# Hybrid Rule-Based Deep Q-Network Agent for LUX AI

---

**Algorithm 1** Overall Agent Architecture

---

1: **Initialize Components:**
2: Create Q-network (neural network for action evaluation)
3: Create target network (stable copy of Q-network)
4: Create replay buffer (stores past experiences)
5: Set hyperparameters: learning rate, discount factor, exploration rates
6: Initialize tracking systems: relic memory, productive tiles, unit assignments
7:
8: **for** each game **do**
9:    Reset game-specific memory
10:    **for** each match (3 matches per game) **do**
11:      **for** each time step **do**
12:        Observe current state and update relic positions
13:        **for** each active unit **do**
14:          Generate random number $u$ from [0,1]
15:          **if** $u <$ DQN probability **then**
16:            Select action using DQN policy (Algorithm 3)
17:          **else**
18:            Select action using rule-based policy (Algorithm 4)
19:          **end if**
20:          Augment with attack action if appropriate (Algorithm 5)
21:        **end for**
22:        Execute all unit actions
23:        Observe rewards and next state
24:        Update productive tile statistics (Algorithm 6)
25:        Calculate shaped rewards for each unit (Algorithm 7)
26:        Store transitions in replay buffer
27:        **if** training step **then**
28:          Perform network update (Algorithm 8)
29:        **end if**
30:      **end for**
31:      Reset match-specific memory
32:    **end for**
33:    Reset all memories for next game
34: **end for**

---

**Algorithm 2** Neural Network Architecture

1: **Q-Network Structure:**
2: Input layer: Game state features (position, energy, distances, etc.)
3: Hidden layer 1: 256 neurons with ReLU activation
4: Hidden layer 2: 256 neurons with ReLU activation
5: Hidden layer 3: 256 neurons with ReLU activation
6: Hidden layer 4: 128 neurons with ReLU activation
7: Output layer: 5 neurons (Q-values for each action)
8:    Actions: 0=Stay, 1=Up, 2=Right, 3=Down, 4=Left
9:
10: **Target Network:**
11: Same architecture as Q-network
12: Updated periodically (every $N$ training steps) for learning stability

---

**Algorithm 3** DQN Action Selection (Epsilon-Greedy)

**Require:** Current state, Q-network, exploration rate
1:
2: Generate random number $u$ from [0,1]
3: **if** $u <$ exploration rate **then**
4:    **return** random action   *// Exploration*
5: **else**
6:    Compute Q-values for all actions using Q-network
7:    **return** action with highest Q-value   *// Exploitation*
8: **end if**

**Algorithm 4** Rule-Based Policy

---

**Require:** Unit ID, time step, game state, relic positions, productive tiles

1:
2: Extract unit position, energy level, and active status
3: **if** unit inactive OR energy below safety threshold **then**
4:     **return** Stay action
5: **end if**
6:
7: Determine number of active units
8: Calculate unit's rank among active units
9: Compute exploration ratio based on game phase:
10:     Early matches: higher exploration (find new relics)
11:     Late matches: lower exploration (focus on collection)
12:
13: Assign unit role based on rank and exploration ratio
14:
15: **if** unit is COLLECTOR **then**
16:     **Priority 1: Stay on productive tiles**
17:     **if** current position is productive tile **then**
18:         **return** Stay action
19:     **end if**
20:
21:     **Priority 2: Move to known productive tiles**
22:     **if** productive tiles exist **then**
23:         Find nearest productive tile
24:         **return** direction toward nearest productive tile
25:     **end if**
26:
27:     **Priority 3: Systematic grid exploration**
28:     **if** unit has no assignment **then**
29:         Assign unit to position in 5×5 grid around relics
30:     **end if**
31:     Calculate target position from assignment
32:     **return** direction toward assigned grid position
33: **else**
34:     *// Unit is EXPLORER*
35:     **Systematic quadrant search**
36:     **if** unit has no exploration target **then**
37:         Assign quadrant based on unit ID (4 quadrants total)
38:         Generate random target within assigned quadrant
39:     **end if**
40:
41:     **if** target reached OR stuck at current target **then**
42:         Generate new random target in quadrant
43:     **end if**
44:     **return** direction toward exploration target
45: **end if**

---

**Algorithm 5** Attack Action Selection

---

**Require:** Unit position, game state, unit energy, attack range

1:
2: Get opponent positions and ally positions
3: **if** energy insufficient OR no opponents visible **then**
4:     **return** no attack
5: **end if**
6:
7: Initialize best score $= -\infty$ and best target $=$ none
8: **for** each possible target location within attack range **do**
9:     Calculate direct hit score (enemies at exact target)
10:    Calculate splash score (enemies adjacent to target)
11:    Calculate friendly fire penalty (allies at or near target)
12:    Calculate distance penalty (slight preference for closer targets)
13:
14:    Total score $=$ direct $+$ splash $-$ friendly fire $-$ distance
15:
16:    **if** score $>$ best score **then**
17:       Update best score and best target
18:    **end if**
19: **end for**
20:
21: **if** best score $\geq$ minimum threshold **then**
22:    **return** attack toward best target
23: **else**
24:    **return** no attack
25: **end if**

