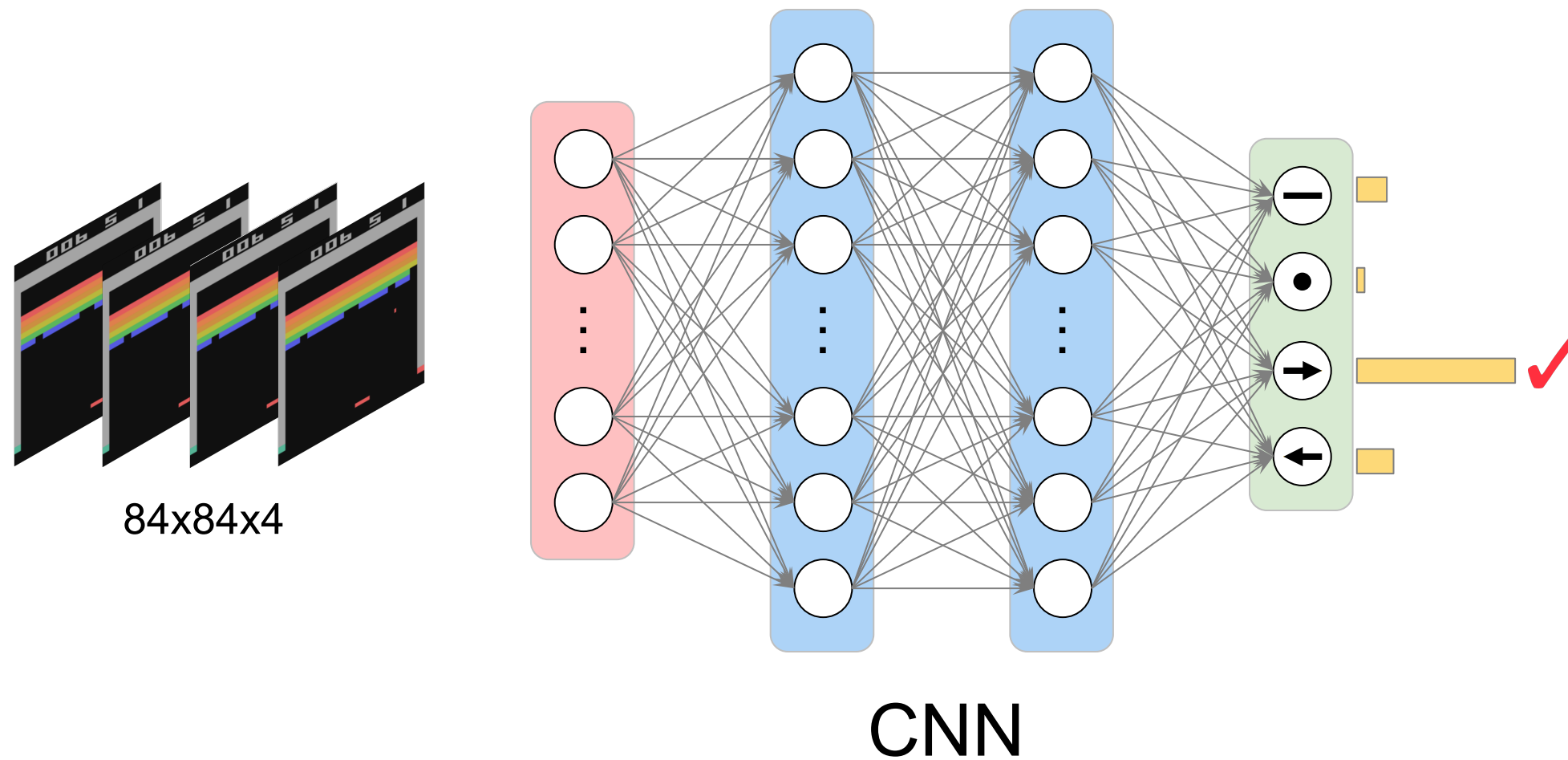
Future reward:  $R_t = r_t + r_{t+1} + \dots + r_n$

Discounted Future reward:  $R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} \dots + \gamma^{n-t} r_n$

Bellman equation:  $Q(s, a) = r + \gamma \max_{a'} Q(s', a')$

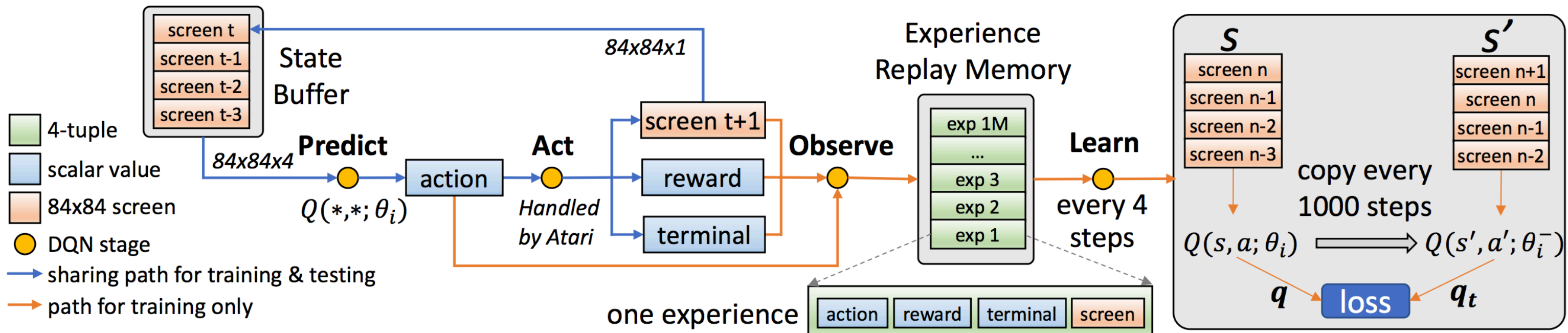
Why deep neural network?  $s$  is too complex:  $256^{84 \times 84 \times 4}$

