

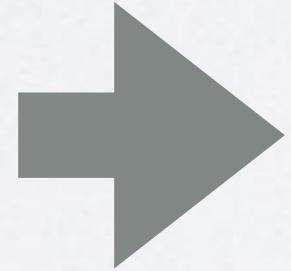
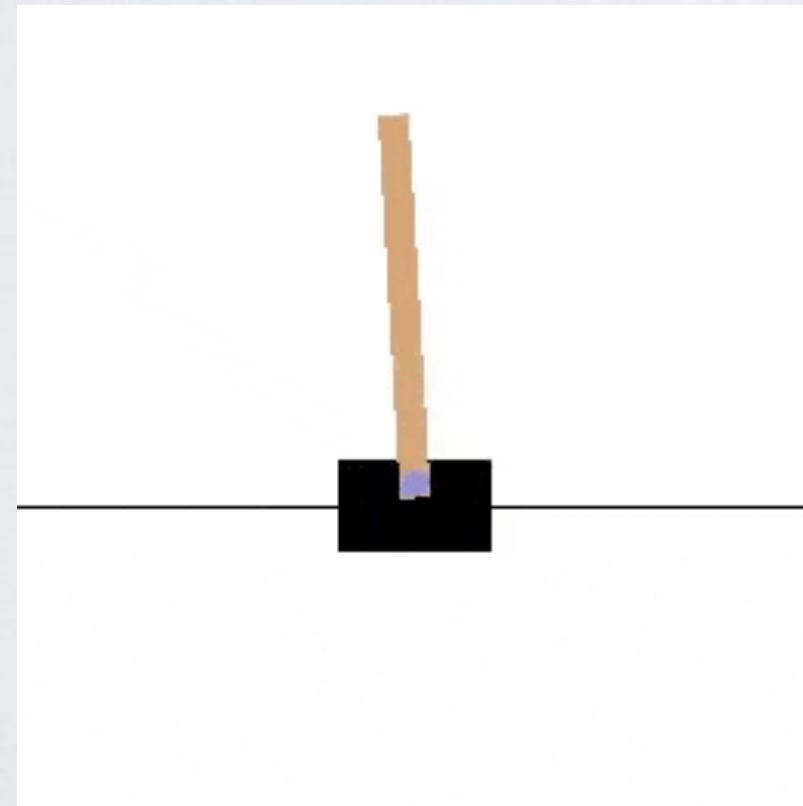
ENGINEERING REINFORCEMENT LEARNING

Deepest
2018. 7. 28 Hosting
Sungwon, Lyu
lyusungwon@dm.snu.ac.kr

REINFORCEMENT LEARNING

- **Value function approximation**

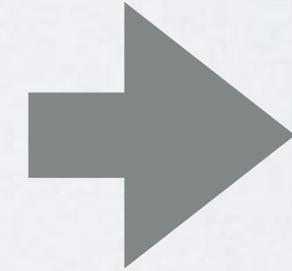
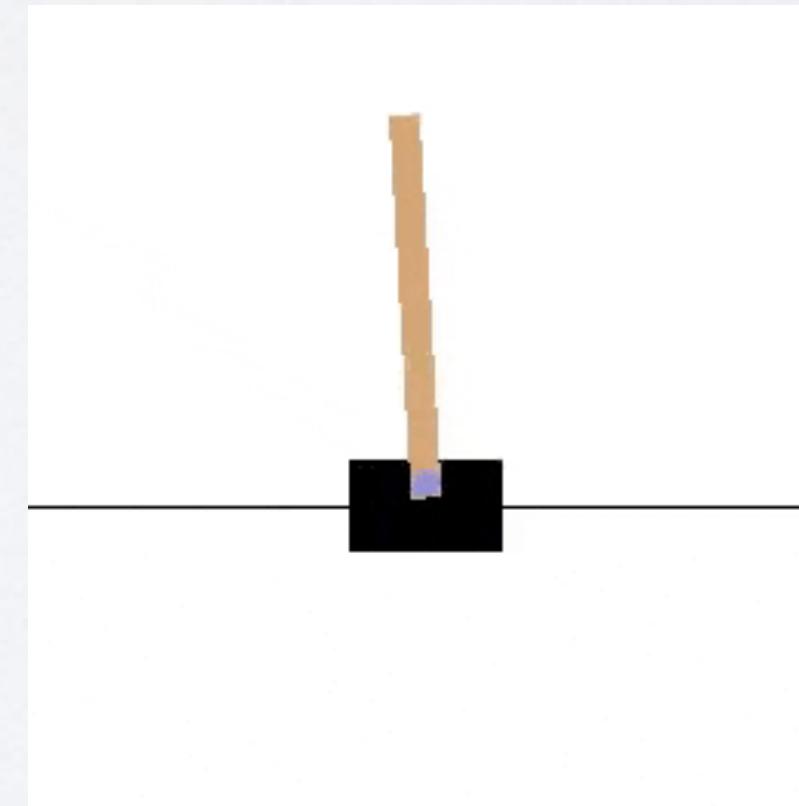
- SARSA, DQN
- $Q(S, A) = \text{Value}$



stay: 0.72
left: 0.86
right: 0.21

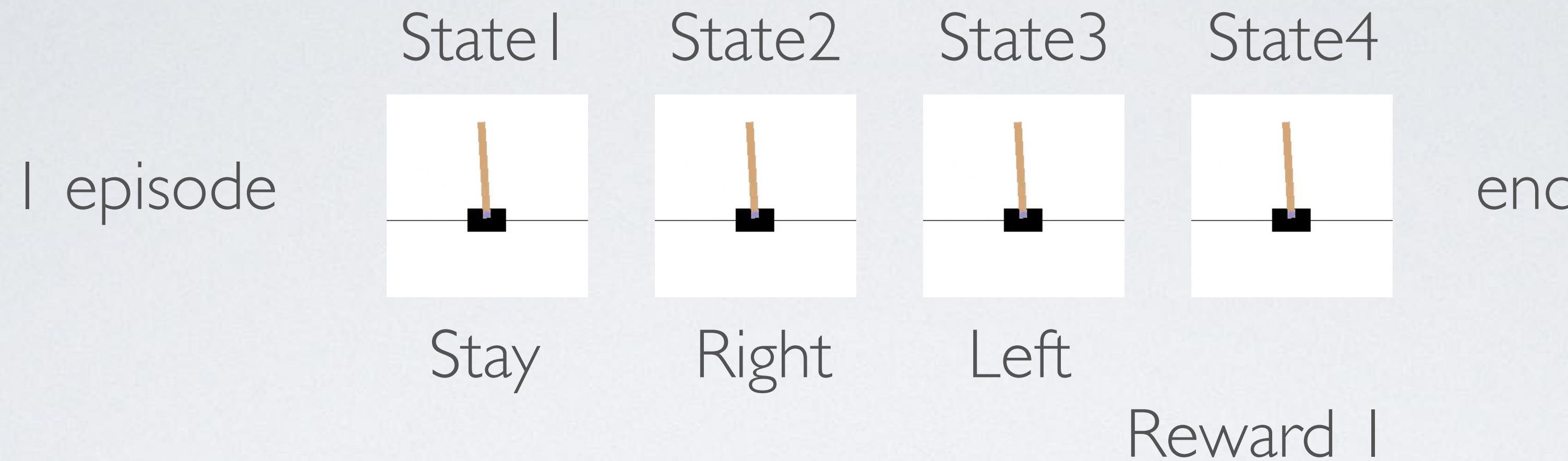
- **Policy gradient**

- REINFORCE, A2C, DDPG, A3C
- $\pi(S) = \text{Probability}$



stay: 36%
left: 45%
right: 19%

REINFORCEMENT LEARNING



- **Temporal Difference**

- Learn from incomplete episode
- Save (S2, Right, Reward 0, S3, Left)
 - $Q(\text{State1}, \text{Stay}) = 0 + 0.9 * Q(\text{State2}, \text{A})$
 - $Q(\text{State2}, \text{Right}) = 0 + 0.9 * Q(\text{State3}, \text{A})$
 - $Q(\text{State3}, \text{Left}) = 1$

- **Monte-Carlo method**

- Learn from complete episode
- Save Entire episode
 - $Q(\text{State3}, \text{Left}) = 1$
 - $Q(\text{State2}, \text{Right}) = 0.9$
 - $Q(\text{State1}, \text{Stay}) = 0.81$

