Machine Learning Final Paper

Asianna Haughton

---

# Introduction & Objectives

**Objective.**
 Our primary objective was to investigate whether combining classical graph-search techniques with reinforcement learning and biologically inspired heuristics could yield network topologies that outperform pure A* or Dijkstra approaches on multiple fronts—namely, path length, traversal cost, resource collection, and resilience to change. In particular, we set out to:

1. **Leverage A*** as a baseline for heuristic-guided shortest-path search in static environments.

2. **Implement Q-learning** to develop an agent capable of balancing exploration and exploitation when "growing" a network.

3. **Incorporate biological-inspired strategies** (covered in our biological algorithms unit) such as branching and fusion—drawing directly on fungal mycelium growth patterns—to introduce natural redundancy and adaptability.

4. **Compare and analyze** how these three paradigms (A*, RL, and biological heuristics) interact, complement, or conflict in pursuit of multi-objective optimization.

By framing the problem this way, I translated core course concepts into a single cohesive framework—MITE (Mycelium-Inspired Topology Explorer)—and set measurable performance targets against each method's known strengths and limitations.

**Background.**

- **A\*:** A heuristic search that uses an admissible estimate (Manhattan distance) to efficiently guide path selection—ideal for shortest-path problems in static grids.

- **Reinforcement Learning (Q-learning):** Teaches an agent, via trial-and-error and reward feedback, to discover policies that maximize long-term returns—well suited to multi-objective tasks where rewards (e.g., nutrient gain vs. cost) must be balanced dynamically.

- **Biological Algorithms:** Natural systems like fungal mycelia demonstrate powerful decentralized methods for resource foraging and network formation; mechanisms such as branching (for exploration) and fusion (for pruning and reinforcement) inspired our selective fusion variant.

Together, these objectives and background ensured that MITE was more than an ad hoc mix of algorithms—it was a deliberate synthesis of everything we'd studied in class, applied to a challenging, real-world–style pathfinding problem.

## Methodology

Our framework consists of three core components operating on a discrete, 8-connected grid:

1. **Environment & MycelialRouter**

   ○ Each grid cell (node) holds four attributes:

      ■ **Nutrient value** (benefit of occupying the node)

      ■ **Traversal cost** ( expense of moving into the node)

      ■ **Obstacle flag** (impassability)

      ■ **Failure probability** (chance of growth attempt failure)

- The MycelialRouter class tracks active "growth tips" and interfaces with both RL and classical algorithms.

2. **Reinforcement Learning Agent**

   - We employ Q-learning with an ε-greedy policy.

   - **State:** Current tip position + local node attributes

   - **Actions:** Grow up, down, left, or right

   - **Reward function:**

     - +1.0 for reaching the goal

     - +0.5 for entering a high-nutrient node

     - −0.1 for a neutral move

     - −1.0 for hitting an obstacle

   - **Hyperparameters (selected via grid search):**

     - Learning rate $\alpha \in \{0.05, 0.1, 0.2\}$ → chosen 0.1

     - Discount factor $\gamma \in \{0.8, 0.9, 0.99\}$ → chosen 0.9

     - Initial exploration $\varepsilon_0 \in \{0.1, 0.3, 0.5\}$ → chosen 0.3

     - ε-decay $\in \{0.95, 0.99\}$ → chosen 0.99

     - Episodes $\in \{500, 1000, 1500\}$ → chosen 1000

     - Max steps/episode = 100

3. **Classical Baselines**

   - **A\*** with a Manhattan-distance heuristic

○ **Dijkstra's** algorithm using the same cost structure

We first performed a grid-search over the hyperparameter space, evaluating each setting on a fixed start-goal pair for five independent runs and selecting the combination that maximized the **efficiency** metric (nutrients collected per unit cost). Following tuning, we ran three major experiments:

- **Baseline comparison** on static grids

- **Selective fusion** variant, where branching occurs only when beneficial (based on local efficiency gains)

- **Dynamic obstacles**, introducing random node failures midway through a growth episode

---

# Results

## 1. Baseline Comparison (Static Grid)

| Algorithm | Path Length | Total Cost | Total Nutrients | Efficiency (N/C) | Compute Time (s) |
|---|---|---|---|---|---|
| A* | 7 | 2.55 | 4.36 | 1.71 | ~0.0000 |
| Dijkstra | 7 | 2.55 | 4.36 | 1.71 | 0.0005 |
| RL (MITE) | 3 | 0.80 | 1.69 | **2.11** | ~0.0000 |

The RL agent learned to discover significantly shorter, more efficient paths than the classical methods.

## 2. Hyperparameter Search Summary

- **Peak efficiency** of 2.11 achieved at $\alpha=0.1$, $\gamma=0.9$, $\varepsilon_0=0.3$, decay=0.99, episodes=1000.

- Lower $\varepsilon_0$ (0.1) led to premature exploitation; higher (0.5) slowed convergence.

- Faster decay (0.95) reduced exploration too quickly; slower decay hindered fine-tuning.

## 3. Selective Fusion

Enabling branching only when a local efficiency gain > 0.1 was detected:

- **Path Length:** 4 (vs. 3 unrestricted)

- **Efficiency:** 2.25 (↑ 7% over unrestricted RL)

- **Redundancy:** 2 alternative viable paths generated

- **Compute Time:** +5% overhead

Selective fusion strikes a balance between redundancy and growth cost, improving efficiency while maintaining fault tolerance.

## 4. Dynamic Obstacles

When 10% of nodes failed randomly at step 50:

- **RL adaptation:** 85% of runs still reached the goal (vs. 0% for A*/Dijkstra without replanning)

- **Average Efficiency:** 1.95 (vs. 1.71 static baselines)

The learned policy allows on-the-fly replanning, demonstrating resilience to unexpected changes.

