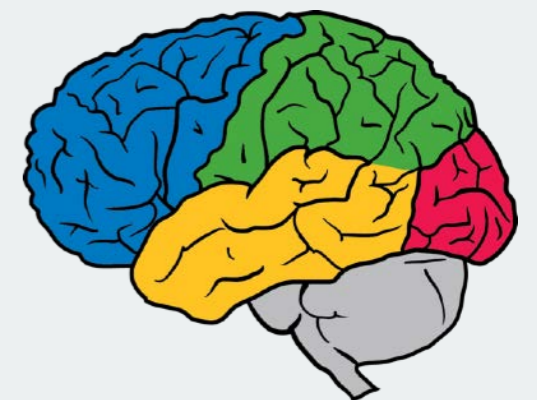
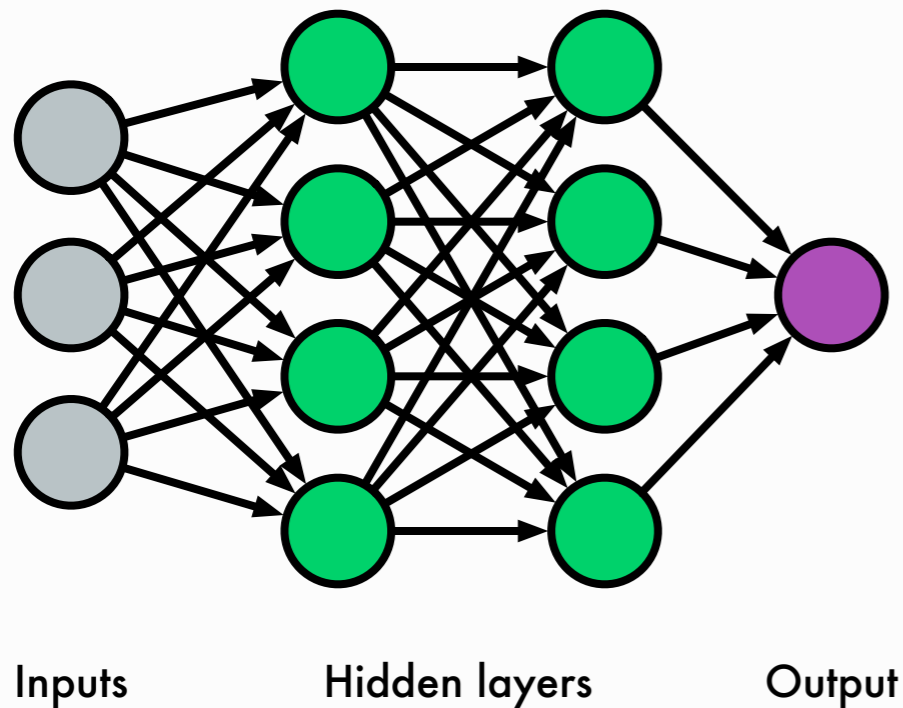


Deep Reinforcement Learning and the Atari 2600

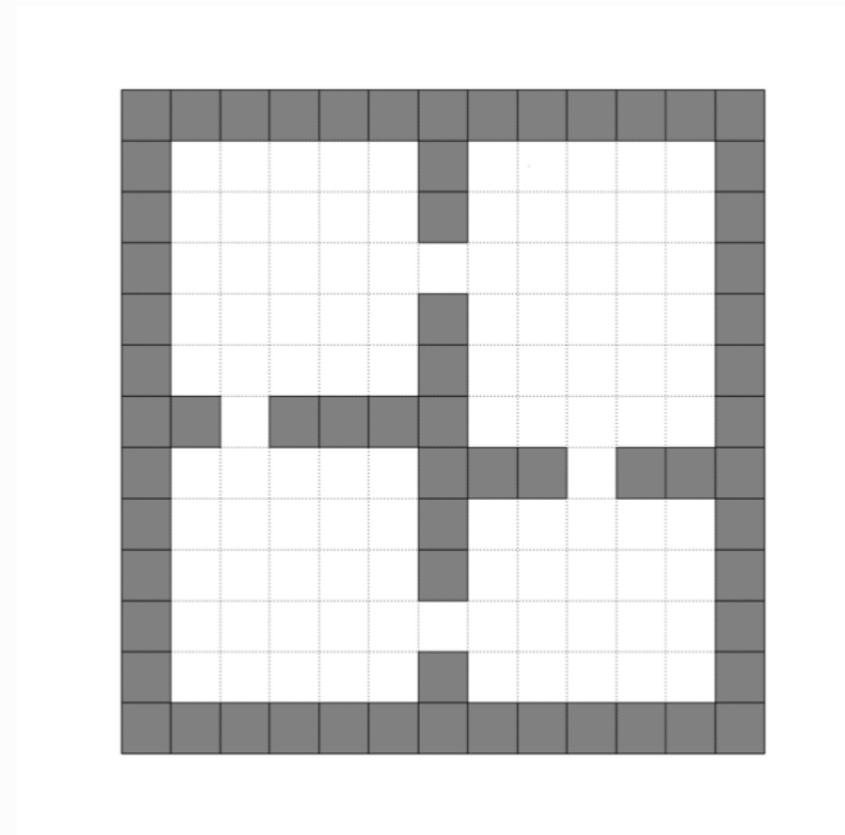
MARC G. BELLEMARE
Google Brain



DEEP



REINFORCEMENT LEARNING



Backgammon (Tesauro, 1994)

Elevator control (Crites and Barto, 1996)

Helicopter control (Ng et al., 2003)

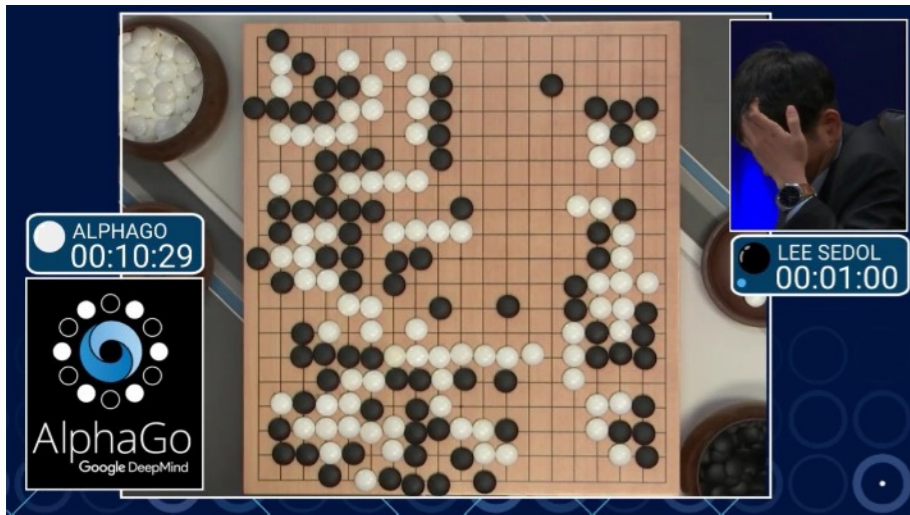
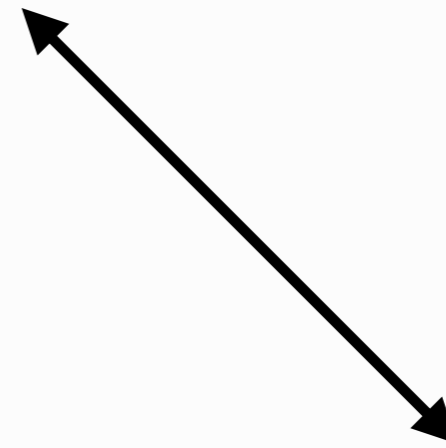
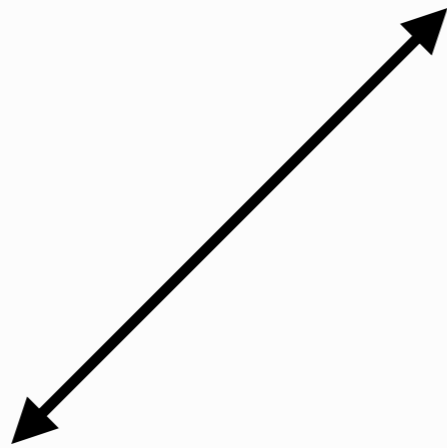
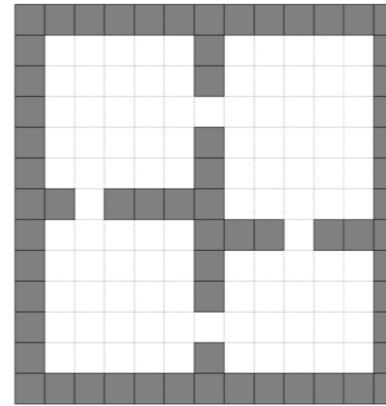
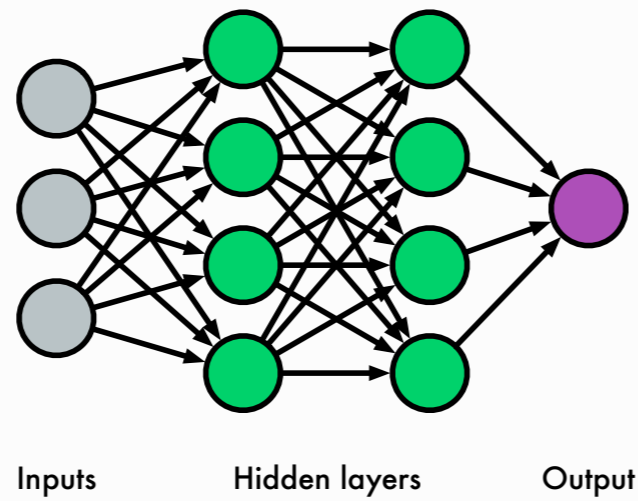
Transfer in tic-tac-toe (Rivest and Precup, 2003; Bellemare et al., unpublished)

...

Atari game-playing (Mnih et al., 2015)

Go (Silver et al., 2016)

Deep RL



Require: $\eta_t^{(x,a)} = \sum_{k=1}^{\infty} p_{t,k} \gamma^{k-1} r_{z_k}$ for each (x, a)

- 1: Sample transition (x_t, a_t, r_t, x_{t+1})
- 2: # Compute distributional Bellman target
- 3: **if** Categorical SARSA **then**
- 4: $a^* \sim \pi(\cdot | x_{t+1})$
- 5: **else if** Categorical Q-learning **then**
- 6: $a^* \leftarrow \arg \max_a \mathbb{E}_{R \sim \eta_t^{(x_{t+1}, a)}} [R]$
- 7: **end if**
- 8: $\hat{\eta}_*^{(x_t, a_t)} \leftarrow (f_{r_t, \gamma}) \# \eta_t^{(x_{t+1}, a^*)}$
- 9: # Project target onto support
- 10: $\hat{\eta}_t^{(x_t, a_t)} \leftarrow \Pi_{\mathcal{C}} \hat{\eta}_*^{(x_t, a_t)}$
- 11: # Compute KL Loss
- 12: Find gradient of $\text{KL}(\hat{\eta}_t^{(x_t, a_t)} || \eta_t^{(x_t, a_t)})$
- 13: Use gradient to generate new estimate

Challenge domains

Algorithms

SOME CHALLENGES IN DEEP RL

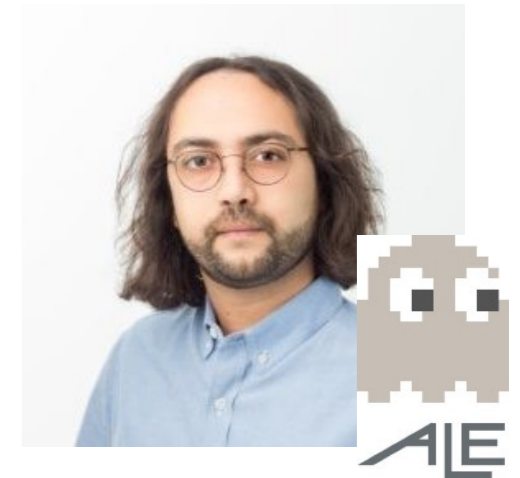
- Training data is **not i.i.d.**
- Hard to get **confidence intervals**
- Potential divergence issues (e.g. van Hasselt et al., 2015)
- State space often unknown
- **Model-free methods reign supreme** (simulation a plus)



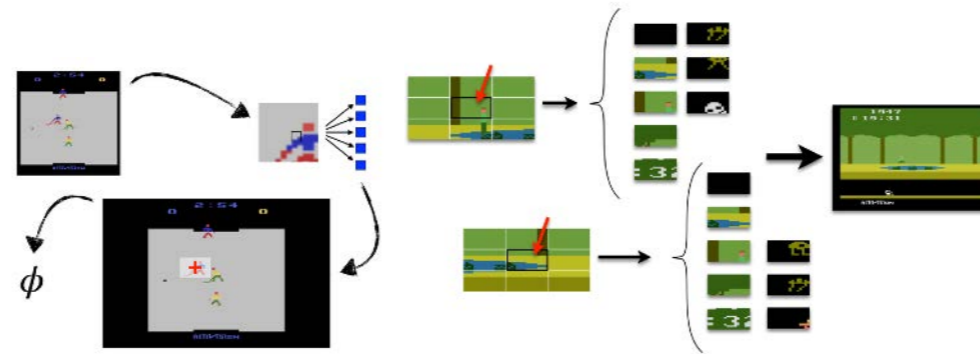
Diuk *et al.*, 2008



Barbados RL Workshop, 2008



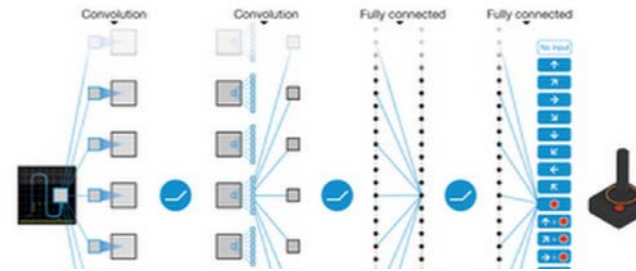
Naddaf, 2008



Bellemare *et al.*, 2011-2013



MGB, Naddaf,
Veness, Bowling 2013



Mnih *et al.*, 2013, 2015

Hausknecht *et al.* 2013

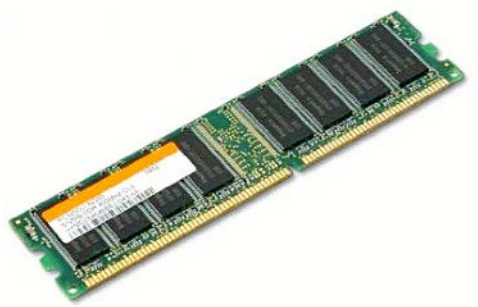
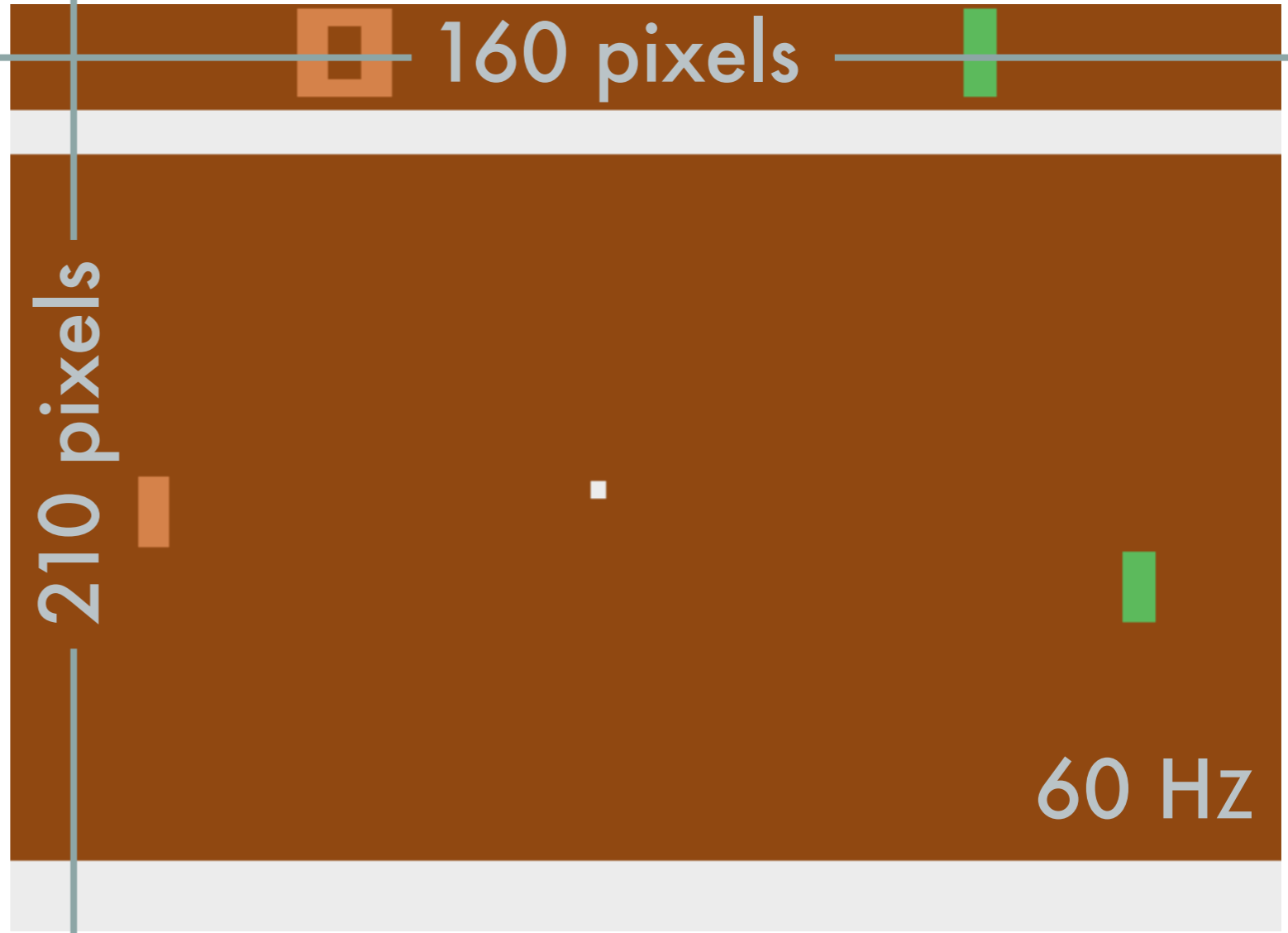
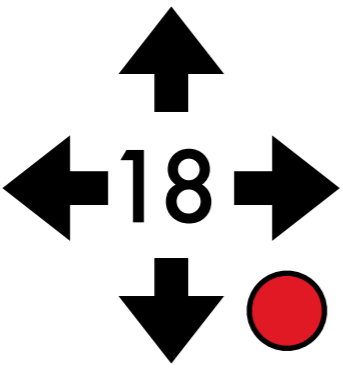
Lipovetsky *et al.* 2015

...

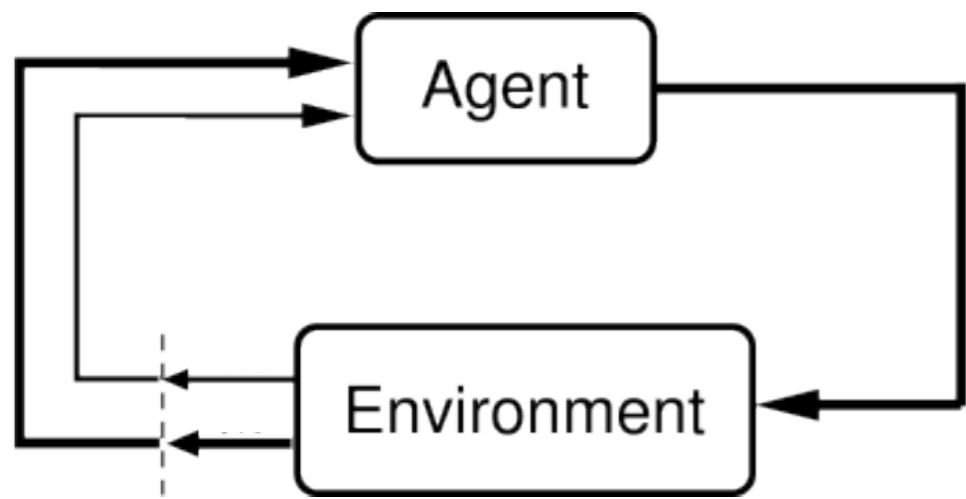




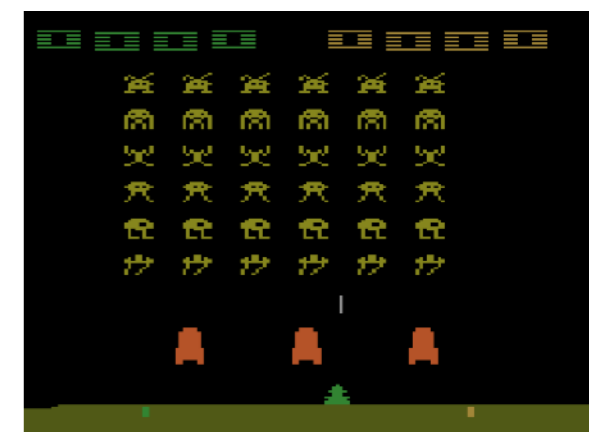
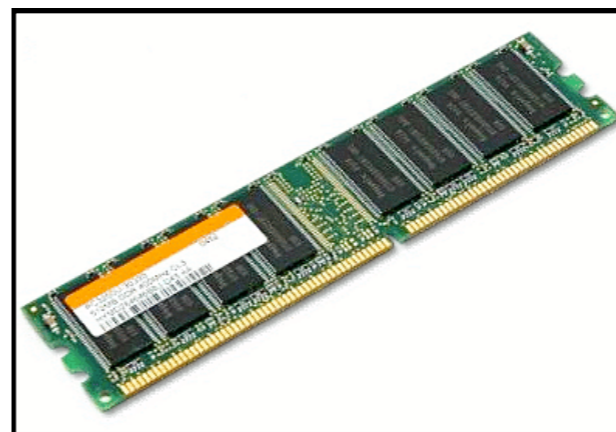
1977



128 BYTES

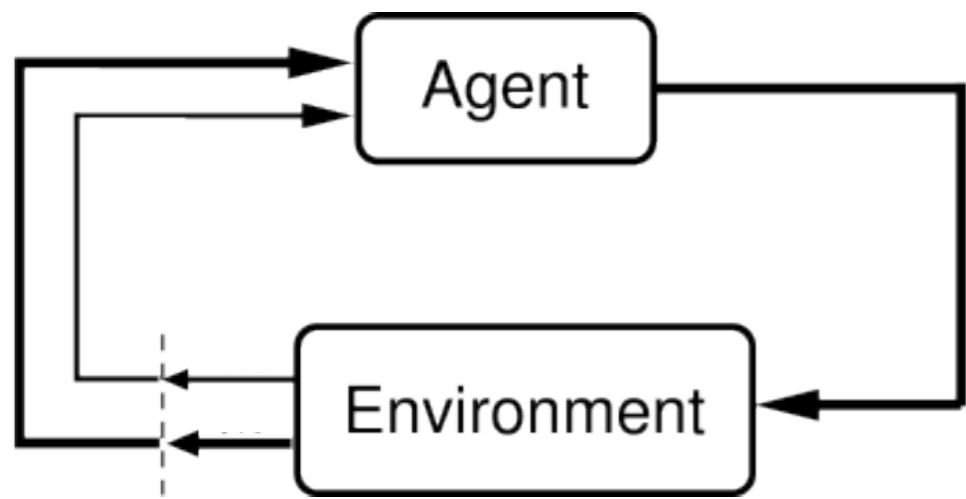


Observation

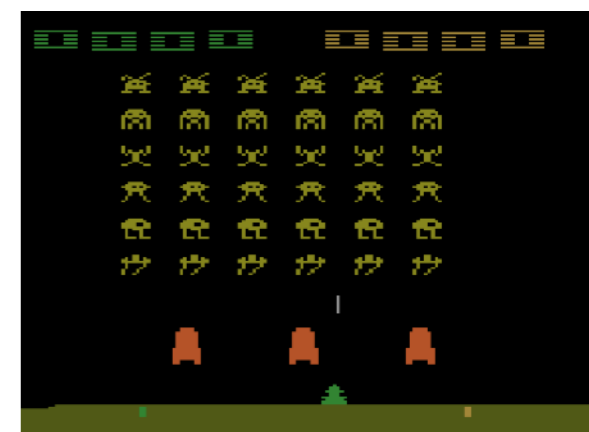
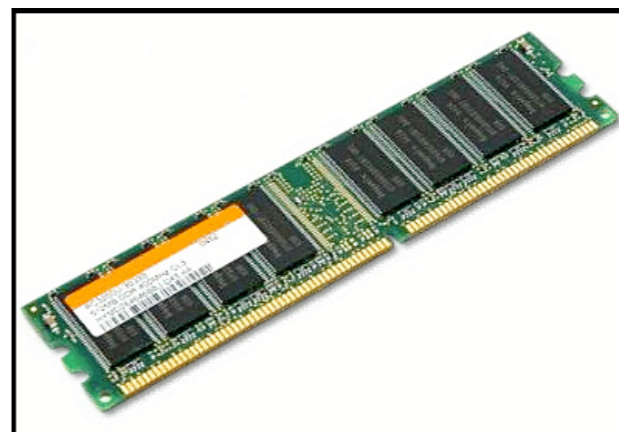


Action





Observation

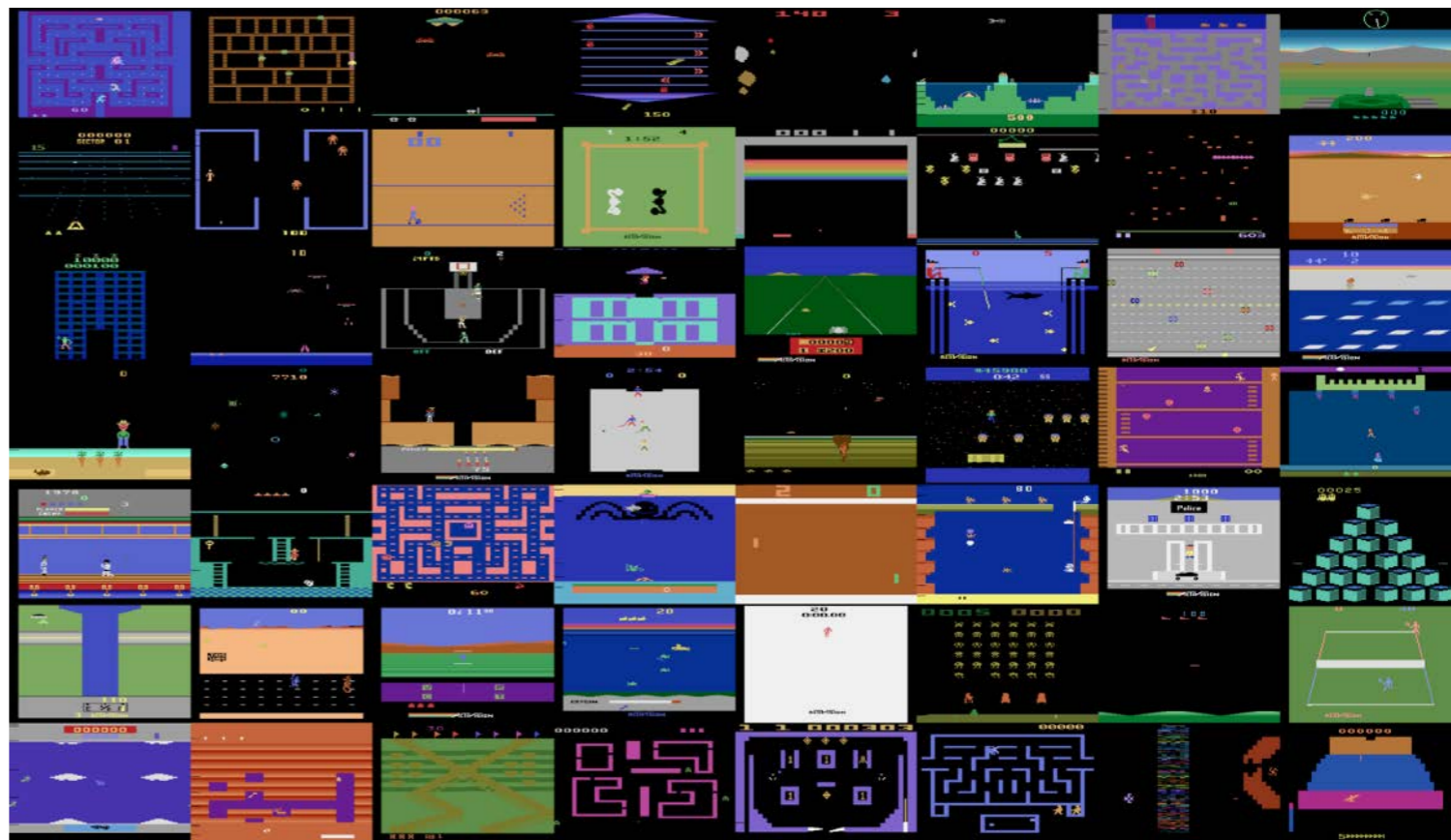
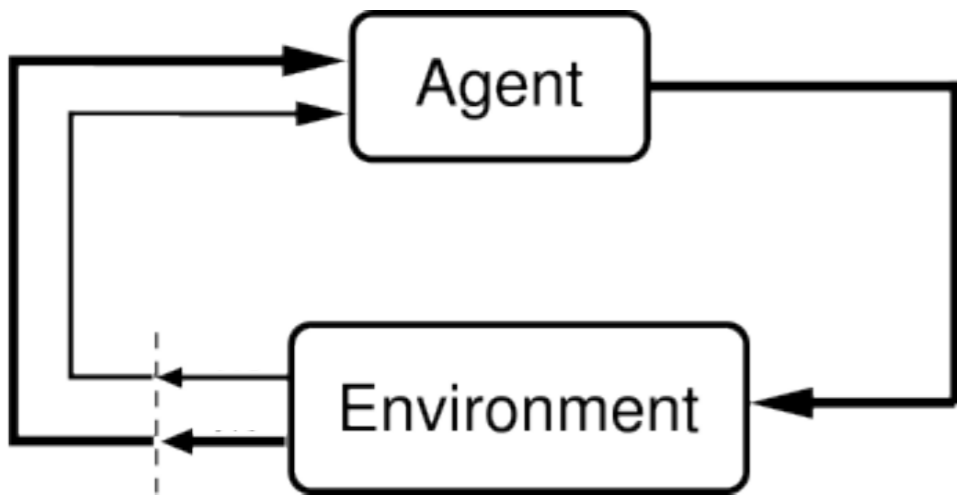


