

Four - DQN Extensions

Placing the blame, not all experience is made equal, the space between bias and variance.

Eligibility Traces

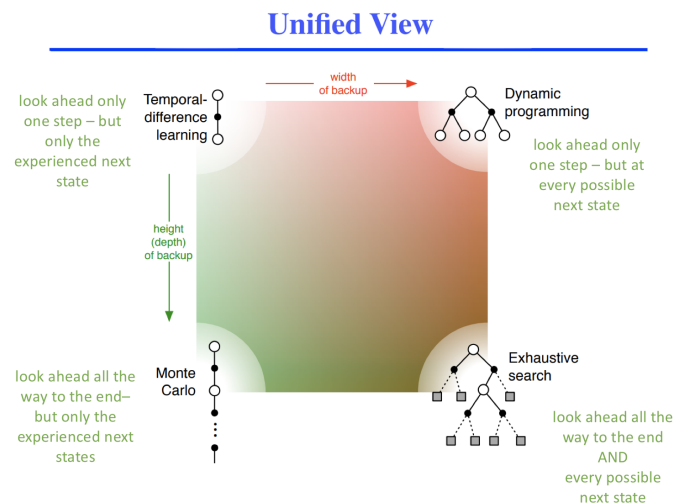


Figure 1: A unified view - [The Long-term of AI & Temporal-Difference Learning – Rich Sutton](#)

Family of methods between Temporal Difference & Monte Carlo

Eligibility traces allow us to **assign TD errors** to different states

- can be useful with delayed rewards or non-Markov environments
- requires more computation
- squeezes more out of data

Allow us to tradeoff between bias and variance

The space between TD and MC

In between TD and MC exist a family of approximation methods known as **n-step returns**

$$TD(\lambda)$$

The family of algorithms between TD and MC is known as $TD(\lambda)$ - weight each return by λ^{n-1} - normalize using $(1 - \lambda)$

$$TD(\lambda) = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} R_t^n$$

$\lambda = 0 \rightarrow TD(0)$

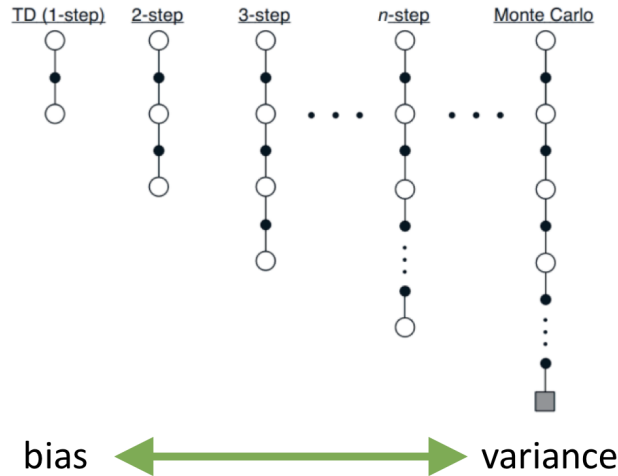


Figure 2: Sutton & Barto

$\lambda = 1 \rightarrow$ Monte Carlo

$TD(\lambda)$ and n-step returns are the same thing

Forward and backward view

We can look at eligibility traces from two perspectives

The **forward** view is helpful for understanding the theory

The **backward** view can be put into practice

We can decompose return into **complex backups** - looking forward to future returns - can use a combination of experience based and model based backups

$$R_t = \frac{1}{2}R_t^2 + \frac{1}{2}R_t^4$$

$$R_t = \frac{1}{2}TD + \frac{1}{2}MC$$

The backward view

The backward view approximates the forward view - forward view is not practical (requires knowledge of the future)

It requires an additional variable in our agents memory - **eligibility trace** $e_t(s)$

