

facebook

Artificial Intelligence Research

Exploration-Exploitation in Reinforcement Learning

Part 3 – Scaling up Exploration to DeepRL

Mohammad Ghavamzadeh, Alessandro Lazaric and Matteo Pirotta

Facebook AI Research

Outline

1 Optimistic Exploration in Deep RL

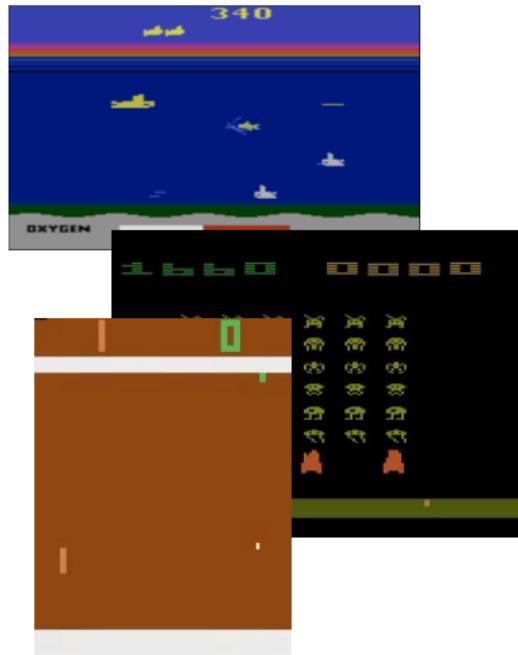
2 Random Exploration in Deep RL

Website

<https://rlgammazero.github.io>

Exploration in DeepRL

these are easy



this is hard, *almost* impossible

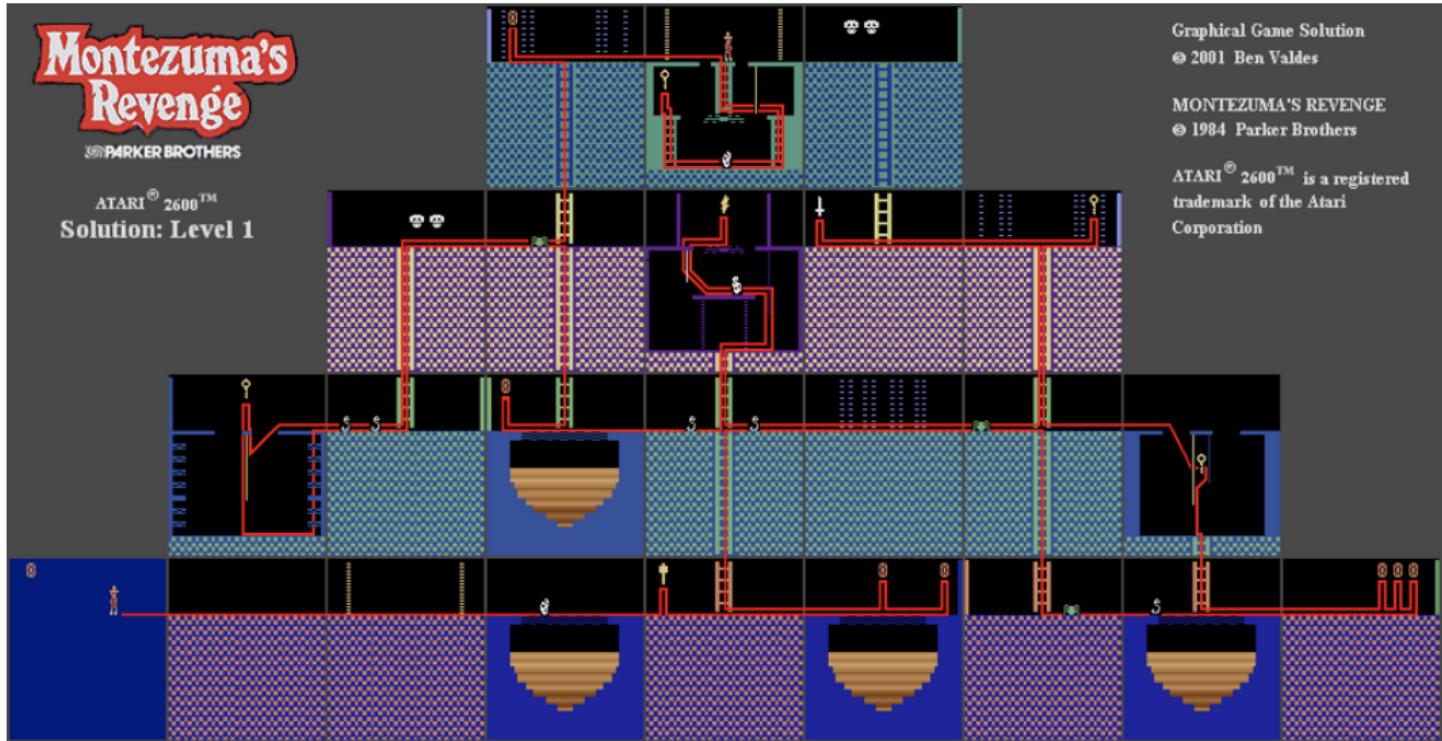


Why?

Random exploration sometimes work!
PONG GIF

Montezuma with random actions!
[Link](#)

Montezuma's Revenge: Level 1



The Four Ingredients Recipe

- 1 Build accurate estimators
- 2 Evaluate the uncertainty of the prediction
- 3 Define a mechanism to combine estimation and uncertainty
- 4 Execute the best policy

The Four Ingredients Recipe

Optimism in face of uncertainty

- 1 Build accurate estimators

$$\widehat{M}_k \Rightarrow V_{\widehat{M}_k}^\pi$$

- 2 Evaluate the uncertainty of the estimators

$$B_{hk}^r(s, a) := \left[\widehat{r}_{hk}(s, a) - \beta_{hk}^r(s, a), \widehat{r}_{hk}(s, a) + \beta_{hk}^r(s, a) \right]$$

$$B_{hk}^p(s, a) := \left\{ p(\cdot|s, a) \in \Delta(\mathcal{S}) : \|p(\cdot|s, a) - \widehat{p}_{hk}(\cdot|s, a)\|_1 \leq \beta_{hk}^p(s, a) \right\}$$

- 3 Define a mechanism to combine estimation and uncertainty

$$\pi_k = \arg \max_{\pi} \max_{M \in \mathcal{M}_k} V_M^\pi$$

The Four Ingredients Recipe

Posterior Sampling

- 1 Build accurate estimators
- 2 Evaluate the uncertainty of the estimators

$\forall \Theta, \quad \mathbb{P}(M^* \in \Theta | H_t, \mu_1) = \mu_t(\Theta) \quad \mu_t \text{ updated using Bayes' rule}$

- 3 Define a mechanism to combine estimation and uncertainty

$$\pi_k = \arg \max_{\pi} V_{\widetilde{M}_k}^{\pi}, \quad \widetilde{M}_k \sim \mu_k$$

“Practical” Limitations

Optimism in face of uncertainty

- Confidence intervals

$$\beta_t^r(s, a) \propto \sqrt{\frac{\log(N_t(s, a)/\delta)}{N_t(s, a)}} \quad \beta_t^p(s, a) \propto \sqrt{\frac{S \log(N_t(s, a)/\delta)}{N_t(s, a)}}$$

- Solving

$$\pi_t = \arg \max_{\pi} \max_{M \in \mathcal{M}_t} V_M^\pi$$

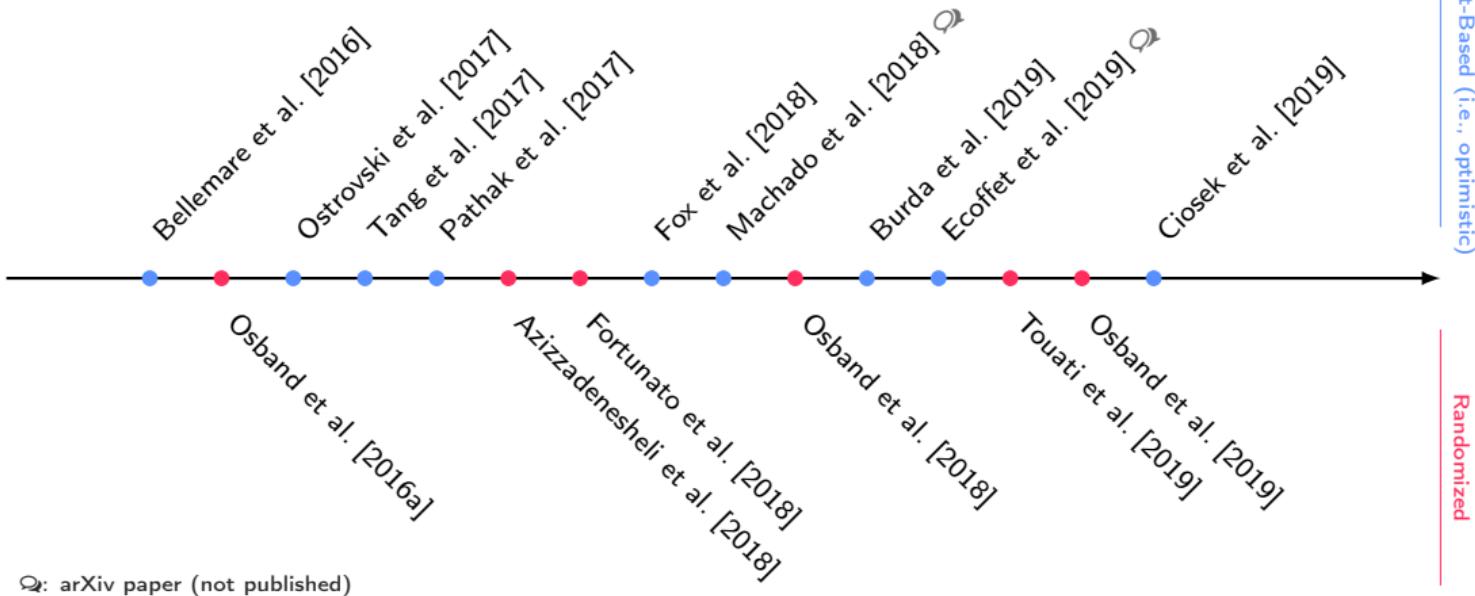
Posterior sampling

- Posterior (dynamics for any state-action pair)

$$\text{Dirichlet}\left(N_t(s'_1|s, a), N_t(s'_2|s, a), \dots, N_t(s'_S|s, a)\right)$$

- Update/sample from a unstructured/non-conjugate posteriors

History: Exploration in DeepRL



Outline

1 Optimistic Exploration in Deep RL

2 Random Exploration in Deep RL

Count-based Exploration

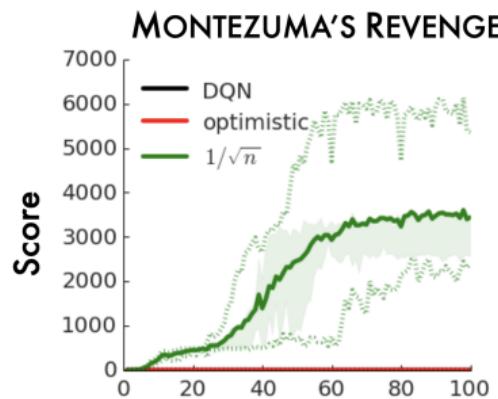
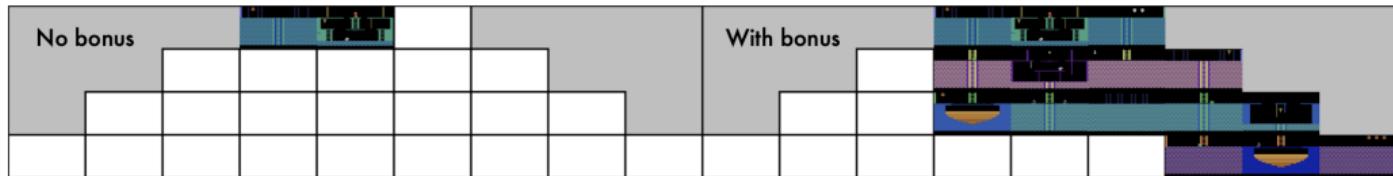
General Scheme

- 1 Estimate a “proxy” for the *number of visits* $\tilde{N}(s_t)$
- 2 Add an *exploration bonus* to the rewards

$$\tilde{r}_t^+ = r_t + \textcolor{red}{c} \sqrt{\frac{1}{\tilde{N}(s_t)}}$$

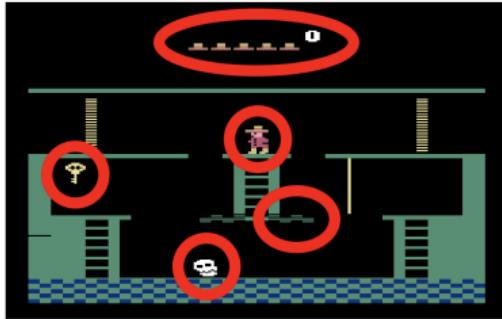
- 3 Run *any DeepRL* algorithm on $\mathcal{D}_t = \{(s_i, a_i, \tilde{r}_i^+, s_{i+1})\}$

Does it work?



