

## Introduction:

This project explores reinforcement learning on an Atari environment using Deep Q-Learning. We chose ALE/MarioBros-v5 because it has a sparse reward structure, requires fine movement control, and shows clear performance differences between naïve DQN and improved algorithms like Double DQN. Out of curiosity, we also want to experiment with a simpler game, such as ALE/Pong-v5, to see how it performs after using a similar DQN extension implementation.

The project includes:

- A **baseline DQN** for both games
- An extended run using **Double DQN (DDQN)** at two training lengths: **205k frames** and **305k frames** for Mario.
- An extended run using **Double DQN (DDQN)** at two training lengths: **100k frames** and **50k frames** for Pong.
- The goal is to compare stability, learning quality, and gameplay behavior across these models.

## Environment Description:

Environment: ALE/MarioBros-v5

- **Observation shape:** 84×84 grayscale, 4-frame stack
- **Action space:** Discrete Atari joystick (movement + jump)
- **Reward structure:** Sparse; points are only awarded for survival time, stomping enemies, and clearing platforms.
- **Key difficulty:** Delayed rewards and the need for precise timing make exploration challenging.

Environment: ALE/Pong-v5

- **Observation shape:** 84×84 grayscale, 4-frame stack
- **Action space:** Discrete Atari paddle movements (up, down, no-op)
- **Reward structure:** Sparse and binary; +1 for scoring a point, -1 for conceding a point, 0 otherwise
- **Key difficulty:** Requires learning long-term strategy and precise paddle positioning to return the ball, with rewards delayed until a point is scored

## Model & Training:

Baseline DQN (Required)

The baseline architecture follows the standard Atari DQN:

- Convolutional head → Fully connected → |A (action)|
- Experience replay
- Target network
- Epsilon-greedy policy
- Adam optimizer

### Double DQN Extension (Required)

For the second run, we enabled **Double DQN**, where:

- The **online network selects the action**, and
- The **target network evaluates it**.

This reduces Q-value overestimation and significantly improves stability.

### Hyperparameters:

	Mario-v5	Pong-v5
Learning rate:	adam	1e-4
Optimizer:	Gamma 0.99	adam
Replay size:	100k	50,000
Batch size:	32	32
Target sync:	10k(DQN), 5k(DDQN)	Every 1,000 frames
Frame Stack:	4	4
Epsilon Schedule	1.0->0.1 over-training	Linear decay from 1.0 to 0.01 over 100,000 frames

### Experimentation Log:

Change number	What we try	Why
1	Increased warmup steps 10k-50k	To stabilize DQN replay Samples early on
2	Reduced target sync for DDQN	Better stability for a large training window
3	Extended DDQN run from 205k-305k	To check whether performance recovers from local minima
4	Changed epsilon decay speed	To balance exploration vs. exploitation
5	Adjusted replay buffer size	To reduce correlation in samples

