

DEEP REINFORCEMENT LEARNING WITH DOUBLE Q-LEARNING

HADO VAN HASSELT, ARTHUR GUEZ, AND DAVID SILVER
GOOGLE DEEPMIND

ABSTRACT. The popular Q-learning algorithm is known to overestimate action values under certain conditions. It was not previously known whether, in practice, such overestimations are common, whether this harms performance, and whether they can generally be prevented. In this paper, we answer all these questions affirmatively. In particular, we first show that the recent DQN algorithm, which combines Q-learning with a deep neural network, suffers from substantial overestimations in some games in the Atari 2600 domain. We then show that the idea behind the Double Q-learning algorithm, which was introduced in a tabular setting, can be generalized to work with large-scale function approximation. We propose a specific adaptation to the DQN algorithm and show that the resulting algorithm not only reduces the observed overestimations, as hypothesized, but that this also leads to much better performance on several games.

The goal of reinforcement learning (Sutton and Barto, 1998) is to learn good policies for sequential decision problems, by optimizing a cumulative future reward signal. Q-learning (Watkins, 1989) is one of the most popular reinforcement learning algorithms, but it is known to sometimes learn unrealistically high action values because it includes a maximization step over estimated action values, which tends to prefer overestimated to underestimated values.

In previous work, overestimations have been attributed to insufficiently flexible function approximation (Thrun and Schwartz, 1993) and noise (van Hasselt, 2010, 2011). In this paper, we unify these views and show overestimations can occur when the action values are inaccurate, irrespective of the source of approximation error. Of course, imprecise value estimates are the norm during learning, which indicates that overestimations may be much more common than previously appreciated.

It is an open question whether, if the overestimations do occur, this negatively affects performance in practice. Overoptimistic value estimates are not necessarily a problem in and of themselves. If all values would be uniformly higher then the relative action preferences are preserved and we would not expect the resulting policy to be any worse. Furthermore, it is known that sometimes it is good to be optimistic: optimism in the face of uncertainty is a well-known exploration technique (Kaelbling et al., 1996). If, however, the overestimations are not uniform and not concentrated at states about which we wish to learn more, then they might negatively affect the quality of the resulting policy. Thrun and Schwartz (1993) give specific examples in which this leads to suboptimal policies, even asymptotically.

To test whether overestimations occur in practice and at scale, we investigate the performance of the recent DQN algorithm (Mnih et al., 2015). DQN combines Q-learning with a flexible

deep neural network and was tested on a varied and large set of deterministic Atari 2600 games, reaching human-level performance on many games. In some ways, this setting is a best-case scenario for Q-learning, because the deep neural network provides flexible function approximation with the potential for a low asymptotic approximation error, and the determinism of the environments prevents the harmful effects of noise. Perhaps surprisingly, we show that even in this comparatively favorable setting DQN sometimes substantially overestimates the values of the actions.

We show that the idea behind the Double Q-learning algorithm (van Hasselt, 2010), which was first proposed in a tabular setting, can be generalized to work with arbitrary function approximation, including deep neural networks. We use this to construct a new algorithm we call Double DQN. We then show that this algorithm not only yields more accurate value estimates, but leads to much higher scores on several games. This demonstrates that the overestimations of DQN were indeed leading to poorer policies and that it is beneficial to reduce them. In addition, by improving upon DQN we obtain state-of-the-art results on the Atari domain.

BACKGROUND

To solve sequential decision problems we can learn estimates for the optimal value of each action, defined as the expected sum of future rewards when taking that action and following the optimal policy thereafter. Under a given policy π , the true value of an action a in a state s is

$$Q_{\pi}(s, a) \equiv \mathbb{E}[R_1 + \gamma R_2 + \dots \mid S_0 = s, A_0 = a, \pi],$$

where $\gamma \in [0, 1]$ is a discount factor that trades off the importance of immediate and later rewards. The optimal value is then $Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$. An optimal policy is easily derived from the optimal values by selecting the highest-valued action in each state.

Estimates for the optimal action values can be learned using Q-learning (Watkins, 1989), a form of temporal difference learning (Sutton, 1988). Most interesting problems are too large to learn all action values in all states separately. Instead, we can learn a parameterized value function $Q(s, a; \theta_t)$. The standard Q-learning update for the parameters after taking action A_t in state S_t and observing the immediate reward R_{t+1} and resulting state S_{t+1} is then

$$(1) \quad \theta_{t+1} = \theta_t + \alpha(Y_t^Q - Q(S_t, A_t; \theta_t)) \nabla_{\theta_t} Q(S_t, A_t; \theta_t).$$

where α is a scalar step size and the target Y_t^Q is defined as

$$(2) \quad Y_t^Q \equiv R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t).$$

This update resembles stochastic gradient descent, updating the current value $Q(S_t, A_t; \theta_t)$ towards a target value Y_t^Q .

Deep Q Networks. A deep Q network (DQN) is a multi-layered neural network that for a given state s outputs a vector of action values $Q(s, \cdot; \theta)$, where θ are the parameters of the network. For an n -dimensional state space and an action space containing m

actions, the neural network is a function from \mathbb{R}^n to \mathbb{R}^m . Two important ingredients of the DQN algorithm as proposed by Mnih et al. (2015) are the use of a target network, and the use of experience replay. The target network, with parameters θ^- , is the same as the online network except that its parameters are copied every τ steps from the online network, so that then $\theta_t^- = \theta_t$, and kept fixed on all other steps. The target used by DQN is then

$$(3) \quad Y_t^{\text{DQN}} \equiv R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t^-).$$

For the experience replay (Lin, 1992), observed transitions are stored for some time and sampled uniformly from this memory bank to update the network. Both the target network and the experience replay dramatically improve the performance of the algorithm (Mnih et al., 2015).

Double Q-learning. The max operator in standard Q-learning and DQN, in (2) and (3), uses the same values both to select and to evaluate an action. This makes it more likely to select overestimated values, resulting in overoptimistic value estimates. To prevent this, we can decouple the selection from the evaluation. This is the idea behind Double Q-learning (van Hasselt, 2010).

In the original Double Q-learning algorithm, two value functions are learned by assigning each experience randomly to update one of the two value functions, such that there are two sets of weights, θ and θ' . For each update, one set of weights is used to determine the greedy policy and the other to determine its value. For a clear comparison, we can first untangle the selection and evaluation in Q-learning and rewrite its target (2) as

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

The Double Q-learning error can then be written as

$$(4) \quad Y_t^{\text{DoubleQ}} \equiv R_{t+1} + \gamma Q(S_{t+1}, \arg\max_a Q(S_{t+1}, a; \theta_t); \theta_t').$$

Notice that the selection of the action, in the $\arg\max$, is still due to the online weights θ_t . This means that, as in Q-learning, we are still estimating the value of the greedy policy according to the current values, as defined by θ_t . However, we use the second set of weights θ_t' to fairly evaluate the value of this policy. This second set of weights can be updated symmetrically by switching the roles of θ and θ' .

