# Deep Q-Learning with Atari Boxing

**Github Link:** https://github.com/tanv99/LLM-Agents-Deep-Q-Learning-with-Atari-Game

## Overview

This project implements a Deep Q-Network (DQN) agent to play Atari Boxing using Stable-Baselines3. Five independent experiments tested different hyperparameters, each training for 500,000 steps (~1,120 episodes).

**Key Results:**

- Best Performance: Experiment 4 (Moderate Epsilon) - 26.25 ± 10.81 test reward
- Worst Performance: Experiment 5 (Boltzmann Policy) - -3.50 ± 12.27 test reward
- Total Training: 5,596 episodes, ~5.8 hours

---

## Game Overview: Atari Boxing-v5

**Environment Specifications:**

- Raw Observation: 210×160×3 RGB images
- Preprocessed: 84×84×4 grayscale stacked frames
- Action Space: 18 discrete actions (movement + punch combinations)
- Rewards: +1 (punch landed), -1 (punch received), 0 (no contact)
- Episode Length: ~447 steps average

**Why Deep Q-Learning is Required:**

Traditional Q-table size: $256^{(84×84×4)} × 18 ≈ 10^{67,000}$ entries - computationally impossible.

Solution: CNN function approximator maps 84×84×4 input → 18 Q-values.

*See Notebook: "Environment Specification and Analysis" cell*

---

## 1. Baseline Performance

**Configuration:**

- Training Steps: 500,000
- Learning Rate ($α$): 1e-4

- Discount Factor (γ): 0.99
- Epsilon Decay: 1.0 → 0.01 (over 80% of training)
- Replay Buffer: 100,000
- Batch Size: 32
- Target Update: Every 1,000 steps

**Results:**

- Test Reward: 22.60 ± 8.48
- Episodes: 1,118
- Training Time: 61.3 minutes

*See Notebook: "Experiment 1: Baseline" cell*

---

# 2. Environment Analysis

**States:**

- Preprocessed state space: 84×84×4 tensor
- Theoretical size: $256^{(84×84×4)} ≈ 10^{67,000}$ possible states

**Actions:**

- 18 discrete actions (8 directions × 2 punch states)

**Q-Table Size:**

- Would require: $10^{67,000} × 18$ entries (impossible)
- DQN Solution: CNN with ~1.2M parameters outputs 18 Q-values

*See Notebook: "Environment Specification and Analysis" cell*

---

# 3. Reward Structure

**Rewards Used:**

- +1: Agent lands punch
- -1: Agent receives punch
- 0: No contact

**Rationale:**

- Native environment rewards - proven effective
- No artificial reward shaping to avoid bias
- Encourages learning genuine winning strategies

- Sparse signal requires extensive exploration (hence 500K steps)

*See Notebook: "Reward Structure" section*

---

# 4. Bellman Equation Parameters

**Alpha (Learning Rate): 1e-4**

- Standard for Adam optimizer with neural networks
- Balances learning speed vs. stability

**Experiment 2 - Higher LR ($\alpha$ = 2.5e-4):**

- Result: 1.15 ± 5.96 test reward
- Conclusion: 2.5× higher LR caused instability and poor performance
- Initial test at 10× (1e-3) caused severe training failure

**Gamma (Discount Factor): 0.99**

- Values rewards ~100 steps into future
- Appropriate for 447-step average episode length

**Experiment 3 - Lower Gamma ($\gamma$ = 0.8):**

- Result: 15.75 ± 14.19 test reward
- Conclusion: Too myopic (only 5-step lookahead), inferior strategic planning

**Finding:** Baseline parameters ($\alpha$=1e-4, $\gamma$=0.99) performed best.

*See Notebook: "Experiment 2: Moderate Learning Rate" and "Experiment 3: Gamma = 0.8"*

---

# 5. Policy Exploration

**Alternative Policy Tested: Boltzmann/Softmax Exploration**

**Epsilon-Greedy (Baseline):**

- $\epsilon$ probability: random action
- (1-$\epsilon$) probability: argmax $Q(s,a)$

**Boltzmann Policy (Experiment 5):**

- $P(a|s) = \exp(Q(s,a)/\tau) / \Sigma \exp(Q(s,a')/\tau)$
- Temperature decay: $\tau = 2.0 \rightarrow 0.1$

**Results:**

- Test Reward: -3.50 ± 12.27 (worst performing)
- High variance indicates unstable learning

**Why It Failed:** The Boltzmann policy struggled because sparse rewards meant Q-values remained similar across actions during early training, making temperature-based exploration less effective than uniform random sampling. Epsilon-greedy's discrete exploration better suited Boxing's 18-action space with infrequent reward signals.

