

Deep Learning Memory Maze Project - Complete Student Code Analysis

Course: Deep Reinforcement Learning **Project:** Memory Maze Navigation with Deep Q-Networks and Policy Gradient Methods **Analysis Date:** November 25, 2025 **Total Files Analyzed:** 14 Python files + 3 configuration files

Executive Summary

This project implements multiple Deep Reinforcement Learning algorithms to solve the **Memory Maze** environment - a challenging 3D navigation task requiring long-term memory. The codebase consists of **two parallel implementations** by different students:

1. **Main Folder Implementation** - Focus on memory architectures (VisualDQN, DRQN, GraphDQN)
2. **Yinsu Folder Implementation** - Production-quality PPO and DQN with advanced features

Overall Quality: Mixed - Ranges from educational prototypes (C+) to production-ready code (A-)

Part I: Project Overview

1.1 The Memory Maze Challenge

Environment: DeepMind's Memory Maze (2022 research environment)

Task Description: - Agent navigates randomly generated 3D mazes - Must collect colored objects in sequence - Border color indicates current target - Correct target: +1 reward, new target assigned - Episode length: 1,000 steps (9×9 maze) - **Critical challenge:** Must remember object locations over time

Why Memory Matters:

Without Memory (Random Agent):

- Finds each target by chance
- Expected reward: 5-10 targets per episode
- Inefficient wandering

With Memory (LSTM/Graph):

- Explores once, builds mental map
- Follows optimal paths
- Expected reward: 20-30+ targets per episode

Environment Variants:

Size	Environment ID	Objects	Steps	Human Mean	Human Max
9x9	MemoryMaze-9x9-v0	3	1,000	26.4	34.8
11x11	MemoryMaze-11x11-v0	4	2,000	44.3	58.0
13x13	MemoryMaze-13x13-v0	5	3,000	55.5	74.5
15x15	MemoryMaze-15x15-v0	6	4,000	67.7	87.7

1.2 Project Structure

```
/home/daisy/Desktop/dl/
    └── Main Folder (Educational Focus)
        ├── agents/
        │   ├── visual_dqn.py          # Baseline CNN-DQN
        │   ├── drqn.py                # LSTM-based DQN
        │   └── graph_dqn.py          # Graph memory (BROKEN)
        ├── configs/
        │   ├── dqn_config.json
        │   └── drqn_config.json
        ├── test_env.py               # Training script
        ├── evaluate_model.py        # Evaluation script
        └── image_generation.py     # Dataset generation

    └── yinsu/3dmaze/ (Production Focus)
        ├── algorithms/
        │   ├── dqn.py                # DQN algorithm (Double + Dueling)
        │   └── ppo.py                # PPO algorithm
        ├── models/
        │   ├── agent.py              # Actor-Critic LSTM (PPO)
        │   └── dqn_agent.py          # DQN model (Dueling support)
        ├── environments/
        │   └── wrappers.py          # Observation preprocessing
        ├── utils/
        │   ├── replay_buffer.py      # Storage for PPO & DQN
        │   └── logger.py             # TensorBoard logging
        ├── configs/
        │   └── config.yaml           # Centralized config
        ├── main_train.py            # Unified training
        └── main_evaluate.py         # Evaluation with videos
```

Part II: Main Folder Implementation

2.1 Agent Architectures

2.1.1 VisualDQN ([agents/visual_dqn.py](#))

Type: Basic Deep Q-Network with CNN vision module

Architecture:

```
Input: (64, 64, 3) RGB image
      ↓
Vision Module:
  Conv2D(3→16, k=3, s=1) → AvgPool(2×2) → ReLU
  Conv2D(16→32, k=3, s=1, p=1) → AvgPool(2×2) → ReLU
  Conv2D(32→32, k=3, s=1, p=1) → ReLU
  Flatten → ~4,000 features
  ↓
Target Color Extraction:
  pixel[0,0] → (R, G, B) - 3 values
  ↓
Concatenate [vision_features, target_color]
  ↓
Head Network:
  Linear(features+3 → 512) → ReLU
  Linear(512 → 128) → ReLU
  Linear(128 → 6) # Q-values for 6 actions
```

Key Innovation (Line 39-41):

```
target_color = x[:, :, 0, 0] # Top-left pixel encodes target
```

The Memory Maze environment cleverly encodes the target color in pixel (0,0), which the network extracts and uses as additional input.

Strengths: - ✓ Simple, understandable baseline - ✓ Smart use of target color hint - ✓ Proper centering of inputs ($x - 0.5$)

Limitations: - ✗ No memory mechanism - ✗ Cannot recall previously seen objects - ✗ Expected performance: 5-10 targets

Training Config:

```
{
    "num_episodes": 8000,
    "batch_size": 8,
    "eps_start": 0.7,
    "eps_decay": 300,
    "gamma": 0.999
}
```

2.1.2 DRQN ([agents/drqn.py](#))

Type: Deep Recurrent Q-Network with LSTM memory

Architecture:

```
Input: (64, 64, 3) RGB image
      ↓
Vision Module:
  Conv2D(3→16, k=3, p=1) + BatchNorm2d + ReLU
  Conv2D(16→32, k=3, p=1) + BatchNorm2d + AvgPool + ReLU
  Conv2D(32→32, k=3, p=1) + BatchNorm2d
  Flatten → 16,384 features
  ↓
  Concatenate [vision_features, target_color (3)]
  Input size: 16,387
  ↓
LSTM:
  Hidden size: 256
  Layers: 3
  Dropout: 0.3
  ↓
Linear Output:
  LSTM_out(256) → 6 Q-values
```

Memory Management (Lines 46-49):

```
def reset_memory(self):
    """Reset LSTM states at episode start"""
    self.prev_hidden = torch.zeros(lstm_layers, hidden_size, device=device)
    self.prev_state = torch.zeros(lstm_layers, hidden_size, device=device)
```

Forward Pass with Memory (Lines 51-64):

```
def forward(self, x):
    vision_features = self.vision_module(x)
    features = torch.cat((vision_features, target_color), dim=1)

    # LSTM processes sequence
    lstm_out, (h, c) = self.lstm(features, (self.prev_hidden, self.prev_state))

    # Detach to prevent backprop through time
    self.prev_hidden = h.detach()
    self.prev_state = c.detach()

    return self.linear_out(lstm_out)
```

Why `.detach()` (Line 61-62)? - Prevents gradients from flowing back through entire episode history - Truncated backpropagation through time (TBPTT) - Reduces memory usage and training time