# Background: Deep Q(uality)-Network (DQN)

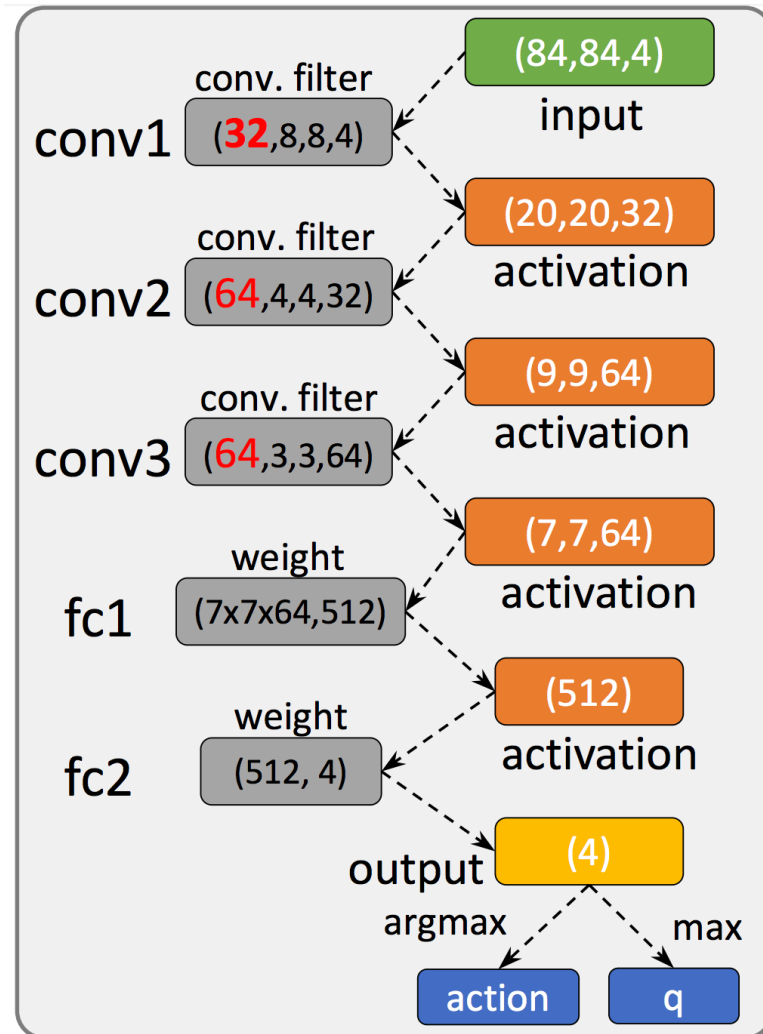


**Objective:**  
maximize the  
total game  
reward

# Background: Deep Q(uality)-Network (DQN)



# Background: Deep Q(uality)-Network (DQN)



$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim ER} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right]$$

a random sample from the experience replay memory:  $(s, a, r, s') \sim ER$

$s = \{sn_{t-3}, sn_{t-2}, sn_{t-1}, sn_t\} \rightarrow DQN(\theta_i) \rightarrow \begin{cases} output_t = \{1.91, 1.72, \mathbf{1.95}, 1.83\} \\ q = \max(output_t) = 1.95 \\ a = \operatorname{argmax}(output_t) = 2 \text{ (right)} \end{cases}$

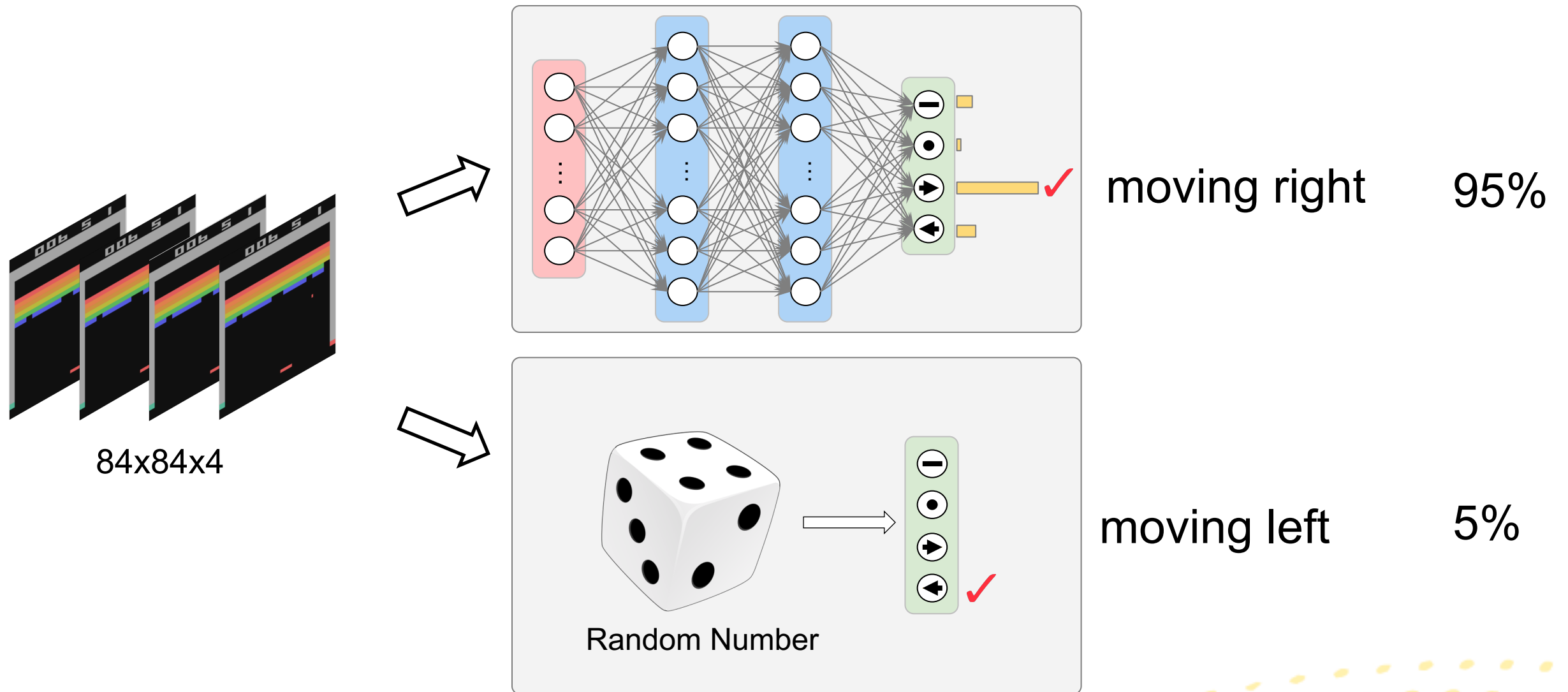
$r = 1$        $\gamma = 0.99$

$s' = \{sn_{t-2}, sn_{t-1}, sn_t, sn_{t+1}\} \rightarrow DQN(\theta_i^-) \rightarrow \begin{cases} output_{t+1} = \{1.35, 1.39, 1.32, \mathbf{1.43}\} \\ q_t = r + \gamma * \max(output_{t+1}) \\ = 1 + 0.99 * 1.43 = 2.4175 \end{cases}$

$$loss = (q - q_t)^2 = q_{diff}^2 = (1.95 - 2.4175)^2 = 0.4675^2 = 0.2186$$



# Background: Exploration and Exploitation dilemma



# Challenges

- Long-time blind training process (understand the model)
  - What strategies are really learned?
  - When are those strategy learned?
  - Which part of the neural network learned those strategies?
- Proper choice of different hyper-parameters (improve the model)
  - E.g., the random rate for the tradeoff between exploration and exploitation

# Contribution: DQNViz

- Visual Analytics System (DQNViz):
  - Effective visual summary
  - Efficient (movement/reward) pattern mining
- Improve DQN training by optimizing the random actions
  - Pattern detection algorithm based on DQNViz



