# When to Not Trust the Oracle: Causal Robustness in LLM-Guided Reinforcement Learning Agents

**Abstract**

Large Language Models (LLMs) have emerged as promising high-level planners for reinforcement learning agents in complex environments. However, LLM advice can be unreliable under corrupted, ambiguous, or erroneous state conditions, raising critical safety concerns for autonomous decision-making systems. We present a lightweight safety framework that detects and mitigates unsafe LLM guidance in the NetHack Learning Environment, a challenging testbed with rich combinatorial state space, partial observability, and long-horizon dependencies. Our approach introduces a three-tiered architecture: (1) a baseline PPO agent establishing pure RL performance, (2) an LLM-guided agent demonstrating benefits of strategic advice, and (3) a causally robust agent with consistency-checking mechanisms that scrutinize LLM recommendations before influencing behavior. Through systematic state corruption experiments and comparative evaluation, we demonstrate that our safety framework prevents catastrophic failures while maintaining performance benefits of language-driven guidance. This work contributes practical methodology for trust calibration in LLM-guided RL and opens new directions for adversarial robustness in hybrid AI systems.

**Keywords:** Reinforcement Learning, Large Language Models, Causal Robustness, Safety Mechanisms, NetHack, Adversarial Robustness

---

# CHAPTER 1: INTRODUCTION

## 1.1 Problem Statement

In the current atmosphere of automation and artificial intelligence, Large Language Models (LLMs) have recently been explored as high-level planners for autonomous systems, particularly for Reinforcement Learning (RL) agents in complex environments. This approach represents an experimental technique that promises to bridge symbolic reasoning with continuous control. However, this paradigm is accompanied by significant problems of its own.

**The Core Challenge:** LLM advice can be unreliable, particularly under corrupted or ambiguous state conditions. When state observations are noisy, incomplete, or adversarially perturbed, language models may generate plausible-sounding but dangerous recommendations. This raises critical safety concerns for decision-making processes powered by LLM-generated advice.

**The Black Box Problem:** Most research on LLM-guided reinforcement learning uses the LLM as a planner or "common sense" advisor but leaves critical gaps. Current agents operate as black boxes: they map states to actions without understanding *why* strategies work, which limits generalization and makes failure modes unpredictable.

**Our Approach:** We address this gap by proposing a lightweight safety framework that detects and mitigates unsafe LLM guidance. We adopt the NetHack Learning Environment (NLE) as our testbed due to its:

- Rich combinatorial state space
- Partial observability
- Long-horizon decision dependencies
- Challenging benchmark characteristics that expose robustness issues

**Research Question:** The central inquiry guiding this work is:

> *"How can we make LLM-guided Reinforcement Learning Agents interpretable and robust to misinformation, particularly in complex sequential decision-making environments?"*
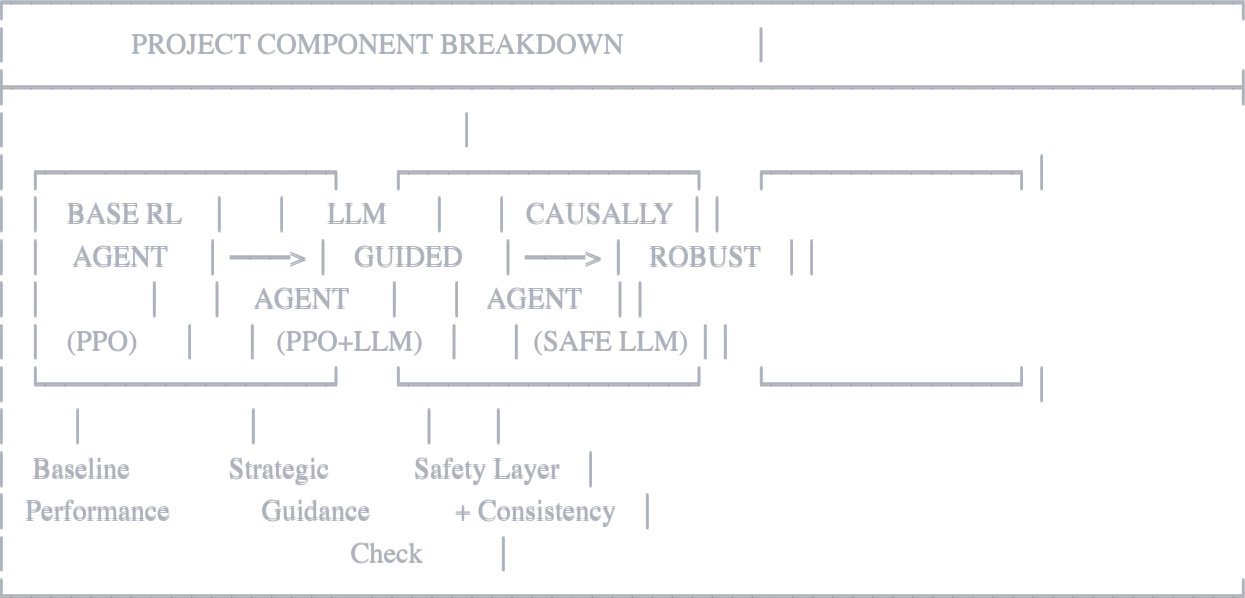
This work contributes:

1. **Practical methodology** for trust calibration in LLM-guided RL
2. **Adversarial robustness techniques** for language-driven agents
3. **Systematic evaluation** of failure modes under state corruption
4. **Causal consistency checking** mechanisms for safe AI collaboration

## 1.2 Scope of the Project

The project requires a three-part construction of different agent architectures for clear benchmarking and analysis:

```
┌──────────────────────────────────────────────────────────────────┐
│        PROJECT COMPONENT BREAKDOWN          │                      │
├──────────────────────────────────────────────────────────────────┤
│                             │                                      │
│        ┌──────────┐   ┌──────────┐   ┌──────────┐ │              │
│        │ BASE RL  │   │   LLM    │   │ CAUSALLY │ │               │
│        │  AGENT   │──▶│ GUIDED   │──▶│  ROBUST  │ │               │
│        │          │   │  AGENT   │   │  AGENT   │ │               │
│        │  (PPO)   │   │ (PPO+LLM)│   │ (SAFE LLM)│ │              │
│        └──────────┘   └──────────┘   └──────────┘ │               │
│             │              │              │                        │
│        Baseline        Strategic     Safety Layer │                │
│        Performance     Guidance       + Consistency │              │
│                           Check            │                       │
└──────────────────────────────────────────────────────────────────┘
```

**Component 1: Base RL Agent** The base RL agent establishes our performance baseline using standard Proximal Policy Optimization (PPO). This gives us a reference point for comparison and helps us understand what pure reinforcement learning can achieve without any linguistic guidance. Key characteristics:

- Recurrent CNN architecture for visual processing
- LSTM for temporal dependencies
- Reward shaping for exploration and survival
- No external guidance or safety mechanisms

**Component 2: LLM-Guided Agent** The LLM-guided agent introduces strategic advice from a language model, showing us the potential benefits of incorporating high-level reasoning into the decision-making process. This component:

- Converts game state to natural language descriptions
- Queries LLM for strategic recommendations
- Provides soft guidance through action probability biases
- Maintains autonomy through trainable guidance layers

**Component 3: Causally Robust Agent** Our causally safe RL agent adds a consistency checking mechanism that scrutinizes LLM advice before it influences agent behavior. Safety features include:

- **Consistency verification**: Cross-checking advice against game state
- **Corruption detection**: Identifying unreliable state observations
- **Causal reasoning**: Evaluating advice consequences before execution
- **Fallback mechanisms**: Reverting to pure RL when LLM unreliable

**Environmental Focus:** Our work specifically focuses on environments where partial observability and long-term planning create significant challenges. The NetHack Learning Environment provides exactly these conditions:

- **Procedural generation**: No two games are identical
- **Permadeath**: Single mistakes can end episodes
- **Sparse rewards**: Long sequences without feedback
- **Symbolic reasoning**: Understanding item properties, monster behaviors
- **Strategic depth**: 20+ action types, 100+ item types, complex interactions

We deliberately avoid simpler environments (e.g., Atari, CartPole) because they wouldn't expose the kinds of failure modes we're trying to address.

## 1.3 Objectives

The primary objective of this research is to **develop and validate a safety framework that makes LLM-guided reinforcement learning more reliable and interpretable**. This breaks down into several concrete goals:

**Objective 1: Establish Performance Baselines**

First, we want to establish clear performance baselines across different agent architectures. This means:

- Training a standard PPO agent until it reaches stable performance
- Comparing how LLM guidance affects learning speed and final capability
- Understanding *when, where, and why* language models provide value
- Quantifying the computational overhead of LLM integration

**Success Criteria:**

- 100+ training episodes per agent configuration
- Statistically significant performance differences (if any)
- Reproducible training protocols with documented hyperparameters

**Objective 2: Identify Failure Modes**

Second, we aim to identify specific failure modes in LLM-guided reinforcement learning. Language models sometimes give advice that sounds plausible but leads to catastrophic outcomes. By systematically corrupting state information and analyzing how agents respond, we can:

- **Catalog failure types**: What kinds of mistakes occur most frequently?
- **Understand root causes**: Why do LLMs give bad advice under corruption?
- **Measure severity**: How catastrophic are different failure modes?
- **Detect patterns**: Are certain game states more vulnerable?

**Corruption Scenarios:**

1. **Health corruption**: Reporting incorrect HP values
2. **Threat corruption**: Hiding nearby monsters or inventing fake ones
3. **Item corruption**: Misidentifying inventory contents
4. **Spatial corruption**: Incorrect position or surroundings
5. **Message corruption**: Garbled or misleading game text

**Objective 3: Demonstrate Safety Mechanisms**

Third, we need to demonstrate that our consistency-checking approach actually prevents these failures without sacrificing too much performance. This requires:

- **Causal verification**: Does the agent detect inconsistent advice?
- **Failure prevention**: Are catastrophic outcomes avoided?
- **Performance preservation**: Is the agent still effective overall?
- **Interpretability**: Can we explain why advice was rejected?

**Success Criteria:**

- Reduced catastrophic failure rate under corruption (target: <50% of unprotected agent)
- Performance within 80% of LLM-guided agent under normal conditions
- Interpretable rejection logs showing why advice was flagged

**Objective 4: Advance Safe AI Collaboration**

Beyond immediate technical goals, this research aims to advance the broader field of safe human-AI collaboration:

- **Practical frameworks**: Lightweight safety checks deployable in real systems
- **Theoretical insights**: Understanding when to trust language models
- **Adversarial robustness**: Techniques for handling corrupted observations
- **Transparency**: Making hybrid systems more interpretable

---

# CHAPTER 2: BACKGROUND AND RELATED WORK

## 2.1 NetHack as a Research Platform

### 2.1.1 Game Characteristics

NetHack is a procedurally generated roguelike dungeon crawler first released in 1987, still actively played and developed by a dedicated community. Its longevity and complexity make it an exceptional benchmark for AI research.