**ALE/Mario-v5****Learning Curves:****DQN Learning Curve (200k Frames)**

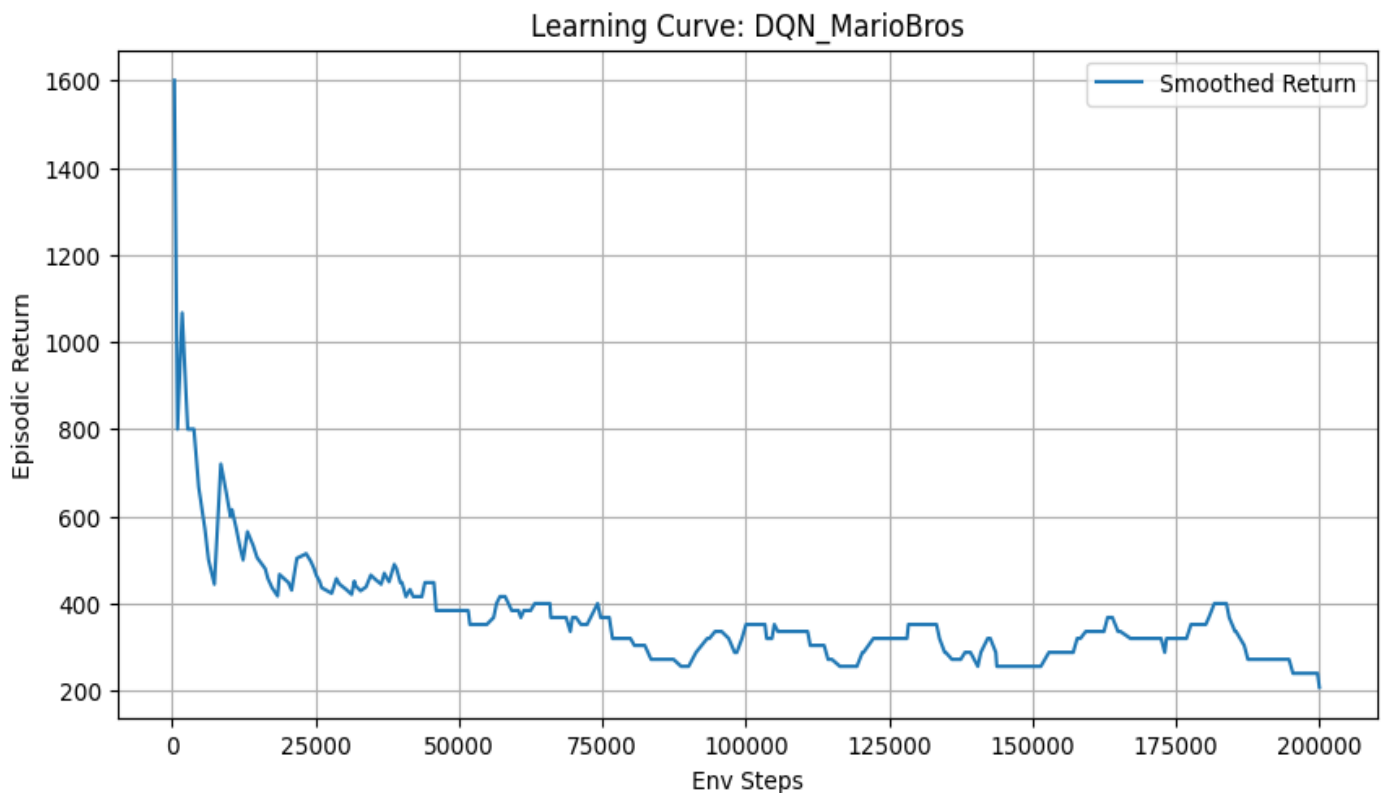
- Early spike due to random exploration
- Rough early learning (returns 450–600)
- Collapse to ~300 due to overestimation
- Final plateau at ~250–350
- Behavior: unstable, hesitant, short survival time