Part I: effective visual summary  
and efficient pattern mining

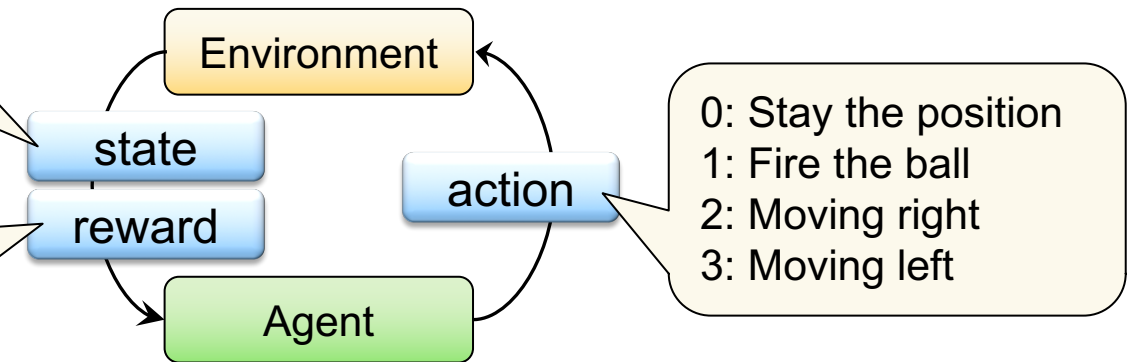
# DQNViz: The Breakout Game



7 points  
4 points  
1 point

Position and speed of the ball, the paddle, state of the bricks, etc. → 4 consecutive screens

0: did not hit anything  
1: hit bottom two rows of bricks  
4: hit middle two rows of bricks  
7: hit top two rows of bricks



$a_1, s_1, r_1, a_2, s_2, r_2, \dots, a_n, s_n, r_n$

Example Data:

Action sequence: 0000012222333312312311000

State sequence: (84x84x3)(84x84x3)...(84x84x3)

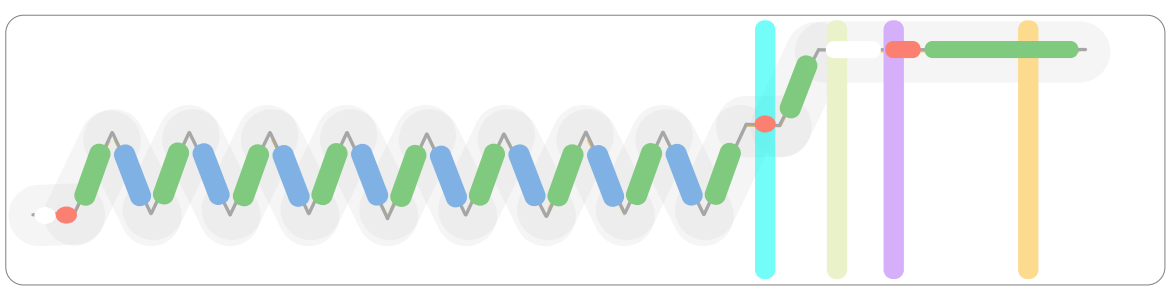
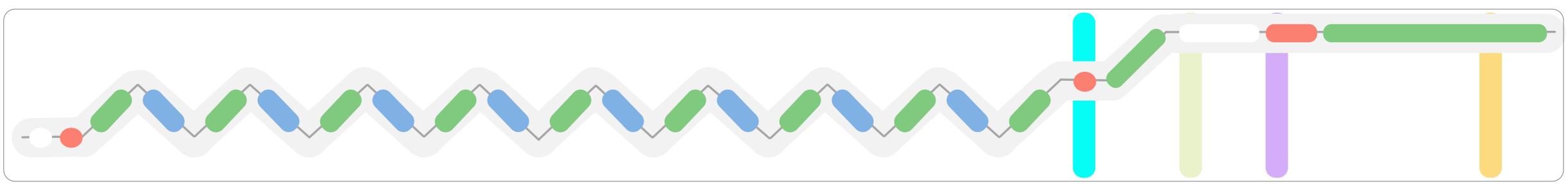
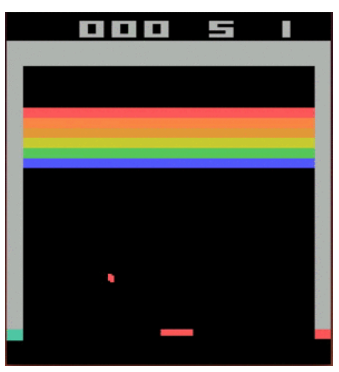
Reward sequence: 100000004000000400000700