* figures from [Bellemare et al., 2016]

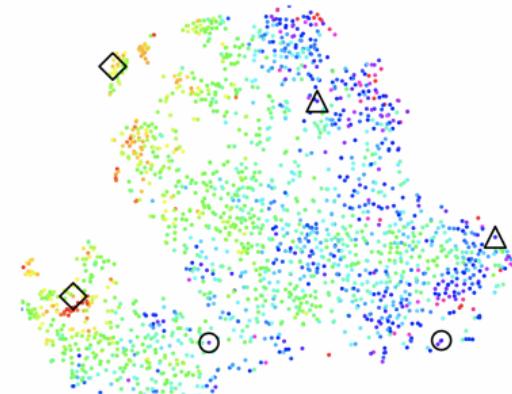
What to Count?



We never see the same state twice (or it is very unlikely)!

How difficult is to learn a state representation?

[Sun et al., 2019]



O: $V = 6.27$

3210

O: $V = 6.14$

3100

△: $V = 6.17$

3545

△: $V = 6.16$

2900

◇: $V = 4.44$

1960

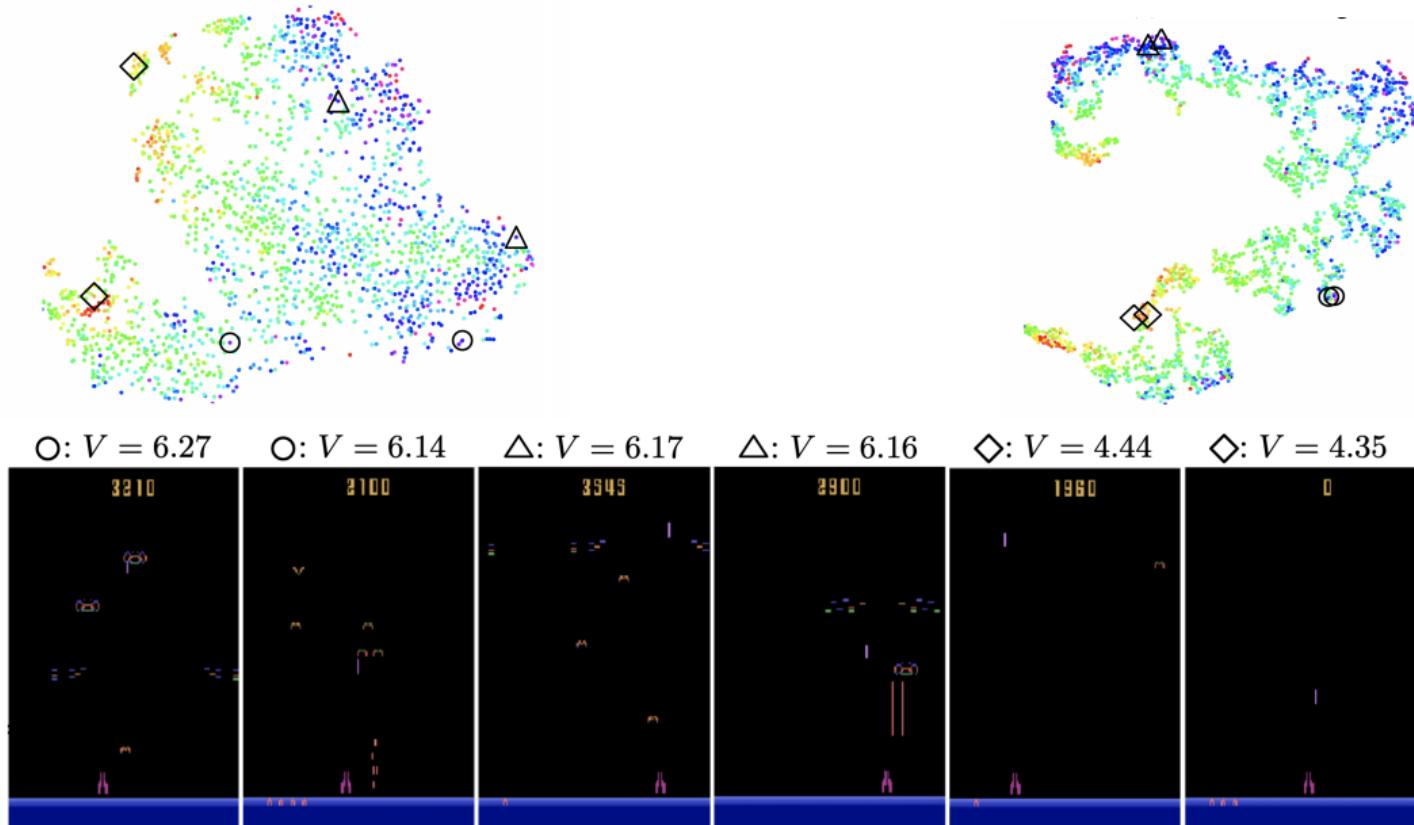
◇: $V = 4.35$

0



How difficult is to learn a state representation?

[Sun et al., 2019]



Count-based Exploration

[Tang et al., 2017]

Algorithm 1: Count-based exploration through static hashing, using SimHash

- 1 Define state preprocessor $g : \mathcal{S} \rightarrow \mathbb{R}^D$
 - 2 (In case of SimHash) Initialize $A \in \mathbb{R}^{k \times D}$ with entries drawn i.i.d. from the standard Gaussian distribution $\mathcal{N}(0, 1)$
 - 3 Initialize a hash table with values $n(\cdot) \equiv 0$
 - 4 **for** each iteration j **do**
 - 5 Collect a set of state-action samples $\{(s_m, a_m)\}_{m=0}^M$ with policy π
 - 6 Compute hash codes through any LSH method, e.g., for SimHash, $\phi(s_m) = \text{sgn}(Ag(s_m))$
 - 7 Update the hash table counts $\forall m : 0 \leq m \leq M$ as $n(\phi(s_m)) \leftarrow n(\phi(s_m)) + 1$
 - 8 Update the policy π using rewards $\left\{ r(s_m, a_m) + \frac{\beta}{\sqrt{n(\phi(s_m))}} \right\}_{m=0}^M$ with any RL algorithm
-

- Use *locality-sensitive hashing* to discretize the input
 - Encode the state into a k -dim vector by random project
small $k =$ more hash collisions
 - Use the sign to discretize
small $\phi(s) \in \{-1, 1\}^k$
- *Count on discrete hashed-states*

Count-based Exploration

[Tang et al., 2017]

Algorithm 1: Count-based exploration through static hashing, using SimHash

- 1 Define state preprocessor $g : \mathcal{S} \rightarrow \mathbb{R}^D$
 - 2 (In case of SimHash) Initialize $A \in \mathbb{R}^{k \times D}$ with entries drawn i.i.d. from the standard Gaussian distribution $\mathcal{N}(0, 1)$
 - 3 Initialize a hash table with values $n(\cdot) \equiv 0$
 - 4 **for** each iteration j **do**
 - 5 Collect a set of state-action samples $\{(s_m, a_m)\}_{m=0}^M$ with policy π
 - 6 Compute hash codes through any LSH method, e.g., for SimHash, $\phi(s_m) = \text{sgn}(Ag(s_m))$
 - 7 Update the hash table counts $\forall m : 0 \leq m \leq M$ as $n(\phi(s_m)) \leftarrow n(\phi(s_m)) + 1$
 - 8 Update the policy π using rewards $\left\{ r(s_m, a_m) + \frac{\beta}{\sqrt{n(\phi(s_m))}} \right\}_{m=0}^M$ with any RL algorithm
-

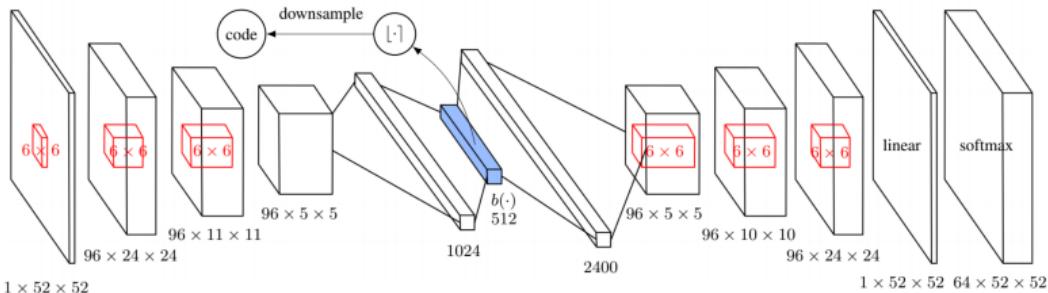
- Use *locality-sensitive hashing* to discretize the input
 - Encode the state into a k -dim vector by random project
small $k =$ more hash collisions
 - Use the sign to discretize
small $\phi(s) \in \{-1, 1\}^k$
- *Count on discrete hashed-states*

👎 Difficult to define a good hashing function

Count-based Exploration

[Tang et al., 2017]

Improve counts by learning a compression



$$L(\{s_n\}_{n=1}^N) = -\frac{1}{N} \sum_{n=1}^N \left[\log p(s_n) - \frac{\lambda}{K} \sum_{i=1}^D \min \left\{ (1 - b_i(s_n))^2, b_i(s_n)^2 \right\} \right]$$

- Entropy loss for the auto-encoder
- “Binarization” loss for the “projection”

Count-based Exploration

[Tang et al., 2017]

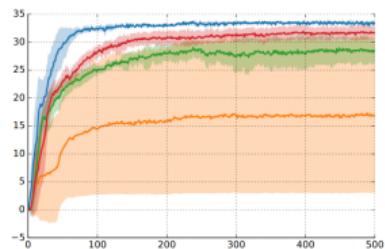
Algorithm 2: Count-based exploration using learned hash codes

- 1 Define state preprocessor $g : \mathcal{S} \rightarrow \{0, 1\}^D$ as the binary code resulting from the autoencoder (AE)
 - 2 Initialize $A \in \mathbb{R}^{k \times D}$ with entries drawn i.i.d. from the standard Gaussian distribution $\mathcal{N}(0, 1)$
 - 3 Initialize a hash table with values $n(\cdot) \equiv 0$
 - 4 **for** each iteration j **do**
 - 5 Collect a set of state-action samples $\{(s_m, a_m)\}_{m=0}^M$ with policy π
 - 6 Add the state samples $\{s_m\}_{m=0}^M$ to a FIFO replay pool \mathcal{R}
 - 7 **if** $j \bmod j_{\text{update}} = 0$ **then**
 - 8 Update the AE loss function in Eq. (3) using samples drawn from the replay pool
 $\{s_n\}_{n=1}^N \sim \mathcal{R}$, for example using stochastic gradient descent
 - 9 Compute $g(s_m) = \lfloor b(s_m) \rfloor$, the D -dim rounded hash code for s_m learned by the AE
 - 10 Project $g(s_m)$ to a lower dimension k via SimHash as $\phi(s_m) = \text{sgn}(Ag(s_m))$
 - 11 Update the hash table counts $\forall m : 0 \leq m \leq M$ as $n(\phi(s_m)) \leftarrow n(\phi(s_m)) + 1$
 - 12 Update the policy π using rewards $\left\{ r(s_m, a_m) + \frac{\beta}{\sqrt{n(\phi(s_m))}} \right\}_{m=0}^M$ with any RL algorithm
-

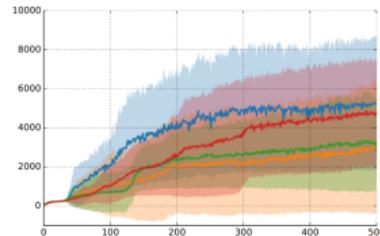
- Use all past history to update the AE
- AE should not be updated too often
we need stable codes!