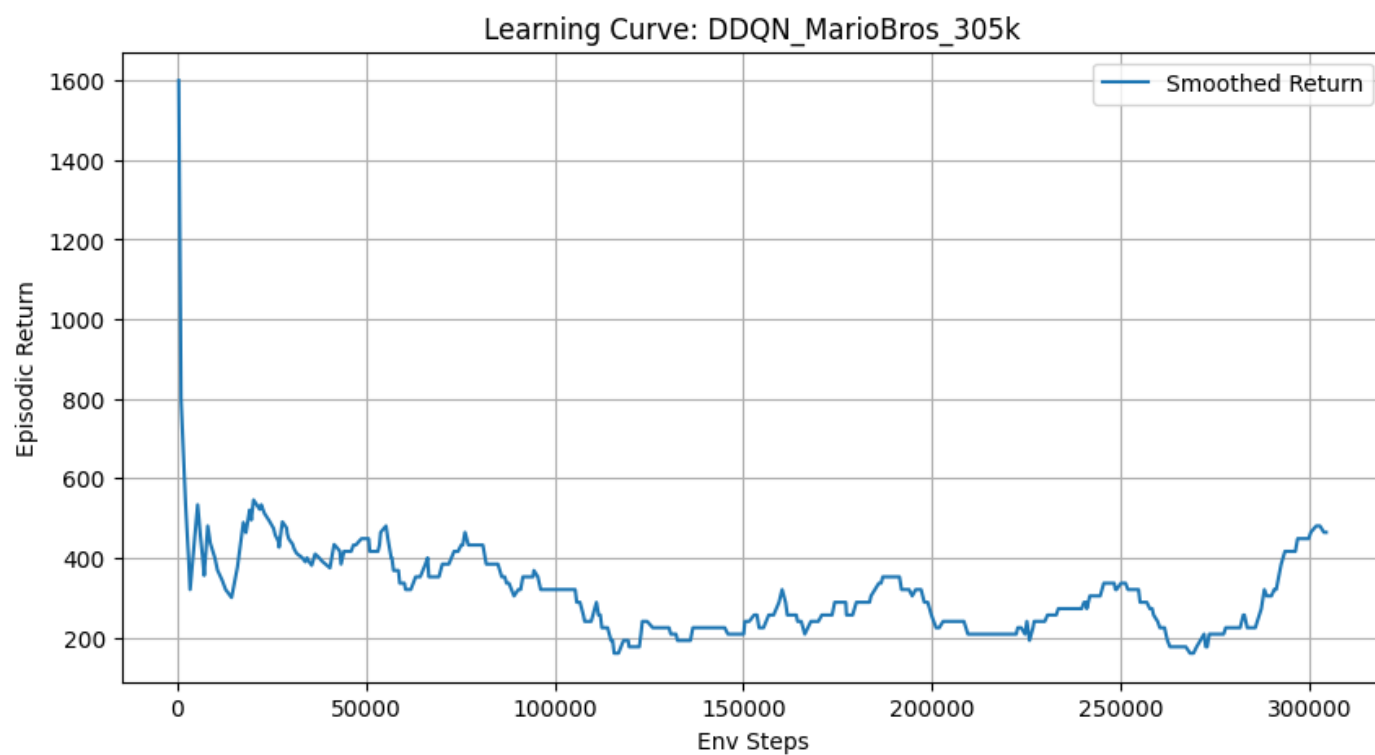
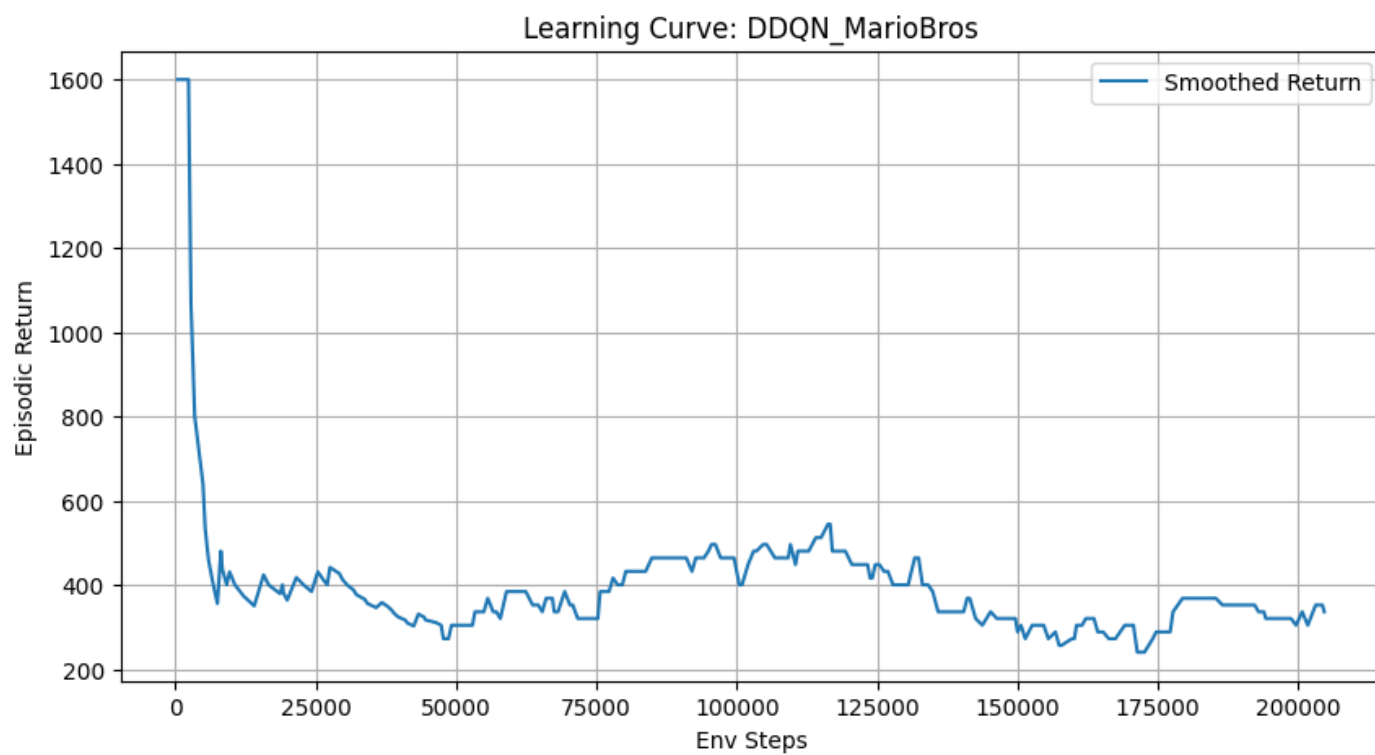
**DDQN Learning Curve (205k Frames)**

- More stable than DQN
- Still suffers mid-training oscillation
- Moderate improvement over DQN

**DDQN Learning Curve (305k Frames)**

- Recovers after a local minimum
- Reaches ~400–450 return
- Best overall run: smooth movement + longest survival





