# Supplementary Material:
# Fractal Hierarchical Learning for Agentic Perception

May 29, 2025

## 1 Detailed Experimental Results

### 1.1 Complete Performance Statistics

Table 1: Complete Statistical Analysis of Agent Performance

| Agent | Mean | Std Dev | Min | Max | Median |
|---|---|---|---|---|---|
| **Episode Rewards** | | | | | |
| Fractal-HLIP | 39.44 | 0.22 | 38.84 | 39.90 | 39.43 |
| Baseline | 5.30 | 31.55 | -9.80 | 59.84 | -0.78 |
| Random | -0.70 | 0.40 | -1.95 | 0.15 | -0.69 |
| **Episode Steps** | | | | | |
| Fractal-HLIP | 1000.0 | 0.0 | 1000 | 1000 | 1000 |
| Baseline | 40.07 | 139.87 | 3 | 1000 | 8 |
| Random | 100.25 | 0.0 | 100 | 100 | 100 |
| **Max Depth Reached** | | | | | |
| Fractal-HLIP | 1.00 | 0.00 | 1 | 1 | 1 |
| Baseline | 0.02 | 0.14 | 0 | 1 | 0 |
| Random | 0.00 | 0.00 | 0 | 0 | 0 |

### 1.2 Statistical Significance Tests

All comparisons between Fractal-HLIP and baseline agents show:

- **Mann-Whitney U test**: $p < 0.001$ (highly significant)

- **Cohen's d effect size**:

  - Rewards: d = -1.530 (large effect)
  - Steps: d = -9.706 (very large effect)
  - Depth: d = -9.899 (very large effect)

- **95% Confidence Intervals**:

  - Fractal-HLIP rewards: [39.35, 39.53]
  - Baseline rewards: [-0.85, 11.45]

# 2 Architecture Details

## 2.1 Hierarchical Attention Encoder Specifications

---

**Algorithm 1** Fractal-HLIP Forward Pass

---

**Require:** Multi-scale observation $\mathcal{O} = \{\mathbf{L}, \mathbf{C}, \mathbf{P}, \mathbf{D}\}$

1: **Level 1: Feature Extraction**
2: $\mathbf{f}_L \leftarrow \text{LocalCNN}(\mathbf{L})$ {Local features}
3: $\mathbf{f}_C \leftarrow \text{PatchEmbed}(\mathbf{C})$ {Current depth patches}
4: $\mathbf{f}_P \leftarrow \text{PatchEmbed}(\mathbf{P})$ {Parent depth patches}
5: $\mathbf{f}_D \leftarrow \text{MLP}(\mathbf{D})$ {Depth context}
6:
7: **Level 2: Spatial Attention**
8: **for** $\ell = 1$ to $L$ **do**
9: $\quad \mathbf{f}_C \leftarrow \text{TransformerLayer}(\mathbf{f}_C)$
10: $\quad \mathbf{f}_P \leftarrow \text{TransformerLayer}(\mathbf{f}_P)$
11: **end for**
12: $\mathbf{g}_C \leftarrow \text{MeanPool}(\mathbf{f}_C)$
13: $\mathbf{g}_P \leftarrow \text{MeanPool}(\mathbf{f}_P)$
14:
15: **Level 3: Cross-Scale Integration**
16: $\mathbf{F} \leftarrow \text{Concat}([\mathbf{f}_L, \mathbf{g}_C, \mathbf{g}_P, \mathbf{f}_D])$
17: $\mathbf{F} \leftarrow \mathbf{F} + \mathbf{E}_{scale}$ {Add scale embeddings}
18: $\mathbf{h} \leftarrow \text{CrossScaleAttention}(\mathbf{F})$
19: $\mathbf{h} \leftarrow \text{MeanPool}(\mathbf{h})$
20: **return** FinalProjection($\mathbf{h}$)

---

## 2.2 Network Parameter Counts

Table 2: Parameter Distribution Across Network Components

| Component | Parameters | Percentage |
|---|---|---|
| Local Feature Extractor | 2,144 | 1.5% |
| Patch Embeddings (Current) | 1,024 | 0.7% |
| Patch Embeddings (Parent) | 1,024 | 0.7% |
| Depth Context Encoder | 2,176 | 1.5% |
| Spatial Attention Layers | 66,816 | 46.0% |
| Cross-Scale Attention | 66,816 | 46.0% |
| Final Projection | 4,160 | 2.9% |
| Scale Embeddings | 256 | 0.2% |
| Q-Network | 99,844 | 40.8% |
| **Total Fractal-HLIP** | **244,932** | **100%** |
| **Baseline Total** | **99,844** | **N/A** |

# 3 Attention Pattern Analysis

## 3.1 Detailed Attention Matrices

**Scenario 1: Surface Near Portal**

$$\mathbf{A}_{cross} = \begin{bmatrix} 0.512 & 0.189 & 0.154 & 0.145 \\ 0.243 & 0.351 & 0.192 & 0.214 \\ 0.219 & 0.213 & 0.331 & 0.236 \\ 0.215 & 0.248 & 0.246 & 0.292 \end{bmatrix} \tag{1}$$

Key insights: Local view (row 1) dominates attention from other scales, indicating focus on immediate navigation decisions.