REINFORCEMENT LEARNING

- **DQN**

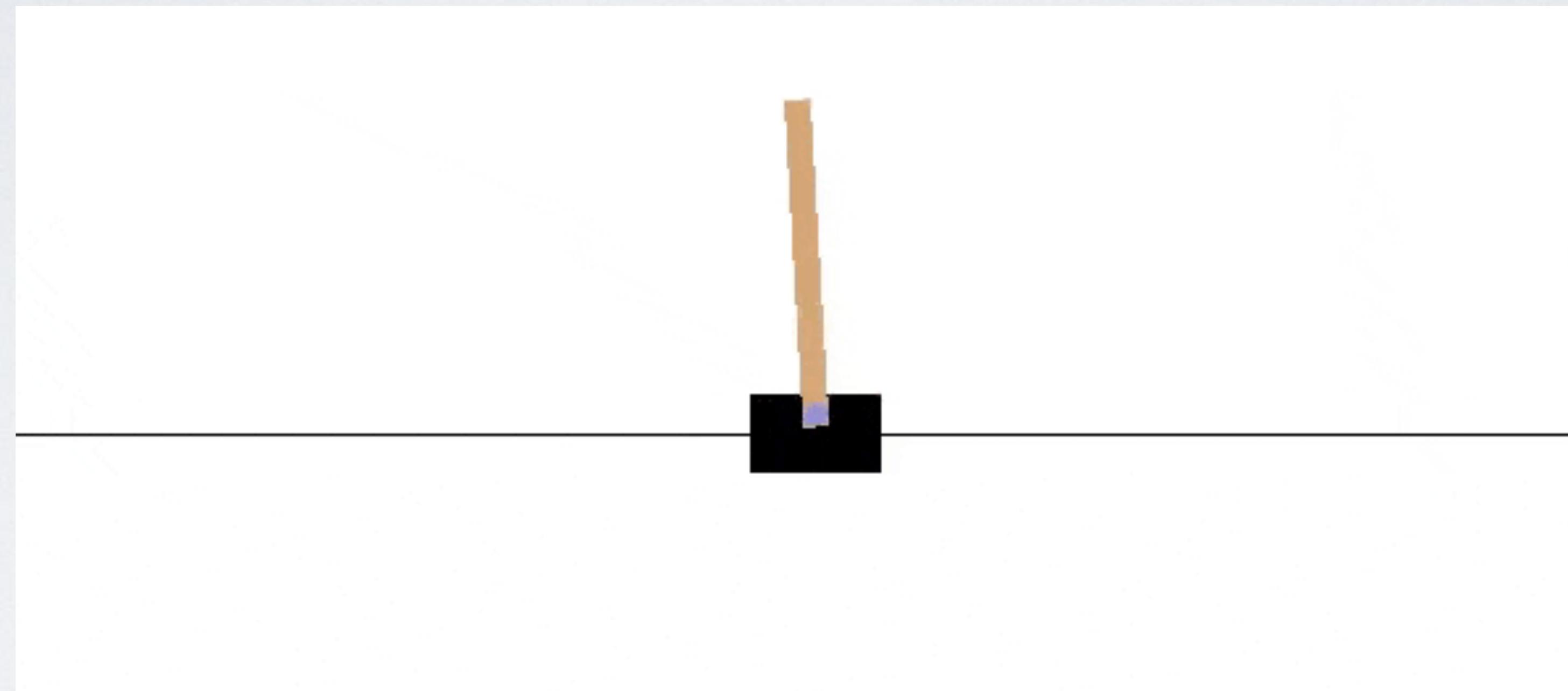
- Value function approximation & TD
- Epsilon decay
- $Q(S,A)$ as neural network (CNN)
- Off-policy (target network)
- Replay Memory

$$Y_t^{DQN} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t^-)$$

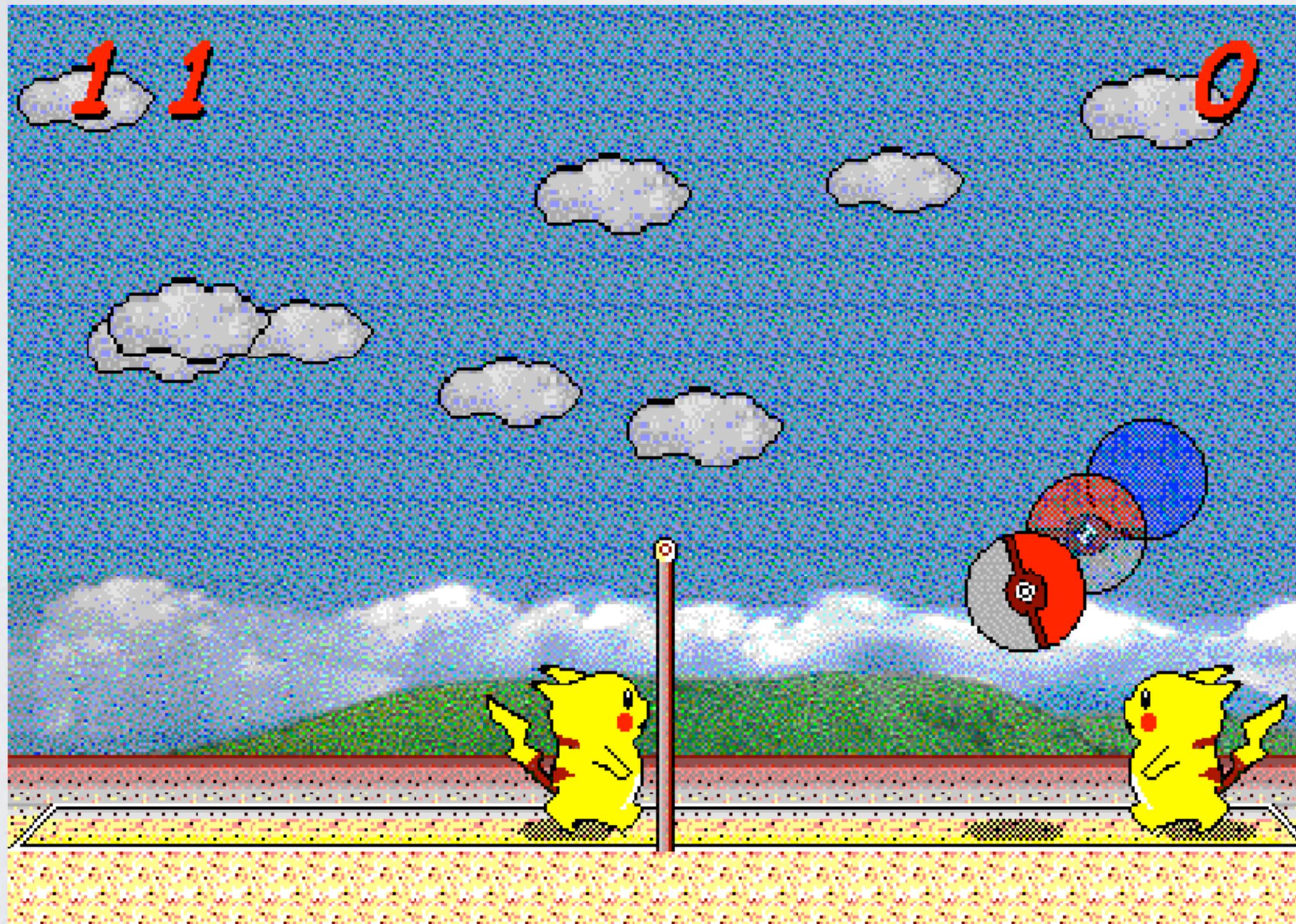
$$J(\theta) = (Q(S_t, a_t; \theta_t) - Y_t^{DQN})^2$$