Count-based Exploration

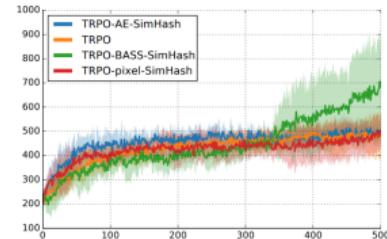
[Tang et al., 2017]



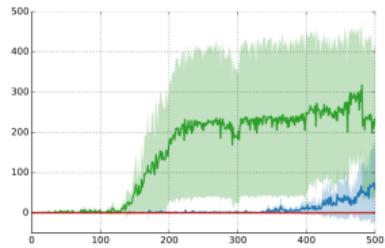
(a) Freeway



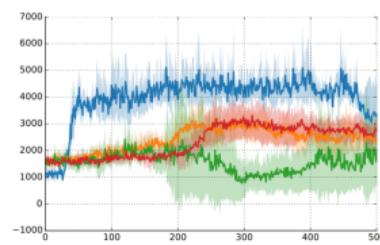
(b) Frostbite



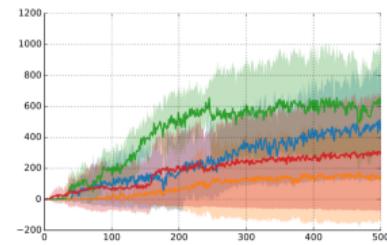
(c) Gravitar



(d) Montezuma's Revenge



(e) Solaris



(f) Venture

Count-based Exploration

[Bellemare et al., 2016, Ostrovski et al., 2017]

- Density estimation over a countable set \mathcal{X}

$$\rho_n(x) = \rho(x|x_1, \dots, x_n) \approx \mathbb{P}[X_{n+1} = x|x_1, \dots, x_n]$$

- Recording probability

$$\rho'_n(x) = \rho(x|x_1, \dots, x_n, x) \approx \mathbb{P}[X_{n+2} = x|x_1, \dots, x_n, X_{n+1} = x]$$

- Pseudo “local” and “total” counts $\tilde{N}_n(x)$ and $\tilde{N}_n(x)$ s.t.

probability of x after observing a new occurrence of x

$$\frac{\tilde{N}_n(x)}{\tilde{n}} = \rho_n(x); \quad \frac{\tilde{N}_n(x) + 1}{\tilde{n} + 1} = \rho'_n(x) \Rightarrow \tilde{N}_n(x) = \frac{\rho_n(x)(1 - \rho'_n(x))}{\rho'_n(x) - \rho_n(x)} = \tilde{n}\rho_n(x)$$

Count-based Exploration

[Bellemare et al., 2016, Ostrovski et al., 2017]

- Density estimation over a countable set \mathcal{X}

$$\rho_n(x) = \rho(x|x_1, \dots, x_n) \approx \mathbb{P}[X_{n+1} = x|x_1, \dots, x_n]$$

- Recording probability

$$\rho'_n(x) = \rho(x|x_1, \dots, x_n, x) \approx \mathbb{P}[X_{n+2} = x|x_1, \dots, x_n, X_{n+1} = x]$$

- Pseudo “local” and “total” counts $\tilde{N}_n(x)$ and $\tilde{N}_n(x)$ s.t.

probability of x after observing a new occurrence of x

$$\frac{\tilde{N}_n(x)}{\tilde{n}} = \rho_n(x); \quad \frac{\tilde{N}_n(x) + 1}{\tilde{n} + 1} = \rho'_n(x) \Rightarrow \tilde{N}_n(x) = \frac{\rho_n(x)(1 - \rho'_n(x))}{\rho'_n(x) - \rho_n(x)} = \tilde{n}\rho_n(x)$$

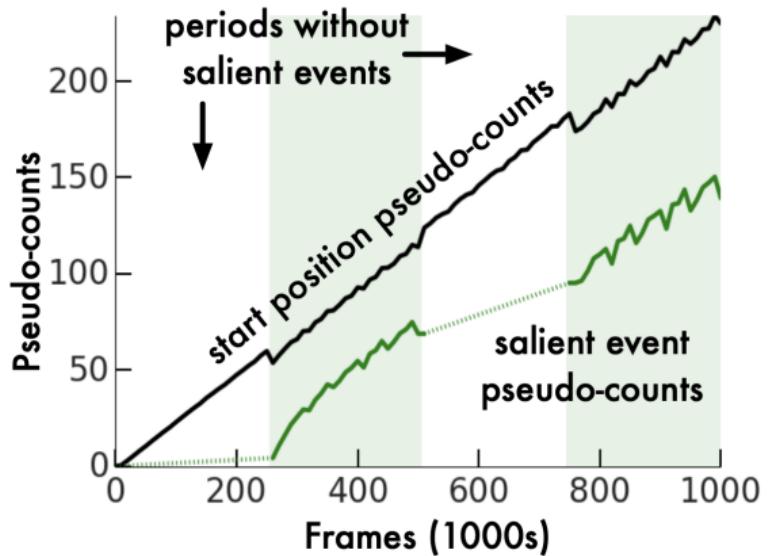
- 👍 Any density estimation algorithm (accurate for images) |

e.g., CTS [Bellemare et al., 2014] or PixelCNN [van den Oord et al., 2016]

- 👎 Density estimation in continuous spaces is hard |

Count-based Exploration

[Bellemare et al., 2016, Ostrovski et al., 2017]



Count-based Exploration

Belle-mare et al. [2016], Ostrovski et al. [2017]

Montezuma!

Prediction-based Exploration

Burda et al. [2019]

What we need is to know how accurate are our predictions.

uncertainty about the model parameters

proxy: prediction error

Sources of prediction errors

- 1 Amount of data 
- 2 Stochasticity (e.g., noisy-TV) 
- 3 Model misspecification 
- 4 Learning dynamics 

Prediction-based Exploration

[Burda et al., 2019]

- Randomly initialize two instances of the same NN (target θ_* and prediction θ_0)

$$f_{\theta_*} : \mathcal{S} \rightarrow \mathbb{R}; \quad f_{\theta} : \mathcal{S} \rightarrow \mathbb{R}$$

- Train the prediction network minimizing loss w.r.t. the target network

$$\theta_n = \arg \min_{\theta} \sum_{t=1}^n \left(f_{\theta}(s_t) - f_{\theta_*}(s_t) \right)^2$$

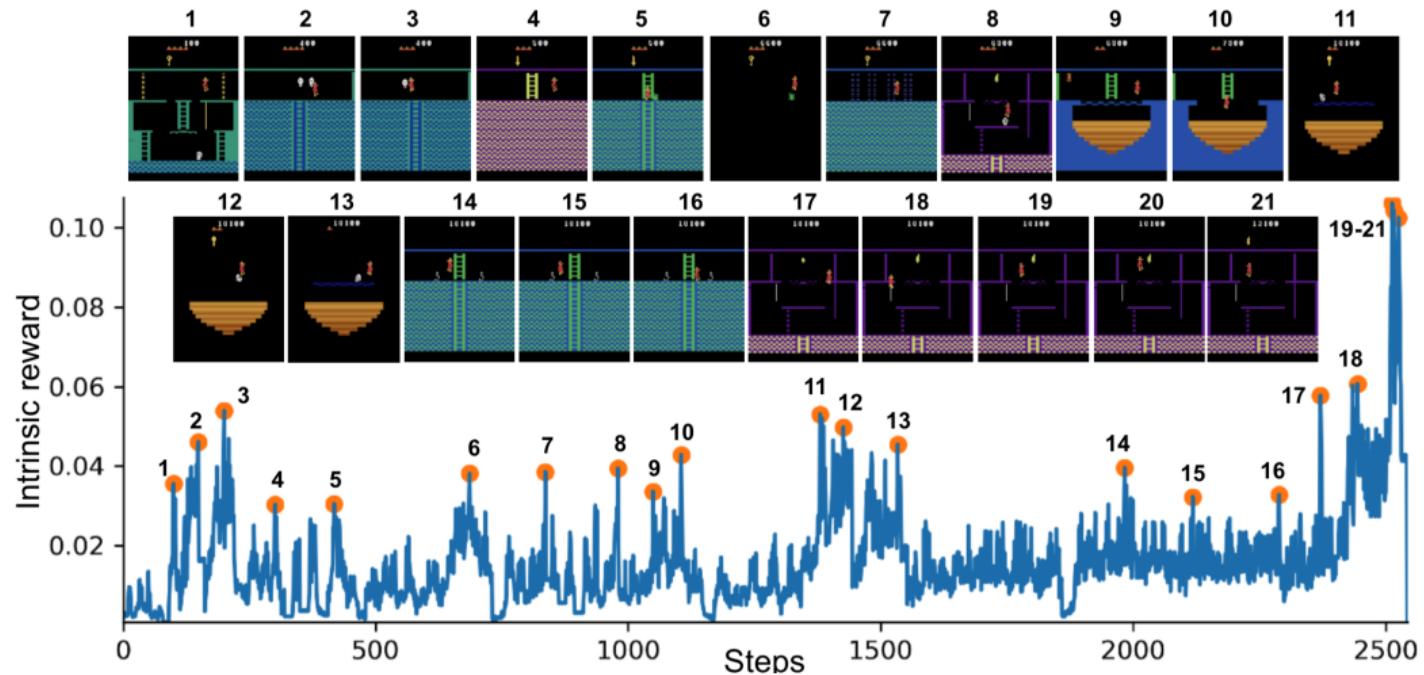
- Build “intrinsic” reward

$$r_t^I = \left| f_{\theta}(s_t) - f_{\theta_*}(s_t) \right|$$

- thumbs up No influence from stochastic transitions
- thumbs up No model misspecification (f_{θ} can exactly predict f_{θ_*})
- thumbs up Influence of learning dynamics can be reduced

Prediction-based Exploration

Burda et al. [2019]



Prediction-based Exploration

[Burda et al., 2019]

General architecture

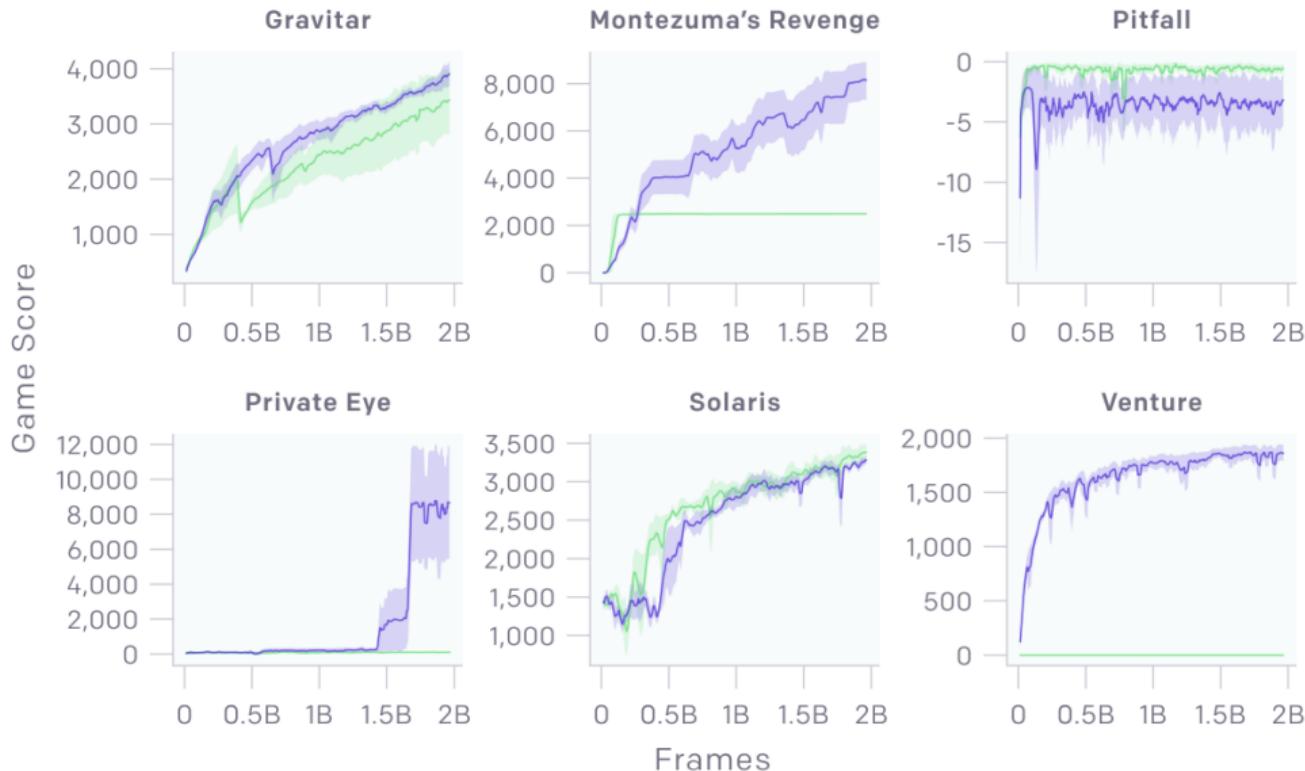
- Separate extrinsic r_t^E and intrinsic reward r_t^I
- PPO with two heads to estimate V^I and V^E
- Greedy policy w.r.t. $V^I + cV^E$

“Tricks”

- Rewards should be in the same range
- Use different discount factors for intrinsic and extrinsic rewards

Prediction-based Exploration

Burda et al [2019]



Prediction-based Exploration

[Burda et al., 2019]