Strengths: - ✓ LSTM can remember maze layout - ✓ BatchNorm for stable training - ✓ Proper memory reset between episodes - ✓ 3-layer deep LSTM for complex patterns

Expected Performance: 15-25 targets per episode

Training Config:

```
{  
    "num_episodes": 2000, // Fewer needed due to memory  
    "batch_size": 128,  
    "eps_start": 0.9,  
    "eps_decay": 200,  
    "lstm_hidden_size": 256,  
    "lstm_num_layers": 3,  
    "dropout": 0.3  
}
```

2.1.3 GraphDQN ([agents/graph_dqn.py](#))

Type: Graph-based spatial memory (INCOMPLETE - BROKEN)

Intended Concept: Build a graph where: - **Nodes:** Encoded observations (128 nodes × 8 features each) - **Edges:** Action transitions between states (6 action types) - **Similarity:** Track which states are related

Architecture (Lines 28-31):

```
self.mem_graph_nodes = torch.zeros((batch_size, 128, 8))  
self.mem_graph_edges = torch.zeros((batch_size, 6, 128, 128))  
self.mem_similarity = torch.full((batch_size, 128, 128), float('inf'))
```

Critical Bugs:

1. Line 118 - Undefined Variable:

```
mem_flat = self.memory.view(...) # ✗ self.memory doesn't exist
```

1. Line 119 - Missing Encoder:

```
mem_encoded = self.mem_encoder(mem_flat) # ✗ self.mem_encoder never created
```

1. Line 92 - Incomplete Logic:

```
# TODO add edge from previous node to new node
```

1. Design Flaw: The `add_mem_state()` method (lines 60-99) merges similar states but never actually adds edges, defeating the purpose of a graph.

Status:  NOT FUNCTIONAL - Will crash if used

What It SHOULD Do: 1. Encode each observation into compact representation 2. Store as graph node 3. Add edge when taking action from one state to another 4. Use graph structure for path planning 5. Query graph to find shortest paths to targets

Potential if Fixed: - Could outperform LSTM on large mazes - Explicit spatial reasoning vs implicit LSTM memory - More interpretable (can visualize graph)

2.2 Training Script (`test_env.py`)

2.2.1 Core Components

File: 324 lines of DQN training logic

Replay Buffer (Lines 57-69):

```
class ReplayBuffer:
    def __init__(self, capacity):
        self.memory = deque([], maxlen=capacity) # Circular buffer

    def push(self, *args):
        self.memory.append(Transition(*args))

    def sample(self, batch_size):
        return random.sample(self.memory, batch_size)
```

Epsilon-Greedy Exploration (Lines 84-108):

```
def select_action(states, step, policy_net, env, device):
    eps = EPS_END + (EPS_START - EPS_END) * np.exp(-step / EPS_DECAY)

    exploit_mask = sample > eps # Take best action
    explore_mask = ~exploit_mask # Random action

    actions[exploit_mask] = policy_net(states).argmax()
    actions[explore_mask] = env.action_space.sample()

    return actions
```

Vectorized: Handles 8 parallel environments simultaneously

Model Optimization (Lines 110-149):

```

def optimize_model(memory, policy_net, optimizer, criterion, device):
    # Sample batch
    batch = Transition(*zip(*transitions))

    # Compute Q(s, a)
    state_action_values = policy_net(states).gather(1, actions)

    # Compute target: r + γ max Q(s', a')
    with torch.no_grad():
        next_state_values = target_net(next_states).max(1).values
    expected_values = rewards + GAMMA * next_state_values

    # Loss and optimization
    loss = F.smooth_l1_loss(state_action_values, expected_values)
    optimizer.zero_grad()
    loss.backward()
    torch.nn.utils.clip_grad_value_(policy_net.parameters(), 100)
    optimizer.step()

```

Training Loop (Lines 151-214):

```

For each episode:
    Reset environment and agent memory

    For each step (1000):
        1. Select action (epsilon-greedy)
        2. Execute in all 8 environments
        3. Store transitions in replay buffer
        4. Optimize policy network
        5. Update progress bar

    Update target network (copy policy weights)
    Save episode rewards

```

2.2.2 Critical Bug (Line 137)

```
next_state_values = target_net(next_states).max(1).values
```

Problem: `target_net` is not in function scope - it's a global variable. Should be passed as parameter.

Impact: Code works but has poor design - hard to test and refactor.

2.3 Evaluation Script ([evaluate_model.py](#))

Purpose: Test trained models without training updates

Key Differences from Training: - No epsilon-greedy (always greedy) - No gradient updates - No replay buffer - Optional image saving

Issues Found:

1. Line 120 - Config Ignored:

```
NUM_ENVS = config.get('num_environments', NUM_ENVS) # Line 117
env = gym.make_vec(..., num_envs=1, ...) # Line 120 - Hardcoded!
```

1. Line 138 - Security Risk:

```
torch.load(weight_path, map_location=device)
# Should be: torch.load(..., weights_only=True)
```

1. Line 109 - Wrong Message:

```
print(f"Training on device: {device}") # Should say "Evaluating"
```

2.4 Image Generation ([image_generation.py](#))

Purpose: Generate supervised learning dataset from synthetic mazes

2.4.1 Maze Generation (Lines 11-26)

Algorithm: Recursive Backtracking (DFS-based)

```
def generate_maze(width, height):
    maze = np.ones((height, width)) # All walls

    def carve(x, y):
        maze[y, x] = 0 # Mark as path
        dirs = [(1,0), (-1,0), (0,1), (0,-1)]
        random.shuffle(dirs)

        for dx, dy in dirs:
            nx, ny = x + 2*dx, y + 2*dy # Jump 2 cells
            if maze[ny, nx] == 1: # Unvisited
                maze[y+dy, x+dx] = 0 # Carve path
                carve(nx, ny) # Recurse

    carve(1, 1)
    return maze
```

2.4.2 Ray-Casting Renderer (Lines 31-79)

Technique: First-person 3D rendering from 2D maze

```

def render_first_person(maze, px, py, angle, fov=100, width=200, height=150):
    For each pixel column:
        1. Calculate ray angle
        2. March ray forward in small steps
        3. Check if ray hits wall
        4. Calculate wall height based on distance
        5. Shade based on distance
        6. Render sky + wall + floor

```

Output: 200×150 RGB image simulating 3D view

2.4.3 Dataset Generation (Lines 119-161)

```

For each free cell in maze:
    For each direction (N, S, E, W):
        1. Render first-person view
        2. Check movability: [forward, left, right]
        3. Create sample dict with image + labels
        4. Randomly split to train/val/test (80/10/10%)
        5. Save as pickle file