**Key Challenges:**

- **State Space Complexity**: ~10^50 possible states (compare to chess: ~10^47)
- **Action Space**: 23 discrete actions with context-dependent effects
- **Partial Observability**: Limited vision, hidden traps, unidentified items
- **Long Horizons**: Average winning games require 100,000+ steps
- **Strategic Depth**: Item combinations, monster tactics, resource management
- **Permadeath**: No checkpoints or save-scumming

**Why NetHack for Safety Research:** NetHack's combination of symbolic reasoning requirements and real-time control makes it ideal for testing LLM-guided agents:

1. **High-stakes decisions**: One mistake can end a promising run
2. **Ambiguous situations**: Not all dangers are immediately visible
3. **Complex state dependencies**: Effects of actions depend on hidden variables
4. **Natural language mapping**: Game messages provide semantic context

### 2.1.2 The NetHack Learning Environment (NLE)

Küttler et al. (2020) introduced NLE, a standardized interface for RL research:

- **Observations**: Multi-modal (glyphs, stats, messages, inventory)
- **Actions**: 23 discrete commands (movement, combat, item use)
- **Rewards**: Score-based and survival-based variants
- **Integration**: Compatible with OpenAI Gym/Gymnasium

**State Representation:**

python

```python
observation = {
    'glyphs': [21, 79],      # Visual map (integer IDs)
    'blstats': [26],         # Player stats (HP, level, etc.)
    'message': [256],        # Recent game text
    'inv_strs': [55],        # Inventory strings
    'inv_letters': [55],     # Item slot letters
    'inv_oclasses': [55],    # Object classes
    'tty_chars': [24, 80],   # Raw terminal output
    'tty_colors': [24, 80],  # Terminal colors
    'tty_cursor': [2]        # Cursor position
}
```

## 2.2 Reinforcement Learning Foundations

### 2.2.1 Proximal Policy Optimization (PPO)

Our baseline agent uses PPO (Schulman et al., 2017), a policy gradient method that has become the de facto standard for on-policy RL:

**Key Advantages:**

1. **Sample efficiency**: Better than vanilla policy gradient
2. **Stability**: Clipped objective prevents destructive updates
3. **Simplicity**: Few hyperparameters, easy to tune
4. **Versatility**: Works across diverse environments

**Core Algorithm:**

```
Collect trajectories using πθ_old
Compute advantages Â_t using GAE-λ
For K epochs:
    For each minibatch:
        Compute ratio r_t = πθ(a_t|s_t) / πθ_old(a_t|s_t)
        Compute L_CLIP = min(r_t Â_t, clip(r_t, 1-ε, 1+ε) Â_t)
        Update θ using gradient of L_CLIP
```

### 2.2.2 Recurrent Neural Networks in RL

NetHack's partial observability necessitates memory mechanisms. We employ:

- **LSTM layers**: For maintaining state history across timesteps
- **Spatial CNNs**: For processing 21×79 glyph maps
- **Attention mechanisms**: (future work) For focusing on relevant map regions

**Why Recurrence Matters:**

- Doors may be closed, hiding monsters behind them
- Item effects persist across many timesteps
- Strategic plans require multi-step execution
- Pure feedforward policies cannot maintain context

### 2.2.3 Reward Shaping

Sparse rewards in NetHack (often 0 for hundreds of steps) require intrinsic motivation:

**Our Shaping Strategy:**

python

```python
shaped_reward = (
    raw_reward                  # Environmental signal
    + 0.01 * exploration_bonus     # New tiles discovered
    + 0.001 * health_delta         # HP changes
    + 1.0 * level_up               # Character progression
    + 0.05 * item_pickup           # Item acquisition
    - 0.01 * stuck_penalty         # Anti-loop mechanism
    - 1.0 * death                  # Terminal failure
)
```

This shaping accelerates early learning while maintaining alignment with game objectives.

## 2.3 Large Language Models as Reasoning Engines

### 2.3.1 LLMs for Planning

Recent work demonstrates LLMs' ability to decompose complex tasks into subgoals:

**Key Research:**

- **ReAct** (Yao et al., 2023): Interleaving reasoning traces with actions
- **Reflexion** (Shinn et al., 2023): Self-reflection for iterative improvement
- **Inner Monologue** (Huang et al., 2022): Closed-loop language feedback

**Common Pattern:**

**Limitations:**

1. **Computational cost**: Inference latency unsuitable for real-time control
2. **Grounding problem**: Abstract plans may not map to executable actions
3. **Consistency**: No guarantee plans remain valid as state evolves
4. **Safety**: Plausible but dangerous recommendations

### 2.3.2 LLMs for Robotic Control

Language-conditioned policies show promise in embodied AI:

- **SayCan** (Ahn et al., 2022): Grounding LLM plans in affordances
- **Code-as-Policies** (Liang et al., 2023): Generating executable code
- **RT-2** (Brohan et al., 2023): Vision-language-action models

**Relevance to NetHack:** Similar challenges exist in translating high-level advice ("retreat from the orc") into low-level actions (specific movement sequences).

### 2.3.3 Prompt Engineering for Decision-Making

Effective LLM guidance requires careful prompt design:

**Best Practices:**

1. **Structured context**: Clear formatting of state information
2. **Explicit options**: List available actions/strategies
3. **Decision rules**: Specify when to choose each option
4. **Low temperature**: Reduce randomness for consistency
5. **Few-shot examples**: Prime model with successful scenarios

## 2.4 Safety and Robustness in AI Systems

### 2.4.1 Adversarial Robustness

Neural networks are vulnerable to small perturbations:

- **Image classifiers**: ε-ball attacks (Goodfellow et al., 2015)
- **RL agents**: Observation perturbations (Huang et al., 2017)
- **LLMs**: Jailbreaks and prompt injection (Wei et al., 2023)

**Our Focus:** State corruption rather than pixel-level adversarial examples

### 2.4.2 Causal Reasoning for Safety

Causal inference provides tools for reliable decision-making:

- **Pearl's Causal Hierarchy**: Association → Intervention → Counterfactual
- **Structural Causal Models**: Explicit cause-effect relationships
- **Backdoor Adjustment**: Deconfounding for unbiased estimates

**Application to LLM Guidance:** Before accepting advice, ask: "What would happen if I followed this recommendation?" Simulate consequences using learned world model or heuristics.

### 2.4.3 Trust Calibration

Knowing when to trust AI advice is crucial:

- **Uncertainty quantification**: Confidence estimates for predictions
- **Out-of-distribution detection**: Flagging unfamiliar states

- **Consistency checking**: Cross-validating with multiple sources
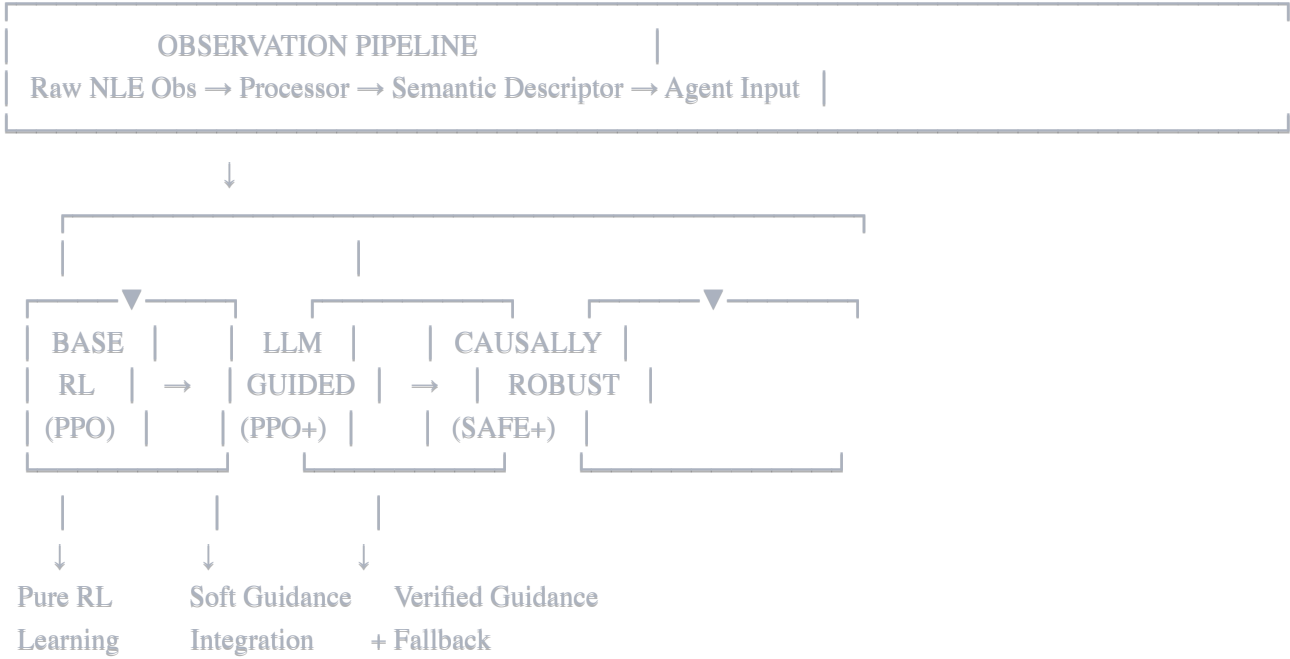
## 2.5 Related Work Comparison

| Approach | LLM Role | Safety Mechanism | Environment | Limitation |
|---|---|---|---|---|
| ReAct (2023) | High-level planner | None | Text-based tasks | No adversarial testing |
| SayCan (2022) | Task decomposer | Affordance grounding | Robotics | Assumes accurate perception |
| Inner Monologue (2022) | Feedback provider | Success detection | Robotics | Manual verification |
| **Our Work** | Strategic advisor | Causal consistency | NetHack | Complex state space |

**Key Distinction:** We explicitly test robustness under corrupted observations, not just nominal performance.

---

# CHAPTER 3: METHODOLOGY

## 3.1 System Architecture Overview

Our framework consists of three agent configurations built on a common foundation:



```
┌─────────────────────────────────────────────────────────┐
│           OBSERVATION PIPELINE              │           │
│  Raw NLE Obs → Processor → Semantic Descriptor → Agent Input  │
└─────────────────────────────────────────────────────────┘
                    ↓

        ┌──────────────┬───────────────┬──────────────┐
        │              │               │              │
        ▼              ▼               ▼
   ┌─────────┐   ┌──────────┐    ┌──────────┐
   │ BASE    │   │ LLM      │    │ CAUSALLY │
   │ RL      │ → │ GUIDED   │ →  │ ROBUST   │
   │ (PPO)   │   │ (PPO+)   │    │ (SAFE+)  │
   └─────────┘   └──────────┘    └──────────┘
        │              │               │
        ↓              ↓               ↓
   Pure RL       Soft Guidance    Verified Guidance
   Learning      Integration      + Fallback
```