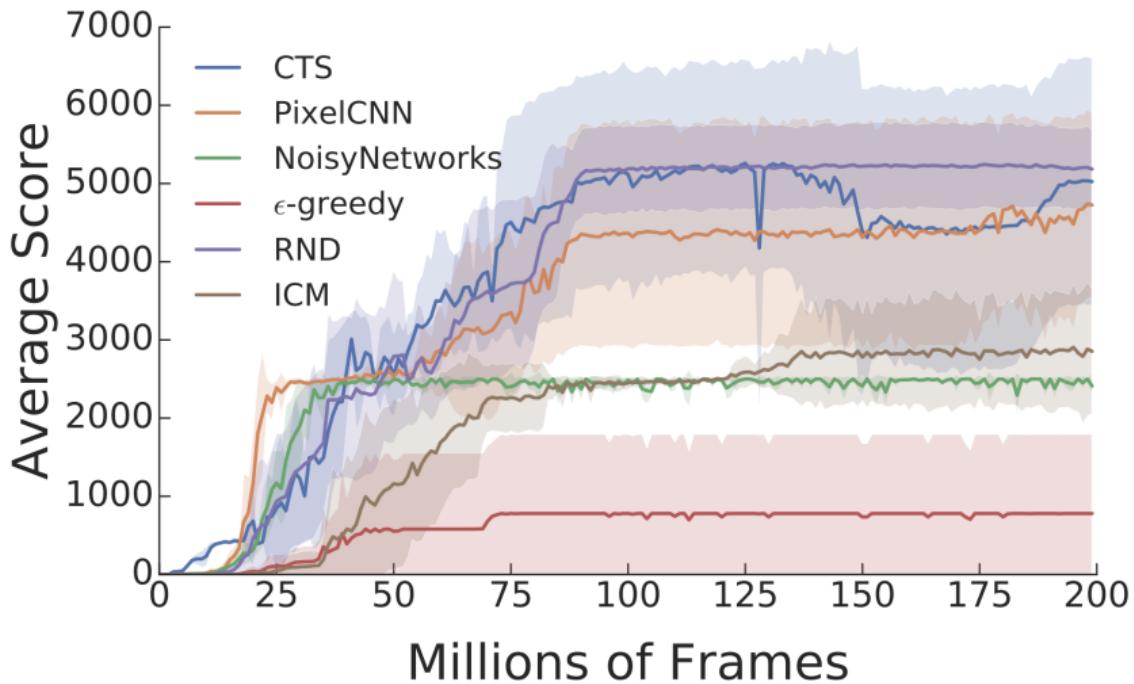
Montezuma!

finds 22 out of the 24 rooms on the first level

Comparison

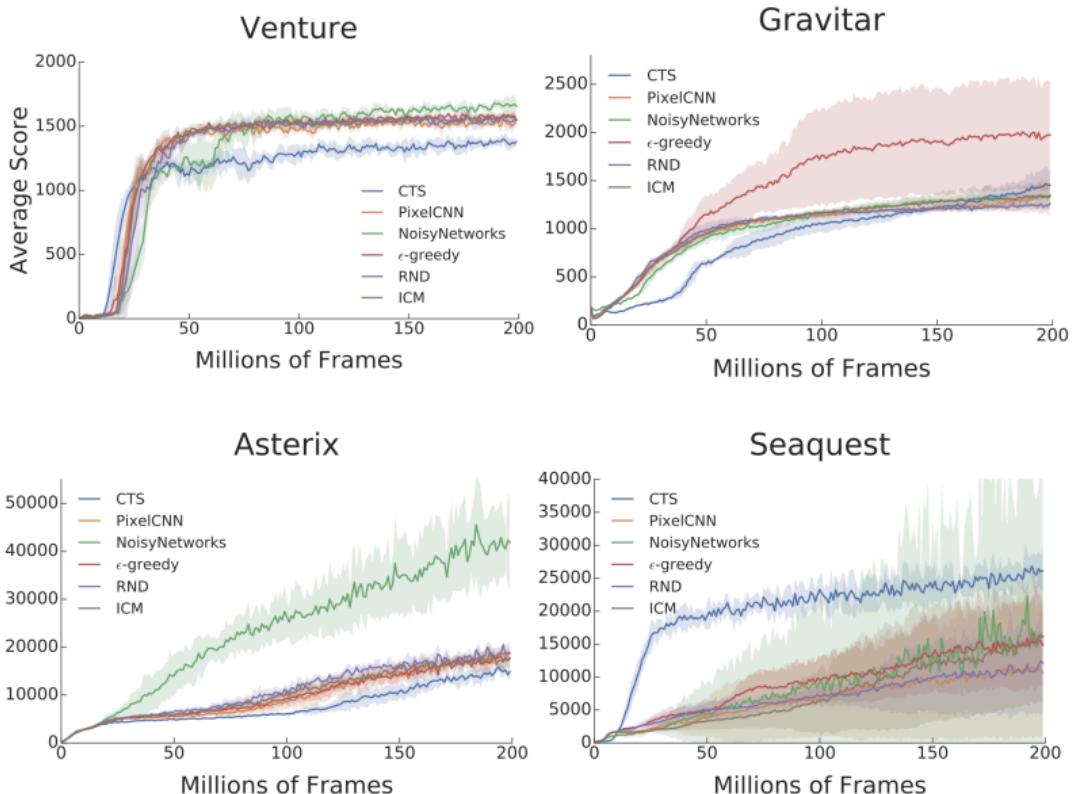
[Taïga et al., 2019]

Montezuma's Revenge



Comparison: not all problems require same amount of exploration

[Taïga et al., 2019]



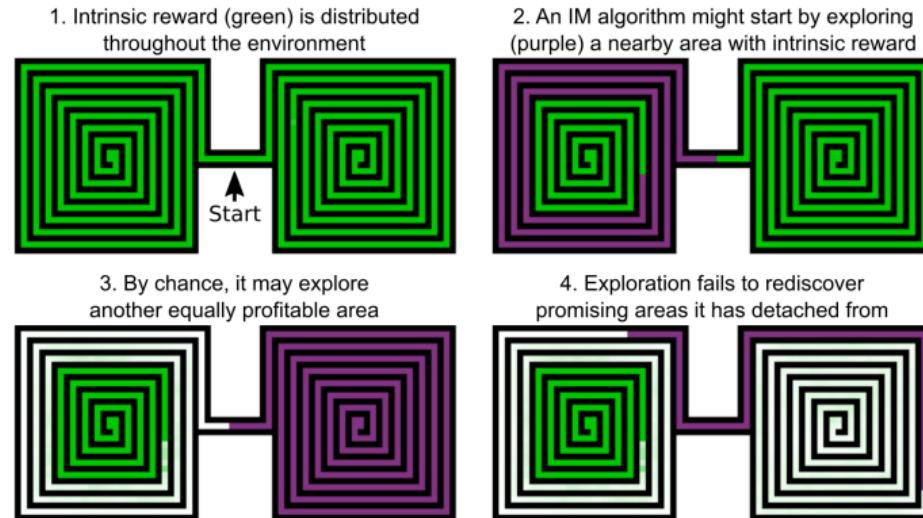
Go-Explore

[Ecoffet et al., 2019]

Issue of *intrinsic motivated* algorithms: *detachment problem*

- Forget about promising areas they have visited
- They do not return to them for further exploration

Green areas indicate intrinsic reward, white indicates areas where no intrinsic reward remains, and **purple** areas indicate where the algorithm is currently exploring.



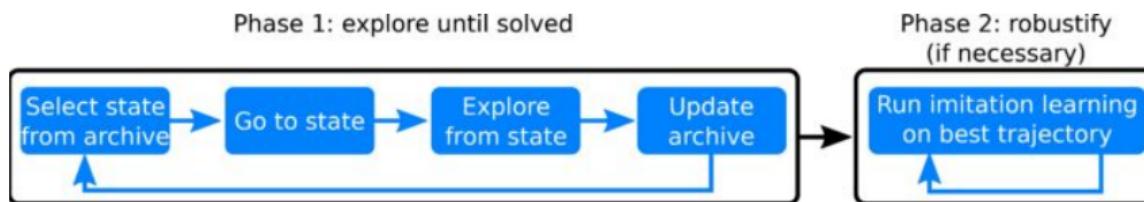
*mainly due to model-free nature

Go-Explore: Structure

[Ecoffet et al., 2019]

1 Exploration

- Select a promising state
- Go to a state
- Explore locally (e.g., randomly)
- Store observations
- Repeat



* similar in spirit to [Lim and Auer, 2012]

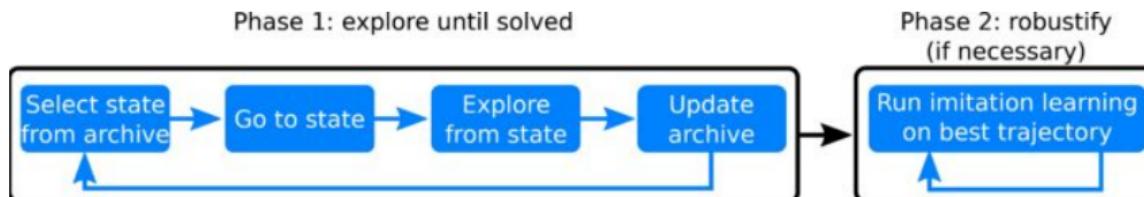
Go-Explore: Structure

[Ecoffet et al., 2019]

1 Exploration

- Select a promising state
- Go to a state
- Explore locally (e.g., randomly)
- Store observations
- Repeat

- Builds an archive of observed states in latent space
- Archive is sorted by relevance (e.g., IM, novelty)
- Knows a policy to reach an observed state (e.g., by replaying)



* similar in spirit to [Lim and Auer, 2012]

Go-Explore: Structure

[Ecoffet et al., 2019]

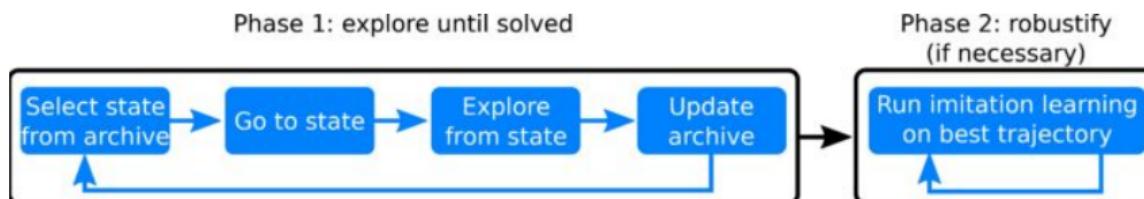
1 Exploration

- Select a promising state
- Go to a state
- Explore locally (e.g., randomly)
- Store observations
- Repeat

- Builds an archive of observed states in latent space
- Archive is sorted by relevance (e.g., IM, novelty)
- Knows a policy to reach an observed state (e.g., by replaying)

2 Robustification

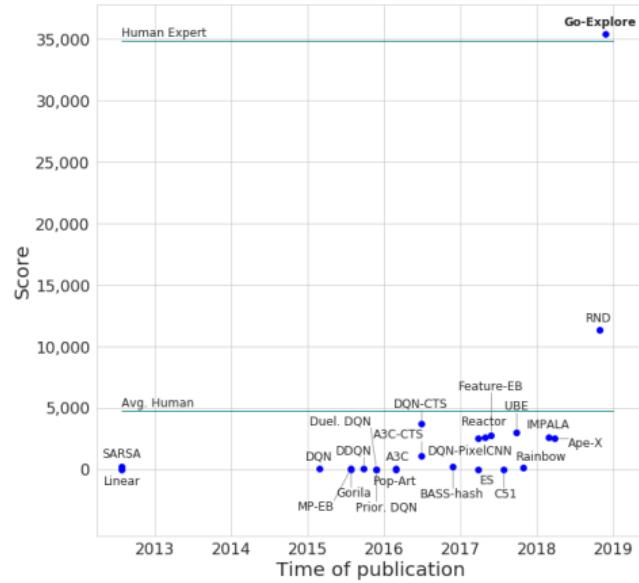
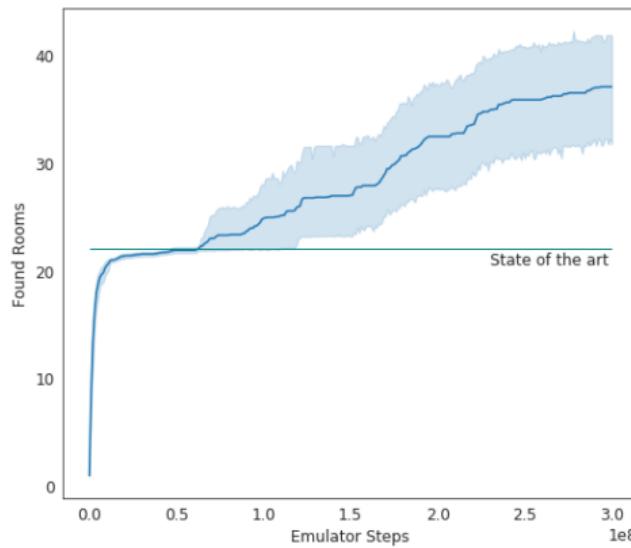
- Against noise
- Imitation learning on best trajectories