```

Bug (Line 174):

```
maze = generate_maze(11, 11)b # ❌ Extra 'b' character
```

Part III: Yinsu's Implementation

3.1 Overall Architecture

Design Philosophy: Production-quality, modular, extensible

Key Differentiators: - Unified framework for multiple algorithms - YAML-based configuration - TensorBoard logging - Comprehensive experiment tracking - Clean separation of concerns

3.2 Algorithm Implementations

3.2.1 DQN ([algorithms/dqn.py](#))

Class: 54 lines of clean, professional code

Features:

1. **Double DQN (Lines 28-34):**

```

if self.use_double_dqn:
    # Policy network selects action
    next_actions = self.policy_net(next_obs).argmax(1).unsqueeze(1)
    # Target network evaluates it
    next_q_values = self.target_net(next_obs).gather(1, next_actions)
else:
    # Standard DQN
    next_q_values = self.target_net(next_obs).max(1)[0].unsqueeze(1)

```

Why Double DQN? - Standard DQN overestimates Q-values (max bias) - Double DQN decouples action selection from evaluation - Significantly more stable learning

1. Smooth L1 Loss (Line 40):

```
loss = F.smooth_l1_loss(current_q, target_q)
```

- Also called Huber loss
- Less sensitive to outliers than MSE
- More stable gradients

• Gradient Clipping (Lines 45-47):

```

for param in self.policy_net.parameters():
    if param.grad is not None:
        param.grad.data.clamp_(-1, 1)

```

Mathematical Foundation:

Bellman Equation:

$$Q(s, a) = r + \gamma \max_{\{a'\}} Q(s', a')$$

Double DQN Update:

$$Q(s, a) = r + \gamma Q_{\text{target}}(s', \text{argmax}_{\{a'\}} Q_{\text{policy}}(s', a'))$$

3.2.2 PPO ([algorithms/ppo.py](#))

Class: 63 lines of textbook PPO implementation

Core Algorithm:

1. Advantage Normalization (Lines 22-24):

```

advantages = rollouts.returns - rollouts.value_preds
advantages = (advantages - advantages.mean()) / (advantages.std() + 1e-5)

```

- Standard practice in PPO
 - Reduces variance, improves stability
- Clipped Surrogate Objective (Lines 40-45):**

```

ratio = exp(new_log_prob - old_log_prob)
surrl1 = ratio * advantage
surrl2 = clip(ratio, 1-ε, 1+ε) * advantage
policy_loss = -min(surrl1, surrl2).mean()

```

Why Clipping? - Prevents too large policy updates - Maintains trust region - More stable than vanilla policy gradient

1. Multi-epoch Updates (Lines 30-60):

```

for epoch in range(ppo_epochs): # 4 epochs
    for mini_batch in batch_generator: # 4 mini-batches
        # Compute losses
        # Update network

```

1. Combined Loss (Line 51):

```
Total = policy_loss + 0.5×value_loss - 0.01×entropy
```

- Policy loss: Improve action selection
- Value loss: Improve state value prediction
- Entropy: Encourage exploration

Mathematical Foundation:

```

PPO Objective:
L^CLIP(θ) = E_t[min(r_t(θ)A_t, clip(r_t(θ), 1-ε, 1+ε)A_t)]

where:
r_t(θ) = π_θ(a_t|s_t) / π_θ_old(a_t|s_t) # Probability ratio
A_t = advantage estimate
ε = 0.1 # Clip parameter

```

3.3 Model Architectures

3.3.1 Actor-Critic RNN ([models/agent.py](#))

Purpose: PPO agent with LSTM memory

Architecture:

```

Input: (84, 84, 4) - 4 grayscale frames stacked
      ↓
CNN Base:
  Conv2D(4→32, k=8, s=4) # Reduces 84×84 to 20×20
  ReLU
  Conv2D(32→64, k=4, s=2) # Reduces to 9×9
  ReLU
  Conv2D(64→64, k=3, s=1) # Reduces to 7×7
  ReLU
  Flatten → 3,136 features
  ↓
LSTMCell(3136 → 256 hidden)
  ↓
    ↗ Actor: Linear(256 → num_actions)
    ↗ Output: Categorical(logits)

    ↗ Critic: Linear(256 → 1)
    ↗ Output: V(s)

```

Key Design Choices:

1. **Atari-DQN CNN Architecture:**
2. Standard from "Playing Atari with Deep RL" (Mnih et al., 2013)
3. Proven effective for visual RL
4. Kernel sizes (8, 4, 3) with strides (4, 2, 1)
5. **LSTMCell vs LSTM:**
6. LSTMCell: Manual state management (more control)
7. LSTM: Automatic state handling
8. Here: LSTMCell for explicit hidden state passing
9. **Shared CNN Base:**
10. Same vision features for actor and critic
11. More efficient than separate networks
12. Enables shared representations

Forward Pass (Lines 38-53):

```

def forward(self, obs, hidden, cell):
    # Normalize and permute
    obs = obs.permute(0, 3, 1, 2).float() / 255.0

    # Extract visual features
    features = self.cnn_base(obs)

    # Update recurrent state
    hidden, cell = self.lstm(features, (hidden, cell))

    # Actor output
    action_dist = Categorical(logits=self.actor_head(hidden))

    # Critic output
    value = self.critic_head(hidden)

    return action_dist, value, (hidden, cell)

```

3.3.2 DQN Agent ([models/dqn_agent.py](#))

Two Modes: Standard DQN or Dueling DQN

Standard Architecture:

```

CNN Base (same as PPO) → 3,136 features
↓
Linear(3136 → 512) → ReLU
Linear(512 → num_actions)

```

Dueling Architecture (Lines 29-40):

```

CNN Base → 3,136 features
↓
Value Stream:
  Linear(3136 → 512) → ReLU → Linear(512 → 1)
  Outputs: V(s)
Advantage Stream:
  Linear(3136 → 512) → ReLU → Linear(512 → num_actions)
  Outputs: A(s,a) for each action

Combine (Line 60):
Q(s,a) = V(s) + [A(s,a) - mean(A(s,a))]

```

Why Dueling? (Wang et al., 2016)

1. **Separates state value from action advantage:**
2. $V(s)$: How good is this state?
3. $A(s,a)$: How much better is action a ?
4. **Better generalization:**
5. Many states have similar values
6. Value stream learns state quality