At each step we decay the trace according to

$$e_t(s) = \gamma \lambda e_{t-1}(s)$$

Unless we visited that state, in which case we accumulate more eligibility

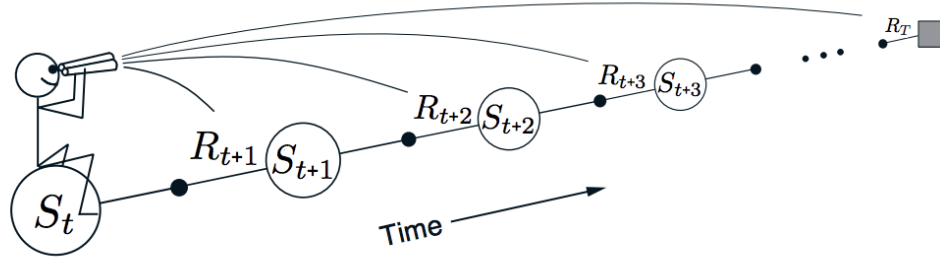


Figure 12.4: The forward view. We decide how to update each state by looking forward to future rewards and states.

Figure 3: Sutton & Barto

$$e_t(s) = \gamma \lambda e_{t-1}(s) + 1$$

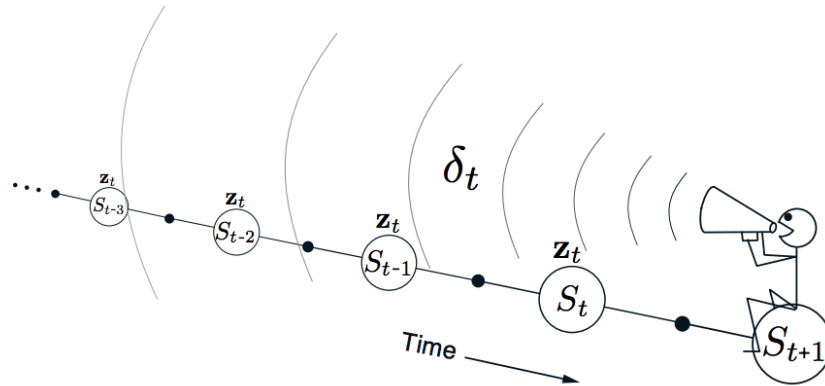


Figure 12.5: The backward or mechanistic view. Each update depends on the current TD error combined with the current eligibility traces of past events.

Figure 4: Sutton & Barto

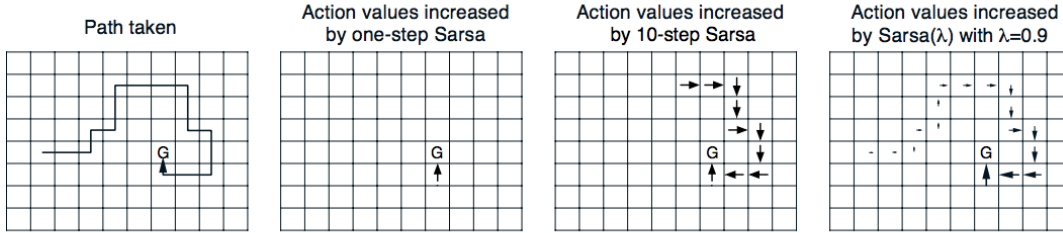
Traces in a grid world

- one step method would only update the last $Q(s, a)$
- n-step method would update all $Q(s, a)$ equally
- eligibility traces updates based on how recently each $Q(s, a)$ was experienced

Prioritized experience replay

Reference = Schaul et. al (2016) Prioritized Experience Replay [TODO link](#)

Example 12.1: Traces in Gridworld The use of eligibility traces can substantially increase the efficiency of control algorithms over one-step methods and even over n -step methods. The reason for this is illustrated by the gridworld example below.



The first panel shows the path taken by an agent in a single episode. The initial estimated values were zero, and all rewards were zero except for a positive reward at the goal location marked by G. The arrows in the other panels show, for various algorithms, which action-values would be increased, and by how much, upon reaching the goal. A one-step method would increment only the last action value, whereas an n -step method would equally increment the last n action's values, and an eligibility trace method would update all the action values up to the beginning of the episode to different degrees, fading with recency. The fading strategy is often the best tradeoff, strongly learning how to reach the goal from the right, yet not as strongly learning the roundabout path to the goal from the left that was taken in this episode. ■

