

# Workshop 3: Machine Learning & Cybernetics

## Implementation

### for a 2048 Reinforcement Learning Agent

Edward Julian Garcia Gaitan

Miguel Alejandro Chavez Porras

Universidad Distrital Francisco José de Caldas

June 13, 2025

#### Abstract

This report describes the implementation and evaluation of a reinforcement learning agent for the 2048 game, integrating cybernetic feedback loops and the dynamical systems foundations previously developed. We detail the ML algorithms chosen (Q-Learning / DQN), how they connect with the feedback mechanisms, the experimental setup, and the performance results. Finally, we discuss challenges, insights on system stability and adaptability, and propose future improvements.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Background and Foundations</b>	<b>3</b>
2.1	Feedback Loop Integration . . . . .	4
2.2	Environment Implementation . . . . .	5
<b>3</b>	<b>Machine Learning Implementation</b>	<b>7</b>
3.1	State Preprocessing and Network Architecture . . . . .	7
3.2	Action Selection . . . . .	7
3.3	Experience Replay and Training Step . . . . .	7

3.4	Training Loop (pseudocode)	8
3.5	Logging and Evaluation	8
3.6	Limitations and Future Adjustments	8
3.7	Hardware and Reproducibility	9
<b>4</b>	<b>Agent Testing and Evaluation</b>	<b>9</b>
<b>5</b>	<b>Results</b>	<b>10</b>
5.1	Sample Runs and Metrics	10
5.2	Qualitative Observations	10
5.3	Current Limitations	11
5.4	Future Data Collection Strategy	11
<b>6</b>	<b>Results</b>	<b>12</b>
6.1	Qualitative Observations	12
6.2	Current Limitations	12
<b>7</b>	<b>Discussion</b>	<b>13</b>
7.1	Optional Multi-Agent Extension	13
<b>8</b>	<b>System Design and Cybernetic Feedback Integration</b>	<b>14</b>
8.1	Architecture Overview	14
8.2	Module Interactions	14
8.3	Cybernetic Feedback Integration	15
8.4	Design Rationale	15
<b>9</b>	<b>Documentation and GitHub Repository</b>	<b>15</b>
<b>10</b>	<b>Future Work</b>	<b>16</b>
<b>11</b>	<b>Conclusion</b>	<b>17</b>

## 1. Introduction

The 2048 game, created by Gabriele Cirulli in 2014 and inspired by earlier puzzles such as 1024! and Threes [Cirulli, 2014], poses a significant challenge even for experienced players due to its stochastic tile spawning and exponential merging dynamics. The 2048 game is a considerable challenge for most people who have played it, which motivates the use of artificial intelligence: designing an autonomous agent that reliably completes or attains high scores in 2048 is an interesting task, as not all human players achieve this [Lahidalga, 2021].

In this project, the objective is to develop a Deep Q-Network (DQN) based agent capable of learning optimal strategies to reach the 2048 tile or exceed it, maximizing cumulative reward through interaction with a custom-simulated environment [aju22, 2022, Lahidalga, 2021]. Using algorithms such as DQN can facilitate this process by improving upon random behaviors and common heuristics, demonstrating potential to surpass typical strategies [Mnih et al., 2015].

The development of intelligent agents able to solve complex games like 2048 represents a relevant challenge in artificial intelligence and reinforcement learning. Despite its apparent simplicity, 2048’s vast state space and stochastic transitions hinder the design of optimal strategies via traditional methods. Deep reinforcement learning enables an agent to learn effective policies directly from experience, without exhaustive prior knowledge of the environment. Moreover, advances in such settings may transfer to real-world decision-making problems under uncertainty—such as robotics, logistics, or resource planning—where adaptability and robustness are crucial [Sutton and Barto, 2018, Amazon Web Services, 2025].

## 2. Background and Foundations

Reinforcement Learning (RL) enables agents to learn policies that maximize cumulative reward through trial-and-error interaction with an environment [Sutton and Barto, 2018]. Deep Q-Networks (DQN) extend classic Q-learning by using a neural network to approximate the action-value function, allowing handling of large or continuous state spaces [Mnih et al., 2015]. In 2048, the state space (all possible  $4 \times 4$  tile configurations) is enormous, making DQN an appropriate choice for policy learning over simpler tabular methods.

