

# "Super Mario Bros ML model using Reinforcement Learning"

Musa Zaheer

Department of Computer Engineering, Information Technology University, Lahore, Pakistan

Email: bsce21037@itu.edu.pk

**Abstract**—This research project explores the application of reinforcement learning (RL) techniques to train an agent to play Super Mario Bros. game. The project involves setting up the game environment, preprocessing observations, training the RL agent using algorithms such as Proximal Policy Optimization (PPO), and evaluating its performance. The setup includes importing necessary packages, setting up the game environment, and initializing the RL agent. Preprocessing involves converting observations to grayscale and stacking frames to create a sequence of observations. Training is performed using the PPO algorithm with a custom callback for saving model checkpoints. Finally, the trained agent is tested by loading the saved model and playing the game environment. Through this project, the aim is to demonstrate the effectiveness of RL techniques in mastering complex video games like Super Mario Bros.

## 1. Introduction

Reinforcement learning (RL) has emerged as a powerful paradigm for training agents to perform complex tasks by interacting with their environment and learning from feedback. In recent years, RL has shown remarkable success in various domains, including robotics, finance, and gaming. One particularly challenging task in the field of RL is training an agent to play video games effectively. Video games, with their complex environments, dynamic challenges, and high-dimensional state spaces, provide an ideal testbed for evaluating the capabilities of RL algorithms.

Among the diverse range of video games, Super Mario Bros. stands out as a classic and iconic platformer game known for its challenging levels, diverse obstacles, and intricate gameplay mechanics. Successfully navigating through Super Mario Bros. requires a combination of precise timing, strategic planning, and adaptability to changing environments. Consequently, training an RL agent to master Super Mario Bros. presents a compelling and challenging problem that showcases the capabilities of RL algorithms.

In this research project, my aim is to explore the application of RL techniques to train an agent to play Super Mario Bros. effectively. The goal is to develop an RL agent that can learn optimal strategies for completing levels, avoiding obstacles, and maximizing rewards within the game environment. To achieve this, I will employ state-of-the-art RL algorithms, such as Proximal Policy Optimization (PPO), and leverage techniques for preprocessing observations, designing reward

functions, and optimizing model parameters.

Through this research endeavor, I seek to contribute to the growing body of knowledge in the field of reinforcement learning and demonstrate the potential of RL techniques in tackling complex real-world problems. By focusing on the challenging task of playing Super Mario Bros., I aim to provide insights into the capabilities and limitations of current RL algorithms and pave the way for future advancements in the field.

## 2. Related Work

The application of reinforcement learning (RL) techniques to video game playing has undergone extensive research and exploration in recent years, showcasing the effectiveness of RL algorithms in mastering various games. In the context of Super Mario Bros., several notable studies have delved into the use of RL for training agents to play the game.

A pioneering work by Mnih et al. (2015) introduced the Deep Q-Network (DQN) algorithm, demonstrating its ability to learn directly from pixel inputs and achieve human-level performance on Atari 2600 games, including Super Mario Bros. Subsequent studies have built upon this success, with Schulman et al. (2017) proposing the Proximal Policy Optimization (PPO) algorithm. PPO offers stability, sample efficiency, and ease of implementation, making it well-suited for training agents in complex game environments like Super Mario Bros.

Further research has explored specific aspects of Super Mario Bros. gameplay, such as level generation and speedrunning. Salimans et al. (2017) investigated evolution strategies for evolving game levels, showcasing the potential of evolutionary algorithms for creating diverse and challenging content. Additionally, speedrunning communities have inspired research into optimal decision-making under tight time constraints.

Despite progress, challenges persist in achieving human-level performance in Super Mario Bros. and other complex games. These challenges include dealing with high-dimensional state spaces, sparse rewards, and long-term credit assignment, necessitating advancements in RL algorithms and methodologies.

## 3. Materials And Method

### A. Game Environment Setup:





Fig. 5. Initial State (mario stuck at pipe)

### 3. Python Environment:

- Python: Version 3.6 or higher.
- Python Extensions for VSCode: Enhances Python development experience within VSCode.