*See Notebook: "Experiment 5: Boltzmann Policy" cell*

---

# 6. Exploration Parameters

**Baseline Epsilon Schedule:**

- Starting ε: 1.0 (100% random)
- Final ε: 0.01 (1% exploration)
- Decay: Linear over 80% of training (400,000 steps)

**Rationale:** Start with full exploration, gradually transition to exploitation, maintain 1% exploration to avoid local optima.

**Experiment 4 - Modified Exploration:**

- Final ε: 0.05 (5% vs 1%)
- Decay: Over 60% of training (300,000 steps)

**Results:**

- Test Reward: 26.25 ± 10.81 (**BEST** performance)
- Faster convergence + higher residual exploration = superior results

**Epsilon at Max Steps:**

- Baseline: ε = 0.01 after 400,000 steps, remains constant
- Experiment 4: ε = 0.05 after 300,000 steps, remains constant

**Conclusion:** 5% residual exploration prevents premature convergence better than 1%.

*See Notebook: "Experiment 4: Faster Epsilon Decay" cell*

---

# 7. Performance Metrics

**Average Steps Per Episode:**

- Experiment 1 (Baseline): 447.2 steps
- Experiment 2 (Higher LR): 447.2 steps
- Experiment 3 (γ=0.8): 444.8 steps
- Experiment 4 (Moderate ε): 447.2 steps
- Experiment 5 (Boltzmann): 447.2 steps

**Analysis:** Episode length remarkably consistent (~445-447 steps) across all experiments, indicating performance differences lie in reward accumulation (punches landed vs. received), not survival duration.

*See Notebook: "Load and Analyze Results" section*

---

# 8. Q-Learning Classification

**Question:** "Does Q-learning use value-based or policy-based iteration? Explain in detail"

Q-learning is fundamentally a **value-based method**. This classification stems from what the algorithm directly learns and optimizes. Value-based methods learn a value function—$V(s)$ for states or $Q(s,a)$ for state-action pairs—and derive the policy implicitly from these values. Q-learning specifically learns $Q(s,a)$, the expected cumulative discounted reward from taking action 'a' in state 's'.

The Bellman equation $Q^*(s,a) = E[r + γ \cdot max\_a' \, Q^*(s',a')]$ defines the optimal value function, where the max operator implicitly represents the greedy policy $π^*(s) = argmax\_a \, Q^*(s,a)$. Crucially, this policy is never explicitly stored or directly optimized—it emerges as a byproduct of value learning.

In Deep Q-Learning, the neural network parameters θ minimize temporal difference error: $L(θ) = E[(r + γ \cdot max\_a' \, Q(s',a';θ^-) - Q(s,a;θ))^2]$. This loss function measures how well Q-values satisfy the Bellman equation, confirming the value-based nature.

**Contrast with Policy-Based:**

Policy-based methods (REINFORCE, PPO) parameterize the policy directly as $π(a|s;θ)$ and optimize θ to maximize expected returns through policy gradient ascent. They never explicitly learn Q-values.

**Why This Matters:**

- Value-based: More sample-efficient for discrete actions, deterministic environments
- Policy-based: Better for continuous actions, stochastic environments

---

# 9. Q-Learning vs LLM Agents

**Question:** "How does Deep Q-Learning differ from LLM-based agents? Explain in detail"

Deep Q-Learning and LLM-based agents represent fundamentally different approaches to sequential decision-making:

**1. Learning Mechanism:**

DQN learns through trial-and-error environmental interaction. The Boxing agent starts with random weights, takes 500,000 game actions, receives +1/-1 rewards for punches, and gradually learns Q(s,a) through Bellman updates. This is purely empirical—no prior knowledge, just data-driven pattern recognition.