game start

game over



Actions: 0133223322332233223322332233223322332233223322331333000113333333333

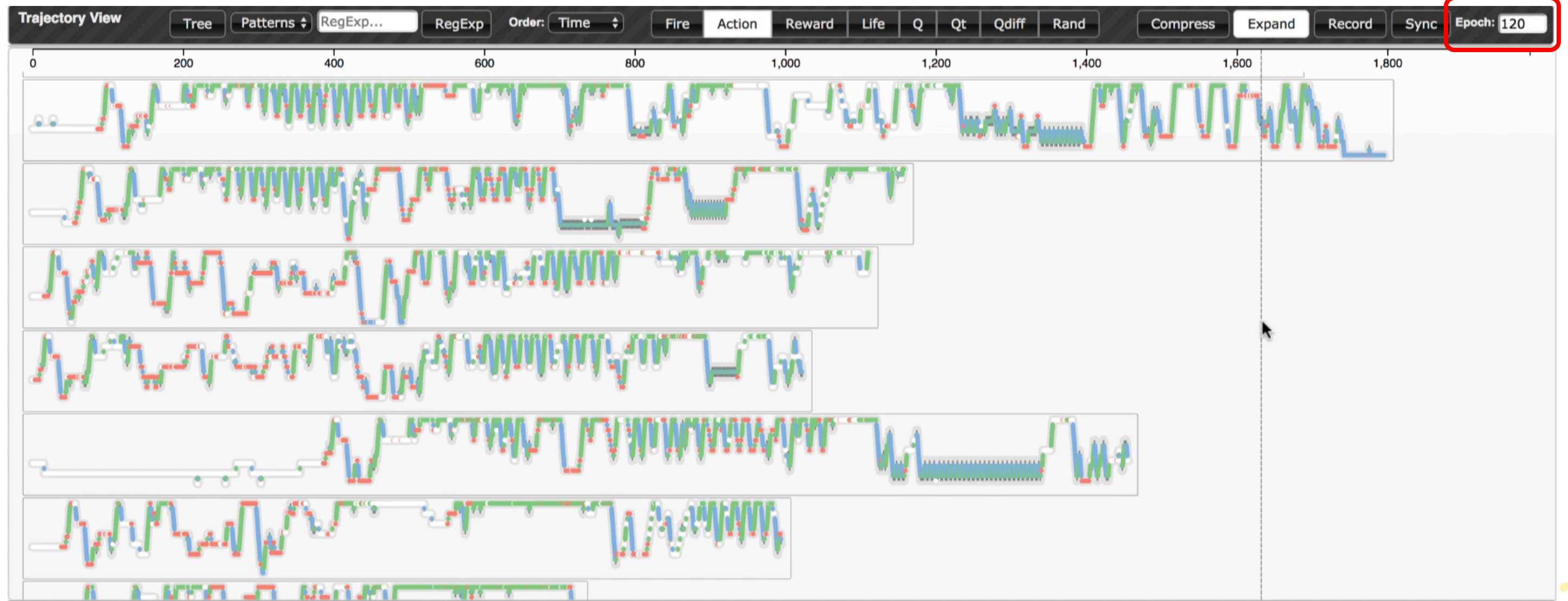


random action    1 point    4 points    7 points

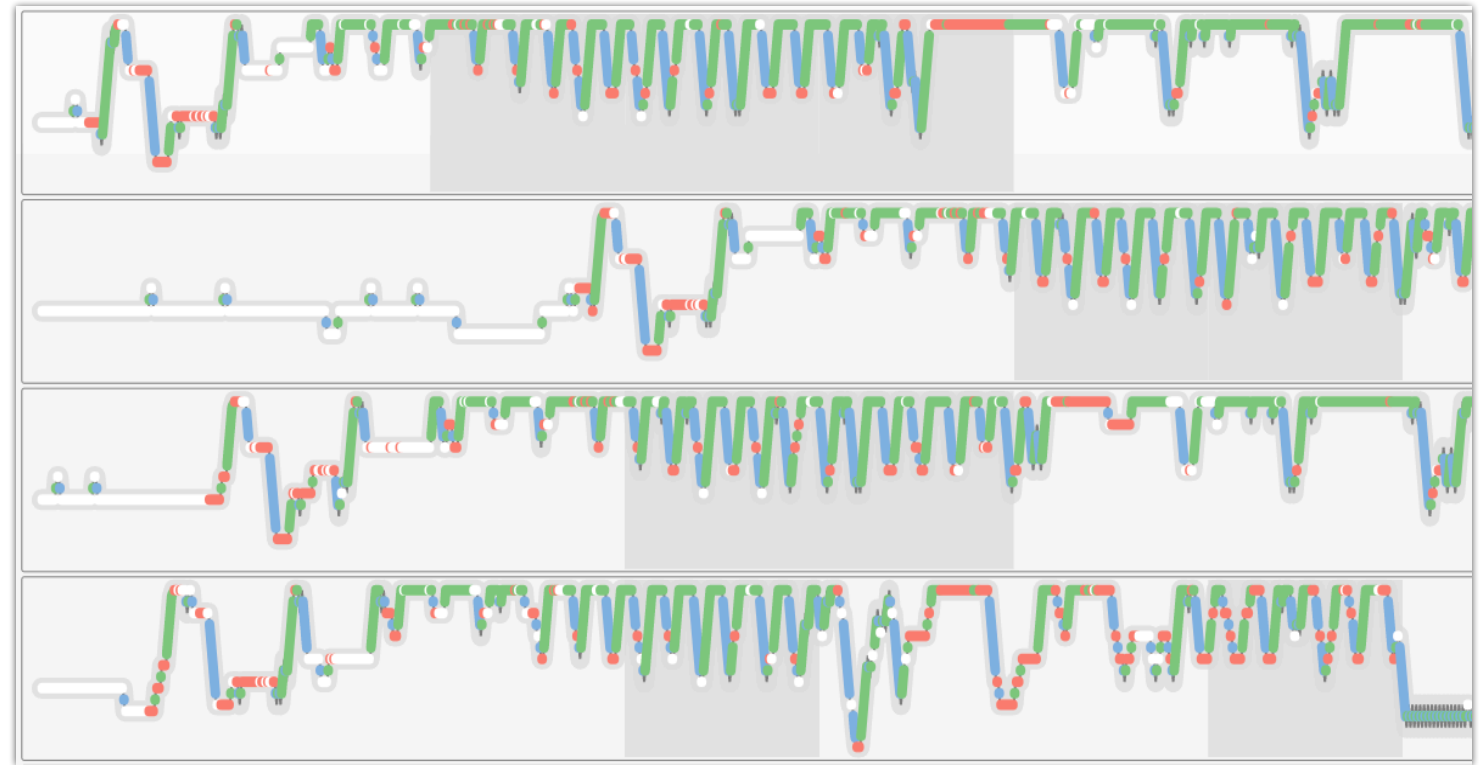
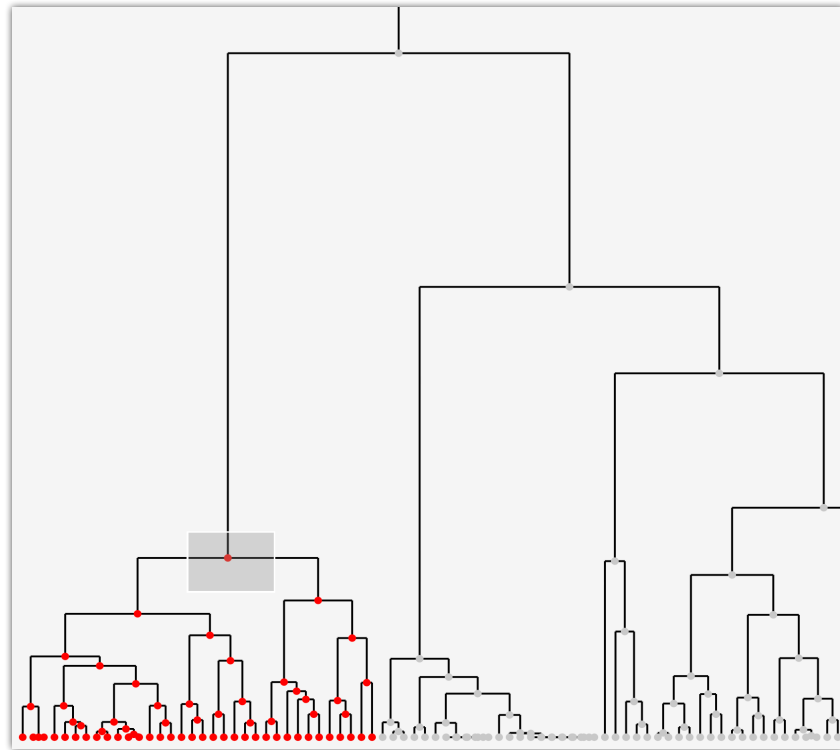
○ Noop (stay)    ● Firing the ball    ● Moving left    ● Moving right

# DQNViz: The Episode/Trajectory View

25000 game steps (actions)