7. Advantage stream focuses on action differences

8. Faster learning:

9. Value updates benefit all actions

10. More sample efficient

Example:

```
State: Middle of empty corridor
V(s) = 10 (decent state)
```

Actions:	Forward	Left	Right	Back
A(s,a):	+0.5	-0.2	-0.2	-1.0
Q(s,a) = V + A:	10.5	9.8	9.8	9.0

Best action: Forward
But network learns $V(s)$ applies to all actions

3.4 Utilities

3.4.1 Replay Buffers ([utils/replay_buffer.py](#))

Two Classes for Two Paradigms:

RolloutStorage (PPO - On-Policy)

Storage:

```
obs: (num_steps+1, num_processes, C, H, W) # uint8
actions: (num_steps, num_processes, 1)
rewards: (num_steps, num_processes, 1)
value_preds: (num_steps+1, num_processes, 1)
returns: (num_steps+1, num_processes, 1)
recurrent_hidden_states: (num_steps+1, num_processes, hidden_dim)
recurrent_cell_states: (num_steps+1, num_processes, hidden_dim)
masks: (num_steps+1, num_processes, 1)
```

GAE Computation (Lines 56-68):

```
def compute_returns(self, next_value, use_gae, gamma, gae_lambda):
    if use_gae:
        gae = 0
        for step in reversed(range(num_steps)):
            δ = rewards[step] + γ*values[step+1]*mask - values[step]
            gae = δ + γ*λ*mask*gae
            returns[step] = gae + values[step]
```

GAE (Generalized Advantage Estimation):

```
Advantage = Σ(γλ)^t δ_t  
where δ_t = r_t + γV(s_{t+1}) - V(s_t)
```

Benefits:

- Reduces variance (vs Monte Carlo)
- Reduces bias (vs TD)
- λ controls bias-variance tradeoff

Batch Generator (Lines 70-90): - Flattens all rollouts - Shuffles randomly - Yields mini-batches for PPO epochs

ReplayBuffer (DQN - Off-Policy)

Efficient Circular Buffer:

```
class ReplayBuffer:  
    def __init__(self, buffer_size, obs_shape, device):  
        # Preallocate numpy arrays  
        self.obs = np.zeros((buffer_size, *obs_shape), dtype=np.uint8)  
        self.actions = np.zeros((buffer_size, 1), dtype=np.int64)  
        self.rewards = np.zeros((buffer_size, 1), dtype=np.float32)  
        self.next_obs = np.zeros((buffer_size, *obs_shape), dtype=np.uint8)  
        self.dones = np.zeros((buffer_size, 1), dtype=np.float32)  
  
        self.ptr = 0 # Current position  
        self.size = 0 # Current size
```

Memory Efficiency:

```
100,000 transitions:  
- obs: 100k × 84 × 84 × 4 × 1 byte = 2.8 GB  
- actions: 100k × 1 × 8 bytes = 0.8 MB  
- rewards: 100k × 1 × 4 bytes = 0.4 MB  
- Total: ~2.8 GB
```

```
If using float32 for obs: 11.2 GB!  
uint8 saves 75% memory
```

3.4.2 Logger ([utils/logger.py](#))

Design Pattern: Abstract factory with inheritance

Base Class:

```
class BaseLogger:  
    def log_scalar(self, tag, value, step): pass  
    def log_hyperparams(self, params): pass  
    def close(self): pass
```

TensorBoardLogger:

```

class TensorBoardLogger(BaseLogger):
    def __init__(self, log_dir):
        self.writer = SummaryWriter(log_dir)

    def log_scalar(self, tag, value, step):
        self.writer.add_scalar(tag, value, step)

    def log_hyperparams(self, hparams):
        # Save to text file for reproducibility
        with open(f"{log_dir}/hyperparams.txt", 'w') as f:
            for k, v in hparams.items():
                f.write(f"{k}: {v}\n")

```

Usage:

```

logger = get_logger("tensorboard", "logs/exp1")
logger.log_scalar("reward", 10.5, step=1000)
logger.log_hyperparams(config)

```

TensorBoard Visualization:

```

tensorboard --logdir logs/
# View at http://localhost:6006

```

3.5 Environment Wrappers

3.5.1 Pipeline Overview

Raw Observation → GrayScale → Resize → FrameStack → Model

```

MemoryMaze: (64, 64, 3) RGB
    ↓ GrayScaleObservation
(64, 64, 1) Grayscale
    ↓ ResizeObservation(84, 84)
(84, 84, 1) Resized
    ↓ FrameStack(k=4)
(84, 84, 4) Stacked frames
    ↓ Model Input

```

3.5.2 GrayScaleObservation (Lines 13-26)

```

class GrayScaleObservation(gym.ObservationWrapper):
    def observation(self, obs):
        gray = cv2.cvtColor(obs, cv2.COLOR_RGB2GRAY)
        return np.expand_dims(gray, axis=-1)

```

Benefits: - 3x less data to process - 3x less memory - Faster training - Color rarely needed for navigation

3.5.3 ResizeObservation (Lines 29-50)

```
class ResizeObservation(gym.ObservationWrapper):
    def observation(self, obs):
        return cv2.resize(obs, (84, 84), interpolation=cv2.INTER_AREA)
```

Why 84x84? 1. Standard from Atari DQN paper 2. Divisible by strides (4, 2, 1) 3. Good balance of detail vs computation 4. After Conv2D($k=8, s=4$): 20x20 5. After Conv2D($k=4, s=2$): 9x9 6. After Conv2D($k=3, s=1$): 7x7

3.5.4 FrameStack (Lines 53-85)

```
class FrameStack(gym.ObservationWrapper):
    def __init__(self, env, k=4):
        self.k = k
        self.frames = deque([], maxlen=k)

    def reset(self):
        obs = self.env.reset()
        for _ in range(self.k):
            self.frames.append(obs)
        return self._get_observation()

    def _get_observation(self):
        return np.concatenate(list(self.frames), axis=2)
```

Why Stack Frames?