LLM agents leverage pre-training on massive text corpora before task-specific fine-tuning. An LLM could reason "In boxing, maintain defensive positioning and counter-punch when openings appear" using linguistic knowledge from training data describing sports strategies. Task behavior is refined through Reinforcement Learning from Human Feedback (RLHF), where human preferences replace environmental rewards.

**2. State Representation:**

DQN processes raw pixels (210×160×3 images) through CNN layers extracting spatial features: opponent position, punch trajectories, movement patterns.

LLM agents operate on discrete text tokens, using transformer embeddings. For Boxing, pixels must be translated to text descriptions ("opponent approaching from left") via vision encoders before language model reasoning applies.

**3. Action Spaces:**

DQN outputs 18 Q-values for discrete game controls (directional movement + punch).

LLM agents generate tokens from 50K+ vocabulary, requiring parsing into executable actions through structured prompting or fine-tuning.

**4. Planning Horizon:**

DQN learns multi-step value through $\gamma^n$ weighting—a punch at timestep t affects Q-values n steps earlier via recursive Bellman updates.

LLM agents perform autoregressive token generation with attention mechanisms, explicitly generating reasoning chains ("Step 1: Block, Step 2: Counter-punch") but limited to 3-10 step lookahead typically.

**5. Knowledge Source:**

DQN learns task-specific policies from scratch—needs 500K Boxing interactions to discover basic strategies.

LLM agents use pre-trained world models—could reason about boxing from pre-training, requiring fewer task-specific examples.

**Example:**

- DQN learns "punch when opponent close" through 500K environmental interactions
- LLM could reason "Attack when in range" from pre-training, then refine via RLHF

---

# 10.  Expected Lifetime Value

**Question:** "What is meant by expected lifetime value in the Bellman equation? Explain in detail"

The expected lifetime value represents the estimated cumulative discounted reward from a given state through episode end. Mathematically:

```
Q(s,a) = E[R_t + γR_{t+1} + γ²R_{t+2} + ... | S_t=s, A_t=a]
```

**Components:**

**"Expected"** - Averages over possible future trajectories weighted by probabilities. In deterministic Boxing, reduces to single trajectory, but framework remains.

**"Lifetime"** - Sum extends until episode termination (~447 steps in Boxing). Includes all future rewards from current action onward.

**"Value"** - Cumulative reward modified by discount factor $\gamma$. With $\gamma=0.99$, a reward 100 steps ahead contributes $0.99^{100} \approx 0.366$ to current Q-value.

**Bellman Decomposition:**

The recursive structure breaks lifetime value into immediate + future components:

```
Q(s,a) = R(s,a) + γ·E[max_a' Q(s',a')]
```

**Concrete Boxing Example:**

State: Opponent approaching with fist raised

Q(block) = +5:

- Immediate reward: 0 (no punch landed)
- Future value: ~+5 (good defensive position → counter-punch opportunities)

Q(aggressive punch) = -2:

- Immediate: Possible +1 if successful
- Future: Negative (likely counter-punch -1 + bad positioning)

The expected lifetime value captures this multi-step strategic reasoning in a single Q-estimate.

**In Deep Q-Learning:**

Neural network approximates expected lifetime value. TD error $[r + \gamma \cdot max_{a'} Q(s',a';\theta^-) - Q(s,a;\theta)]^2$ measures how well current Q-estimates satisfy Bellman equation. Minimizing this error through gradient descent trains network to predict cumulative future rewards.

# 11.  RL for LLM Agents

**Question:** "How might RL concepts apply to LLM-based agents? Explain in detail"

Reinforcement learning principles from Deep Q-Learning directly inform modern LLM agent training through several key applications:

**1. Reward Modeling (RLHF):**

DQN learns Q(s,a) from environmental rewards (+1/-1 for punches). LLM agents use Reinforcement Learning from Human Feedback (RLHF):

- Train reward model $R_\theta(s,a)$ on human preference comparisons
- Use learned rewards to fine-tune language model via PPO
- Analogous to DQN learning from environmental rewards

**ChatGPT Training Example:**

1. Supervised learning (like initializing DQN)
2. RLHF fine-tuning (like Q-learning updates)
3. Result: Agent maximizes human satisfaction (like cumulative reward)

**2. Exploration-Exploitation:**

DQN: $\varepsilon$-greedy balances trying new actions vs. exploiting learned policy

LLM: Temperature sampling

- High temperature $\rightarrow$ diverse outputs (exploration)
- Low temperature $\rightarrow$ deterministic responses (exploitation)

**3. Credit Assignment:**

DQN: Which action caused delayed reward? Bellman updates propagate value backwards.