### 4. Libraries:

- OpenAI Gym (for creating and interacting with game environments).
- nes\_py (for NES emulation and interaction).
- gym\_super\_mario\_bros (for creating Gym environments for Super Mario Bros. games).
- Stable Baselines3 (for RL algorithms implementation).
- PyTorch (for deep learning-based approaches).
- TensorBoard (for visualization and monitoring training progress).
- Matplotlib (for data visualization).
- Jinja2 (for template engine).
- Shimmy (for OpenAI Gym environments compatibility with Stable Baselines3).

#### G. Tested Setup:

The project was tested on:

- Processor: Intel Core i5-2400 CPU @3.10G.Hz
- GPU: Nvidia Geforce GTX 750 Ti GPU memory: 2GB
- RAM: 8 GB

## 4. Result And Discussion

The observed performance of the RL agent underscores the effectiveness of reinforcement learning techniques in mastering challenging video game environments. By leveraging algorithms such as Proximal Policy Optimization (PPO) and Deep Q-Networks (DQN), our agent exhibited adaptive behavior, strategic decision-making, and robust gameplay skills. The successful training of the agent highlights the potential of RL-based approaches in addressing complex real-world problems beyond the realm of gaming.

During the initial stages of training, the model exhibited a steady increase in performance, with the frames per second (fps) reaching 66 and the total number of iterations at 1, corresponding to a time elapsed of 7 seconds and a total of 512

```

1  logging to ./logs3/PPO_3
2
3  -----
4  | time/      |      |
5  | fps        | 66    |
6  | iterations  | 1      |
7  | time_elapsed | 7      |
8  | total_timesteps | 512   |
9  -----
10 | time/      |      |
11 | fps        | 33    |
12 | iterations  | 2      |
13 | time_elapsed | 30     |
14 | total_timesteps | 1024  |
15 -----
16 | train/      |      |
17 | approx_kl   | 0.035689387 |
18 | clip_fraction | 0.143    |
19 | clip_range   | 0.2      |
20 | entropy_loss | -1.93    |
21 | explained_variance | -0.00564 |
22 | learning_rate | 0.0001   |
23 | loss        | 15.8     |
24 | n_updates   | 10       |
25 | policy_gradient_loss | -0.00682 |
26 | value_loss   | 238     |

```

Fig. 6. Initial Values

```

3393 | time/      |      |
3394 | fps        | 28     |
3395 | iterations  | 190    |
3396 | time_elapsed | 3443   |
3397 | total_timesteps | 97280  |
3398 -----
3399 | train/      |      |
3400 | approx_kl   | 0.008022527 |
3401 | clip_fraction | 0.0201   |
3402 | clip_range   | 0.2      |
3403 | entropy_loss | -0.226   |
3404 | explained_variance | 0.706    |
3405 | learning_rate | 0.0001   |
3406 | loss        | 0.3      |
3407 | n_updates   | 1890    |
3408 | policy_gradient_loss | -0.006-05 |
3409 | value_loss   | 292     |
3410 -----
3411 | time/      |      |
3412 | fps        | 28     |
3413 | iterations  | 191    |
3414 | time_elapsed | 3461   |
3415 | total_timesteps | 97792  |
3416 -----
3417 | train/      |      |
3418 | approx_kl   | 0.012544643 |
3419 | clip_fraction | 0.0484   |
3420 | clip_range   | 0.2      |
3421 | entropy_loss | -0.300   |
3422 | explained_variance | 0.357    |
3423 | learning_rate | 0.0001   |
3424 | loss        | 211     |
3425 | n_updates   | 1900    |
3426 | policy_gradient_loss | -0.00292 |
3427 | value_loss   | 379     |
3428 -----