Figure 5: Sutton & Barto

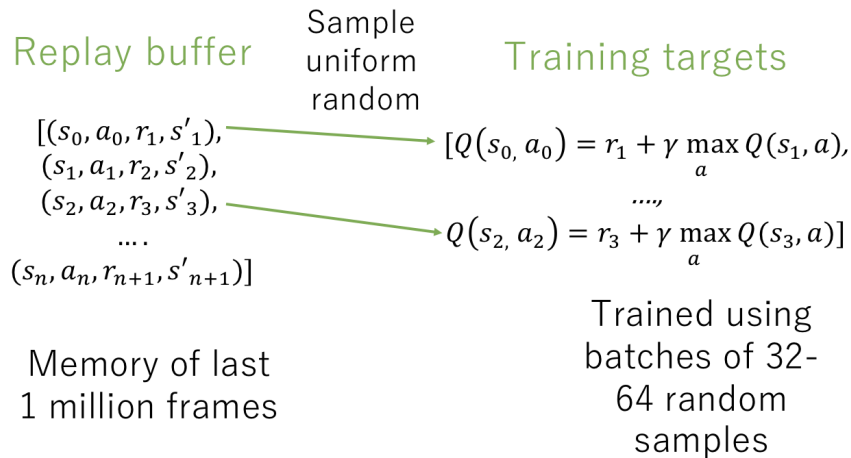


Figure 6: fig

Naive experience replay

Experience replay liberates agents from learning from experience as sampled. Prioritized experience replay liberates agents from learning from experience in the same frequency as experienced.

But some experience is more useful for learning than others. Prioritized experience replay quantifies the value of experience using the temporal difference error. Evidence of more frequent replay of high temporal difference error experiences has been found in the hippocampus of rodents.

$$error = r + \gamma Q(s', a) - Q(s, a)$$

The temporal difference error can be thought of a measure of surprise - a high temporal difference error means that the experienced transition gave a different reward than the value function expected.

Non-random sampling introduces two problems

1. loss of diversity - we will only sample from high TD error experiences
2. introduce bias - non-independent sampling

Schaul et. al (2016) solves these problems by

1. loss of diversity -> make the prioritization stochastic
2. correct bias -> use importance sampling

Stochastic prioritization

Greedy sampling of experience from memory is simply taking the highest priority samples from memory. New experiences are given the maximal priority in order to guarantee replay.

A binary heap structure is used for the priority queue, which is $\mathcal{O}(1)$ to find the highest priority sample and $\mathcal{O}(\log N)$ to update priorities after learning.

Greedy sampling has the following problems

- low TD errors on the first visit mean experience will not be replayed
- noise in the environment rewards effects the rate of sampling (which compounds with noise in the value function approximation)
- focus on a small subset of experience, making the value function prone to overfitting

These problems are overcome using a stochastic sampling procedure. The probability of sampling experience $P(i)$ is based on the probabilities:

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

Where α controls the strength of prioritization. Schaul et. al propose two methods for how we assign the priorities p_i :

1. directly based on the temporal difference error - $p_i = \delta + \epsilon$, where ϵ is a small constant to incentive revisiting with zero error experience
2. rank based - $p_i = \frac{1}{rank(i)}$, where the rank is based on the temporal difference error δ

The rank based will be more robust and less sensitive to outliers.

TODO - at the bottom of page 4

DDQN

Reference = TODO

DDQN = Double Deep Q-Network - first introduced in a tabular setting in 2010 - reintroduced in the context of DQN in 2016