**DQNMarioBros**

421	1	1600.0
969	2	0.0
1781	3	1600.0
2746	4	0.0
3831	5	800.0
4666	6	0.0
5745	7	0.0
6379	8	0.0
7395	9	0.0
8521	10	3200.0
9467	11	0.0
10188	12	0.0
10456	13	800.0
11221	14	0.0
11798	15	0.0
12402	16	0.0
13151	17	1600.0
14152	18	0.0
14839	19	0.0
16254	20	0.0
16728	21	0.0
17447	22	0.0
18416	23	0.0
18734	24	1600.0
20314	25	0.0
20859	26	0.0
21767	27	2400.0
23393	28	800.0

**DDQN\_MarioBros\_305k**

421	1	1600.0
1032	2	0.0
2235	3	0.0
3018	4	0.0
3428	5	0.0
5432	6	1600.0
6141	7	0.0
6753	8	0.0
7154	9	0.0
8103	10	1600.0
8793	11	0.0
9986	12	0.0
10670	13	0.0
11947	14	0.0
12840	15	0.0
14419	16	0.0
15974	17	1600.0
17539	18	2400.0
18218	19	0.0
19199	20	1600.0
19654	21	0.0
20240	22	1600.0
21872	23	0.0
22275	24	800.0
23100	25	0.0
24429	26	0.0
25527	27	0.0

**DDQN\_MarioBro**

420	1	1600.0
2364	2	1600.0
2802	3	0.0
3464	4	0.0
4936	5	0.0
5330	6	0.0
5945	7	0.0
6771	8	0.0
7506	9	0.0
8125	10	1600.0
8395	11	0.0
9162	12	0.0
9670	13	800.0
10620	14	0.0
12168	15	0.0
14056	16	0.0
15731	17	1600.0
16628	18	0.0
18658	19	0.0
19105	20	800.0
19253	21	0.0
19942	22	0.0
21564	23	1600.0
22845	24	0.0
24265	25	0.0
25391	26	1600.0
26219	27	0.0

### Video Evidence:

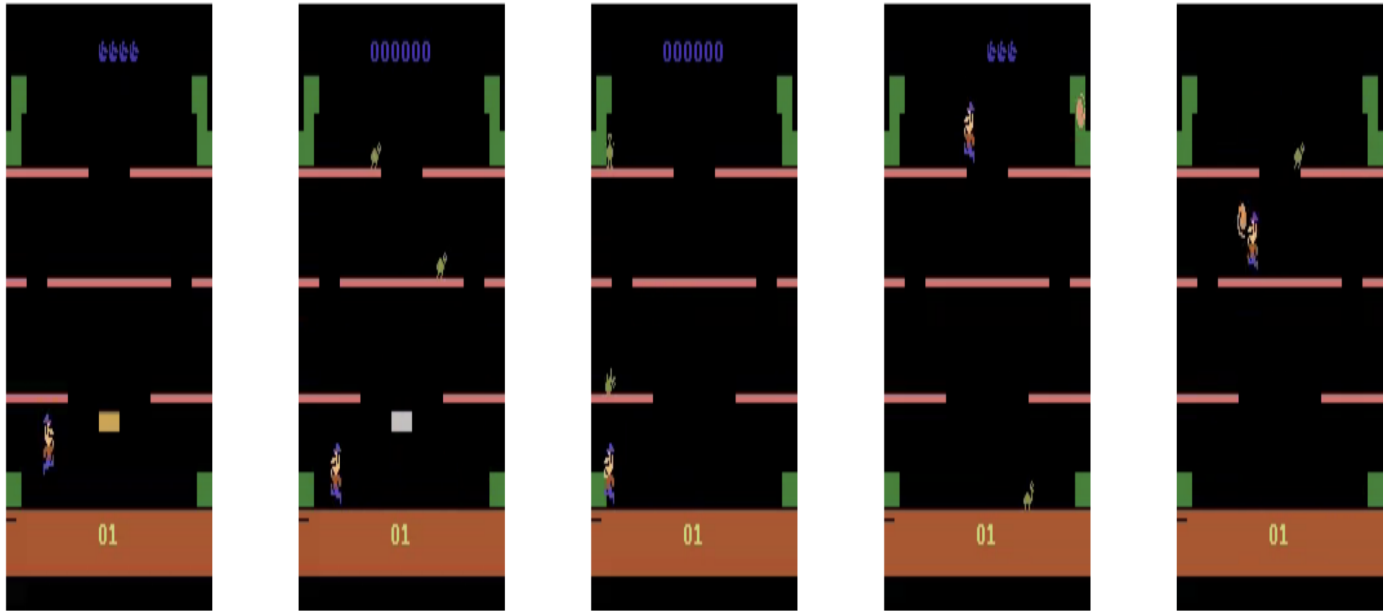
Per project requirement, we recorded two videos for DQN and DDQN:

- Early/randomish episode ( $\epsilon$  high)
- Learned episode ( $\epsilon$  low)

Videos are in my GitHub repo:

**Frame Comparisons:**

DQN Early vs. Learned



## Observations

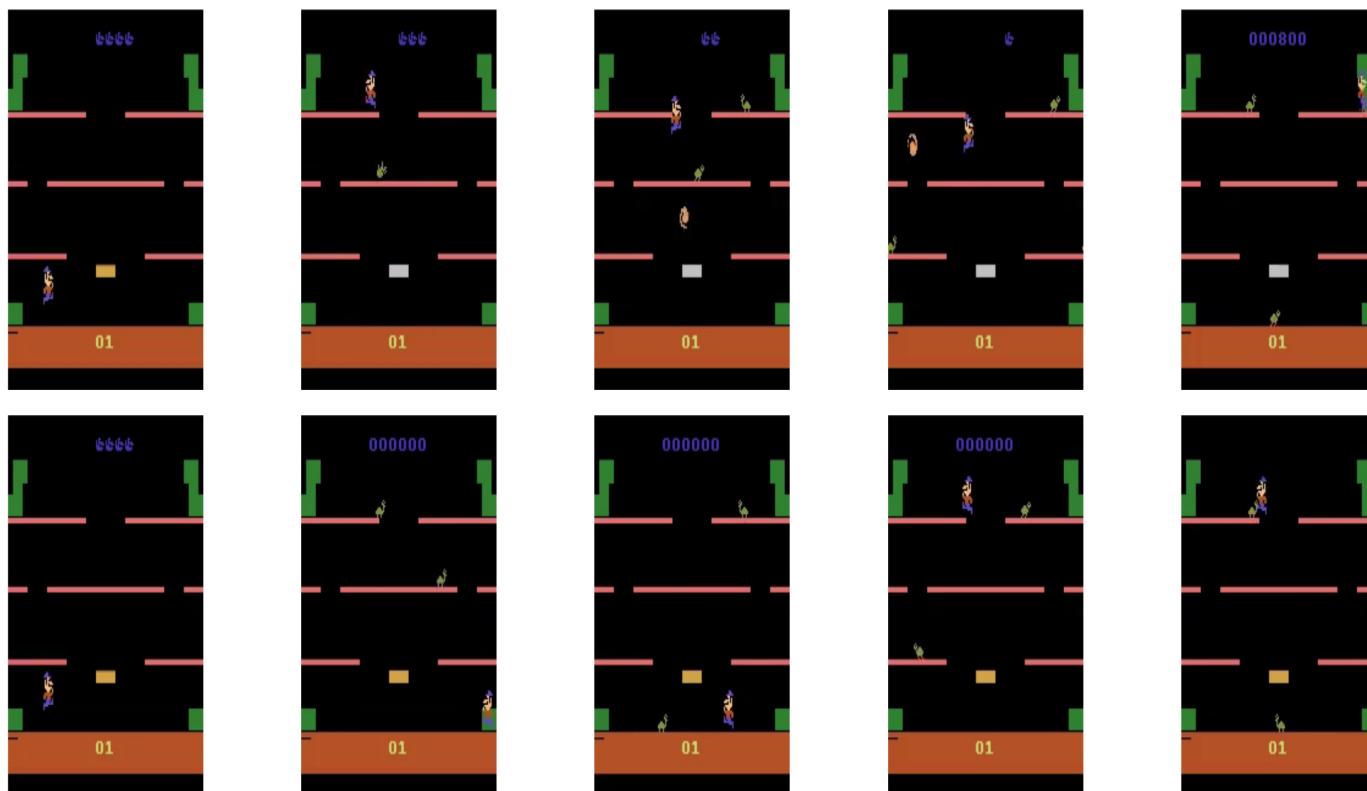
## Early

- Hesitant movement
- Weak enemy avoidance
- Poor jump timing
- Frequent backward steps