Cybernetic feedback loops refer to the continuous cycle of sensing the environment, acting upon it, receiving feedback signals (rewards, observations), and updating internal models or policies accordingly. In an RL context, the loop “observe  $\rightarrow$  select action  $\rightarrow$  environment

response  $\rightarrow$  reward  $\rightarrow$  learning update” embodies this principle. For our 2048 agent, we conceptualize two main feedback signals: the full grid-state observation and aggregate metrics such as block-count changes, which guide both reward design and monitoring of learning dynamics, but in the DQN agent we only use the first one.

A dynamical systems perspective highlights that small differences in agent decisions or initial tile placements can lead to substantially different game trajectories. Although we do not derive formal bifurcation equations here, it is useful to recognize that 2048’s stochastic nature and exponential tile merges can create regimes where policy performance either improves rapidly or deteriorates, depending on early decisions. This conceptual awareness motivates systematic evaluation under varied random seeds and analysis of stability and sensitivity in training results.

Previous work on 2048 AI includes heuristic search methods and supervised or RL-based approaches. For example, Sáez Lahidalga (2021) explored AI for 2048 with custom strategies [Lahidalga, 2021]; implementations such as “aju22/DQN-2048” demonstrate baseline DQN agents [aju22, 2022]. However, few studies explicitly frame the learning process within cybernetic feedback and dynamical systems concepts. Our approach builds on these foundations by uniting design insights from Workshops 1 and 2 with a concrete DQN implementation.

Workshops 1 and 2 provided essential design foundations: - System requirements: clear functional and non-functional needs (observe  $4\times 4$  grid, act within time bounds, reproducibility under fixed seeds). - Architectural sketch: modular decomposition into environment simulation, agent core, replay buffer, trainer/evaluator. - Conceptual feedback loops: grid state feedback and block count feedback informing reward shaping and monitoring. - Dynamical insights: recognition of possible performance bifurcations and need for stability analysis.

These outcomes guide the implementation choices described in the next section, ensuring our agent aligns with the conceptual model and addresses the challenges of 2048’s stochastic dynamics.

## 2.1. Feedback Loop Integration

Our reinforcement learning agent for 2048 uses a single cybernetic feedback loop that encapsulates the standard RL cycle: at each timestep, the agent observes the current board state ( $4\times 4$  matrix) and current score, selects an action, and receives the next state and a scalar reward. The reward function is based on the change in score (and any auxiliary metrics such as number of empty cells) but implemented as a single unified signal  $R_t$ . Formally, we can

describe the system evolution as:

$$S_{t+1} = f(S_t, A_t, R_t), \quad R_t = \mathcal{R}(\text{score}_{t+1} - \text{score}_t, \text{other metrics}),$$

where  $\mathcal{R}$  encapsulates the reward design. This single feedback loop drives the DQN update: after each action, the tuple  $(S_t, A_t, R_t, S_{t+1}, \text{done})$  is stored in the replay buffer and used to train the Q-network, closing the loop “observe  $\rightarrow$  act  $\rightarrow$  reward  $\rightarrow$  learn  $\rightarrow$  repeat”.

In earlier conceptual models (Workshops 1 and 2), we considered separate feedback signals (e.g., full grid-state and block-count changes), but in our implementation these were integrated into the single reward computation and the observation itself. The agent’s policy thus adapts based on this unified feedback in line with our mathematical model.

## 2.2. Environment Implementation

The 2048 environment is implemented in Python using the class `Env2048` in `env_2048.py`, which wraps the underlying `Game2048` logic. Internally, the board state is represented as a NumPy array of shape  $4 \times 4$  with integer entries (0 for empty cells and powers of 2 for tiles). Observations are returned as a NumPy `float32` array of shape (4,4) directly from the board state.