MOTIVATION



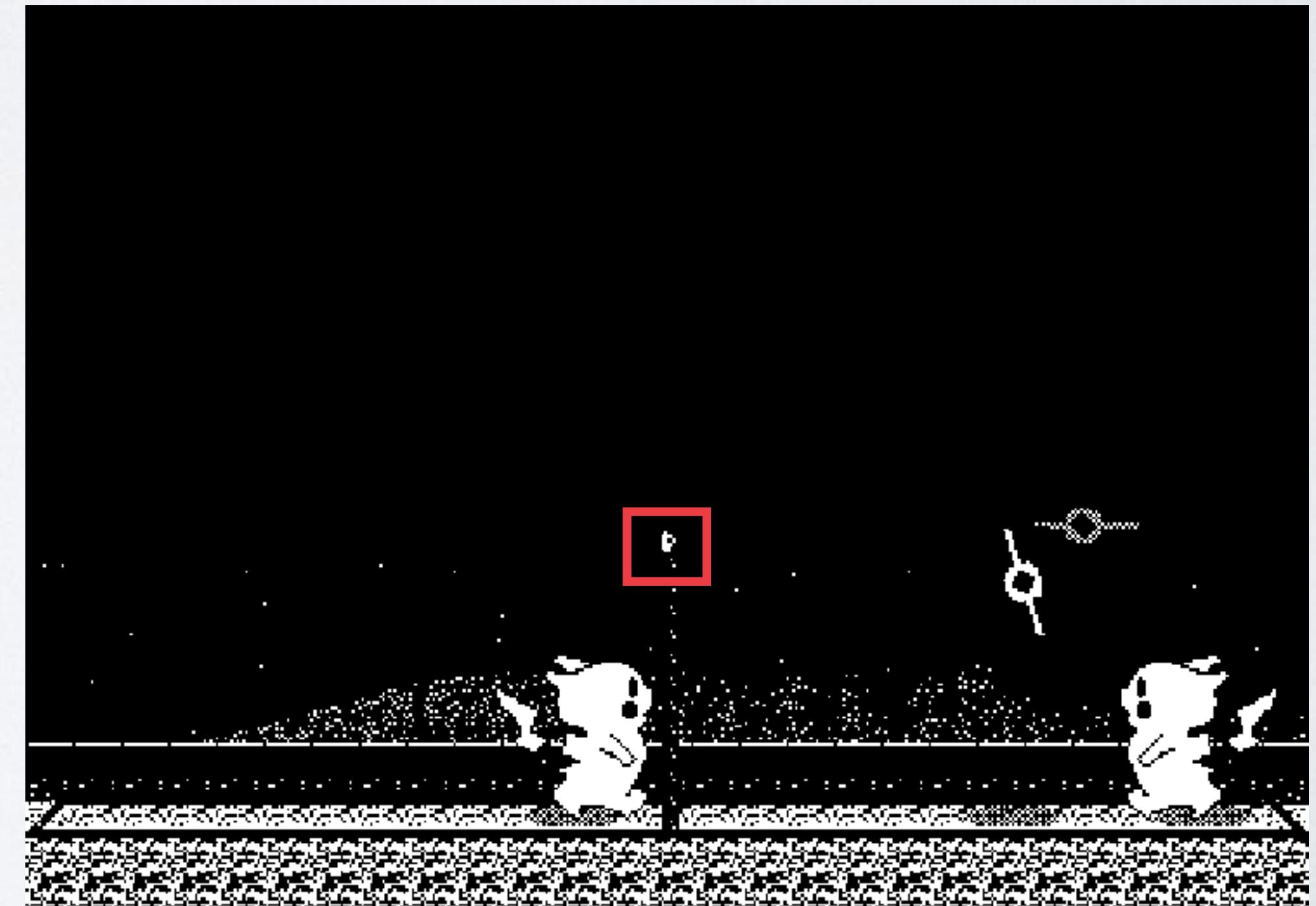
PIKACHU - VOLLEYBALL



- Window only?
 - No Linux
 - No Server
- No environment toolkit
 - Screen, key-input, timing
 - Start, end, state, action, reward
- Noisy background
 - Clouds
 - Scores

VERSION 1.0

- Window only?
 - Wine
- No environment toolkit
 - Screenshot (mss package)
 - Template matching
 - Key input (pynput package)
 - Feature engineering (Start / end / restart)
- Noisy background
 - => Color filtering



Fail

DQN-EXTENSION



DQN-EXTENSION

- **Double DQN**

- Prevent over-optimization

$$Y_t^{DQN} = R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t^-)$$

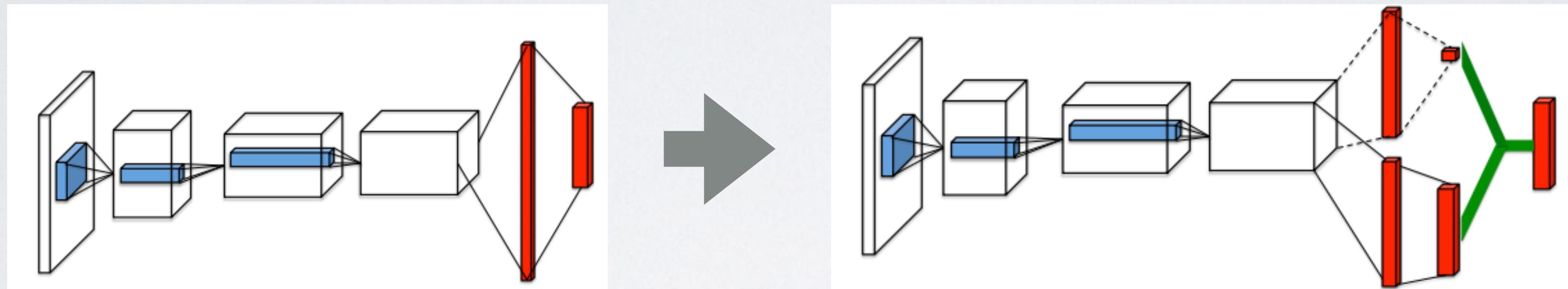
$$Y_t^{DoubleQ} = R_{t+1} + \gamma Q(S_{t+1}, \text{argmax}_a Q(S_{t+1}, a; \theta_t); \theta_t')$$

- Implementation: | min cut
- => DDQN

DQN-EXTENSION

- **Dueling DQN**

- Action is not important for some state



$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + (A(s, a; \theta, \alpha) - \frac{1}{\|\mathcal{A}\|} \sum_{a'} A(s, a'; \theta, \alpha))$$

- Implementation: 5min cut
- => Dueling DQN

DQN-EXTENSION

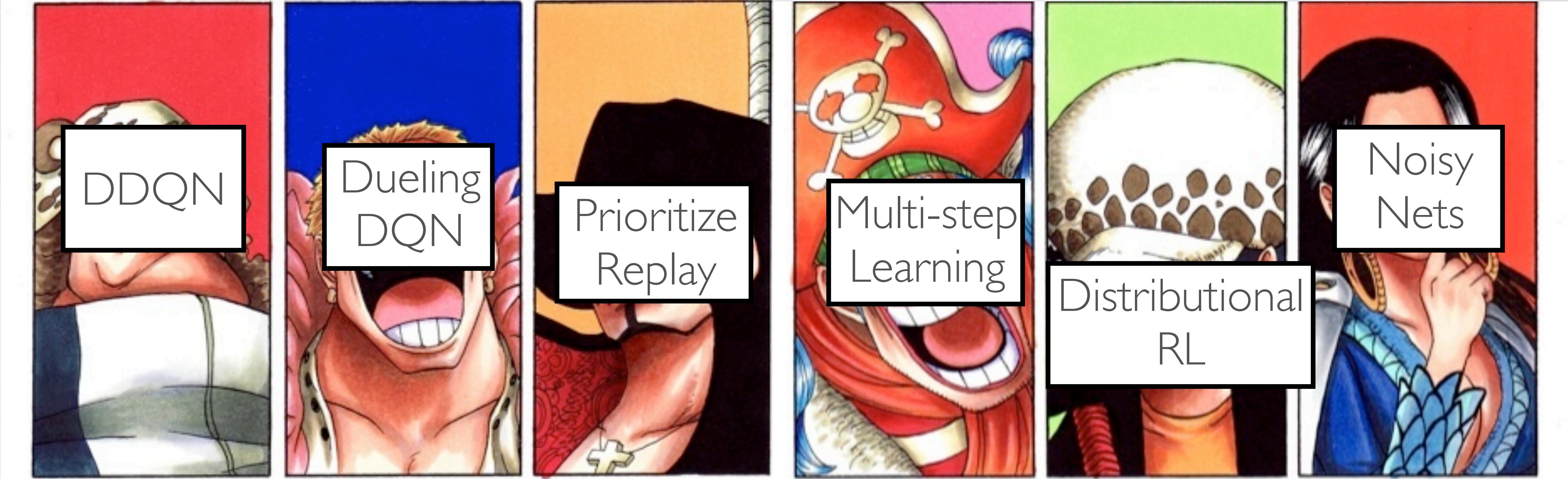
- **Prioritize replay**

- Not all experience are equal
- Weighted Sample according to ‘Surprise’

$$\delta_j = R_j + \gamma_j Q_{target}(S_j, \text{argmax}_a Q(S_j, a)) - Q(S_{j-1}, A_{j-1})$$

- Implementation: 30min cut
- => Prioritized DQN

MORE DQN-EXTENSION



- => Rainbow

Source:

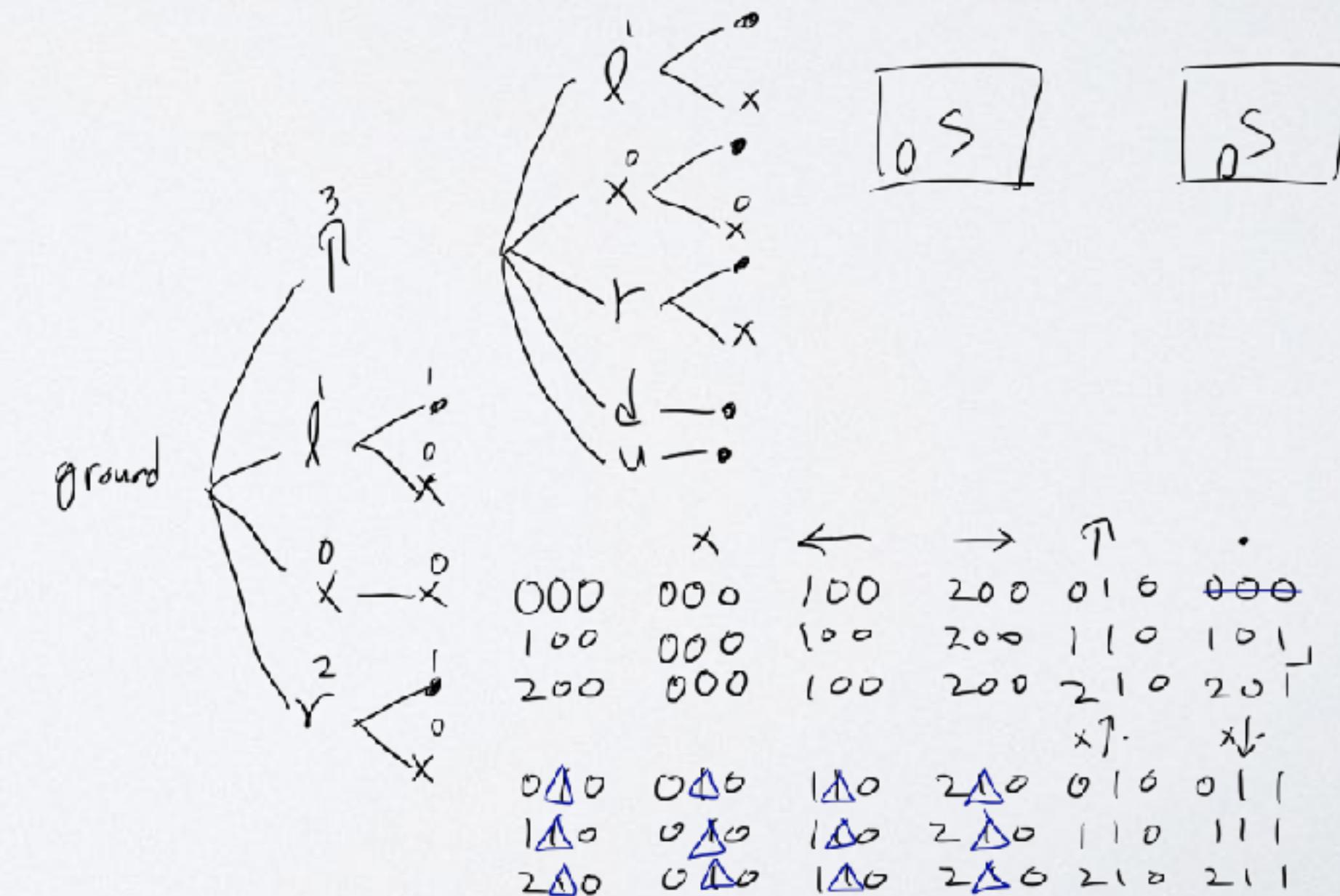
Bellemare, Marc G., Will Dabney, and Rémi Munos. "A distributional perspective on reinforcement learning." *arXiv preprint arXiv:1707.06887* (2017).

Fortunato, Meire, et al. "Noisy networks for exploration." *arXiv preprint arXiv:1706.10295* (2017).

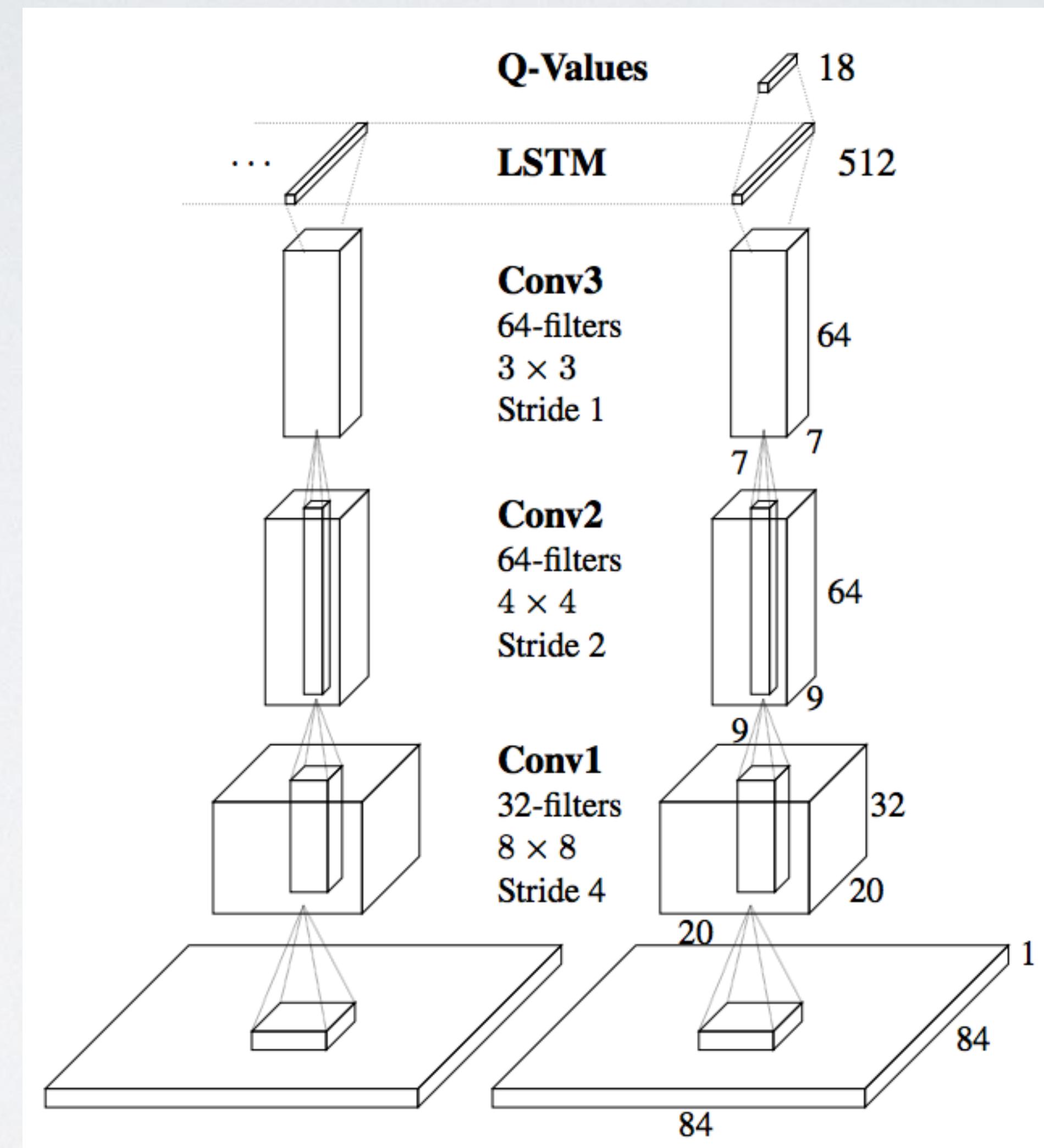
Hessel, Matteo, et al. "Rainbow: Combining improvements in deep reinforcement learning." *arXiv preprint arXiv:1710.02298*(2017).

VERSION 2.0

- Pytorch
- Prioritized Dueling DDQN
- Model
 - c16(3,1)-c32(3,1)-c64(3,1)-c64(3,1)-c64(3,1) => c16(8,4)-c32(4,2)-c32(3,1)
- State: $4 * 140 * 215 \Rightarrow 4 * 84 * 129$
- Action: 9 => 6
- Frame: 0.125 -> 0.2