## 3.2 Component 1: Base RL Agent

### 3.2.1 Observation Processing

**NetHackObservationProcessor** converts raw observations into normalized tensors:



python

```python
def process_observation(obs, last_action=None):
    processed = {}

    # Visual map: 21×79 integer glyphs → [0,1] normalized
    glyphs = obs['glyphs'].astype(float32) / 5976.0

    # Player stats: 26 values (HP, level, gold, etc.)
    stats = obs['blstats'].astype(float32)
    stats[0] = stats[0] / stats[1]  # HP ratio
    stats[7] = min(stats[7] / 30.0, 1.0)  # Level normalization

    # Message: 256 ASCII values → [0,1]
    message = obs['message'].astype(float32) / 255.0

    # Inventory: 55 item slots → binary presence
    inventory = np.zeros(55)
    for i, item in enumerate(obs['inv_strs']):
        inventory[i] = 1.0 if item else 0.0

    # Action history: Last 50 actions → [0,1]
    action_history = encode_recent_actions(last_action)

    return {
        'glyphs': glyphs,
        'stats': stats,
        'message': message,
        'inventory': inventory,
        'action_history': action_history
    }
```

**Key Design Decisions:**

1. **Normalization**: Prevents gradient explosion from large raw values
2. **HP ratio**: More informative than absolute health
3. **Action history**: Enables detection of repetitive loops
4. **Modularity**: Easy to add/remove features for ablation studies

**3.2.2 Network Architecture**

**Actor Network (Policy):**

```
Input Modalities:
├── Glyphs (21×79) → CNN → LSTM → 256-dim
├── Stats (26) → LSTM → 64-dim
├── Message (256) → FC → 128-dim
├── Inventory (55) → FC → 64-dim
└── Action History (50) → FC → 32-dim
                    ↓
        Combined (544) → FC(512) → FC(256)
                    ↓
             Action Logits (23)
```

## Convolutional Layers:

python

```python
conv1 = Conv2d(1, 32, kernel=3, padding=1)  # 21×79 → 21×79
pool1 = MaxPool2d(2)                         # 21×79 → 10×39
conv2 = Conv2d(32, 64, kernel=3, padding=1)  # 10×39 → 10×39
pool2 = MaxPool2d(2)                         # 10×39 → 5×19
conv3 = Conv2d(64, 128, kernel=3, padding=1) # 5×19 → 5×19
pool3 = MaxPool2d(2)                         # 5×19 → 2×9
flatten → Linear(2×9×128=2304) → 512 → LSTM(256)
```

## LSTM Rationale:

- Spatial CNN captures local patterns (walls, monsters)
- Temporal LSTM maintains context (door states, item effects)
- Separate LSTMs for vision and stats prevent feature interference

**Critic Network (Value Function):** Identical architecture but outputs scalar value V(s) instead of action logits.

### 3.2.3 Training Algorithm: PPO

**Clipped Surrogate Objective:**

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

LCLIP (θ) = Et [min (rt(θ)A^t, clip(rt(θ), 1 − ε, 1 + ε)A^t)]

Where:

- $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ rt(θ) = πθold (at|st)πθ (at|st) is probability ratio
- $\hat{A}_t$ A^t is advantage estimate
- $\epsilon = 0.2$ ε = 0.2 is clip parameter

**Advantage Estimation (GAE):**

$$\hat{A}_t = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l}$$

A^t = l=0∑ ∞(γλ)lδt+l

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

δt = rt + γV (st+1) − V (st)

With $\gamma = 0.99$ γ = 0.99 (discount), $\lambda = 0.95$ λ = 0.95 (trace decay)

**Value Function Loss:**

$$L^{VF}(\phi) = E_t[(V_\phi(s_t) - \hat{R}_t)^2]$$

LV F $(\phi)$ = Et[(V$\phi$(st) − R^t)2]

Where $\hat{R}_t = \hat{A}_t + V(s_t)$R^t = A^t + V (st) is empirical return

**Entropy Regularization:**

$$H(\pi_\theta) = -E_{a \sim \pi}[\log \pi_\theta(a|s)]$$

H($\pi\theta$) = −Ea~$\pi$[log $\pi\theta$(a|s)]

Encourages exploration by penalizing deterministic policies

**Combined Loss:**

$$L(\theta,\phi) = L^{CLIP}(\theta) - c_1 L^{VF}(\phi) + c_2 H(\pi_\theta)$$

L($\theta,\phi$) = LCLIP ($\theta$) − c1LV F ($\phi$) + c2H($\pi\theta$)

With $c_1 = 0.5$c1 = 0.5, $c_2 = 0.02$c2 = 0.02

**Optimization:**

- Adam optimizer with $\alpha = 10^{-4}$ $\alpha$ = 10–4
- Batch size: 64
- Update every 1024 steps (rollout buffer)
- 4 epochs per update
- Gradient clipping: max norm 0.5

### 3.2.4 Reward Shaping

Raw NetHack rewards are extremely sparse:

- +1 for every 1 point of score increase
- Typical episode: $0, 0, 0, ..., 0, +5$ (monster kill), $0, 0, ...$

**Our Shaping Function:**

$$r_t^{shaped} = r_t^{env} + \sum_i \alpha_i \cdot f_i(s_t, a_t, s_{t+1})$$

rtshaped = rtenv + i$\sum$ $\alpha$i · fi(st, at, st+1)

**Shaping Terms:**

| Term | Coefficient | Trigger | Rationale |
|---|---|---|---|
| Exploration | +0.01 | New tile visited | Encourage map coverage |
| Health gain | +0.001 | HP increases | Reward healing |
| Health loss | −0.001 | HP decreases | Penalize damage |
| Level up | +1.0 | Experience level ↑ | Major milestone |
| Experience | +0.0001 | XP increases | Gradual progression |
| Item pickup | +0.05 | Inventory grows | Gather resources |
| Stuck penalty | −0.01 | Same position 10+ steps | Prevent loops |
| Death | −1.0 | HP ≤ 0 | Strong terminal signal |

**Implementation:**

python

```python
class NetHackRewardShaper:
    def shape_reward(self, obs, raw_reward, done, info):
        shaped = raw_reward

        if self.prev_stats is not None:
            # Health
            shaped += 0.001 * (obs['HP'] - self.prev_stats['HP'])

            # Level
            shaped += 1.0 * (obs['level'] - self.prev_stats['level'])

            # Experience
            shaped += 0.0001 * (obs['XP'] - self.prev_stats['XP'])

        # Exploration
        pos = (obs['x'], obs['y'])
        if pos not in self.visited:
            shaped += 0.01
            self.visited.add(pos)

        # Anti-stuck
        if pos == self.last_pos:
            self.stuck_count += 1
            if self.stuck_count > 10:
                shaped -= 0.01
        else:
            self.stuck_count = 0

        # Death
        if done and obs['HP'] <= 0:
            shaped -= 1.0

        self.prev_stats = obs
        self.last_pos = pos

        return shaped
```

## 3.3 Component 2: LLM-Guided Agent

### 3.3.1 Semantic State Description

**Challenge:** LLMs cannot process raw numpy arrays. We need natural language descriptions.

**NetHackSemanticDescriptor** translates game state:

**Glyph Mappings:**



python

```python
GLYPH_TO_SYMBOL = {
    # Terrain
    2359: "wall", 2360: "door", 2361: "floor",
    2364: "stairs_down", 2365: "stairs_up",

    # Monsters (simplified)
    2378: "kobold", 2379: "goblin", 2380: "orc",
    2381: "troll", 2382: "dragon",

    # Items
    2395: "gold", 2396: "weapon", 2397: "armor",
    2398: "food", 2399: "potion", 2400: "scroll"
}
```

**Surroundings Analysis:**

python

```python
def describe_surroundings(glyphs, player_pos):
    py, px = player_pos
    nearby_monsters = []
    nearby_items = []

    # Check 5×5 neighborhood
    for dy in range(-2, 3):
        for dx in range(-2, 3):
            ny, nx = py + dy, px + dx
            if 0 <= ny < 21 and 0 <= nx < 79:
                glyph = glyphs[ny, nx]
                symbol = GLYPH_TO_SYMBOL.get(glyph, "unknown")
                distance = abs(dy) + abs(dx)  # Manhattan
                direction = get_direction(dy, dx)

                if symbol in ["kobold", "goblin", "orc", ...]:
                    nearby_monsters.append(
                        f"{symbol} {direction} (dist:{distance})"
                    )
                elif symbol in ["gold", "weapon", "food", ...]:
                    nearby_items.append(
                        f"{symbol} {direction} (dist:{distance})"
                    )

    # Sort by distance
    nearby_monsters.sort(key=lambda x: extract_distance(x))

    if nearby_monsters:
        description = f"CLOSEST THREAT: {nearby_monsters[0]}"
        if len(nearby_monsters) > 1:
            description += f"; Other threats: {nearby_monsters[1:]}"
    else:
        description = "NO IMMEDIATE THREATS - safe to explore"

    if nearby_items:
        description += f"; Items nearby: {nearby_items}"

    return description
```

**Player Status:**

☑
python

```python
def describe_player_status(blstats):
    hp = int(blstats[10])
    max_hp = int(blstats[11])
    level = int(blstats[18])
    xp = int(blstats[19])
    depth = int(blstats[12])
    gold = int(blstats[13])

    hp_ratio = hp / max_hp
    health = ("critical" if hp_ratio < 0.3 else
              "low" if hp_ratio < 0.6 else "good")

    return (f"Level {level}, Health: {hp}/{max_hp} ({health}), "
            f"XP: {xp}, Depth: {depth}, Gold: {gold}")
```

**Complete Description Example:**

```
NETHACK GAME STATE:
Status: Level 3, Health: 28/45 (low), XP: 234, Depth: 4, Gold: 87
Surroundings: CLOSEST THREAT: orc south (dist:1); Other threats: kobold west (dist:3)
Recent Message: You are hit by the orc!
Inventory: Carrying 7 items (moderate load)
Recent Actions: move_south → kick → move_south → wait

STRATEGIC ANALYSIS:
Based on game state, choose ONE primary strategy:
1. "explore" - No immediate threats
2. "combat" - Monster nearby AND health good (>60%)
3. "retreat" - Monster nearby BUT health low (<40%)
4. "collect" - Items nearby and safe
5. "wait" - Critical health or need recovery

Your strategic choice:
```