* similar in spirit to [Lim and Auer, 2012]

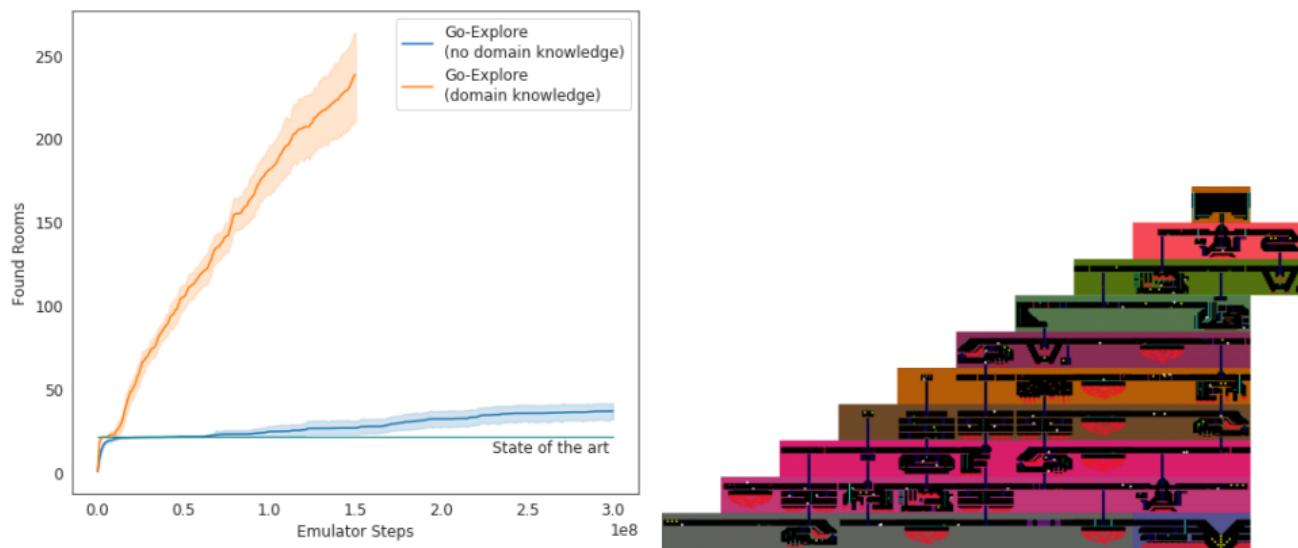
Go-Explore on Montezuma

[Ecoffet et al., 2019]



Go-Explore on Montezuma

[Ecoffet et al., 2019]



Using domain knowledge for the state representation, Phase 1 of Go-Explore finds a 238 rooms, solves over 9 levels on average

Image Credit: Wikimedia Foundation

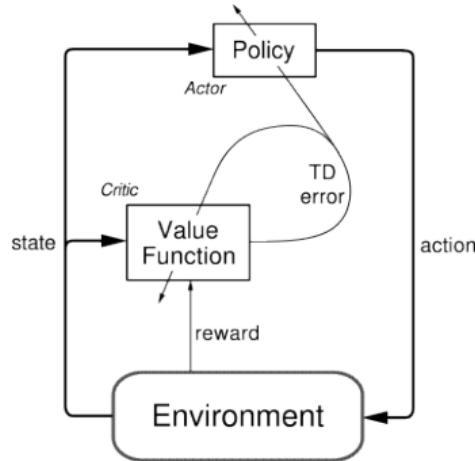
Go-Explore on Pitfall

[Ecoffet et al., 2019]

Pitfall!

Optimistic Actor-Critic

[Ciosek et al., 2019]



- Actor: decides which action to take
⇒ π
- Critic: estimates the goodness of an action in a state
⇒ Q

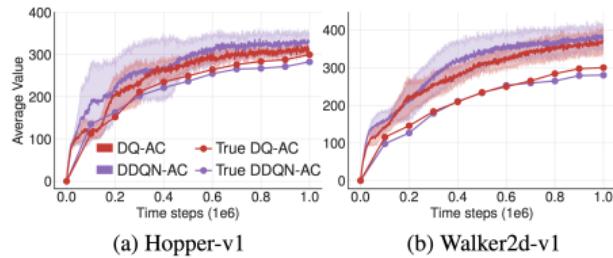
Policy performance: $J^\pi = \mathbb{E}_\pi \left[\sum_t \gamma^t r_t \right]$

Policy gradient:

$$\nabla J^\pi = \mathbb{E}_{s \sim d^\pi} \left[\sum_a \nabla \pi(s, a) Q^\pi(s, a) \right]$$

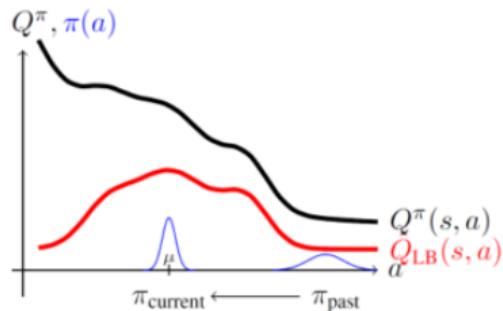
Optimistic Actor-Critic

[Ciosek et al., 2019]



* image from [Fujimoto et al., 2018]

👎 Issues with LB and greedy

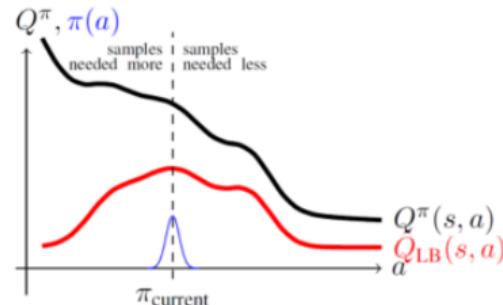


(a) Pessimistic underexploration

To limit overestimation (i.e., positive bias)

- use of two *identical* Q-functions
- train them independent

$$Q_{LB} = \min\{Q_1, Q_2\}$$



(b) Directional uninformedness

Optimistic Actor-Critic

[Ciosek et al., 2019]

- Build an upper-bound to the true Q -value (*optimistic Q-value*)

$$\hat{Q}(s, a) = \mu_Q(s, a) + \beta \sigma_Q(s, a)$$

with

$$\mu_Q(s, a) = \frac{1}{2} \left(Q^1(s, a) + Q^2(s, a) \right), \quad \sigma_Q(s, a) = \sqrt{\sum_{i \in \{1, 2\}} \frac{1}{2} (Q^i(s, a) - \mu_Q(s, a))^2}$$

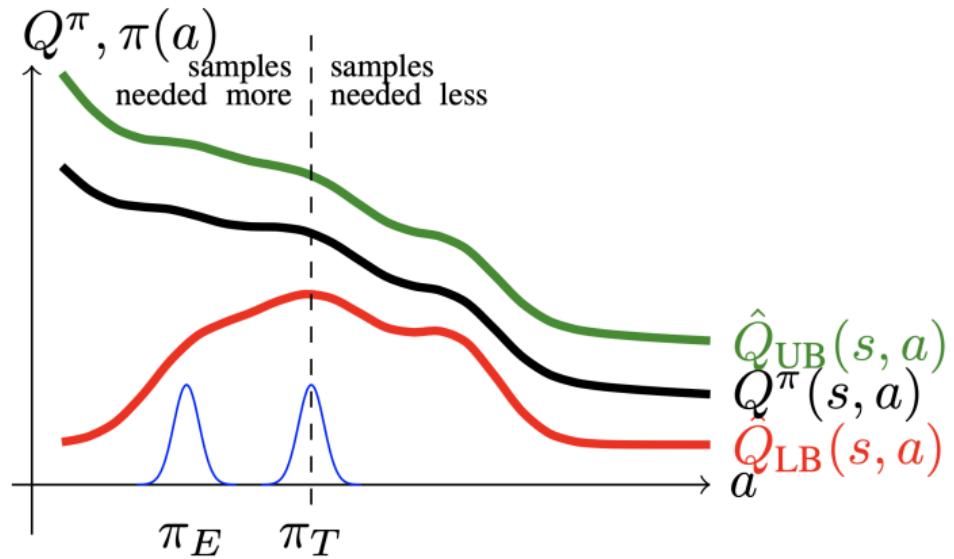
- Exploration policy* by soft-update (for stability): $\pi_E = \mathcal{N}(\mu_E, \Sigma_E)$

$$\begin{aligned} (\mu_E, \sigma_E) &= \arg \max_{\mu, \Sigma} \mathbb{E}_{a \sim \mathcal{N}(\mu, \Sigma)} [\hat{Q}(s, a)] \\ \text{s.t. } &KL(\mathcal{N}(\mu, \Sigma), \mathcal{N}(\mu_T, \Sigma_T)) \leq \delta \end{aligned}$$

* for computational efficiency, they fit a linear model on \hat{Q}

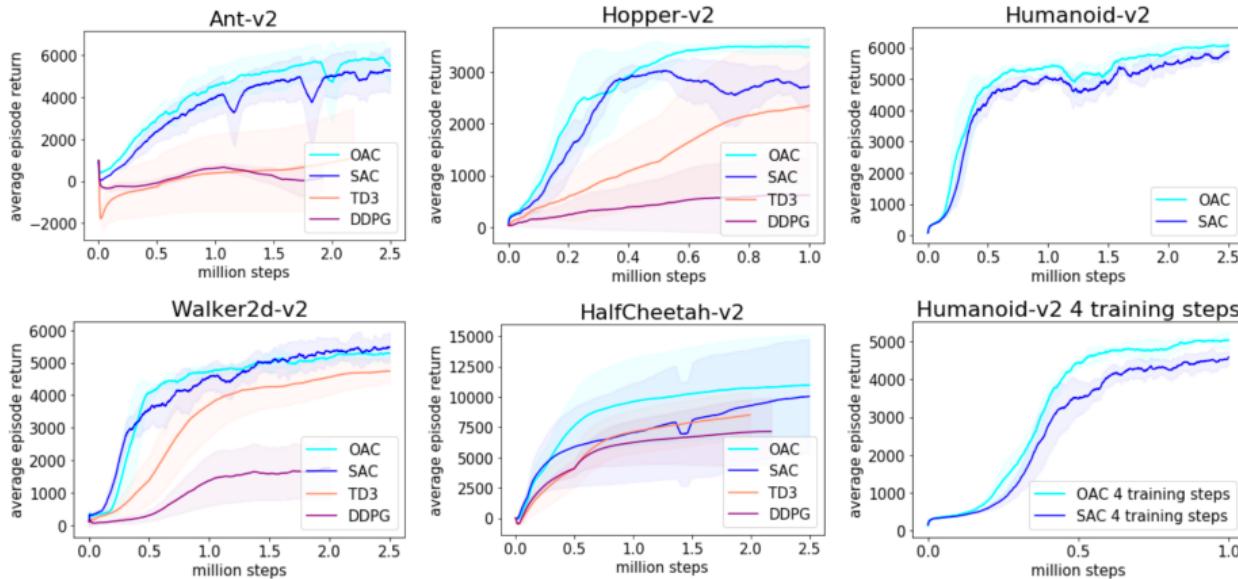
Optimistic Actor-Critic

[Ciosek et al., 2019]



Optimistic Actor-Critic: Experiments

[Ciosek et al., 2019]



Outline

1 Optimistic Exploration in Deep RL

2 Random Exploration in Deep RL

Randomized Exploration

General Scheme

- 1 Estimate the parameters θ for either policy or value function
- 2 Add randomness to the parameters $\tilde{\theta} = \theta + \text{noise}$
- 3 Run the corresponding (greedy) policy

Remark: changing weights induces a consistent, and potentially very complex, state-dependent change in policy over multiple time steps

- ⇒ long-term exploration
- ⇒ no dithering

Randomized Exploration

General Scheme

- 1 Estimate the parameters θ for either policy or value function
- 2 Add randomness to the parameters $\tilde{\theta} = \theta + \text{noise}$
- 3 Run the corresponding (greedy) policy

Remark: changing weights induces a consistent, and potentially very complex, state-dependent change in policy over multiple time steps

- ⇒ long-term exploration
- ⇒ no dithering

⚠ The randomness needs to represent “uncertainty”

