# HW5: Machine Learning Paper Review

Submitted by: malikap@oregonstate.edu

## 1. Non-Technical Background

Human-level control through deep reinforcement learning is a very influential paper in the field of artificial intelligence, specifically in the area of deep learning and reinforcement learning. This paper was published in 2015 by researchers at DeepMind Technologies. The authors of this paper are:

1. Volodymyr Mnih (Lead Author)
2. Koray Kavukcuoglu
3. David Silver
4. Andrei A. Rusu
5. Joel Veness
6. Marc G. Bellemare
7. Alex Graves
8. Martin Riedmiller
9. Andreas K. Fidjeland
10. Georg Ostrovski
11. Stig Petersen
12. Charles Beattie
13. Amir Sadik
14. Ioannis Antonoglou
15. Helen King
16. Dharshan Kumaran
17. Daan Wierstra
18. Shane Legg
19. Demis Hassabis

### 1.1. Who

All the authors are affiliated with DeepMind Technologies, a British artificial intelligence subsidiary of Alphabet Inc. (the parent company of Google). Among the authors, Demis Hassabis and Shane Legg are the co-founders of DeepMind. They along with David Silver and Koray Kavukcuoglu are likely among the principal investigators. The paper does not explicitly mention the roles of each author as student or intern.

## 1.2. Where

This paper was published in the Nature journal on 25th February 2015, and has received a significant number of citations, indicating its significant impact in the field. As of now, it has been cited 27,026 (source: Google Scholar) times. There is not much information on the internet about this paper receiving any awards or nomination, it also has no significant media reports, yet it is a highly impactful paper in the field of Reinforcement Learning and Deep Learning. The trend in number of citations seems to be increasing.

## 1.3. Reproducibility

The authors did not create a new dataset for their research. Instead, they used Atari 2600 games from the Arcade Learning Environment (ALE). ALE provides an environment for AI research with hundreds of Atari 2600 game environments, which is a standard benchmark in reinforcement learning research. The Atari 2600 games used in ALE are freely available online for research purposes.

There is no train/dev/test split. The training involved the agent interacting with the environment (playing the game) and learning from this interaction. The evaluation is performed on the same environment but with the agent's learning mechanisms disabled, to assess how well the trained policy performs.

There is a lot of content available online related to this paper. Here are some YouTube links:
- https://www.youtube.com/watch?v=V1eYniJ0Rnk
- https://www.youtube.com/watch?v=A0OVmImyFEA

The original paper did not release its code. But the core concepts and results of the paper have been reproduced and validated by the AI community. The algorithm description, implementation, and pseudocode is provided in detail in the paper and can be reproduced easily.


# 2. Core

## 2.1. The Problem

This paper aimed to solve a specific problem in the domain of reinforcement learning and control systems. The authors presented an algorithm that could achieve human-level performance in a wide range of highly challenging tasks, specifically playing various Atari 2600 video games. The main challenge was to create a system that could learn directly from high-dimensional sensory inputs (raw pixel data from video games' screen) and make decisions to maximize cumulative rewards (get a high score) in these environments (games).
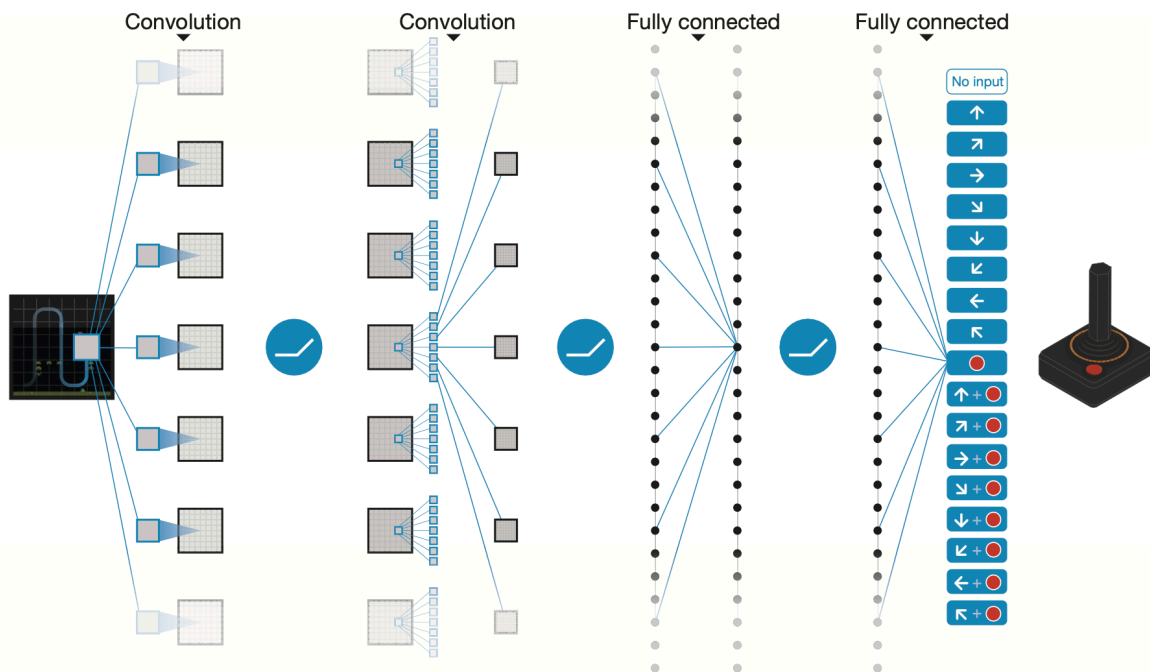
**Figure 1:** The Convolutional Neural network takes in the game's screen as input and outputs the highest expected reward action as dictated by the policy.