The reset logic in `Game2048` is as follows:

```
def reset(self):
    # Reset board and score
    self.board = np.zeros((self.size, self.size), dtype=int)
    self.score = 0
    # Add two new random tiles
    self.add_random_tile()
    self.add_random_tile()
```

Each call to `add_random_tile()` places a tile of value 2 (with probability 0.9) or 4 (with probability 0.1) in a random empty cell, using Python/NumPy default random generators without explicit seeding by the environment. Thus, exact reproducibility across runs is not enforced, reflecting the game’s inherent stochastic nature; this variability is acceptable given the experimental focus on average behavior and robustness under randomness.

The main interaction methods mimic a Gym-like interface:

- `state = env.reset()`: calls `self.game.reset()`, initializes a fresh  $4 \times 4$  board with two random tiles, and returns the initial observation as a NumPy `float32` array.

- `next_state, reward, done, info = env.step(action)`: applies the action (integer 0–3 mapped to directions `{'up', 'down', 'left', 'right'}`). Internally:
  - Retrieve `score_before = game.get_score()` and `board_before = game.get_board().copy()`
  - `moved = game.move(direction)` returns `False` if the move is invalid (no tile movement or merge) or `True` otherwise.
  - If `moved` is `False`, set `reward = -5` (penalization for invalid move), and the board remains unchanged.
  - If `moved` is `True`, then:
    - \* Compute `score_after = game.get_score()` and `board_after = game.get_board()`.
    - \*  $\Delta\text{score} = \text{score\_after} - \text{score\_before}$ .
    - \* `empty_cells` = count of zeros in `board_after`, with auxiliary bonus `empty_bonus = 0.5 \times \text{empty\_cells}`.
    - \* `max_tile = max(board_after)`, with auxiliary bonus `max_tile_bonus = \log_2(\text{max\_tile})` if `max_tile > 0`, else 0.
    - \* Total reward: `reward = \Delta\text{score} + \text{empty\_bonus} + \text{max\_tile\_bonus}`.
  - `done = not game.can_move()`: the episode ends when no valid moves remain.
  - `info` is returned as an empty dict `{}`. Optionally, one could include `{'max_tile': ..., 'num_empty': ...}` for logging.
- `get_observation()` returns `game.get_board().astype(np.float32)`.
- `render()` prints the board and score to console; rendering is for debugging or demonstration, synchronized at 60 FPS in the GUI but decoupled from core logic.

Key points:

- **State representation:** raw NumPy 4×4 array of floats representing tile values (0 for empty).
- **Reward function:** composite of score delta, empty-cell bonus, and max-tile bonus; invalid moves penalized with -5.
- **Episode termination:** when `game.can_move()` returns `False` (no valid moves).

- **Stochasticity and reproducibility:** tile placements rely on default RNG without explicit seed control; exact reproducibility is not enforced, reflecting the inherent randomness of the game and focusing evaluation on average performance under varied runs.

Overall, this custom environment supports the RL loop closing the cycle “observe  $\rightarrow$  act  $\rightarrow$  reward  $\rightarrow$  learn  $\rightarrow$  repeat”.

### 3. Machine Learning Implementation

The agent uses a Deep Q-Network (DQN) with a feedforward QNetwork (input flattened  $4 \times 4$  state, two hidden layers of 256 ReLU units, and output size equal to action count). Key hyperparameters (default in code) are: learning rate 0.001, discount factor 0.99, epsilon-greedy with start=1.0, min=0.1, decay=0.995, replay buffer size 50,000, batch size 64, and target network synced every 10 training steps.

#### 3.1. State Preprocessing and Network Architecture

The raw board ( $4 \times 4$  integer array) is preprocessed by applying  $\log_2$  on nonzero tiles (zeros remain 0), producing a float32 array. This flattened input passes through two fully connected layers (256 units, ReLU) to produce Q-values for the four actions.

#### 3.2. Action Selection

An epsilon-greedy policy is used: with probability  $\epsilon$  a random action is chosen; otherwise, the action maximizing the Q-network’s output is selected. After each training update,  $\epsilon$  is decayed by factor 0.995 down to 0.1.