OVEROPTIMISM DUE TO ESTIMATION ERRORS

Q-learning's overestimations were first investigated by Thrun and Schwartz (1993), who showed that if the action values contain random errors uniformly distributed in an interval $[-\epsilon, \epsilon]$ then each target is overestimated up to $\gamma\epsilon\frac{m-1}{m+1}$, where m is the number of actions. In addition, Thrun and Schwartz give a concrete example in which these overestimations even asymptotically lead to sub-optimal policies, and show the overestimations manifest themselves in a small toy problem when using function approximation. Later van Hasselt (2010) argued that noise in the environment can lead to overestimations even when using tabular representation, and proposed Double Q-learning as a solution.

In this section we demonstrate more generally that estimation errors of any kind can induce an upward bias, regardless of whether these errors are due to environmental noise, function

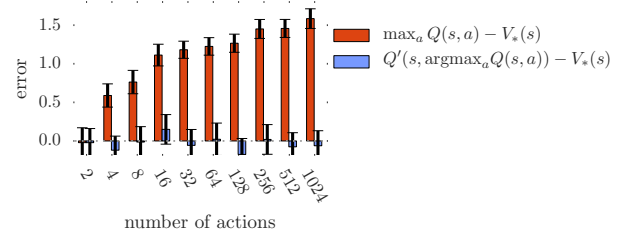
approximation, non-stationarity, or any other source. This is important, because in practice any method will incur some inaccuracies during learning, simply due to the fact that the true values are initially unknown.

The result by Thrun and Schwartz (1993) cited above gives an upper bound to the overestimation for a specific setup, but it is also possible, and potentially more interesting, to derive a lower bound.

Theorem 1. Consider a state s in which all the true optimal action values are equal at $Q_*(s, a) = V_*(s)$ for some $V_*(s)$. Let Q_t be arbitrary value estimates that are on the whole unbiased in the sense that $\sum_a (Q_t(s, a) - V_*(s)) = 0$, but that are not all correct, such that $\frac{1}{m} \sum_a (Q_t(s, a) - V_*(s))^2 = C$ for some $C > 0$, where $m \geq 2$ is the number of actions in s . Under these conditions, $\max_a Q_t(s, a) \geq V_*(s) + \sqrt{\frac{C}{m-1}}$. This lower bound is tight. Under the same conditions, the lower bound on the absolute error of the Double Q-learning estimate is zero. (Proof in appendix.)

Note that we did not need to assume that estimation errors for different actions are independent. This theorem shows that even if the value estimates are on average correct, estimation errors of any source can drive the estimates up and away from the true optimal values.

The lower bound in Theorem 1 decreases with the number of actions. This is an artifact of considering the lower bound, which requires very specific values to be attained. More typically, the overoptimism increases with the number of actions as shown here:



The orange bars show the bias in a single Q-learning update when the action values are $Q(s, a) = V_*(s) + \epsilon_a$ and the errors $\{\epsilon_a\}_{a=1}^m$ are independent standard normal random variables. The second set of action values Q' , used for the blue bars, was generated identically and independently. All bars are the average of 100 repetitions. Q-learning's overestimations here indeed increase with the number of actions. Double Q-learning is unbiased. As another example, if for all actions $Q_*(s, a) = V_*(s)$ and the estimation errors $Q_t(s, a) - V_*(s)$ are uniformly random in $[-1, 1]$, then the overoptimism is $\frac{m-1}{m+1}$. (Proof in appendix.)

We now turn to function approximation and consider a real-valued continuous state space with 10 discrete actions in each state. For simplicity, the true optimal action values in this example depend only on state so that in each state all actions have the same true value. These true values are shown in the left column of plots in Figure 1 (purple lines) and are defined as either $Q_*(s, a) = \sin(s)$ (top row) or $Q_*(s, a) = 2 \exp(-s^2)$ (middle and bottom rows). The left plots also show an approximation for a single action (green lines) as a function of state as well as the samples the estimate is based on (green dots). The estimate is a d -degree polynomial that is fit to the true values at sampled states, where $d = 6$ (top and middle rows) or $d = 9$ (bottom row). The

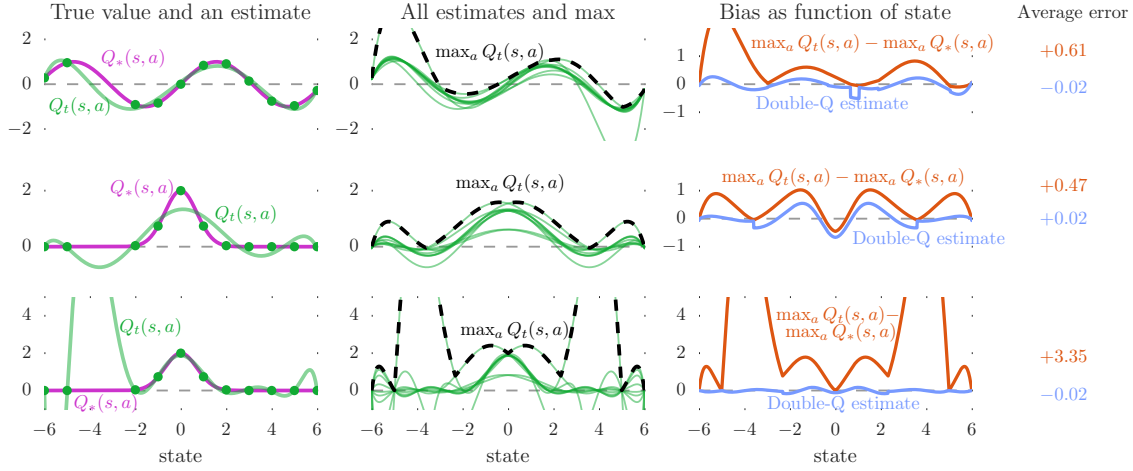


FIGURE 1. Illustration of overestimations during learning. In each state (x-axis), there are 10 actions. The **left column** shows the true values $V_*(s)$ (purple line). All true action values are defined by $Q_*(s, a) = V_*(s)$. The green line shows estimated values $Q(s, a)$ for one action as a function of state, fitted to the true value at several sampled states (green dots). The **middle column** plots shows all the estimated values (green), and the maximum of these values (dashed black). The maximum is higher than the true value (purple, left plot) almost everywhere. The **right column** plots shows the difference in orange. The blue line in the right plots is the estimate used by Double Q-learning with a second set of samples for each. The blue line is much closer to zero, indicating less bias. The three **rows** correspond to different true functions (left, purple) or capacities of the fitted function (left, green). (Details in the text)

samples match the true function exactly: there is no noise and we assume we have **ground truth** for the action value on these sampled states. The approximation is inexact even on the sampled states for the top two rows because the function approximation is insufficiently flexible. In the bottom row, the function is flexible enough to fit the green dots, but this reduces the accuracy in unsampled states. Notice that the sampled states are spaced further apart near the left side of the left plots, resulting in larger estimation errors. In many ways this is a typical learning setting, where at each point in time we only have limited data.