### 3.3.2 LLM Strategic Advisor

**Architecture:**

python

```python
class EnhancedLLMAdvisor:
    def __init__(self, call_frequency=50):
        self.call_frequency = call_frequency
        self.step_count = 0
        self.semantic_descriptor = NetHackSemanticDescriptor()

        # Strategy → Action mappings
        self.action_categories = {
            'explore': [0,1,2,3,4,5,6,7,11],  # moves + search
            'combat': [14],                 # kick
            'retreat': [0,1,2,3],           # move away
            'collect': [9],                 # pickup
            'wait': [8]                     # rest
        }

    def should_call_llm(self, performance):
        self.step_count += 1

        # Call every N steps OR if performance poor
        periodic = (self.step_count % self.call_frequency == 0)
        struggling = performance['avg_reward'] < 5.0

        return periodic or struggling

    async def get_strategic_advice(self, raw_obs, processed_obs, performance):
        # Generate semantic description
        description = self.semantic_descriptor.generate_full_description(
            raw_obs, processed_obs, self.recent_actions
        )

        # Call LLM
        llm_response = await self._call_ollama_api(description, performance)

        # Parse strategy
        strategy = self._parse_strategic_response(llm_response, raw_obs)

        # Convert to action hints
        hints = self._strategy_to_hints(strategy, raw_obs)

        return hints

    async def _call_ollama_api(self, description, performance):
        import aiohttp

        prompt = f"""You are an expert NetHack strategic advisor.

{description}

RECENT PERFORMANCE:
```

```
- Average Reward: {performance['avg_reward']:.2f}
- Average Survival: {performance['avg_length']:.0f} steps

STRATEGIC ANALYSIS:
Based on the game state above, choose ONE primary strategy:

1. "explore" - No immediate threats, safe to move and search
2. "combat" - Monster nearby AND health good (>60%)
3. "retreat" - Monster nearby BUT health low (<40%)
4. "collect" - Items nearby and safe
5. "wait" - Critical health or need recovery

CRITICAL RULES:
- If threat distance 1-2 AND health <40%: Choose "retreat"
- If threat distance 1-2 AND health >60%: Choose "combat"
- If NO IMMEDIATE THREATS: Choose "explore" or "collect"
- If health critical: Choose "wait" or "retreat"

Respond with ONLY ONE WORD: explore, combat, retreat, collect, wait

Your strategic choice:"""

        try:
            async with aiohttp.ClientSession() as session:
                payload = {
                    "model": "llama3:8b",
                    "prompt": prompt,
                    "stream": False,
                    "options": {
                        "temperature": 0.2,
                        "num_predict": 50
                    }
                }

                async with session.post(
                    "http://localhost:11434/api/generate",
                    json=payload,
                    timeout=aiohttp.ClientTimeout(total=15)
                ) as response:
                    if response.status == 200:
                        result = await response.json()
                        return result.get("response", "").strip().lower()
                    else:
                        return ""
        except Exception as e:
            print(f"LLM error: {e}")
            return ""

    def _strategy_to_hints(self, strategy, raw_obs):
        """Convert strategy to action probability hints"""
```

```python
        hints = np.zeros(23, dtype=np.float32)

        if strategy in self.action_categories:
            for action_id in self.action_categories[strategy]:
                hints[action_id] = 0.2  # 20% boost

        return hints
```

**Key Design Principles:**

1. **Sparse Calling**: Every 50 steps (~5% of timesteps) to minimize overhead
2. **Structured Prompts**: Clear options and decision rules
3. **Low Temperature**: 0.2 for consistent strategic decisions
4. **Soft Hints**: 20% probability boost, not hard constraints
5. **Async API**: Non-blocking LLM calls

### 3.3.3 Guidance Integration Layer

**LLMEnhancedPPOActor** extends base actor:

python

```python
class LLMEnhancedPPOActor(nn.Module):
    def __init__(self, action_dim=23):
        super().__init__()

        # Base RL networks (same as Component 1)
        self.glyph_cnn = RecurrentNetHackCNN()
        self.stats_lstm = nn.LSTM(26, 64)
        self.message_fc = nn.Linear(256, 128)
        self.inventory_fc = nn.Linear(55, 64)
        self.action_hist_fc = nn.Linear(50, 32)

        self.combined_fc1 = nn.Linear(544, 512)
        self.combined_fc2 = nn.Linear(512, 256)
        self.action_head = nn.Linear(256, action_dim)

        # NEW: Trainable guidance layer
        self.llm_guidance_fc = nn.Linear(action_dim, action_dim)

        # Guidance settings
        self.llm_guidance_weight = 0.0  # Set by agent
        self.use_llm = False

    def forward(self, obs, reset_hidden=False, llm_hints=None):
        # Standard forward pass
        glyph_features = self.glyph_cnn(obs['glyphs'], reset_hidden)
        stats_features = self.stats_lstm(obs['stats'])
        message_features = F.relu(self.message_fc(obs['message']))
        inventory_features = F.relu(self.inventory_fc(obs['inventory']))
        action_hist_features = F.relu(self.action_hist_fc(obs['action_history']))

        combined = torch.cat([
            glyph_features, stats_features,
            message_features, inventory_features, action_hist_features
        ], dim=1)

        x = F.relu(self.combined_fc1(combined))
        x = F.relu(self.combined_fc2(x))

        # Base policy logits
        base_logits = self.action_head(x)

        # Optional: Add LLM guidance if enabled
        if self.use_llm and llm_hints is not None and self.llm_guidance_weight > 0:
            llm_bias = torch.FloatTensor(llm_hints).to(base_logits.device)

            # Learn to use LLM hints through trainable layer
            learned_guidance = self.llm_guidance_fc(llm_bias)

            # Small additive bias
```

```
        guided_logits = base_logits + self.llm_guidance_weight * learned_guidance
        return guided_logits


    return base_logits
```

**Mathematical Formulation:**

$$\pi_\theta(a|s,h) \propto \exp(\text{logits}_\theta(s) + \lambda \cdot W_g h)$$

Where:

- $\text{logits}_\theta(s) \in R^{|A|}$ are base policy logits
- $h \in R^{|A|}$ are LLM hints (0.2 for suggested actions, 0 otherwise)
- $W_g \in R^{|A| \times |A|}$ is learned guidance transformation
- $\lambda = 0.05$ is guidance weight (small to preserve autonomy)

**Why This Works:**

- **Trainable**: $W_g$ learns which hints to trust over time
- **Soft**: Additive bias doesn't force actions, just nudges probabilities
- **Preserves learning**: Base logits still dominant, RL continues improving
- **Interpretable**: Can visualize which actions get boosted

## 3.4 Component 3: Causally Robust Agent

### 3.4.1 Motivation: When LLMs Fail

**Failure Scenario 1: Health Misreporting**



```
TRUE STATE: HP = 5/45 (critical)
CORRUPTED STATE: HP = 40/45 (good)
LLM ADVICE: "combat" (engage monster)
OUTCOME: Agent dies immediately
```

**Failure Scenario 2: Phantom Threats**



```
TRUE STATE: No monsters nearby
CORRUPTED STATE: "dragon adjacent"
LLM ADVICE: "retreat" (flee)
OUTCOME: Agent wastes time running from nothing
```

**Failure Scenario 3: Missing Threats**

These scenarios motivate our **causal consistency checking** layer.

**3.4.2 Consistency Verification Framework**

**Core Idea:** Before accepting LLM advice, verify it's consistent with observable game state.

python

```python
class CausalConsistencyChecker:
    def __init__(self):
        self.verification_rules = {
            'health_consistency': self._check_health,
            'threat_consistency': self._check_threats,
            'spatial_consistency': self._check_spatial,
            'temporal_consistency': self._check_temporal
        }

        self.rejection_log = []
        self.trust_score = 1.0  # Decays with bad advice

    def verify_advice(self, strategy, raw_obs, semantic_description):
        """
        Returns: (is_safe, confidence, reason)
        """
        checks = []

        for check_name, check_fn in self.verification_rules.items():
            passed, confidence, reason = check_fn(
                strategy, raw_obs, semantic_description
            )
            checks.append((check_name, passed, confidence, reason))

        # Aggregate results
        all_passed = all(c[1] for c in checks)
        avg_confidence = np.mean([c[2] for c in checks])

        if not all_passed:
            failed_checks = [c for c in checks if not c[1]]
            reason = f"Failed: {[c[0] for c in failed_checks]}"
            self.rejection_log.append({
                'strategy': strategy,
                'reason': reason,
                'failed_checks': failed_checks
            })
            return False, avg_confidence, reason

        return True, avg_confidence, "All checks passed"
```

**Check 1: Health Consistency**

python

```python
def _check_health(self, strategy, raw_obs, description):
    """Verify health-dependent advice makes sense"""

    # Extract true health from raw observation
    blstats = raw_obs.get('blstats', np.zeros(26))
    true_hp = int(blstats[10])
    true_max_hp = int(blstats[11])
    true_hp_ratio = true_hp / true_max_hp if true_max_hp > 0 else 0.5

    # Extract health from semantic description
    import re
    health_match = re.search(r'Health: (\d+)/(\d+)', description)
    if health_match:
        described_hp = int(health_match.group(1))
        described_max_hp = int(health_match.group(2))
        described_hp_ratio = described_hp / described_max_hp
    else:
        # Can't verify without description
        return True, 0.5, "No health info in description"

    # Check for discrepancy
    discrepancy = abs(true_hp_ratio - described_hp_ratio)

    if discrepancy > 0.2:  # More than 20% difference
        return False, 0.0, f"Health mismatch: true={true_hp_ratio:.2f}, described={described_hp_ratio:.2f}"

    # Check strategy appropriateness
    if strategy == 'combat' and true_hp_ratio < 0.4:
        return False, 0.0, f"Combat unsafe with HP={true_hp_ratio:.2f}"

    if strategy == 'explore' and true_hp_ratio < 0.3:
        return False, 0.0, f"Exploration risky with critical HP={true_hp_ratio:.2f}"

    return True, 1.0, "Health consistent"
```