DRQN

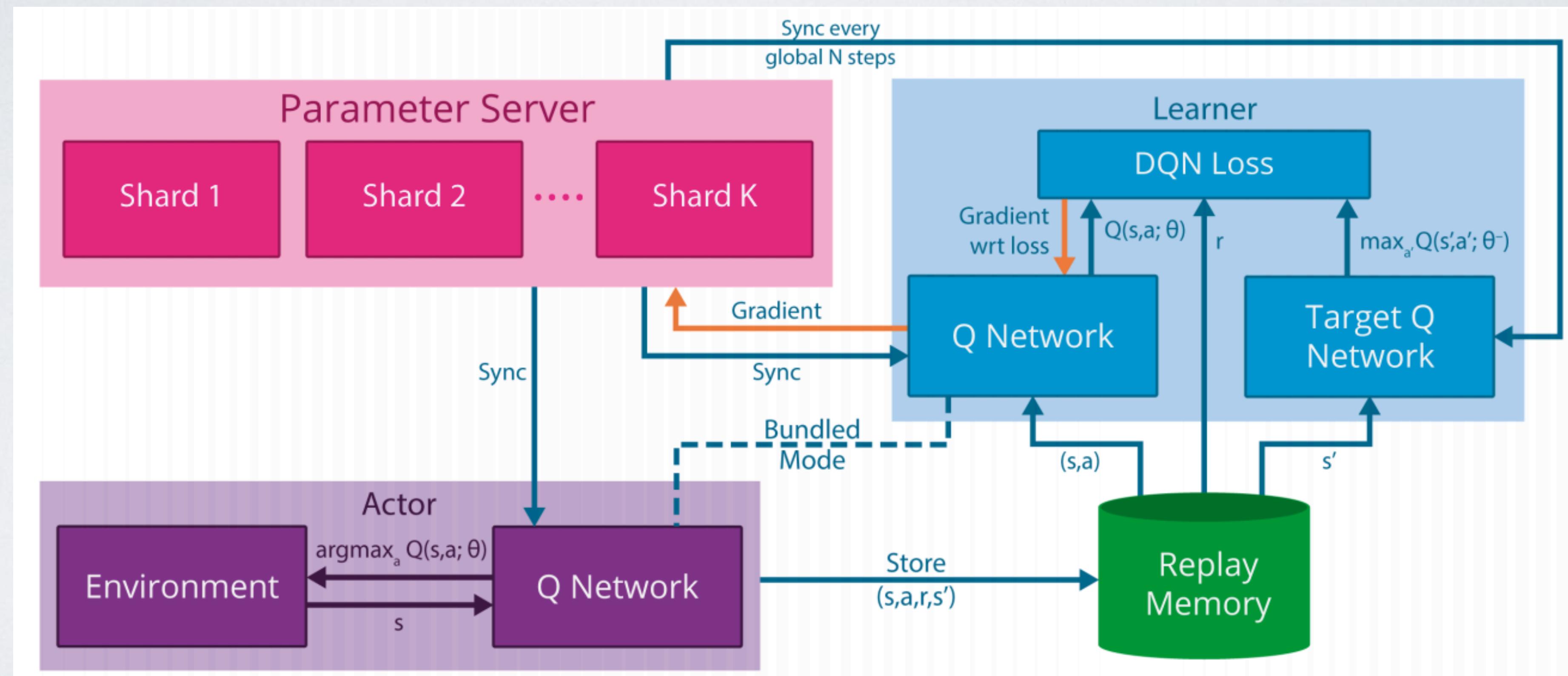


XVFB

- XVFN(X Virtual Frame Buffer)
 - Display server: operate graphics on memory
 - How?
 - sudo apt-get install xvfb
 - xvfb :99 -ac -screen 0 1280x1024x24 > /dev/null &
 - echo "export DISPLAY=:99" >> ~/.bashrc
 - DISPLAY=:99 python actor.py
- xdotool
 - sudo apt-get install xdotool
 - DISPLAY=:99 xdotool search --name “PIKA”
 - DISPLAY=:99 xdotool key --window 12582913 KP_Enter

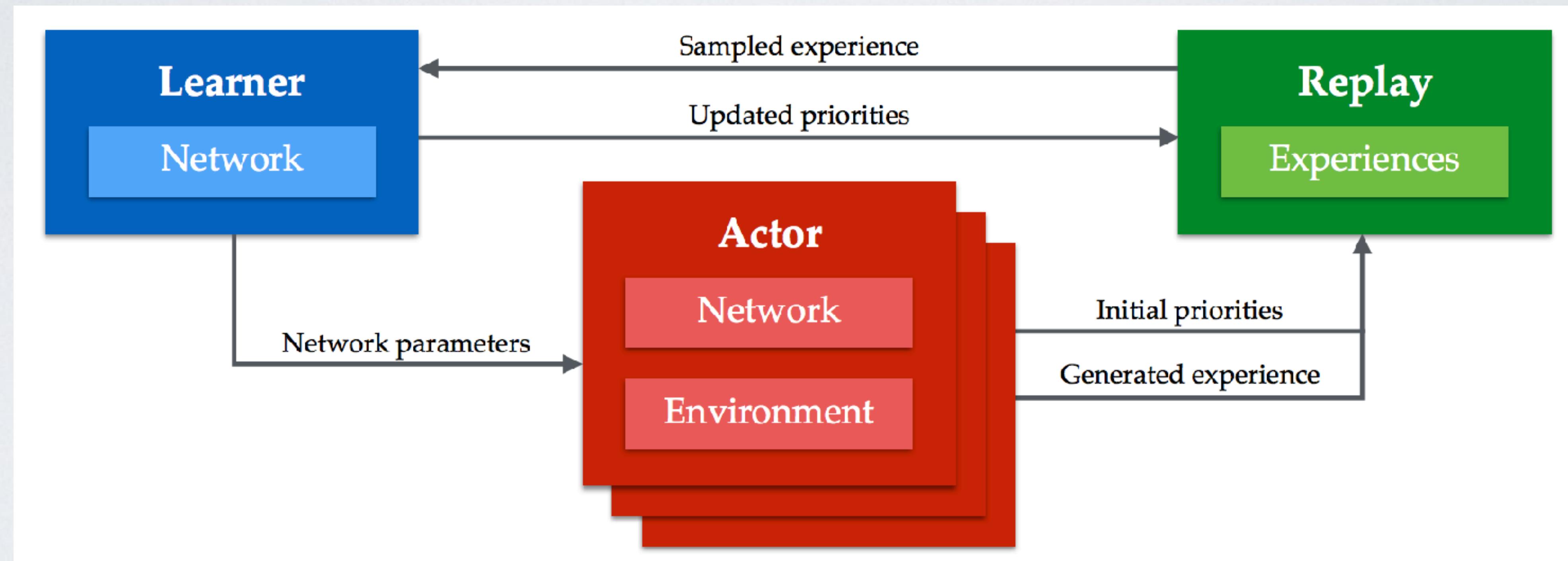
GORILA

- Gorila (General Reinforcement Learning Architecture)



APE-X

- Ape-X DQN / DPG



STATE-OF-THE-ART BENCHMARK

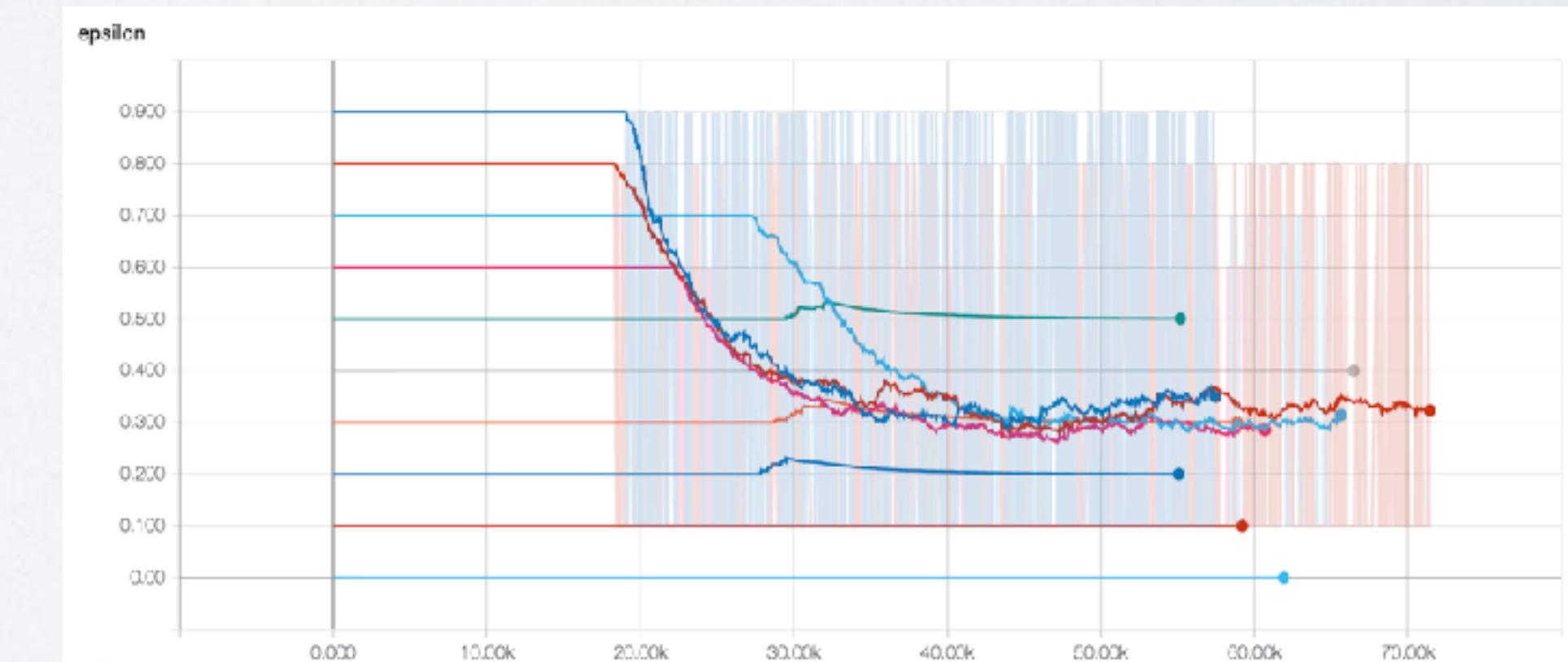
Algorithm	Frames	Bowling		Venture	
		%	raw	%	raw
RUDDER	12M	62.10	108.55	96.55	1,147
rainbow	200M	5.01	30	0.46	5.5
Prioritized DDQN	200M	28.71	62.6	72.67	863
Noisy DQN	200M	39.39	77.3	0	0
Dueling DDQN	200M	30.81	65.5	41.85	497
DQN	200M	19.84	50.4	13.73	163
Distributional DQN	200M	37.06	74.1	93.22	1,107
DDQN	200M	32.7	68.1	8.25	98
Ape-X DQN	22,800M	-17.6	4	152.67	1,813
Random	-	0	23.1	0	0
Human	-	100	160.7	100	1,187

Table 1: Results of RUDDER and other methods when learning the Atari games Bowling and Venture. Normalized human-percentage and raw scores over 200 testing-games with no-op starting condition: A3C scores are not reported, as not available for no-op starting condition. Scores for other methods were taken from previous publications [8, 39]. The RUDDER model is chosen based only on its training loss over 12M frames.