The middle column of plots in Figure 1 shows estimated action value functions for all 10 actions (green lines), as functions of state, along with the maximum action value in each state (black dashed line). Although the true value function is the same for all actions, the approximations differ because we have supplied different sets of sampled states.¹ The maximum is often higher than the ground truth shown in purple on the left. This is confirmed in the right plots, which shows the difference between the black and purple curves in orange. The orange line is almost always positive, indicating an upward bias. The right plots also show the estimates from Double Q-learning in blue², which are on average much closer to zero. This demonstrates that Double Q-learning indeed can successfully reduce the overoptimism of Q-learning.

The different rows in Figure 1 show variations of the same experiment. The difference between the top and middle rows is the true value function, demonstrating that overestimations

are not an artifact of a specific true value function. The difference between the middle and bottom rows is the flexibility of the function approximation. In the left-middle plot, the estimates are even incorrect for some of the sampled states because the function is insufficiently flexible. The function in the bottom-left plot is more flexible but this causes higher estimation errors for unseen states, resulting in higher overestimations. This is important because flexible parametric function approximators are often employed in reinforcement learning (see, e.g., Tesauro 1995; Sallans and Hinton 2004; Riedmiller 2005; Mnih et al. 2015).

In contrast to van Hasselt (2010) we did not use a statistical argument to find overestimations, the process to obtain Figure 1 is fully deterministic. In contrast to Thrun and Schwartz (1993), we did not rely on inflexible function approximation with irreducible asymptotic errors; the bottom row shows that a function that is flexible enough to cover all samples leads to high overestimations.

In the examples above, overestimations occur even when assuming we have samples of the *true* action value at certain states. The value estimates can further deteriorate if we bootstrap off of action values that are already overoptimistic, since this causes overestimations to propagate throughout our estimates. Although *uniformly* overestimating values might not hurt the resulting policy, in practice overestimation errors will differ for different states and actions. Overestimation combined with bootstrapping then has the pernicious effect of propagating the wrong relative information about which states are more valuable than others, directly affecting the quality of the learned policies.

The overestimations should not be confused with optimism in the face of uncertainty (Sutton, 1990; Agrawal, 1995; Kaelbling et al., 1996; Auer et al., 2002; Brafman and Tennenholtz, 2003; Szita and Lőrincz, 2008; Strehl et al., 2009), where an exploration bonus is given to states or actions with uncertain values. Conversely, the overestimations discussed here occur only

¹Each action-value function is fit with a different subset of integer states. States -6 and 6 are always included to avoid extrapolations, and for each action two adjacent integers are missing: for action a_1 states -5 and -4 are not sampled, for a_2 states -4 and -3 are not sampled, and so on. This causes the estimated values to differ.

²We arbitrarily used the samples of action a_{i+5} (for $i \leq 5$) or a_{i-5} (for $i > 5$) as the second set of samples for the double estimator of action a_i .

after updating, resulting in overoptimism in the face of apparent certainty. This was already observed by Thrun and Schwartz (1993), who noted that, in contrast to optimism in the face of uncertainty, these overestimations actually can impede learning an optimal policy. We will see this negative effect on policy quality confirmed later in the experiments as well: when we reduce the overestimations using Double Q-learning, the policies improve.

DOUBLE DQN

A main goal of this paper is to investigate whether the overestimations of Q-learning occur in practice and, when they do occur, if they hurt performance. To test these hypotheses, in the next section we analyze the performance of the DQN algorithm and compare it to the Double DQN algorithm that we will now construct. The idea of Double Q-learning is to decompose the max operation in the target into action selection and action evaluation. Although not fully decoupled, the target network in the DQN architecture provides a natural candidate for the second value function, without having to introduce additional networks. We therefore propose to evaluate the greedy policy according to the online network, but using the target network to estimate its value. In reference to both Double Q-learning and DQN, we refer to the resulting learning algorithm as Double DQN. The update used by Double DQN is the same as for DQN, but replacing the target Y_t^{DQN} with

$$Y_t^{\text{DoubleDQN}} \equiv R_{t+1} + \gamma Q(S_{t+1}, \arg\max_a Q(S_{t+1}, a; \theta_t), \theta_t^-).$$

In comparison to Double Q-learning (4), the weights of the second network θ_t' are replaced with the weights of the target network θ_t^- for the evaluation of the current greedy policy. The update to the target network stays unchanged from DQN, and remains a periodic copy of the online network.

This version of Double DQN is perhaps the minimal possible change to DQN towards Double Q-learning. The goal is to get most of the benefit of Double Q-learning, while keeping the rest of the DQN algorithm intact for a fair comparison, and with minimal computational overhead.

EMPIRICAL RESULTS

In this section, we analyze the overestimations of DQN and show that Double DQN improves over DQN both in terms of value accuracy and in terms of policy quality. To further test the robustness of the approach we additionally evaluate the algorithms with random starts generated from expert human trajectories, as proposed by Nair et al. (2015).

Our testbed contains 49 Atari 2600 games, using the Arcade Learning Environment (Bellemare et al., 2013). The goal is for a single algorithm, with a fixed set of hyperparameters, to learn to play each of the games separately from interaction given only the screen pixels as input. This is a demanding testbed: not only are the inputs high-dimensional, the game visuals and game mechanics vary substantially between games. Good solutions must therefore rely heavily on the learning algorithm — it is not practically feasible to overfit the domain by relying only on tuning.

We closely follow the experimental setting and network architecture outlined by Mnih et al. (2015). Briefly, the network architecture is a convolutional neural network (Fukushima, 1988; LeCun et al., 1998) with 3 convolution layers and a fully-connected

hidden layer (approximately 1.5M parameters in total). The network takes the last four frames as input and outputs the action value of each action. On each game, the network is trained on a single GPU for 200M frames, or approximately 1 week.

Results on overoptimism. Figure 2 shows examples of DQN’s overestimations in six Atari games. DQN and Double DQN were both trained under the exact conditions described by Mnih et al. (2015). DQN is consistently and sometimes vastly overoptimistic about the value of the current greedy policy, as can be seen by comparing the orange learning curves in the top row of plots to the straight orange lines, which represent the actual discounted value of the best learned policy. More precisely, the (averaged) value estimates are computed regularly during training with full evaluation phases of length $T = 125,000$ steps as

$$\frac{1}{T} \sum_{t=1}^T \arg\max_a Q(S_t, a; \theta).$$

The ground truth averaged values are obtained by running the best learned policies for several episodes and computing the actual cumulative rewards. Without overestimations we would expect these quantities to match up (i.e., the curve to match the straight line at the right of each plot). Instead, the learning curves of DQN consistently end up much higher than the true values. The learning curves for Double DQN, shown in blue, are much closer to the blue straight line representing the true value of the final policy. Note that the blue straight line is often higher than the orange straight line. This indicates that Double DQN does not just produce more accurate value estimates but also better policies.