DDQN aims to overcome the **maximization bias** of Q-Learning

Maximization bias

Imagine a state where $Q(s, a) = 0$ for all a

Our estimates are normally distributed above and below 0

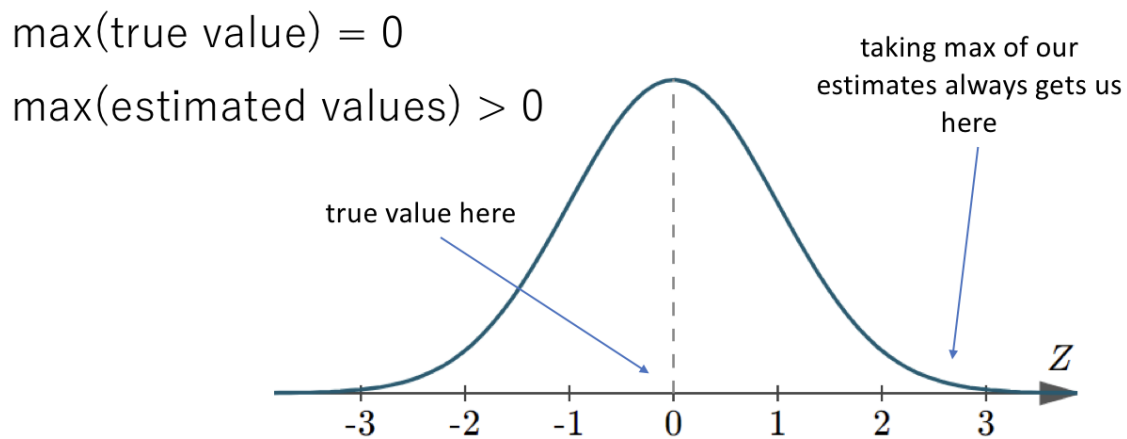


Figure 7: Taking the maximum across our estimates makes our estimate positively biased

The DDQN modification to DQN makes use of the target network as a different function to approximate $Q(s,a)$

Original DQN target

$$r + \gamma \max_a Q(s, a; \theta^-)$$

DDQN target

$$r + \gamma Q(s', \operatorname{argmax}_a Q(s', a; \theta); \theta^-)$$

- select the action according to the online network
- quantify the value that action using the target network

Distributional Q-Learning

ref = TODO

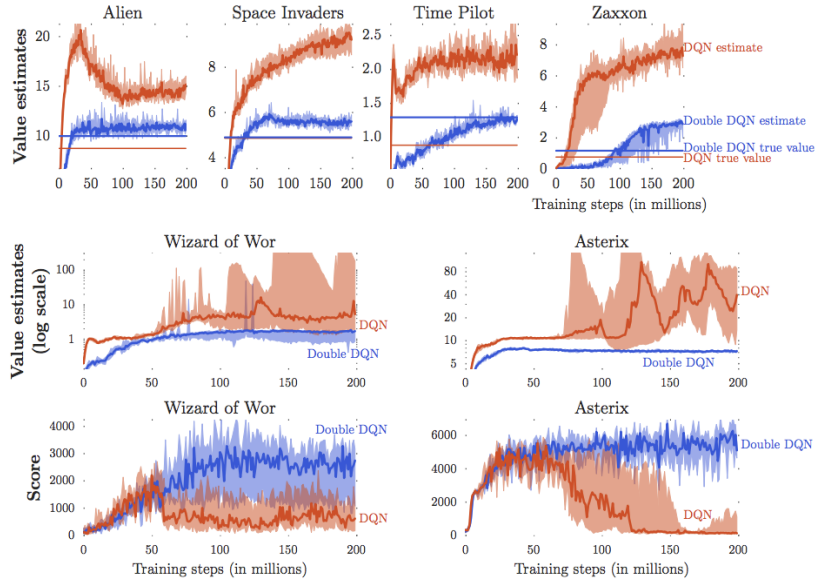


Figure 3: The **top** and **middle** rows show value estimates by DQN (orange) and Double DQN (blue) on six Atari games. The results are obtained by running DQN and Double DQN with 6 different random seeds with the hyper-parameters employed by Mnih et al. (2015). The darker line shows the median over seeds and we average the two extreme values to obtain the shaded area (i.e., 10% and 90% quantiles with linear interpolation). The straight horizontal orange (for DQN) and blue (for Double DQN) lines in the top row are computed by running the corresponding agents after learning concluded, and averaging the actual discounted return obtained from each visited state. These straight lines would match the learning curves at the right side of the plots if there is no bias. The **middle** row shows the value estimates (in log scale) for two games in which DQN’s overoptimism is quite extreme. The **bottom** row shows the detrimental effect of this on the score achieved by the agent as it is evaluated during training: the scores drop when the overestimations begin. Learning with Double DQN is much more stable.

