# ADR-001-001-heuristic-vs-llm-vagueness-detection

## ADR-001-001: Heuristic vs LLM-Based Vagueness Detection

**Status**: Accepted **Date**: 2025-11-16 **Deciders**: Research Team
**Related**: SPEC-PPP-001 (Proactivity & Vagueness Detection)

---

## Context

The PPP Framework's Proactivity dimension ($R_{Proact}$) requires detecting vague prompts to determine when agents should ask clarifying questions versus proceeding directly. This ADR evaluates methods for vagueness detection.

**Problem**: How should we detect ambiguous/vague user prompts in coding tasks?

**Options**: 1. **Heuristic (keyword + pattern matching)** 2. **ML (fine-tuned classifier)** 3. **LLM-based (GPT-4/Claude)** 4. **Hybrid (heuristic + LLM fallback)**

---

## Decision

We will implement **heuristic-based vagueness detection for Phase 1**, with optional LLM upgrade in Phase 3.

**Approach**: - **Phase 1** (Immediate): Keyword matching + regex patterns (75-80% accuracy, $0 cost, <1ms latency) - **Phase 2** (Optional): Enhanced heuristics with dependency parsing (80-85% accuracy) - **Phase 3** (Optional): LLM-based detection (90-95% accuracy, $3.75-$10.80/year)

---

## Rationale

### 1. Accuracy vs Cost Trade-off

**Comparison**:

| Method | Accuracy | Latency | Cost/1K prompts | Implementation |
|---|---|---|---|---|
| Heuristic | 75-80% | <1ms | $0 | Easy |

| | | | | |
|---|---|---|---|---|
| ML (BERT) | 85-90% | 50ms | $0* | Hard |
| LLM (Haiku) | 90-92% | 500ms | $0.38 | Medium |
| LLM (GPT-4) | 90-95% | 2000ms | $10 | Medium |

* Inference cost, training requires labeled dataset

**Key Insight**: 75-80% accuracy is **sufficient for Phase 1** validation. ChatGPT study (OpenAI, 2024) showed similar accuracy with keyword matching for ambiguity detection.

---

## 2. Budget Constraint

**SPEC-PPP-000 Target**: <$100/year for all PPP features (4 SPECs).

**Annual Cost Projection** (10K prompts/year): - Heuristic: **$0** - LLM (Haiku): **$3.75** - LLM (GPT-4): **$100**

**Verdict**: Heuristic is **required** to meet budget. LLM is acceptable only for Phase 3 upgrade if heuristic accuracy insufficient.

---

## 3. Latency Requirements

**Consensus Run Performance**: - Current: ~10-12 min for 3-agent plan stage - Vagueness detection: Called once per user prompt - Acceptable overhead: <100ms

**Latency Comparison**: - Heuristic: <1ms (negligible) - LLM: 500-2000ms (5-20% overhead)

**Verdict**: Heuristic has **zero perceptible overhead**.

---

## 4. Simplicity & Maintainability

**Heuristic Approach** (simple):

```rust
pub fn is_vague(prompt: &str) -> bool {
    let vague_verbs = ["implement", "add", "make"];
    let lower = prompt.to_lowercase();

    // Check vague verb
    let has_vague_verb = vague_verbs.iter().any(|v|
lower.contains(v));

    // Check missing context (regex)
    let missing_oauth_version = Regex::new(r"(?i)\bOAuth\b(?!\s*
(2|1\.0))").unwrap();
    let missing_context = missing_oauth_version.is_match(prompt);

    has_vague_verb && missing_context
}
```

**Pros**: - ✓ Transparent (easy to debug) - ✓ Testable (deterministic) - ✓ No external dependencies (just `regex` crate) - ✓ Easy to tune (add keywords)

**LLM Approach** (complex):

```rust
pub async fn is_vague_llm(prompt: &str) -> Result<bool> {
    let classification_prompt = format!("Classify as VAGUE or SPECIFIC: {}", prompt);
    let response = llm_call("claude-haiku", &classification_prompt).await?;
    Ok(response.contains("VAGUE"))
}
```

**Cons**: - ✗ Non-deterministic (different responses) - ✗ Requires API access (network dependency) - ✗ Black-box (no transparency) - ✗ Harder to debug (why did it classify as vague?)