1. **Velocity Information:** `` Frame t-3: Object at position (10, 20) Frame t-2: Object at position (12, 20) Frame t-1: Object at position (14, 20) Frame t: Object at position (16, 20)

Inference: Object moving right at 2 pixels/frame ````

1. Markov Property:

2. Single frame: Non-Markov (can't infer velocity)
3. Stacked frames: Approximately Markov
4. Network can learn motion patterns

5. Standard Practice:

6. Atari DQN: 4 frames
7. Proven effective for visual RL

3.6 Main Training Script

3.6.1 Unified Framework ([main_train.py](#))

Key Innovation: Single script handles both PPO and DQN

Command-Line Interface:

```

# Train DQN only
python main_train.py --algorithm DQN

# Train PPO only
python main_train.py --algorithm PPO

# Train both sequentially
python main_train.py --algorithm ALL

```

3.6.2 PPO Training Loop (Lines 99-195)

High-Level Flow:

```

def train_ppo(config, device, experiment_dir):
    # 1. Setup
    envs = create_vectorized_env(num_envs=8)
    agent = ActorCriticRNN(...)
    ppo_algo = PPO(agent, config)
    rollouts = RolloutStorage(...)

    # 2. Training loop
    for update in range(num_updates):
        # Collect rollouts (128 steps × 8 envs)
        for step in range(128):
            action_dist, value, new_hidden = agent(obs, hidden, cell)
            action = action_dist.sample()
            obs, reward, done, info = envs.step(action)
            rollouts.insert(obs, hidden, action, log_prob, value, reward, mask)

        # Compute returns with GAE
        rollouts.compute_returns(next_value, use_gae=True, ...)

        # Update policy (4 epochs × 4 mini-batches)
        value_loss, policy_loss, entropy = ppo_algo.update(rollouts)

        # Logging and checkpointing
        if update % log_interval == 0:
            logger.log_scalar("reward", mean_reward, step)
            logger.log_scalar("value_loss", value_loss, step)

```

Parallel Environments:

```

8 processes collecting data simultaneously:

Env 0: [s0, a0, r0, s1] → [s1, a1, r1, s2] → ...
Env 1: [s0, a0, r0, s1] → [s1, a1, r1, s2] → ...
...
Env 7: [s0, a0, r0, s1] → [s1, a1, r1, s2] → ...

After 128 steps: 128 × 8 = 1,024 transitions collected

```

3.6.3 DQN Training Loop (Lines 197-277)

High-Level Flow:

```

def train_dqn(config, device, experiment_dir):
    # 1. Setup
    env = create_env() # Single environment
    policy_net = DQNAgent(use_dueling=True)
    target_net = DQNAgent(use_dueling=True)
    target_net.load_state_dict(policy_net.state_dict())
    replay_buffer = ReplayBuffer(100000, ...)

    # 2. Training loop
    for step in range(10_000_000):
        # Epsilon-greedy action selection
        epsilon = calculate_epsilon(step)
        if random.random() < epsilon:
            action = env.action_space.sample()
        else:
            action = policy_net(obs).argmax()

        # Execute and store
        next_obs, reward, done, info = env.step(action)
        replay_buffer.push(obs, action, reward, next_obs, done)

        # Learn from replay buffer
        if step > learning_starts and step % train_freq == 0:
            batch = replay_buffer.sample(32)
            loss = dqn_algo.update(batch)

        # Update target network
        if step % target_update_freq == 0:
            dqn_algo.update_target_network()

```

Epsilon Decay (Lines 279-286):

```

def calculate_epsilon(step):
    start, end, decay = 1.0, 0.01, 200000
    fraction = min(1.0, step / decay)
    return start + fraction * (end - start)

```

Linear decay: $1.0 \rightarrow 0.01$ over 200k steps

3.7 Configuration Management

File: `configs/config.yaml` (141 lines)

Complete Configuration:

```

algorithm_to_run: "DQN" # or "PPO"

env:
  name: "MemoryMaze-9x9-v0"
  wrappers:
    grayscale: {enabled: true}
    resize: {enabled: true, shape: [84, 84]}
    frame_stack: {enabled: true, k: 4}

model:
  recurrent_hidden_dim: 256

ppo:
  learning_rate: 0.00025
  gamma: 0.99
  use_gae: true
  gae_lambda: 0.95
  clip_epsilon: 0.1
  ppo_epochs: 4
  num_mini_batches: 4
  value_loss_coef: 0.5
  entropy_coef: 0.01
  max_grad_norm: 0.5
  num_processes: 8
  num_steps_per_update: 128

dqn:
  dueling: true
  double: true
  learning_rate: 0.0001
  gamma: 0.99
  buffer_size: 100000
  batch_size: 32
  train_frequency: 4
  target_update_frequency: 1000
  learning_starts: 10000
  epsilon_start: 1.0
  epsilon_end: 0.01
  epsilon_decay_steps: 200000

training:
  device: "auto"
  seed: 42
  num_env_steps: 10000000
  output_dir: "experiment_results/"
  logger_type: "tensorboard"
  log_interval_steps: 10000
  save_interval_steps: 100000

live_visualization:
  enabled: false
  eval_interval_steps: 50000

evaluation:
  record_video: true
  video_folder: "videos"
  num_episodes: 10

```

Advantages: - ✓ Single source of truth - ✓ Easy to modify hyperparameters - ✓ Experiment reproducibility - ✓ No hardcoded values

Part IV: Comparative Analysis

4.1 Side-by-Side Comparison

Feature	Main Folder	Yinsu Folder
Code Quality	C+	A-
Lines of Code	~500 (training)	~350 (training)
Algorithms	DQN, DRQN, GraphDQN*	PPO, DQN (Double+Dueling)
Architecture	Script-based	Modular packages
Configuration	JSON (partial)	YAML (comprehensive)
Logging	Print statements	TensorBoard
Experiment Tracking	Manual	Automated timestamps
Checkpointing	Basic	Robust + recovery
Error Handling	Minimal	Better
Documentation	Comments	Docstrings + comments
Vectorization	8 parallel envs (DQN)	8 parallel (PPO only)
Memory Efficiency	Good (uint8)	Excellent (uint8 + preallocated)
Evaluation	Basic metrics	Comprehensive stats + video
Code Reusability	Low	High
Maintainability	Medium	High
Extensibility	Low	High

*GraphDQN is broken

4.2 Algorithm Comparison

4.2.1 Feature Matrix

Algorithm	Memory	On/Off-Policy	Sample Efficiency	Stability	Complexity
VisualDQN	None	Off-policy	Medium	Medium	Low
DRQN	LSTM	Off-policy	High	Medium	Medium
GraphDQN*	Graph	Off-policy	Highest	?	High
PPO (Yinsu)	LSTM	On-policy	Low	High	High
DQN (Yinsu)	None	Off-policy	High	High	Medium