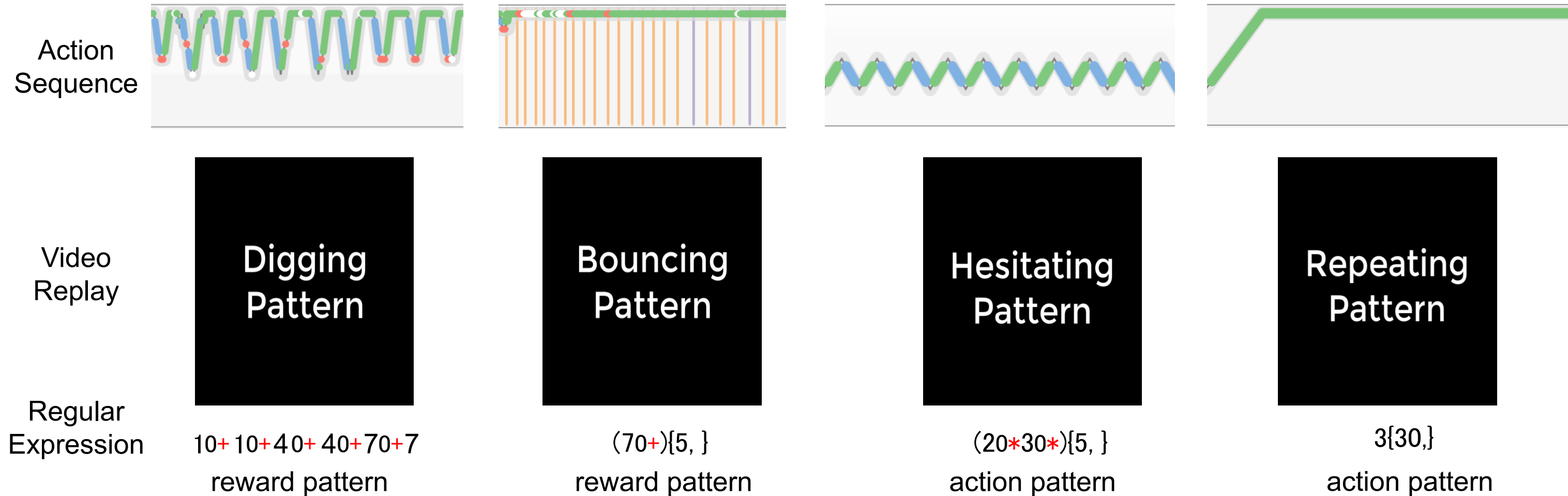


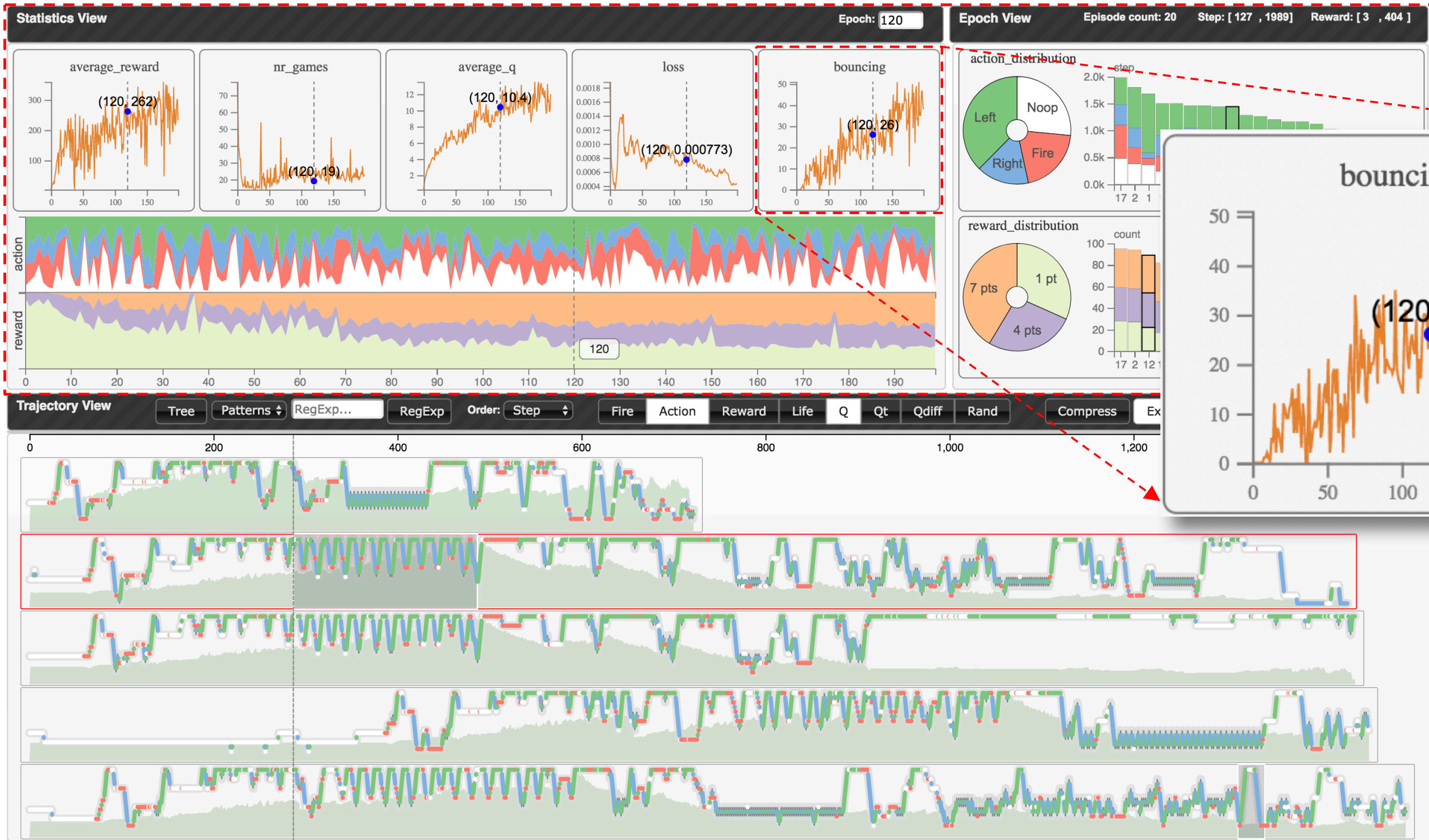
# DQNViz: Pattern Mining



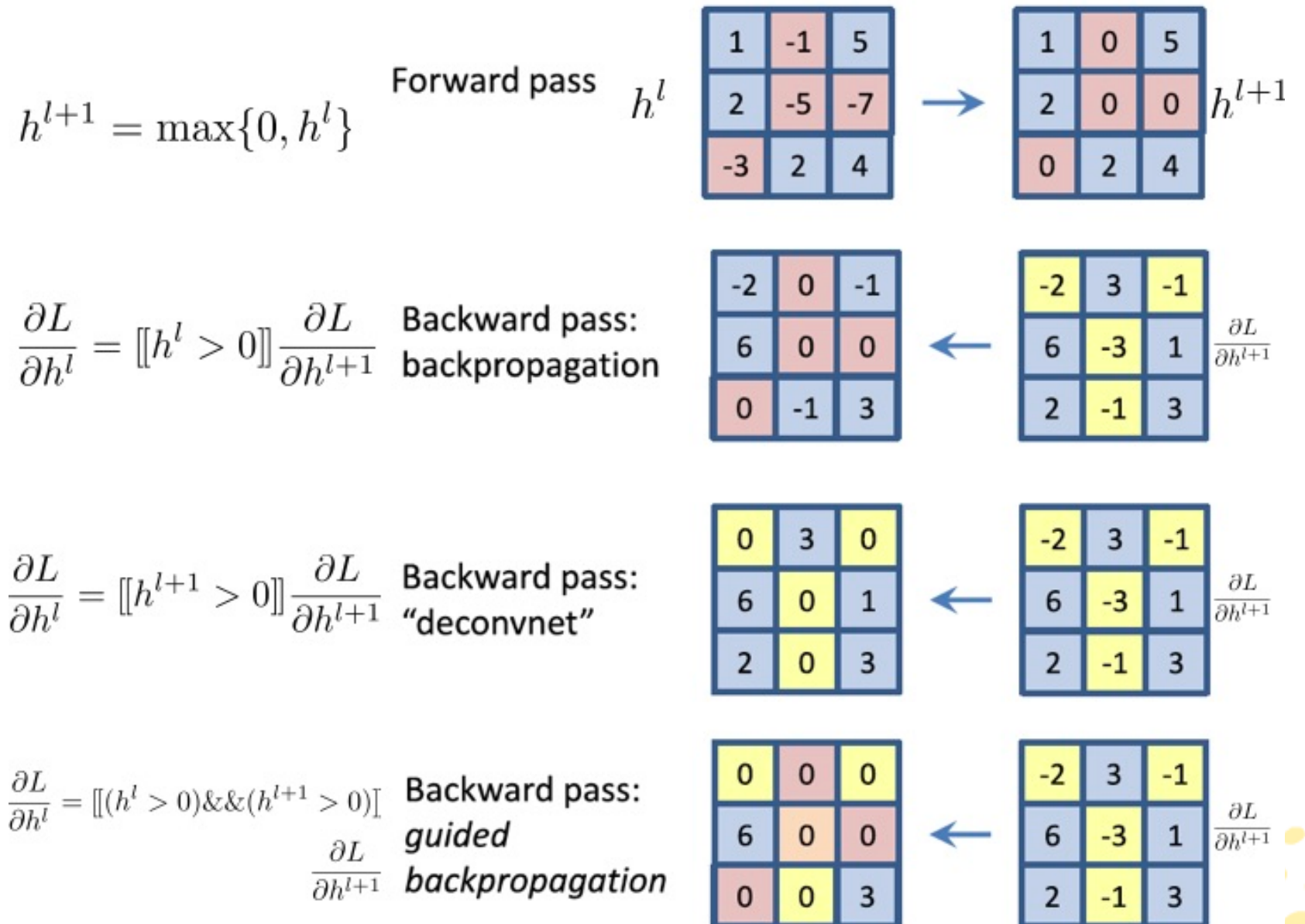
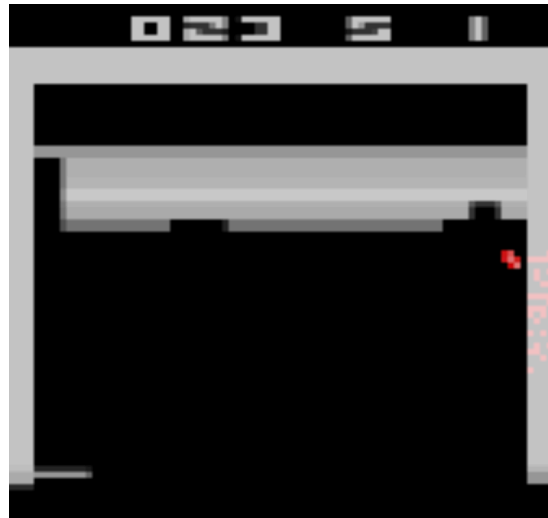


# DQNViz: Pattern Mining





# Guided Back-propagation



Saliency map from filter  $i$

4	3	1
2	3	0
0	2	0



The 3rd max = 3



1	0	0
0	0	0
0	0	0

Binary saliency map  
(activated pixel = 1)

Saliency map from filter  $j$

0	2	3
1	3	4
4	5	6



The 3rd max = 4



0	0	0
0	0	0
0	1	1

Binary saliency map  
(activated pixel = 2)

Saliency map from filter  $k$

1	2	1
2	2	2
2	2	2



The 3rd max = 2



0	0	0
0	0	0
0	0	0

Binary saliency map  
(activated pixel = 0)