Action



Reward



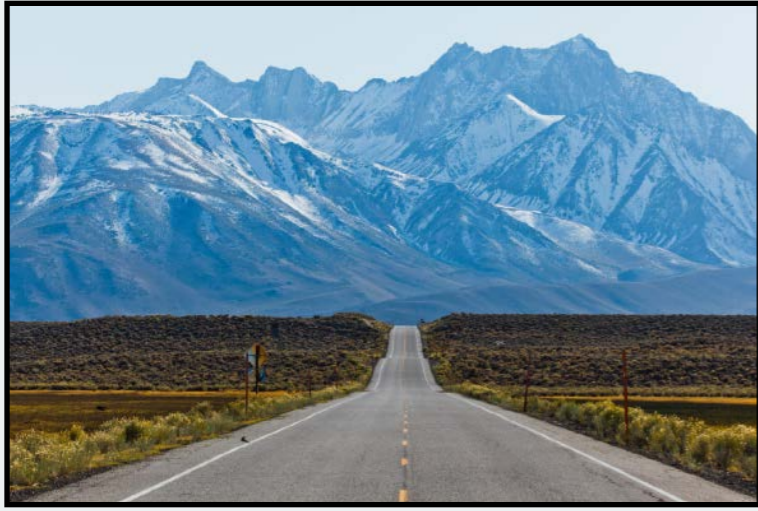




GENERAL COMPETENCY

```
// Returns the vector of the minimal set of actions needed to play  
// the game.  
ActionVect ALEInterface::getMinimalActionSet() {  
    if (!romSettings.get()) {  
        throw std::runtime_error("ROM not set");  
    }  
    return romSettings->getMinimalActionSet();  
}  
  
// Returns the frame number since the loading of the ROM  
int ALEInterface::getFrameNumber() {  
    return environment->getFrameNumber();  
}  
  
// Returns the frame number since the start of the current episode  
int ALEInterface::getEpisodeFrameNumber() const {  
    return environment->getEpisodeFrameNumber();  
}  
  
// Returns the current game screen  
ALEScreen& ALEInterface::getScreen() {  
    return environment->getScreen();  
}
```





```

// Returns the vector of the minimal set of actions needed to play
// the game.
ActionVect ALEInterface::getMinimalActionSet() {
    if (!IronSettings.get()){
        throw std::runtime_error("ROM not set");
    }
    return IronSettings->getMinimalActionSet();
}

// Returns the frame number since the loading of the ROM
int ALEInterface::getFrameNumber() {
    return environment->getFrameNumber();
}

// Returns the frame number since the start of the current episode
int ALEInterface::getEpisodeFrameNumber() const {
    return environment->getEpisodeFrameNumber();
}

// Returns the current game screen
ALEScreen& ALEInterface::getScreen() {
    return environment->getScreen();
}

```

```

// Returns the vector of the minimal set of actions needed to play
// the game.
ActionVect ALEInterface::getMinimalActionSet() {
    if (!IronSettings.get()){
        throw std::runtime_error("ROM not set");
    }
    return IronSettings->getMinimalActionSet();
}

// Returns the frame number since the loading of the ROM
int ALEInterface::getFrameNumber() {
    return environment->getFrameNumber();
}

// Returns the frame number since the start of the current episode
int ALEInterface::getEpisodeFrameNumber() const {
    return environment->getEpisodeFrameNumber();
}

// Returns the current game screen
ALEScreen& ALEInterface::getScreen() {
    return environment->getScreen();
}

```

```

// Returns the vector of the minimal set of actions needed to play
// the game.
ActionVect ALEInterface::getMinimalActionSet() {
    if (!IronSettings.get()){
        throw std::runtime_error("ROM not set");
    }
    return IronSettings->getMinimalActionSet();
}

// Returns the frame number since the loading of the ROM
int ALEInterface::getFrameNumber() {
    return environment->getFrameNumber();
}

// Returns the frame number since the start of the current episode
int ALEInterface::getEpisodeFrameNumber() const {
    return environment->getEpisodeFrameNumber();
}

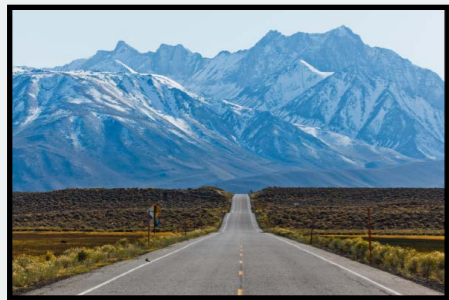
// Returns the current game screen
ALEScreen& ALEInterface::getScreen() {
    return environment->getScreen();
}

```



NARROW COMPETENCY

Diverse



```
/* Returns the vector of the initial set of actions needed to play
 * the game.
 */
virtual Action* getInitialActions() {
    if (!m_initialized) {
        show_status_bar_error("RM not set");
        return m_settings->getInitialActions();
    }
}

/* Returns the frame number since the loading of the RM
 * at A2Interface::getFrameNumber().
 */
return environment->getFrameNumber();

/* Returns the frame number since the start of the current episode
 * at A2Interface::getEpisodeFrameNumber().
 */
return environment->getEpisodeFrameNumber();

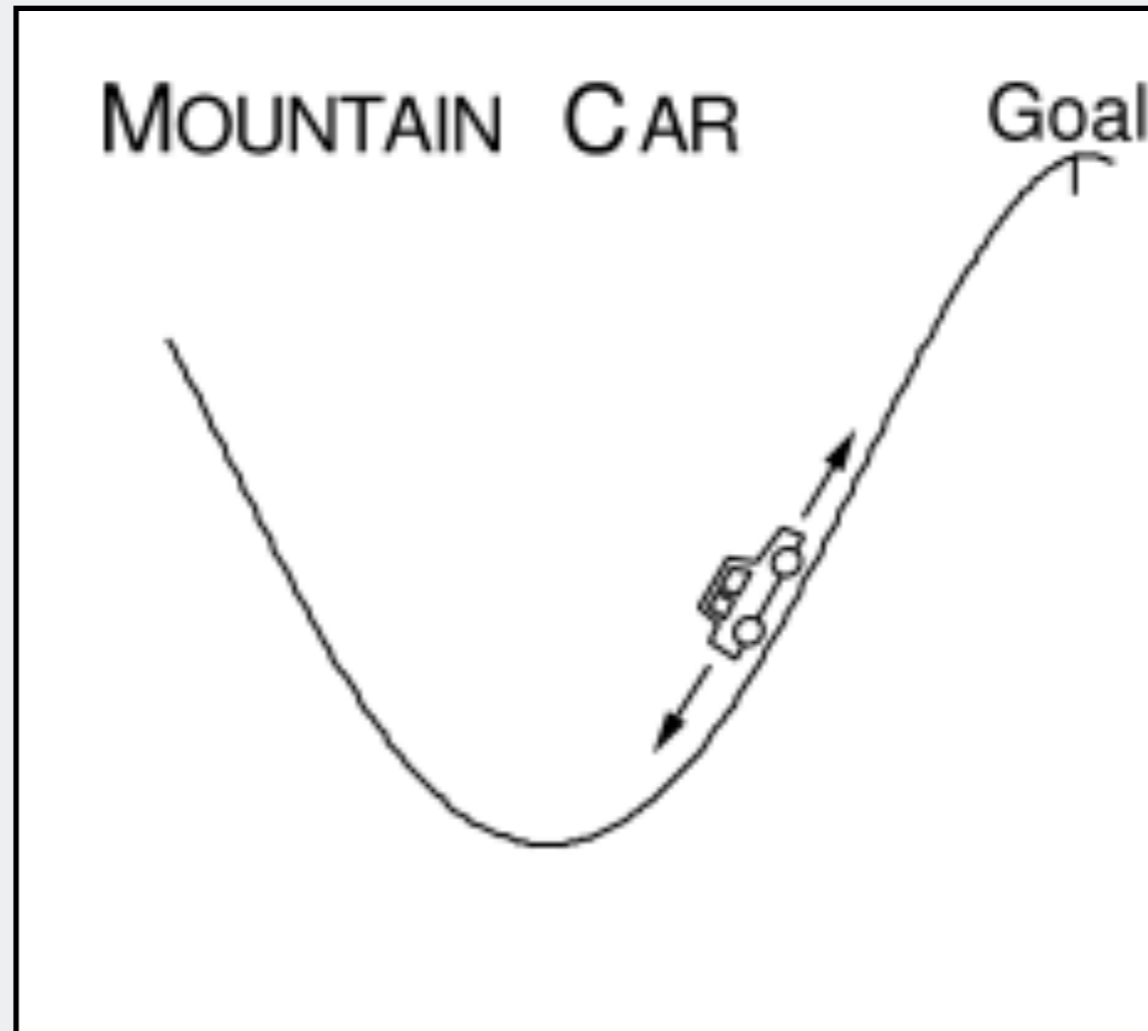
/* Returns the current game screen
 * at A2Interface::getScreen().
 */
```



Interesting

Independent

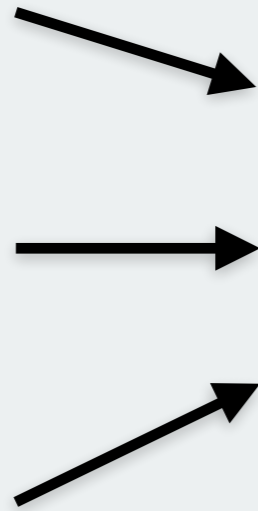
Diverse



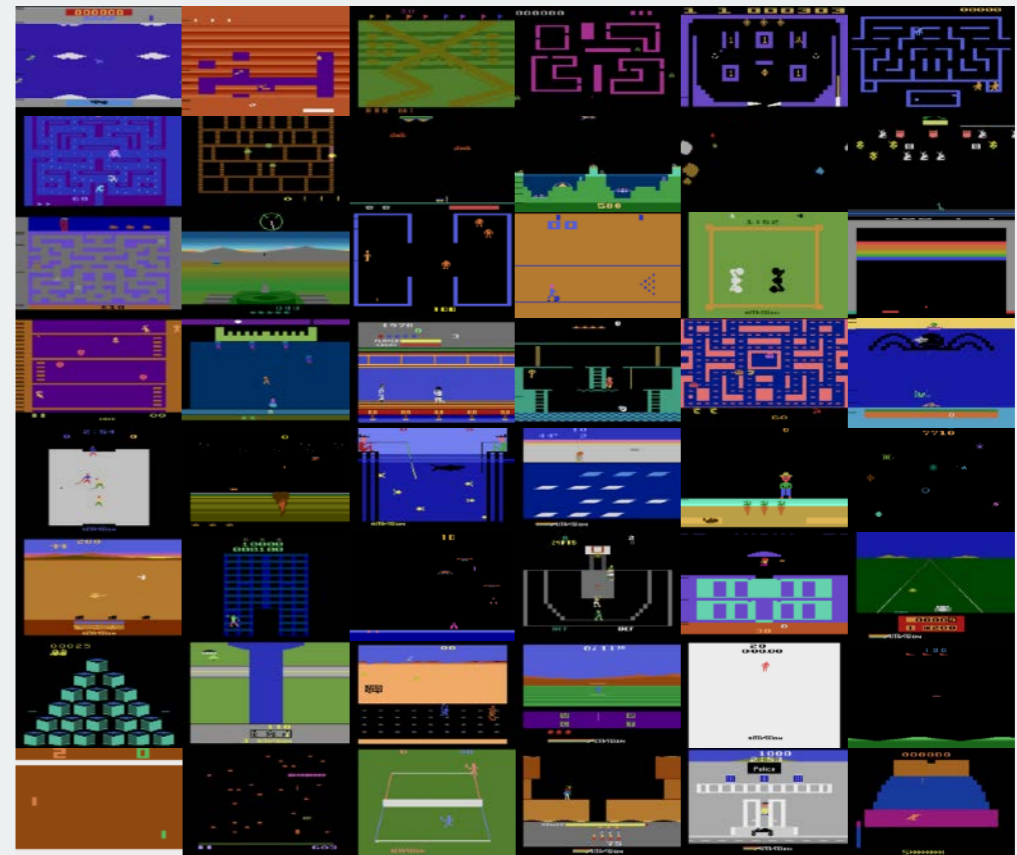
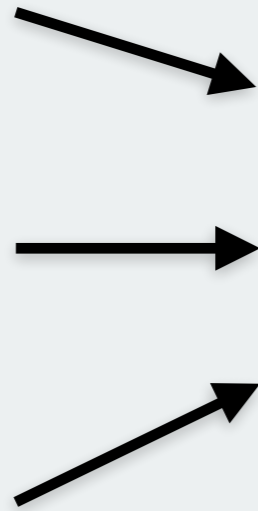
“We argue that reinforcement learning is particularly vulnerable to **environment overfitting** and propose as a remedy generalized methodologies [...]”

Whiteson et al., 2011

Diverse



Diverse



Interesting

(to people)

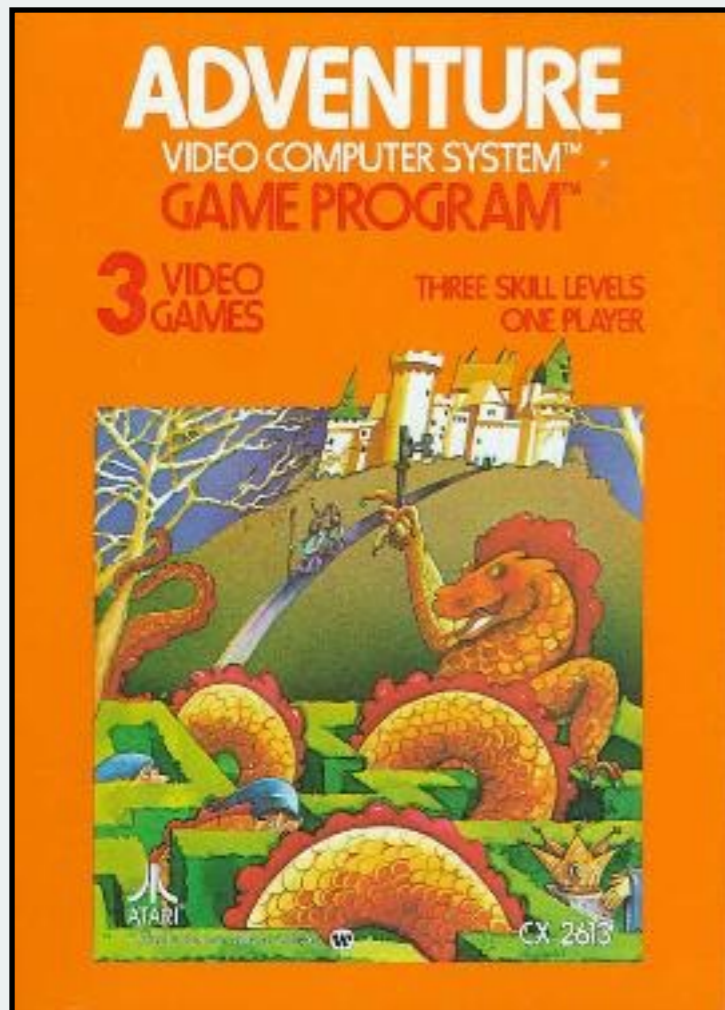


CC Wikipedia

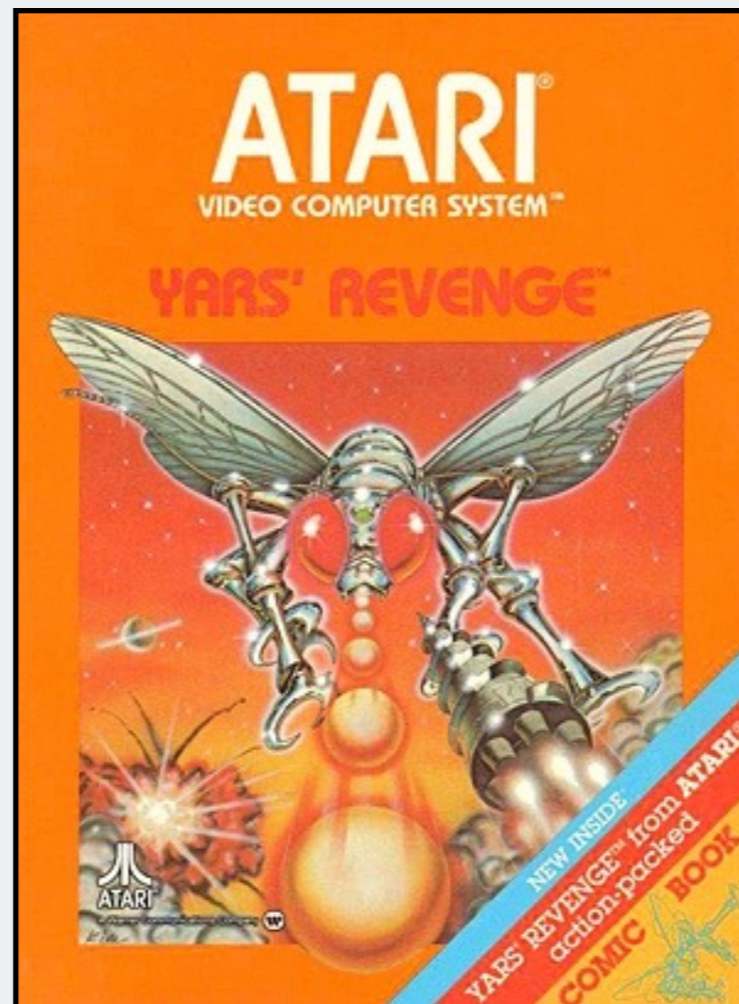


Independent

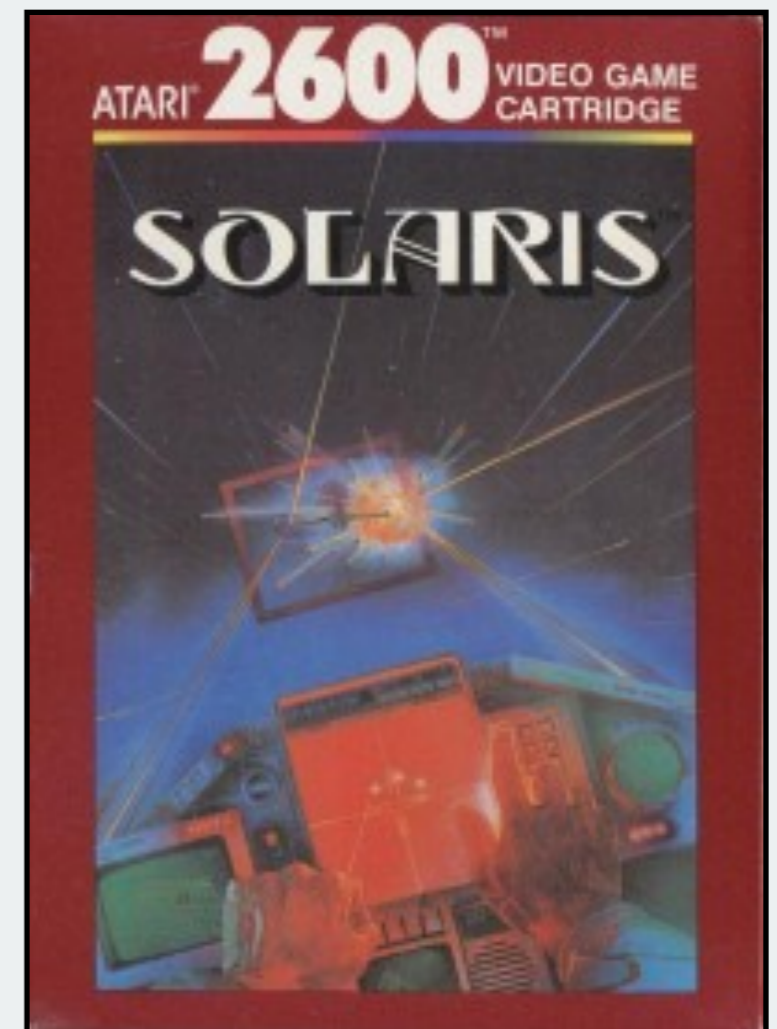
(by and for people)



1979



1982

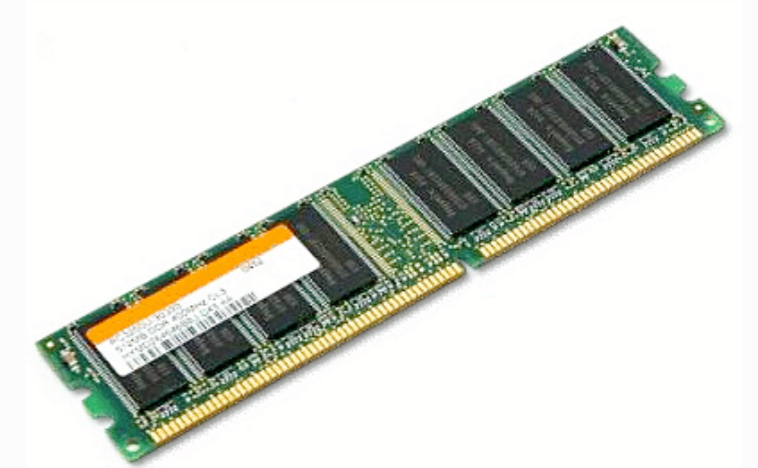
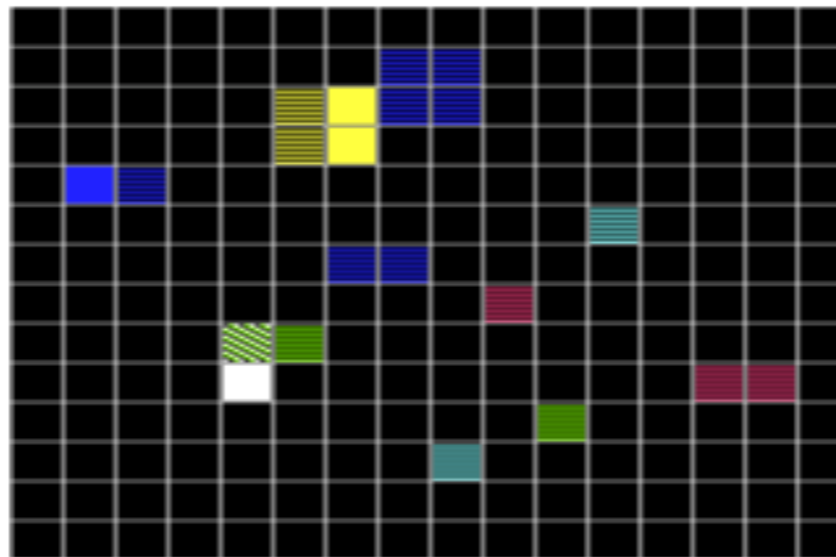


1986

EARLY ATTEMPTS (2010-2013)



BASS/Basic Features



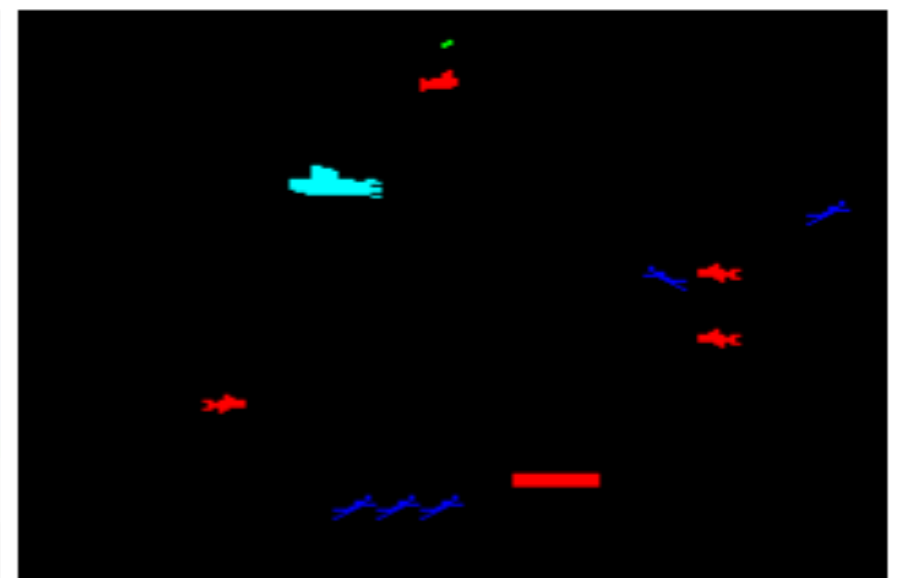
RAM



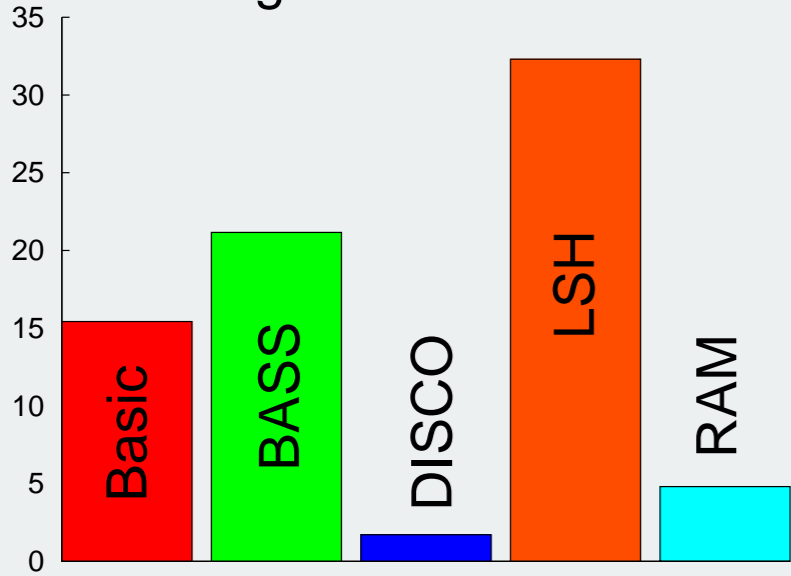
LSH on pixels



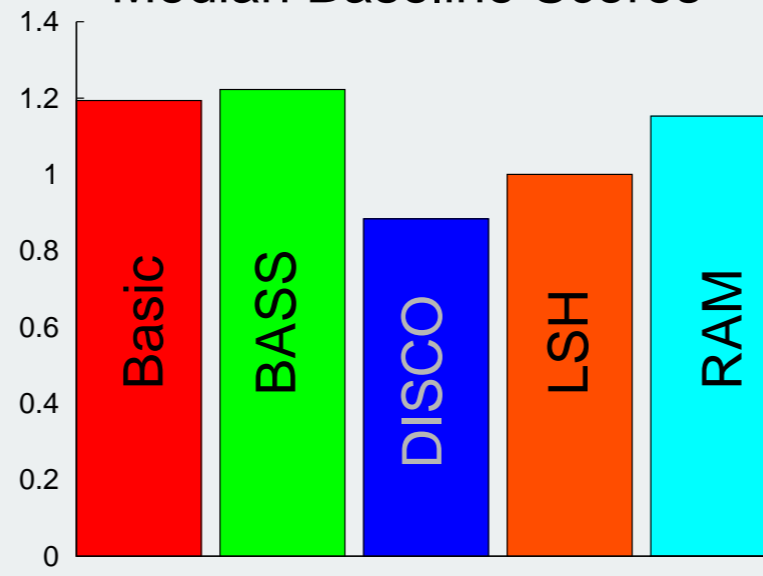
DISCO: Object Detection



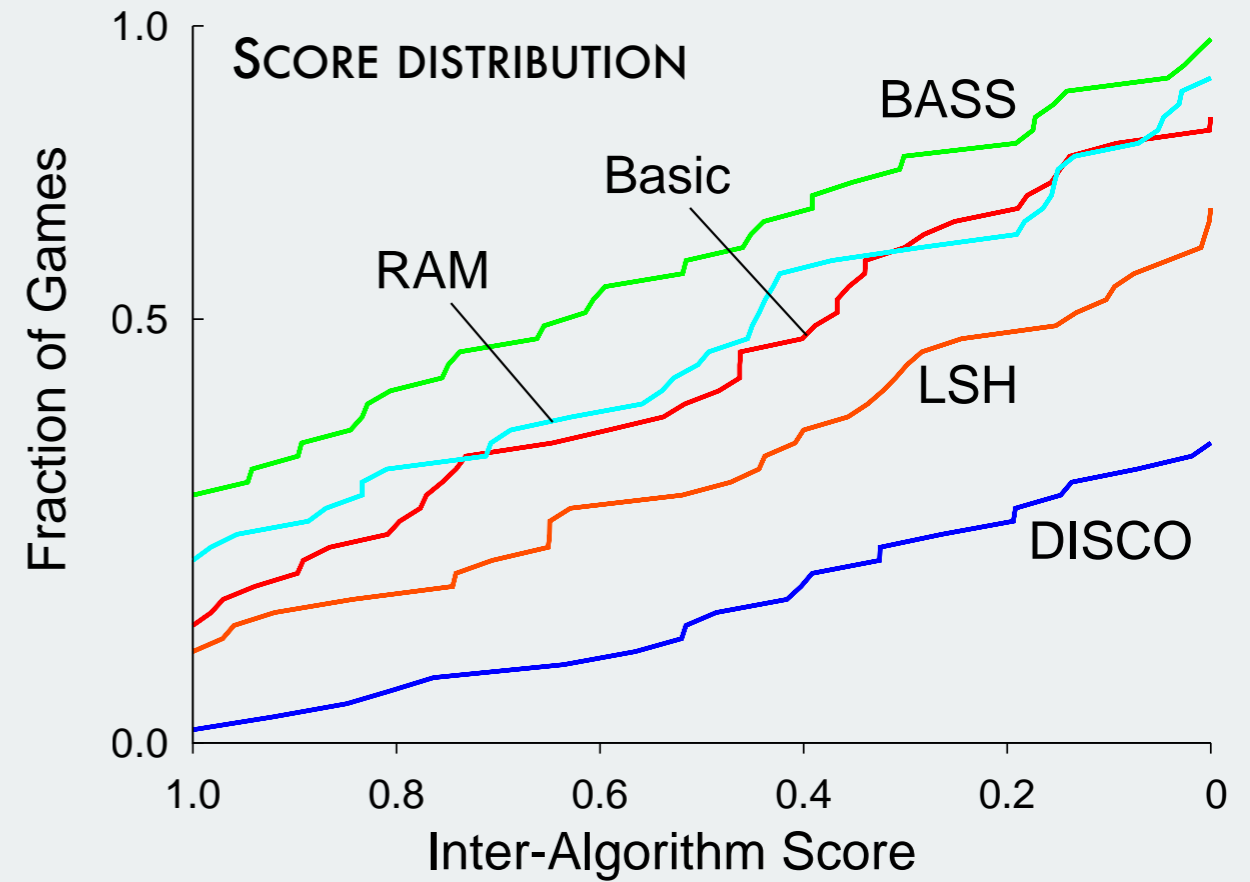
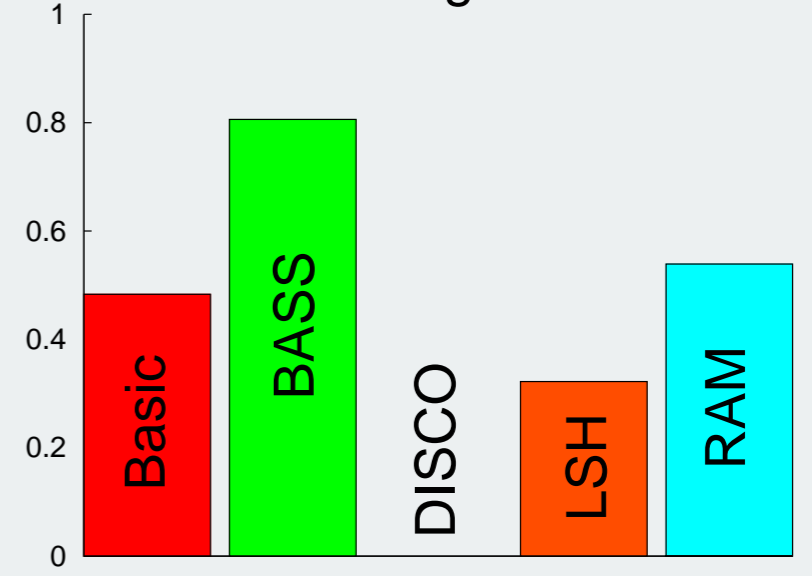
Average Baseline Scores