Figure 8: fig

Beyond the expectation

All the reinforcement learning we have seen focuses on the expectation (i.e. the mean)

$$Q(s, a) = \mathbb{E}[G_t] = \mathbb{E}[r + \gamma Q(s', a)]$$

In 2017 DeepMind introduced the idea of the value distribution

State of the art results on Atari (at the time - Rainbow is currently SOTA)



The expectation of 7.5 min will never occur in reality!

Figure 9: fig

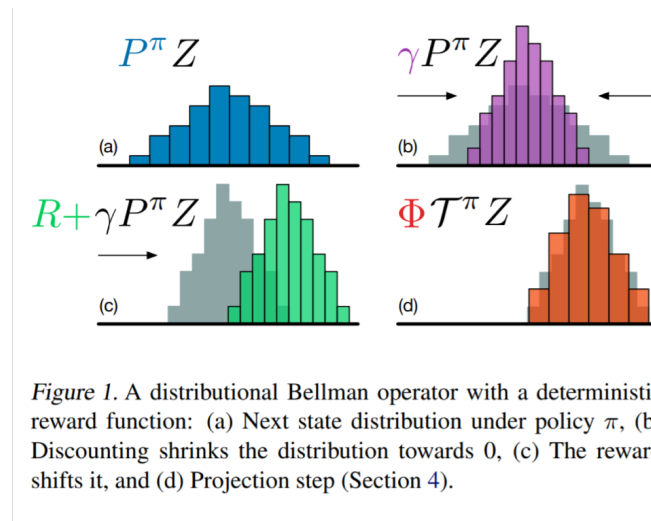


Figure 10: Bellamare et. al 2017

Rainbow

All the various improvements to DQN address different issues

	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
C51	701%	178%	40	50
UNREAL [†]	880%	250%	-	-

Figure 6. Mean and median scores across 57 Atari games, measured as percentages of human baseline (H.B., [Nair et al., 2015](#)).

[†] The UNREAL results are not altogether comparable, as they were generated in the asynchronous setting with per-game hyperparameter tuning ([Jaderberg et al., 2017](#)).

Figure 11: Bellamare et. al 2017

- DDQN - overestimation bias
- prioritized experience replay - sample efficiency
- dueling - generalize across actions
- multi-step bootstrap targets - bias variance tradeoff
- distributional Q-learning - learn categorical distribution of $Q(s, a)$
- noisy DQN - stochastic layers for exploration