This paper features an artificial intelligence agent capable of learning optimal policies directly from high dimensional sensory inputs using end-to-end deep reinforcement learning.

The problem was not entirely new. Reinforcement learning (RL) had been explored in various forms before. Before this paper, some of the notable approaches in Reinforcement Learning relied on handcrafted features and were limited to environments with fully observed, low-dimensional state spaces. These methods were not able to generalize to more complex, high-dimensional environments, like video games or real-world scenarios, and were tailored to specific tasks. The authors wanted to build a system that could learn to excel in various tasks (in this case, different video games) without changing the architecture or learning algorithms for each new task.

The previous methods in reinforcement learning struggled with high-dimensional sensory inputs (like images and videos). The authors wanted to build a system which is also capable of learning successful policies directly from raw visual inputs as encountered in natural tasks, such as video frames from games. Instead of relying on handcrafted features or domain knowledge, they wanted to develop an approach that could learn directly from raw data. This meant creating an algorithm that could understand and interpret the input (game screens), execute actions, and learn from the outcomes of these actions.

With the deep learning revolution around 2012, there was increased interest in using deep neural networks as function approximators in Reinforcement Learning. Model-based deep RL algorithms were used to estimate a forward model of the environment dynamics using supervised learning. On the other hand, model-free algorithms were used to learn a policy directly without modeling the environment dynamics. These methods represented advances in handling high-dimensional data but still faced challenges in terms of generalization and computational efficiency.

## 2.2. Importance of the problem

This problem is very important for advancing AI knowledge. Demonstrating that an AI system can learn and adapt to a wide array of tasks without task-specific tuning is important for developing more general AI systems. Success in this area would lead to potential application far beyond video games, including robotics and automated systems in various industries.

The problem is inherently challenging because it involves learning directly from high-dimensional sensory inputs (raw pixel data), requiring the processing of complex visual environments. Each game presents its own unique rules and objectives, making the task of generalizing across different games very difficult.

**Complexity of High-Dimensional Sensory Inputs:** Learning directly from high-dimensional sensory inputs, such as raw pixel data from video games, is inherently challenging. Traditional AI and machine learning methods have struggled with this level of complexity without significant pre-processing or feature extraction.

**Generalization Across Diverse Tasks:** Creating an algorithm that could not only learn but also generalize its learning (without any modification to the algorithm) across a wide variety of tasks (in this case, different Atari 2600 games) was a very hard challenge. This is because each game has different rules, objectives, and visual representations, and therefore required a versatile and adaptive approach.

**Lack of Structured Training Data:** Unlike supervised learning, reinforcement learning does not rely on labeled training data. The agent must learn from interaction, which can be an inefficient and complex process, especially in environments where feedback (rewards) is sparse or delayed.