Median Baseline Scores

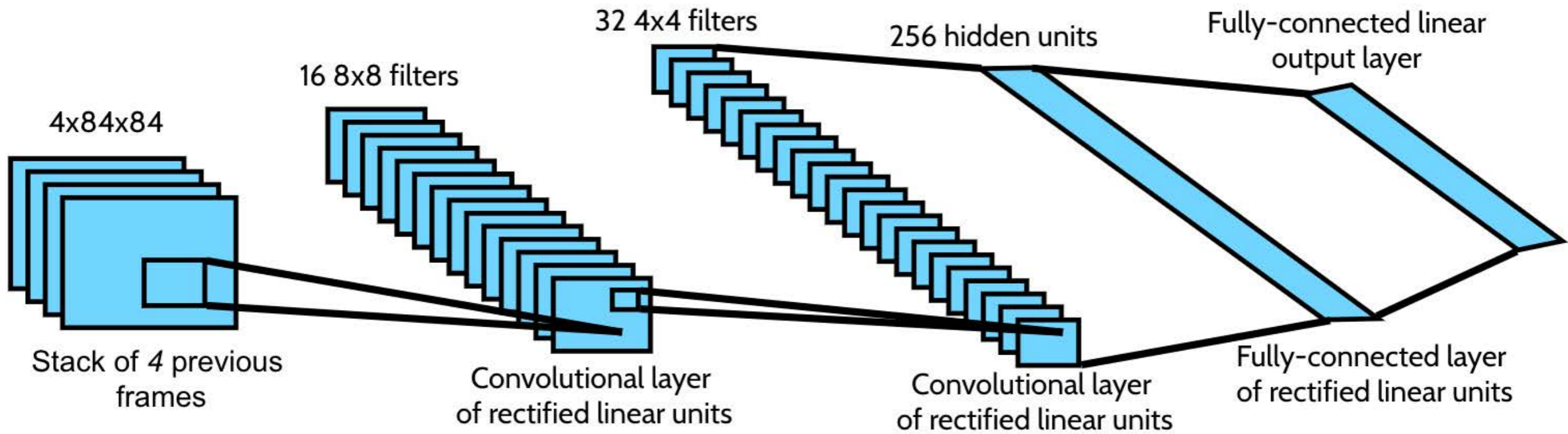


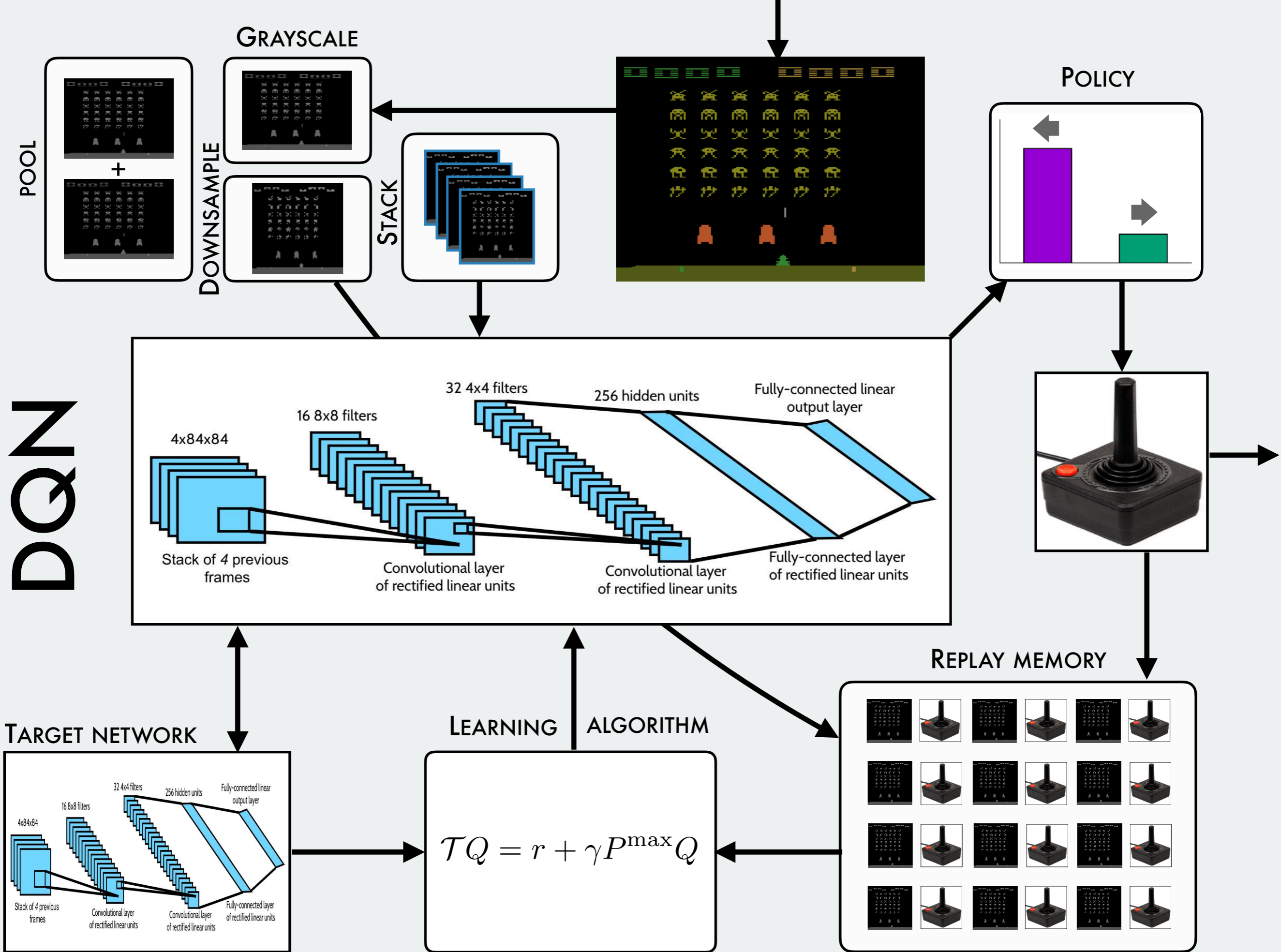
Median Inter-Algorithm Scores



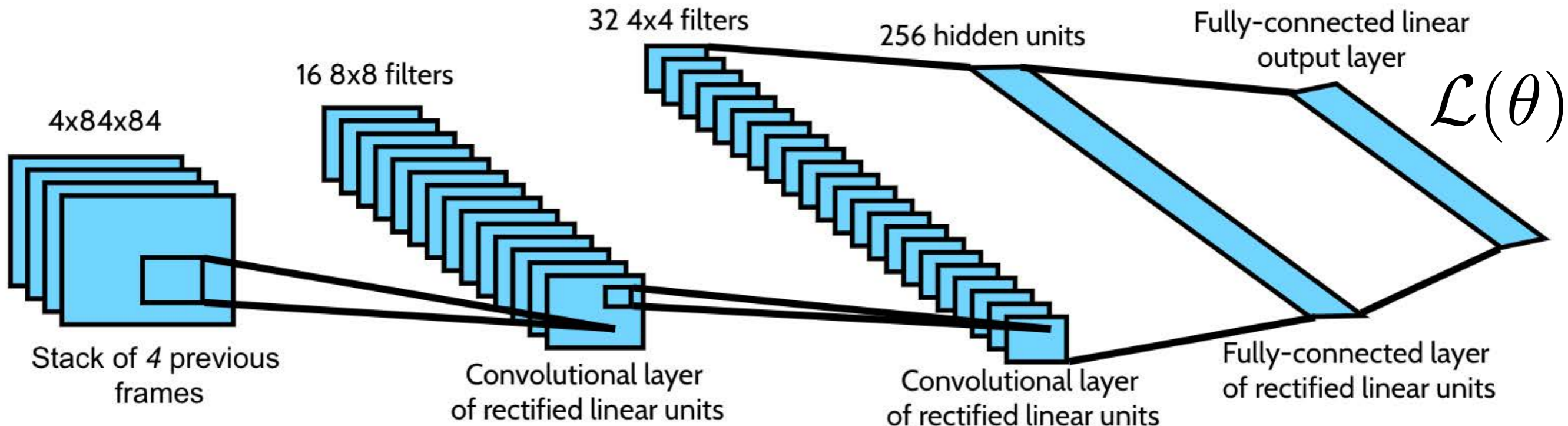
DEEP Q-NETWORKS (DQN)

Mnih et al., 2015





$$\nabla_{\theta} \mathcal{L}(\theta)$$

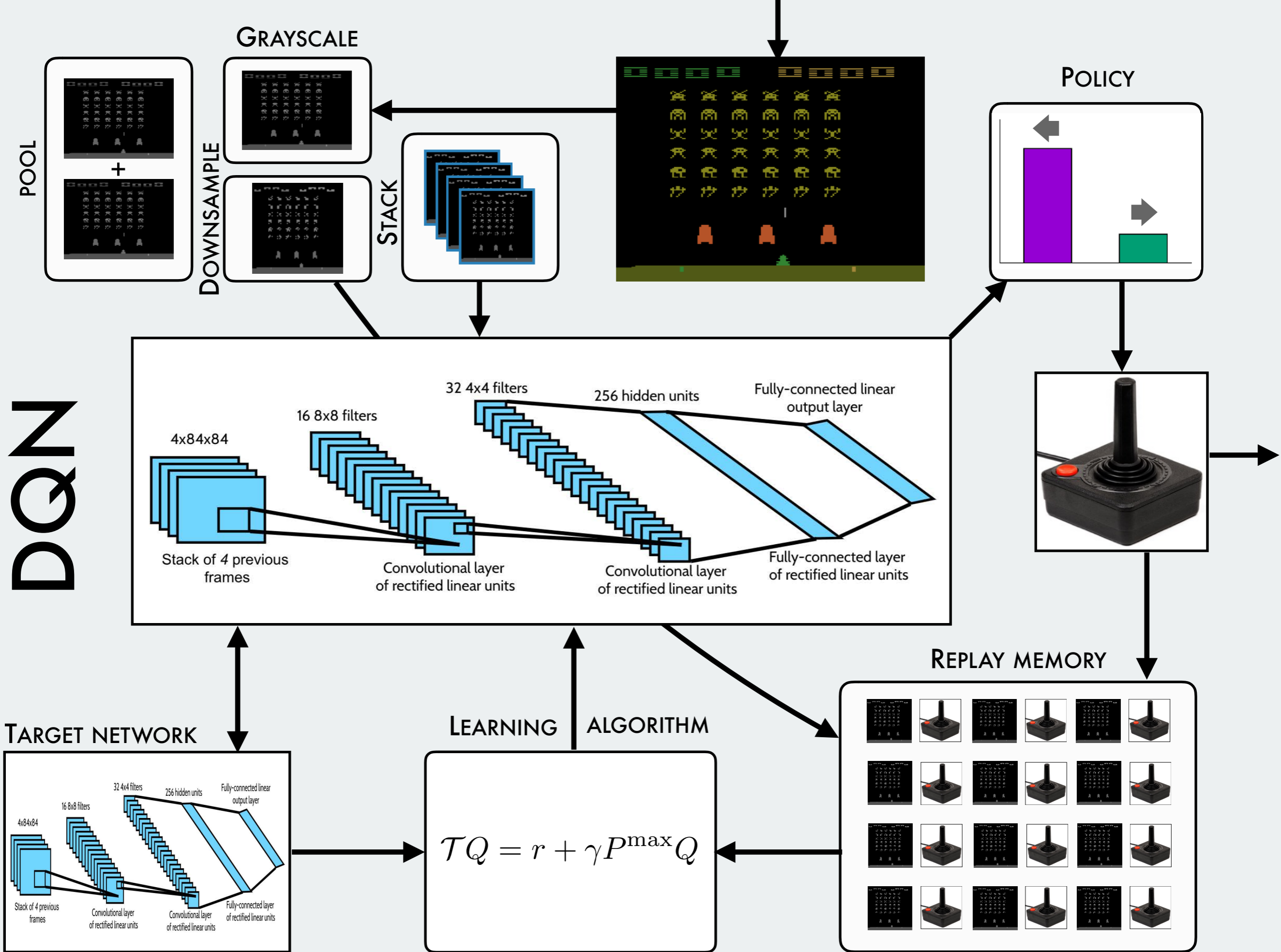


$$\mathcal{L}(\theta) = \frac{1}{2} \left(r + \gamma \max_{a' \in \mathcal{A}} Q_{\tilde{\theta}}(x', a') - Q_{\theta}(x, a) \right)^2$$

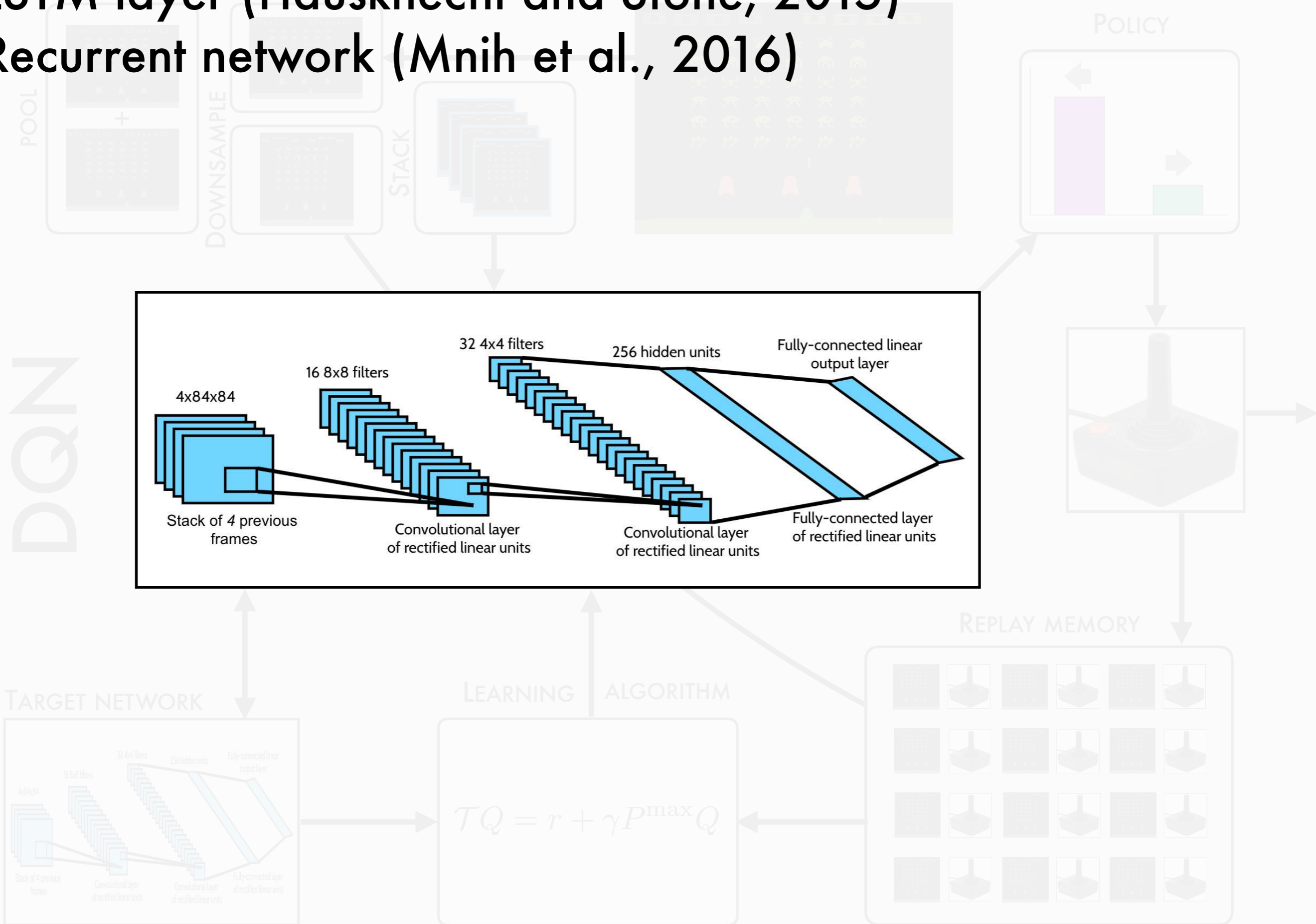
LEARNING ALGORITHM

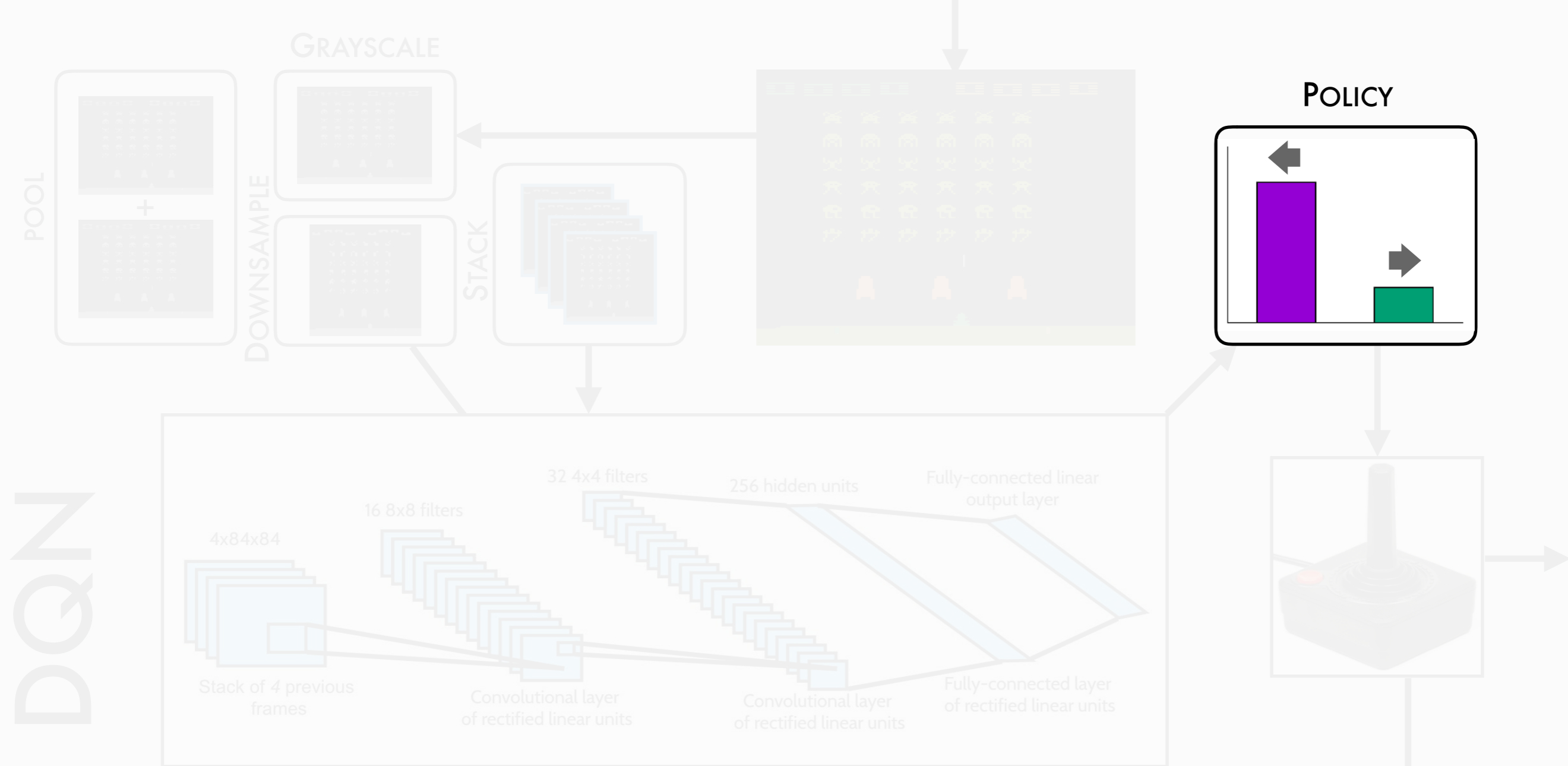
$$\mathcal{T}Q = r + \gamma P^{\max} Q$$





LSTM layer (Hausknecht and Stone, 2015) Recurrent network (Mnih et al., 2016)





Dueling networks (Wang et al., 2016)

Asynchronous actor-critic (Mnih et al., 2016)

Bootstrap Q-functions (Osband et al., 2016)

Noisy networks (Fortunato et al.; Plappert et al., 2017)

LEARNING ALGORITHM

$$TQ = r + \gamma P^{\max} Q$$

TARGET NETWORK

REPLAY MEMORY

Stack of 4 previous frames

Convolutional layer of rectified linear units

Convolutional layer of rectified linear units

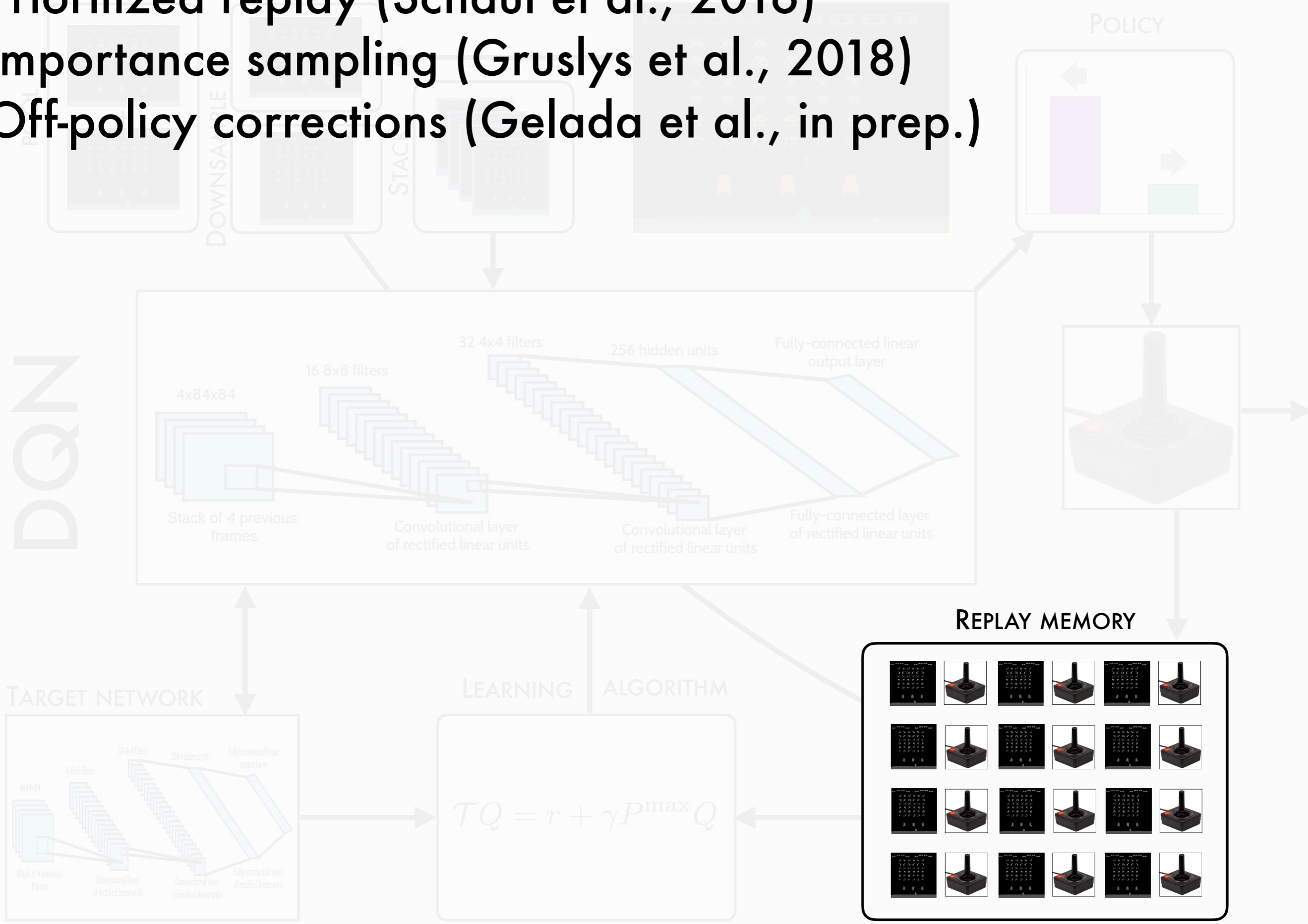
Fully-connected layer of rectified linear units

256 hidden units

Fully-connected linear output layer

Prioritized replay (Schaul et al., 2016)
Importance sampling (Gruslys et al., 2018)
Off-policy corrections (Gelada et al., in prep.)

DQN



Double Q-learning (van Hasselt et al., 2015)

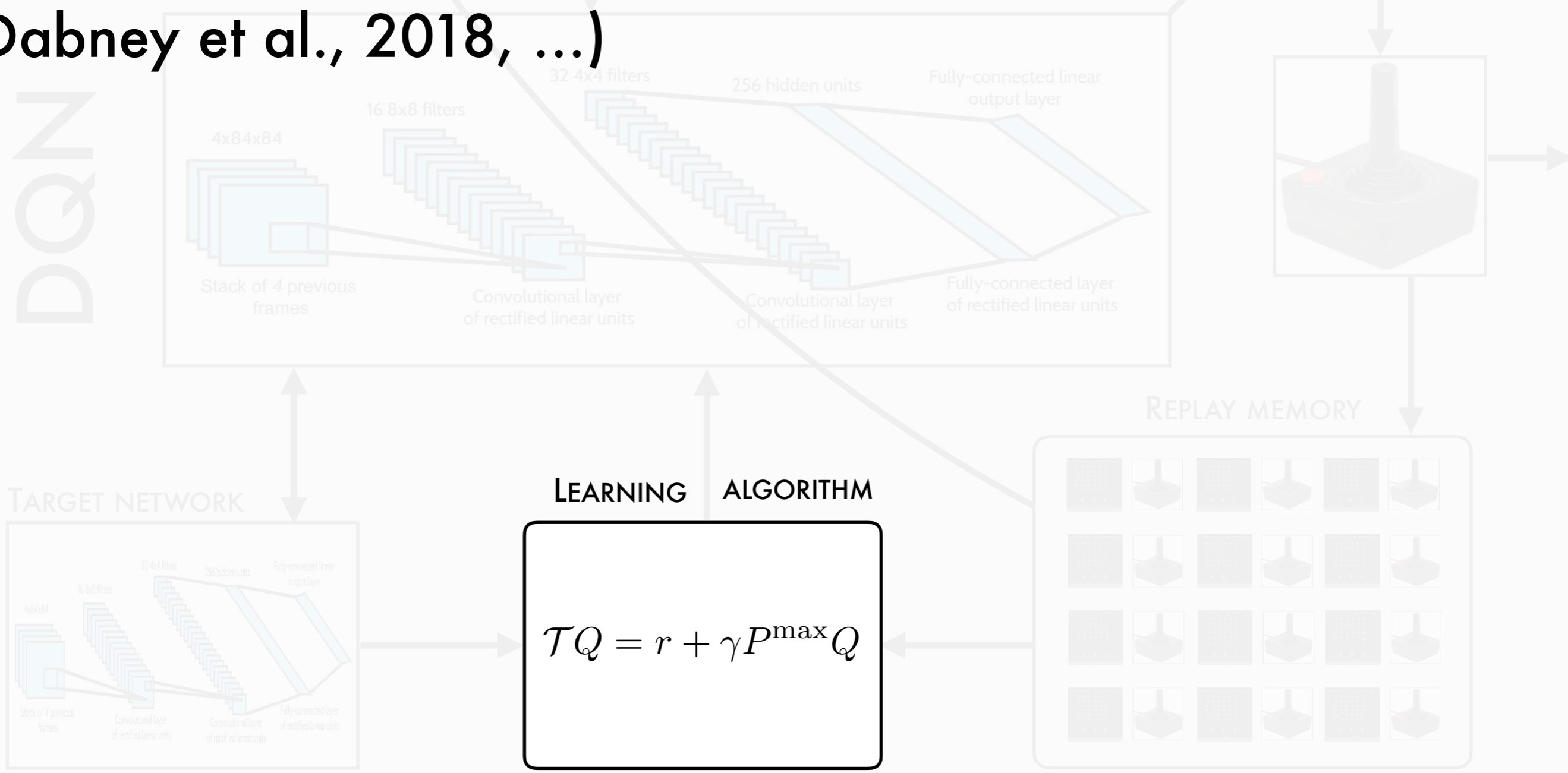
Advantage learning (Bellemare et al., 2015)

Q(λ) (Harutyunyan et al., 2016)

Retrace (Munos et al., 2016)

Distributional methods (Bellemare et al., 2017,

Dabney et al., 2018, ...)