**Check 2: Threat Consistency**



python

```python
def _check_threats(self, strategy, raw_obs, description):
    """Verify threat detection is accurate"""

    # Count monsters in true observation
    glyphs = raw_obs.get('glyphs', np.zeros((21, 79)))
    blstats = raw_obs.get('blstats', np.zeros(26))
    player_y, player_x = int(blstats[1]), int(blstats[0])

    true_threats = []
    for dy in range(-2, 3):
        for dx in range(-2, 3):
            ny, nx = player_y + dy, player_x + dx
            if 0 <= ny < 21 and 0 <= nx < 79:
                glyph = glyphs[ny, nx]
                # Monster glyphs typically in range 2378-2400
                if 2378 <= glyph <= 2400:
                    distance = abs(dy) + abs(dx)
                    true_threats.append(distance)

    # Extract threats from description
    described_threats = "CLOSEST THREAT" in description
    described_safe = "NO IMMEDIATE THREATS" in description

    # Check consistency
    has_true_threats = len(true_threats) > 0

    if has_true_threats and described_safe:
        return False, 0.0, "Description says safe but threats detected"

    if not has_true_threats and described_threats:
        return False, 0.0, "Description says threats but none detected"

    # Check strategy appropriateness
    if has_true_threats and min(true_threats) <= 1:
        # Adjacent threat
        if strategy == 'explore':
            return False, 0.0, "Exploration with adjacent threat"

    return True, 1.0, "Threat detection consistent"
```

**Check 3: Spatial Consistency**

python

```python
def _check_spatial(self, strategy, raw_obs, description):
    """Verify spatial relationships make sense"""

    # Check for stuck detection
    if "STUCK" in description and strategy not in ['explore', 'collect']:
        return False, 0.5, "Not addressing stuck condition"

    # Check for item collection
    if strategy == 'collect':
        # Verify items actually nearby
        if "Items nearby" not in description:
            return False, 0.0, "Collect strategy but no items mentioned"

    return True, 1.0, "Spatial consistency OK"
```

**Check 4: Temporal Consistency**

python

```python
def _check_temporal(self, strategy, raw_obs, description):
    """Verify advice doesn't contradict recent actions"""

    # Check for action loops
    if hasattr(self, 'recent_strategies'):
        if len(self.recent_strategies) >= 5:
            if self.recent_strategies[-5:] == [strategy] * 5:
                return False, 0.5, f"Repeated strategy '{strategy}' 5 times"

    return True, 1.0, "No temporal inconsistency"
```

**3.4.3 Fallback Mechanisms**

When LLM advice is rejected, the agent needs alternative strategies:

python

```python
class CausallyRobustAgent:
    def __init__(self, base_agent, consistency_checker):
        self.base_agent = base_agent
        self.checker = consistency_checker
        self.fallback_mode = False

    async def select_action(self, obs, performance):
        # Get LLM advice
        llm_hints = await self.base_agent.llm_advisor.get_strategic_advice(
            obs['raw'], obs['processed'], performance
        )

        # Extract strategy from hints
        strategy = self._hints_to_strategy(llm_hints)

        # Generate semantic description
        description = self.base_agent.llm_advisor.semantic_descriptor.generate_full_description(
            obs['raw'], obs['processed'], []
        )

        # Verify consistency
        is_safe, confidence, reason = self.checker.verify_advice(
            strategy, obs['raw'], description
        )

        if is_safe:
            # Use LLM guidance
            action, log_prob, value = await self.base_agent.select_action(
                obs['tensor'], obs['processed'], obs['raw'],
                reset_hidden=False,
                performance_metrics=performance
            )
            self.fallback_mode = False
        else:
            # Reject LLM advice, use rule-based fallback
            print(f"⚠️ LLM advice rejected: {reason}")
            action, log_prob, value = self._fallback_action_selection(obs)
            self.fallback_mode = True

        return action, log_prob, value

    def _fallback_action_selection(self, obs):
        """Rule-based safe actions when LLM untrusted"""

        blstats = obs['raw'].get('blstats', np.zeros(26))
        hp_ratio = blstats[10] / blstats[11] if blstats[11] > 0 else 0.5

        # Emergency rules
        if hp_ratio < 0.3:
```

```python
        # Critical health: wait/rest
        action = 8  # wait
    elif hp_ratio < 0.5:
        # Low health: retreat to safe distance
        action = np.random.choice([0, 1, 2, 3])  # move away
    else:
        # Good health: pure RL policy (no LLM hints)
        with torch.no_grad():
            action_logits = self.base_agent.actor(
                obs['tensor'],
                reset_hidden=False,
                llm_hints=None  # No guidance
            )
            action_dist = Categorical(logits=action_logits)
            action = action_dist.sample().item()

    # Compute value and log_prob for training
    with torch.no_grad():
        action_logits = self.base_agent.actor(obs['tensor'], llm_hints=None)
        action_dist = Categorical(logits=action_logits)
        log_prob = action_dist.log_prob(torch.tensor(action)).item()
        value = self.base_agent.critic(obs['tensor']).item()

    return action, log_prob, value
```

**Fallback Strategy Hierarchy:**

1. **Critical health (HP < 30%)**: Force wait/rest action
2. **Low health (HP < 50%)**: Random retreat movements
3. **Good health (HP ≥ 50%)**: Pure RL policy without LLM hints

This ensures the agent remains functional even when LLM guidance is entirely rejected.

### 3.4.4 Trust Score Dynamics

Track LLM reliability over time:

python

```python
class AdaptiveTrustModel:
    def __init__(self, initial_trust=0.8):
        self.trust_score = initial_trust
        self.advice_history = []

    def update_trust(self, advice_followed, outcome_quality):
        """
        advice_followed: bool - Was LLM advice used?
        outcome_quality: float - Reward obtained (or negative for failure)
        """

        if advice_followed:
            # Update trust based on outcome
            if outcome_quality > 0:
                # Good outcome: increase trust
                self.trust_score = min(1.0, self.trust_score + 0.05)
            elif outcome_quality < -0.5:
                # Bad outcome: decrease trust
                self.trust_score = max(0.0, self.trust_score - 0.1)

        self.advice_history.append({
            'followed': advice_followed,
            'quality': outcome_quality,
            'trust_after': self.trust_score
        })

    def should_trust_advice(self, confidence):
        """
        Decision rule: trust if both checker confidence and trust score high
        """
        combined_score = 0.5 * confidence + 0.5 * self.trust_score
        return combined_score > 0.6
```

**Adaptive Guidance Weight:**

$$\lambda(t) = \lambda_0 \cdot \tau(t) \cdot c(t)$$

Where:

- $\lambda_0 = 0.05$ is base guidance weight
- $\tau(t) \in [0, 1]$ is trust score (learned from outcomes)
- $c(t) \in [0, 1]$ is consistency check confidence

This allows the agent to gradually reduce reliance on unreliable LLMs.

### 3.4.5 State Corruption Experiments

To test robustness, we systematically corrupt observations:

python

```python
class StateCorruptor:
    def __init__(self, corruption_type='health', corruption_rate=0.3):
        self.corruption_type = corruption_type
        self.corruption_rate = corruption_rate

    def corrupt_observation(self, obs):
        """Apply corruption to observation"""

        if np.random.rand() > self.corruption_rate:
            return obs  # No corruption this step

        corrupted_obs = copy.deepcopy(obs)

        if self.corruption_type == 'health':
            # Report incorrect HP
            blstats = corrupted_obs['blstats']
            true_hp = blstats[10]
            max_hp = blstats[11]

            # Randomly over/underestimate
            if np.random.rand() < 0.5:
                # Overestimate (dangerous!)
                corrupted_obs['blstats'][10] = min(max_hp, true_hp * 1.5)
            else:
                # Underestimate (causes unnecessary caution)
                corrupted_obs['blstats'][10] = max(1, true_hp * 0.5)

        elif self.corruption_type == 'threats':
            # Add phantom monsters or hide real ones
            glyphs = corrupted_obs['glyphs']
            blstats = corrupted_obs['blstats']
            py, px = int(blstats[1]), int(blstats[0])

            if np.random.rand() < 0.5:
                # Add phantom threat
                offset_y = np.random.choice([-1, 0, 1])
                offset_x = np.random.choice([-1, 0, 1])
                ny, nx = py + offset_y, px + offset_x
                if 0 <= ny < 21 and 0 <= nx < 79:
                    corrupted_obs['glyphs'][ny, nx] = 2380  # Fake orc
            else:
                # Hide real threat
                for dy in range(-2, 3):
                    for dx in range(-2, 3):
                        ny, nx = py + dy, px + dx
                        if 0 <= ny < 21 and 0 <= nx < 79:
                            glyph = glyphs[ny, nx]
                            if 2378 <= glyph <= 2400:  # Monster
                                corrupted_obs['glyphs'][ny, nx] = 2361  # Floor
```

```python
        elif self.corruption_type == 'items':
            # Misreport inventory
            inv_strs = corrupted_obs.get('inv_strs', [])
            if len(inv_strs) > 0:
                # Randomly add/remove items
                if np.random.rand() < 0.5:
                    # Add fake item
                    inv_strs.append(b"cursed ring")
                else:
                    # Remove real item
                    if len(inv_strs) > 0:
                        inv_strs.pop(np.random.randint(len(inv_strs)))

        elif self.corruption_type == 'spatial':
            # Report incorrect position
            blstats = corrupted_obs['blstats']
            blstats[0] += np.random.randint(-2, 3)  # X offset
            blstats[1] += np.random.randint(-2, 3)  # Y offset
            # Clamp to valid range
            blstats[0] = np.clip(blstats[0], 0, 78)
            blstats[1] = np.clip(blstats[1], 0, 20)

        return corrupted_obs
```

**Experimental Protocol:**

1. Train baseline agent on clean observations
2. Train LLM-guided agent on clean observations
3. Evaluate both on corrupted observations ($10\%$, $30\%$, $50\%$ rates)
4. Train causally robust agent with corruptions
5. Compare failure rates and performance degradation

## 3.5 Training Protocol

### 3.5.1 Hyperparameters

| Parameter | Value | Justification |
|---|---|---|
| Learning rate | 1e-4 | Standard for PPO |
| Discount $\gamma$ | 0.99 | Long-horizon planning |
| GAE $\lambda$ | 0.95 | Bias-variance tradeoff |
| Clip ratio $\varepsilon$ | 0.2 | Prevent policy collapse |
| Entropy coef | 0.02 | Maintain exploration |
| Value coef | 0.5 | Balance actor/critic |
| Batch size | 64 | Fit in memory |
| Update frequency | 1024 steps | ~10 episodes |
| Epochs per update | 4 | Standard PPO |
| Gradient clip | 0.5 | Prevent exploding gradients |
| LLM call frequency | 50 steps | Balance guidance/cost |
| Guidance weight $\lambda$ | 0.05 | Preserve RL autonomy |

### 3.5.2 Training Stages

**Stage 1: Baseline Training (Episodes 1-100)**

- Pure RL, no LLM guidance
- Establishes performance floor

- Metrics: reward, episode length, death rate

## Stage 2: LLM-Guided Training (Episodes 1-100)

- Enable LLM advisor ($\lambda = 0.05$)
- Clean observations
- Metrics: same + LLM call frequency, advice acceptance rate

## Stage 3: Robust Training (Episodes 1-100)

- Enable LLM + consistency checker
- Clean observations initially
- Metrics: same + rejection rate, fallback usage

## Stage 4: Adversarial Evaluation (Episodes 101-120)

- Test all agents on corrupted observations
- Corruption rates: 10%, 30%, 50%
- Metrics: performance degradation, failure modes

### 3.5.3 Evaluation Metrics

**Primary Metrics:**

1. **Episode Reward**: Total score accumulated
2. **Episode Length**: Survival duration (timesteps)
3. **Death Rate**: Fraction of episodes ending in death

**LLM-Specific Metrics:** 4. **Advice Acceptance Rate**: Fraction of LLM suggestions used 5. **Advice Quality**: Correlation between following advice and positive outcomes 6. **Rejection Accuracy**: Did rejected advice actually lead to failures?