LLM: Which tokens contributed to human preference? Advantage function $A(s,a) = Q(s,a) - V(s)$ in PPO serves similar role.

**4. Multi-Step Reasoning:**

DQN: Temporal difference learning captures delayed consequences through $\gamma^n$ discounting

LLM: Chain-of-thought prompting decomposes complex queries into sequential sub-problems

**Concrete Application - Customer Service Chatbot:**

- Initial supervised fine-tuning on human conversations (like pre-training)
- Reward model $R_\theta$ scores responses for helpfulness (replaces +1/-1 punches)

- PPO fine-tunes to maximize expected $R_\theta$ (like optimizing Q-values)
- Temperature sampling explores diverse responses during training
- Low temperature at deployment exploits learned preferences

---

# 12.  Planning in RL vs LLM Agents

**Question:** "How does planning differ in traditional RL vs LLM agents? Include examples and frameworks"

**Traditional RL Planning (Model-Based):**

Uses explicit environment model to simulate future trajectories before selecting actions.

**AlphaGo Example:**

- Builds search tree of possible move sequences
- Simulates millions of game continuations
- Evaluates leaf nodes with value network $V(s)$
- Selects action with highest visit count after extensive search

**Process:**

1. Expand tree with legal moves
2. Simulate continuations using policy network $\pi(a|s)$
3. Evaluate positions with value network
4. Backpropagate values through tree
5. Output: Best move after 800K simulations

**Characteristics:**

- Explicit forward simulation
- Requires known/learned environment dynamics
- Computationally expensive (minutes per decision)

**LLM Agent Planning (Implicit Reasoning):**

Generates intermediate reasoning steps as natural language before final answers.

**Chain-of-Thought Example:**

Question: "Best boxing strategy when health is low?"

LLM Planning:
1. "First, prioritize defense to avoid further damage"
2. "Look for counter-punch opportunities"
3. "Maintain distance until opening appears"

Answer: "Play defensively and wait for openings to counter"

**Characteristics:**

- Sequential token generation (autoregressive)
- Planning embedded in text
- No explicit world model—knowledge in transformer weights
- Moderate computation (forward passes)

**Key Differences:**

| Aspect | RL Planning | LLM Planning |
|--------|-------------|--------------|
| Method | Tree search/rollouts | Text generation |
| Model | Explicit T(s'\|s,a) | Implicit in weights |
| Evaluation | V(s) or Q(s,a) | Language likelihood |
| Depth | 100+ steps possible | Typically 3-10 steps |
| Cost | High (millions of sims) | Moderate |
| Domain | Task-specific | Broadly applicable |

**Chess Move Comparison:**

AlphaZero:

- 40 moves lookahead
- 800K simulations
- Minutes of computation
- Exhaustive search within depth

GPT-4:

- ~3 moves lookahead
- Reasoning in text
- Seconds of computation
- Approximate but incorporates broader strategy

---

# 13. Q-Learning Algorithm

**Question:** "Explain the Q-learning algorithm. Include pseudocode and mathematics"

Q-learning learns the optimal action-value function $Q^*(s,a)$ = expected cumulative discounted reward from state s taking action a.