1. Distributional reinforcement learning

2. Exploration with pseudo-counts



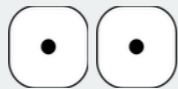
x

$E R(x) :$

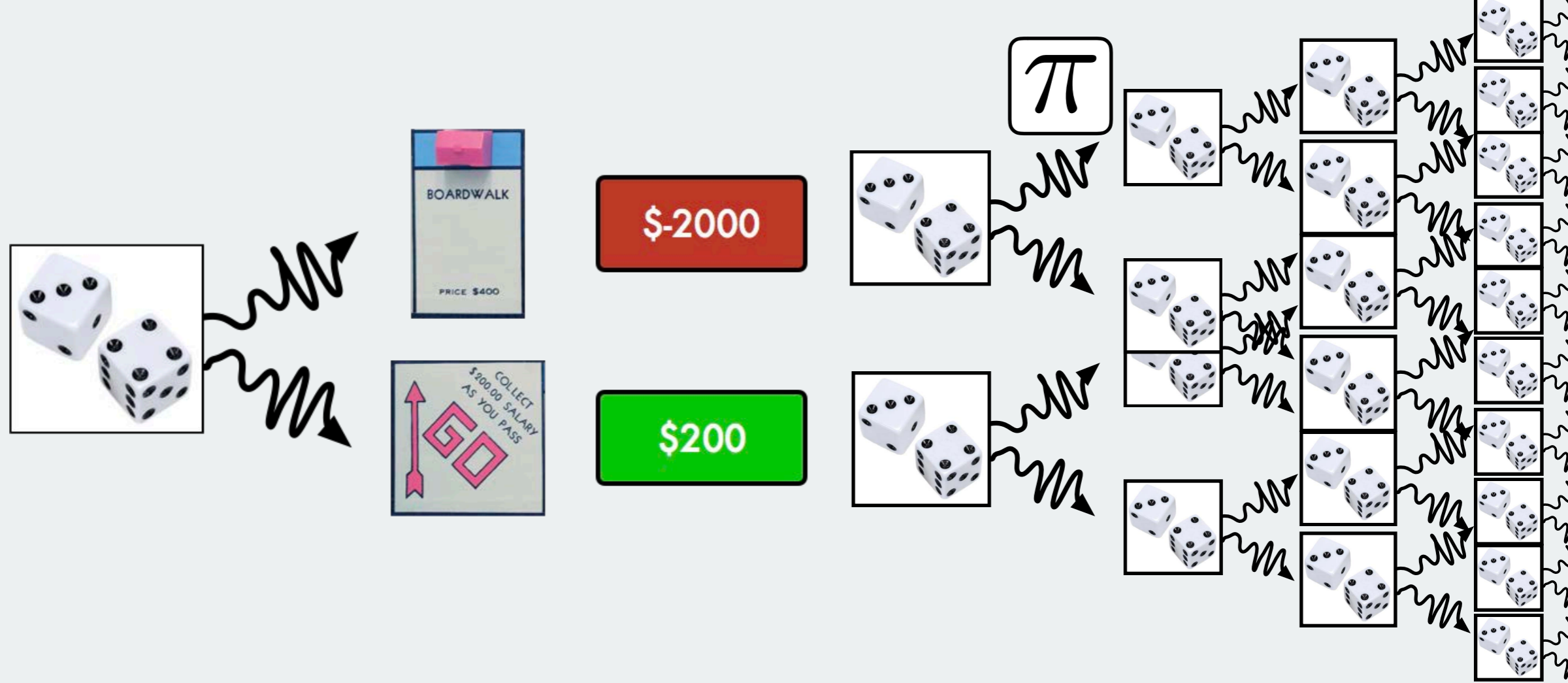
$$\frac{1}{36} (-2000)$$

$R(x)$

$$+ \frac{35}{36} (200)$$



139



$$R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) \dots$$

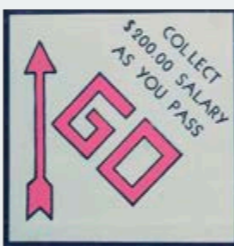
BELLMAN EQUATION

$$V^\pi(x) = \mathbf{E} R(x) + \gamma \mathbf{E}_{x' \sim P^\pi} V^\pi(x')$$

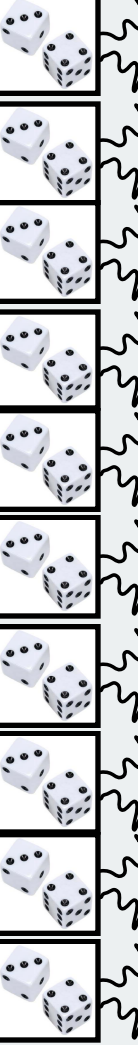
GROUND TRUTH



-\$2000

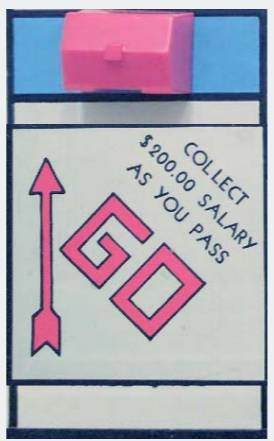


\$200



$$V^\pi(x) = \mathbf{E} R(x) + \gamma \mathbf{E}_{x' \sim P^\pi} V^\pi(x')$$

IMPLIED MODEL



\$139

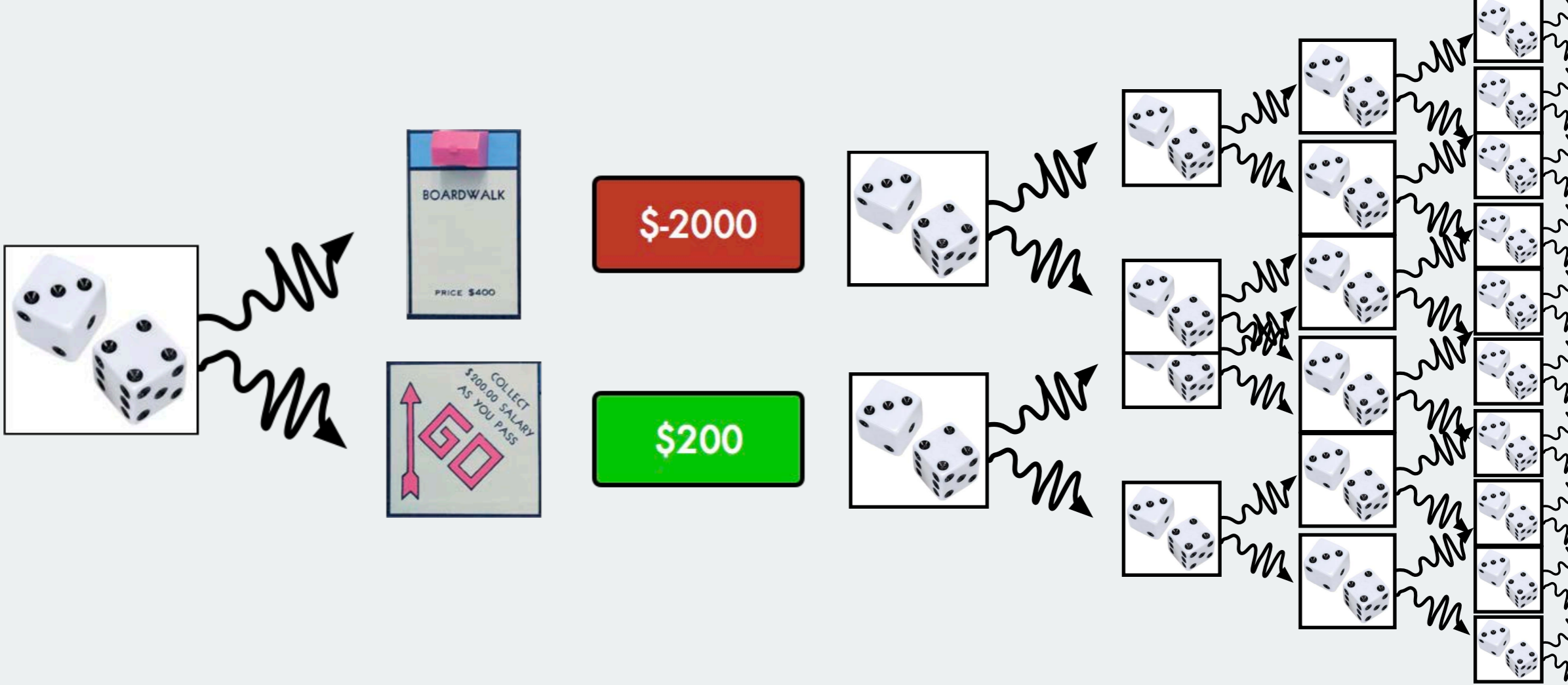


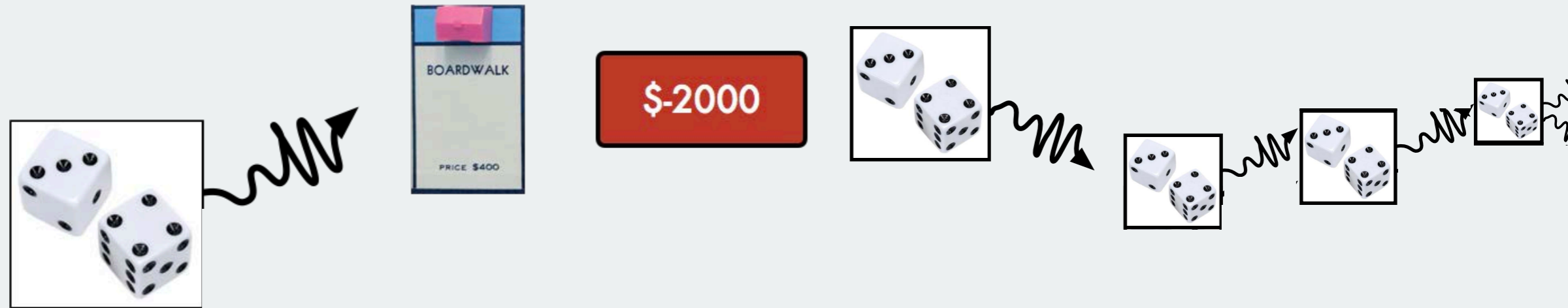
SCHOKNECHT 2002; PARR ET AL., 2008; SUTTON ET AL. 2008

IMPLIED MODEL



$$\begin{aligned} V^\pi(x) &= \mathbf{E} [R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \dots \mid x_0 = x, a_t \sim \pi] \\ &= \mathbf{E} [R(x_0)] + \gamma \mathbf{E} [R(x_1)] + \gamma^2 \mathbf{E} [R(x_2)] + \dots \end{aligned}$$





$$V^\pi(x) = \mathbf{E} [R(x_0) + \gamma R(x_1) + \gamma^2 R(x_2) + \dots]$$

BELLMAN EQUATION

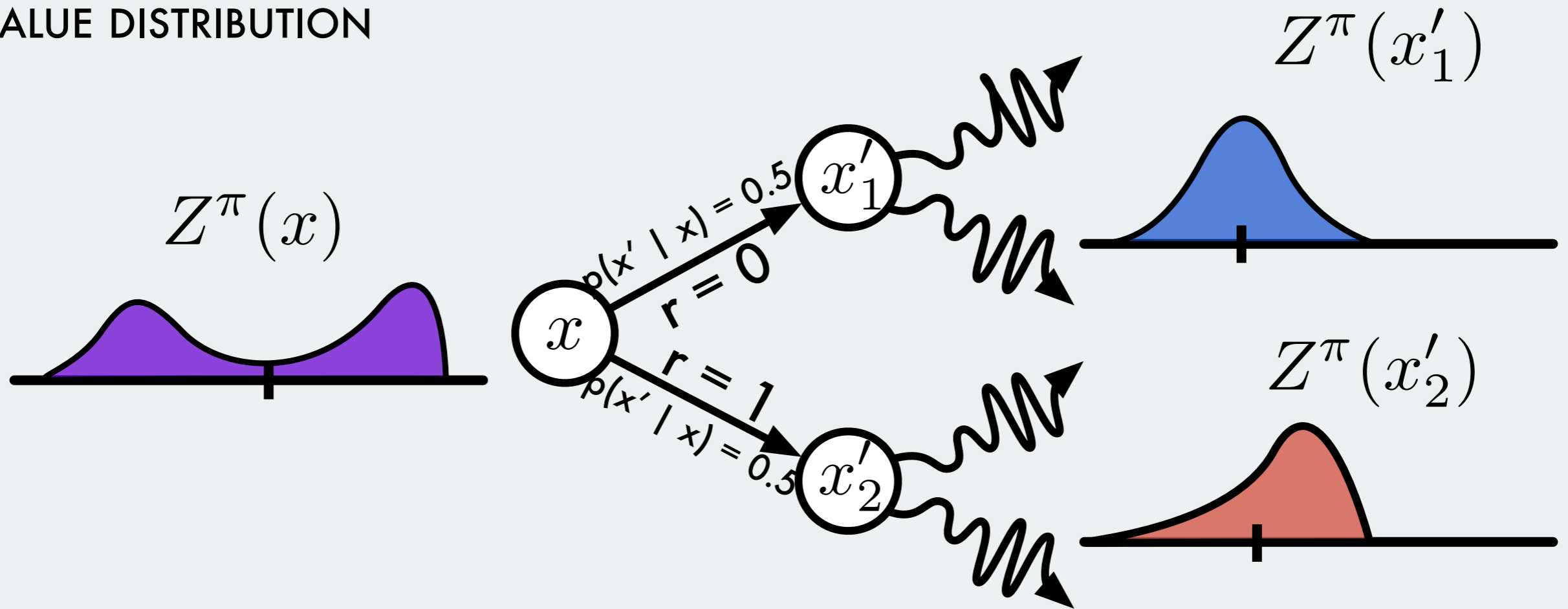
$$\cancel{\mathbf{E}} [Z^\pi(x)] = \cancel{\mathbf{E}} R(x) + \gamma \cancel{\mathbf{E}} [Z^\pi(X')]$$

DISTRIBUTIONAL BELLMAN EQUATION

$$Z^\pi(x) \stackrel{D}{=} R(x) + \gamma Z^\pi(X')$$

$$Z^\pi(x) \stackrel{D}{=} R(x) + \gamma Z^\pi(X')$$

↑
VALUE DISTRIBUTION



Bellman (1957): Bellman equation for **mean**

Sobel (1982): ... for **variance**

Engel (2003): ... for **Bayesian uncertainty**

Azar et al. (2011), Lattimore & Hutter (2012): ... for **higher moments**

Morimura et al. (2010, 2010b): ... for **densities**

\$150

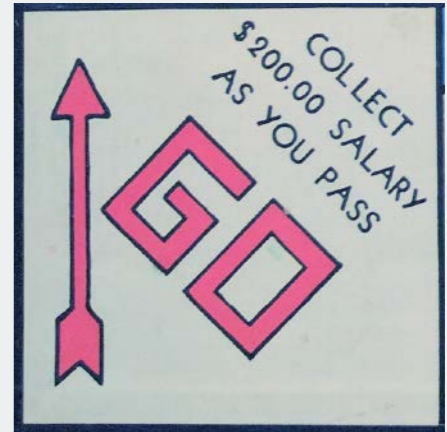


+

\$300



\$200



\$450



+

\$300



\$200



\$450

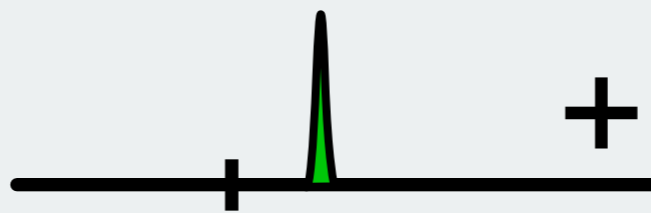
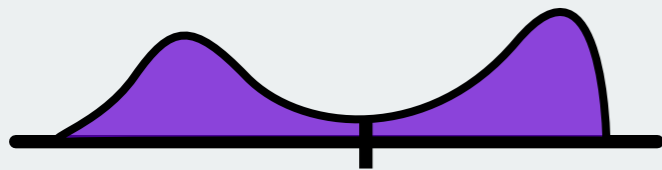
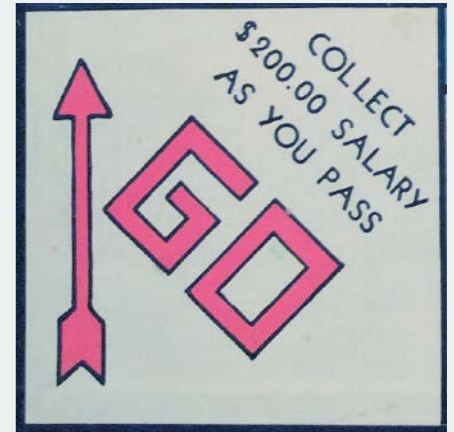


+

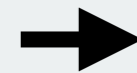
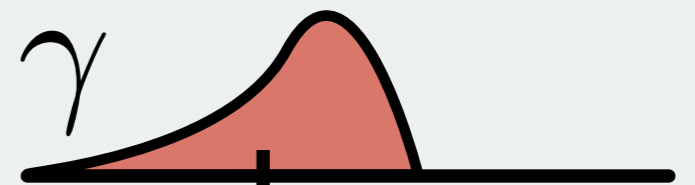
γ \$300



\$200



+



\$450

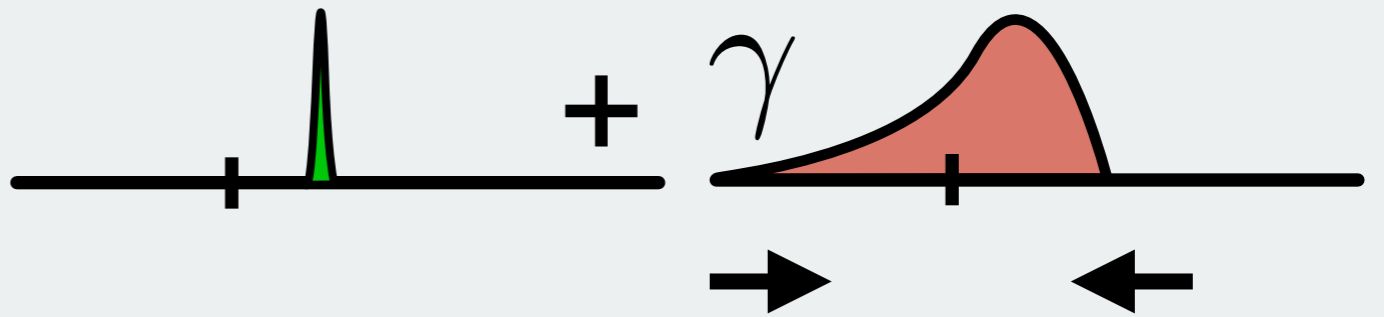
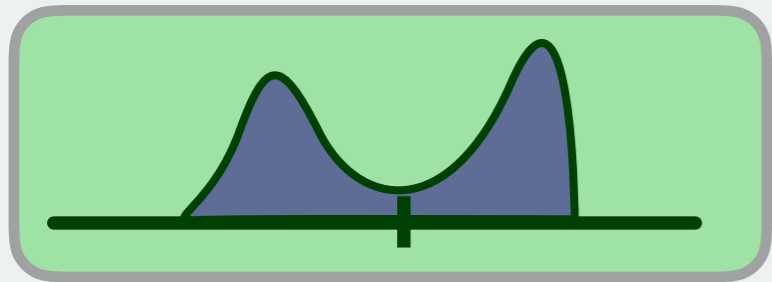
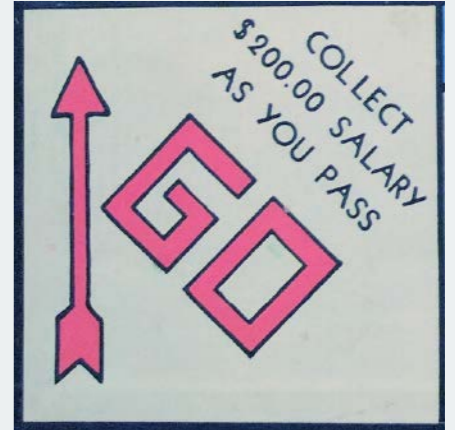


+

γ \$300



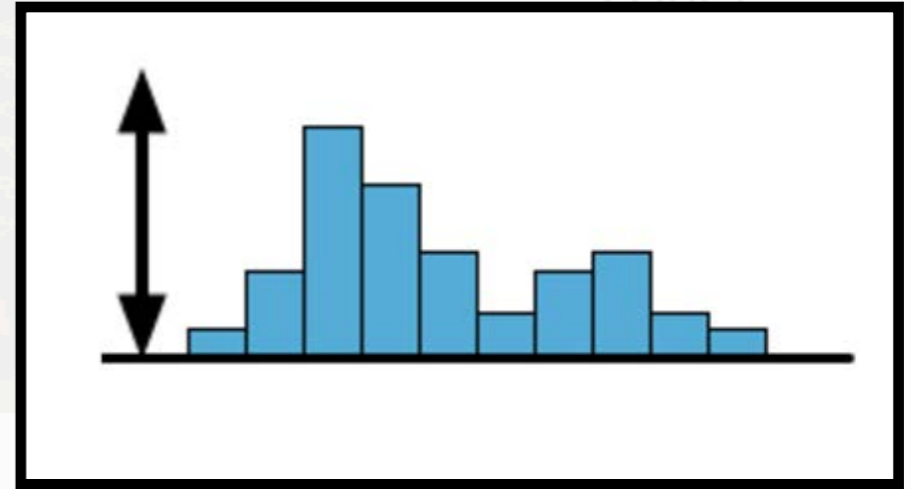
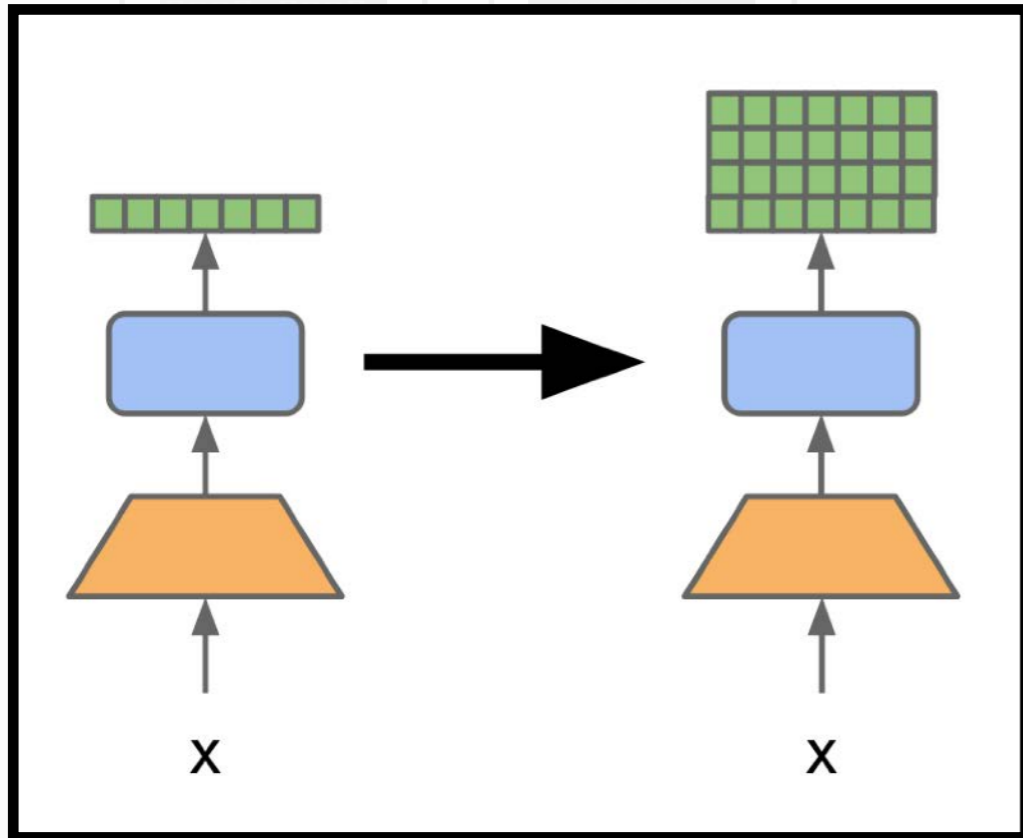
\$200



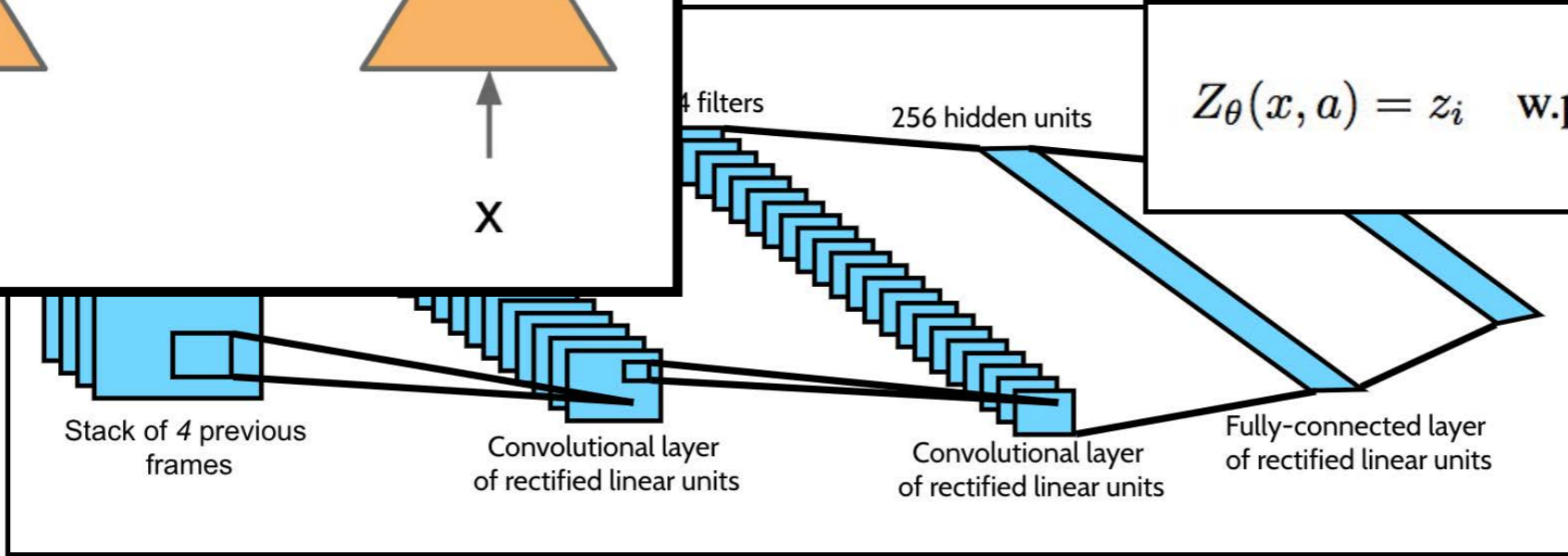
$Q(x, a)$

$p(z | x, a)$

Discrete distribution

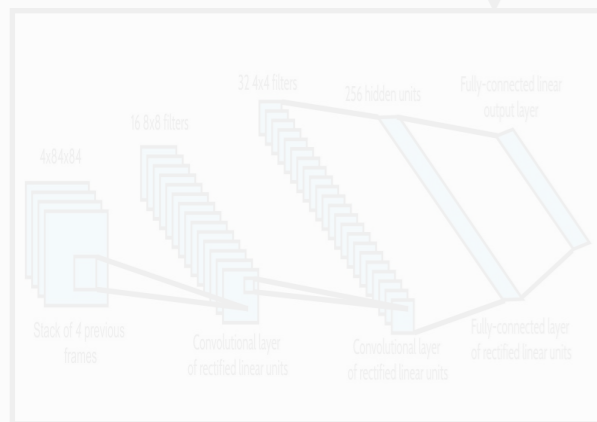


$$Z_{\theta}(x, a) = z_i \quad \text{w.p.} \quad p_i(x, a) := \frac{e^{\theta_i(x, a)}}{\sum_j e^{\theta_j(x, a)}}.$$