U

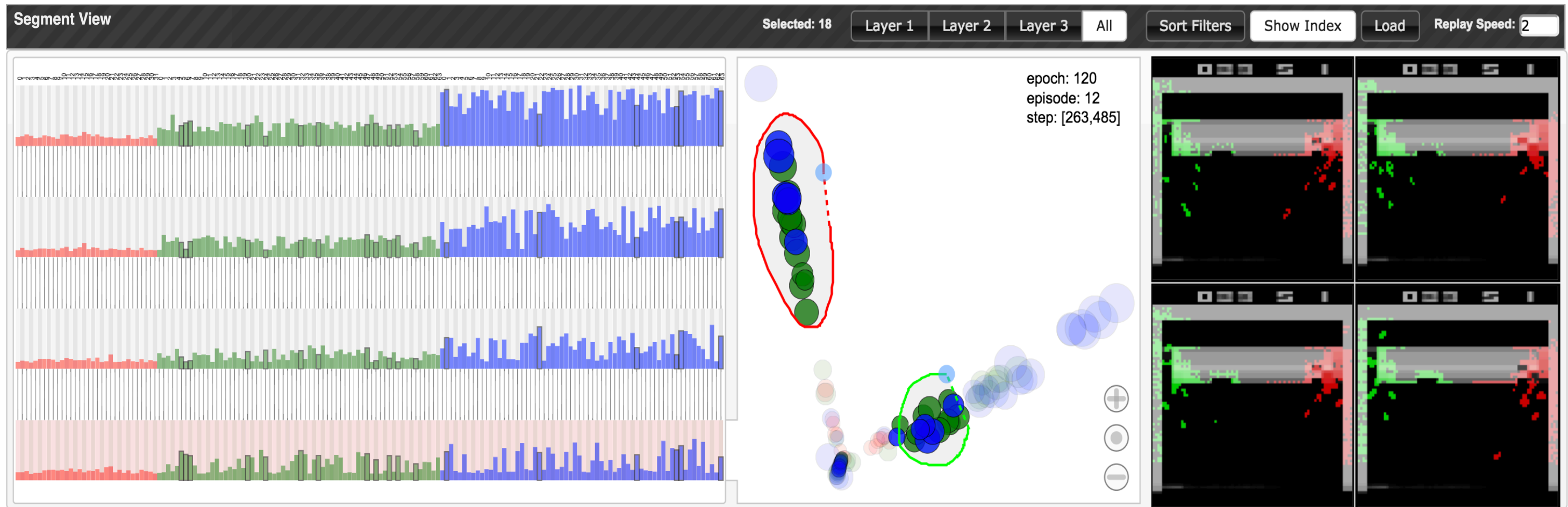
U

=

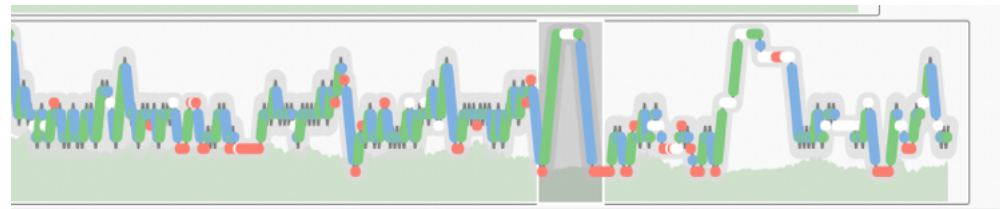
1	0	0
0	0	0
0	1	1

Final saliency map

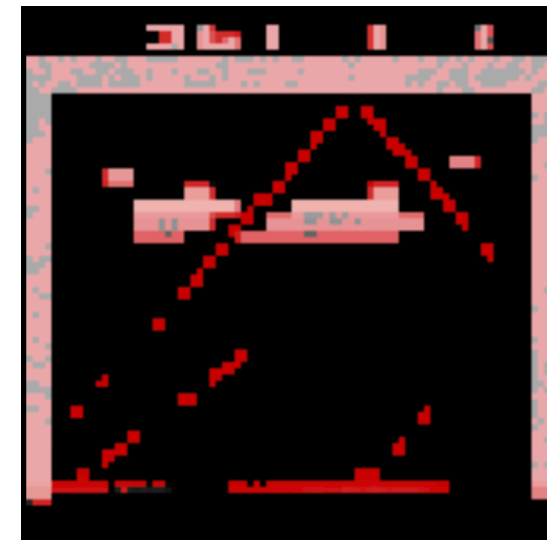
# Guided Back-propagation



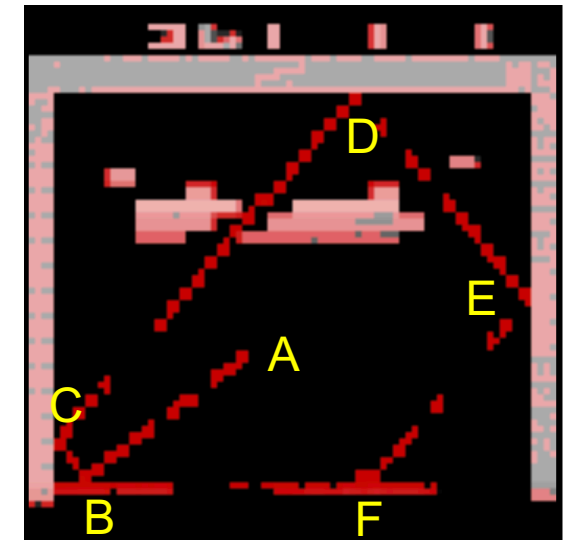
# Guided Back-propagation



Epoch 0



Epoch 10

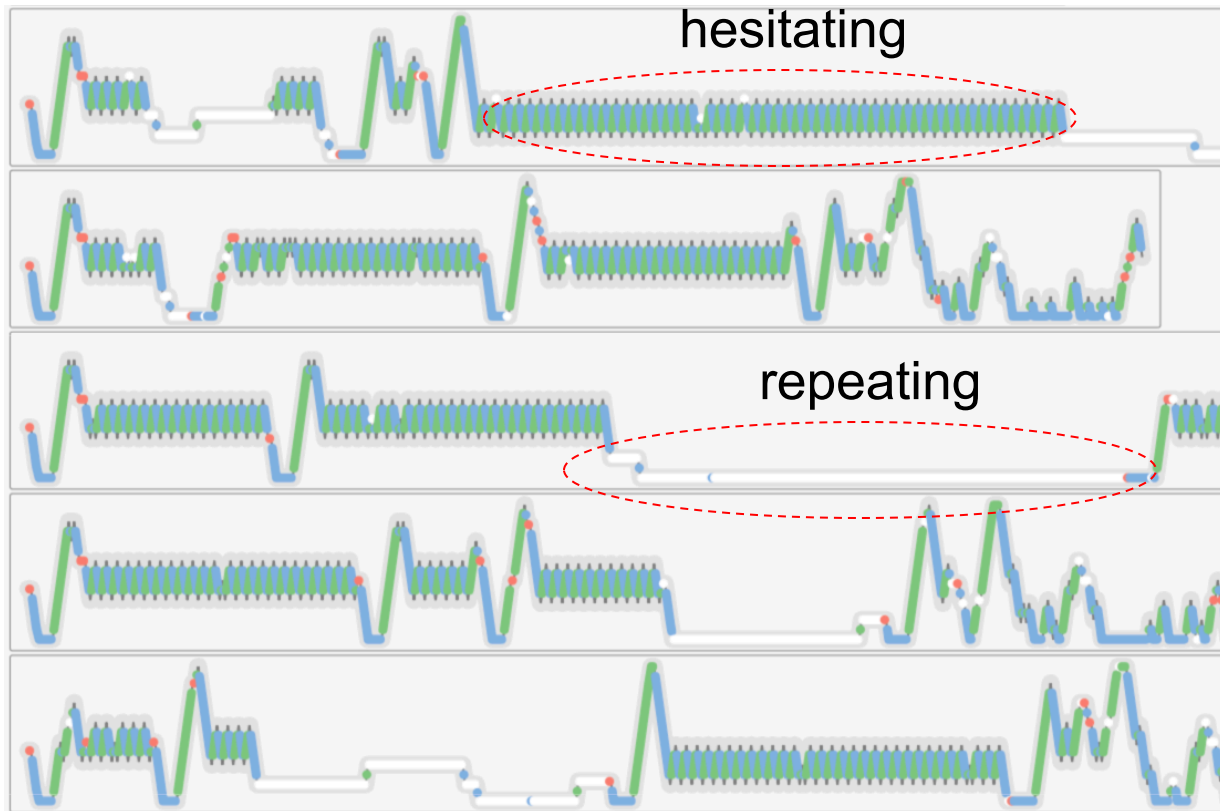


Epoch 120

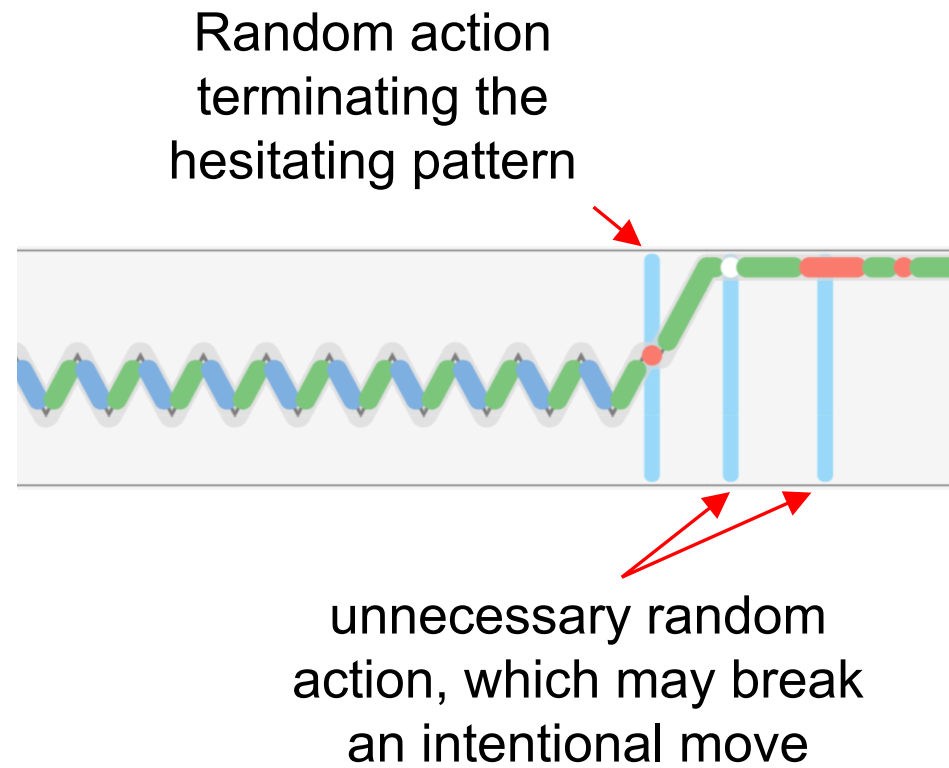


Part II: Optimize random actions to improve the DQN model training.