**Mathematical Foundation:**

Bellman Optimality Equation:

```
Q*(s,a) = E[R(s,a) + γ·max_a' Q*(s',a')]
Update Rule:
Q(s,a) ← Q(s,a) + α[r + γ·max_a' Q(s',a') - Q(s,a)]
                   └_____┘
                                 TD Target
```

**Components:**

- **TD Target:** r + γ·max_a' Q(s',a') (observed reward + discounted best next action)
- **TD Error:** δ = Target - Q(s,a) (prediction error)
- **Learning Rate:** α controls update magnitude

**Pseudocode:**

```
# Initialize
Initialize Q(s,a) for all states, actions
epsilon = max_epsilon

for episode in range(total_episodes):
    s = environment.reset()
    done = False

    while not done:
        # Action selection (ε-greedy)
        if random() < epsilon:
            a = random_action()
        else:
            a = argmax_a Q(s,a)

        # Execute action
        s', r, done = environment.step(a)

        # Q-learning update (Bellman equation)
        td_target = r + gamma * max_a' Q(s',a')
        td_error = td_target - Q(s,a)
        Q(s,a) = Q(s,a) + alpha * td_error

        s = s'
```

Deep Q-Learning Extensions:

1. Neural Network: Q(s,a;θ) replaces Q-table
2. Experience Replay: Sample random transitions from buffer D

3. Target Network: Use $\theta^-$ for stable TD targets

Loss Function: $L(\theta) = E[(r + \gamma \cdot max\_a' Q(s',a';\theta^-) - Q(s,a;\theta))^2]$

See Notebook: "Helper Functions" and preprocessing cells

---

# 14. LLM Integration

Question: "How could you integrate Deep Q-Learning with LLM systems? Describe architectures and applications"

Architecture 1: Hierarchical (LLM Strategy + DQN Tactics)

High-Level (LLM): Input: "Navigate robot to pick up red block" Output: ["Move to block", "Grasp object", "Lift upward"]

Low-Level (DQN):
  State: Joint angles + camera pixels Actions: Motor torques Policy: $\pi(a|s, goal)$ conditioned on LLM sub-goal

**Application:** Robotic manipulation - LLM decomposes complex instructions, DQN learns precise motor control.

### Architecture 2: LLM as Advisor

```python
class LLMAdvisedDQN:
    def select_action(self, state):
        # Consult LLM periodically
        if turn % 10 == 0:
            advice = llm.generate(f"Game state: {state}. Best strategy?")
            self.strategy = advice

        # DQN uses LLM context
        q_values = dqn(state, self.strategy)
        return argmax(q_values)
```

**Application:** Trading bot - LLM analyzes market news for regime, DQN executes trades within regime.

### Architecture 3: LLM Reward Shaping

```python
def get_reward(state, action, next_state, env_reward):
    base_reward = env_reward  # +1/-1 from environment

    # LLM evaluates strategic quality
    description = state_to_text(state, action)
```

```
    llm_eval = llm.rate_move(description)  # -1 to +1

    return base_reward + 0.1 * llm_eval
```

**Application:** Provides dense reward signals for sparse environments, accelerates learning.

**Practical Examples:**

1. **Autonomous Driving:** LLM interprets traffic signs/situations, DQN controls vehicle
2. **Personal Assistant:** LLM handles language understanding, DQN optimizes UI/UX through learned preferences

---

# 15.  Code Attribution

**Original Implementation (My Work):**

- 5-experiment design and configurations
- Multi-day training system with Google Drive persistence
- `train_experiment()` function and pipeline
- `BoltzmannExplorationCallback` and `train_boltzmann_experiment()`
- Results aggregation and visualization code
- All markdown documentation

**Adapted Code (Modified):**

**AtariPreprocessing:**

- Source: Mnih et al. (2015) preprocessing pipeline
- Modifications: Adapted for Gymnasium API, custom wrapper implementation
- Functions: Grayscale conversion, 84×84 resize, frame skipping

**FrameStack:**

- Source: Standard frame stacking pattern
- Modifications: Custom 4-frame implementation
- Function: Temporal information (velocity/motion)

**EpisodeTracker:**

- Source: Stable-Baselines3 BaseCallback interface
- Modifications: Original tracking logic implementation

**Third-Party Libraries (Unmodified):**

- Stable-Baselines3: DQN implementation, training loops (MIT)
- Gymnasium: Environment interface (MIT)

- ALE: Boxing-v5 emulation (GPL-2.0)
- OpenCV, NumPy, Matplotlib: Utilities

**Academic References:**

- Mnih et al. (2015) - DQN algorithm
- Sutton & Barto (2018) - Q-learning theory
- Stable-Baselines3 docs - Library usage