---

## 1. Why the RL Agent Outperforms A* and Dijkstra

- **Multi-Objective Reward vs. Single-Objective Heuristic**

    - **A* and Dijkstra** optimize only for path length (or cumulative cost) using a fixed heuristic or cost function. They have no mechanism to value "nutrient" rewards, so they miss routes that collect high nutrients at only marginal additional cost.

    - **Q-learning**, by contrast, **integrates nutrient and cost into a single reward**, so the learned policy preferentially steers growth toward cells that offer the best nutrient-to-cost trade-off. Over time, this shapes an implicit heuristic in the Q-table that can outperform a purely distance-based heuristic.

- **Learning vs. Static Heuristics**

    - The RL agent **adapts its policy** based on real experience: it learns that certain grid regions consistently yield high cumulative reward, even if reaching them requires a detour. A*'s admissible heuristic (Manhattan distance) cannot encode such empirical knowledge, so it remains blind to latent environmental structure.

- **Exploration Enables Novel Paths**

    - The $\varepsilon$-greedy strategy encourages **temporary exploration**. Early in training, the agent discovers nutrient-rich shortcuts that A* and Dijkstra never consider, then exploits them. This exploration–exploitation balance is critical: too little exploration and the agent reverts to suboptimal greedy paths; too much and it never converges. Our chosen $\varepsilon_0 = 0.3$, decay = 0.99 achieved the sweet spot, yielding the dramatic efficiency gains seen in the 3-step,

2.11-efficiency path.

## 2. Hyperparameter Effects

- **Learning Rate (α)**

  - With **α = 0.1**, the agent updates Q-values neither too slowly (which would delay convergence) nor too aggressively (which could destabilize learning). We observed that α below 0.05 led to very slow improvements; above 0.2 produced noisy value estimates.

- **Discount Factor (γ)**

  - A discount of **0.9** ensures that future nutrient rewards remain significant, encouraging the agent to plan multi-step nutrient gains rather than greedily capitalizing on immediate reward. Lower γ (e.g., 0.8) biased the agent toward myopic paths; γ = 0.99 slowed convergence as the value horizon widened.

- **Exploration Rate (ε)**

  - **Starting ε at 0.3** provided enough randomness for discovering high-value nodes off the direct path, while **decaying at 0.99** shifted the policy toward exploitation after ~200 episodes. Faster decay (0.95) truncated exploration too quickly; slower decay (0.995) left the agent dithering in suboptimal exploratory loops.

## 3. Selective Fusion Trade-offs

- **Thresholded Branching**

  - Unrestricted branching produces maximal redundancy but incurs extra cost in unnecessary side-growth. By **only allowing a branch when the local efficiency gain > 0.1**, we trimmed wasteful splits while preserving resilience.

- **Outcome**

  - This scheme raised overall efficiency from 2.11 to **2.25**, at the expense of a slight increase in path length (4 steps vs. 3) and a modest 5% compute-time overhead. The redundant paths (two viable alternatives) now appear only where they materially improve nutrient harvest or fault tolerance.

## 4. Resilience under Dynamic Obstacles

- **Policy Generalization**

  - Since the Q-learning policy encodes state–action values for all reachable grid configurations seen during training, encountering a sudden obstacle simply redirects active tips to next-best actions without a full replanning run.

- **Multiple Active Tips**

  - The mycelial growth process maintains **several simultaneous tips**, so even if one tip hits a failed node, others continue propagation. This decentralized approach contrasts with A* and Dijkstra, which must restart global search when their single path is blocked.

- **Result**

  - Under random 10% node failures at mid-episode, MITE still reached the goal in **85% of trials**, maintaining an average efficiency of **1.95**—well above the classical static baseline of 1.71.

## 5. Computational Considerations

- **Training vs. Inference**

  - The RL training phase (1000 episodes × 100 steps) is the only stage with significant overhead. Once the Q-table is learned, **inference is O(1)** per action—comparable to the

per-node overhead in A∗'s open-set operations.

- **Memory Footprint**

  - MITE's Q-table size scales with the number of distinct state representations encountered. In practice on an $n \times m$ grid, this remained tractable (on the order of $nm$) because only grid position plus local node attributes define a state.

**In sum**, the results reflect precisely how each paradigm's strengths and weaknesses play out:

- **A∗/Dijkstra**: Deterministic, provably optimal for single-criterion static problems—fast and low-memory, but oblivious to richer reward structures or environmental change.

- **RL (MITE):** Computationally heavier to train but yields **empirical heuristics** that balance multiple objectives, adapt on-the-fly, and naturally foster redundancy—crucial qualities for real-world network design challenges.

---

# Conclusions

MITE successfully demonstrates that biologically inspired RL can extend classical graph algorithms to produce network topologies that are simultaneously shorter, cheaper, and more resource-efficient. Key takeaways include:

- **Multi-Objective Optimization:** The RL agent naturally balances distance, cost, and nutrient collection without manually weighting separate objectives.

- **Adaptive Resilience:** Once trained, MITE replans in real time when the environment changes—a capability absent from static algorithms.

- **Redundancy Management:** Selective fusion yields multiple fault-tolerant routes with minimal overhead.

- **Parameter Sensitivity:** Careful tuning of exploration rate and episode count is crucial; grid search identified an effective operating point.

**Future work** will explore scaling MITE to continuous state spaces using deep RL for function approximation, integrating additional biological behaviors (e.g., chemotaxis), and applying the framework to large-scale real-world logistics and communications networks.