More extreme overestimations are shown in the middle two plots, where DQN is highly unstable on the games Asterix and Wizard of Wor. Notice the log scale for the values on the y -axis. The bottom two plots show the corresponding scores for these two games. Notice that the increases in value estimates for DQN in the middle plots coincide with decreasing scores in bottom plots. Again, this indicates that the overestimations are harming the quality of the resulting policies. If seen in isolation, one might perhaps have been tempted to think the observed instability is related to inherent instability problems of off-policy learning with function approximation (Baird, 1995; Tsitsiklis and Van Roy, 1997; Sutton et al., 2008; Maei, 2011; Sutton et al., 2015). However, we see that learning is much more stable with Double DQN, suggesting that the cause for these instabilities is in fact Q-learning’s overoptimism. Figure 2 only shows a few examples, but overestimations were observed for DQN in all 49 tested Atari games, albeit in varying amounts.

Quality of the learned policies. Overoptimism does not always adversely affect the quality of the learned policy. For example, DQN achieves optimal behavior in Pong despite slightly overestimating of the policy value. Nevertheless, reducing overestimations can significantly benefit the stability of the learning; we see clear examples of this in Figure 2. We now assess more generally how much Double DQN helps in terms of policy quality by evaluating on all 49 games that DQN was tested on.

As described by Mnih et al. (2015) we start each evaluation episode by executing a special action that does not affect the environment, a so-called no-op action, up to 30 times to provide

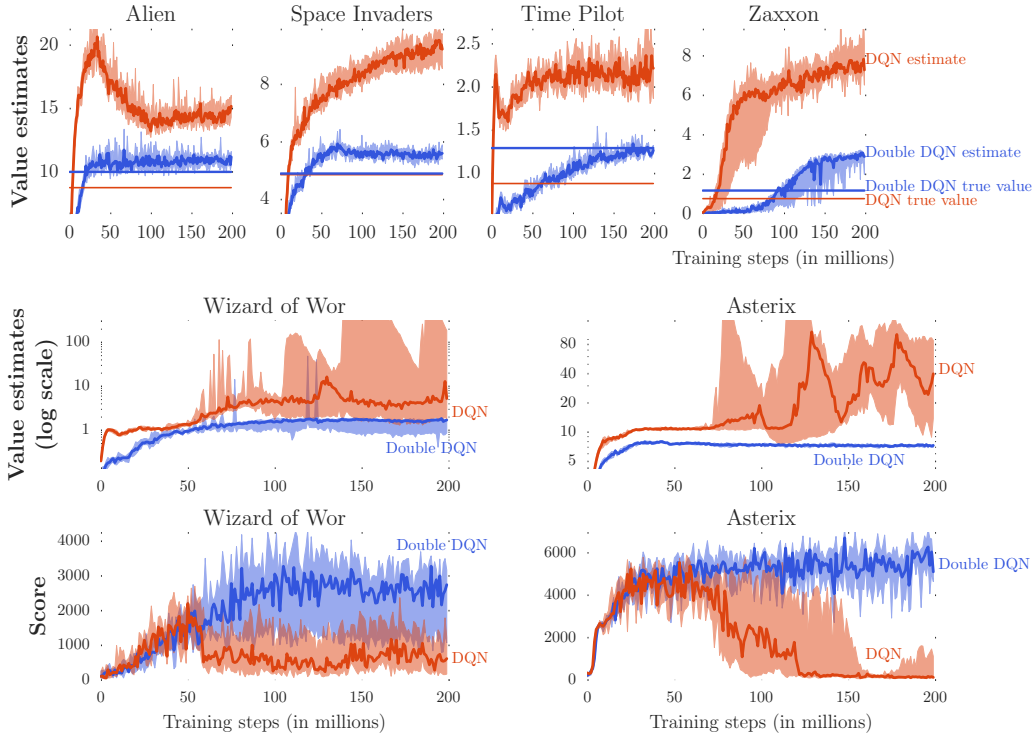


FIGURE 2. The **top** and **middle** rows show value estimates by DQN (orange) and Double DQN (blue) on six Atari games. The results are obtained by running DQN and Double DQN with 6 different random seeds with the hyper-parameters employed by Mnih et al. (2015). The darker line shows the median over 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.

different starting points for the agent. Some exploration during evaluation provides additional randomization. For Double DQN we used the exact same hyper-parameters as for DQN, to allow for a controlled experiment focused just on reducing overestimations. The learned policies are evaluated for 5 mins of emulator time (18,000 frames) with an ϵ -greedy policy where $\epsilon = 0.05$. The scores are averaged over 100 episodes. The only difference between Double DQN and DQN is the target, using $Y_t^{\text{DoubleDQN}}$ rather than Y_t^{DQN} . This evaluation is somewhat adversarial, as the used hyper-parameters were tuned for DQN but not for Double DQN.

To obtain summary statistics across games, we normalize the score for each game as follows:

$$(5) \quad \text{score}_{\text{normalized}} = \frac{\text{score}_{\text{agent}} - \text{score}_{\text{random}}}{\text{score}_{\text{human}} - \text{score}_{\text{random}}}.$$

The ‘random’ and ‘human’ scores are the same as used by Mnih et al. (2015), and are given in the appendix.

Table 1, under **no ops**, shows that on the whole Double DQN clearly improves over DQN. A detailed comparison (in appendix) shows that there are several games in which Double DQN greatly improves upon DQN. Noteworthy examples include Road Runner (from 233% to 617%), Asterix (from 70% to 180%), Zaxxon (from 54% to 111%), and Double Dunk (from 17% to 397%).

	no ops		random starts	
	DQN	Double DQN	DQN	Double DQN
Median	93.5%	114.7%	56.6%	86.8%
Mean	241.1%	330.3%	146.0%	883.3%

TABLE 1. Summary statistics of normalized performance.

In the table we have not included the Gorila algorithm (Nair et al., 2015), which is a massively distributed version of DQN, because the architecture and infrastructure is sufficiently different to make a direct comparison unclear. For reference, we note that Gorila obtained median and mean normalized scores of 96% and 495%, respectively.

Robustness to Human starts. One concern with the previous evaluation is that in deterministic games with a unique starting point the learner could potentially learn to remember sequences of actions without much need to generalize. While successful, the solution would not be particularly robust. By testing the agents from various starting points, we can test whether the found solutions generalize well, and as such provide a challenging testbed for the learned policies (Nair et al., 2015).