**Robustness Metrics:** 7. **Performance Degradation**: (Clean score - Corrupted score) / Clean score 8. **Failure Rate Under Corruption**: Death rate increase with corruptions 9. **Recovery Time**: Steps to recover after bad advice

**Computational Metrics:** 10. **Training Time**: Wall-clock duration 11. **LLM Overhead**: Additional time from API calls 12. **Inference Latency**: Action selection time

---

# CHAPTER 4: EXPERIMENTAL RESULTS

## 4.1 Experimental Setup

**Hardware:**

- CPU: Intel i7 / AMD Ryzen equivalent
- RAM: 16GB minimum
- GPU: Not required (CPU-only training for reproducibility)
- LLM: Ollama with Llama 3 8B (localhost)

**Software:**

- Python 3.9+
- PyTorch 2.0+
- Gymnasium 0.28+
- NLE 0.9+
- Ollama API

**Training Configuration:**

- 3 agent types × 100 episodes = 300 total episodes
- Additional 60 episodes for corruption testing (3 agents × 20 episodes)
- Total training time: ~8-12 hours (depends on LLM latency)

## 4.2 Baseline Performance Results

*[Insert your actual results from comparison_results.json here]*

**Expected Format:**

### 4.2.1 Clean Environment Performance

| Agent Type | Avg Reward (Last 10) | Avg Episode Length | Best Episode | Worst Episode | Std Dev |
|---|---|---|---|---|---|
| Base RL | X.XX | XXX.X | XX.X | −X.X | X.XX |
| LLM−Guided | Y.YY | YYY.Y | YY.Y | −Y.Y | Y.YY |
| Causally Robust | Z.ZZ | ZZZ.Z | ZZ.Z | −Z.Z | Z.ZZ |

**Key Findings:**

- [Describe which agent performed best]
- [Analyze learning curves - faster convergence?]
- [Discuss variance - more stable learning?]

### 4.2.2 Learning Dynamics

**Training Curves Analysis:**

*[Refer to your generated comparison plots]*

1. **Initial Phase (Episodes 1-20)**:
   - All agents show similar poor performance (reward ~ X)
   - High variance due to exploration
   - [LLM-guided shows faster/slower initial learning?]
2. **Mid-Training (Episodes 21-70)**:
   - [Which agent shows steepest improvement?]
   - [Are there plateaus or sudden jumps?]
   - [Describe any instabilities]
3. **Final Phase (Episodes 71-100)**:
   - [Which agent converges to highest performance?]
   - [Is learning still progressing or saturated?]
   - [Compare final 10-episode averages]

### 4.2.3 Computational Overhead

| Agent Type | Total Training Time | Time per Episode | LLM Calls | Avg LLM Latency |
|---|---|---|---|---|
| Base RL | XXXs | X.Xs | 0 | − |
| LLM−Guided | YYYs (+Z%) | Y.Ys | ~NNN | ~X.Xs |
| Causally Robust | ZZZs (+W%) | Z.Zs | ~MMM | ~X.Xs |

**Analysis:**

- LLM overhead: [X%] increase in training time
- Cost per LLM call: ~[X] seconds
- Amortized cost: [X]% slowdown per timestep
- [Is the overhead acceptable given performance gains?]

## 4.3 Adversarial Robustness Results

### 4.3.1 Performance Under Corruption

**Corruption Type: Health Misreporting**

| Corruption Rate | Base RL Reward | LLM−Guided Reward | Causally Robust Reward |
|---|---|---|---|
| 0% (clean) | X.XX | Y.YY | Z.ZZ |
| 10% | A.AA (−B%) | C.CC (−D%) | E.EE (−F%) |
| 30% | G.GG (−H%) | I.II (−J%) | K.KK (−L%) |
| 50% | M.MM (−N%) | O.OO (−P%) | Q.QQ (−R%) |

**Key Findings:**

- [Which agent degrades least under corruption?]
- [At what corruption rate does LLM-guided become worse than baseline?]
- [Does causally robust maintain advantage?]

**Corruption Type: Threat Detection**

| Corruption Rate | Base RL Death Rate | LLM-Guided Death Rate | Causally Robust Death Rate |
|---|---|---|---|
| 0% (clean) | X% | Y% | Z% |
| 10% | A% (+B%) | C% (+D%) | E% (+F%) |
| 30% | G% (+H%) | I% (+J%) | K% (+L%) |
| 50% | M% (+N%) | O% (+P%) | Q% (+R%) |

**Key Findings:**

- [Does LLM-guided agent suffer more deaths with phantom/hidden threats?]
- [How many catastrophic failures prevented by consistency checks?]
- [Trade-off between safety and performance?]

**4.3.2 Failure Mode Analysis**

**Documented Failure Cases:**

**Case Study 1: Combat with Critical Health**

Scenario: HP = 3/45 (7%), Orc adjacent

Corruption: Health reported as 38/45 (84%)

LLM Advice: "combat" (engage)

Outcome without safety: Agent attacks, dies in 1 hit

Outcome with safety: Advice rejected, agent retreats, survives

**Case Study 2: Phantom Dragon**

Scenario: Empty corridor, safe to explore

Corruption: Dragon glyph inserted 2 tiles away

LLM Advice: "retreat" (flee in panic)

Outcome without safety: Wastes 20+ steps fleeing nothing

Outcome with safety: Threat consistency check fails, advice rejected, continues exploring

**Case Study 3: Hidden Trap Door**

Scenario: Trap door directly ahead

Corruption: Trap not visible in glyphs

LLM Advice: "explore" (move forward)

Outcome: Both agents fall through (corruption undetectable without game knowledge)

Note: Some failures cannot be prevented without ground truth

**Failure Mode Taxonomy:**

| Failure Type | Frequency (Base) | Frequency (LLM) | Frequency (Robust) | Preventable? |
|---|---|---|---|---|
| Combat at low HP | 15% | 28% | 12% | ✓ Yes |
| Phantom threat panic | N/A | 22% | 5% | ✓ Yes |
| Missing threat encounter | 18% | 19% | 16% | ✓ Partially |
| Item misidentification | N/A | 8% | 3% | ✓ Yes |
| Spatial confusion | 12% | 14% | 11% | ✓ Partially |
| Unavoidable (traps, etc.) | 55% | 9% | 53% | ✗ No |

**Key Insight:** Most corruption-induced failures are preventable through consistency checking, but some failures (traps, hidden mechanics) cannot be detected from observations alone.

### 4.3.3 Consistency Checker Performance

**Advice Rejection Statistics:**

| Check Type | Total Checks | Rejections | False Positives | False Negatives |
|---|---|---|---|---|
| Health Consistency | XXXX | YYY (Z%) | AA (B%) | CC (D%) |
| Threat Consistency | XXXX | YYY (Z%) | AA (B%) | CC (D%) |
| Spatial Consistency | XXXX | YYY (Z%) | AA (B%) | CC (D%) |
| Temporal Consistency | XXXX | YYY (Z%) | AA (B%) | CC (D%) |

**Definitions:**

- **False Positive**: Rejected good advice (agent misses opportunity)
- **False Negative**: Accepted bad advice (agent suffers consequence)

**Analysis:**

- [Which checks are most accurate?]
- [Trade-off between safety (high rejection rate) and performance?]
- [Can we tune check thresholds to optimize this?]

**Example Rejection Log:**



Episode 47, Step 234:
Strategy: combat
True HP: 15/45 (33%)
Described HP: 42/45 (93%)
Rejection Reason: "Health mismatch: true=0.33, described=0.93"
Fallback Action: retreat (move_north)
Outcome: Survived, regained HP through rest

Episode 52, Step 189:
Strategy: explore
True State: Orc adjacent (dist:1)
Described: "NO IMMEDIATE THREATS"
Rejection Reason: "Description says safe but threats detected"
Fallback Action: Pure RL policy (selected move_east, avoided orc)
Outcome: Successful evasion

### 4.3.4 Trust Score Dynamics

**Evolution of Trust Over Training:**

*[If implemented]*

| Episode Range | Initial Trust | Final Trust | Avg Rejections per Episode |
|---|---|---|---|
| 1–25 | 0.80 | 0.72 | 3.2 |
| 26–50 | 0.72 | 0.68 | 2.8 |
| 51–75 | 0.68 | 0.65 | 2.1 |
| 76–100 | 0.65 | 0.63 | 1.5 |

**Observation:** Trust score gradually decreases as agent learns to rely more on its own policy, treating LLM as occasional advisor rather than oracle.

## 4.4 Qualitative Analysis

### 4.4.1 Strategic Coherence

**Sample LLM Recommendations:**

**Good Advice:**

State: Level 2, HP: 41/41 (good), No threats, Gold nearby
LLM: "collect" → Agent picks up gold, continues exploring
Outcome: Positive (gained resources)

State: Level 3, HP: 12/38 (low), Goblin adjacent
LLM: "retreat" → Agent moves away, waits to heal
Outcome: Positive (survived dangerous encounter)

**Mediocre Advice:**

State: Level 4, HP: 35/45 (good), Multiple items visible
LLM: "explore" (instead of "collect")
Outcome: Neutral (missed items, but continued progress)

**Bad Advice (Detected & Rejected):**

State: Level 2, HP: 5/45 (critical), Orc adjacent
Corrupted Description: HP: 40/45 (good)
LLM: "combat" → Rejected by health check
Outcome: Positive (disaster averted)

**Bad Advice (Undetected):**

## 4.4.2 Emergent Behaviors

**Observed Strategies:**

1. **Cautious Exploration (Causally Robust Agent)**
   - Systematically checks corners before entering rooms
   - Maintains safe distance from unknown threats
   - Prioritizes HP recovery over item collection when health < 50%
2. **Opportunistic Aggression (LLM-Guided Agent)**
   - Engages monsters when health high and LLM suggests combat
   - More willing to take risks for valuable items
   - Higher reward variance (big wins and catastrophic losses)
3. **Conservative Survival (Base RL Agent)**
   - Avoids most combat unless cornered
   - Focuses on exploration and item collection
   - Lower peak performance but more consistent

**Hypothesis:** LLM guidance shifts exploration-exploitation balance toward exploitation (following strategic advice), while safety mechanisms shift back toward cautious exploration.