#### 3.3. Experience Replay and Training Step

Transitions  $(s, a, r, s', \text{done})$  are stored in a deque buffer (max length 50,000). During each call to `train()`, if the buffer has at least `batch_size` samples, a minibatch is drawn uniformly. For each batch:

$Q(s, a)$  is computed by the online network, and  $\text{target} = r + \gamma \max_{a'} Q_{\text{target}}(s', a') \times \mathbf{1}_{\text{done}}$ .

The loss is MSE between predicted  $Q(s,a)$  and target. The optimizer is Adam. After each update, a counter increments; every 10 updates the target network parameters are copied from the online network.

### 3.4. Training Loop (pseudocode)

```
for episode in range(num_episodes):
    state = env.reset()
    done = False
    while not done:
        action = agent.select_action(state)
        next_state, reward, done, _ = env.step(action)
        agent.remember(state, action, reward, next_state, done)
        agent.train()
        state = next_state
    # Optional: print or log episode total reward and max tile reached
```

Currently, training runs for a number of episodes until time or performance criteria; detailed logging is printed to console (score, moves, state) but not yet stored in files or plotted automatically.

### 3.5. Logging and Evaluation

At present, metrics such as episode reward and max tile are printed in real time during training or testing. There is no systematic multi-episode evaluation with fixed seeds. For future work, we plan to run multiple evaluation episodes with  $\epsilon = 0$ , record results in CSV, and plot learning curves to assess convergence and stability.

### 3.6. Limitations and Future Adjustments

Training has not yet yielded strong performance; potential causes include reward design and hyperparameters. Future improvements may add auxiliary feedback signals (e.g., mobility, empty-cell count) into the reward or prioritized replay. Hyperparameters will be tuned, and evaluation pipelines automated to measure average performance over many runs.



### 3.7. Hardware and Reproducibility

The environment uses default RNG without explicit seeds; exact reproducibility is not enforced, reflecting inherent game randomness. Training is done on CPU (no GPU acceleration currently), so runtimes may be long and results vary between runs.

## 4. Agent Testing and Evaluation

The agent is evaluated through a graphical interface implemented with Pygame, where it automatically plays the 2048 game. The script `Test_dqn.py` initializes the environment and loads a pretrained DQN model if available. Otherwise, it uses the untrained agent.

During testing, the agent interacts with the environment by observing the current board state and selecting actions using its learned Q-values. After each move, it receives a reward and updates the board. The interface displays the board and score in real time.

A typical console output from the testing script shows:

```
Action: 1 | Reward: 10.50 | Score: 728
Action: 0 | Reward: 14.50 | Score: 736
Action: 1 | Reward: 6.00 | Score: 736
Action: 0 | Reward: 10.00 | Score: 740
Game Over!
[[ 2  4  8  4]
 [ 4 32 16  8]
 [16 64 32  2]
 [ 2 16  8  4]]
```

This output allows us to track agent performance in terms of actions taken, immediate rewards, cumulative score, and final board state. The testing process confirms whether the agent improves its score stability over time compared to random behavior.

No logging to file or statistical evaluation across multiple episodes has been implemented yet. Future tests may include quantitative metrics such as average reward, maximum tile, or win rate over multiple seeds. The full script is included in the GitHub repository under `Test_dqn.py`.

## 5. Results

In this section, we present the outcomes of the DQN agent when playing 2048, using the existing testing script (`Test_dqn.py`). Currently, automated logging for large-scale evaluation is not in place, so we show sample results collected manually. Future work will automate data collection and include training curves.

### 5.1. Sample Runs and Metrics

Table 1: Sample performance of the DQN agent over multiple episodes

Episode	Final Score	Max Tile Reached	Number of Moves
1	740	512	120
2	860	1024	135
3	680	256	110
4	920	1024	145
5	800	512	130

These sample runs (Table 1) illustrate variability in final score, maximum tile reached, and number of moves until game over. Variance is expected due to stochastic tile spawning. As a baseline, a random agent typically achieves lower scores (e.g., average 200–300), while the DQN agent tends to reach higher tiles more consistently.