We obtained 100 starting points sampled for each game from a human expert’s trajectory, as proposed by Nair et al. (2015). We start an evaluation episode from each of these starting points and run the emulator for up to 108,000 frames (30 mins at 60Hz

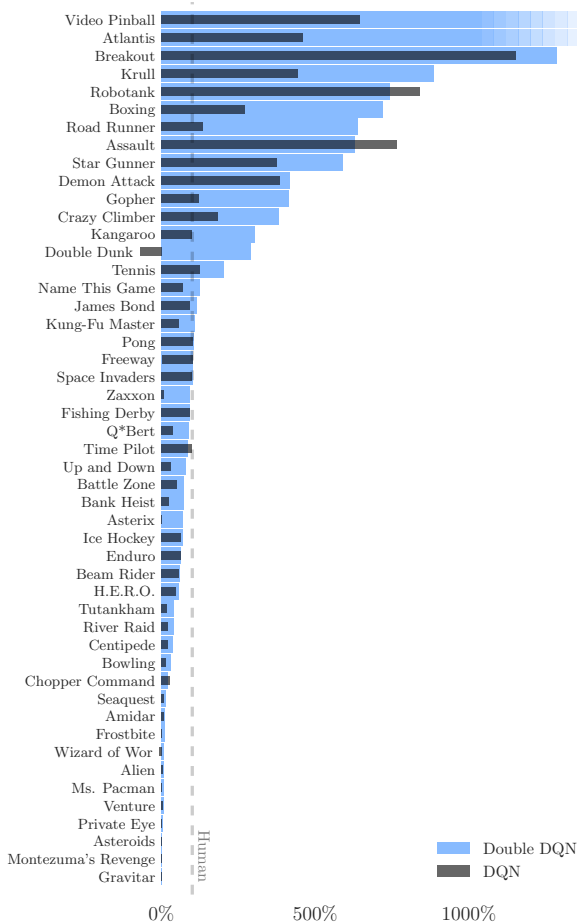


FIGURE 3. Normalized scores of DQN and Double DQN on 49 Atari 2600 games, tested for 100 episodes per game with random starts.

including the trajectory before the starting point). Each agent is only evaluated on the rewards accumulated after the starting point.

Table 1, under **random starts**, reports summary statistics for this evaluation. Double DQN obtains clearly higher median and mean scores. Again we propose not to compare directly to Gorila DQN (Nair et al., 2015), but for completeness note it obtained a median of 78% and a mean of 259%. Detailed results are available in Figure 3 and in the appendix. On several games the improvements are striking, in some cases bringing scores much closer to human, or even surpassing these. This includes Asterix (from -1% to 69%), Bank Heist (from 26% to 72%), Double Dunk (from -71% to 290%), Q-Bert (from 37% to 88%), Up and Down (from 30% to 80%), and Zaxxon (from 9% to 93%). Only on two games, Assault and Robotank, we see a notable decrease in performance, although in both cases the performance stays well above human level.

Double DQN appears more robust to this more challenging evaluation, suggesting that appropriate generalizations occur and that the found solutions do not exploit the determinism of the environments. This is appealing, as it indicates progress towards finding general solutions rather than a deterministic sequence of steps that would be less robust.

DISCUSSION

This paper has five contributions. First, we have shown why Q-learning can be overoptimistic in large-scale problems, even if these are deterministic, due to the inherent estimation errors of learning. Second, by analyzing the value estimates on Atari games we have shown that these overestimations are more common and severe in practice than previously acknowledged. Third, we have shown that Double Q-learning can be used at scale to successfully reduce this overoptimism, resulting in more stable and reliable learning. Fourth, we have proposed a specific implementation called Double DQN, that uses the existing architecture and deep neural network of the DQN algorithm without requiring additional networks or parameters. Finally, we have shown that Double DQN finds better policies, obtaining new state-of-the-art results on the Atari 2600 domain.

ACKNOWLEDGMENTS

We would like to thank Tom Schaul, Volodymyr Mnih, Marc Bellemare, Thomas Degris, and Richard Sutton for helpful comments on earlier versions of this document and the ideas presented therein, and everyone at Google DeepMind for technical and non-technical support and for helping create a constructive research environment.

REFERENCES

- R. Agrawal. Sample mean based index policies with $O(\log n)$ regret for the multi-armed bandit problem. *Advances in Applied Probability*, pages 1054–1078, 1995.
- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- L. Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning: Proceedings of the Twelfth International Conference*, pages 30–37, 1995.
- M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *J. Artif. Intell. Res. (JAIR)*, 47:253–279, 2013.
- R. I. Brafman and M. Tennenholtz. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *The Journal of Machine Learning Research*, 3:213–231, 2003.
- K. Fukushima. Neocognitron: A hierarchical neural network capable of visual pattern recognition. *Neural networks*, 1(2):119–130, 1988.
- L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4:237–285, 1996.
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- L. Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3):293–321, 1992.
- H. R. Maei. *Gradient temporal-difference learning algorithms*. PhD thesis, University of Alberta, 2011.
- V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- A. Nair, P. Srinivasan, S. Blackwell, C. Alcicek, R. Fearon, A. D. Maria, V. Panneershelvam, M. Suleyman, C. Beattie, S. Petersen, S. Legg, V. Mnih, K. Kavukcuoglu, and D. Silver. Massively parallel methods

- for deep reinforcement learning. In *Deep Learning Workshop, ICML*, 2015.
- M. Riedmiller. Neural fitted Q iteration - first experiences with a data efficient neural reinforcement learning method. In J. Gama, R. Camacho, P. Brazdil, A. Jorge, and L. Torgo, editors, *Proceedings of the 16th European Conference on Machine Learning (ECML'05)*, pages 317–328. Springer, 2005.
- B. Sallans and G. E. Hinton. Reinforcement learning with factored states and actions. *The Journal of Machine Learning Research*, 5: 1063–1088, 2004.
- A. L. Strehl, L. Li, and M. L. Littman. Reinforcement learning in finite MDPs: PAC analysis. *The Journal of Machine Learning Research*, 10:2413–2444, 2009.
- R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44, 1988.
- R. S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Proceedings of the seventh international conference on machine learning*, pages 216–224, 1990.
- R. S. Sutton and A. G. Barto. *Introduction to reinforcement learning*. MIT Press, 1998.
- R. S. Sutton, C. Szepesvári, and H. R. Maei. A convergent $O(n)$ algorithm for off-policy temporal-difference learning with linear function approximation. *Advances in Neural Information Processing Systems 21 (NIPS-08)*, 21:1609–1616, 2008.
- R. S. Sutton, A. R. Mahmood, and M. White. An emphatic approach to the problem of off-policy temporal-difference learning. *arXiv preprint arXiv:1503.04269*, 2015.
- I. Szita and A. Lőrincz. The many faces of optimism: a unifying approach. In *Proceedings of the 25th international conference on Machine learning*, pages 1048–1055. ACM, 2008.
- G. Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- S. Thrun and A. Schwartz. Issues in using function approximation for reinforcement learning. In M. Mozer, P. Smolensky, D. Touretzky, J. Elman, and A. Weigend, editors, *Proceedings of the 1993 Connectionist Models Summer School*, Hillsdale, NJ, 1993. Lawrence Erlbaum.
- J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690, 1997.
- H. van Hasselt. Double Q-learning. *Advances in Neural Information Processing Systems*, 23:2613–2621, 2010.
- H. van Hasselt. *Insights in Reinforcement Learning*. PhD thesis, Utrecht University, 2011.
- C. J. C. H. Watkins. *Learning from delayed rewards*. PhD thesis, University of Cambridge England, 1989.