Least-Squares Value Iteration

Estimate Q^* from samples

Input: Dataset $\mathcal{D}_k = (s_{hi}, a_{hi}, r_{hi})_{h=1, i=1}^{H, k}$

Set $\hat{Q}_{H+1}(s, a) = 0$

for $h = H, \dots, 1$ **do** // backward induction

 Compute

$$y_{hi} = r_{hi} + \max_{a \in \mathcal{A}} \hat{Q}_{h+1,k}(s_{h+1,i}, a) = r_{hi} + \hat{V}_{h+1,k}(s_{h+1,i}), \quad i = 1, \dots, k$$

 Build regression dataset $\mathcal{D}_h^{\text{reg}} = \{\phi_h(s_{hi}, a_{hi}), y_{hi}\}_i$

 Compute

$$\hat{\theta}_{hk} = \arg \min_{\theta} \left\{ \mathcal{L}(\theta, \mathcal{D}^{\text{reg}}) := \frac{1}{k} \sum_{i=1}^k (y_{hi} - Q_{hk}(s_{hi}, a_{hi} | \theta))^2 \right\}$$

end

return $\{\hat{\theta}_{hk}\}_{h=1}^H$

Optimize \mathcal{L} by *gradient descent*

Randomization on LSVI

How to force exploration

- Perturb observed rewards
- Perturb parameters (e.g., based on posterior uncertainty)

Randomized Value Function (RVF) [Osband et al., 2019, 2018, Azizzadenesheli et al., 2018, Lipton et al., 2018, Touati et al., 2019, Osband et al., 2019]

RVF: Reward Perturbation

[Osband et al., 2018, 2019]

Input: Dataset $\mathcal{D}_k = (s_{hi}, a_{hi}, r_{hi})_{h=1, i=1}^{H, k}$

Set $\widehat{Q}_{H+1}(s, a) = 0$

for $h = H, \dots, 1$ **do** // backward induction

Perturb rewards

$$\tilde{r}_{hi} = r_{hi} + \omega_{hi}, \quad \omega_{hi} \sim \mathcal{N}(0, \sigma^2)$$

Compute

$$\tilde{y}_{hi} = \tilde{r}_{hi} + \max_{a \in \mathcal{A}} \widehat{Q}_{h+1, k}(s_{h+1, i}, a) = \tilde{r}_{hi} + \widehat{V}_{h+1, k}(s_{h+1, i}), \quad i = 1, \dots, k$$

Build regression dataset $\widetilde{\mathcal{D}}_h^{\text{reg}} = \{(s_{hi}, a_{hi}), \tilde{y}_{hi}\}_i$

Sample θ^p from prior

Compute

$$\widehat{\theta}_{hk} = \arg \min_{\theta} \left\{ \mathcal{L}^B(\theta, \theta^p, \widetilde{\mathcal{D}}^{\text{reg}}) := \frac{1}{k} \sum_{i=1}^k (\tilde{y}_{hi} - Q_{hk}(s_{hi}, a_{hi} | \theta))^2 + \mathcal{R}(\theta, \theta^p) \right\}$$

end

return $\{\widehat{\theta}_{hk}\}_{h=1}^H$

RVF: Perturb Parameters

[Osband et al., 2016b]

Input: Dataset $\mathcal{D}_k = (s_{hi}, a_{hi}, r_{hi})_{h=1, i=1}^{H, k}$

Set $\bar{Q}_{H+1}(s, a) = 0$

for $h = H, \dots, 1$ **do** // backward induction
 Compute

$$\bar{y}_{hi} = r_{hi} + \max_{a \in \mathcal{A}} \bar{Q}_{h+1,k}(s_{h+1,i}, a) = r_{hi} + \bar{V}_{h+1,k}(s_{h+1,i}), \quad i = 1, \dots, k$$

Bootstrapping randomized estimates

Build regression dataset $\bar{\mathcal{D}}_h^{\text{reg}} = \{\phi_h(s_{hi}, a_{hi}), \bar{y}_{hi}\}_i$

Sample θ^p from prior

Compute

$$\hat{\theta}_{hk} = \arg \min_{\theta} \left\{ \mathcal{L}^B(\theta, \theta^p, \bar{\mathcal{D}}^{\text{reg}}) := \frac{1}{k} \sum_{i=1}^k (\bar{y}_{hi} - Q_{hk}(s_{hi}, a_{hi} | \theta))^2 + \mathcal{R}(\theta, \theta^p) \right\}$$

Sample $\xi_{hk} \sim \mathcal{N}(0, \Sigma_{hk}^{-1})$

Set $\bar{\theta}_{hk} = \hat{\theta}_{hk} + \xi_{hk}$

end

return $\{\bar{\theta}_{hk}\}_{h=1}^H$

$$\hat{\theta} = \mathbb{E}[\theta | \mathcal{D}^{\text{reg}}, \text{prior}], \quad \Sigma^{-1} = \text{Cov}[\theta | \mathcal{D}^{\text{reg}}, \text{prior}]$$

RLSVI as Regression on Perturbed Data

[Osband et al., 2018, 2019]

Bayesian Linear Regression (posterior structure)

- True parameter is $\theta^* \in \mathbb{R}^d \Rightarrow$ we want to estimate it
- Assume *Gaussian prior* $\mathcal{N}(\bar{\theta}, \lambda I)$
- Dataset $\mathcal{D} = (x_i, y_i)_{i=1}^N$, where

$$y_i = x_i^\top \theta^* + \epsilon_i \quad , \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

- Solve $\min_{\theta} \mathcal{L}^B(\theta, \bar{\theta}, \mathcal{D})$
- Conditional *posterior* μ_p

$$\theta^* | \mathcal{D} \sim \mu_p = \mathcal{N}\left(\underbrace{\Sigma^{-1} \left(\frac{1}{\sigma^2} X^\top y + \frac{1}{\lambda} \bar{\theta} \right)}_{\widehat{\theta}}, \Sigma^{-1} \right)$$

$$\Sigma = \frac{1}{\sigma^2} X^\top X + \frac{1}{\lambda} I$$

RLSVI as Regression on Perturbed Data

[Osband et al., 2018, 2019]

Target perturbation

■ Compute

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{\sigma^2} \sum_{i=1}^N (y_i + \omega_i - x_i^\top \theta)^2 + \frac{1}{\lambda} \|\tilde{\theta} - \theta\|_2^2$$

$\Rightarrow \hat{\theta} \sim \mu_p$

👉 Computational generation of posterior samples for linear Bayesian regression

i.e., we can sample μ_p by fitting a least-squares estimate

RLSVI as Regression on Perturbed Data

[Osband et al., 2018, 2019]

Target perturbation

■ Compute

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{\sigma^2} \sum_{i=1}^N (y_i + \omega_i - x_i^\top \theta)^2 + \frac{1}{\lambda} \|\tilde{\theta} - \theta\|_2^2$$

$\Rightarrow \hat{\theta} \sim \mu_p$

✍ Computational generation of posterior samples for linear Bayesian regression

i.e., we can sample μ_p by fitting a least-squares estimate

For linear models,

poster sampling = regularized least-squares on perturbed data

✍ For tabular MDPs, $x_i = e_{s,a}$ and $\theta = Q$

backward induction on randomized rewards = RLSVI (see Part 2)

RVF: issues

Reward Perturbation

- Minimize least-squares problem for any reward structure
e.g., by gradient descent
- Not so easy to define the magnitude of the reward perturbation

Posterior Sampling

- Posterior variance
 - easy for linear model
 - hard (almost impossible) for generic models
- A lot of approximate schemas for computing the posterior

Randomized Prior

[Osband et al., 2018]

Algorithm 1 Randomized prior functions for ensemble posterior.

Require: Data $\mathcal{D} \subseteq \{(x, y) | x \in \mathcal{X}, y \in \mathcal{Y}\}$, loss function \mathcal{L} , neural model $f_\theta: \mathcal{X} \rightarrow \mathcal{Y}$, Ensemble size $K \in \mathbb{N}$, noise procedure `data_noise`, distribution over priors $\mathcal{P} \subseteq \{\mathbb{P}(p) | p: \mathcal{X} \rightarrow \mathcal{Y}\}$.

- 1: **for** $k = 1, \dots, K$ **do**
- 2: | initialize $\theta_k \sim$ Glorot initialization [23].
- 3: | form $\mathcal{D}_k = \text{data_noise}(\mathcal{D})$ (e.g. Gaussian noise or bootstrap sampling [50]).
- 4: | sample prior function $p_k \sim \mathcal{P}$.
- 5: | optimize $\nabla_{\theta|\theta=\theta_k} \mathcal{L}(f_\theta + p_k; \mathcal{D}_k)$ via ADAM [28].
- 6: **return** ensemble $\{f_{\theta_k} + p_k\}_{k=1}^K$.

$$\mathcal{L}_\gamma(\theta; \theta^-, p, \mathcal{D}) := \sum_{t \in \mathcal{D}} \left(r_t + \gamma \max_{a' \in \mathcal{A}} \overbrace{(f_{\theta^-} + p)(s'_t, a')}^{\text{target } Q} - \overbrace{(f_\theta + p)(s_t, a_t)}^{\text{online } Q} \right)^2$$

Posterior Distribution for Deep Neural Networks

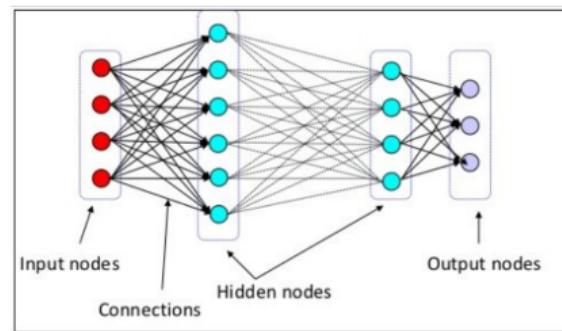
Bayesian DQN [Azizzadenesheli et al., 2018]

⚠ Same tools as in linear bandit

- 1 Bayesian linear regression with given feature $\phi(s) \in \mathbb{R}^d$ and given target vector for each action y_a

$$\mu_a = (\Phi_a^\top \Phi_a)^{-1} \Phi_a^\top y_a \quad \Sigma_a = \Phi_a^\top \Phi_a$$

- 2 Draw a weight vector at random $w_a \sim \mathcal{N}(\mu_a, \Sigma_a^{-1})$
- 3 Run the corresponding (greedy) policy
 $a_t = \arg \max_a Q(s_t, a) := \arg \max_a w_a^\top \phi(s_t)$
- 4 Train ϕ with standard NN to estimate Q



Posterior Distribution for Deep Neural Networks

Bayesian DQN [Azizzadenesheli et al., 2018]

Game	BDQN	DDQN	DDQN ⁺	Bootstrap	NoisyNet	CTS	Pixel	Reactor	Human	SC	SC ⁺	Step
Amidar	5.52k	0.99k	0.7k	1.27k	1.5k	1.03k	0.62k	1.18k	1.7k	22.9M	4.4M	100M
Alien	3k	2.9k	2.9k	2.44k	2.9k	1.9k	1.7k	3.5k	6.9k	-	36.27M	100M
Assault	8.84k	2.23k	5.02k	8.05k	3.1k	2.88k	1.25k	3.5k	1.5k	1.6M	24.3M	100M
Asteroids	14.1k	0.56k	0.93k	1.03k	2.1k	3.95k	0.9k	1.75k	13.1k	58.2M	9.7M	100M
Asterix	58.4k	11k	15.15k	19.7k	11.0	9.55k	1.4k	6.2k	8.5k	3.6M	5.7M	100M
BeamRider	8.7k	4.2k	7.6k	23.4k	14.7k	7.0k	3k	3.8k	5.8k	4.0M	8.1M	70M
BattleZone	65.2k	23.2k	24.7k	36.7k	11.9k	7.97k	10k	45k	38k	25.1M	14.9M	50M
Atlantis	3.24M	39.7k	64.76k	99.4k	7.9k	1.8M	40k	9.5M	29k	3.3M	5.1M	40M
DemonAttack	11.1k	3.8k	9.7k	82.6k	26.7k	39.3k	1.3k	7k	3.4k	2.0M	19.9M	40M
Centipede	7.3k	6.4k	4.1k	4.55k	3.35k	5.4k	1.8k	3.5k	12k	-	4.2M	40M
BankHeist	0.72k	0.34k	0.72k	1.21k	0.64k	1.3k	0.42k	1.1k	0.72k	2.1M	10.1M	40M
CrazyClimber	124k	84k	102k	138k	121k	112.9k	75k	119k	35.4k	0.12M	2.1M	40M
ChopperCmd	72.5k	0.5k	4.6k	4.1k	5.3k	5.1k	2.5k	4.8k	9.9k	4.4M	2.2M	40M
Enduro	1.12k	0.38k	0.32k	1.59k	0.91k	0.69k	0.19k	2.49k	0.31k	0.82M	0.8M	30M
Pong	21	18.82	21	20.9	21	20.8	17	20	9.3	1.2M	2.4M	5M