**Consulted (No Code Copied):**

- ykteh93/Deep_Reinforcement_Learning-Atari
- shubhlohiya/playing-atari-with-deep-RL
- Stable-Baselines3 examples

*Complete citations in notebook: "Code Attribution and Citations" section*

# 16. Licensing

**License:** MIT License (see final notebook cell)

**Key Points:**

1. Original code released under MIT
2. Compatible with all dependencies:
   - Stable-Baselines3: MIT
   - Gymnasium: MIT
   - ALE: GPL-2.0 (used as library, no modifications)
3. Academic citations separate from code license
4. Permits free use, modification, distribution

**GPL Compliance:**

ALE is GPL-2.0 licensed, but this project uses it as an unmodified library dependency, which is permitted under GPL without requiring derivative works to be GPL-licensed.

*Complete license text in notebook: "MIT License" cell*

# Conclusion

This Deep Q-Learning implementation successfully demonstrates reinforcement learning fundamentals while systematically exploring hyperparameter effects. The baseline agent achieved competent Boxing performance (22.60 test reward) through 500,000 environment interactions, learning effective strategies from sparse rewards.

**Key Experimental Findings:**

1. **Learning Rate:** Higher α (2.5e-4) destabilized learning, degrading performance to 1.15 reward

2. **Discount Factor:** Lower γ (0.8) produced myopic strategies (15.75 reward), insufficient for multi-step planning
3. **Exploration Schedule:** Moderate decay (60% duration, 5% final ε) outperformed baseline (80%, 1%), achieving best results (26.25 reward)
4. **Policy Comparison:** Boltzmann exploration catastrophically failed (-3.50 reward) versus epsilon-greedy for sparse rewards

**Validated Theoretical Predictions:**

Results align with Mnih et al. (2015) recommendations: α=1e-4 and γ=0.99 are well-calibrated for Atari environments. Experiments provide empirical validation while revealing that moderate exploration schedules offer incremental improvements.

**RL and LLM Agent Synergies:**

The assignment connects classical RL to modern LLM agents through multiple lenses. RLHF mirrors Q-learning's reward-driven optimization, temperature sampling parallels epsilon-greedy exploration, and chain-of-thought reasoning relates to multi-step value propagation. Hybrid architectures combining LLM strategic planning with DQN tactical execution point toward future AI systems leveraging both pre-trained knowledge and task-specific learning.

**Professional Standards:**

The implementation demonstrates:

- **Technical Competence:** Successful DQN implementation with appropriate tools
- **Scientific Rigor:** Systematic 5-experiment design testing distinct hypotheses
- **Problem-Solving:** Google Drive persistence enabling multi-day training
- **Communication:** Clear documentation of methodology, results, and implications

**Broader Implications:**

This work illustrates that while DQN excels at precise reactive control through environmental interaction (learning pixel→action mappings for Boxing), LLM agents provide complementary capabilities in linguistic reasoning and knowledge generalization. The future of AI agents likely involves synthesizing these approaches: LLM world models guiding DQN-based control, creating systems that reason abstractly while acting optimally in specific environments.

The findings confirm that Deep Q-Learning remains a powerful approach for sequential decision-making in well-defined environments with discrete actions and visual observations. The systematic experimental methodology provides a template for investigating algorithm behavior across parameter spaces. As a portfolio piece, this project evidences capability for independent research, rigorous experimentation, and professional software engineering suitable for machine learning research or engineering roles.

**Final Result:** Moderate epsilon-greedy exploration (60% decay to 5% final) with standard DQN hyperparameters (α=1e-4, γ=0.99) yields optimal Atari Boxing performance.

# References

**Academic Papers:**

1. Mnih, V., et al. (2015). "Human-level control through deep reinforcement learning." *Nature*, 518(7540), 529-533.
2. Sutton, R. S., & Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.

**Software:**

- Stable-Baselines3: https://github.com/DLR-RM/stable-baselines3 (MIT)
- Gymnasium: https://github.com/Farama-Foundation/Gymnasium (MIT)
- ALE: https://github.com/Farama-Foundation/Arcade-Learning-Environment (GPL-2.0)