C51

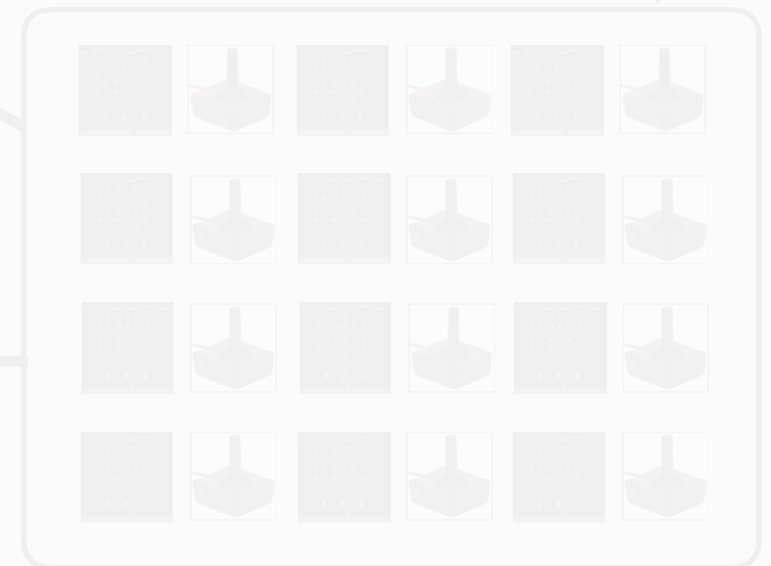
TARGET NETWORK

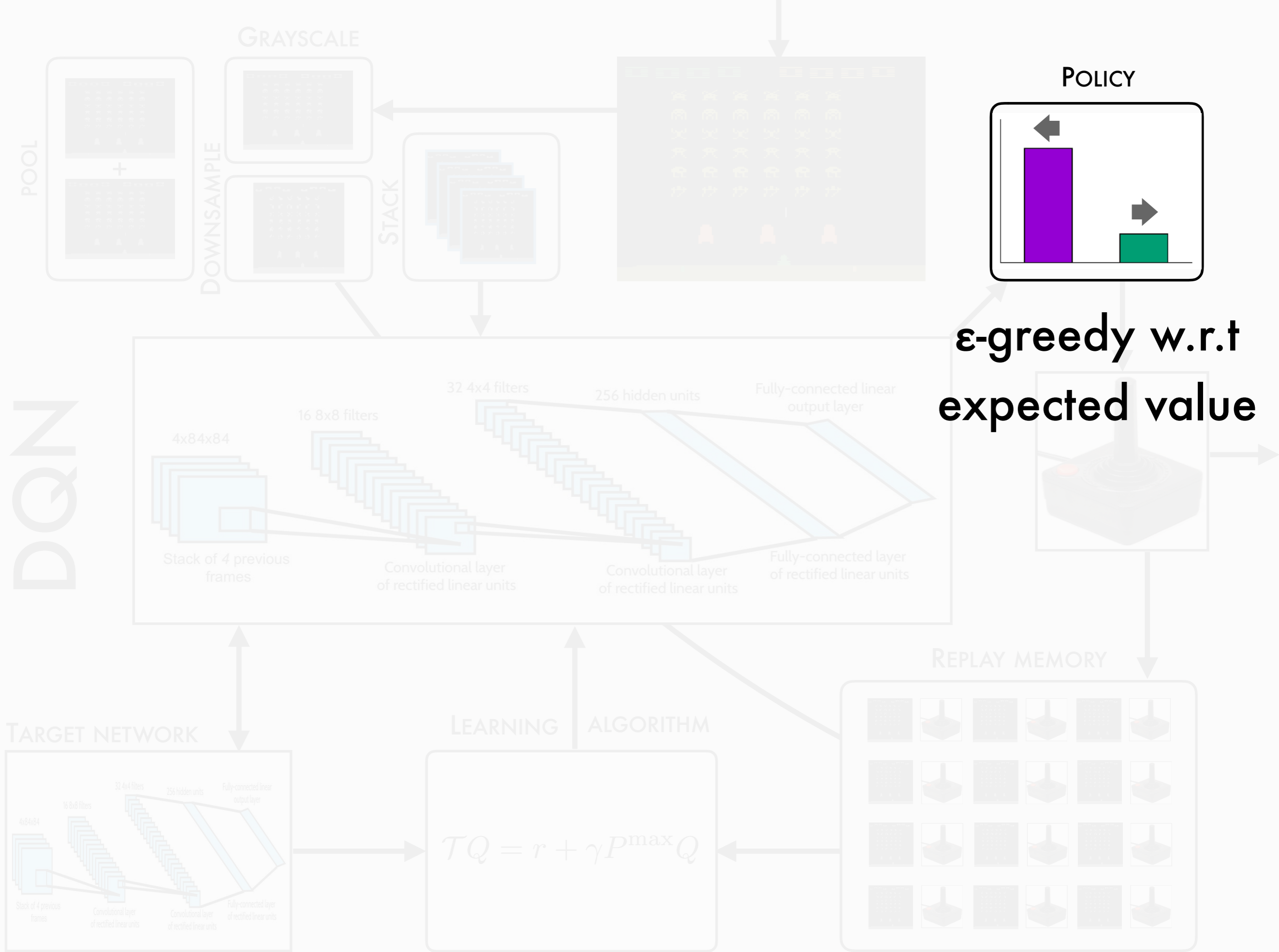


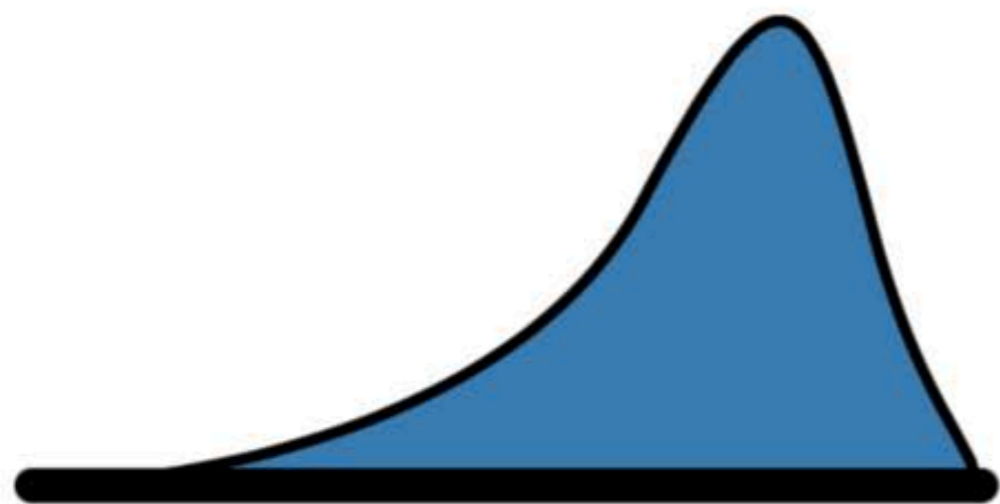
LEARNING ALGORITHM

$$\mathcal{T}Q = r + \gamma P^{\max} Q$$

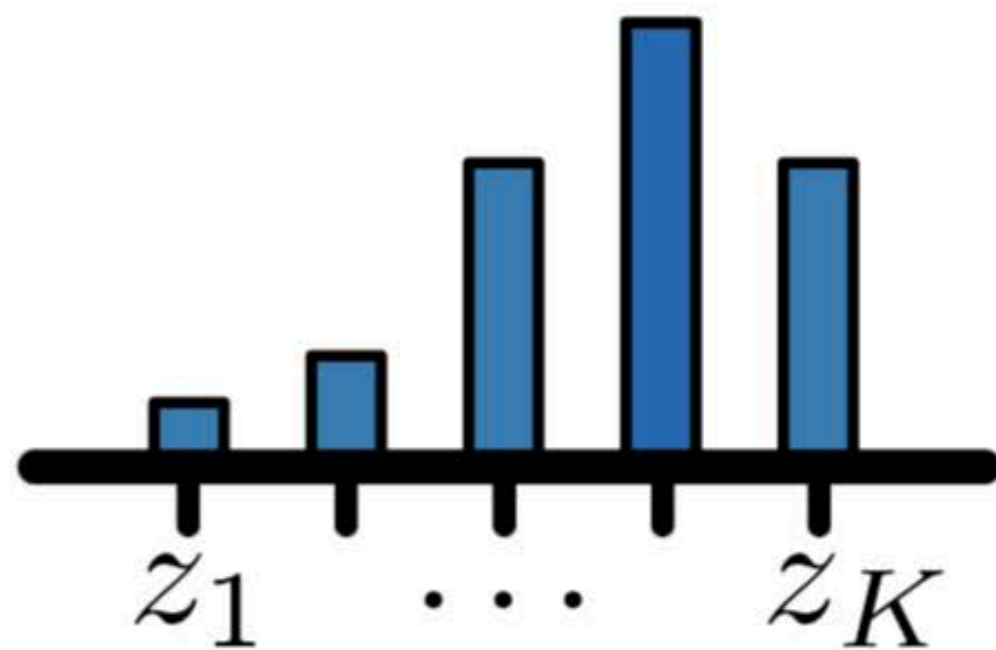
REPLAY MEMORY



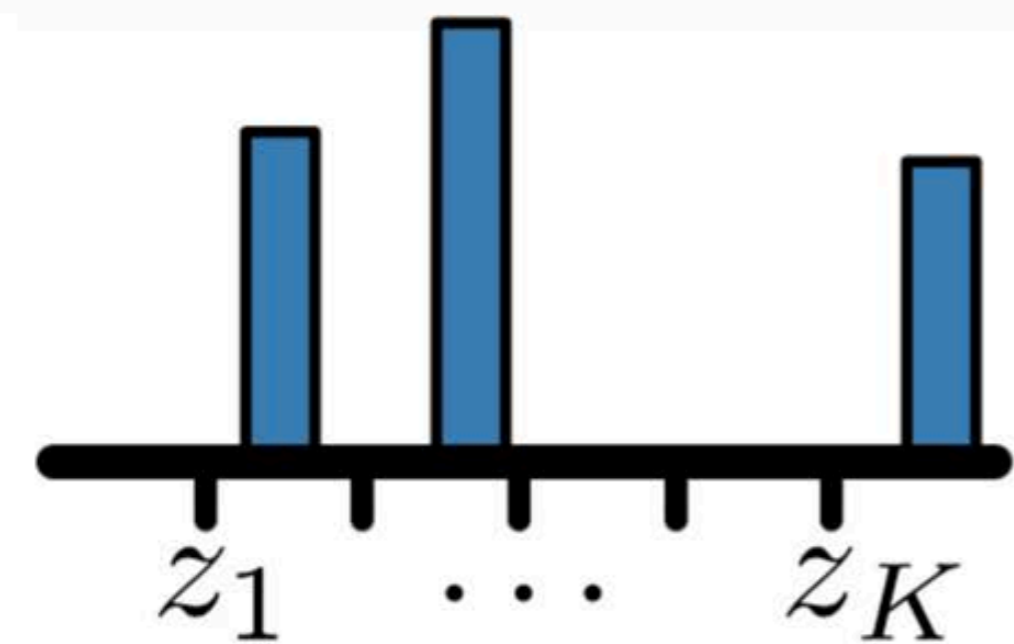




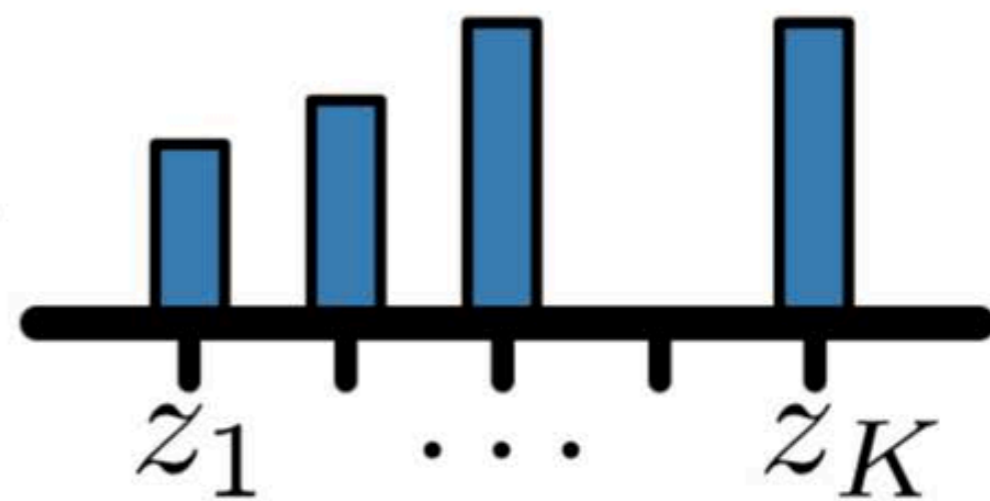
\approx



APPROXIMATION



$\xrightarrow{\Pi_c}$



PROJECTION

$$\mathcal{T}^\pi Z = R + \gamma P^\pi Z$$

1. From x, a , **sample** a transition:

$$r, X', A' \sim R(x, a), P(\cdot | x, a), \pi(\cdot | X')$$

2. Compute **sample backup**

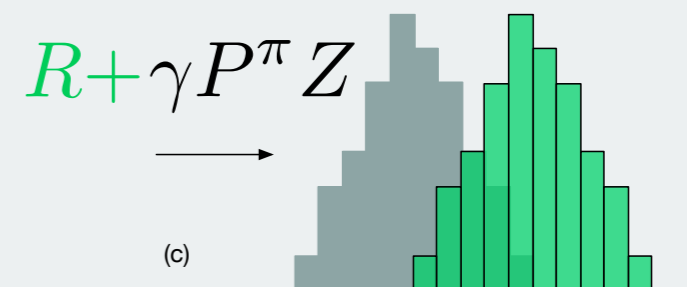
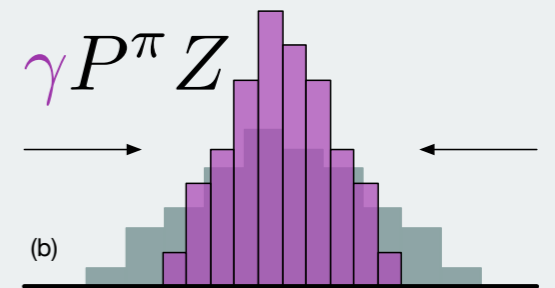
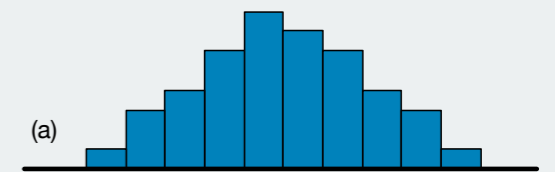
$$\hat{\mathcal{T}}^\pi Z(x, a) := r + \gamma Z(X', A')$$

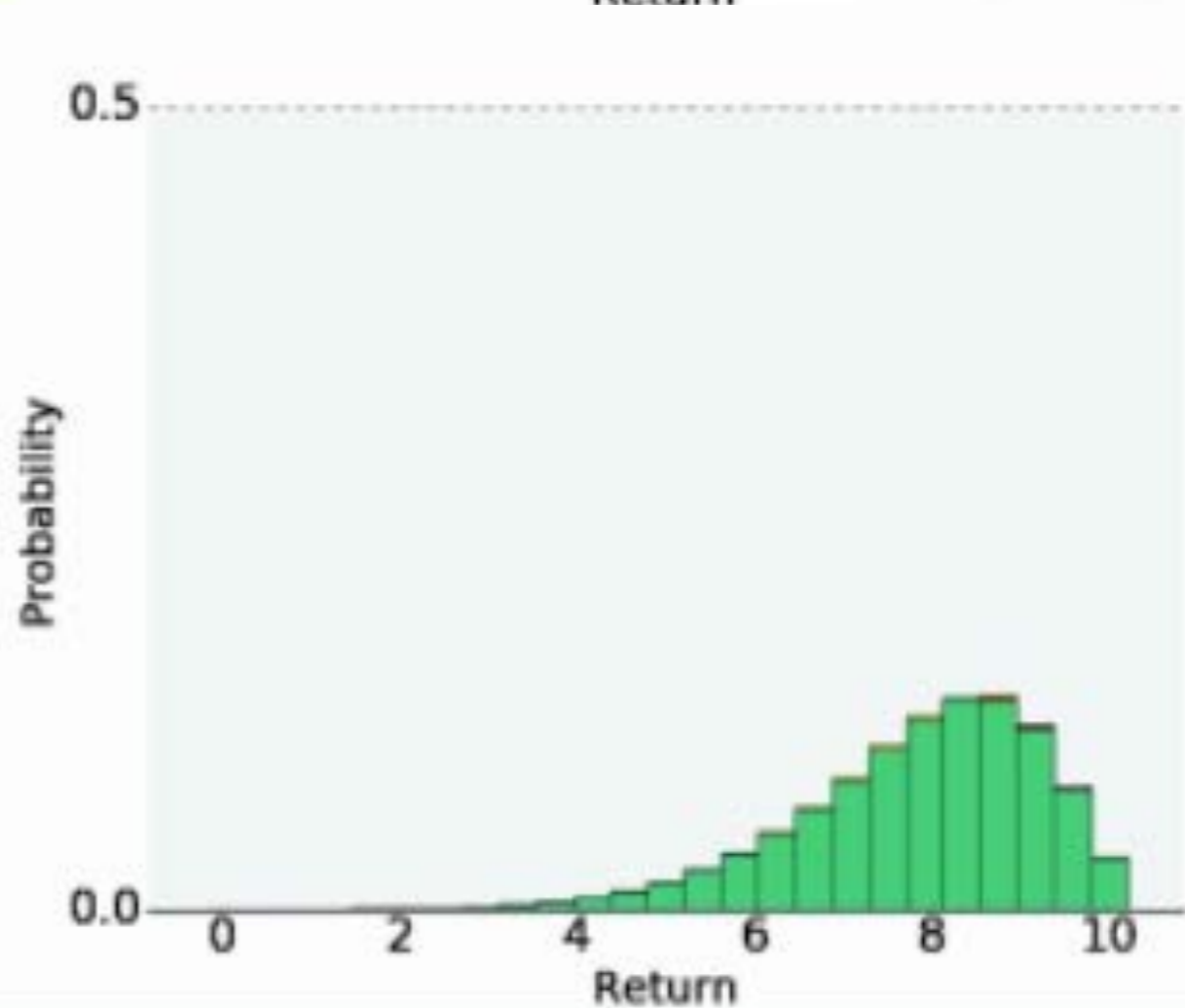
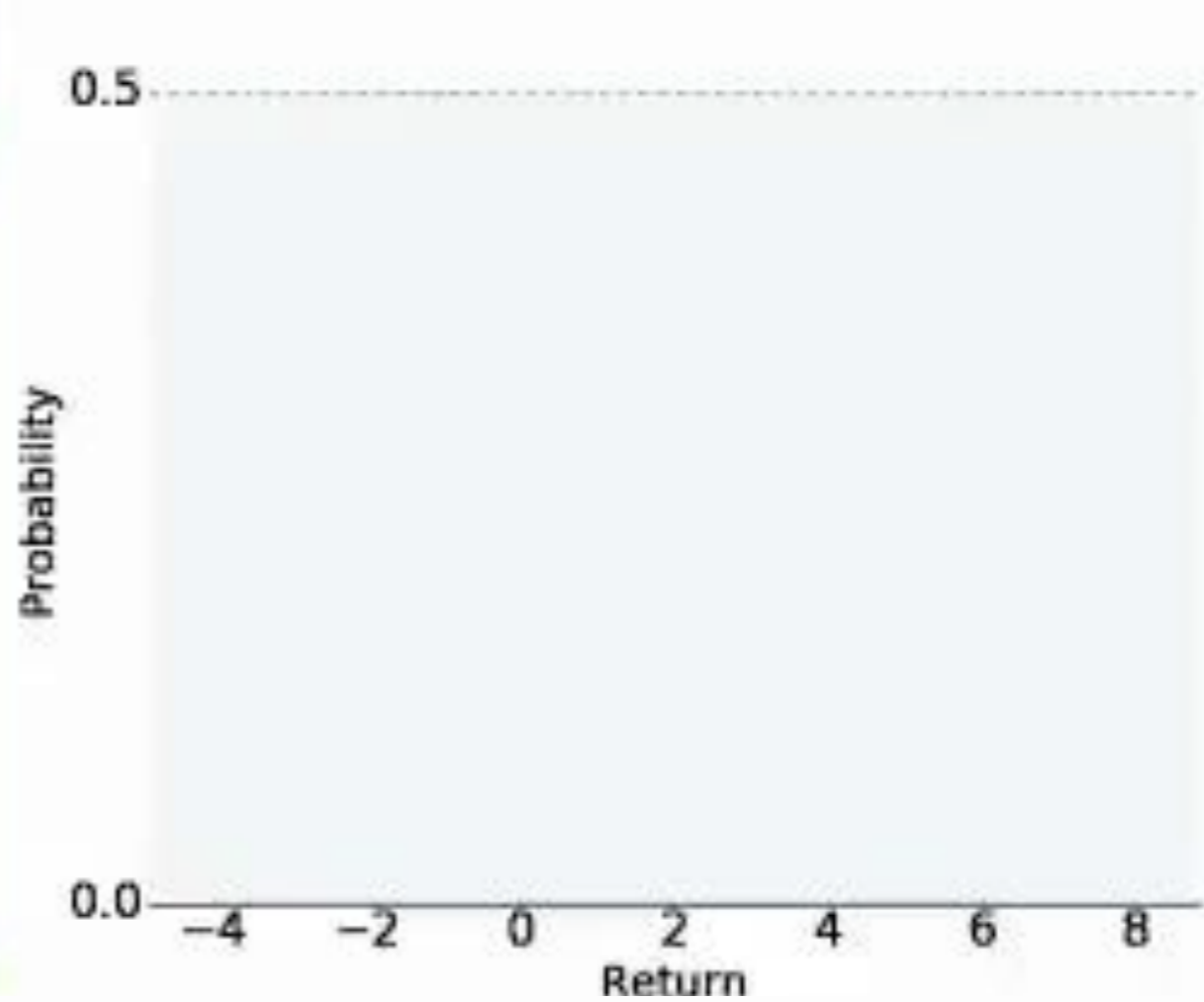
3. **Project** onto approximation support:

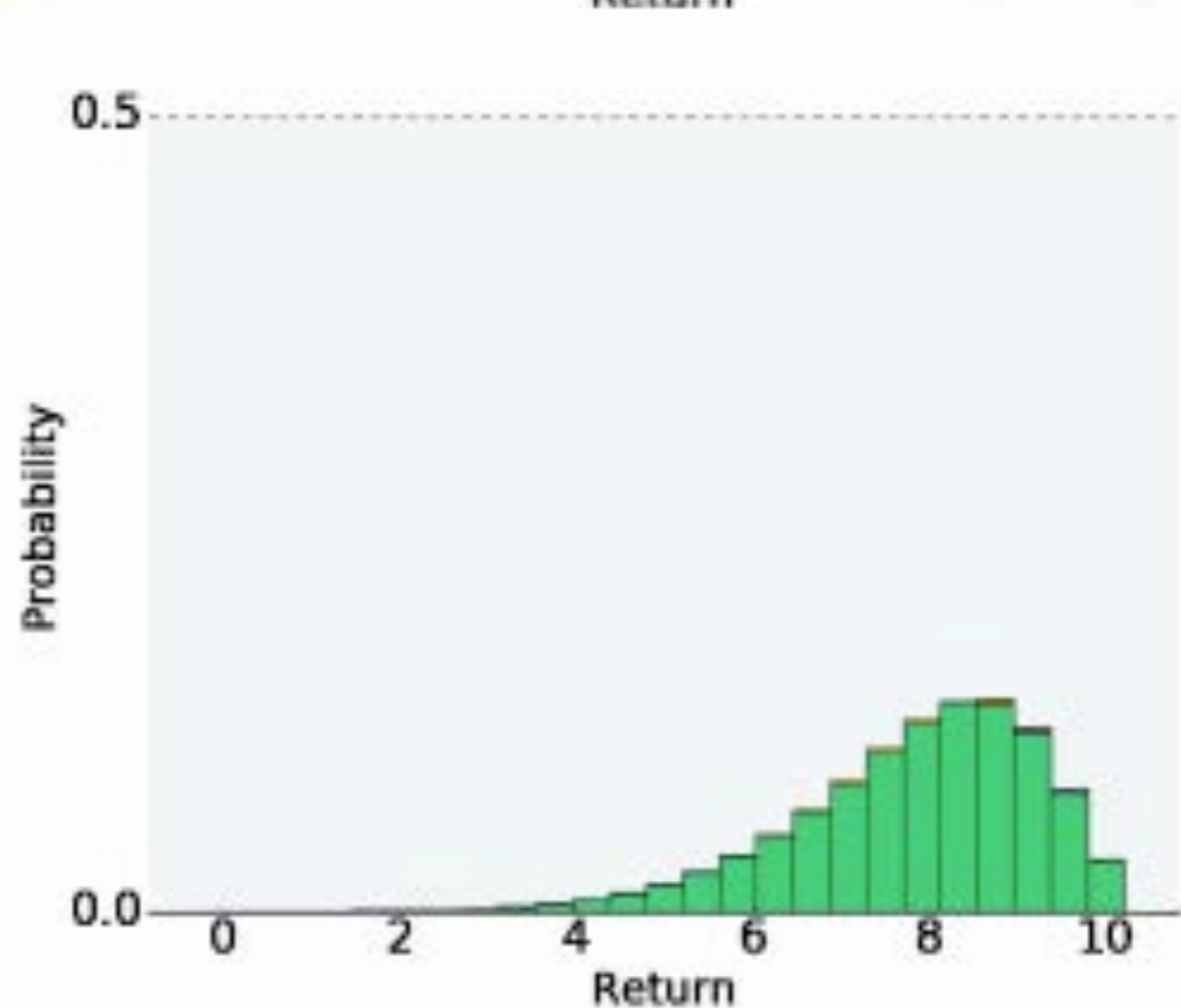
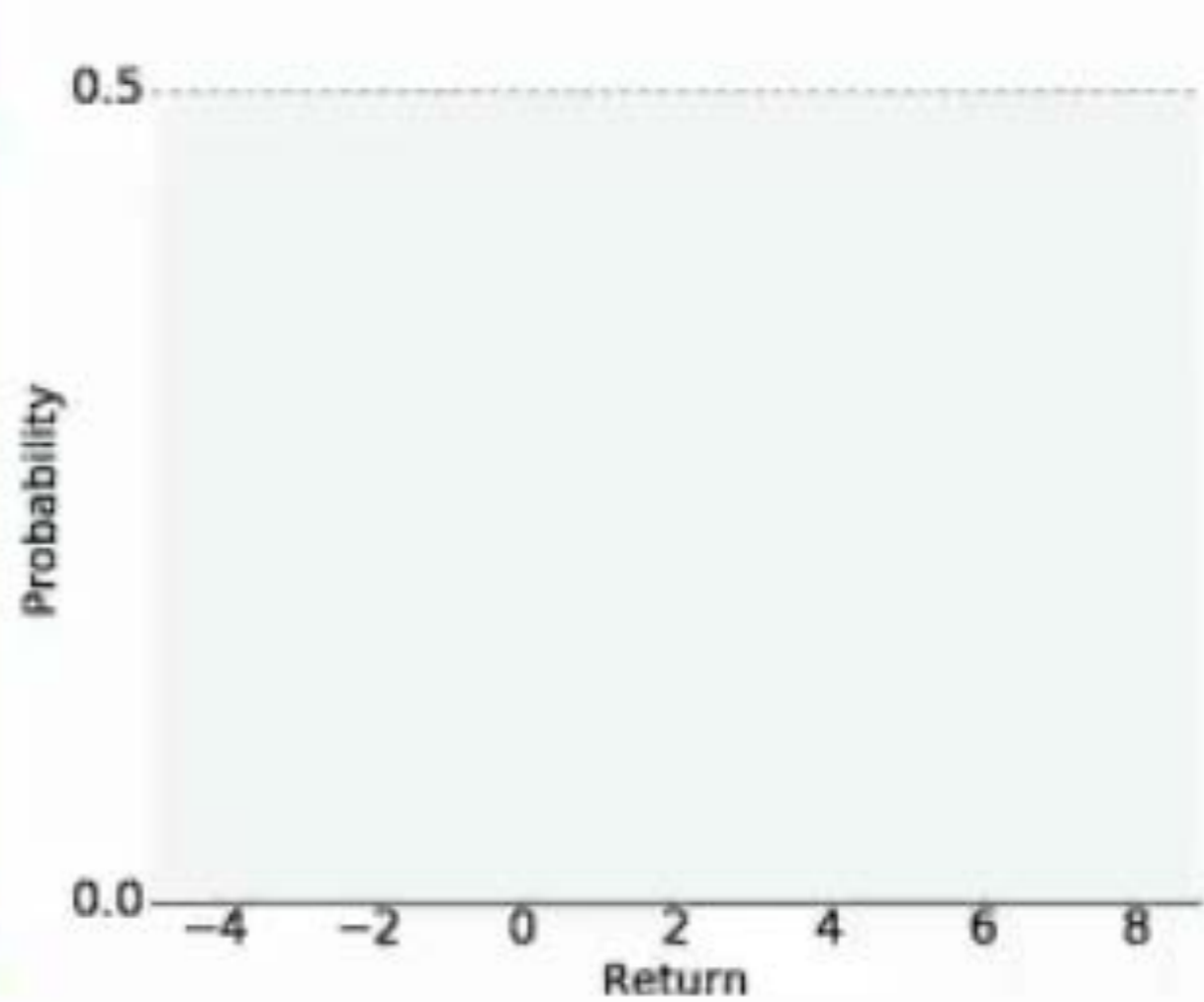
$$\Phi \hat{\mathcal{T}}^\pi Z(x, a)$$

4. Update towards projection, i.e.
take a **KL-minimizing step**

$P^\pi Z$





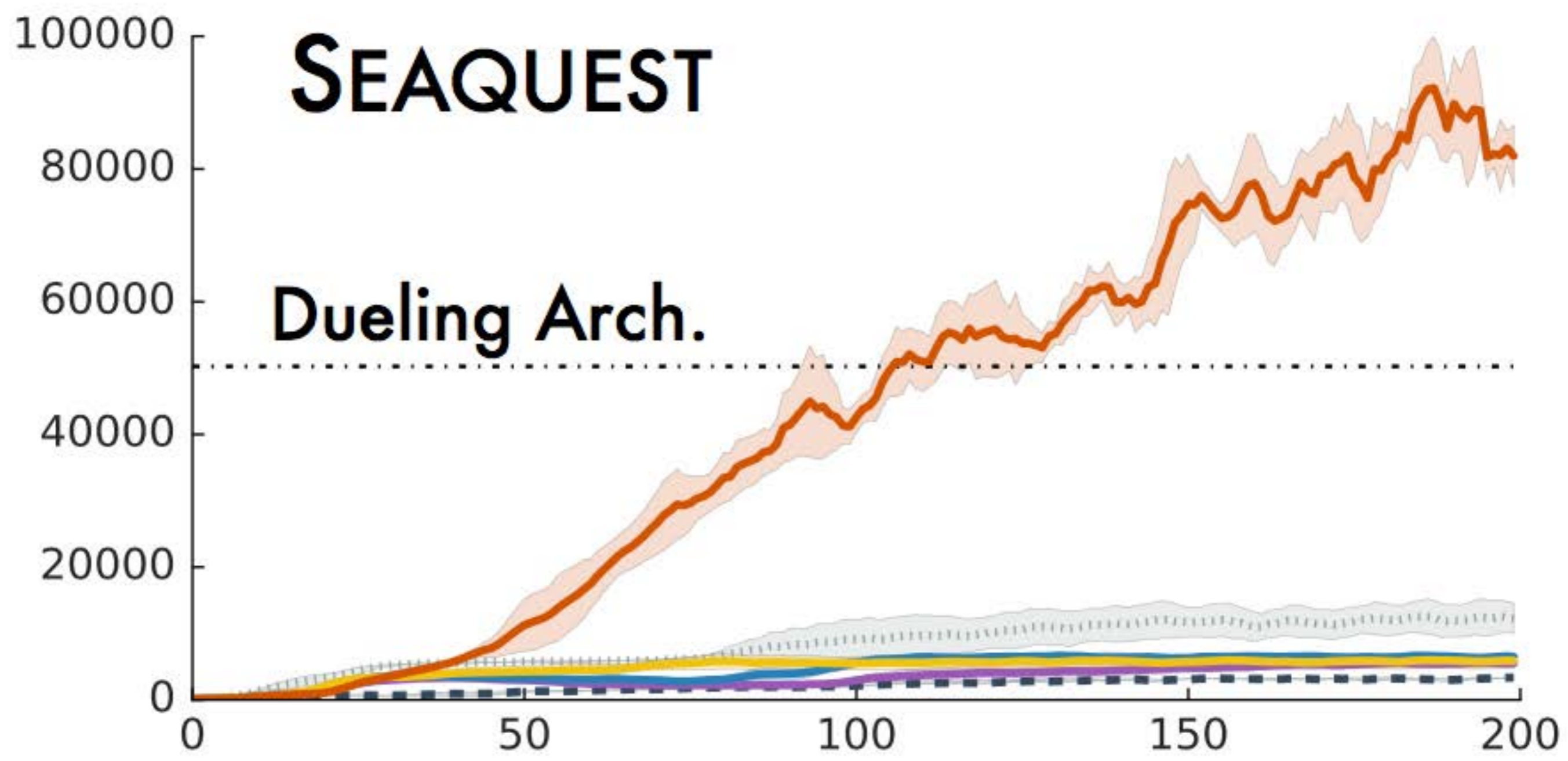


	Mean	Median	> H.B.	> DQN
DQN	228%	79%	24	0
DDQN	307%	118%	33	43
DUEL.	373%	151%	37	50
PRIOR.	434%	124%	39	48
PR. DUEL.	592%	172%	39	44