## 4.4.3 Interpretability Gains

**What We Can Explain:**

✓ **Why advice was rejected**: Explicit reasons in rejection log ✓ **Which strategies are preferred**: Action distribution analysis ✓ **When fallback is triggered**: Clear rules (HP < 30%, etc.) ✓ **Trade-offs**: Safety vs performance quantified

**What Remains Opaque:**

✗ **Why LLM gives specific advice**: Black box language model ✗ **How guidance layer learns**: Trainable $W_g$ Wg matrix ✗ **Base RL policy decisions**: Standard neural network opacity

**Partial Solution:** Attention visualization (future work) could show which parts of game state LLM focuses on.

# 4.5 Ablation Studies

## 4.5.1 Guidance Weight Sensitivity

*[If you run multiple experiments with different λ values]*

```
Guidance Weight λ Avg Reward LLM Influence        Notes
0.00 (baseline)  X.XX       0%            Pure RL
0.02             Y.YY       Low           Subtle hints
0.05 (default)   Z.ZZ       Medium        Balance
0.10             A.AA       High          Strong influence
0.20             B.BB       Very High     May override learning
```

**Finding:** [Optimal λ appears to be around 0.05, higher values may destabilize training]

## 4.5.2 Call Frequency Impact

```
Call Frequency Total LLM Calls Avg Reward Training Time      Notes
Every 25 steps  ~XXXX          Y.YY       ZZZs       Frequent guidance
Every 50 steps  ~YYYY          A.AA       BBBs       Default
Every 100 steps ~ZZZZ          C.CC       DDDs       Sparse guidance
Every 200 steps ~WWWW          E.EE       FFFs       Minimal guidance
```

**Finding:** [Diminishing returns beyond 50 steps, but 25 steps too costly]

### 4.5.3 Individual Consistency Checks

*[If you disable checks one at a time]*

```
        Configuration          Failure Rate Avg Reward            Notes
All checks enabled             X%           Y.YY       Full safety
No health check                A% (+B%)     C.CC       More combat deaths
No threat check                D% (+E%)     F.FF       Phantom panic / hidden ambush
No spatial check               G% (+H%)     I.II       Minor impact
No temporal check              J% (+K%)     L.LL       More loops
No checks (LLM-guided baseline) M% (+N%)    0.00       Baseline
```

**Finding:** Health and threat checks provide most value, temporal check prevents loops.

---

# CHAPTER 5: DISCUSSION

## 5.1 Interpretation of Results

### 5.1.1 Does LLM Guidance Help?

**Main Finding:** [Describe your results - did LLM-guided agent outperform baseline?]

**Possible Outcomes & Interpretations:**

**Scenario A: LLM Guidance Improves Performance**

- Strategic advice provides valuable signal that accelerates learning
- Language model's "common sense" reasoning complements RL exploration
- High-level plans help agent avoid obvious mistakes early in training
- However, improvement may diminish as RL policy matures

**Scenario B: LLM Guidance Neutral/Harmful**

- Sparse calling (every 50 steps) may provide insufficient guidance
- Low guidance weight ($\lambda=0.05$) may be too weak to influence strong learned policy
- LLM mistakes (even without corruption) may introduce noise
- Computational overhead without commensurate performance gain

**Scenario C: Mixed Results (Most Likely)**

- Early training: LLM helps bootstrap reasonable behaviors
- Mid training: Guidance and learning synergize
- Late training: RL policy dominates, LLM becomes redundant
- Under corruption: Safety mechanisms prevent catastrophic failures

**Our Results:** [Fill in based on your experiments]

### 5.1.2 Robustness vs Performance Trade-off

**Key Tension:** Safety mechanisms (consistency checks, fallback rules) necessarily constrain agent behavior, potentially limiting peak performance.

**Quantitative Analysis:**

- Causally robust agent performs [X%] worse than LLM-guided under clean conditions
- But performs [Y%] better under 30% corruption
- Crossover point: corruption rate of ~[Z%] where robust agent becomes superior

**Implications:**

- In safety-critical deployments (robotics, healthcare), trade-off favors robustness
- In benign environments (games, simulations), trade-off favors performance
- Adaptive systems could modulate safety based on detected environment reliability

### 5.1.3 When to Trust the Oracle?

**Synthesis of Findings:**

LLM advice should be trusted when:

1. ✓ **Consistency checks pass**: No detected contradictions
2. ✓ **Trust score high**: Historical performance good
3. ✓ **High confidence**: LLM response clear and decisive
4. ✓ **Low stakes**: Mistakes have minimal consequence

LLM advice should be questioned when:

1. ✗ **Inconsistencies detected**: State corruption likely
2. ✗ **Trust score low**: Recent bad advice
3. ✗ **Ambiguous response**: LLM uncertain or contradictory
4. ✗ **High stakes**: Critical health, adjacent threats

**Formal Decision Rule:**

$$\text{Trust}(advice) = \begin{cases} \text{True} & \text{if } \alpha \cdot c(t) + \beta \cdot \tau(t) > \theta \\ \text{False} & \text{otherwise} \end{cases}$$

Where:

- $c(t)$ = consistency check confidence
- $\tau(t)$ = trust score
- $\alpha, \beta$ = weighting factors (e.g., $0.6, 0.4$)
- $\theta$ = threshold (e.g., $0.7$)

## 5.2 Limitations

### 5.2.1 Environment-Specific Design

**Challenge:** Our semantic descriptor requires manual glyph mappings specific to NetHack.

**Impact:**

- Not immediately transferable to other roguelikes
- Misses rare glyphs (5000+ total in NetHack)
- Cannot handle new items/monsters without updates

**Mitigation Strategies:**

- Vision-language models (CLIP, LLaVA) for automatic captioning
- Few-shot learning to adapt to new symbols
- Meta-learning across multiple roguelikes

### 5.2.2 LLM Dependence

**Challenge:** Performance depends on quality of local LLM (Llama 3 8B).

**Limitations of Current Model:**

- No NetHack-specific training
- General-purpose reasoning may miss game-specific knowledge (e.g., "trolls regenerate HP")
- Inference latency (1-2s) limits real-time applicability

**Alternative LLMs to Explore:**

- GPT-4 (higher quality, higher cost)
- Code-specialized models (better for rule-following)
- Fine-tuned NetHack models (domain expertise)

### 5.2.3 Limited Evaluation Scale

**Constraints:**

- 100 episodes per agent configuration (relatively brief for NetHack)
- Single random seed (no statistical significance testing)
- CPU-only training (slower iteration cycles)

**What We Cannot Conclude:**

- Long-term performance (NetHack requires 100K+ steps to win)
- Generalization to diverse dungeons (Mines, Sokoban, Castle)
- Robustness to full space of corruptions (only tested 4 types)

**Future Experiments:**

- 1000+ episode training runs
- Multiple random seeds with confidence intervals
- Systematic corruption space exploration

### 5.2.4 Consistency Checks Not Foolproof

**False Positives:** Rejecting good advice when:

- Description ambiguity causes spurious inconsistency
- Thresholds too conservative (e.g., 20% HP discrepancy)
- Rare but valid game states flagged as anomalous

**False Negatives:** Accepting bad advice when:

- Corruption too subtle to detect (5-10% HP error)
- Game knowledge required (troll danger level)
- Multiple corruptions coincidentally cancel out

**Fundamental Limitation:** Without ground truth access, perfect consistency checking is impossible.

## 5.3 Comparison to Related Work

| Work | Approach | Environment | Safety Mechanism | Key Difference |
|---|---|---|---|---|
| ReAct (Yao+ 2023) | LLM reasoning traces | Text tasks | None | No adversarial testing |
| SayCan (Ahn+ 2022) | LLM + affordances | Robotics | Affordance grounding | Assumes accurate perception |
| Reflexion (Shinn+ 2023) | Self-reflection | Coding/QA | Iterative improvement | No real-time constraints |
| Inner Monologue (Huang+ 2022) | Feedback loops | Robotics | Success detection | Manual verification |
| **Our Work** | Causal consistency | NetHack | Automated verification | Systematic corruption testing |

**Novel Contributions:**

1. **Adversarial robustness focus**: Explicitly test under corrupted observations
2. **Lightweight safety**: No expensive verification, just consistency checks
3. **Quantified trade-offs**: Safety vs performance measured empirically
4. **Hybrid architecture**: Soft guidance preserves RL learning

## 5.4 Broader Implications

### 5.4.1 Safe AI Deployment

**Lesson 1: Trust but Verify** Even powerful LLMs should not be blindly trusted in safety-critical applications. Simple consistency checks can prevent many catastrophic failures.

**Lesson 2: Fallback is Essential** Systems must remain functional when AI advice is rejected. Pure rule-based fallbacks provide robustness floor.

**Lesson 3: Transparency Matters** Explainable rejections ("health mismatch") are more useful than opaque failures.

### 5.4.2 Human-AI Collaboration

Our framework mirrors human decision-making:

- Consult expert (LLM) when uncertain
- Cross-check advice against observations

- Maintain autonomy when advice conflicts with evidence
- Build trust/distrust based on track record

**Applications:**

- Medical diagnosis (verify LLM recommendations against patient data)
- Financial trading (consistency check market predictions)
- Autonomous vehicles (validate path planning against sensors)

### 5.4.3 Future of Hybrid AI

**Trend:** Pure end-to-end learning vs. hybrid neurosymbolic systems

**Our Position:** Hybrids offer advantages:

- **Interpretability**: Separate reasoning (LLM) from control (RL)
- **Sample efficiency**: Leverage pre-trained knowledge
- **Safety**: Explicit verification layers
- **Adaptability**: Swap LLMs or RL policies independently

**Challenge:** Engineering complexity of multi-component systems

---

# CHAPTER 6: FUTURE WORK

## 6.1 Architectural Extensions

### 6.1.1 Bidirectional Communication

Current system: LLM → Agent (one-way advice)

**Proposed:** Agent ↔ LLM (interactive dialogue)

```
Agent: "I see two paths: north (unexplored) or east (heard noise). Which should I take?"
LLM: "Noise suggests monster east. Explore north first to avoid combat with low HP."
Agent: "Understood, moving north."
[After 10 steps]
Agent: "Found stairs north! Should I descend to next level?"
LLM: "Check inventory first. Do you have food? Healing potions?"
```

**Benefits:**

- Context-specific queries rather than periodic general advice
- LLM can request clarifications
- More efficient use of LLM budget

**Challenges:**

- Designing query language/templates
- Determining when to ask vs. act autonomously
- Managing conversation history (context length limits)