APPENDIX

Theorem 1. Consider a state s in which all the true optimal action values are equal at $Q_*(s, a) = V_*(s)$ for some $V_*(s)$. Let Q_t be arbitrary value estimates that are on the whole unbiased in the sense that $\sum_a (Q_t(s, a) - V_*(s)) = 0$, but that are not all zero, such that $\frac{1}{m} \sum_a (Q_t(s, a) - V_*(s))^2 = C$ for some $C > 0$, where $m \geq 2$ is the number of actions in s . Under these conditions, $\max_a Q_t(s, a) \geq V_*(s) + \sqrt{\frac{C}{m-1}}$. This lower bound is tight. Under the same conditions, the lower bound on the absolute error of the Double Q-learning estimate is zero.

Proof of Theorem 1. Define the errors for each action a as $\epsilon_a = Q_t(s, a) - V_*(s)$. Suppose that there exists a setting of $\{\epsilon_a\}$ such that $\max_a \epsilon_a < \sqrt{\frac{C}{m-1}}$. Let $\{\epsilon_i^+\}$ be the set of positive ϵ of size n , and $\{\epsilon_j^-\}$ the set of strictly negative ϵ of size $m - n$ (such that $\{\epsilon\} = \{\epsilon_i^+\} \cup \{\epsilon_j^-\}$). If $n = m$, then $\sum_a \epsilon_a = 0 \implies \epsilon_a = 0 \forall a$, which contradicts $\sum_a \epsilon_a^2 = mC$. Hence, it must be that $n \leq m - 1$. Then, $\sum_{i=1}^n \epsilon_i^+ \leq n \max_i \epsilon_i^+ < n \sqrt{\frac{C}{m-1}}$, and therefore (using the constraint $\sum_a \epsilon_a = 0$) we also have that $\sum_{j=1}^{m-n} |\epsilon_j^-| < n \sqrt{\frac{C}{m-1}}$. This implies $\max_j |\epsilon_j^-| < n \sqrt{\frac{C}{m-1}}$. By Hölder’s inequality, then

$$\begin{aligned} \sum_{j=1}^{m-n} (\epsilon_j^-)^2 &\leq \sum_{j=1}^{m-n} |\epsilon_j^-| \cdot \max_j |\epsilon_j^-| \\ &< n \sqrt{\frac{C}{m-1}} n \sqrt{\frac{C}{m-1}}. \end{aligned}$$

We can now combine these relations to compute an upper-bound on the sum of squares for all ϵ_a :

$$\begin{aligned} \sum_{a=1}^m (\epsilon_a)^2 &= \sum_{i=1}^n (\epsilon_i^+)^2 + \sum_{j=1}^{m-n} (\epsilon_j^-)^2 \\ &< n \frac{C}{m-1} + n \sqrt{\frac{C}{m-1}} n \sqrt{\frac{C}{m-1}} \\ &= C \frac{n(n+1)}{m-1} \\ &\leq mC. \end{aligned}$$

This contradicts the assumption that $\sum_{a=1}^m \epsilon_a^2 < mC$, and therefore $\max_a \epsilon_a \geq \sqrt{\frac{C}{m-1}}$ for all settings of ϵ that satisfy the constraints. We can check that the lower-bound is tight by setting $\epsilon_a = \sqrt{\frac{C}{m-1}}$ for $a = 1, \dots, m-1$ and $\epsilon_m = -\sqrt{(m-1)C}$. This verifies $\sum_a \epsilon_a^2 = mC$ and $\sum_a \epsilon_a = 0$.

The only tight lower bound on the absolute error for Double Q-learning $|Q'_t(s, \arg\max_a Q_t(s, a)) - V_*(s)|$ is zero. This can be seen by because we can have

$$Q_t(s, a_1) = V_*(s) + \sqrt{C \frac{m-1}{m}},$$

and

$$Q_t(s, a_i) = V_*(s) - \sqrt{C \frac{1}{m(m-1)}}, \text{ for } i > 1.$$

Then the conditions of the theorem hold. If then, furthermore, we have $Q'_t(s, a_1) = V_*(s)$ then the error is zero. The remaining action values $Q'_t(s, a_i)$, for $i > 1$, are arbitrary. \square

Theorem 2. Consider a state s in which all the true optimal action values are equal at $Q_*(s, a) = V_*(s)$. Suppose that the estimation errors

$Q_t(s, a) - Q_*(s, a)$ are independently distributed uniformly randomly in $[-1, 1]$. Then,

$$\mathbb{E} \left[\max_a Q_t(s, a) - V_*(s) \right] = \frac{m-1}{m+1}$$

Proof. Define $\epsilon_a = Q_t(s, a) - Q_*(s, a)$; this is a uniform random variable in $[-1, 1]$. The probability that $\max_a Q_t(s, a) \leq x$ for some x is equal to the probability that $\epsilon_a \leq x$ for all a simultaneously. Because the estimation errors are independent, we can derive

$$\begin{aligned} P(\max_a \epsilon_a \leq x) &= P(X_1 \leq x \wedge X_2 \leq x \wedge \dots \wedge X_m \leq x) \\ &= \prod_{a=1}^m P(\epsilon_a \leq x). \end{aligned}$$

The function $P(\epsilon_a \leq x)$ is the cumulative distribution function (CDF) of ϵ_a , which here is simply defined as

$$P(\epsilon_a \leq x) = \begin{cases} 0 & \text{if } x \leq -1 \\ \frac{1+x}{2} & \text{if } x \in (-1, 1) \\ 1 & \text{if } x \geq 1 \end{cases}$$

This implies that

$$\begin{aligned} P(\max_a \epsilon_a \leq x) &= \prod_{a=1}^m P(\epsilon_a \leq x) \\ &= \begin{cases} 0 & \text{if } x \leq -1 \\ \left(\frac{1+x}{2}\right)^m & \text{if } x \in (-1, 1) \\ 1 & \text{if } x \geq 1 \end{cases} \end{aligned}$$

This gives us the CDF of the random variable $\max_a \epsilon_a$. Its expectation can be written as an integral

$$\mathbb{E} \left[\max_a \epsilon_a \right] = \int_{-1}^1 x f_{\max}(x) dx,$$

where f_{\max} is the probability density function of this variable, defined as the derivative of the CDF: $f_{\max}(x) = \frac{d}{dx} P(\max_a \epsilon_a \leq x)$, so that for $x \in [-1, 1]$ we have $f_{\max}(x) = \frac{m}{2} \left(\frac{1+x}{2}\right)^{m-1}$. Evaluating the integral yields

$$\begin{aligned} \mathbb{E} \left[\max_a \epsilon_a \right] &= \int_{-1}^1 x f_{\max}(x) dx \\ &= \left[\left(\frac{x+1}{2}\right)^m \frac{mx-1}{m+1} \right]_{-1}^1 \\ &= \frac{m-1}{m+1}. \end{aligned} \quad \square$$

EXPERIMENTAL DETAILS FOR THE ATARI 2600 DOMAIN

We selected the 49 games to match the list used by Mnih et al. (2015), see Tables below for the full list. Each agent step is composed of four frames (the last selected action is repeated during these frames) and reward values (obtained from the Arcade Learning Environment (Bellemare et al., 2013)) are clipped between -1 and 1.