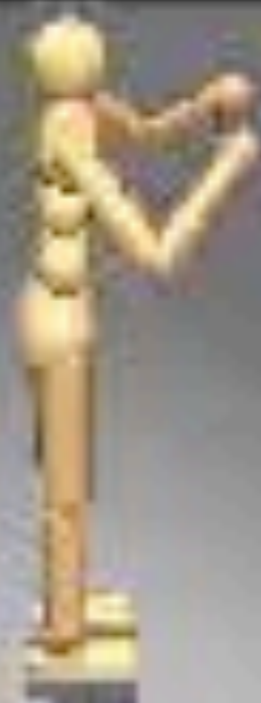


SEAQUEST

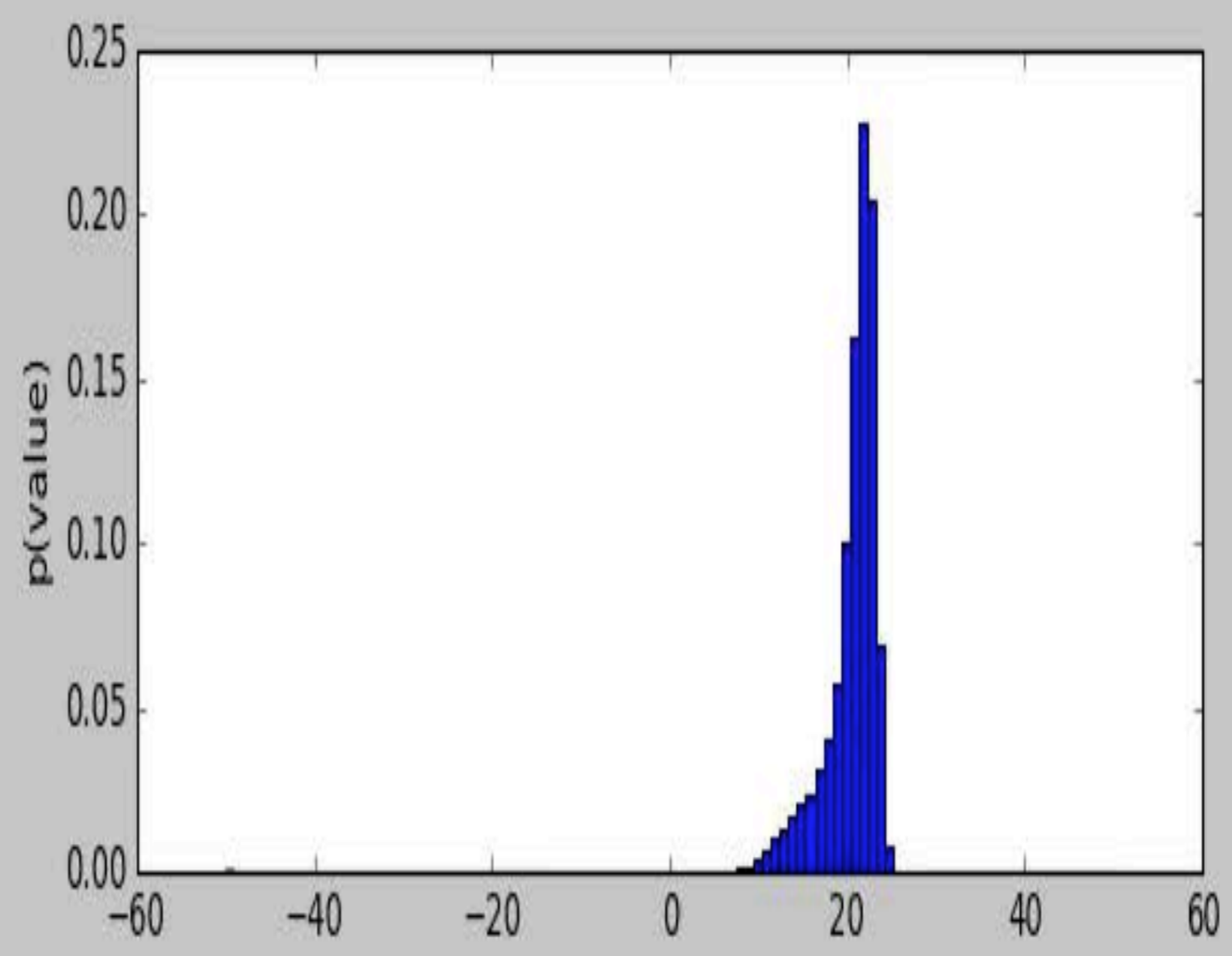
Dueling Arch.



Training Frames (millions)

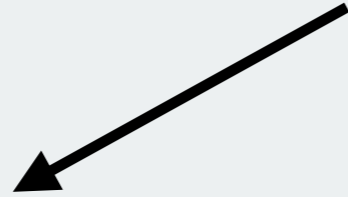


Time	00:00.00
Reward	0.002
Cumulative	0.002



DISTRIBUTIONAL PERSPECTIVE

BELLEMARE ET AL., 2016



QUANTILE REGRESSION

DABNEY ET AL., 2017

IMPLICIT QUANTILE NETWORKS

DABNEY ET AL., 2018

DISTRIBUTIONAL ACTOR-CRITIC

BARTH-MARON ET AL., 2018

GENERATIVE QUANTILE NETWORKS

OSTROVSKI ET AL., 2018

THE RAINBOW AGENT

HESSEL ET AL., 2018

CRAMER DISTANCE

BELLEMARE ET AL., 2017

ANALYSIS OF C51

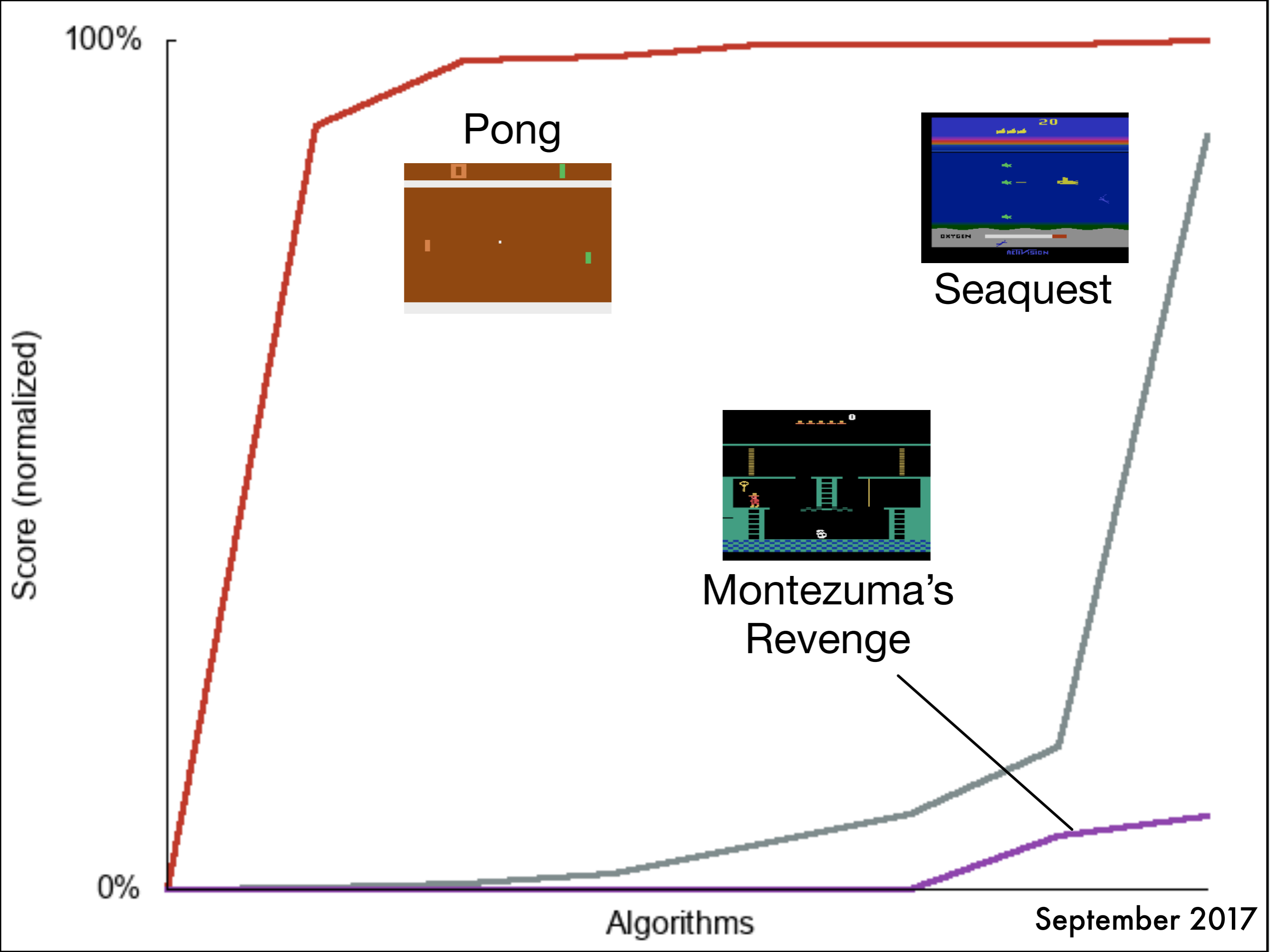
ROWLAND ET AL., 2017

"AS EXPECTED"

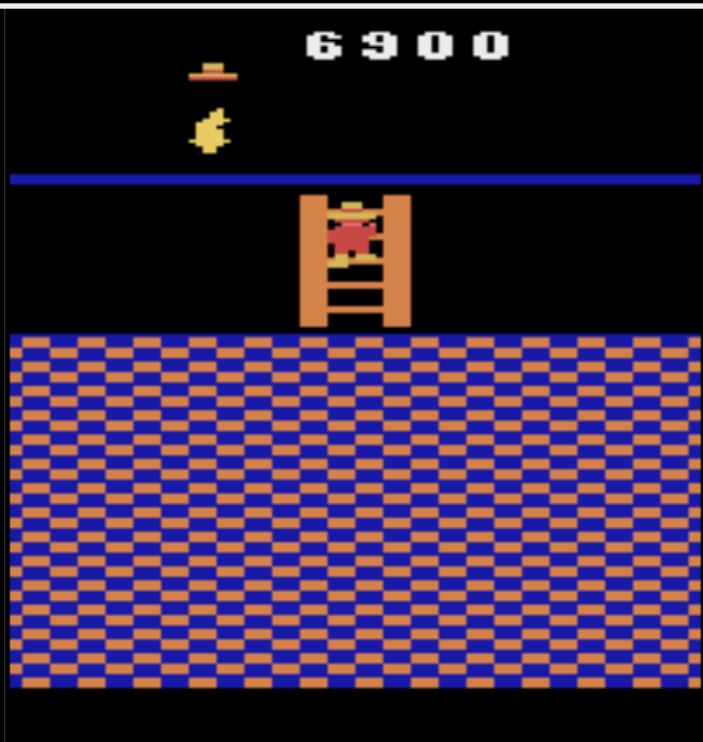
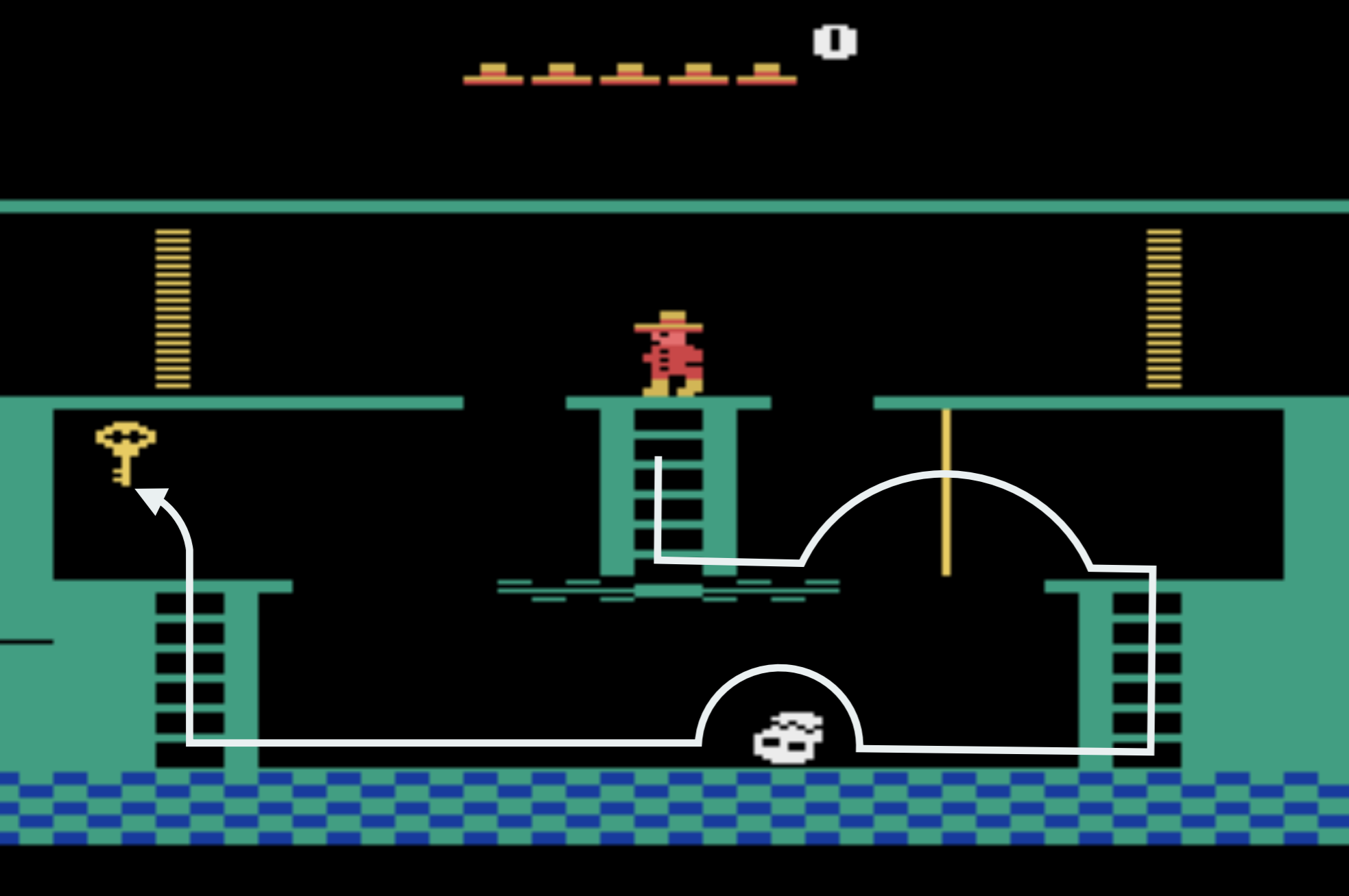
LYLE ET AL., 2018

PGMRL WORKSHOP, ICML

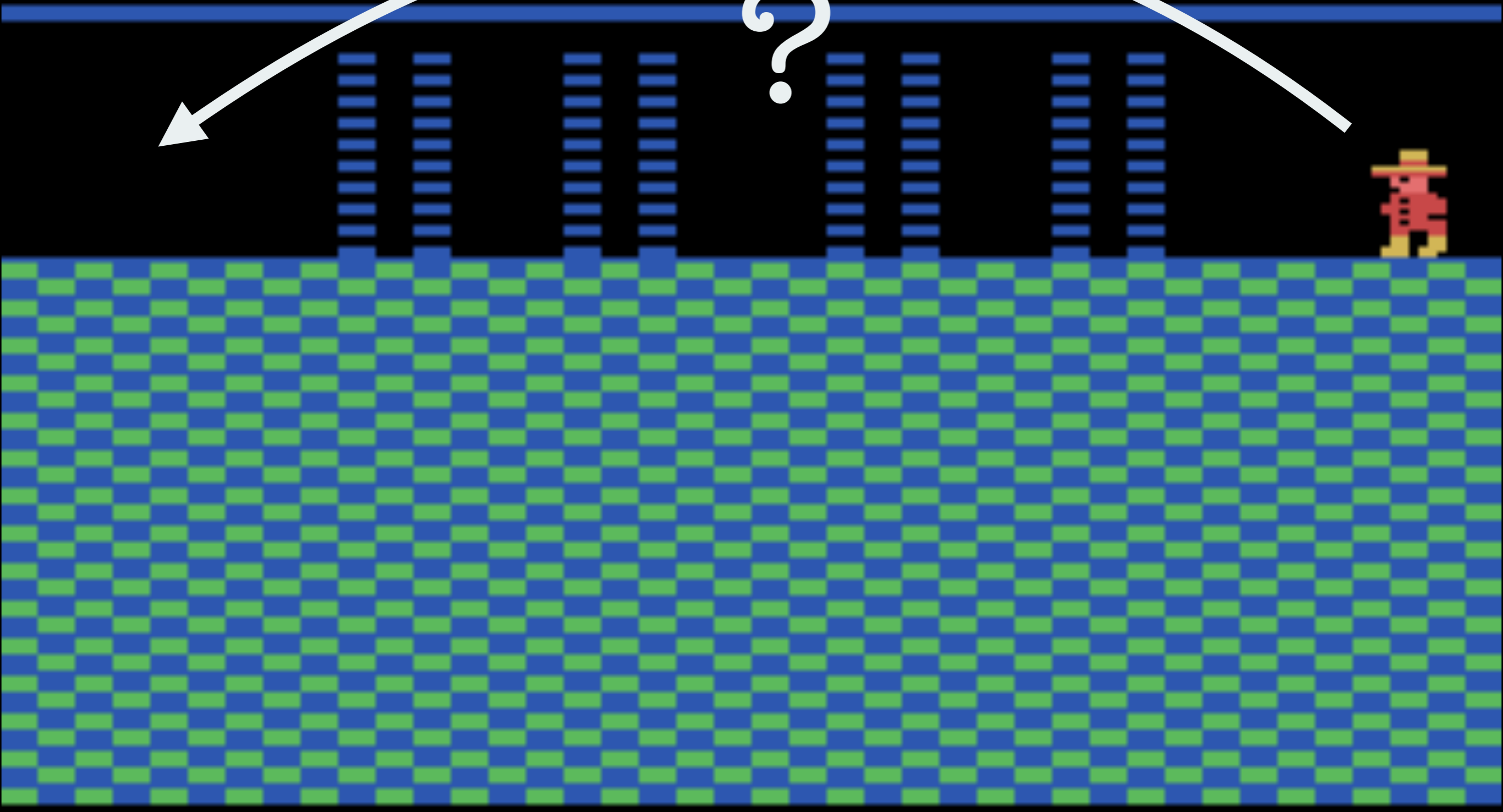
1. Distributional reinforcement learning
- 2. Exploration with pseudo-counts**







3600



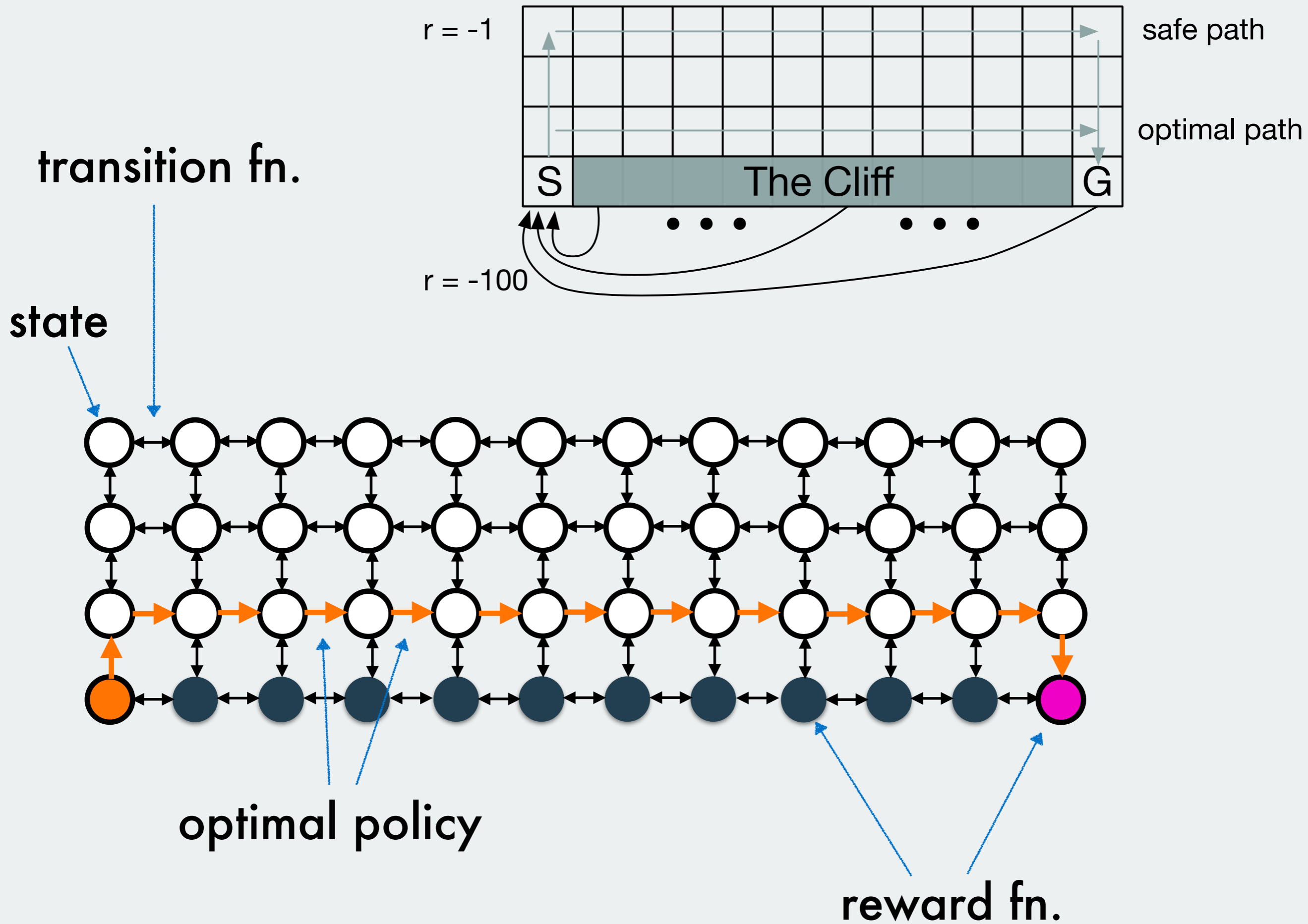
EXPLORATION

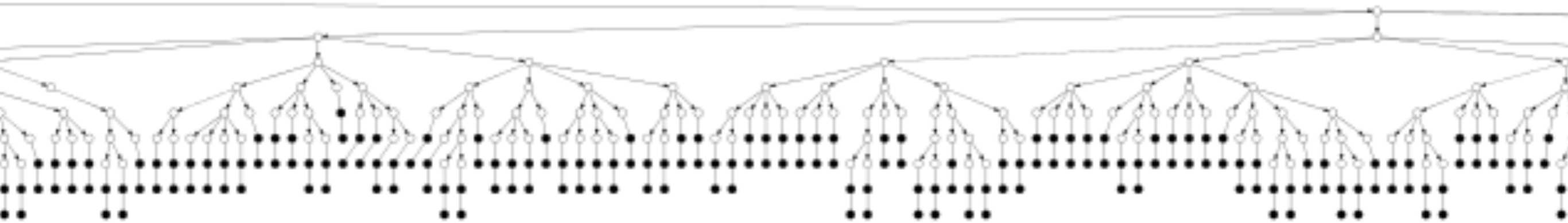
Bellman equation:

$$Q(x, a) = r(x, a) + \gamma \mathbf{E}_{x' \sim P} \max_{a' \in \mathcal{A}} Q(x', a')$$

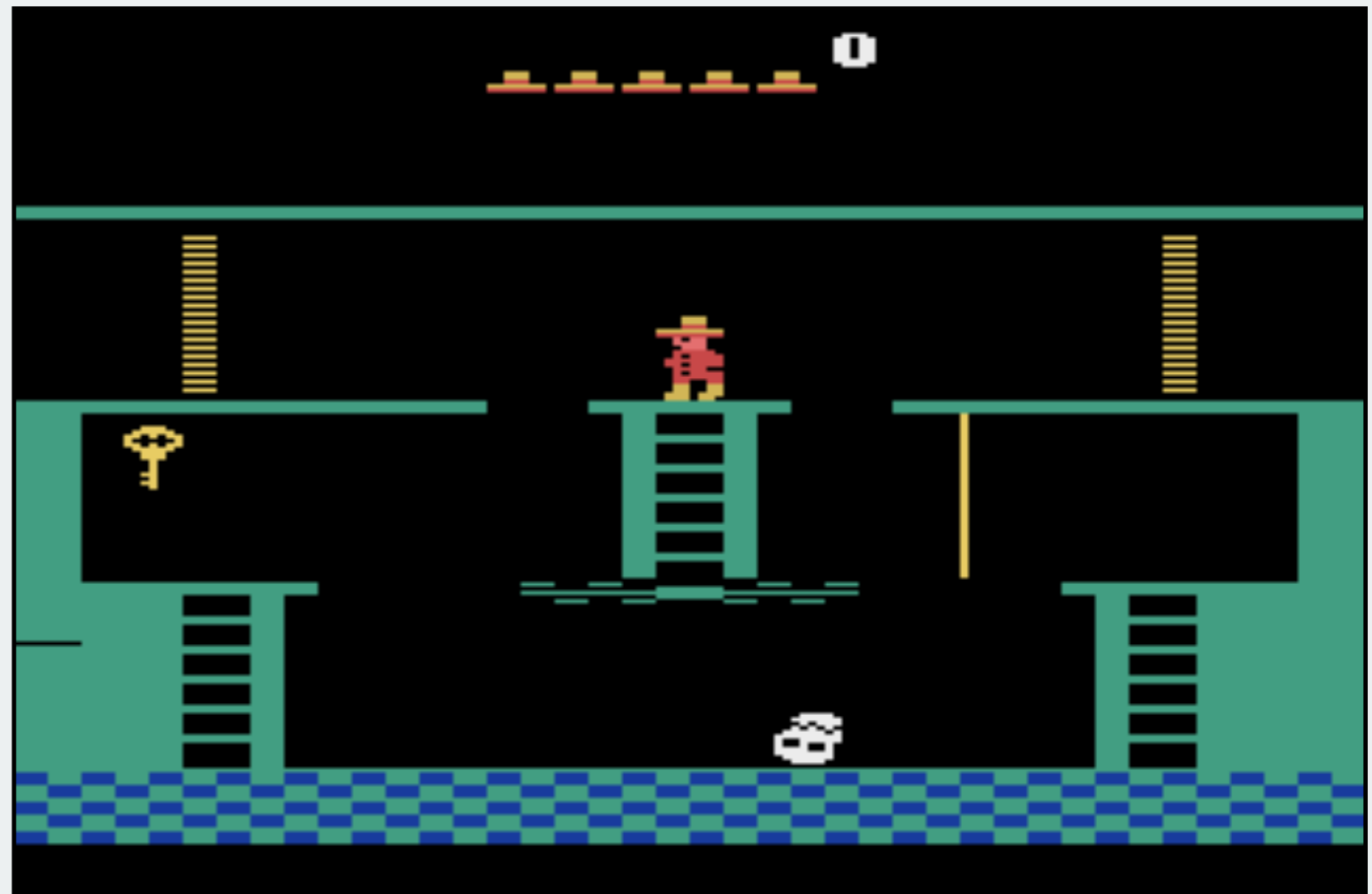
Exploration bonus (Strehl and Littman, 2008):

$$Q(x, a) = \hat{r}(x, a) + \gamma \mathbf{E}_{x' \sim \hat{P}} \max_{a' \in \mathcal{A}} Q(x', a') + \frac{\beta}{\sqrt{N(x, a)}}$$

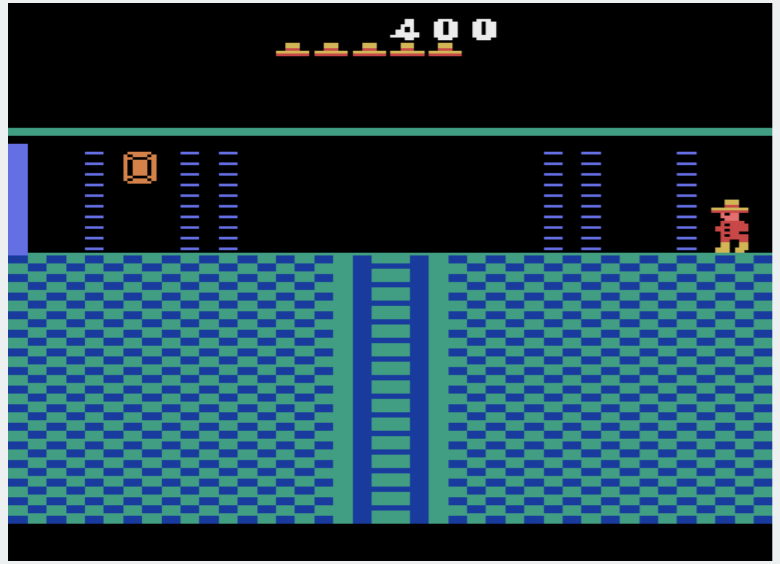
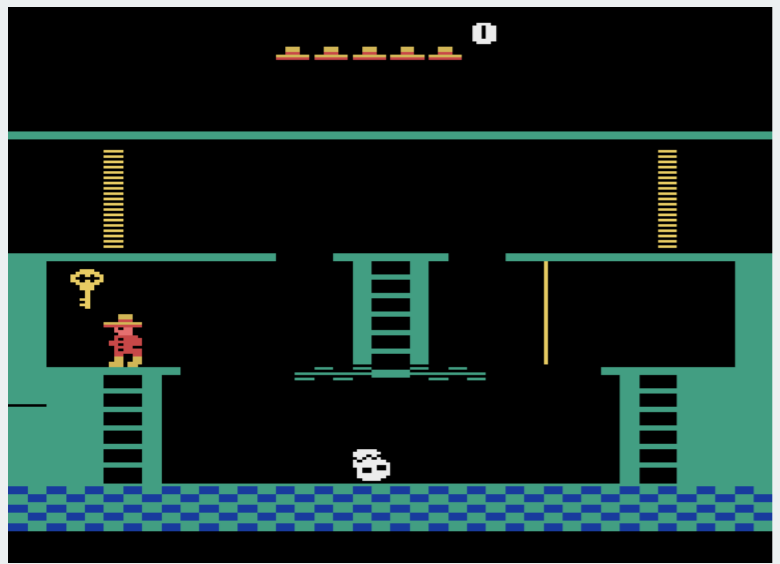