VERSION 3.0

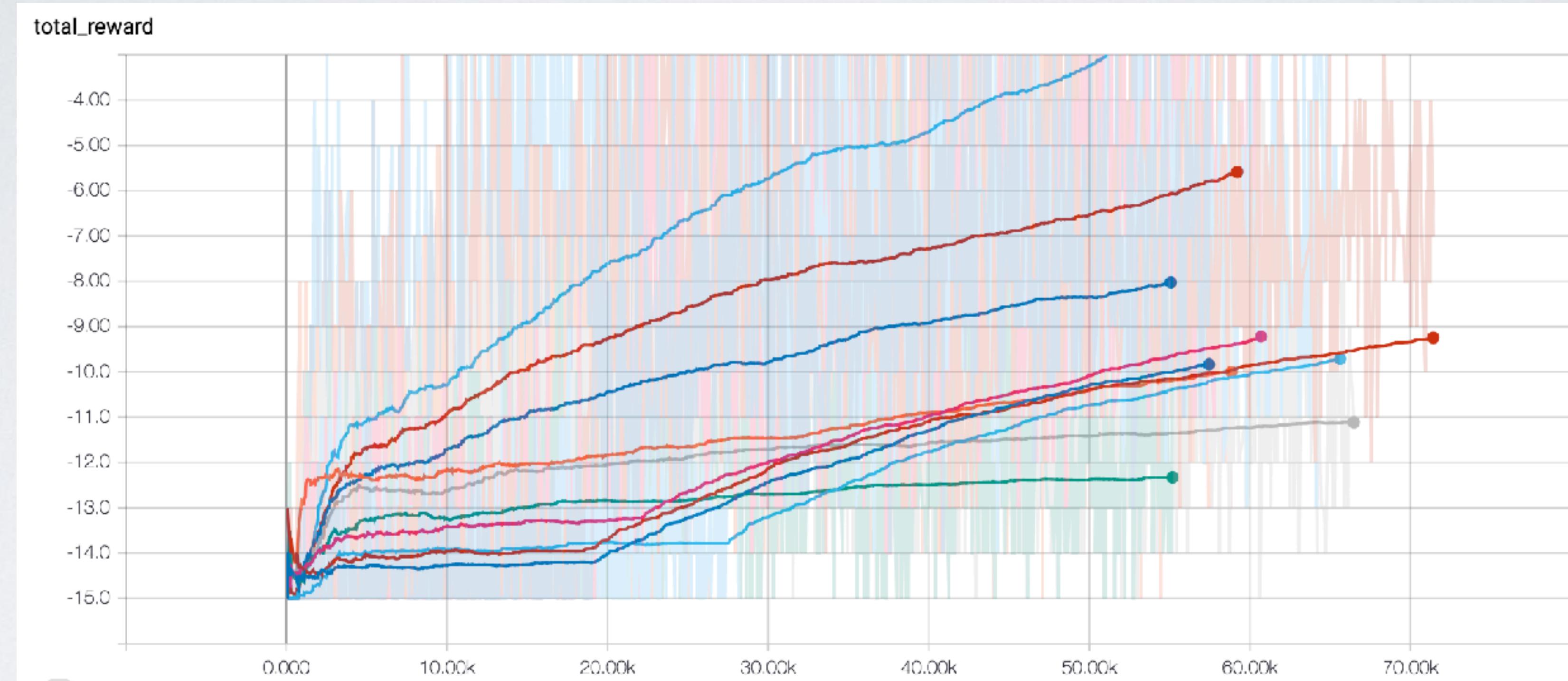
- Multi-thread => Multi-process
 - Communicate through file
 - 10 actors and 1 learner
 - Different epsilon per actor
- Run on server
 - Learner on gpu / actors on cpu
- Memory allocation
 - 40000 experience replay for learner
 - 4000 per actor

```
(0) 0: watch* "ubuntu"
(1) 1: python "ubuntu"
(2) 2: python "ubuntu"
(3) 3: pika.exe "ubuntu"
(4) 4: python "ubuntu"
(5) 5: pika.exe "ubuntu"
(6) 6: python- "ubuntu"
(7) 7: pika.exe "ub top - 15:09:00 up 93 days, 4:26, 169 users, load average: 50.22, 48.66, 49.02
(8) 8: python "ubur Tasks: 844 total, 20 running, 816 sleeping, 7 stopped, 1 zombie
(9) 9: pika.exe "ub %Cpu(s): 44.7 us, 22.5 sy, 0.0 ni, 32.8 id, 0.0 wa, 0.0 hi, 0.0 si, 0.0 st
(a) 10: python "ub KiB Mem : 13191706+total, 4057884 free, 39507568 used, 88351616 buff/cache
(b) 11: pika.exe "ub KiB Swap: 7812092 total, 6264740 free, 1547352 used. 91039424 avail Mem
(c) 12: python "ub
(d) 13: pika.exe "u PID USER PR NI VIRT RES SHR S %CPU %MEM TIME+ COMMAND
(e) 14: python "ub 794 sungwon 20 0 5977376 373592 92196 R 247.2 0.3 1743:46 python
(f) 15: pika.exe "u 18044 sungwon 20 0 6228932 495424 90796 R 231.4 0.4 12978:23 python
(g) 16: python "ub 22715 sungwon 20 0 6074112 346256 92000 R 224.8 0.3 123:54.15 python
(h) 17: pika.exe "u 2416 sungwon 20 0 6181128 452900 92476 S 200.0 0.3 6584:00 python
(i) 18: python "ub 28457 sungwon 20 0 6309544 581176 91480 S 185.8 0.4 18758:23 python
(j) 19: pika.exe "u 32247 sungwon 20 0 5989560 323812 92628 S 135.6 0.2 841:51.38 python
(k) 20: python "ub 26905 sungwon 20 0 6145956 414600 90480 R 126.7 0.3 15082:39 python
(l) 21: pika.exe "u 6343 sungwon 20 0 6202300 537452 91748 R 126.1 0.4 13580:58 python
(m) 22: python "ub 19807 sungwon 20 0 6018528 353980 92104 R 120.1 0.3 120:10.37 python
9179 jehyuk 20 0 322536 203288 11580 R 100.0 0.2 2:40.32 conda
33154 sungwon 20 0 50.014g 9.379g 293940 R 99.7 7.5 17381:21 python
16642 sungwon 20 0 45.744g 3.007g 411340 R 92.4 2.4 3589:46 python
19405 sungwon 20 0 48892 13996 2900 R 86.8 0.0 16283:33 wineserver
28198 sungwon 20 0 2803528 19764 15920 R 65.3 0.0 4910:32 pika.exe
22563 sungwon 20 0 2803592 19944 16104 S 64.4 0.0 31:53.45 pika.exe
32148 sungwon 20 0 2803528 19840 16000 R 64.0 0.0 236:30.18 pika.exe
[0] 0:[tmux]* 1:python
[0] <a.exe 20:python 21:pika.exe 22:python 23:top*> "ubuntu" 15:08 27-Jul-18
```

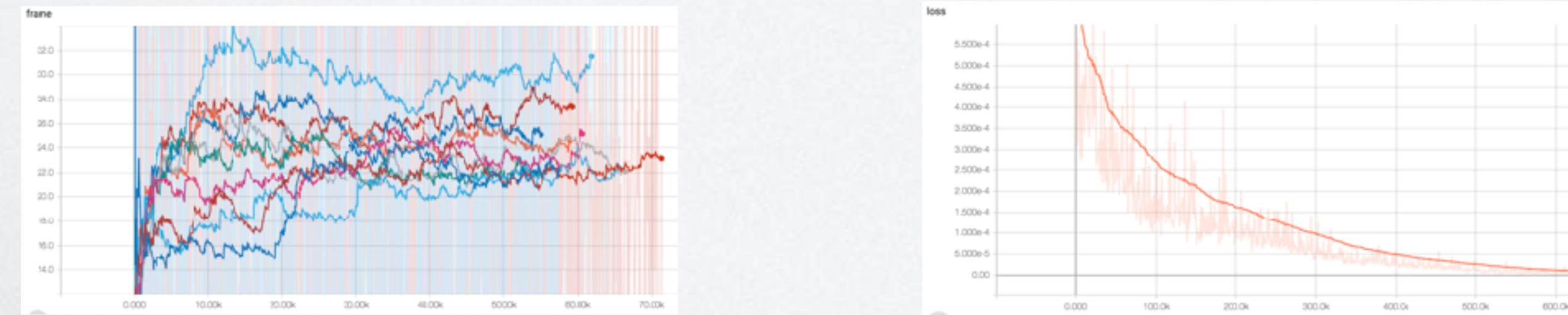


RESULT I

- Total reward (15 - Score difference)