# Improve Random Actions



An unsuccessful training due to the hesitating and repeating pattern

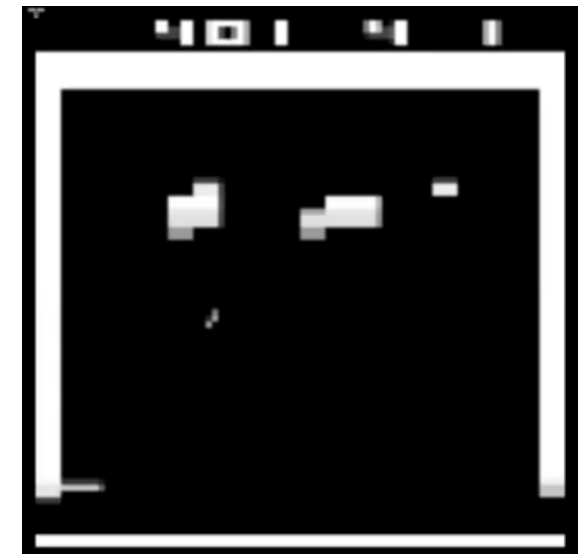
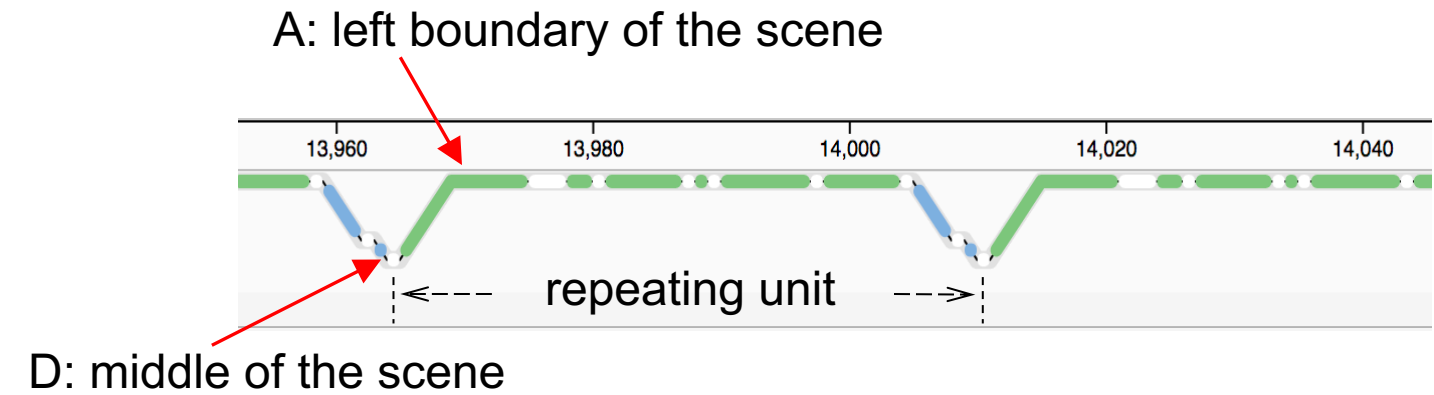
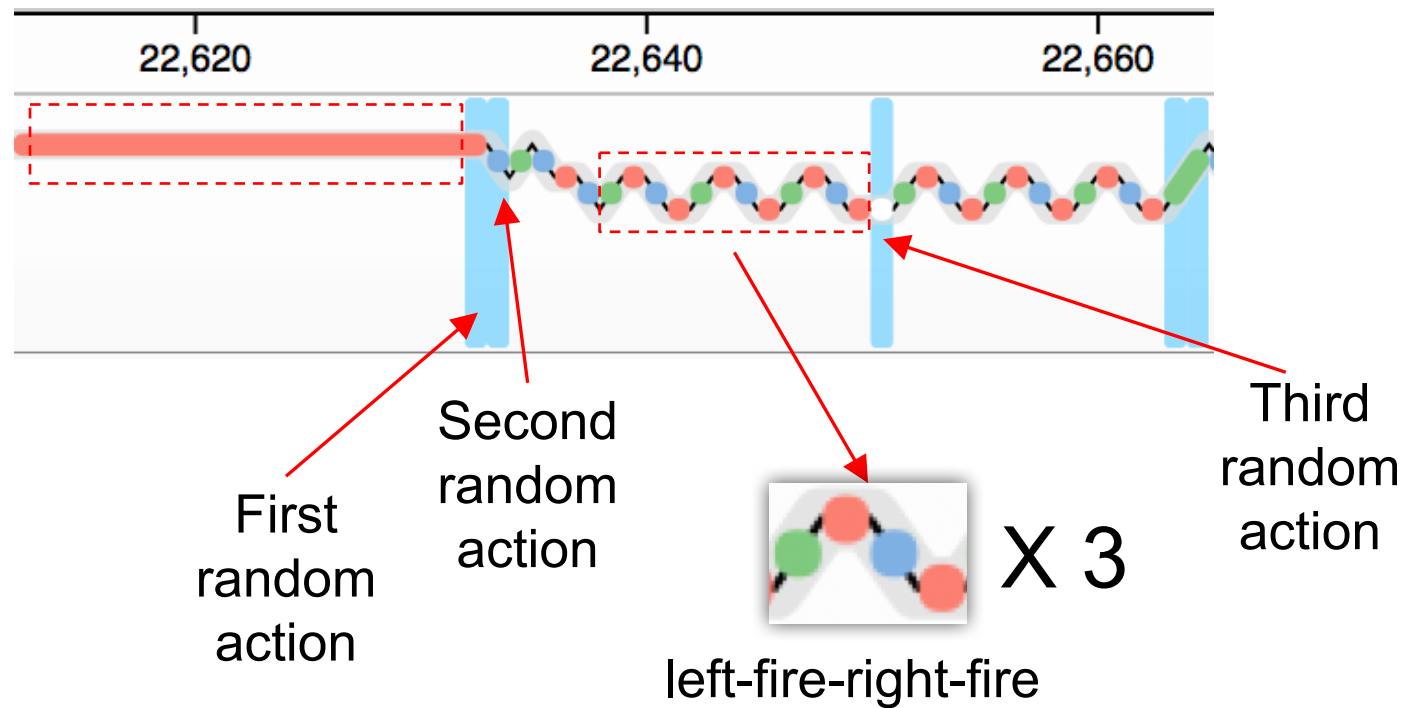


Random actions break the agent's intentional move and cause a life-loss.



# Improve Random Actions

- Noop (stay)
- Firing the ball
- Moving left
- Moving right



A-B-C-D-C-B-A

# Quantitative Evaluation

Result of averaging over 10 runs

Algorithm	# of Steps	# of episodes	Total rewards	# or random steps
Random rate 5%	25,000	16.6	4198.6	1269.4
Our RegExp Alg.	25,000	11.4	4899.2	503
Random rate 2%	25,000	9.9	3780.8	492.1

# Conclusion

- We present DQNViz, a visual analytics system that provides effective multi-level visual summaries of the large multi-faceted data generated from DQN trainings.
- Based on our analysis of the training data, we identified typical movement and reward patterns of the agent, and those patterns have helped in controlling the random actions of the DQN model.

# Thanks! Questions?

This work was supported in part by US Department of Energy Los Alamos National Laboratory contract 47145, UT-Battelle LLC contract 4000159447, NSF grants IIS-1250752, IIS1065025, and US Department of Energy grants DE-SC0007444, DEDC0012495, program manager Lucy Nowell.