Rainbow combines these improvements

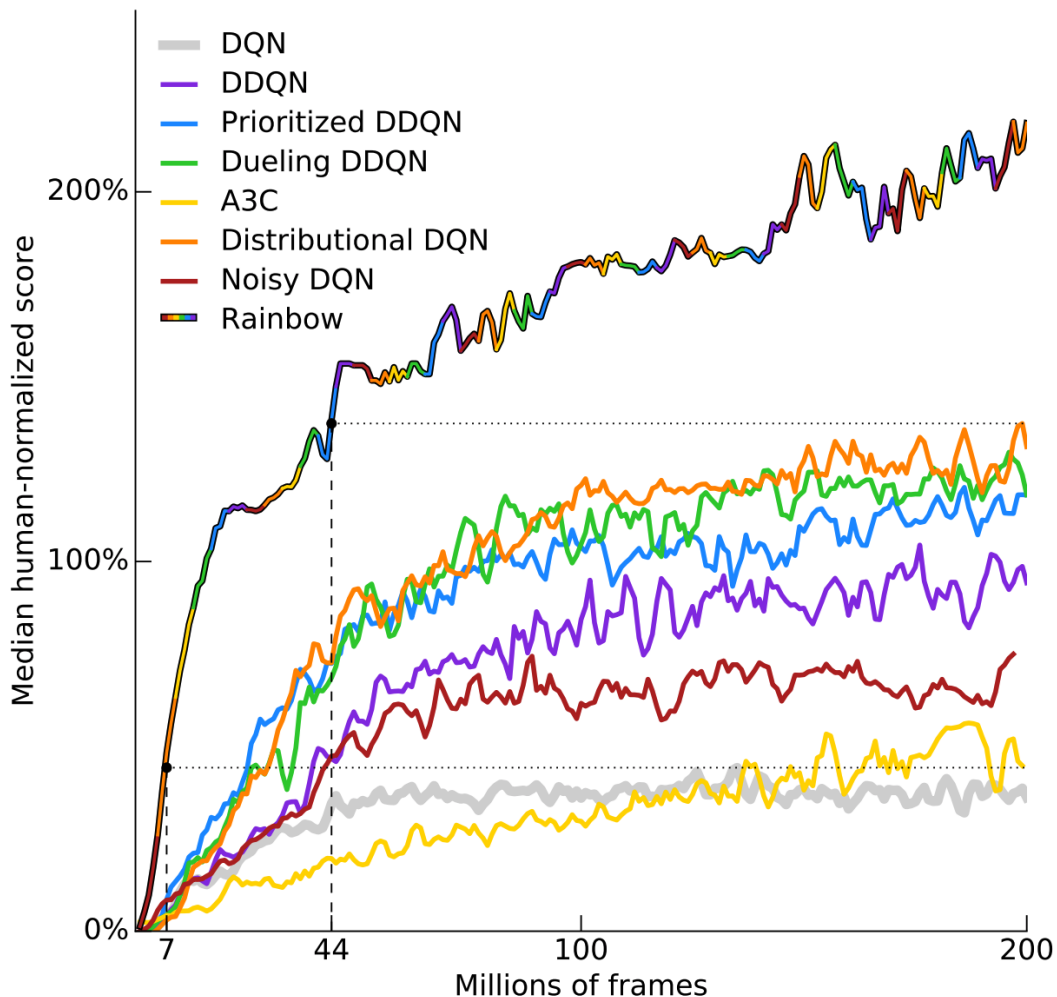


Figure 12: Median human-normalized performance across 57 Atari games. We compare our integrated agent (rainbowcolored) to DQN (grey) and six published baselines. Note that we match DQN’s best performance after 7M frames, surpass any baseline within 44M frames, and reach substantially improved final performance. Curves are smoothed with a moving average over 5 points.

Evaluation Methodology. We evaluated all agents on 57 Atari 2600 games from the arcade learning environment (Bellemare et al. 2013). We follow the training and evaluation procedures of Mnih et al. (2015) and van Hasselt et al. (2016). The average scores of the agent are evaluated during training, every 1M steps in the environment, by suspending learning and evaluating the latest agent for 500K frames. Episodes are truncated at 108K frames (or 30 minutes of simulated play), as in van Hasselt et al. (2016).

Figure 13: fig

Parameter	Value
Min history to start learning	80K frames
Adam learning rate	0.0000625
Exploration ϵ	0.0
Noisy Nets σ_0	0.5
Target Network Period	32K frames
Adam ϵ	1.5×10^{-4}
Prioritization type	proportional
Prioritization exponent ω	0.5
Prioritization importance sampling β	$0.4 \rightarrow 1.0$
Multi-step returns n	3
Distributional atoms	51
Distributional min/max values	$[-10, 10]$

Table 1: Rainbow hyper-parameters

Figure 14: fig

Learning speed. As in the original DQN setup, we ran each agent on a single GPU. The 7M frames required to match DQN’s final performance correspond to less than 10 hours of wall-clock time. A full run of 200M frames corresponds to approximately 10 days, and this varies by less than 20% between all of the discussed variants. The litera-

Agent	no-ops	human starts
DQN	79%	68%
DDQN (*)	117%	110%
Prioritized DDQN (*)	140%	128%
Dueling DDQN (*)	151%	117%
A3C (*)	-	116%
Noisy DQN	118%	102%
Distributional DQN	164%	125%
Rainbow	223%	153%

Table 2: Median normalized scores of the best agent snapshots for Rainbow and baselines. For methods marked with an asterisk, the scores come from the corresponding publication. DQN’s scores comes from the dueling networks paper, since DQN’s paper did not report scores for all 57 games. The others scores come from our own implementations.

Figure 15: fig