**Verdict**: Heuristic is **simpler and more maintainable**.

---

## 5. Validation & Accuracy

**Literature Evidence**: - Haptik (2023): 75-80% accuracy with keyword-based ambiguity detection in chatbots - ChatGPT acknowledgment: "guesses" instead of asking clarifying questions (OpenAI, 2024) - TrustNLP (2025): 85-90% accuracy with dependency parsing + keywords

**Phase 1 Validation Plan**: 1. Create labeled dataset (100 prompts: 50 vague, 50 specific) 2. Run heuristic detector, measure accuracy 3. Target: >75% accuracy, >0.75 F1 score

**Upgrade Trigger**: If accuracy <75% after 100 production runs → implement Phase 2 (enhanced heuristics).

---

## 6. Domain-Specific Patterns

**Coding Task Vagueness Indicators** (from research):

**Vague Verbs** (lack specificity): - "implement", "add", "make", "create", "do", "build", "fix", "update"

**Missing Context** (common in coding tasks): - "OAuth" without version → "OAuth 2.0" - "database" without type → "PostgreSQL" - "authentication" without method → "JWT" - "API" without version/endpoint → "REST API v2"

**Ambiguous Quantifiers**: - "some", "a few", "better", "faster", "good"

**Heuristic Pattern Library**:

```rust
lazy_static! {
    static ref MISSING_OAUTH_VERSION: Regex =
        Regex::new(r"(?i)\bOAuth\b(?!\s*(2|1\.0))").unwrap();
    static ref MISSING_DB_TYPE: Regex =
        Regex::new(r"(?i)\bdatabase\b(?!\s*(SQL|NoSQL|PostgreSQL|MySQL))").unwrap();
    static ref MISSING_AUTH_TYPE: Regex =
        Regex::new(r"(?i)\bauth\b(?!\s*(JWT|OAuth|SAML))").unwrap();
```

```
    }
```

**Verdict**: Domain-specific patterns improve heuristic accuracy to 75-80% without ML/LLM.

---

# Comparison to Alternatives

## Alternative 1: ML-Based (Fine-Tuned BERT)

**Approach**: Train BERT classifier on labeled dataset (vague/specific).

**Pros**: - ✓ High accuracy (85-90%) - ✓ Context-aware (understands semantics) - ✓ Local inference (no API calls)

**Cons**: - ✗ Training complexity (requires labeled dataset) - ✗ Slow (50-500ms on CPU) - ✗ Large model files (100MB-1GB) - ✗ Requires GPU for real-time use

**Verdict**: ✗ Reject - Over-engineering for Phase 1. Consider for Phase 2 if heuristic accuracy insufficient.

---

## Alternative 2: LLM-Based (Claude Haiku)

**Approach**: Use Claude Haiku for vagueness classification.

**Pros**: - ✓ High accuracy (90-92%) - ✓ Context-aware (domain knowledge) - ✓ Easy to implement (API call) - ✓ Provides reasoning (explainable)

**Cons**: - ✗ Slow (300-1000ms latency) - ✗ Costs $0.000375/prompt ($3.75/year for 10K prompts) - ✗ Non-deterministic (different responses) - ✗ Requires API access (network dependency)

**Verdict**: ⚠ Phase 3 upgrade if heuristic accuracy <80% in production.

---

## Alternative 3: LLM-Based (GPT-4)

**Approach**: Use GPT-4 for vagueness classification.

**Pros**: - ✓ Highest accuracy (90-95%) - ✓ Best context understanding

**Cons**: - ✗ Very slow (1-2 seconds latency) - ✗ Expensive ($0.01/prompt = $100/year for 10K prompts) - ✗ Exceeds budget constraint ($100/year for ALL PPP features)

**Verdict**: ✗ Reject - Too expensive, violates budget constraint.

---

## Alternative 4: Hybrid (Heuristic + LLM Fallback)

**Approach**: Use heuristic first, LLM only if uncertain (score 0.4-0.6).

**Workflow**:

```rust
    pub async fn is_vague_hybrid(prompt: &str) -> Result<bool> {
```

```
    let heuristic_score = vagueness_score(prompt);

    // Confident cases (70%)
    if heuristic_score < 0.3 || heuristic_score > 0.7 {
        return Ok(heuristic_score > 0.5);
    }

    // Uncertain cases (30%): use LLM
    is_vague_llm(prompt).await
}
```

**Cost Reduction**: - Confident cases (70%): Free (heuristic) - Uncertain cases (30%): $0.000375 (LLM) - **Effective cost**: $0.000375 × 0.3 = **$0.0001125/prompt** - **70% cost reduction** vs pure LLM

**Verdict**: ⚠ Phase 3 optimization if LLM costs exceed budget.

---

## Decision Matrix

**Scoring** (0-10 scale, higher better):

| Criterion | Weight | Heuristic | ML (BERT) | LLM (Haiku) | Hybr |
|-----------|--------|-----------|-----------|-------------|------|
| **Accuracy** | 0.3 | 7 (75-80%) | 8.5 (85-90%) | 9 (90-92%) | 8.5 (8 90%) |
| **Latency** | 0.25 | 10 (<1ms) | 7 (50ms) | 5 (500ms) | 8 (15( avg) |
| **Cost** | 0.2 | 10 ($0) | 10 ($0) | 7 ($3.75/yr) | 9 ($1.1: |
| **Simplicity** | 0.15 | 9 (regex) | 4 (training) | 8 (API) | 6 (log |
| **Maintainability** | 0.1 | 8 (patterns) | 5 (model) | 9 (prompts) | 7 (2 syster |
| **Weighted Score** | - | **8.63** | **7.48** | **7.70** | **7.98** |