---

**Algorithm 6** Productive Tile Detection

**Require:** Current state, team points, previous points, position statistics

 1: 
 2: Calculate points gained this step = current points − previous points
 3: Find all unit positions near known relics (within distance 3)
 4: 
 5: **for** each position near relics **do**
 6:     Increment occupancy counter for this position
 7: **end for**
 8: 
 9: **if** points gained > 0 AND units are near relics **then**
10:     Distribute gained points among near-relic positions
11:     **for** each position that was occupied **do**
12:         Update total points earned at this position
13:         Calculate confidence = total points / times occupied
14: 
15:         **if** confidence > 0.5 AND occupied at least 3 times **then**
16:             Mark position as "productive tile"
17:         **end if**
18:     **end for**
19: **end if**
20: 
21: **Periodic cleanup (every 100 steps):**
22: Remove tiles from productive set if confidence drops below 0.7

**Algorithm 7** Reward Shaping Function

**Require:** State, action, team reward, productive tiles, statistics

1:
2: Initialize shaped reward = 3.0 × team reward
3: Extract current position and next position from states
4:
5: **Component 1: Productive Tile Bonuses**
6: **if** currently on productive tile **then**
7:     Add large bonus scaled by tile confidence
8:     **if** stayed on same tile (didn't move) **then**
9:         Add additional staying bonus
10:     **end if**
11: **end if**
12: **if** moved to a productive tile **then**
13:     Add bonus for reaching productive tile
14: **end if**
15:
16: **Component 2: Distance Improvement**
17: **if** productive tiles are known **then**
18:     Calculate distance improvement toward nearest productive tile
19:     Add reward proportional to improvement
20: **else**
21:     Calculate distance improvement toward nearest relic
22:     Add smaller reward proportional to improvement
23: **end if**
24:
25: **Component 3: Energy Management**
26: Calculate energy change
27: **if** energy increased **then**
28:     Add small bonus for energy recovery
29: **end if**
30: **if** energy below threshold AND not on productive tile **then**
31:     Add penalty for risky low energy state
32: **end if**
33:
34: **Component 4: Movement Encouragement**
35: **if** unit moved AND not leaving productive tile **then**
36:     Add small bonus to encourage exploration
37: **end if**
38:
39: **Component 5: Survival Bonus**
40: Add small constant bonus for staying alive
41:
42: **Component 6: Leaving Productive Tile Penalty**
43: **if** left productive tile without going to another productive tile **then**
44:     Add penalty scaled by confidence of abandoned tile
45: **end if**
46:
47: **return** shaped reward

**Algorithm 8** Training Procedure (Double DQN)
***
**Require:** Replay buffer, batch size, Q-network, target network
 1:
 2: **if** not enough experiences in buffer **then**
 3:    **return**   *// Skip training*
 4: **end if**
 5:
 6: Sample random batch of experiences from replay buffer
 7: Each experience contains: (state, action, reward, next state, done flag)
 8:
 9: **Step 1: Compute current Q-value predictions**
10: Pass states through Q-network
11: Extract Q-values for actions that were actually taken
12:
13: **Step 2: Compute target Q-values (Double DQN method)**
14: Use Q-network to SELECT best actions for next states
15: Use target network to EVALUATE those actions
16: Calculate targets = reward + discount × evaluated Q-value
17: If episode ended, target = reward only (no future value)
18:
19: **Step 3: Calculate loss and update**
20: Compute mean squared error between predictions and targets
21: Calculate gradients of loss with respect to network parameters
22: Clip gradients to prevent instability
23: Update network parameters using Adam optimizer
24:
25: **Step 4: Decay exploration rate**
26: Reduce exploration rate: new rate = current rate × decay factor
27: Ensure rate doesn't go below minimum value
28:
29: **Step 5: Periodically update target network**
30: **if** training step is multiple of update frequency **then**
31:    Copy Q-network parameters to target network
32: **end if**
***

# Key Concepts Explained

### Hybrid Policy

The agent combines two decision-making strategies:

- **Rule-based policy:** Hand-crafted rules for reliable behavior

- **DQN policy:** Learned through experience via neural networks

- A probability parameter controls which policy is used

### Q-Network

A neural network that estimates the "quality" (Q-value) of each action in a given state. Higher Q-values indicate better expected future rewards.

### Experience Replay

Past experiences are stored in a buffer and randomly sampled during training. This breaks correlations between consecutive experiences and improves learning stability.

### Target Network

A slowly-updated copy of the Q-network used for computing learning targets. This prevents the "moving target" problem and stabilizes training.

### Reward Shaping

The basic team reward is augmented with additional signals to guide learning:

- Bonuses for productive tiles

- Distance improvement rewards

- Energy management signals

- Exploration encouragement

### Productive Tile Detection

Statistical tracking identifies which positions near relics actually generate points. Confidence increases as more evidence accumulates.