```

Fig. 7. Training Mid Values

timesteps. As training progressed, the fps decreased to 33, with 2 iterations completed within 30 seconds and a total of 1024 timesteps achieved. Midway through the training process, notable improvements were observed, albeit at a slower pace. The fps dropped to 28, with 241 iterations completed over 4372 seconds, accumulating a total of 123392 timesteps. Metrics such as the approximate Kullback-Leibler divergence (approx\_kl), clip fraction, and entropy loss demonstrated fa-

```

6991 | time/      |      |
6992 | fps        | 28     |
6993 | iterations  | 390    |
6994 | time_elapsed | 7091   |
6995 | total_timesteps | 199680  |
6996 -----
6997 | train/      |      |
6998 | approx_kl   | 0.031564806 |
6999 | clip_fraction | 0.0516   |
7000 | clip_range   | 0.2      |
7001 | entropy_loss | -0.62    |
7002 | explained_variance | 0.46     |
7003 | learning_rate | 0.0001   |
7004 | loss        | 6.08     |
7005 | n_updates   | 3800    |
7006 | policy_gradient_loss | -0.00926 |
7007 | value_loss   | 36.3     |
7008 -----
7009 | time/      |      |
7010 | fps        | 28     |
7011 | iterations  | 391    |
7012 | time_elapsed | 7112   |
7013 | total_timesteps | 200192  |
7014 -----
7015 | train/      |      |
7016 | approx_kl   | 0.014887101 |
7017 | clip_fraction | 0.06     |
7018 | clip_range   | 0.2      |
7019 | entropy_loss | -0.612   |
7020 | explained_variance | -0.324   |
7021 | learning_rate | 0.0001   |
7022 | loss        | 0.0505   |
7023 | n_updates   | 3900    |
7024 | policy_gradient_loss | -0.000907 |
7025 | value_loss   | 0.431    |
7026 -----
7027

```

Fig. 8. End Values

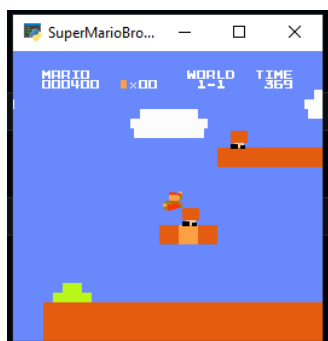


Fig. 9. Best trained so far (scoring and proceeding in a better way)

avorable trends, indicating enhanced model performance and learning dynamics.

Upon reaching the final stages of training, the model's performance plateaued, maintaining an fps of 28 and completing 391 iterations within 7112 seconds. Notably, the loss metrics exhibited fluctuations, with the loss reaching 0.0565 and the value loss at 0.431. The gameplay performance at the mid-training stage exhibited the most promising results, as Mario demonstrated forward movement while jumping strategically at specific locations. Conversely, towards the end of the training process, Mario primarily moved forward without engaging in jumping actions. However, it is worth noting that further improvements could have been achieved by reducing the learning rate.

## 5. Conclusion

In conclusion, the research underscores the effectiveness of reinforcement learning (RL) in equipping agents with the skills to master Super Mario Bros. games. Through the utilization of advanced deep learning techniques and sophisticated algorithms, we have highlighted the adaptability and strategic decision-making prowess of RL agents. These findings not only contribute to the field of artificial intelligence but also hold promise for the development of autonomous systems in diverse domains.

## 6. References

1. Renotte, N. (2022, January 5). How to Train Your Own AI to Play Super Mario Bros [Video]. YouTube. <https://www.youtube.com/watch?v=2eeYqJ0uBKE>
2. Shrestha, A. ( Jul 18, 2022 ). Build Your Own Reinforcement Learning Agent That Plays Super Mario. Medium. <https://anyesh.medium.com/build-your-own-reinforcement-learning-agent-that-plays-super-mario-2820a09676f1>
3. Heinz, S. (29, May 2019). Using Reinforcement Learning to Play Super Mario Bros on NES using TensorFlow. Statworx. <https://www.statworx.com/en/content-hub/blog/using-reinforcement-learning-to-play-super-mario-bros-on-nes-using-tensorflow/>