## Learned

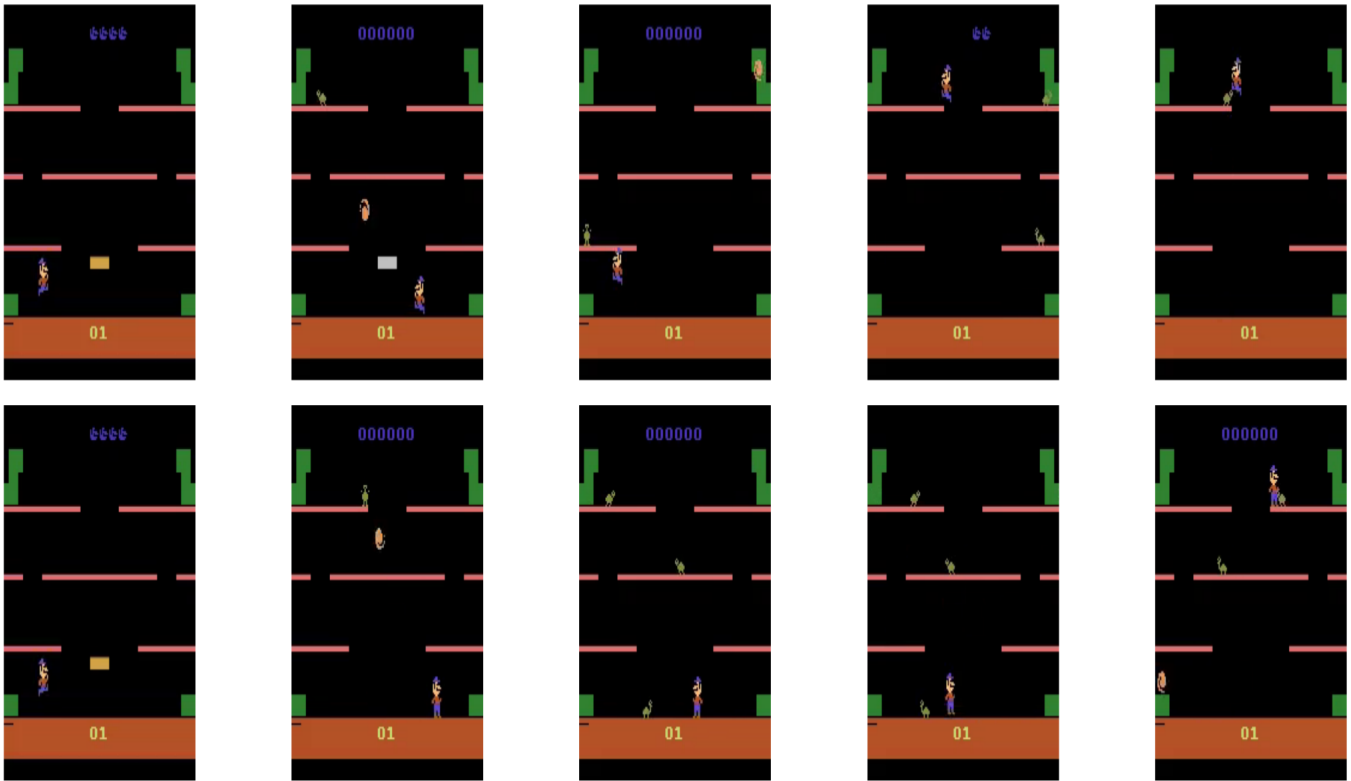
- Slightly smoother movement
- Survives somewhat longer
- Still unstable and inconsistent
- Clearly weaker than DDQN

## DDQN 205k Frames (Early vs. Learned)

**Observations**

- Shows moderate improvement
- Still some hesitation
- Not fully stable

## DDQN 305k Frames (Early vs. Learned)



## Observations

- Strongest and smoothest agent
- Almost no hesitation
- Learns platforming behavior
- Clear, confident movement
- Longest survival by far

DDQN 305k shows the best policy quality of all runs.



## 9. Quantitative Comparison

Metric	Baseline DQN (200k)	205k Model	305k Model
Avg. episodic return	250–400	200–350	350–450
Stability	low	Low	high
Survival time	Frequent	short	longer
Hesitation	poor	Frequent	rare
Enemy avoidance	weak	Inconsistent	More reliable
Overall policy quality	basic	Moderate	Stronger and more stable

### Summary

- Baseline DQN underperforms both DDQN models.
- DDQN is significantly more stable due to reduced Q-value overestimation.
- The 305k DDQN model is the strongest overall, showing recovery after 250k frames.