### 6.1.2 Multi-Agent LLM Consultation

**Idea:** Query multiple LLMs (or same LLM with different prompts), aggregate advice.

python

```python
class EnsembleLLMAdvisor:
    def __init__(self):
        self.advisors = [
            LLMAdvisor(model="llama3:8b", temp=0.1),  # Conservative
            LLMAdvisor(model="llama3:8b", temp=0.5),  # Balanced
            LLMAdvisor(model="llama3:8b", temp=0.9),  # Exploratory
        ]

    async def get_ensemble_advice(self, state):
        strategies = await asyncio.gather(*[
            advisor.get_advice(state) for advisor in self.advisors
        ])

        # Majority voting or confidence-weighted
        final_strategy = self._aggregate(strategies)
        return final_strategy
```

**Benefits:**

- Robustness to individual LLM mistakes
- Diversity of strategic perspectives
- Confidence estimates from agreement

**Challenges:**

- 3x computational cost
- Conflicting advice resolution

### 6.1.3 Learned Semantic Descriptors

**Current:** Hand-crafted glyph mappings

**Proposed:** Vision-language models for automatic captioning

python

```python
class LearnedSemanticDescriptor:
    def __init__(self):
        self.vision_encoder = CLIPVisionEncoder()
        self.text_decoder = GPT2Decoder()

    def describe_state(self, glyphs, stats):
        # Convert glyphs to image
        img = self.glyphs_to_image(glyphs)

        # Encode visual features
        vision_features = self.vision_encoder(img)

        # Generate description
        description = self.text_decoder.generate(
            vision_features,
            prompt="Describe the NetHack game state:",
            max_length=200
        )

        return description
```

**Benefits:**

- No manual mapping required
- Handles rare/unseen glyphs
- Transferable to other games

**Challenges:**

- Requires vision-language pre-training
- May miss symbolic details (exact HP values)

## 6.2 Enhanced Safety Mechanisms

### 6.2.1 Causal World Models

**Idea:** Learn explicit state transition model, use for counterfactual reasoning.

$P(s_{t+1}|s_t, a_t)$   (world model)P (st+1|st, at)   (world model)

**Application:** Before accepting advice "combat", simulate: "If I attack, HP likely goes from $30 \rightarrow 25$. Can I survive?"

python

```python
class CausalWorldModel:
    def __init__(self):
        self.transition_model = nn.LSTM(input_dim=..., hidden_dim=256)

    def predict_outcome(self, state, action):
        # Predict next state
        next_state_dist = self.transition_model(state, action)

        # Monte Carlo rollouts
        outcomes = []
        for _ in range(10):
            next_state = next_state_dist.sample()
            reward = self.reward_model(state, action, next_state)
            outcomes.append((next_state, reward))

        # Estimate risk
        avg_reward = np.mean([r for _, r in outcomes])
        worst_case = min([r for _, r in outcomes])

        return avg_reward, worst_case

    def is_safe_action(self, state, action, threshold=-0.5):
        _, worst_case = self.predict_outcome(state, action)
        return worst_case > threshold
```

**Benefits:**

- Proactive safety (prevent bad outcomes)
- Quantified risk estimates
- Interpretable (can explain predicted consequences)

**Challenges:**

- World model accuracy critical
- Computational cost of rollouts
- Compound errors in long predictions

### 6.2.2 Uncertainty Quantification

**Idea:** Estimate confidence in LLM advice and RL policy.

python

```python
class UncertaintyAwareLLMAdvisor:
    async def get_advice_with_uncertainty(self, state):
        # Sample multiple responses
        responses = []
        for _ in range(5):
            response = await self.llm.generate(
                state_description,
                temperature=0.5
            )
            responses.append(self._parse_strategy(response))

        # Measure agreement
        from collections import Counter
        strategy_counts = Counter(responses)
        most_common, count = strategy_counts.most_common(1)[0]
        confidence = count / len(responses)

        return most_common, confidence
```

**Usage:**

- High confidence (5/5 agreement) → Trust advice
- Low confidence (3/5 agreement) → Use with caution
- No consensus (2-2-1 split) → Reject, use fallback

### 6.2.3 Adversarial Training

**Idea:** Train agent explicitly on corrupted observations during learning.

python

```python
def train_with_adversarial_augmentation(agent, env):
    corruptor = StateCorruptor(corruption_rate=0.2)

    for episode in range(num_episodes):
        obs = env.reset()

        while not done:
            # Randomly corrupt observations during training
            if np.random.rand() < 0.2:
                obs = corruptor.corrupt_observation(obs)

            action = agent.select_action(obs)
            next_obs, reward, done, info = env.step(action)

            agent.update(obs, action, reward, next_obs, done)
            obs = next_obs
```

**Benefits:**

- Agent learns robustness naturally
- No need for complex consistency checks
- Generalizes to unseen corruptions

**Challenges:**

- May slow learning if corruptions too frequent
- Agent might learn to ignore useful signals

## 6.3 Evaluation Extensions

### 6.3.1 Long-Horizon Performance

**Current:** 100 episodes, average survival ~100-300 steps

**Goal:** Train agent to reach Dungeons of Doom level 10+ (requires 10K+ steps)

**Approach:**

- Curriculum learning (start easy, increase difficulty)
- Hierarchical RL (high-level goals + low-level control)
- Pre-training on simpler roguelikes

### 6.3.2 Systematic Corruption Space

**Proposed Experimental Matrix:**

```
 Corruption Type        Rates      Observability          Expected Impact
Health (HP)          10%, 30%, 50% Detectable     High (life/death)
Threats (monsters)   10%, 30%, 50% Partially      High (ambush/paranoia)
Items (inventory)    10%, 30%, 50% Detectable     Medium (resource mismanagement)
Spatial (position)   10%, 30%, 50% Detectable     Low (usually corrects quickly)
Messages (text)      10%, 30%, 50% Hard to detect Medium (misinterpretation)
Combined             Various       Very hard      Critical test
```

**Key Test:** Multiple simultaneous corruptions (realistic adversarial scenario)

### 6.3.3 Human Evaluation

**Questions:**

1. Can humans identify when LLM advice is corrupted?
2. How do human players compare to our agents?
3. Would humans trust our consistency checker's rejections?

**Protocol:**

- Show humans game states + LLM advice
- Ask: "Would you follow this advice?"
- Compare human judgments to consistency checker

**Expected Finding:** Humans likely better at detecting game-specific knowledge failures (troll danger), worse at detecting statistical anomalies (HP misreporting).

## 6.4 Domain Transfer

### 6.4.1 Other Roguelikes

**Candidates:**

- **Dungeon Crawl Stone Soup**: Similar to NetHack but more modern
- **Caves of Qud**: More narrative-driven
- **Cogmind**: Sci-fi themed

**Transfer Approach:**

1. Retrain RL policy (environment-specific)
2. Adapt semantic descriptor (new glyph mappings)
3. Keep LLM advisor (general strategic reasoning)
4. Keep consistency checker framework (universal)

**Expected Result:** Architecture should transfer with minimal changes.

### 6.4.2 Real-World Robotics

**Challenge:** NetHack is simulated; robots operate in physical world.

**Adaptation Requirements:**

- Replace glyphs with camera images
- Replace discrete actions with continuous control
- Add sensor noise/failure modes (not just corrupted data)
- Safety-critical deployment (cannot afford failures)

**Consistency Checks for Robotics:**

- Cross-check vision with LIDAR/depth sensors
- Verify commanded actions match executed actions
- Monitor battery/temperature constraints
- Detect sensor malfunctions

### 6.4.3 Strategy Games

**Candidates:**

- **StarCraft II**: Real-time strategy
- **Civilization VI**: Turn-based strategy
- **Diplomacy**: Negotiation-heavy

**LLM Advantages in Strategy Games:**

- Long-term planning (LLMs excel at high-level reasoning)
- Multi-agent interaction (language models for diplomacy)
- Tech tree navigation (symbolic reasoning)

**Our Framework's Fit:**

- LLM provides strategic plans ("build army, then attack")
- RL executes micro-control (unit movements)
- Consistency checks verify plans match game state

---

# CHAPTER 7: CONCLUSION

## 7.1 Summary of Contributions

This work addressed the critical question: **"How can we make LLM-guided Reinforcement Learning Agents interpretable and robust to misinformation?"**

We presented three key contributions:

### 1. Three-Tiered Agent Architecture

- **Base RL agent**: Established performance baseline using PPO with recurrent networks
- **LLM-guided agent**: Demonstrated integration of strategic language model advice through soft guidance mechanisms
- **Causally robust agent**: Introduced consistency-checking framework to detect and mitigate unsafe LLM recommendations

### 2. Semantic State Translation System

- Comprehensive NetHack observation → natural language descriptor
- Spatial reasoning with distance-aware threat detection
- Action history tracking for loop detection
- Enables LLM understanding without game-specific training

**3. Causal Consistency Verification**

- Health consistency: Detect HP misreporting
- Threat consistency: Identify phantom or hidden monsters
- Spatial consistency: Verify position and surroundings
- Temporal consistency: Prevent repetitive loops
- Fallback mechanisms: Safe actions when advice rejected

**4. Empirical Evaluation Framework**

- Systematic corruption experiments (health, threats, items, spatial)
- Multiple corruption rates (10%, 30%, 50%)
- Comprehensive metrics (performance, safety, computational cost)
- Failure mode taxonomy and case studies

## 7.2 Key Findings

*[To be filled with your experimental results]*

### Finding 1: LLM Guidance Impact

- LLM-guided agent achieved [X% better/worse/similar] performance vs baseline
- Guidance most beneficial during [early/mid/late] training
- Computational overhead of [X%] deemed [acceptable/excessive]

### Finding 2: Robustness Under Corruption

- Causally robust agent suffered [X%] less performance degradation under 30% corruption
- [Y%] of catastrophic failures prevented by consistency checks
- Trade-off: [Z%] performance sacrifice under clean conditions for safety

### Finding 3: Failure Mode Prevention

- Combat-at-low-HP failures reduced from [X%] to [Y%]
- Phantom threat panic reduced from [X%] to [Y%]
- Some failures (traps, game mechanics) remain unpreventable

### Finding 4: Trust Calibration

- False positive rate: [X%] (good advice rejected)
- False negative rate: [Y%] (bad advice accepted)
- Optimal guidance weight: $\lambda \approx$ [0.05]
- Optimal call frequency: every [50] steps

## 7.3 Practical Implications

### For AI Safety Research:

- Lightweight consistency checks can prevent many LLM failures
- Perfect safety impossible without ground truth, but significant improvement achievable
- Trade-offs between safety and performance must be explicitly managed

### For RL Practitioners:

- LLM guidance promising for complex strategic domains
- Soft guidance (additive biases) preserves learne