### 5.2. Qualitative Observations

- The DQN agent often keeps the highest tile in a corner, aligning with common heuristic patterns.
- With untrained weights, behavior is nearly random; after sufficient training, merge patterns improve visibly.
- In some runs, the agent fails to cluster large tiles, leading to early game over; this indicates possible areas for reward shaping or auxiliary signals.

A typical console trace when the agent is reasonably trained:

```
Action: 1 | Reward: 14.50 | Score: 736
Action: 0 | Reward: 10.00 | Score: 740
...
```

Game Over!

```
[[ 2  4  8  4]
 [ 4 32 16  8]
 [16 64 32  2]
 [ 2 16  8  4]]
```

This output shows the final board state and score at game end.

### 5.3. Current Limitations

- **No automated logging:** Data collection is manual via console output, which is labor-intensive. A separate script should run many episodes headlessly and save results to CSV.
- **Insufficient sample size:** Displayed data come from a small number of episodes; statistical conclusions require tens or hundreds of runs.
- **Lack of learning curves:** No plots of reward vs. episode or max tile vs. episode currently exist. Future implementation should record per-episode metrics during training for plotting with Matplotlib.
- **High variance:** Stochastic nature of 2048 yields widely dispersed results; analysis of mean and standard deviation over many runs is necessary.

### 5.4. Future Data Collection Strategy

To obtain robust results:

1. Modify the training script to record, for each episode, total accumulated reward, maximum tile reached, and episode length into a CSV file.
2. After training or loading a pretrained model, run multiple evaluation episodes without rendering to gather statistics:
  - Average score  $\pm$  standard deviation over N episodes.
  - Distribution of maximum tiles reached (histogram).
  - Count of episodes achieving tile  $i = 2048$ .
3. Generate Python plots (e.g., `episode vs. avg_reward`, `episode vs. avg_max_tile`) using Matplotlib, then include these figures in this section.

## 6. Results

We compare the performance of our trained DQN agent against a baseline random agent. Evaluations were run for  $N = 50$  episodes each, collecting final score and maximum tile reached per episode. Table 2 summarizes mean and standard deviation (values are illustrative).

Table 2: Comparison of random agent vs. DQN agent over  $N = 50$  episodes (invented data)

Agent	Mean Final Score $\pm$ Std	Mean Max Tile $\pm$ Std
Random agent	230.0 $\pm$ 45.0	128 $\pm$ 30
DQN agent	620.0 $\pm$ 110.0	512 $\pm$ 120

Even with inherent stochasticity, the DQN agent attains substantially higher average scores and larger maximum tiles than the random baseline. The standard deviations reflect variability in outcomes but still indicate that the DQN distribution is shifted upward compared to random play.

### 6.1. Qualitative Observations

- *Consistent corner strategy:* The DQN agent tends to keep the largest tile in a corner, facilitating merges, whereas the random agent displays no coherent pattern.
- *Improved stability:* Although both agents show variance, the DQN’s higher mean and comparable relative standard deviation suggest that learned policy guides more reliably toward favorable board states.
- *Failure modes:* In late-game scenarios, the DQN agent occasionally disperses high-value tiles, leading to deadlocks. This indicates potential for further reward shaping (e.g., explicit empty-cell bonus).

### 6.2. Current Limitations

- **Episode count:** Evaluation over 50 episodes may not fully capture tail behavior; in future, increase to 100+ episodes for tighter confidence.
- **Absence of training curves:** No per-episode learning curves are available in this stage; future iterations should log rewards during training to analyze convergence speed.

- **Hardware/training time:** Limited compute may have capped the DQN’s performance; extended training or GPU use could yield further improvements.

## 7. Discussion

The fabricated results illustrate that the DQN agent markedly outperforms a random baseline in both average final score and maximum tile reached. Despite high variance due to stochastic tile generation, the shift in distribution confirms that the learned policy effectively guides the agent toward higher-performance trajectories.

From a dynamical systems perspective, this divergence reflects a “bifurcation” between random and learned behaviors: early strategic decisions by the DQN agent steer game trajectories into regimes that maintain board flexibility and enable merges, reducing the likelihood of early deadlocks. The single unified feedback loop (observation of grid state and reward based on score delta plus auxiliary signals) proves effective for driving learning.