Network Architecture. The convolution network used in the experiment is exactly the one proposed by proposed by Mnih et al. (2015), we only provide details here for completeness. The input to the network is a 84x84x4 tensor containing a rescaled, and gray-scale, version of the last four frames. The first convolution layer convolves the input with 32 filters of size 8 (stride 4), the second layer has 64 layers of size 4 (stride 2), the final convolution layer has 64 filters of size 3 (stride 1). This is followed by a fully-connected hidden layer of 512 units. All these layers are separated by Rectifier Linear Units (ReLU). Finally, a fully-connected linear layer projects to the output of the network, i.e., the Q-values. The optimization employed to train the network is RMSProp (with momentum parameter 0.95).

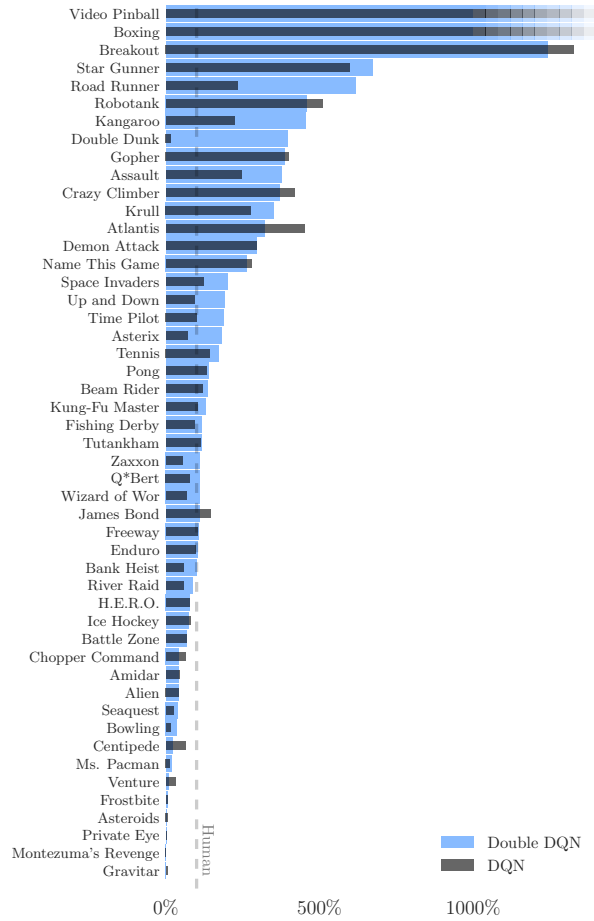


FIGURE 4. The normalized scores of DQN (narrow black bars) and Double DQN (broad blue bars) for the evaluation with no ops from Mnih et al. (2015).

Hyper-parameters. In all experiments, the discount was set to $\gamma = 0.99$, and the learning rate to $\alpha = 0.00025$. The number of steps between target network updates was $\tau = 10,000$. Training is done over 50M steps (i.e., 200M frames). The agent is evaluated every 1M steps, and the best policy across these evaluations is kept as the output of the learning process. The size of the experience replay memory is 1M tuples. The memory gets sampled to update the network every 4 steps with minibatches of size 32. The simple exploration policy used is an ϵ -greedy policy with the ϵ decreasing linearly from 1 to 0.1 over 1M steps.

SUPPLEMENTARY RESULTS IN THE ATARI 2600 DOMAIN

Figure 4 shows the normalized scores for DQN and Double DQN for the evaluation used by Mnih et al. (2015), and as given in Table 3 below. The Tables below provide further detailed results for our experiments in the Atari domain.

Game	Random	Human	DQN	Double DQN
Alien	227.80	6875.40	3069.33	2907.30
Amidar	5.80	1675.80	739.50	702.10
Assault	222.40	1496.40	3358.63	5022.90
Asterix	210.00	8503.30	6011.67	15150.00
Asteroids	719.10	13156.70	1629.33	930.60
Atlantis	12850.00	29028.10	85950.00	64758.00
Bank Heist	14.20	734.40	429.67	728.30
Battle Zone	2360.00	37800.00	26300.00	25730.00
Beam Rider	363.90	5774.70	6845.93	7654.00
Bowling	23.10	154.80	42.40	70.50
Boxing	0.10	4.30	71.83	81.70
Breakout	1.70	31.80	401.20	375.00
Centipede	2090.90	11963.20	8309.40	4139.40
Chopper Command	811.00	9881.80	6686.67	4653.00
Crazy Climber	10780.50	35410.50	114103.33	101874.00
Demon Attack	152.10	3401.30	9711.17	9711.90
Double Dunk	-18.60	-15.50	-18.07	-6.30
Enduro	0.00	309.60	301.77	319.50
Fishing Derby	-91.70	5.50	-0.80	20.30
Freeway	0.00	29.60	30.30	31.80
Frostbite	65.20	4334.70	328.33	241.50
Gopher	257.60	2321.00	8520.00	8215.40
Gravitar	173.00	2672.00	306.67	170.50
H.E.R.O.	1027.00	25762.50	19950.33	20357.00
Ice Hockey	-11.20	0.90	-1.60	-2.40
James Bond	29.00	406.70	576.67	438.00
Kangaroo	52.00	3035.00	6740.00	13651.00
Krull	1598.00	2394.60	3804.67	4396.70
Kung-Fu Master	258.50	22736.20	23270.00	29486.00
Montezuma's Revenge	0.00	4366.70	0.00	0.00
Ms. Pacman	307.30	15693.40	2311.00	3210.00
Name This Game	2292.30	4076.20	7256.67	6997.10
Pong	-20.70	9.30	18.90	21.00
Private Eye	24.90	69571.30	1787.57	670.10
Q*Bert	163.90	13455.00	10595.83	14875.00
River Raid	1338.50	13513.30	8315.67	12015.30
Road Runner	11.50	7845.00	18256.67	48377.00
Robotank	2.20	11.90	51.57	46.70
Seaquest	68.40	20181.80	5286.00	7995.00
Space Invaders	148.00	1652.30	1975.50	3154.60
Star Gunner	664.00	10250.00	57996.67	65188.00
Tennis	-23.80	-8.90	-2.47	1.70
Time Pilot	3568.00	5925.00	5946.67	7964.00
Tutankham	11.40	167.60	186.70	190.60
Up and Down	533.40	9082.00	8456.33	16769.90
Venture	0.00	1187.50	380.00	93.00
Video Pinball	16256.90	17297.60	42684.07	70009.00
Wizard of Wor	563.50	4756.50	3393.33	5204.00
Zaxxon	32.50	9173.30	4976.67	10182.00

TABLE 2. Raw scores for the no-op evaluation condition (5 minutes emulator time). Random, Human and DQN are as given by Mnih et al. (2015).