**Advancing the Field of AI:** Successfully tackling this problem would represent a major advancement in artificial intelligence, particularly in demonstrating the potential of deep learning when combined with reinforcement learning.

**Applications Beyond Gaming:** While the research used video games as a test environment, the underlying principles and techniques are applicable to a wide range of real-world problems, from robotics to automated systems in various industries.

**Understanding and Improving Learning Algorithms:** This research could provide deeper insights into how machines can learn and adapt in complex, dynamic environments, contributing to the broader understanding of both artificial and natural intelligence.

## 2.3. Authors' Approach

The authors tackled the problem using a new and innovative approach that combined deep learning with reinforcement learning. Their new method introduced the development of a Deep Q-Network (DQN).

**Deep Q-Network (DQN):** DQN is a type of algorithm used in reinforcement learning that combines Q-learning with deep neural networks. The goal of DQN is to learn the optimal action-value function (which predicts the expected reward of taking a certain action in a given state) using a neural network.

**Use of Convolutional Neural Networks (CNNs):** The authors made use of CNNs to process and learn from raw pixel data from the video games. CNNs are highly effective in interpreting visual input. This allowed the system to understand and analyze the game screens directly, without the need for feature extraction or manual input processing.

**Role of DNN in Decision Making:** After the CNN processes the input, the extracted features are fed into a Deep Neural Network (DNN) structure, which is responsible for approximating the Q-value function. This part of the network, often composed of fully connected layers, takes the high-level features identified by the CNN and uses them to estimate the value of taking different actions. The DNN ultimately outputs the Q-values, which are used to make decisions (selecting actions that maximize expected reward).

**Integration of Deep Learning with Q-Learning:** The authors used a deep neural network to approximate the Q-function. In reinforcement learning, the Q-function estimates the value of taking a certain action in a given state, considering the expected future rewards. Q-Learning is a form of reinforcement learning where an agent learns a policy, telling it what action to take under what circumstances. It does so by learning the Q-values (quality of actions) for each state-action pair.

$$Q^*(s, a) = \max_\pi \mathbb{E}\left[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots \mid s_t = s, a_t = a, \pi\right]$$

The optimal action-value function, $Q^*(s, a)$, returns the best expected future rewards an agent can achieve by taking action $a$ in state $s$, considering both immediate and future rewards. The immediate reward is $r_t$, and future rewards are factored in with diminishing importance through a discount factor $\gamma$. This factor ensures that rewards expected further in the future are weighted less. The equation seeks the best long-term strategy by maximizing these expected rewards over all possible actions and states, guiding the agent's decisions to achieve the most favorable outcomes in a learning process.

**Experience Replay:** To improve the learning process, the authors implemented an experience replay mechanism. In reinforcement learning, an AI agent learns from its experiences. But if it only learns from what's happening right now, it might miss out on some important lessons from the past. To solve this problem, and increase effectiveness, the agent's experiences at each time step (its observations, actions, rewards, and new observations) were stored in a replay memory. During the training phase, mini-batches of these experiences were randomly sampled from the replay memory. This approach breaks the correlation between consecutive samples and therefore stabilizes the learning algorithm.

**Iterative Policy Learning:** The DQN was trained to select actions by estimating the optimal action-value function. It used a variant of Q-learning, with the key difference being that the action values (Q-values) were approximated using the deep neural network. The policy (mapping from environment state to action) was continuously updated as the network learned from new experiences, improving the agent's performance over time.

**Preprocessing and Frame Skipping:** The authors applied preprocessing to simplify the input state representation. This included converting game frames to grayscale and down-sampling to reduce the computational load. They also implemented frame skipping (only selecting every k-th frame) to speed up the learning process and make it more manageable for the network.

**Reward Clipping:** To make the algorithm more robust across various games with different scoring systems, the authors clipped the rewards to a fixed range. This prevented the scale of the rewards from influencing the learning process.

## 2.4. Remarkable Results

The Deep Q-Learning approach which involved the development of the Deep Q-Network (DQN), resulted in impressive results as discussed below.

**Human-Level Performance:** One of the most important outcomes was that the DQN achieved or exceeded human-level performance in a majority of the Atari 2600 games tested. This is a highly praiseworthy and innovative achievement, as it demonstrated that an AI system could learn and adapt to a wide range of complex tasks to the extent of matching or surpassing human skills.