*Broken implementation

4.2.2 Expected Performance (9x9 Maze)

Algorithm	Training Time	Final Reward	Convergence	Memory Usage
VisualDQN	24-36h	8-12	Slow	3 GB
DRQN	18-30h	18-25	Medium	3 GB
PPO (Yinsu)	12-24h	20-28	Fast	2 GB
DQN-Dueling (Yinsu)	24-40h	22-30	Medium	3 GB

4.3 Code Organization

Main Folder:

Pros:

- + Focuses on different memory mechanisms
- + Educational - easy to understand individual components
- + Good vectorization support

Cons:

- Mixed concerns (training + model in one file)
- Hardcoded hyperparameters
- Poor experiment management
- GraphDQN broken
- No production features

Yinsu Folder:

Pros:

- + Clean separation of concerns
- + Production-ready features
- + Comprehensive configuration
- + Two state-of-the-art algorithms
- + Excellent logging and tracking
- + Extensible architecture

Cons:

- 239 lines of commented code
- Empty README
- No input validation
- Assumes specific gym version

Part V: Comprehensive Code Issues

5.1 Critical Bugs

Main Folder:

1. GraphDQN Completely Broken

2. Location: `agents/graph_dqn.py`
3. Lines 118-119: Undefined variables (`self.memory`, `self.mem_encoder`)
4. Line 92: Incomplete logic (TODO comment)
5. Impact:  Will crash if used

6. Global Variable Abuse

7. Location: `test_env.py:137`
8. `target_net` not in function scope
9. Impact:  Poor design, hard to test

10. Config Ignored

11. Location: `evaluate_model.py:120`
12. NUM_ENVS read from config but hardcoded to 1
13. Impact:  Misleading configuration

14. Syntax Error

15. Location: `image_generation.py:174`
16. `maze = generate_maze(11, 11)b` - extra 'b'
17. Impact:  Script won't run

Yinsu Folder:

1. Commented Code Bloat

2. Location: `main_evaluate.py:1-239`

3. 239 lines of old code versions

4. Impact:  Poor code hygiene

5. Config Path Hardcoded

6. Location: `main_evaluate.py:483`

7. `configs/config.yaml` hardcoded

8. Impact:  Less flexible

5.2 Security Issues

1. Unsafe Pickle Loading

2. Both folders

3. `torch.load()` without `weights_only=True`

4. Impact:  Arbitrary code execution risk

5.3 Code Quality Issues

Main Folder:

1. **No error handling** - Missing try-except blocks

2. **Unused imports** - `Transition` in `evaluate_model.py`

3. **Misleading messages** - "Training" instead of "Evaluating"

4. **No validation** - Config parameters not checked

5. **Inconsistent naming** - `test_env.py` does training

Yinsu Folder:

1. **No README** - Empty file

2. **No version checks** - Assumes gym 0.23.1

3. **No input validation** - Config not validated

4. **Large commented sections** - Should use git

5.4 Performance Issues

1. Image Saving in Loop

2. Location: `test_env.py:172-175`

3. Saves image every step when enabled
 4. Impact: ⚠️ Massive slowdown
- 5. No Memory Profiling**
6. Both folders
 7. No documentation of RAM requirements
 8. Impact: ⚠️ May OOM on smaller machines
-

Part VI: Deep Learning Concepts Applied

6.1 Q-Learning Foundation

Bellman Optimality Equation:

$$Q^*(s, a) = E[r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

where:

$Q^*(s, a)$ = optimal action-value function
 r = immediate reward
 γ = discount factor (0.999)
 s' = next state
 a' = next action

DQN Approximation:

$$Q(s, a; \theta) \approx Q^*(s, a)$$

where θ = neural network parameters

Loss Function:

$$L(\theta) = E[(r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta))^2]$$

where θ^- = target network parameters

6.2 Policy Gradient Foundation

Objective:

$$J(\theta) = E_{\{\pi_\theta\}}[\sum \gamma^t r_t]$$

Maximize expected cumulative reward

REINFORCE Gradient:

$$\nabla_\theta J(\theta) = E_\pi[\nabla_\theta \log \pi_\theta(a|s) Q^\pi(s,a)]$$

PPO Clipped Objective:

$$L^{CLIP}(\theta) = E_t[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) \hat{A}_t)]$$

where:

$$r_t(\theta) = \pi_\theta(a_t|s_t) / \pi_\theta(a_t|s_t)$$

\hat{A}_t = advantage estimate

$$\epsilon = 0.1$$

6.3 Convolutional Feature Learning

Receptive Field Growth:

Layer 0 (Input): 84x84
 Layer 1 (Conv k=8, s=4): 20x20, receptive field = 8x8
 Layer 2 (Conv k=4, s=2): 9x9, receptive field = 18x18
 Layer 3 (Conv k=3, s=1): 7x7, receptive field = 26x26

Feature Hierarchy: - Layer 1: Edges, textures - Layer 2: Shapes, patterns - Layer 3: Objects, spatial relationships

6.4 Recurrent Memory Mechanisms

LSTM Equations:

Forget gate: $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$
 Input gate: $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$
 Output gate: $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$

Candidate: $\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$
 Cell state: $C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$
 Hidden state: $h_t = o_t \odot \tanh(C_t)$

Why LSTM for Memory Maze? - Can remember over 1000 timesteps - Selective forgetting (forget gate) - Protected memory (cell state) - Gradient flow (no vanishing gradient)

6.5 Exploration Strategies

Epsilon-Greedy:

```
a = {  
    argmax_a Q(s,a)      with probability 1-ε  
    random action        with probability ε  
}
```

Epsilon Decay:

Linear: $\epsilon_t = \epsilon_{\text{start}} + (\epsilon_{\text{end}} - \epsilon_{\text{start}}) \times (t / \text{decay_steps})$
Exponential: $\epsilon_t = \epsilon_{\text{end}} + (\epsilon_{\text{start}} - \epsilon_{\text{end}}) \times \exp(-t / \text{decay_steps})$

PPO Entropy Bonus:

Entropy = $-\sum \pi(a|s) \log \pi(a|s)$

High entropy → More exploration
Low entropy → More exploitation

Loss includes: $-\alpha \times \text{Entropy}$ ($\alpha = 0.01$)

Part VII: Training Details

7.1 Hyperparameter Analysis

DQN (Main Folder - VisualDQN)

Batch Size: 8

- Very small for DQN
- Higher variance gradients
- Faster iteration but less stable

Epsilon Decay: 300 steps

- VERY fast decay
- Reaches near-greedy quickly
- May underexplore

Episodes: 8000

- $8000 \times 1000 = 8M$ steps
- Reasonable for 9x9 maze

Gamma: 0.999

- Very high discount
- Values long-term rewards
- Appropriate for 1000-step episodes