Game	Normalized DQN	Normalized Double DQN
Alien	42.75 %	40.31 %
Amidar	43.93 %	41.69 %
Assault	246.17 %	376.81 %
Asterix	69.96 %	180.15 %
Asteroids	7.32 %	1.70 %
Atlantis	451.85 %	320.85 %
Bank Heist	57.69 %	99.15 %
Battle Zone	67.55 %	65.94 %
Beam Rider	119.80 %	134.73 %
Bowling	14.65 %	35.99 %
Boxing	1707.86 %	1942.86 %
Breakout	1327.24 %	1240.20 %
Centipede	62.99 %	20.75 %
Chopper Command	64.78 %	42.36 %
Crazy Climber	419.50 %	369.85 %
Demon Attack	294.20 %	294.22 %
Double Dunk	17.10 %	396.77 %
Enduro	97.47 %	103.20 %
Fishing Derby	93.52 %	115.23 %
Freeway	102.36 %	107.43 %
Frostbite	6.16 %	4.13 %
Gopher	400.43 %	385.66 %
Gravitar	5.35 %	-0.10 %
H.E.R.O.	76.50 %	78.15 %
Ice Hockey	79.34 %	72.73 %
James Bond	145.00 %	108.29 %
Kangaroo	224.20 %	455.88 %
Krull	277.01 %	351.33 %
Kung-Fu Master	102.37 %	130.03 %
Montezuma's Revenge	0.00 %	0.00 %
Ms. Pacman	13.02 %	18.87 %
Name This Game	278.29 %	263.74 %
Pong	132.00 %	139.00 %
Private Eye	2.53 %	0.93 %
Q*Bert	78.49 %	110.68 %
River Raid	57.31 %	87.70 %
Road Runner	232.91 %	617.42 %
Robotank	508.97 %	458.76 %
Seaquest	25.94 %	39.41 %
Space Invaders	121.49 %	199.87 %
Star Gunner	598.09 %	673.11 %
Tennis	143.15 %	171.14 %
Time Pilot	100.92 %	186.51 %
Tutankham	112.23 %	114.72 %
Up and Down	92.68 %	189.93 %
Venture	32.00 %	7.83 %
Video Pinball	2539.36 %	5164.99 %
Wizard of Wor	67.49 %	110.67 %
Zaxxon	54.09 %	111.04 %

TABLE 3. Normalized results for no-op evaluation condition (5 minutes emulator time).

Game	Human	DQN	Double DQN
Alien	6371.30	570.20	736.40
Amidar	1540.40	133.40	188.20
Assault	628.90	3332.30	2774.30
Asterix	7536.00	124.50	5285.00
Asteroids	36517.30	697.10	1274.50
Atlantis	26575.00	76108.00	260556.00
Bank Heist	644.50	176.30	469.80
Battle Zone	33030.00	17560.00	25240.00
Beam Rider	14961.00	8672.40	9107.90
Bowling	146.50	41.20	62.30
Boxing	9.60	25.80	68.50
Breakout	27.90	303.90	338.70
Centipede	10321.90	3773.10	5166.60
Chopper Command	8930.00	3046.00	2483.00
Crazy Climber	32667.00	50992.00	94315.00
Demon Attack	3442.80	12835.20	13943.50
Double Dunk	-14.40	-21.60	-6.40
Enduro	740.20	475.60	475.90
Fishing Derby	5.10	-2.30	-3.40
Freeway	25.60	25.80	26.50
Frostbite	4202.80	157.40	497.30
Gopher	2311.00	2731.80	8742.80
Gravitar	3116.00	216.50	170.00
H.E.R.O.	25839.40	12952.50	15341.40
Ice Hockey	0.50	-3.80	-3.10
James Bond	368.50	348.50	416.00
Kangaroo	2739.00	2696.00	8254.00
Krull	2109.10	3864.00	6130.40
Kung-Fu Master	20786.80	11875.00	22771.00
Montezuma's Revenge	4182.00	50.00	30.00
Ms. Pacman	15375.00	763.50	1401.80
Name This Game	6796.00	5439.90	7871.50
Pong	15.50	16.20	17.70
Private Eye	64169.10	298.20	3385.10
Q*Bert	12085.00	4589.80	10713.30
River Raid	14382.20	4065.30	6579.00
Road Runner	6878.00	9264.00	43884.00
Robotank	8.90	58.50	52.00
Seaquest	40425.80	2793.90	5863.00
Space Invaders	1464.90	1449.70	1495.70
Star Gunner	9528.00	34081.00	53052.00
Tennis	-6.70	-2.30	11.00
Time Pilot	5650.00	5640.00	5375.00
Tutankham	138.30	32.40	63.60
Up and Down	9896.10	3311.30	7988.90
Venture	1039.00	54.00	75.00
Video Pinball	15641.10	20228.10	214282.30
Wizard of Wor	4556.00	246.00	934.00
Zaxxon	8443.00	831.00	7874.00

TABLE 4. Raw scores for the human start condition (30 minutes emulator time). Human and DQN are as given by Nair et al. (2015).

Game	Normalized DQN	Normalized Double DQN
Alien	5.57%	8.28%
Amidar	8.31%	11.89%
Assault	765.04%	627.77%
Asterix	-1.17%	69.27%
Asteroids	-0.06%	1.55%
Atlantis	460.90%	1804.78%
Bank Heist	25.72%	72.28%
Battle Zone	49.56%	74.60%
Beam Rider	56.92%	59.90%
Bowling	14.67%	31.77%
Boxing	270.53%	720.00%
Breakout	1153.44%	1286.26%
Centipede	20.44%	37.37%
Chopper Command	27.53%	20.59%
Crazy Climber	183.73%	381.67%
Demon Attack	385.42%	419.10%
Double Dunk	-71.43%	290.48%
Enduro	64.25%	64.29%
Fishing Derby	92.36%	91.22%
Freeway	100.78%	103.52%
Frostbite	2.23%	10.44%
Gopher	120.49%	413.23%
Gravitar	1.48%	-0.10%
H.E.R.O.	48.06%	57.69%
Ice Hockey	63.25%	69.23%
James Bond	94.11%	113.99%
Kangaroo	98.40%	305.25%
Krull	443.36%	886.79%
Kung-Fu Master	56.59%	109.67%
Montezuma's Revenge	1.20%	0.72%
Ms. Pacman	3.03%	7.26%
Name This Game	69.89%	123.88%
Pong	101.93%	106.08%
Private Eye	0.43%	5.24%
Q*Bert	37.13%	88.49%
River Raid	20.91%	40.18%
Road Runner	134.75%	638.94%
Robotank	840.30%	743.28%
Seaquest	6.75%	14.36%
Space Invaders	98.85%	102.34%
Star Gunner	377.00%	591.02%
Tennis	125.73%	203.51%
Time Pilot	99.52%	86.79%
Tutankham	16.55%	41.13%
Up and Down	29.67%	79.63%
Venture	5.20%	7.22%
Video Pinball	-644.88%	32157.42%
Wizard of Wor	-7.95%	9.28%
Zaxxon	9.49%	93.23%

TABLE 5. Normalized scores for the human start condition (30 minutes emulator time).