Posterior Distribution for Deep Neural Networks

BBQ-Networks [Lipton et al., 2018]

- Uses variational inference to quantify uncertainty
- Uses independent factorized Gaussians as an approximate posterior

MNF-DQN [Touati et al., 2019]

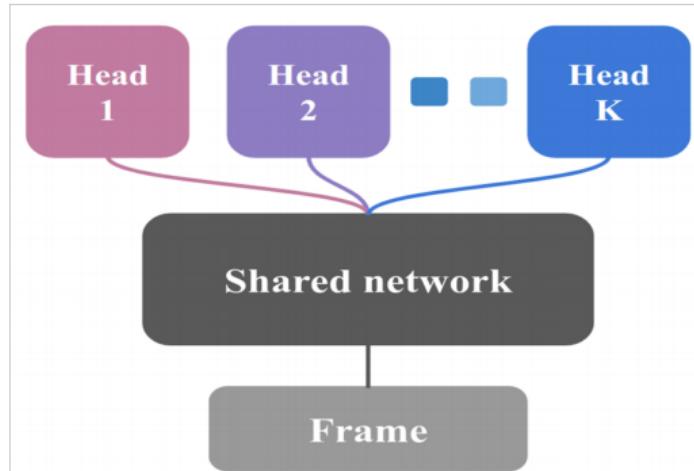
- Leverages recent advances in variational Bayesian NN
- Computationally and statistically efficient
- Uses normalizing multiplicative flows (MNF) in order to account for the uncertainty of estimates for efficient exploration

Bootstrap DQN

[Osband et al., 2016a]

- Define multiple value functions Q_k
- Update functions with different datasets
- Share part of the architecture

another way of approximating a sample from posterior



Bootstrap DQN

[Osband et al., 2016a]

Algorithm 1 Bootstrapped DQN

- 1: **Input:** Value function networks Q with K outputs $\{Q_k\}_{k=1}^K$. Masking distribution M .
- 2: Let B be a replay buffer storing experience for training.
- 3: **for** each episode **do**
- 4: Obtain initial state from environment s_0
- 5: Pick a value function to act using $k \sim \text{Uniform}\{1, \dots, K\}$
- 6: **for** step $t = 1, \dots$ until end of episode **do**
- 7: Pick an action according to $a_t \in \arg \max_a Q_k(s_t, a)$
- 8: Receive state s_{t+1} and reward r_t from environment, having taking action a_t
- 9: Sample bootstrap mask $m_t \sim M$
- 10: Add $(s_t, a_t, r_{t+1}, s_{t+1}, m_t)$ to replay buffer B
- 11: **end for**
- 12: **end for**

- M_t determines the type of bootstrapping strategy

$$g_t^k = m_t^k (y_t^Q - Q_k(s_t, a_t; \theta)) \nabla_\theta Q_k(s_t, a_t, ; \theta)$$

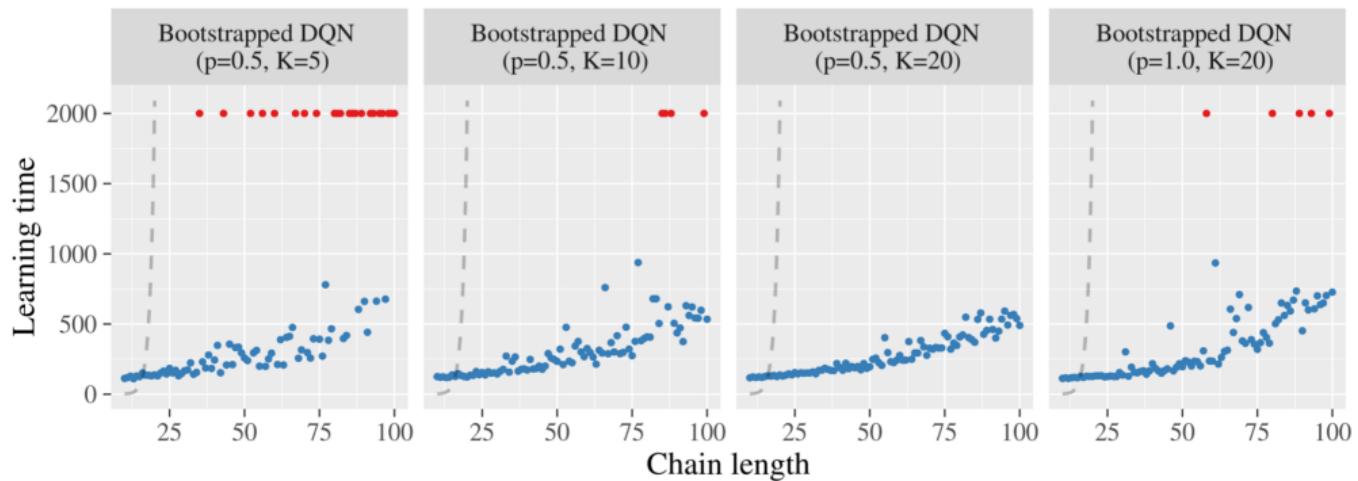
with target $y_t = r_t + \max_a Q(s_{t+1}, a; \theta^-)$

Bootstrap DQN

[Osband et al., 2016a]



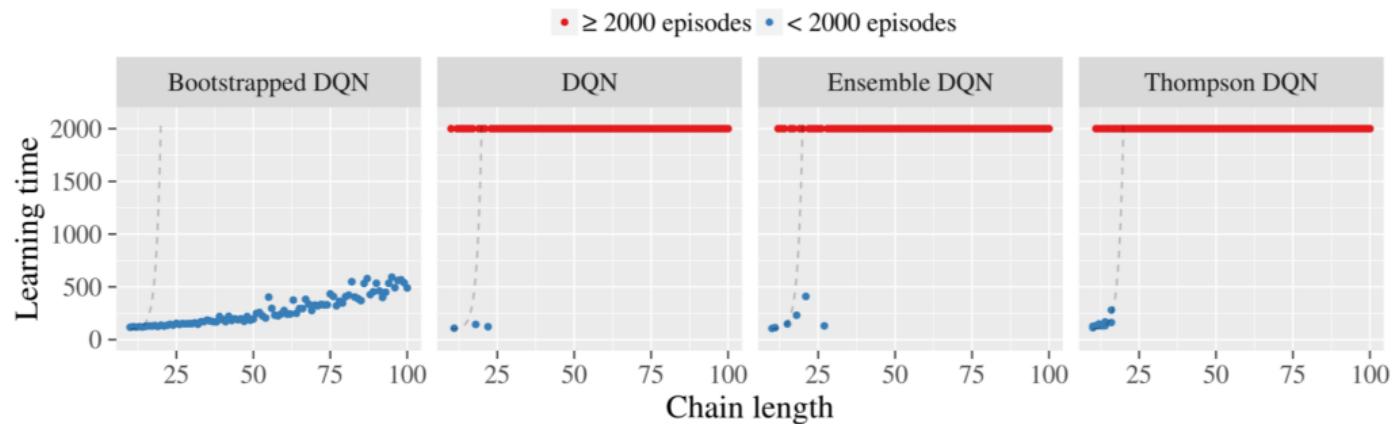
• ≥ 2000 episodes • < 2000 episodes



Masking rule for samples in episode k : $m_k \sim Ber(p)$

Bootstrap DQN

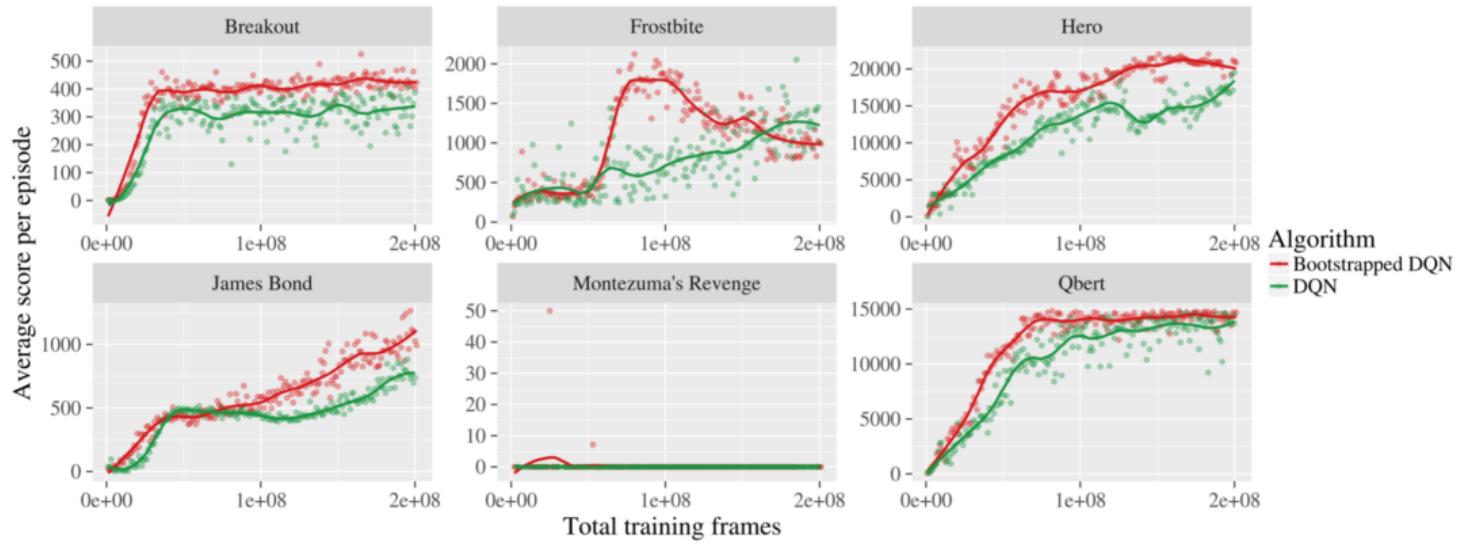
[Osband et al., 2016a]



- Ensemble DQN: ensemble policy?
- Thompson DQN: resample at each step

Bootstrap DQN

[Osband et al., 2016a]



Noisy Networks

[Fortunato et al., 2018]

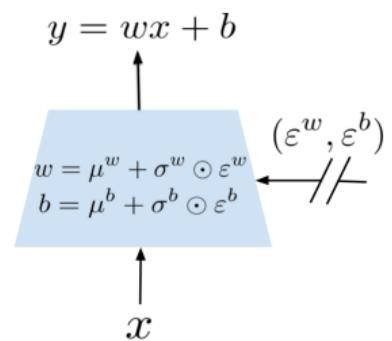
- Normal NN layer $y = wx + b$
- Double the parameters with mean and variance $w \rightarrow \mu^w, \sigma^w$ and $b \rightarrow \mu^b, \sigma^b$
- Whenever a layer is evaluated draw $\varepsilon^w, \varepsilon^b \sim \mathcal{D}$
- Evaluate the “random” layer as

$$y = (\mu^w + \sigma^w \odot \varepsilon^w) + \mu^b + \sigma^b \odot \varepsilon^b$$
- Let $\zeta = (\mu^w, \sigma^w, \mu^b, \sigma^b)$, define the expected loss

$$\bar{L}(\zeta) = \mathbb{E}_\varepsilon [L(\zeta, \varepsilon)]$$

- Gradient estimation

$$\nabla_\zeta \bar{L}(\zeta) = \mathbb{E}_\varepsilon [\nabla_\zeta L(\zeta, \varepsilon)] \approx \frac{1}{n} \sum_{i=1}^n \nabla_\zeta L(\zeta, \varepsilon_i)$$



Noisy Networks

[Fortunato et al., 2018]

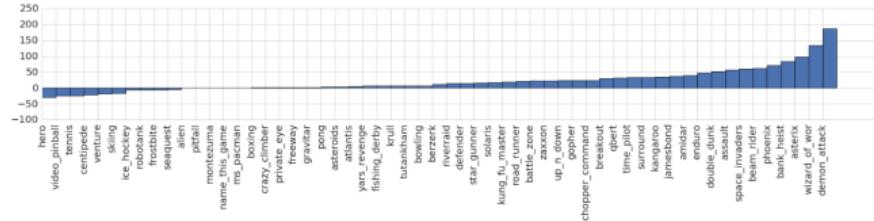
Noise models

- Independent noise $\varepsilon_{i,j}$ for each weight i at layer j
- Factorized noise $\varepsilon_{i,j} = f(\varepsilon_i)f(\varepsilon_j)$ (e.g., $f(x) = \text{sgn}(x)\sqrt{x}$)
- Independent noise for target and online networks