DQN (Main Folder - DRQN)

Batch Size: 128

- Standard size
- More stable gradients
- Better for LSTM

Epsilon: 0.9 → 0.05 over 200 steps

- Faster decay than VisualDQN
- LSTM helps exploration

Episodes: 2000

- Fewer needed (better sample efficiency)
- LSTM learns faster

LSTM: 256 hidden, 3 layers, 0.3 dropout

- Large capacity
- Deep for complex patterns
- Dropout prevents overfitting

DQN (Yinsu - Dueling + Double)

Buffer Size: 100,000

- Standard for Atari
- ~2.8 GB memory

Batch Size: 32

- Standard
- Good balance

Epsilon: 1.0 → 0.01 over 200k steps

- MUCH slower than main folder
- Better exploration
- More appropriate

Learning Starts: 10,000

- Collects diverse data first
- Better initial replay buffer

Train Frequency: 4

- Update every 4 steps
- More sample efficient

Target Update: Every 1000 steps

- Less frequent = more stable
- Standard practice

PPO (Yinsu)

Processes: 8

- Parallel data collection
- 8x faster sampling

Steps per Update: 128

- $128 \times 8 = 1024$ transitions
- Substantial data per update

PPO Epochs: 4

- Reuses data 4 times
- Sample efficient

Mini-batches: 4

- $1024 / 4 = 256$ per batch
- Multiple gradient updates

Clip Epsilon: 0.1

- Standard PPO value
- Prevents large updates

GAE Lambda: 0.95

- Standard value
- Good bias-variance tradeoff

Entropy Coefficient: 0.01

- Moderate exploration bonus

7.2 Expected Training Curves

VisualDQN

Episode 0-1000: Reward ~5 (random)

- Pure exploration
- Learning basic navigation

Episode 1000-3000: Reward 5→10

- Epsilon dropped to ~0.05
- Starts following walls
- Occasional lucky finds

Episode 3000-8000: Reward 10-12

- Greedy policy
- No memory of past sightings
- Performance plateau

DRQN

Episode 0-500: Reward ~5 (random)

- Collecting diverse experiences

Episode 500-1000: Reward 5→15

- LSTM learning patterns
- Starts remembering corridors

Episode 1000-1500: Reward 15→22

- Memory kicking in
- Recalls object locations

Episode 1500-2000: Reward 22→25

- Near-optimal for memory capacity
- Stable performance

PPO (Yinsu)

Updates 0-100: Reward ~5

- Random policy initialization

Updates 100-500: Reward 5→20

- Rapid improvement
- LSTM + parallel envs

Updates 500-1000: Reward 20→26

- Refinement
- Near-human performance

Updates 1000+: Reward 26→28

- Convergence
- Stable plateau

DQN Dueling + Double (Yinsu)

```
Steps 0-10k: Reward ~5
- Learning_starts period
- Pure exploration

Steps 10k-100k: Reward 5→15
- Network learning
- Epsilon still high (>0.5)

Steps 100k-500k: Reward 15→25
- Dueling architecture helping
- Double DQN reducing overestimation

Steps 500k-1M: Reward 25→28
- Fine-tuning
- Near-optimal
```

Part VIII: Recommendations

8.1 Critical Fixes (High Priority)

Main Folder:

1. **Fix or Remove GraphDQN** ```python # Option 1: Fix self.mem_encoder =
 nn.Linear(mem_features*mem_size, hidden_dim)
Option 2: Remove rm agents/graph_dqn.py ````
1. **Fix target_net Scoping** python def optimize_model(..., target_net): # Add parameter ...
2. **Use Safe Model Loading** python torch.load(path, map_location=device, weights_only=True)
3. **Fix Syntax Error** python maze = generate_maze(11, 11) # Remove 'b'

Yinsu Folder:

1. Clean	Up	Commented	Code
bash # Remove lines 1-239 from main_evaluate.py # Remove lines 1-70 from config.yaml			
2. Add README.md	```markdown # Memory Maze Deep RL		
## Installation pip install -r requirements.txt			
## Training	python main_train.py --algorithm DQN		
## Evaluation	python main_evaluate.py --exp-dir results/... ````		
1. Add Config Validation	python def validate_config(config): required = ['env', 'training', 'algorithm_to_run'] for key in required: assert key in config, f"Missing {key}"		

8.2 Performance Improvements (Medium Priority)

1. Implement Prioritized Experience Replay

2. Importance sampling
3. Better sample efficiency
4. Standard in modern DQN

5. Add n-step Returns

6. Better credit assignment
7. Faster learning
8. Common in DQN variants

9. Implement Noisy Networks

- Replace epsilon-greedy
- Parameter space noise
- Often superior exploration

10. Add Distributed Training

- Ape-X DQN
- DD-PPO
- Much faster training

8.3 Code Quality (Low Priority)

1. **Add Type Hints** `python def select_action(states: torch.Tensor, step: int, policy_net: nn.Module) -> torch.Tensor:`
2. **Add Unit Tests** `python def test_replay_buffer(): buffer = ReplayBuffer(100, (84, 84, 4), 'cpu') buffer.push(obs, 0, 1.0, next_obs, False) assert len(buffer) == 1`
3. **Add Logging Framework** `python import logging logger = logging.getLogger(__name__) logger.info(f"Episode {i}: reward={r}")`
4. **Use pathlib** `python from pathlib import Path model_path = Path(output_dir) / "models" / f"agent_{step}.pth"`

8.4 Feature Additions

1. Curiosity-Driven Exploration

- Intrinsic rewards
- Better for sparse rewards
- Memory Maze has sparse rewards

2. Hindsight Experience Replay

- Learn from failures
- Better sample efficiency
- Useful for goal-based tasks

3. Self-Play / Population-Based Training

- Multiple agents
- Diverse strategies
- More robust policies

4. Attention Mechanisms

- Attend to relevant objects
- Better than LSTM for long sequences
- Interpretable

5. Graph Neural Networks

- Complete the GraphDQN idea
 - Explicit spatial reasoning
 - Could outperform LSTM
-

Part IX: Learning Outcomes

9.1 Demonstrated Skills

Main Folder Student(s):

Strengths: - ✓ Understanding of DQN fundamentals - ✓ LSTM integration for memory - ✓ Creative approach (GraphDQN concept) - ✓ Vectorized environment usage - ✓ Smart use of environment features (target color)