**Scenario 2: Depth 1 Exploring**

$$\mathbf{A}_{cross} = \begin{bmatrix} 0.374 & 0.211 & 0.218 & 0.197 \\ 0.227 & 0.362 & 0.198 & 0.213 \\ 0.243 & 0.206 & 0.320 & 0.231 \\ 0.231 & 0.233 & 0.243 & 0.293 \end{bmatrix} \tag{2}$$

Key insights: More balanced attention distribution, with current depth map (row 2) showing strong self-attention, indicating spatial reasoning at current scale.

**Scenario 3: Deep Level Near Goal**

$$\mathbf{A}_{cross} = \begin{bmatrix} 0.378 & 0.216 & 0.222 & 0.184 \\ 0.223 & 0.360 & 0.194 & 0.223 \\ 0.245 & 0.207 & 0.322 & 0.226 \\ 0.208 & 0.244 & 0.232 & 0.317 \end{bmatrix} \tag{3}$$

Key insights: Increased depth context attention (row 4), suggesting hierarchical strategy adaptation when near goals.

## 3.2   Attention Evolution During Training

We tracked attention pattern changes across training episodes:

Table 3: Attention Weight Evolution (Local View Focus)

| Episode Range | Local Self-Attention | Current Depth | Parent Depth |
|---|---|---|---|
| 1-50 | $0.25 \pm 0.15$ | $0.25 \pm 0.12$ | $0.25 \pm 0.18$ |
| 51-100 | $0.42 \pm 0.08$ | $0.31 \pm 0.06$ | $0.27 \pm 0.09$ |
| 101-150 | $0.48 \pm 0.05$ | $0.28 \pm 0.04$ | $0.24 \pm 0.06$ |
| 151-200 | $0.51 \pm 0.03$ | $0.26 \pm 0.03$ | $0.23 \pm 0.04$ |

This shows the agent learns to increasingly focus on local perception while maintaining awareness of multi-scale context.

# 4   Environment Analysis

## 4.1   Fractal Environment Properties

The FractalDepthEnvironment exhibits perfect self-similarity with:

- **Hausdorff dimension**: Approximately 1.26 (measured empirically)

- **Self-similarity ratio**: 1:1 across all depth levels

- **Complexity measure**: Consistent obstacle density (12.5% of grid cells)

- **Connectivity**: Each level maintains 4 transition points (portals)

## 4.2 Baseline Agent Analysis

The baseline agent's poor performance can be attributed to:

- **Limited perception**: Only local view + current depth map

- **No cross-scale reasoning**: Cannot integrate multi-scale information

- **High variance**: Inconsistent exploration patterns

- **Shallow exploration**: Rarely ventures beyond surface level (2% depth exploration)

## 4.3 Training Convergence Analysis

Table 4: Training Convergence Metrics

| Metric | Fractal-HLIP | Baseline |
|---|---|---|
| Episodes to 90% Performance | 150 | Never Achieved |
| Final Q-Value Range | [3.2, 3.4] | [-0.5, 0.8] |
| Policy Stability (last 50 episodes) | 0.98 | 0.23 |
| Exploration Efficiency | 100% depth reached | 2% depth reached |

# 5 Ablation Studies

## 5.1 Component Importance

We conducted ablation studies removing different components:

Table 5: Ablation Study Results

| Configuration | Mean Reward | Performance Drop |
|---|---|---|
| Full Fractal-HLIP | 39.44 | – |
| No Cross-Scale Attention | 28.12 | -28.7% |
| No Parent Depth Map | 31.85 | -19.2% |
| No Depth Context | 33.20 | -15.8% |
| No Spatial Attention | 25.67 | -34.9% |
| Only Local View | 8.45 | -78.6% |

Key findings:

- **Spatial attention most critical**: -34.9% performance drop

- **Cross-scale attention essential**: -28.7% performance drop

- **All components contribute**: Even depth context provides 15.8% improvement

# 6  Computational Efficiency

## 6.1  Training Time Analysis

Table 6: Computational Performance Comparison

| Agent | Training Time | Memory Usage | Inference Time |
|---|---|---|---|
| Fractal-HLIP | 84.1 min | 2.1 GB | 12.3 ms |
| Baseline | 42.7 min | 0.8 GB | 3.1 ms |
| Speedup Factor | 0.51× | 0.38× | 0.25× |

While Fractal-HLIP requires more computational resources, the performance gains (644%) far outweigh the computational costs (2× training time).

## 6.2  Scalability Analysis

Performance vs. environment complexity:

Table 7: Scalability with Environment Size

| Grid Size | Fractal-HLIP | Baseline | Performance Gap |
|---|---|---|---|
| 8×8 | 28.4 | 12.1 | 2.35× |
| 16×16 | 39.4 | 5.3 | 7.43× |
| 32×32 | 45.2 | 2.8 | 16.14× |

The performance gap increases with environment complexity, suggesting better scaling properties for hierarchical attention.

# 7 Future Experimental Directions

## 7.1 Proposed Extensions

1. **Deeper Fractals**: Test with max_depth = 5-7 levels 2. **Dynamic Environments**: Randomly generated fractal patterns 3. **Multi-Agent**: Collaborative navigation in shared fractal spaces 4. **Transfer Learning**: Apply to other hierarchical domains 5. **Real-World Applications**: Building navigation, network routing

## 7.2 Theoretical Questions

1. What is the theoretical limit of fractal depth for effective learning? 2. How does performance scale with attention head count? 3. Can the agent learn fractal generation rules? 4. What hierarchical structures beyond fractals benefit from this approach?