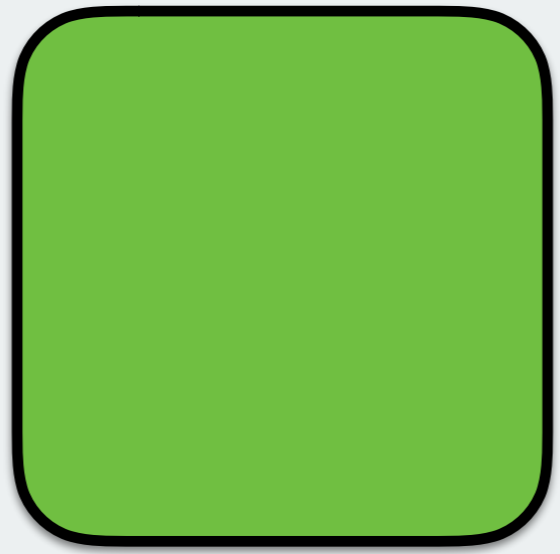
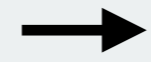
Most observations
experienced *once*:
up to 90% singletons
(Blundell et al., 2016)



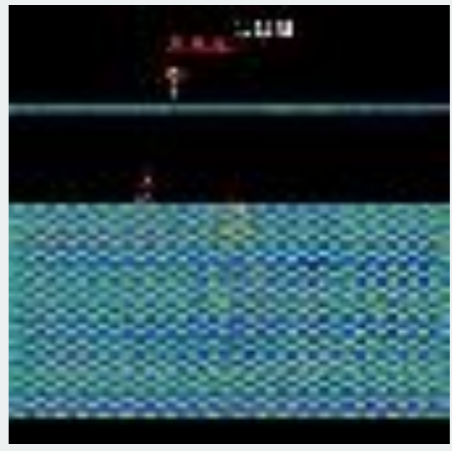
GENERATIVE MODEL



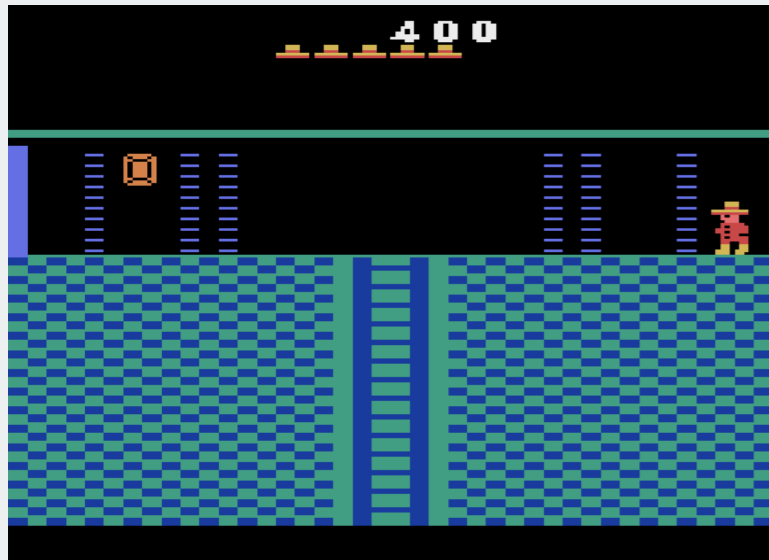
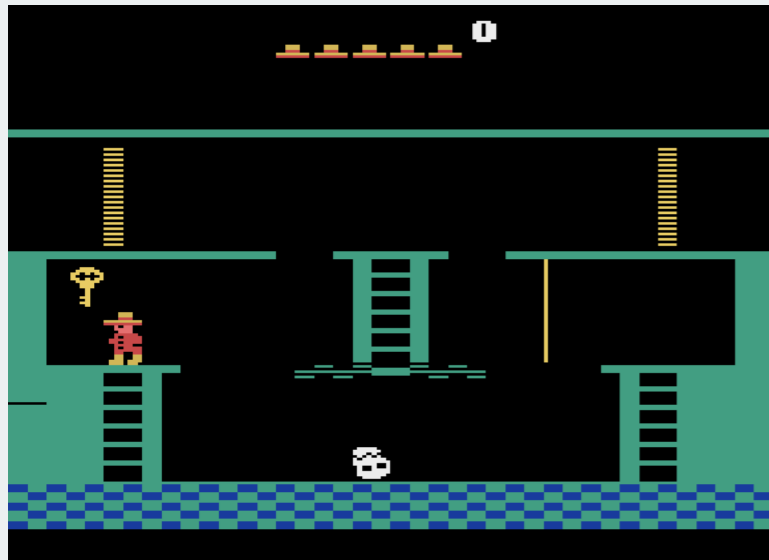
Train



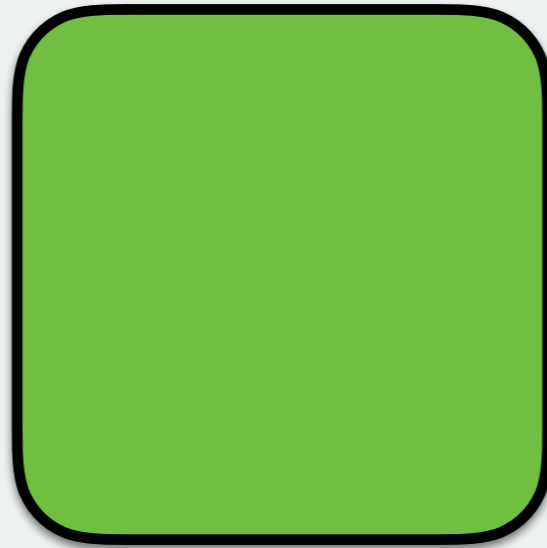
Generate



DENSITY MODEL



Train



Query

$$\rho_n(x)$$

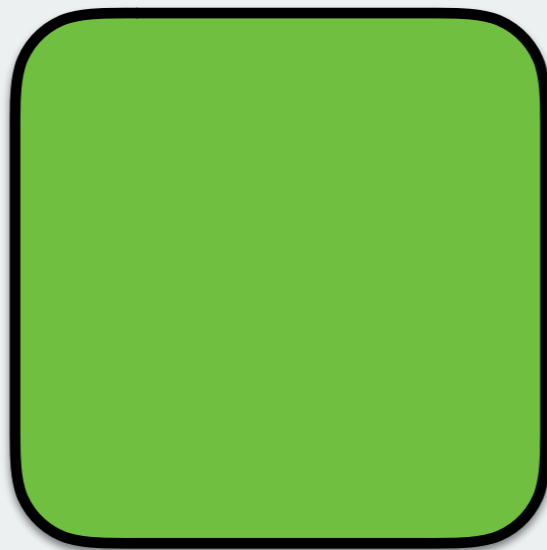
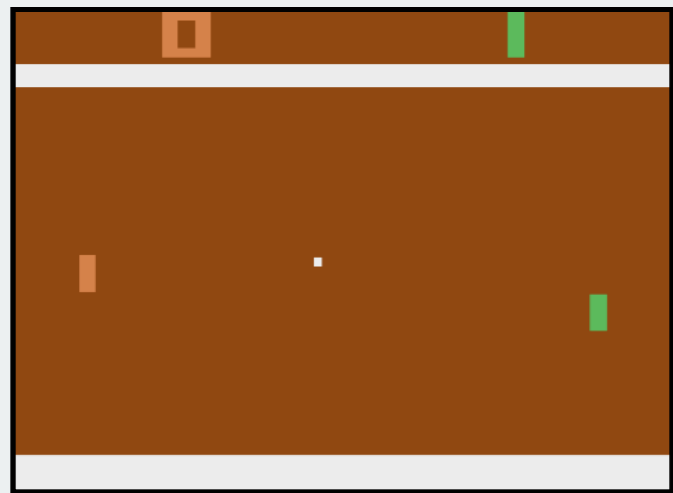
Assumptions

1. density \propto frequency
2. model quality = generalization

$$\rho_n(x) = \frac{\hat{N}_n(x)}{\hat{n}} \quad \rho'_n(x) = \frac{\hat{N}_n(x) + 1}{\hat{n} + 1}$$

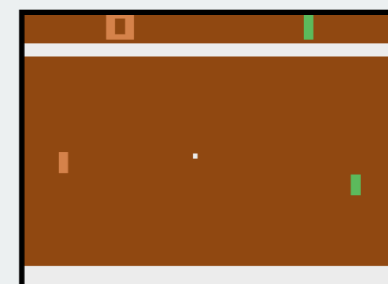


$$\hat{N}_n(x) = \rho_n(x) \frac{1 - \rho'_n(x)}{\rho'_n(x) - \rho_n(x)}$$

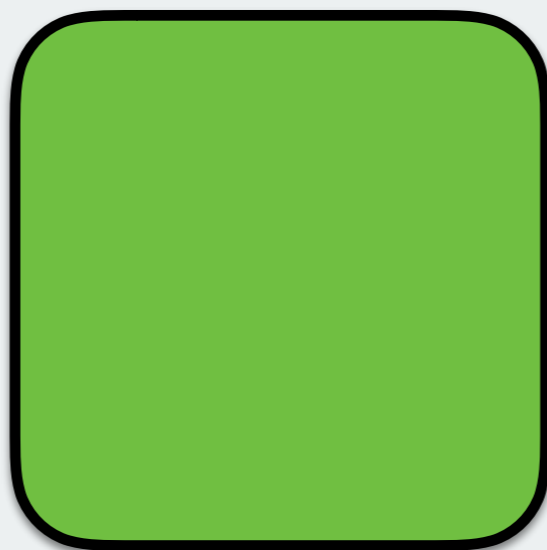
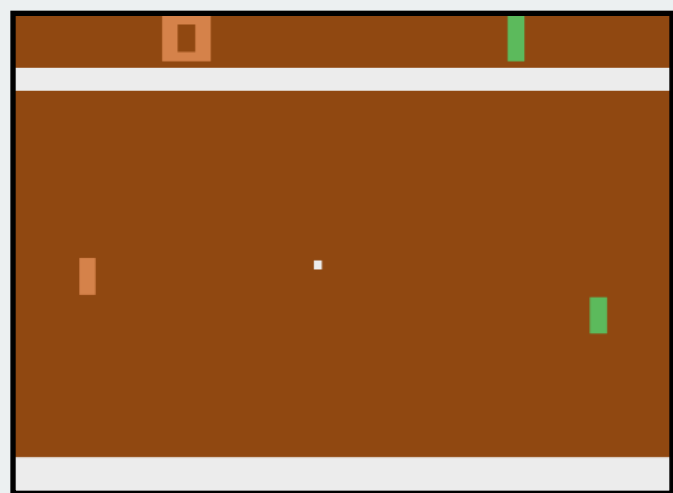


CODING PROBABILITY
(DENSITY)

$$\rho_n(x)$$



TRAIN

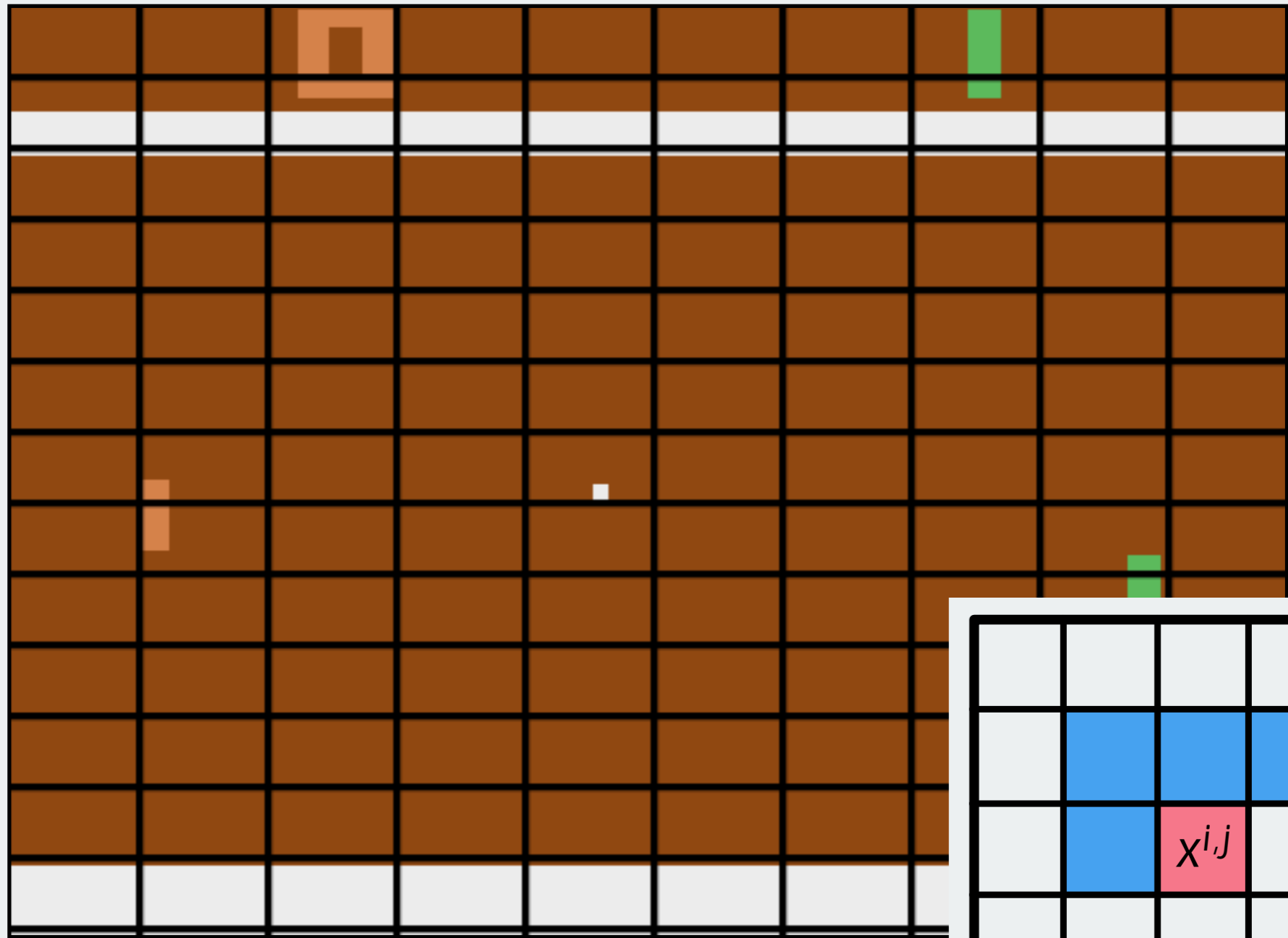


$$\rho'_n(x)$$

RECODING
PROBABILITY

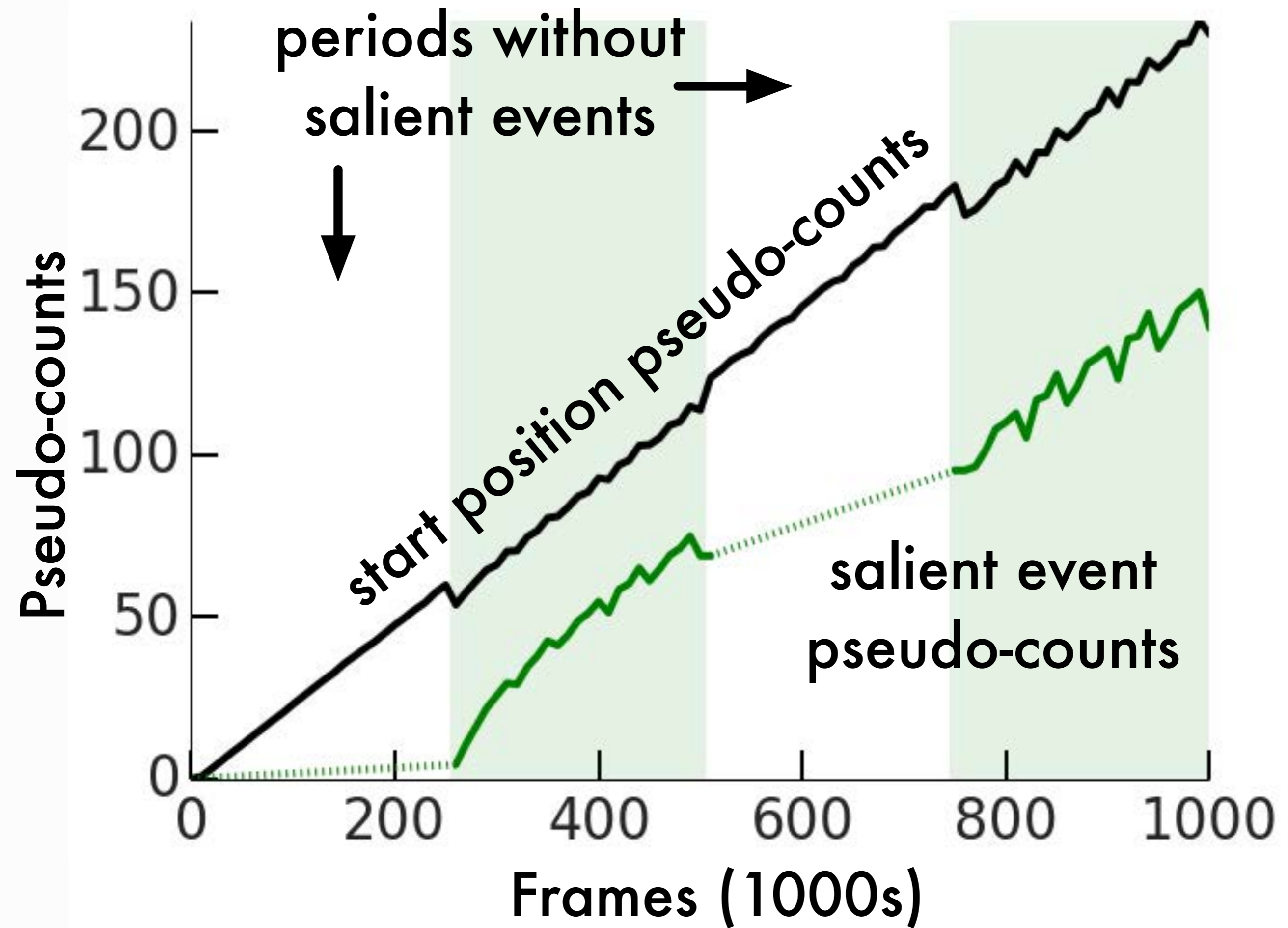
THE "CTS" MODEL

$$\rho_n(x) := \prod_{i=1}^K \rho_n^i(x^i | x^{<i})$$





ACTIVISION



EXPLORATION

Bellman equation:

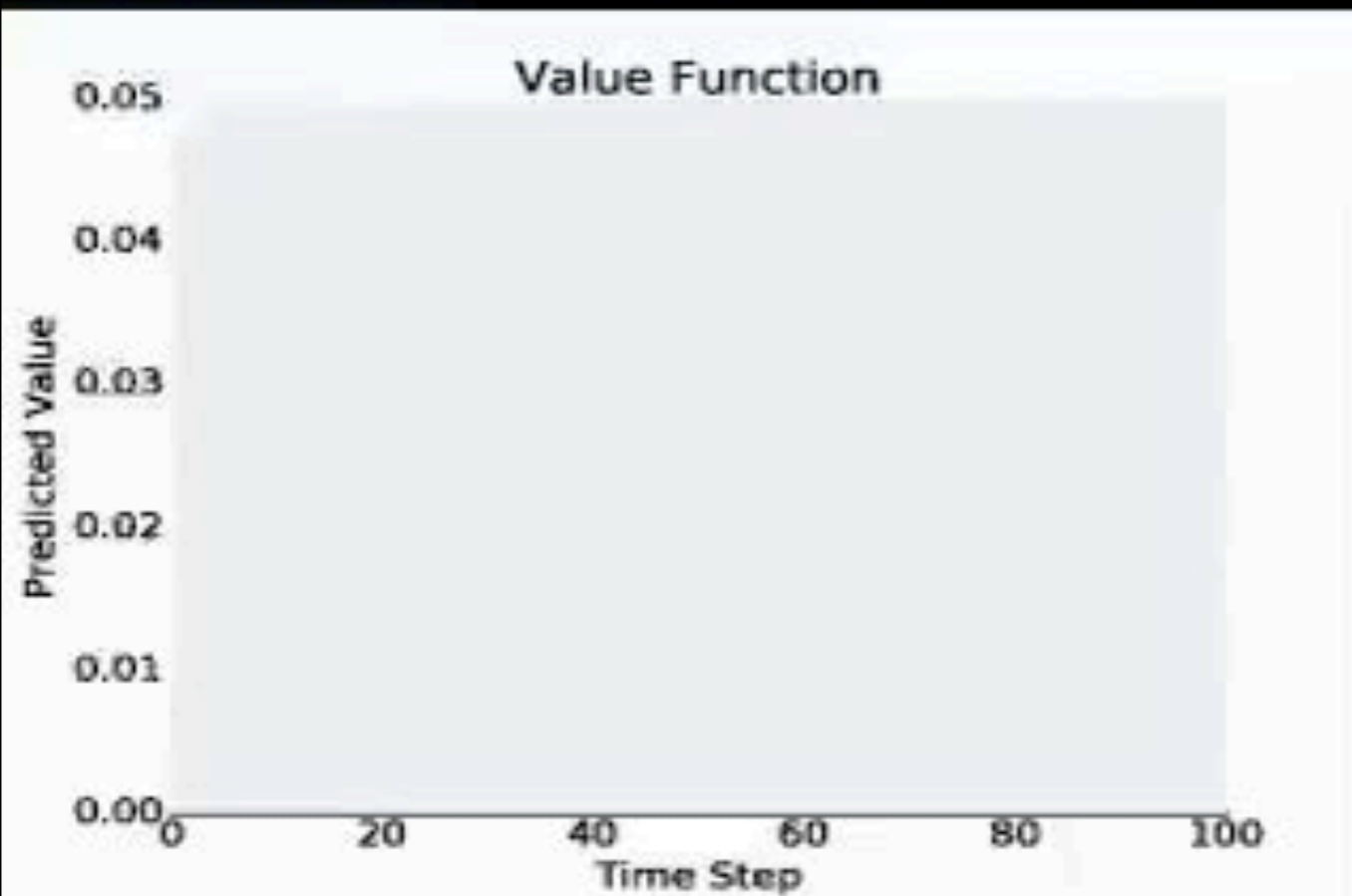
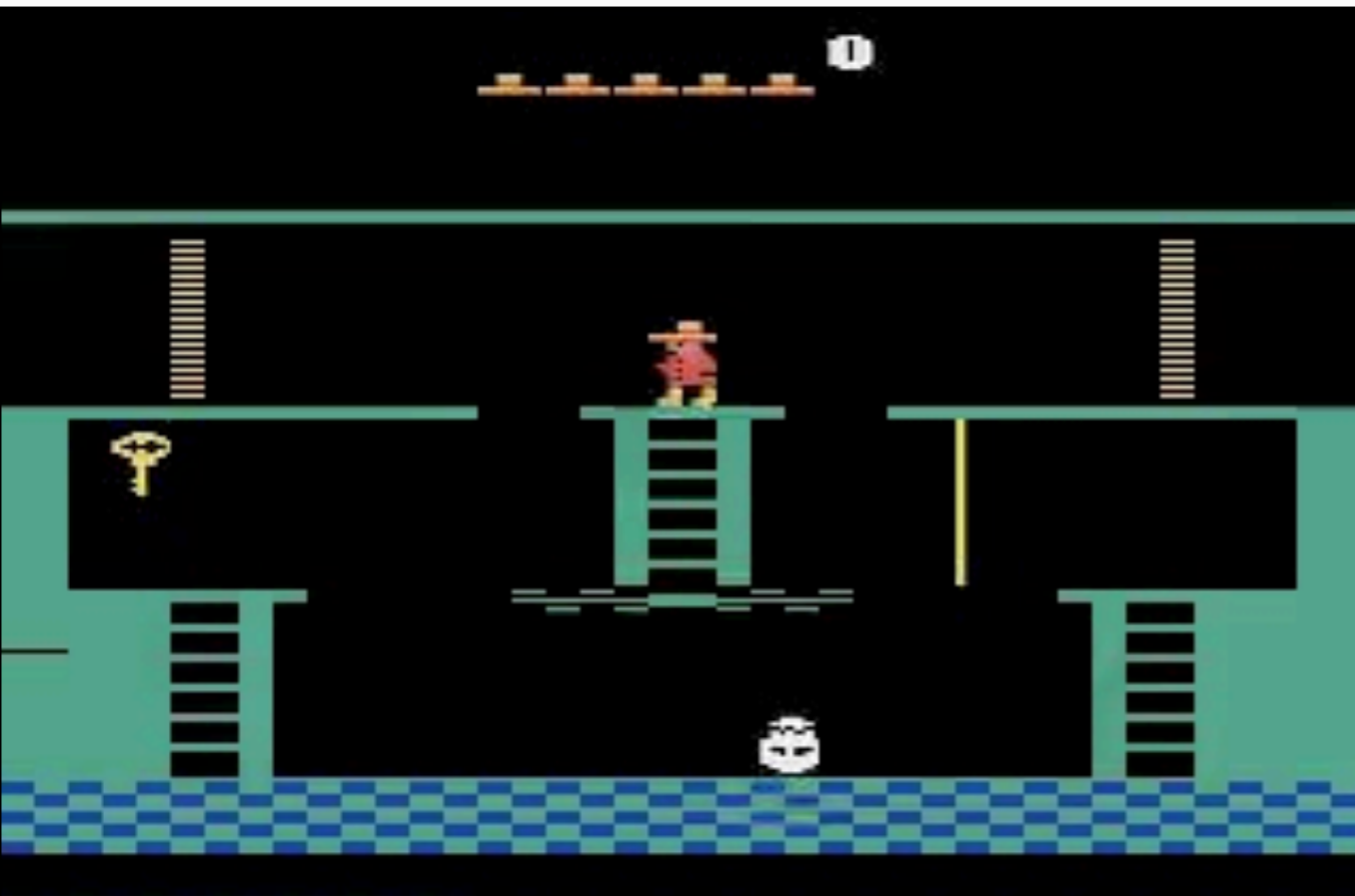
$$Q(x, a) = r(x, a) + \gamma \mathbf{E}_{x' \sim P} \max_{a' \in \mathcal{A}} Q(x', a')$$

Exploration bonus (Strehl and Littman, 2008)

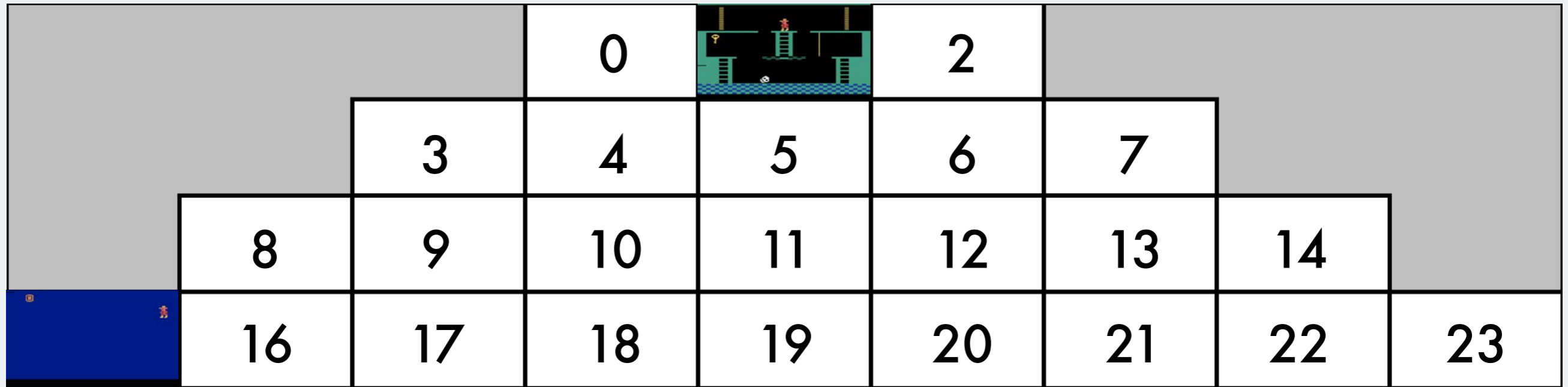
$$Q(x, a) = \hat{r}(x, a) + \gamma \mathbf{E}_{x' \sim \hat{P}} \max_{a' \in \mathcal{A}} Q(x', a') + \frac{\beta}{\sqrt{N(x, a)}}$$

Pseudo-count bonus

$$Q(x, a) = \hat{r}(x, a) + \gamma \mathbf{E}_{x' \sim \hat{P}} \max_{a' \in \mathcal{A}} Q(x', a') + \frac{\beta}{\sqrt{\hat{N}(x)}}$$