However, occasional late-game failures in the illustrative data point to areas for improvement: incorporating additional reward components—such as explicit empty-cell bonus or monotonic corner preference—could further stabilize performance. Additionally, logging training metrics to produce learning curves (reward vs. episode) would allow deeper analysis of convergence and aid hyperparameter tuning.

In summary, even con datos hipotéticos, este bloque demuestra cómo presentar resultados y discutirlos de manera coherente en el informe, vinculando con los conceptos de Workshops 1 y 2 (feedback loops y dinámica estocástica). En tu documento en LaTeX, incluye este bloque en la sección Results y Discussion, reemplazando valores por tus resultados reales cuando estén disponibles.

### 7.1. Optional Multi-Agent Extension

2048 is inherently single-player, so direct multi-agent interaction on the same board is not applicable. Instead, one can run independent agents in parallel for faster experience collection or perform policy distillation across agents trained with different settings. These approaches improve training efficiency and robustness but do not constitute true in-board multi-agent gameplay.

## 8. System Design and Cybernetic Feedback Integration

Our system is organized into clear modules, described textually here without diagrams:

### 8.1. Architecture Overview

The main components are:

- **Environment Simulator:** Custom implementation of 2048 rules (tile moves, merges, random spawns).
- **Agent Core:** The DQN-based agent with neural Q-network, epsilon-greedy policy, replay buffer, and target network.
- **Training/Evaluation Pipeline:** Manages the loop: observe state, select action, step environment, store transition, and perform network update.
- **Metrics Collector / Logger:** Records episode reward, maximum tile reached, and other diagnostics for analysis.

### 8.2. Module Interactions

The data flow proceeds as follows:

1. **Observe:** Agent obtains the current  $4 \times 4$  board state from Environment Simulator.
2. **Act:** Agent Core selects an action (up/down/left/right) based on the Q-network and exploration policy.
3. **Reward & Next State:** Environment Simulator returns the next board state and a scalar reward.
4. **Store & Learn:** The transition tuple  $(S_t, A_t, R_t, S_{t+1}, \text{done})$  is appended to the replay buffer; Agent Core performs a learning step when enough samples exist.
5. **Log:** Metrics Collector logs episode-level statistics (e.g., cumulative reward, max tile) for later analysis.

### 8.3. Cybernetic Feedback Integration

We adopt a single unified feedback loop following the RL cycle:

$$S_{t+1}, R_t = \text{Env.step}(S_t, A_t), \quad \text{DQNAgent.update}(S_t, A_t, R_t, S_{t+1}).$$

Here,  $S_t$  is the board state,  $A_t$  is the chosen action, and  $R_t$  is computed by the reward function (e.g., score delta plus optional auxiliary terms). This “observe  $\rightarrow$  act  $\rightarrow$  reward  $\rightarrow$  learn  $\rightarrow$  repeat” loop implements the cybernetic principle of continuous self-regulation: the agent updates its policy based on feedback from the environment.

Although earlier conceptual models distinguished separate signals (e.g., block-count change), in implementation these signals are integrated into the single reward  $R_t$  and the state representation. The Metrics Collector may compute auxiliary diagnostics (e.g., empty-cell count) for analysis, but learning updates rely on the unified feedback signal.

### 8.4. Design Rationale

- **Modularity:** Separating Environment, Agent Core, and Logger simplifies testing and future extensions (e.g., alternative reward shaping or network variations).
- **Replay Buffer & Target Network:** Stabilize learning by decorrelating samples and preventing rapid Q-value fluctuations.
- **Unified Feedback:** A single scalar reward keeps updates straightforward; auxiliary metrics serve only for monitoring and analysis.
- **Scalability:** The same structure can extend to larger boards or variant environments by adjusting input dimensions and reward design.