**Winner**: **Heuristic (8.63)** for Phase 1

**Hybrid (7.98)** is second-best, viable for Phase 3 if heuristic accuracy insufficient.

---

## Consequences

### Positive

1. ✅ **Zero cost**: No API calls, no model training
2. ✅ **Fast**: <1ms latency (negligible overhead)
3. ✅ **Simple**: Transparent, testable, easy to debug
4. ✅ **Maintainable**: Easy to add new patterns
5. ✅ **Budget-compliant**: Meets <$100/year constraint
6. ✅ **Sufficient accuracy**: 75-80% validated in literature

## Negative

1. ⚠ **Lower accuracy**: 75-80% vs 90-95% with LLM
   - **Mitigation**: Upgrade to Phase 2 (enhanced heuristics) or Phase 3 (LLM) if accuracy insufficient
   - **Impact**: 20-25% false negatives (vague prompts not detected)
2. ⚠ **Pattern maintenance**: Requires manual curation of keywords/patterns
   - **Mitigation**: Quarterly review of missed cases, add patterns as needed
   - **Impact**: Low (10-15 minutes/quarter)
3. ⚠ **Context-insensitive**: Cannot understand complex semantics
   - **Example**: "Implement the same approach as last time" (vague, but heuristic may miss)
   - **Mitigation**: Multi-turn context tracking (future work)

## Neutral

1. ⼮ **Phased upgrade path**: Clear migration to LLM if needed
   - Phase 1 → Phase 2 (dependency parsing)
   - Phase 2 → Phase 3 (hybrid LLM)
   - No breaking changes

---

# Implementation Plan

## Phase 1: Heuristic Foundation (Immediate)

**Timeline**: 2-3 days

**Deliverables**: 1. `vagueness_detector.rs` (keyword + pattern matching) 2. Test suite (20 test cases) 3. Validation dataset (100 prompts) 4. Integration with trajectory logging

**Acceptance Criteria**: - [ ] Accuracy >75% on validation dataset - [ ] Latency <5ms (p95) - [ ] 80%+ test coverage

---

## Phase 2: Enhanced Heuristics (If Needed)

**Timeline**: 1-2 weeks **Trigger**: Accuracy <80% after 100 production runs

**Enhancements**: 1. Dependency parsing (nlprule crate) 2. Domain-specific patterns (technology vocabulary) 3. Multi-part prompt handling

**Target Accuracy**: 80-85%

---

## Phase 3: LLM Upgrade (Optional)

**Timeline**: 1 week **Trigger**: Accuracy <85% OR user feedback indicates poor detection

**Approach**: Hybrid (heuristic + LLM fallback)

**Target**: - Accuracy: 90-95% - Cost: $1.13-$3.75/year (70% reduction vs pure LLM)

# Validation Strategy

## Test Dataset

**Creation**: 1. Collect 100 real prompts from codex-tui logs 2. Manual labeling (2 engineers, resolve disagreements) 3. 50 vague, 50 specific

**Vague Examples**: - "Implement OAuth" (missing version) - "Add authentication" (missing type) - "Fix the bug" (missing context)

**Specific Examples**: - "Implement OAuth2 with Google provider using PKCE" - "Add JWT authentication with HS256 signing" - "Fix null pointer in user_service.rs line 42"

## Metrics

**Accuracy**: (TP + TN) / Total **Precision**: TP / (TP + FP) **Recall**: TP / (TP + FN) **F1 Score**: 2 * (Precision * Recall) / (Precision + Recall)

**Targets**: - Phase 1: Accuracy >75%, F1 >0.75 - Phase 2: Accuracy >80%, F1 >0.80 - Phase 3: Accuracy >90%, F1 >0.90

# References

1. Sun, W., et al. (2025). "Training Proactive and Personalized LLM Agents." arXiv:2511.02208
2. Haptik (2023). "90% User Response Rate to Clarification Questions in Chatbots"
3. OpenAI (2024). "ChatGPT acknowledges guessing instead of asking clarifying questions"
4. TrustNLP (2025). "Ambiguity Detection and Uncertainty Calibration for Question Answering"
5. findings.md (SPEC-PPP-001): Literature review of vagueness detection methods

# Notes

**Why not start with LLM?** - Budget constraint ($100/year for ALL PPP features) - Latency overhead (500-2000ms unacceptable for <1s operation) - Heuristic accuracy (75-80%) is **good enough** for Phase 1 validation

**Why not ML (BERT)?** - Training complexity (requires labeled dataset + GPU) - No accuracy advantage over hybrid LLM approach in Phase 3 - Heuristic → LLM is simpler upgrade path than Heuristic → ML

**Success Criteria**: If Phase 1 achieves >75% accuracy in production, this ADR is validated. If <75%, upgrade to Phase 2 (enhanced heuristics).

**Decision**: **Accepted** (2025-11-16) **Implemented**: Phase 1 (heuristic) in vagueness_detector_poc.rs **Next Review**: After 100 production consensus runs