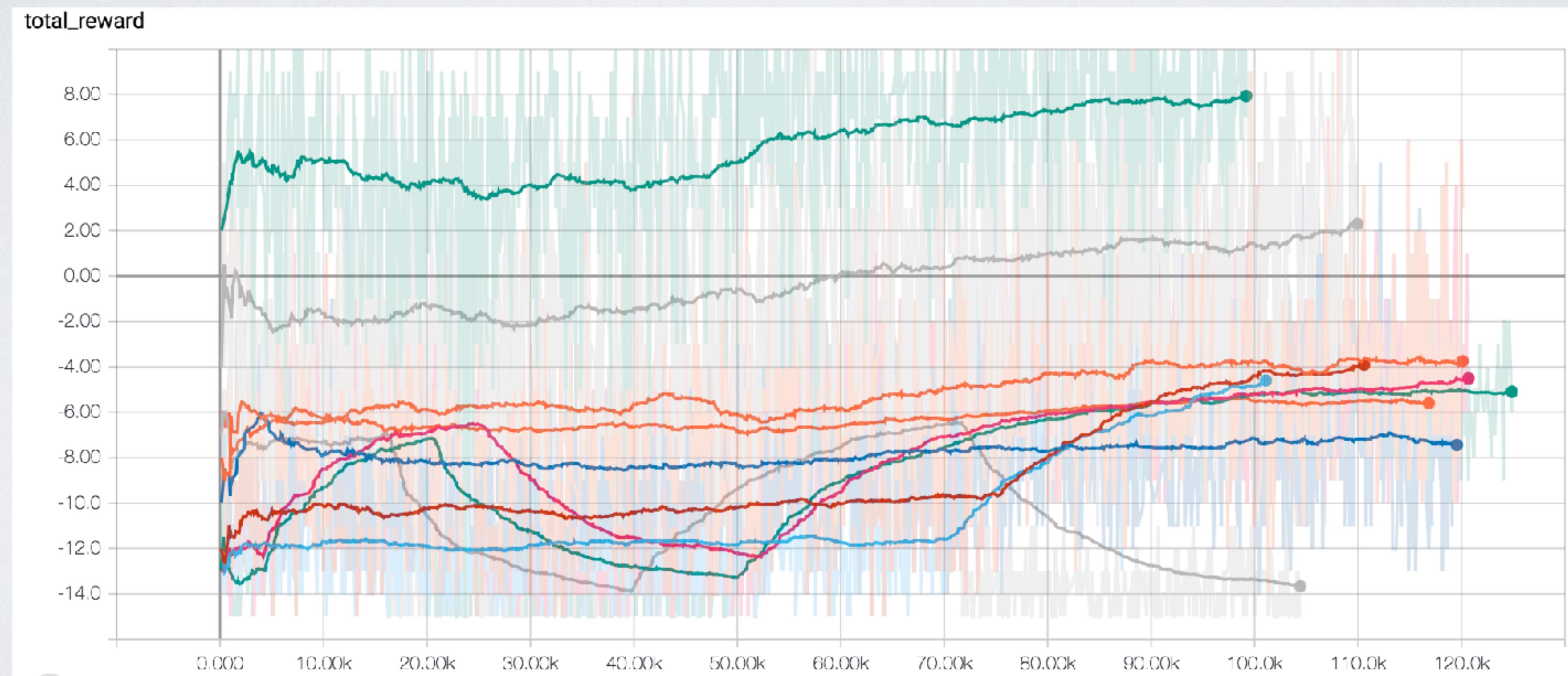


- Frame, loss



RESULT 2

- Total reward (15 - Score difference)