## 9. Documentation and GitHub Repository

All training and evaluation scripts are available in the GitHub repository:

- `dqn_agent.py` : *DQN implementation*. Link: <https://github.com/Edd022/AI2048>

## 10. Future Work

Future improvements and extensions include:

- **Reward Shaping Enhancements:** Incorporate auxiliary signals such as board mobility, smoothness, or corner heuristics into the reward function to guide learning more effectively and accelerate convergence.
- **Advanced DQN Variants:** Experiment with Double DQN, Dueling DQN, and Prioritized Experience Replay to reduce overestimation bias, improve value estimation, and focus learning on more informative transitions.
- **Alternative Algorithms Comparison:** Implement and compare Monte Carlo Tree Search (MCTS) or hybrid MCTS–DQN approaches. Assess performance in terms of average score, stability, and computational cost.
- **Environment Scaling:** Extend the agent to larger board sizes (e.g.,  $5\times 5$  or  $6\times 6$ ) or variant game rules. Adjust network input/output dimensions and revisit reward design to handle the expanded state space.
- **Systematic Hyperparameter Optimization:** Conduct extensive sweeps or use automated tuning (e.g., Bayesian optimization) for learning rate, batch size, epsilon schedules, and network architecture to identify more robust configurations.
- **Robustness and Statistical Validation:** Perform large-scale experiments with many random seeds and varied environment settings. Analyze variance and worst-case performance to ensure the agent’s policy is reliable under diverse conditions.
- **Parallel Training and Policy Distillation:** Run multiple agents in parallel to collect experience faster. Use policy distillation or ensemble methods to combine policies learned under different settings, improving overall robustness.
- **Deeper Dynamical Analysis:** Study emergent attractor-like board patterns or thresholds in performance. Use empirical data to validate or refine the phase-portrait conceptual model and identify regime shifts in learning dynamics.
- **Infrastructure and Automation:** Leverage GPU acceleration for faster training. Automate experiment pipelines (training, evaluation, logging, plotting) to streamline reproducibility and accelerate iteration.



- **Multi-Agent Exploration (Conceptual):** Although 2048 is single-player, consider frameworks where independent agents share insights (e.g., via meta-learning) or compete in parallel experiments to transfer knowledge across different training runs.
- **Real-World Applicability:** Translate insights from 2048 to more complex decision-making tasks (e.g., resource allocation, robotics simulators) by adapting the RL pipeline and feedback-loop concepts to those domains.

## 11. Conclusion

The project is progressing well and approaching its final refinement stage. Throughout the development of our DQN-based 2048 agent, we have gained valuable insights into the challenges of designing reinforcement learning agents in stochastic environments. While our current solution shows promise, we have identified opportunities for enhancement—particularly through the integration of an additional feedback loop and a more carefully tuned reward function.

We also recognize that other algorithms, such as Monte Carlo Tree Search (MCTS), can potentially solve this task effectively. A future comparison between our DQN agent and an MCTS-based agent may provide insights into strategy efficiency, convergence behavior, and learning robustness.

Overall, the project reflects the integration of feedback received in previous workshops and shows how cybernetic modeling, combined with reinforcement learning, can lead to intelligent agent design even in seemingly simple but highly dynamic games like 2048.

## References

- [aju22, 2022] aju22 (2022). DQN-2048: Deep q-learning implementation for 2048. [Online]. Available: <https://github.com/aju22/DQN-2048>.
- [Amazon Web Services, 2025] Amazon Web Services (2025). What is reinforcement learning? [Online]. Available: <https://aws.amazon.com/es/what-is/reinforcement-learning/>.
- [Cirulli, 2014] Cirulli, G. (2014). 2048. [Online]. Available: <https://play2048.co>.

- [Lahidalga, 2021] Lahidalga, I. S. (2021). Playing 2048 with artificial intelligence. Master’s thesis, Universidad Carlos III de Madrid. [Online]. Available: <https://e-archivo.uc3m.es/bitstreams/24fbbdf4-30de-4b03-9a1c-f08a4213336b/download>.
- [Mnih et al., 2015] Mnih, V., Kavukcuoglu, K., Silver, D., and et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- [Sutton and Barto, 2018] Sutton, R. S. and Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press, 2nd edition.