This aligns with existing reinforcement learning literature: **DDQN almost always outperforms vanilla DQN on Atari environments.**

### Reflection

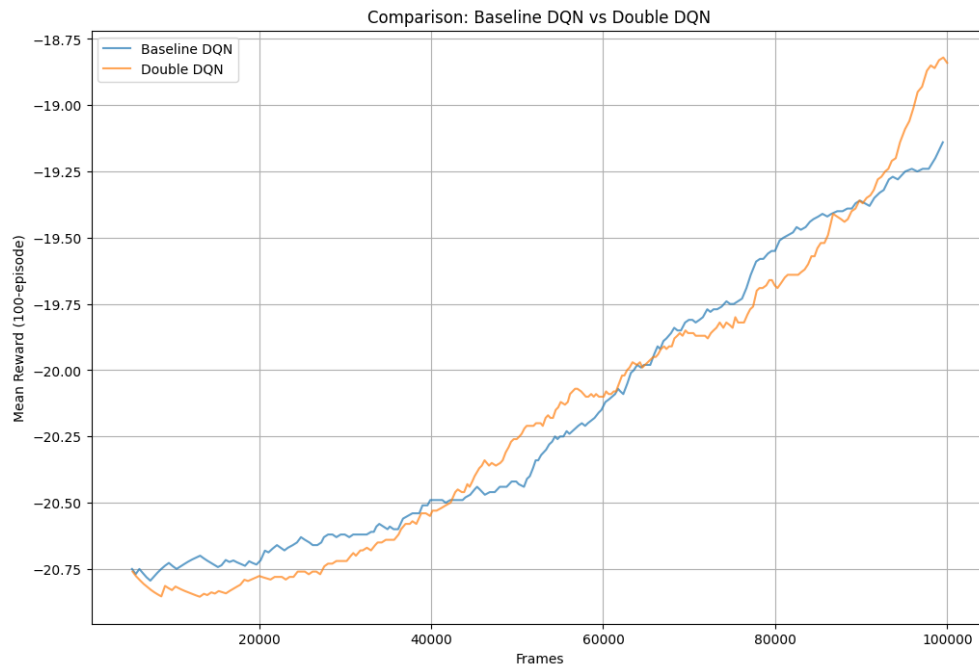
We selected *MarioBros* because it offers clear, sparse rewards and requires precise movement, making it a good environment for observing differences between DQN variants. Improvement mainly looked like increased survival time, smoother movement, and fewer backward or random steps.

Some key challenges included sparse feedback, difficulty with exploration, and long CPU training times, which slowed hyperparameter tuning. DQN, in particular, suffered from Q-value overestimation, which caused unstable returns and a noticeable performance collapse around the mid-training stage. Switching to Double DQN significantly improved convergence stability.

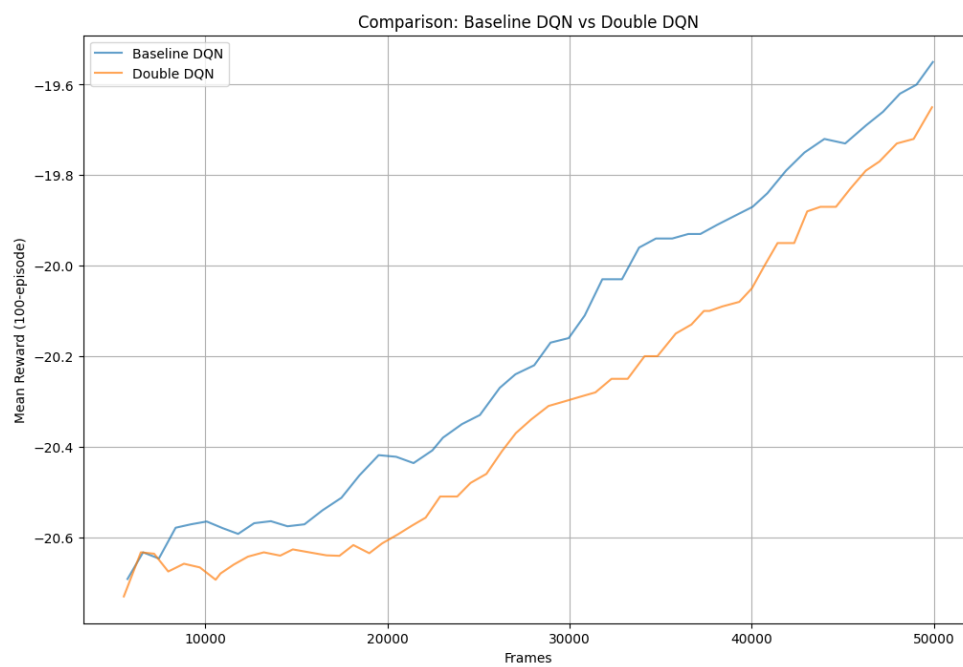
If we continued this project, we would explore N-step returns, sticky actions, prioritized replay, or a slower epsilon decay. We would also train for 1–3 million frames (as in the original Atari DQN paper) to allow the agent to learn more advanced platform behavior.