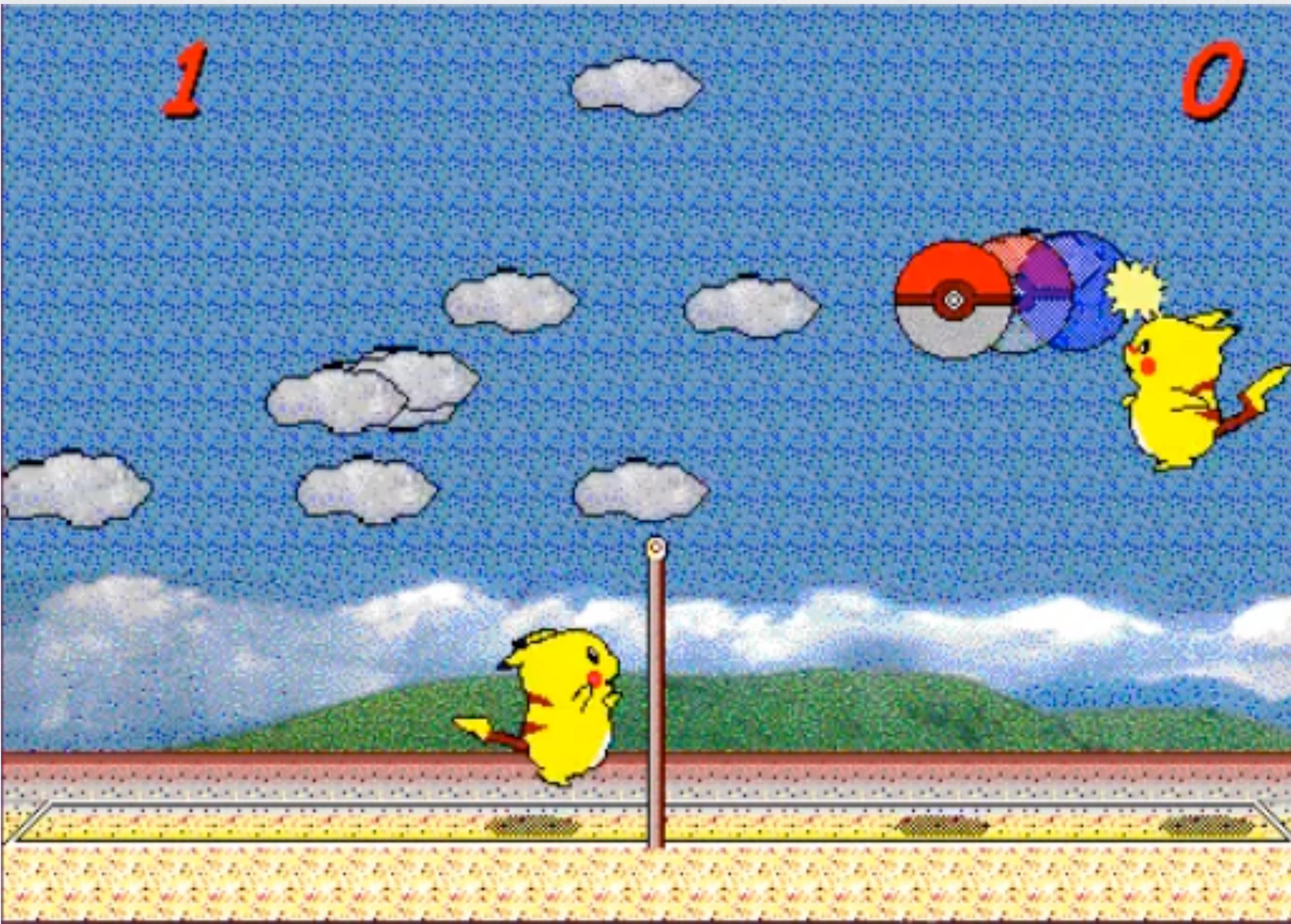
DEMO

FEATURES

- Infinite loop

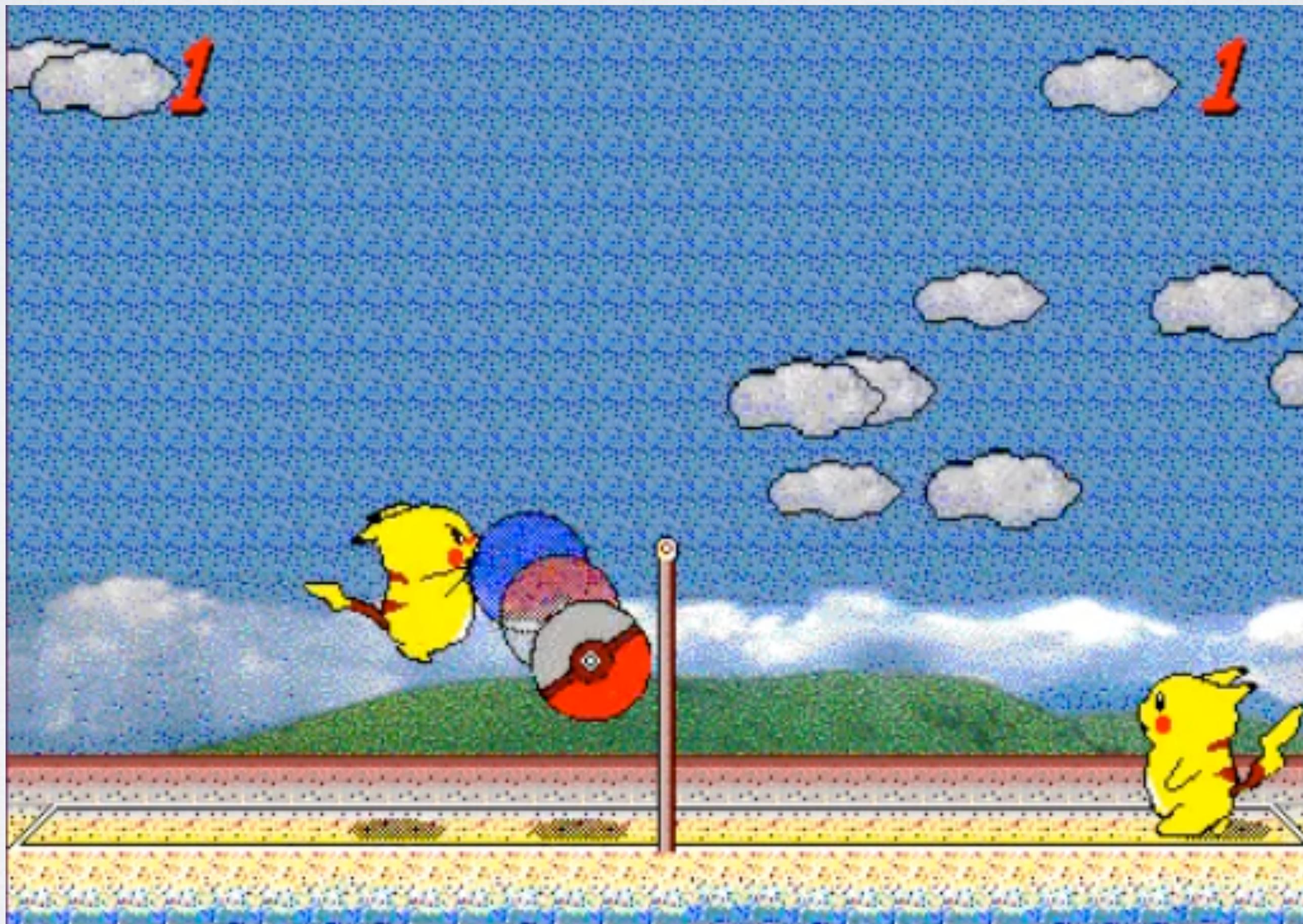
FEATURES

- Mode



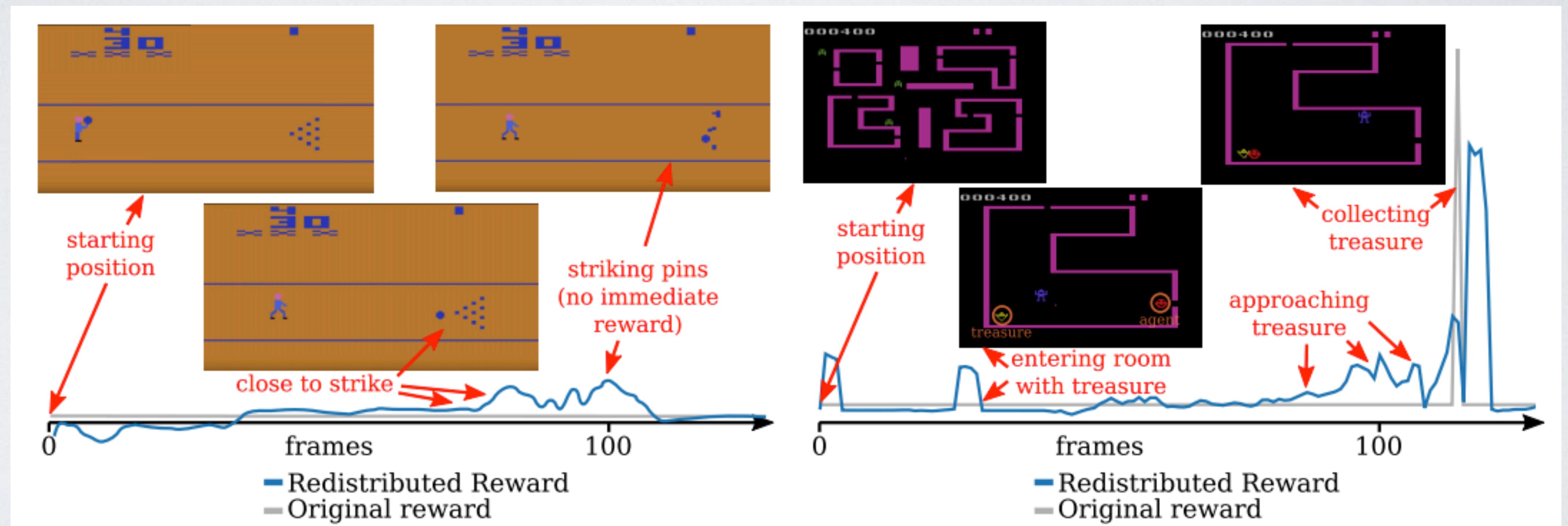
FEATURES

- Suicide



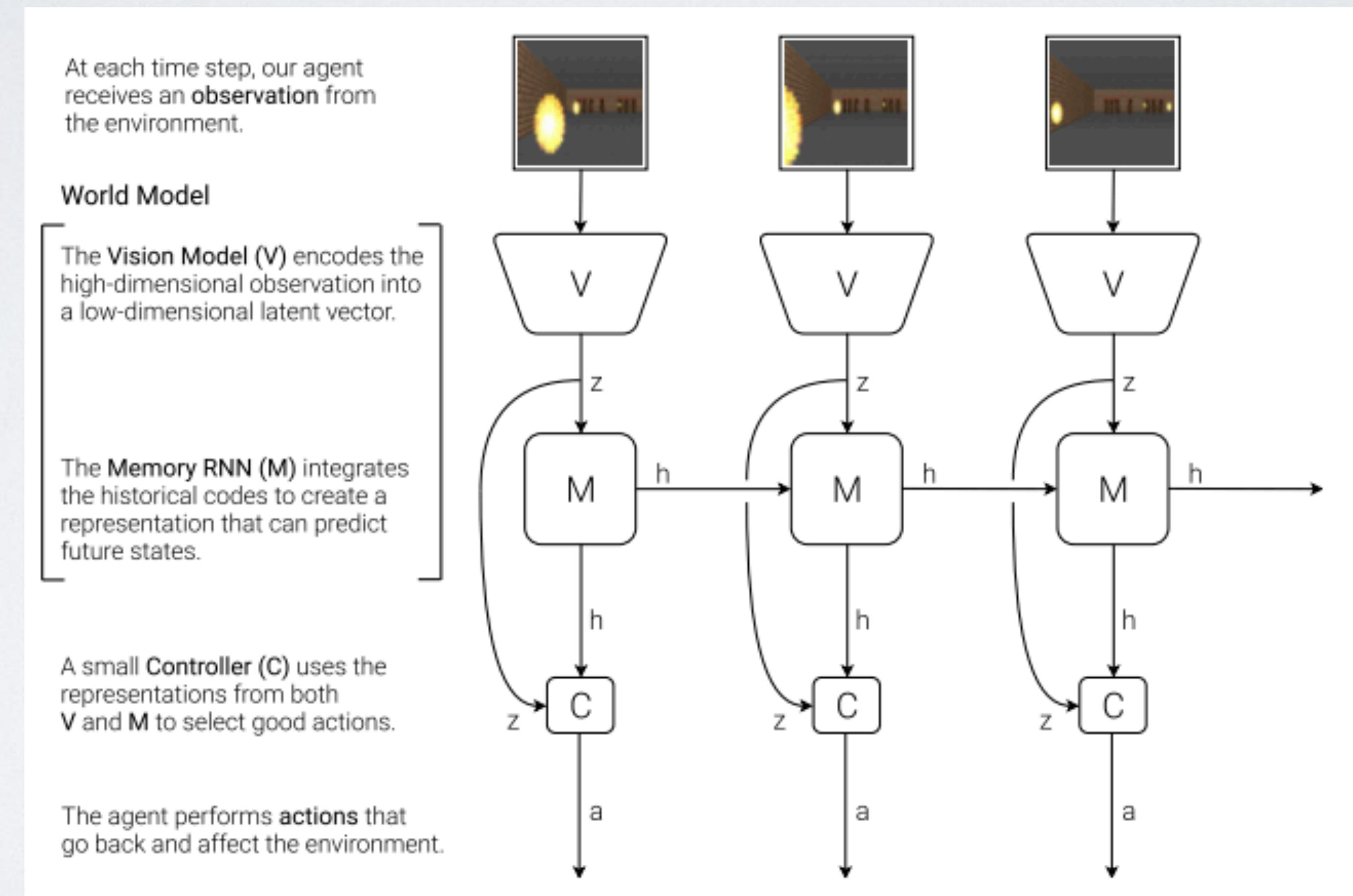
RUDDER

- RUDDER: Return Decomposition for Delayed Rewards



WORLD MODEL

- World Model



CONCLUSION

- Reinforcement learning can be very inefficient
 - Do everything before training
 - Keep the states / action size small as possible
- Parallel training is necessary
 - Model size matters
 - Hardware matters
- Lots of drawbacks in reinforcement learning
 - Opportunity?
 - Why?

SUMMARY

- Reinforcement learning basics
 - Value approximation function
 - Policy gradient
 - Temporal difference
 - Monte Carlo method
 - DQN
- DQN extensions
 - Double DQN
 - Dueling DQN
 - Prioritized experience replay
- Parallel training
 - Gorila
 - Ape-X
- Recent models
 - RUDDER
 - World model
- Engineering Tips