**Learning from Raw Visual Inputs:** The DQN was successful in learning directly from high-dimensional sensory inputs (raw pixel data from game screens) without the need for feature extraction or hand-engineered input representations. This is a highly innovative advancement, showing the potential of deep learning networks in interpreting and making decisions based on complex, unstructured data.

**Generalization Across Tasks:** The DQN was not specifically fine-tuned to any single game but was instead able to learn and adapt across various games with different rules, objectives, and visual styles. This generalization capability is a huge step forward in AI, and shows the potential for creating similar systems in real-world settings.

**Setting a New Standard in AI Research:** The success of the DQN in achieving human-level control in complex tasks set a new standard in the field of AI. It has opened up new avenues for research, particularly in the application of deep learning to decision-making and control problems.

**Influencing Subsequent Research and Applications:** The principles and techniques demonstrated in this research have influenced subsequent work in AI and machine learning. The use of deep learning in combination with reinforcement learning has been explored further for various applications, ranging from robotics to complex system optimization.

# 3. Further

## 3.1. Some of the potential flaws in the authors' method and results

**Computational Resources:** The Deep Q-Network (DQN) required a large amount of computational resources for training, including high-end GPUs and long processing times. This level of resource requirement could be a limitation for broader applications or for researchers with limited access to such computational power.

**Sample Efficiency:** DQN's training process was not very sample efficient. It required a large number of interactions with the environment (gameplay in the case of Atari games) to achieve good performance. This inefficiency can be problematic in real-world scenarios where obtaining such vast amounts of experience is impractical or expensive.

**Generalization Beyond Training Environment:** Although the DQN showed impressive performance and generalization across different games, there's still a question about how well this generalization would extend to tasks significantly different from the training environment, especially real-world tasks that are more complex and less structured than Atari games. The real world environments are often non-deterministic and do not give immediate rewards after each action.

## 3.2. Follow up idea on this line of work

A good follow up to this line of research would be to explore the applicability of the described methods in real-world scenarios. Unlike video game environments, real-world scenarios are often non-deterministic and do not give immediate rewards after each action. For example, In self-driving cars, decisions must be made in real-time with uncertain and dynamic conditions,

like changing traffic or weather. The feedback (safe arrival at a destination) is delayed, and the reward structure (safety, efficiency) is highly complex. Exploring strategies for environments where feedback (rewards) is sparse or delayed could be a great follow up to improve the applicability of this research in real-world scenarios.

### 3.3. Overall Feeling

My overall impression of this paper is that it represents a important milestone in the field of artificial intelligence. The integration of deep learning with reinforcement learning to create a AI system that could learn and excel in a wide array of tasks (achieving human level performance in Atari 2600 video games) is a highly innovative achievement that is significant step forward in the development of more general and adaptable AI systems.

The single take-home message from this paper is the great potential of deep learning techniques when applied with reinforcement learning. By using a Convolutional Neural Network to interpret high-dimensional sensory input and then using a Deep Neural Network to make decisions in complex environments, the Deep Q-Network (DQN) demonstrated that it's possible for an AI system to learn, excel, and adapt to a wide range of video game environments in an end-to-end manner, directly from raw image inputs. This approach not only achieved remarkable results in video games environments but also set a ground for future research for applicability of these techniques in various real-world scenarios.

# 4. Relevance

This paper covered a wide variety of different topic from Deep Learning, CNNs, and Reinforcement Learning. Although some of these topics were not part of the course, the students did have the foundational knowledge (particularly in the area of Deep Learning and Neural Networks) to form a solid base from which they can grasp the fundamental aspects of this paper.

I have already reproduced the results of this paper for my Capstone project for Udacity's Machine Learning Engineer Nanodegree (GitHub: https://github.com/1998apoorvmalik/udacity-machine-learning-engineer-nanodegree). So, I just briefly went through some of my course modules, my Capstone project report, and read the paper again.

# 5. References

1. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., & Hassabis, D. (2015, February 25). Human-level control through deep reinforcement learning. Nature. Retrieved December 3, 2022, from https://www.nature.com/articles/nature14236