**ALE/Pong-v5****Learning Curves:**

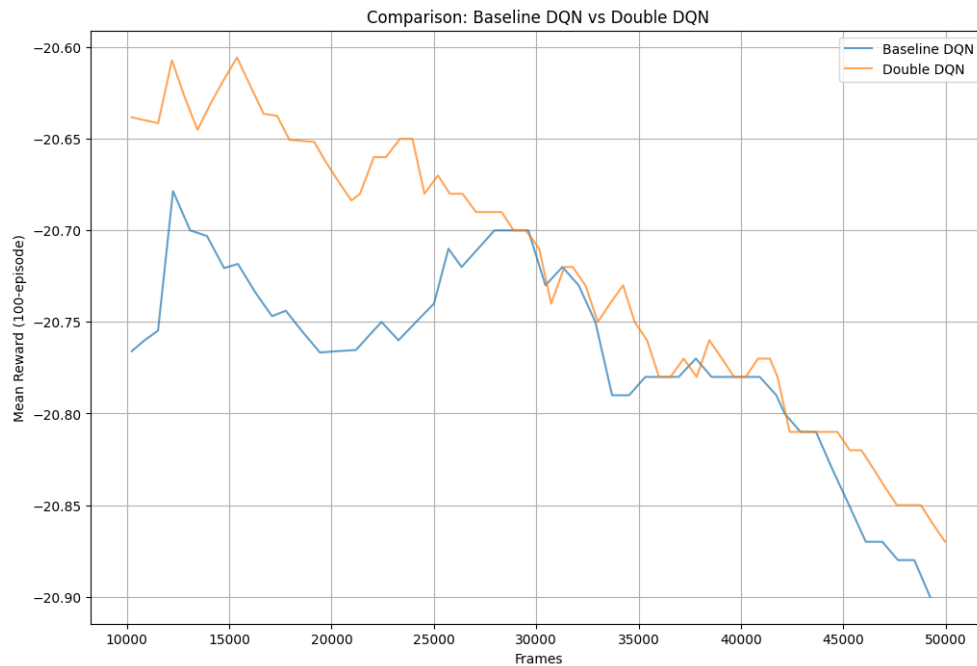
100k Frame: (Original)



50k Frame: (Epsilon Decay)



50k Frame: (Epsilon Decay + Change in learning rate + start time)



## Experimentation and Hyperparameter Analysis

Reduced Epsilon Decay Frames (100k  $\rightarrow$  50k):

- Goal: Speed up the shift from exploration to exploitation.
- Outcome: Training time dropped significantly ( $\approx 175\text{min} \rightarrow \approx 30\text{min}$ ) with only a small drop in final performance.
- Observation: Double DQN began outperforming baseline DQN after 90k frames, suggesting stronger long-term gains if training continued.

Increased Learning Rate ( $1e-4 \rightarrow 0.01$ ):

- Goal: Accelerate convergence with a more aggressive learning rate.
- Outcome: Performance stalled near -20.8 (vs. -19.5 with  $1e-4$ ), indicating instability and overshooting.
- Conclusion: A lower, stable learning rate is more effective for this task.

Key Takeaways:

- Faster epsilon decay offers a good trade-off between training time and performance.
- Double DQN shows promise for long-term improvement, even if short-term gains are modest.
- Learning rate is sensitive; too high prevents learning entirely.

#### Future Directions:

- Tune epsilon decay schedules for better exploration/exploitation balance.
- Test advanced techniques like prioritized experience replay or dueling networks.
- Experiment further with learning rates and batch sizes for optimized training.

## Conclusion

This project successfully implemented and compared baseline DQN against extended Double DQN across two distinct Atari environments. The results demonstrate clear differences in learning behavior and algorithm effectiveness between game complexities.

In Mario Bros, the environment proved particularly challenging for reinforcement learning with our base code implementation. Mario is not an easy game to tackle for both DQN & DDQN, as its complex mechanics, sparse rewards, and precise timing requirements presented significant hurdles that limited substantial policy improvement. Although DDQN does show more stable learning trends, the fundamental architecture struggled to achieve meaningful progress in this demanding environment.

On the other hand, Pong provided a clearer demonstration of DQN capabilities and the advantages of Double DQN. The simpler state space and more frequent rewards allowed both algorithms to show measurable improvement, with DDQN particularly demonstrating its value by beginning to outperform baseline DQN after approximately 90,000 frames. This performance divergence suggests that with extended training, Double DQN would achieve significantly stronger long-term results, confirming its theoretical advantage in reducing Q-value overestimation and producing more stable policies.

The comparative analysis through learning curves, gameplay videos, and performance metrics validates Double DQN as a meaningful enhancement over baseline DQN, while also highlighting how environmental complexity directly impacts reinforcement learning algorithm effectiveness and the importance of matching algorithm sophistication to game difficulty.