Growth Areas: - ⚠ Code organization and modularity - ⚠ Testing and validation - ⚠ Completing complex implementations - ⚠ Production code practices

Grade: B- (Good understanding, needs refinement)

Yinsu:

Strengths: - ✓ Excellent software engineering - ✓ Multiple advanced algorithms - ✓ Production-quality code - ✓ Comprehensive configuration - ✓ Experiment tracking - ✓ Clean architecture

Growth Areas: - ⚠️ Code hygiene (commented code) - ⚠️ Documentation (README) - ⚠️ Input validation

Grade: A- (Excellent, minor issues)

9.2 Deep RL Concepts Covered

Both implementations demonstrate:

1. **Value-Based RL:**

- 2. Q-learning
- 3. Bellman equations
- 4. Temporal difference learning
- 5. Experience replay
- 6. Target networks

7. **Policy-Based RL (Yinsu):**

- 8. Policy gradients
- 9. Actor-Critic
- 10. PPO clipped objective
- 11. Advantage estimation (GAE)

12. **Memory Mechanisms:**

- 13. LSTM for sequential data
- 14. Recurrent state management
- 15. Graph memory (concept)

16. **Neural Network Architectures:**

- 17. CNNs for vision
- 18. Dueling networks
- 19. Actor-Critic split
- 20. Recurrent networks

21. **Training Techniques:**

- 22. Epsilon-greedy exploration
- 23. Gradient clipping
- 24. Batch normalization
- 25. Entropy regularization
- 26. Learning rate scheduling

27. **Engineering Practices:**

- 28. Vectorized environments
- 29. Replay buffers
- 30. Model checkpointing

31. Experiment logging
 32. Configuration management
-

Part X: Conclusion

10.1 Project Summary

This project represents a comprehensive exploration of Deep Reinforcement Learning applied to a challenging memory-based navigation task. The codebase showcases:

Technical Achievements: - Multiple algorithm implementations (DQN, DRQN, PPO) - Advanced features (Dueling, Double DQN, GAE) - Memory mechanisms (LSTM, attempted Graph) - Production-quality engineering (Yinsu)

Educational Value: - Demonstrates evolution from prototype to production - Shows different approaches to same problem - Highlights importance of software engineering in ML - Covers wide range of deep RL concepts

10.2 Overall Assessment

Folder	Code Quality	Completeness	Innovation	Production-Ready
Main	C+	60% (GraphDQN broken)	High (Graph concept)	No
Yinsu	A-	95%	Medium	Yes (with fixes)
Combined	B+	80%	High	Partial

10.3 Best Use Cases

Main Folder Code: - Teaching basic DQN - Demonstrating LSTM memory - Educational examples - Research prototypes

Yinsu Folder Code: - Production deployments - Research baselines - Teaching advanced RL - Starting point for new algorithms

10.4 Final Recommendations

For Students: 1. Study Yinsu's code organization 2. Learn from main folder's creative approaches 3. Fix GraphDQN as learning exercise 4. Add features from recommendation list 5. Run experiments and compare results

For Instructors: 1. Use main folder for teaching basics 2. Use Yinsu folder for best practices 3. Assign GraphDQN fix as project 4. Compare implementations in class 5. Discuss software engineering importance

For Researchers: 1. Build on Yinsu's framework 2. Implement Graph memory properly 3. Try attention mechanisms 4. Scale to larger mazes (13×13, 15×15) 5. Publish results

10.5 Expected Results Summary

9×9 Memory Maze - 1000 Steps:

Algorithm	Expected Reward	Training Time	Memory Usage
Random Baseline	5-7	N/A	N/A
VisualDQN	8-12	24-36h	3 GB
DRQN	18-25	18-30h	3 GB
PPO (Yinsu)	20-28	12-24h	2 GB
DQN-Dueling (Yinsu)	22-30	24-40h	3 GB
Human Baseline	26.4	N/A	N/A
Human Max	34.8	N/A	N/A

Key Insight: Memory (LSTM) is critical - improves performance by 2-3×

10.6 Future Directions

1. Complete GraphDQN Implementation

2. Fix bugs
3. Test on 9×9 maze
4. Compare to LSTM
5. Scale to 15×15

6. Hybrid Approaches

7. Combine Graph + LSTM

8. Attention + LSTM

9. Multi-scale memory

10. Transfer Learning

11. Pre-train on small mazes

12. Fine-tune on large mazes

13. Cross-task transfer

14. Interpretability

15. Visualize LSTM activations

16. Probe what network remembers

17. Attention heatmaps

18. Real-World Applications

19. Robot navigation

20. Autonomous vehicles

21. Game AI

22. Logistics planning

End of Complete Analysis

Appendices

A. File Inventory

Main Folder: - `agents/visual_dqn.py` - 52 lines - `agents/drqn.py` - 67 lines - `agents/graph_dqn.py` - 126 lines (broken) - `test_env.py` - 324 lines - `evaluate_model.py` - 149 lines - `image_generation.py` - 182 lines - `configs/dqn_config.json` - 18 lines - `configs/drqn_config.json` - 23 lines

Yinsu Folder: - `algorithms/dqn.py` - 54 lines - `algorithms/ppo.py` - 63 lines - `models/agent.py` - 57 lines - `models/dqn_agent.py` - 64 lines - `utils/replay_buffer.py` - 126 lines - `utils/logger.py` - 52 lines - `environments/wrappers.py` - 85 lines - `main_train.py` - 348 lines - `main_evaluate.py` - 600 lines (239 commented) - `configs/config.yaml` - 141 lines (70 commented)

Total: ~2,500 lines of student code

B. Dependencies

Core:

- Python 3.8+
- PyTorch 2.0+
- NumPy 1.24+
- Gymnasium 0.23.1

Environment:

- memory-maze
- dm_control
- mujoco

Utilities:

- tensorboard
- opencv-python
- pyyaml
- tqdm

Visualization:

- pygame
- pillow
- imageio

C. Hardware Requirements

Minimum: - CPU: 4 cores - RAM: 8 GB - GPU: Not required (CPU mode works) - Disk: 20 GB

Recommended: - CPU: 8+ cores - RAM: 16 GB - GPU: NVIDIA RTX 2060+ (6GB VRAM) - Disk: 50 GB SSD

Training Time: - CPU only: 3-5 days - GPU (RTX 3080): 12-36 hours

Analysis Date: November 25, 2025 Total Files: 17 Python + 3 Config Total Lines: ~2,500 Analysis

Time: Comprehensive review of all code Analyzer: Claude Code (Sonnet 4.5)