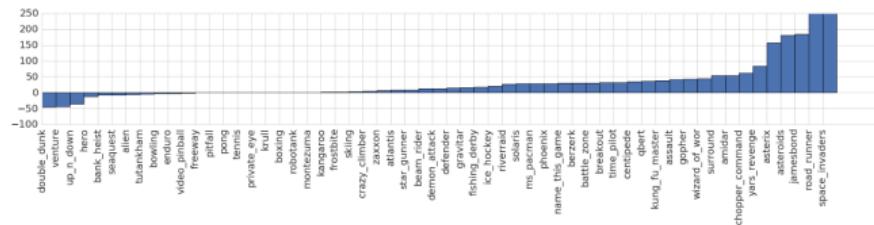
$$y_t = r_t + \max_{a'} Q(s'_t, a'; \varepsilon', \zeta^-); \quad L_t(\zeta, \varepsilon) = (y_t - Q(s_t, a_t; \varepsilon, \zeta))^2$$

Noisy Networks

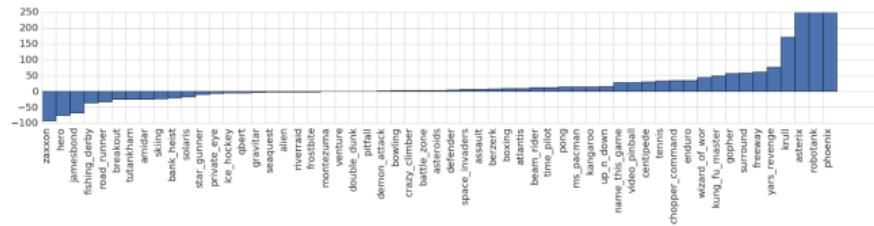
[Fortunato et al., 2018]



(a) Improvement in percentage of NoisyNet-DQN over DQN (Mnih et al., 2015)



(b) Improvement in percentage of NoisyNet-Dueling over Dueling (Wang et al., 2016)



(c) Improvement in percentage of NoisyNet-A3C over A3C (Mnih et al., 2016)

Comparison

[Touati et al., 2019]

Simple Chain domain

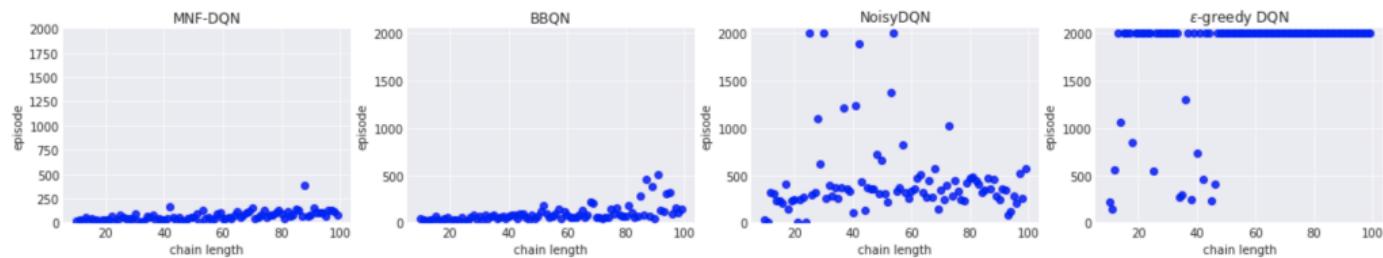
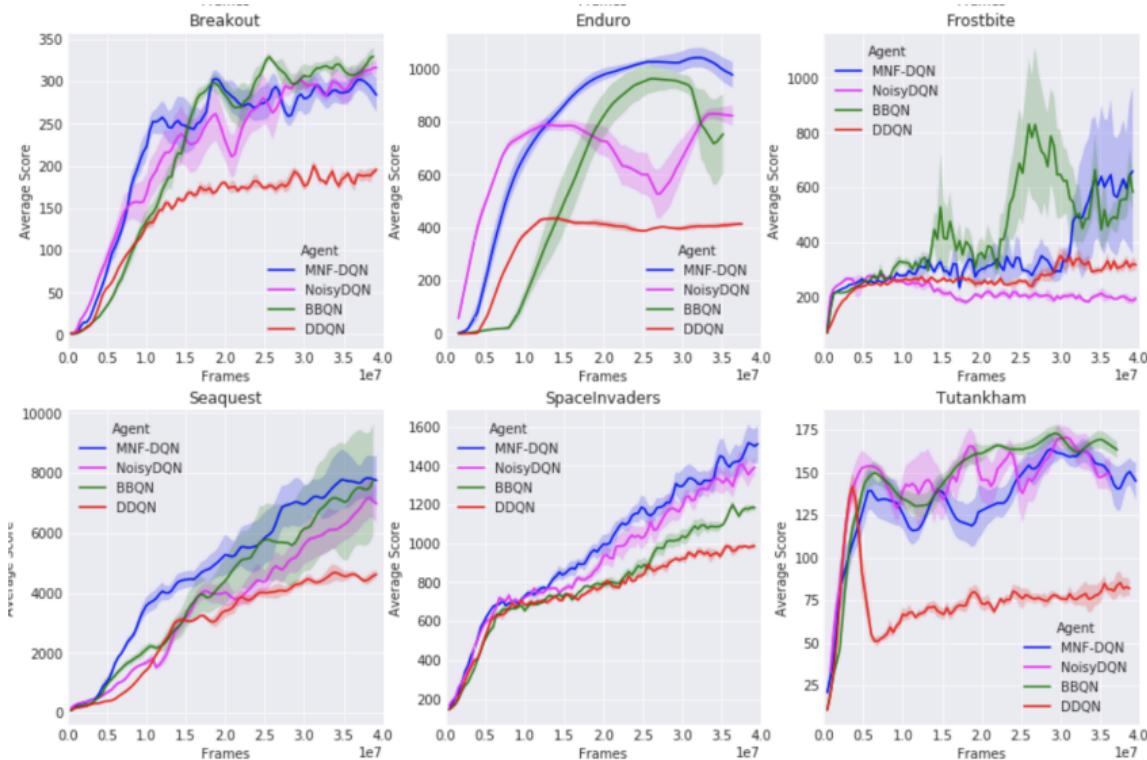


Figure 1: Median number of episodes (max 2000) required to solve the n-chain problem for (figure from left to right) MNF-DQN, BBQN, NoisyDQN and ϵ -greedy DQN. The median is obtained over 10 runs with different seeds. We see that MNF-DQN consistently performs best across different chain lengths.

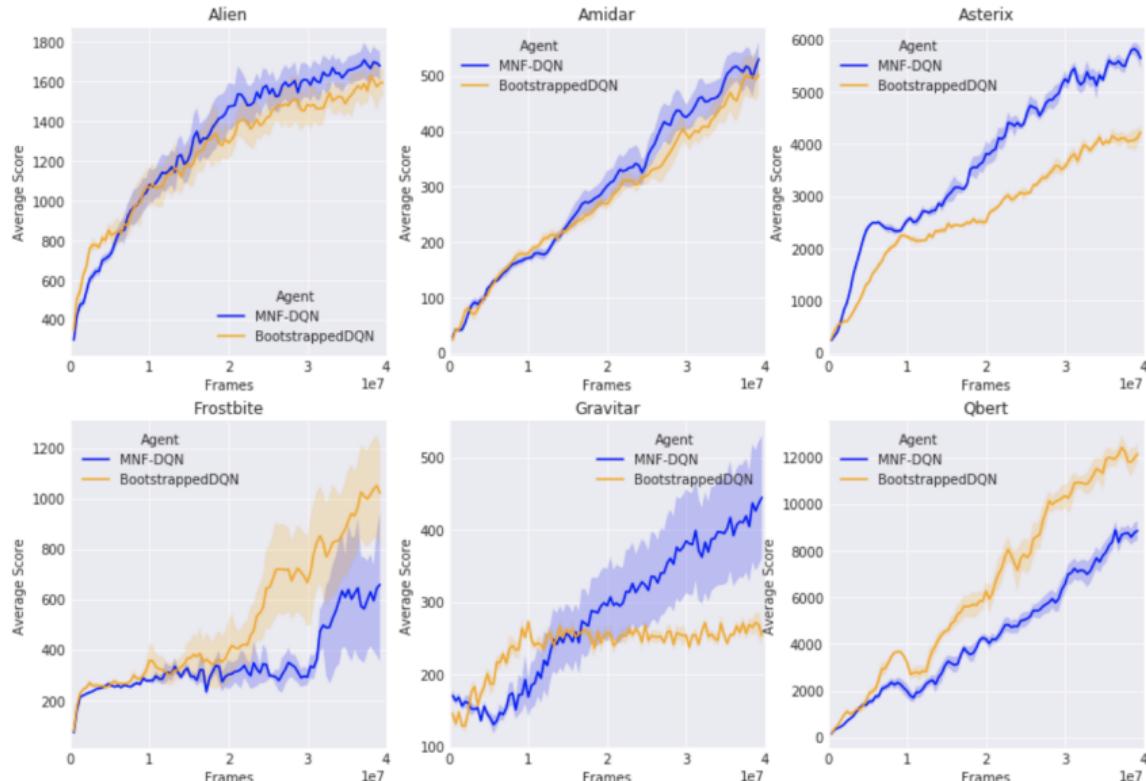
Comparison: Atari

[Touati et al., 2019]



Comparison: Atari

[Touati et al., 2019]



Remarks

- Exploration needs to account for uncertainty in the predictions
- Should account for long-term effect

Exploration at the level of (value/policy/model) parameters

Randomized explorations performs often better than optimism

- Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration through bayesian deep q-networks. In *ITA*, pages 1–9. IEEE, 2018.
- Marc G. Bellemare, Joel Veness, and Erik Talvitie. Skip context tree switching. In *ICML*, volume 32 of *JMLR Workshop and Conference Proceedings*, pages 1458–1466. JMLR.org, 2014.
- Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. In *NIPS*, pages 1471–1479, 2016.
- Yuri Burda, Harrison Edwards, Amos J. Storkey, and Oleg Klimov. Exploration by random network distillation. In *ICLR (Poster)*. OpenReview.net, 2019.
- Kamil Ciosek, Quan Vuong, Robert Loftin, and Katja Hofmann. Better exploration with optimistic actor critic. In *NeurIPS*, pages 1785–1796, 2019.
- Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. Go-explore: a new approach for hard-exploration problems. *CoRR*, abs/1901.10995, 2019.
- Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Matteo Hessel, Ian Osband, Alex Graves, Volodymyr Mnih, Rémi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. In *ICLR (Poster)*. OpenReview.net, 2018.
- Lior Fox, Leshem Choshen, and Yonatan Loewenstein. DORA the explorer: Directed outreaching reinforcement action-selection. In *ICLR (Poster)*. OpenReview.net, 2018.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *ICML*, volume 80 of *Proceedings of Machine Learning Research*, pages 1582–1591. PMLR, 2018.
- Shiau Hong Lim and Peter Auer. Autonomous exploration for navigating in mdps. In *COLT*, volume 23 of *JMLR Proceedings*, pages 40.1–40.24. JMLR.org, 2012.
- Zachary C. Lipton, Xiujun Li, Jianfeng Gao, Lihong Li, Faisal Ahmed, and Li Deng. Bbq-networks: Efficient exploration in deep reinforcement learning for task-oriented dialogue systems. In *AAAI*, pages 5237–5244. AAAI Press, 2018.

- Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. Count-based exploration with the successor representation. *CoRR*, abs/1807.11622, 2018.
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped DQN. In *NIPS*, pages 4026–4034, 2016a.
- Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 2377–2386. JMLR.org, 2016b.
- Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep reinforcement learning. In *NeurIPS*, pages 8626–8638, 2018.
- Ian Osband, Benjamin Van Roy, Daniel J. Russo, and Zheng Wen. Deep exploration via randomized value functions. *Journal of Machine Learning Research*, 20(124):1–62, 2019. URL <http://jmlr.org/papers/v20/18-339.html>.
- Georg Ostrovski, Marc G. Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration with neural density models. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 2721–2730. PMLR, 2017.
- Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 2778–2787. PMLR, 2017.
- Yifan Sun, Yaqi Duan, Hao Gong, and Mengdi Wang. Learning low-dimensional state embeddings and metastable clusters from time series data. In *NeurIPS*, pages 4563–4572, 2019.
- Adrien Ali Taïga, William Fedus, Marlos C. Machado, Aaron C. Courville, and Marc G. Bellemare. Benchmarking bonus-based exploration methods on the arcade learning environment. *CoRR*, abs/1908.02388, 2019.
- Haoran Tang, Rein Houthooft, Davis Foote, Adam Stooke, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. #exploration: A study of count-based exploration for deep reinforcement learning. In *NIPS*, pages 2753–2762, 2017.

- Ahmed Touati, Harsh Satija, Joshua Romoff, Joelle Pineau, and Pascal Vincent. Randomized value functions via multiplicative normalizing flows. In *UAI*, page 156. AUAI Press, 2019.
- Aäron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1747–1756. JMLR.org, 2016.