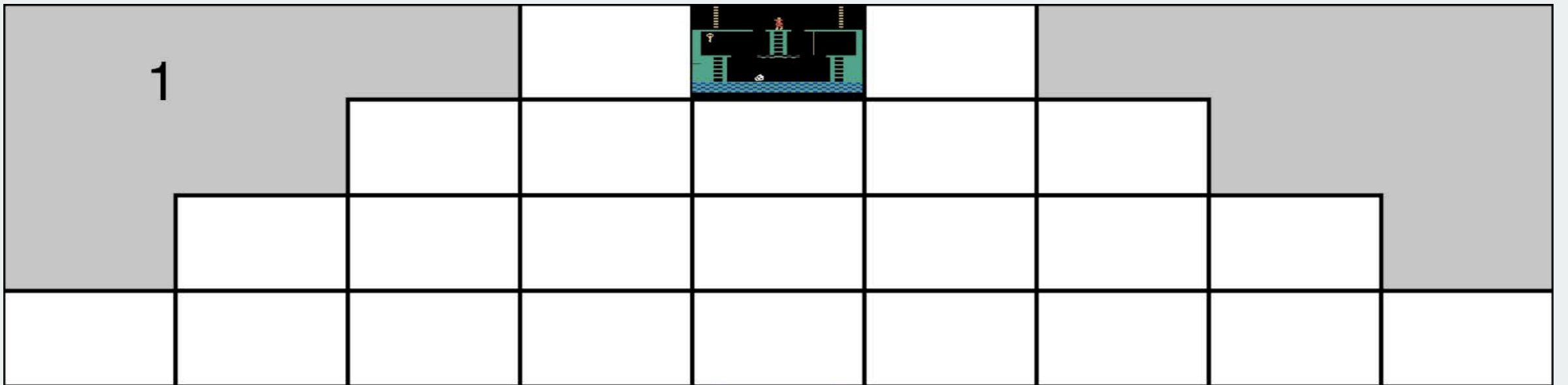


START

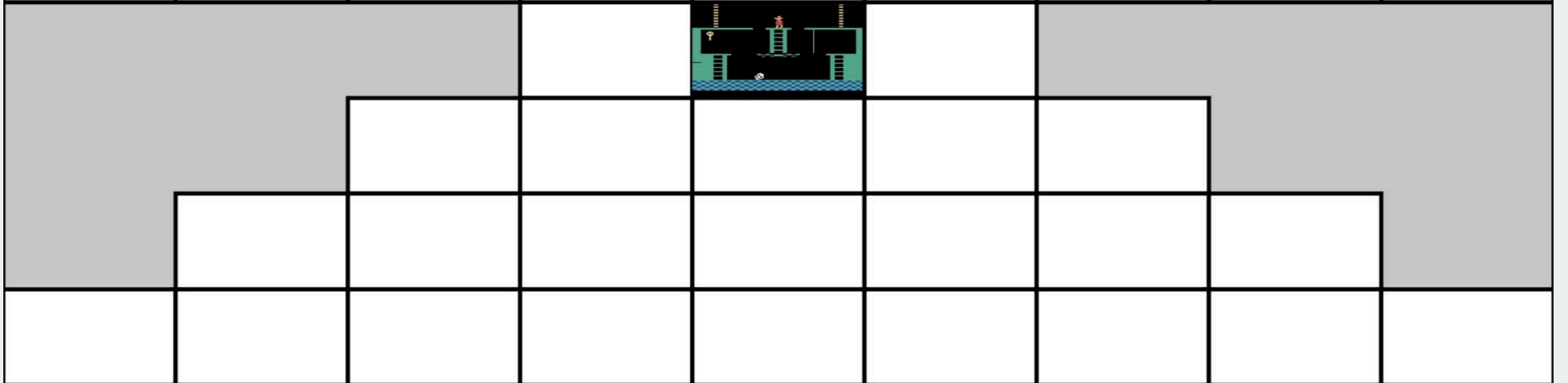


END

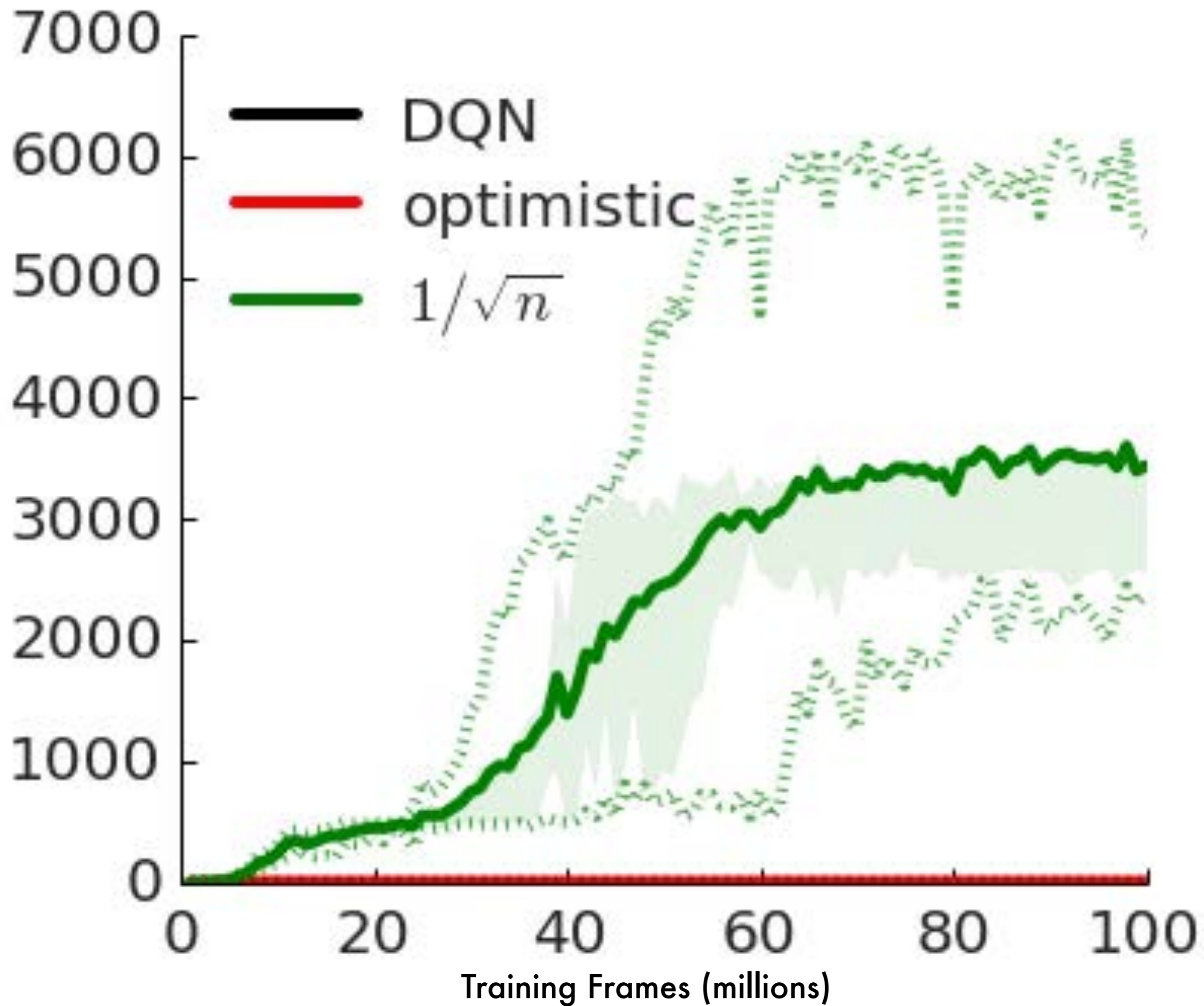
DQN



DQN
+ bonus



Average Score



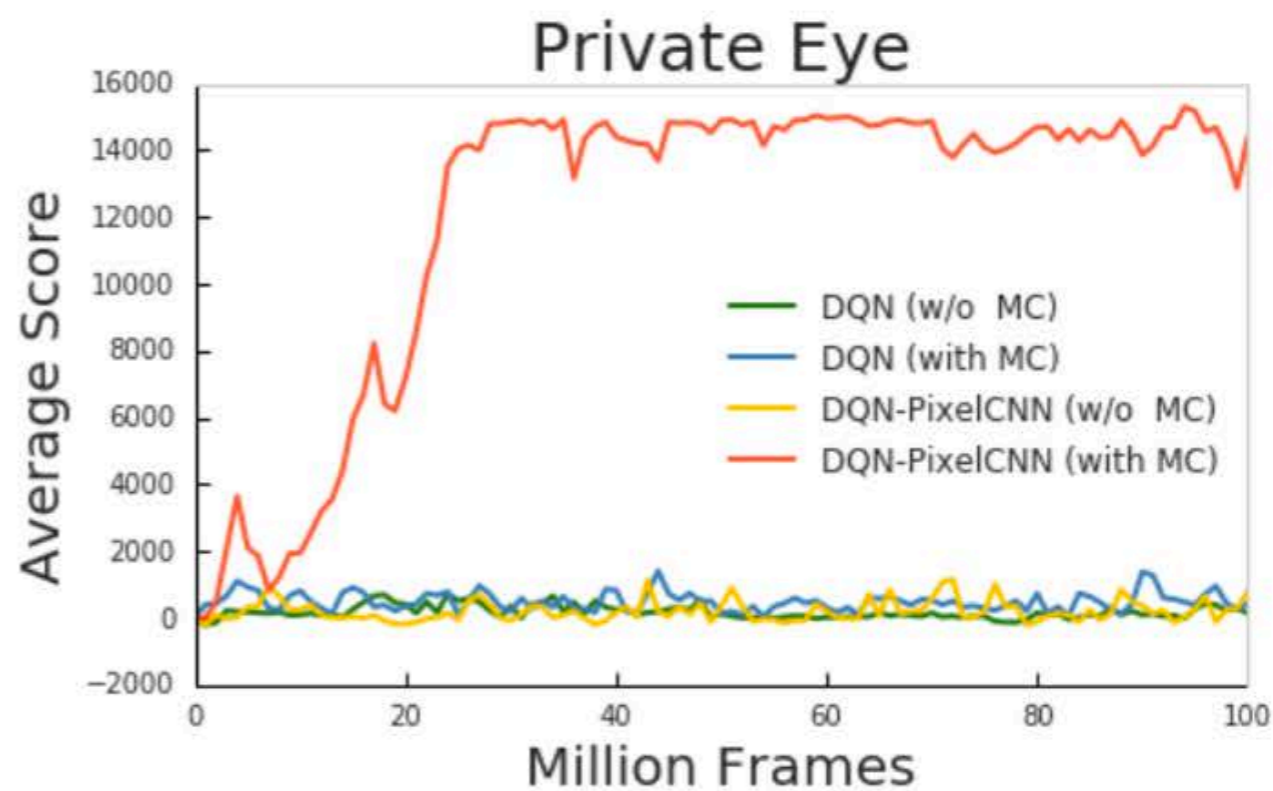
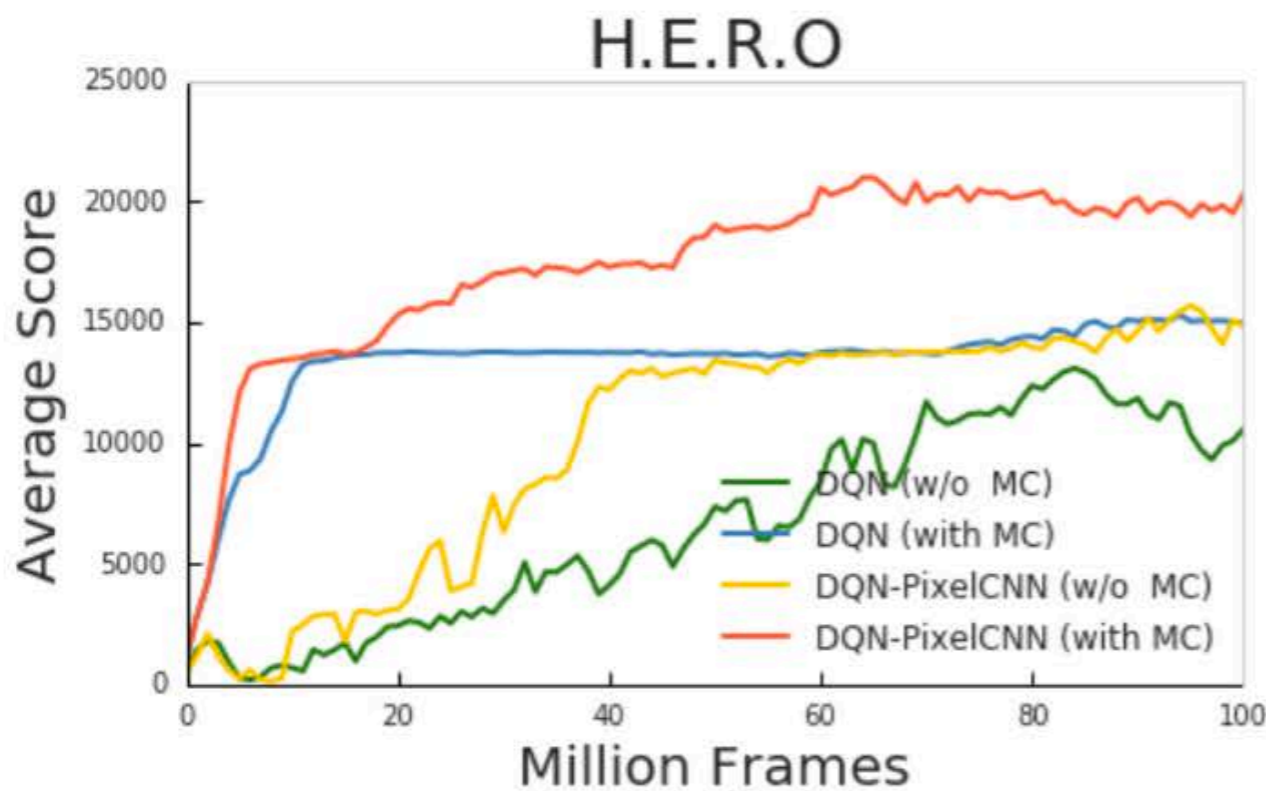
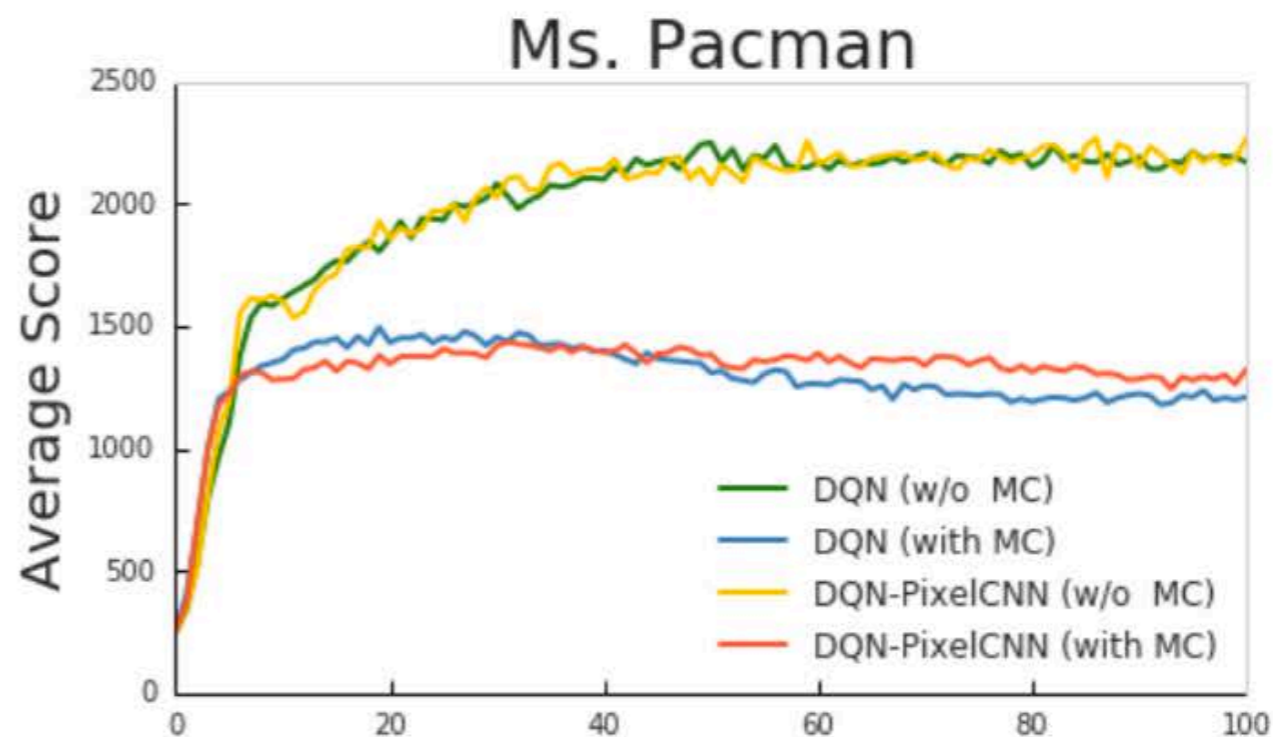
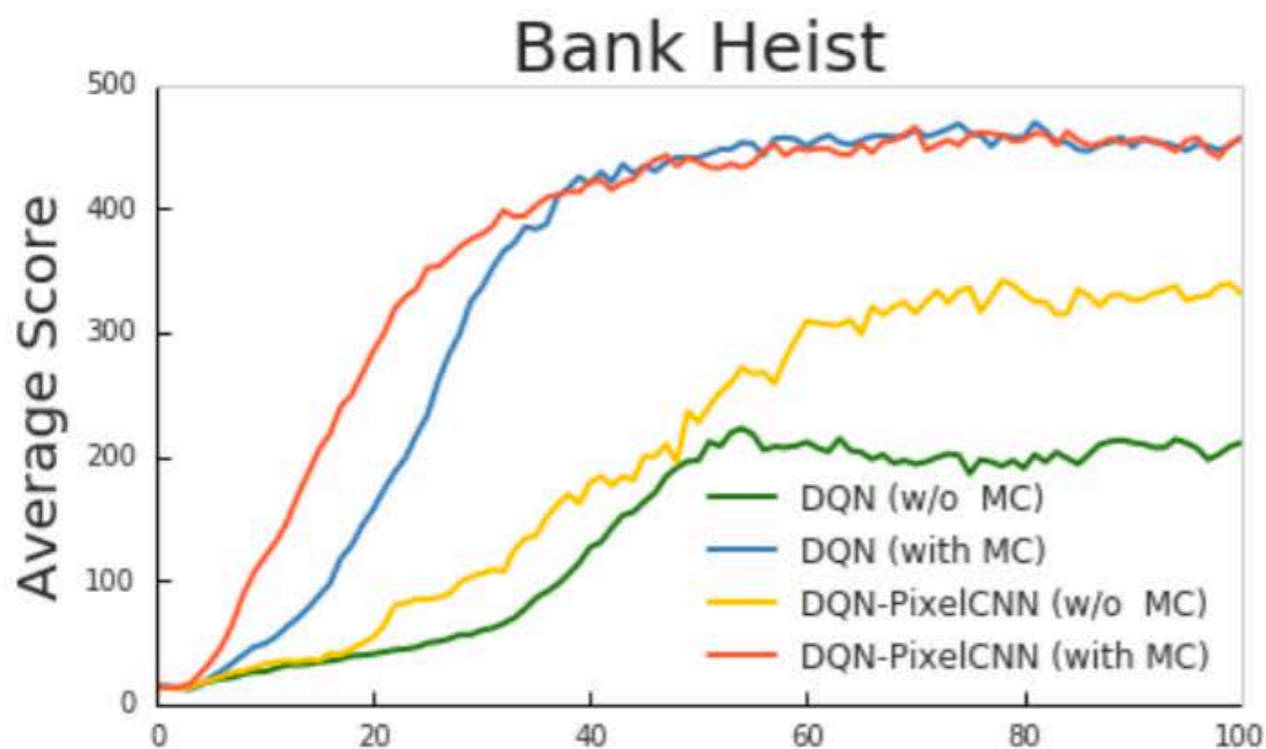
CREDIT ASSIGNMENT ISSUES IN EXPLORATION

- Another important key factor: **fast credit assignment**

$$\tilde{Q}(x_t, a_t) = q \times \left[r_+(x_t, a_t) + \gamma \max_{a' \in \mathcal{A}} Q(x_{t+1}, a') \right] \\ + (1 - q) \times \left[\sum_{i=0}^{\infty} \gamma^i r_+(x_{t+i}, a_{t+i}) \right]$$

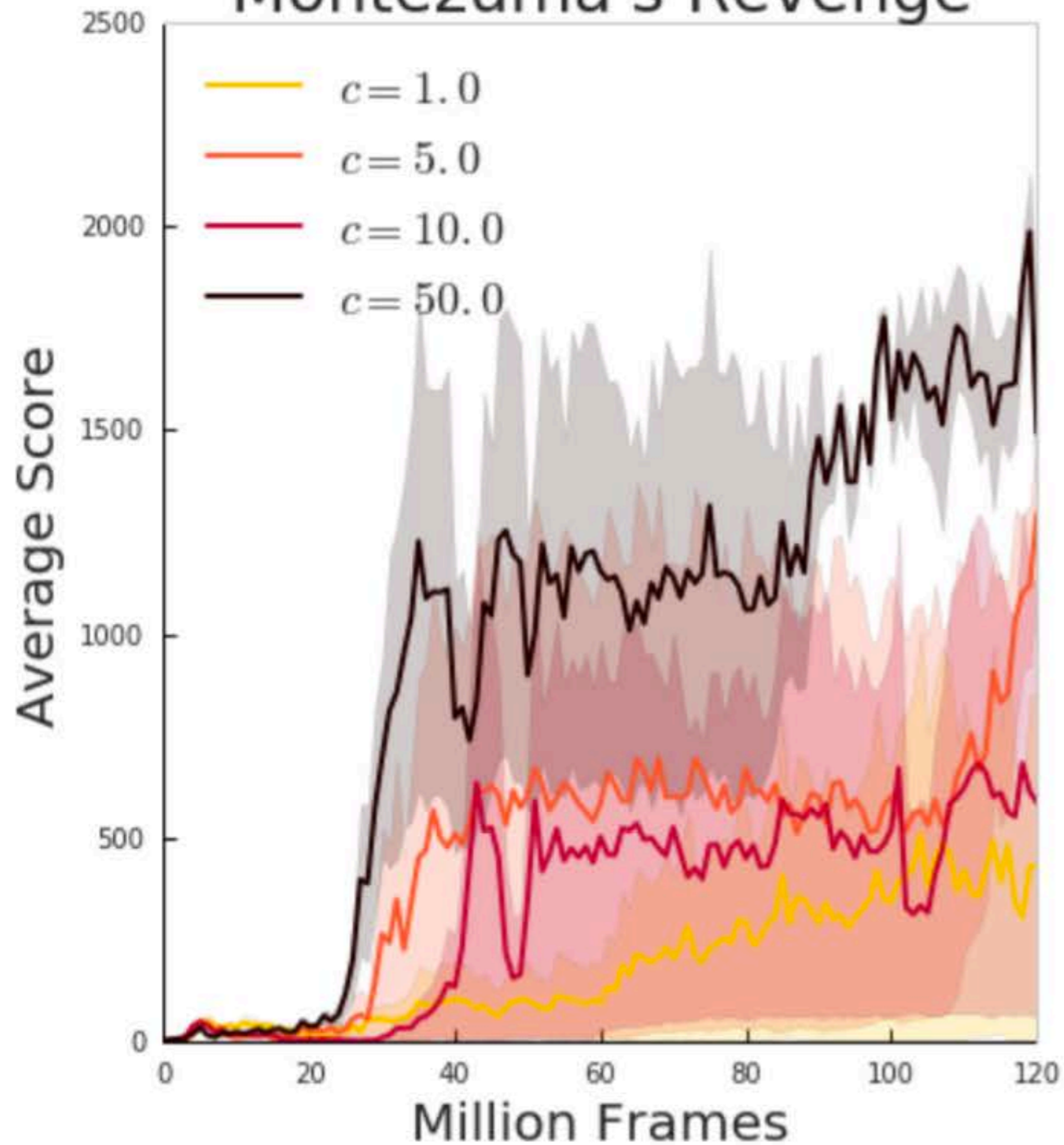
Mixed Monte-Carlo (MMC) update

EFFECT OF MIXED MONTE CARLO UPDATE

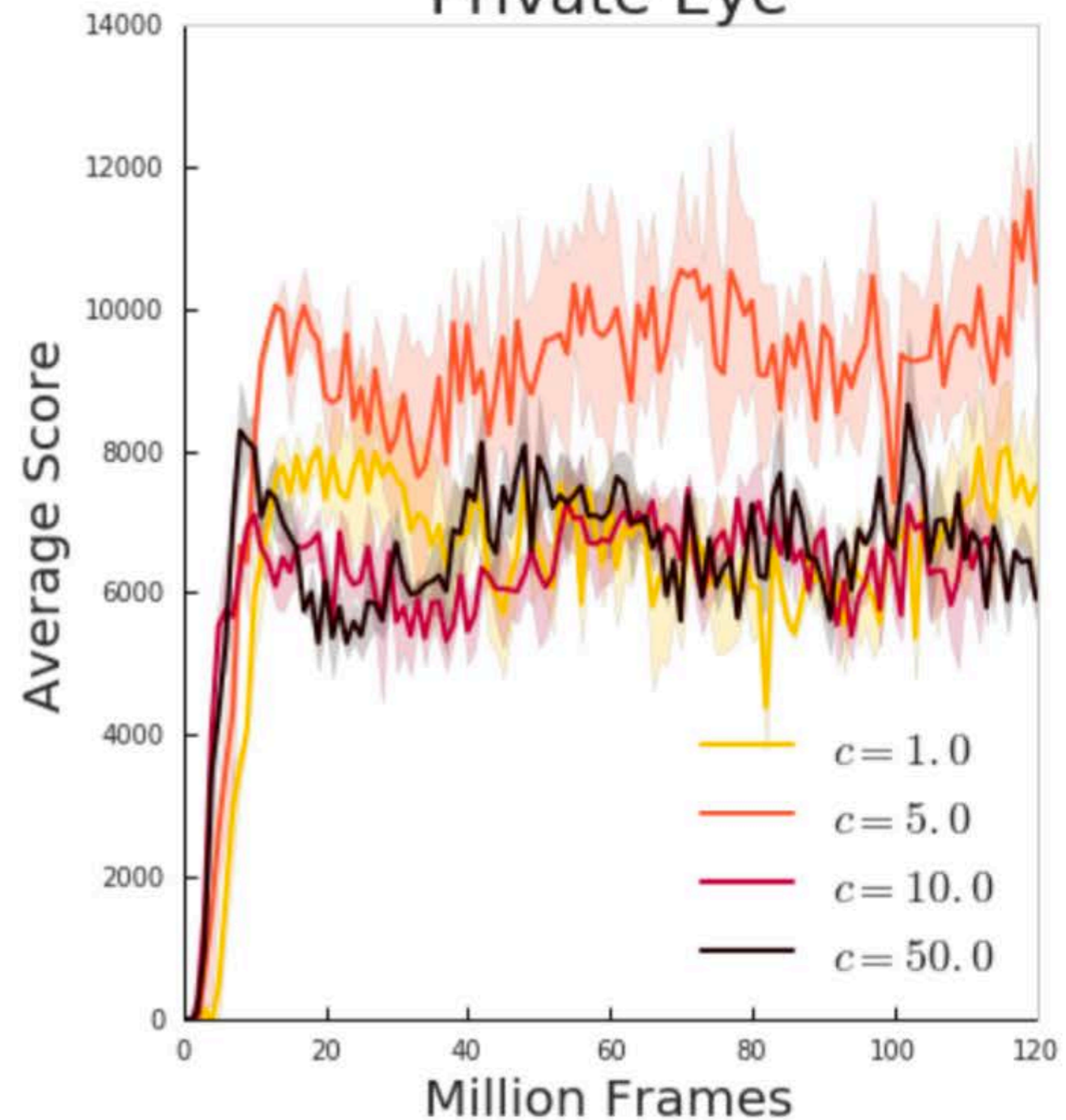


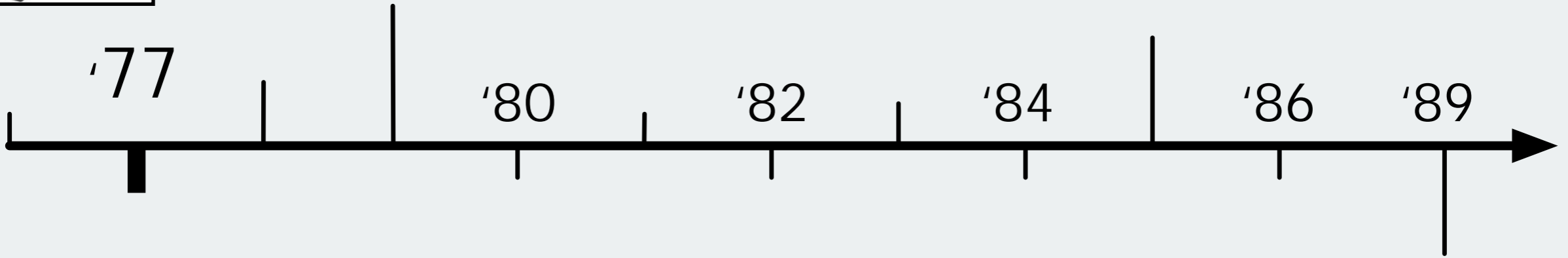
REMOVING EXTRINSIC REWARDS

Montezuma's Revenge



Private Eye





- Super-human agents
- Scoring exploit
- Sub-human agents



■ Breakout '77

■ Space Invaders

■ Asteroids
■ Bowling

■ Freeway
■ Ice Hockey
■ Tennis
□ Ms. Pac-Man
□ Venture

■ Assault
■ Asterix
■ Beam Rider
■ Elevator Action
■ Enduro
■ James Bond
■ Seaquest
■ Up N' Down
■ Bank Heist
■ Krull
□ Frostbite
□ Private Eye
□ Q*Bert

■ Road Runner

'80

'82

'84

'86

'89

■ Pong
■ Video Pinball



■ Battlezone
■ Berzerk
■ Boxing
■ Carnival
■ Centipede
■ Crazy Climber
■ Fishing Derby
■ Phoenix
□ Skiing

■ Amidar
■ Atlantis
■ Chopper Cmd.
■ Demon Attack
■ Gopher
■ Name this Game
■ Pooyan
■ River Raid
■ Star Gunner
■ Time Pilot
■ Yar's Revenge
■ Zaxxon
■ Tutankham
□ Alien
□ Gravitar
□ Journey Escape
□ Kangaroo
□ Pitfall!

■ Kung-Fu Master
□ H.E.R.O.
□ Montezuma's Revenge

□ Solaris

■ Double Dunk

Deep Reinforcement Learning and the Atari 2600

MARC G. BELLEMARE
Google Brain

