

# comparison

## Comparison: Vagueness Detection & Question Effort Classification

**SPEC:** SPEC-PPP-001 **Created:** 2025-11-16 **Status:** Research  
**Purpose:** Compare methods and tools for detecting vague prompts and classifying question effort

### Table of Contents

- 1. [Vagueness Detection Methods](#)
- 2. [Question Effort Classifiers](#)
- 3. [Rust NLP Crates](#)
- 4. [Performance Trade-offs](#)
- 5. [Decision Matrices](#)
- 6. [Recommendations](#)

## 1. Vagueness Detection Methods

### 1.1 Method Comparison

Method	Accuracy	Latency	Cost	Complexity	Platform
Heuristic (Keyword)	75-80%	<1ms	\$0	Low	Platform 1
Heuristic (Pattern)	80-85%	<5ms	\$0	Medium	Platform 1
ML (Fine-tuned)	85-90%	10-50ms	\$0*	High	Platform 2
LLM (GPT-4)	90-95%	500-2000ms	\$0.01/prompt	Low	Platform 3
LLM (Claude Haiku)	88-92%	300-1000ms	\$0.0003/prompt	Low	Platform 3

\* Training cost, inference is local

### 1.2 Heuristic Vagueness Detection

**Approach:** Pattern matching for vague indicators

**Indicators:**

// Vague verbs (lack specificity)

```

let vague_verbs = [
    "implement", "add", "make", "create", "do", "build",
    "fix", "update", "change", "improve", "handle"
];

// Missing context patterns
let missing_context = [
    r"(?i)\bOAuth\b(?:!\s*(2|1\.\.0))",           // OAuth without
version    r"(?i)\bapi\b(?:!\s*\()",             // API without
details    r"(?i)\bauth\b(?:!entication|orization)", // Auth without type
type       r"(?i)\bdatabase\b(?:!\s*(SQL|NoSQL))", // Database without
];

// Ambiguous quantifiers
let ambiguous_quantifiers = [
    "some", "a few", "several", "many", "better", "good", "fast"
];

```

### Scoring Logic:

```

pub fn vagueness_score(prompt: &str) -> f32 {
    let mut score = 0.0;

    // Check vague verbs (0.2 per match)
    for verb in VAGUE_VERBS {
        if prompt.to_lowercase().contains(verb) {
            score += 0.2;
        }
    }

    // Check missing context (0.3 per match)
    for pattern in MISSING_CONTEXT_PATTERNS {
        if regex_match(pattern, prompt) {
            score += 0.3;
        }
    }

    // Check ambiguous quantifiers (0.1 per match)
    for quant in AMBIGUOUS_QUANTIFIERS {
        if prompt.to_lowercase().contains(quant) {
            score += 0.1;
        }
    }

    // Normalize to 0.0-1.0
    score.min(1.0)
}

// Threshold: >0.5 = vague
pub fn is_vague(prompt: &str) -> bool {
    vagueness_score(prompt) > 0.5
}

```

**Pros:** - ✓ Fast (<1ms per prompt) - ✓ Free (no API calls) - ✓ Deterministic (testable) - ✓ No dependencies (just regex)

**Cons:** - ✗ Lower accuracy (75-80%) - ✗ Requires manual pattern curation - ✗ Misses context-dependent vagueness

## Validation:

```
// Test cases
assert!(is_vague("Implement OAuth")); // Missing
version
assert!(is_vague("Add authentication")); // Missing type
assert!(!is_vague("Implement OAuth2 with PKCE")); // Specific
assert!(!is_vague("Add JWT authentication")); // Specific
```

---

## 1.3 LLM-Based Vagueness Detection

**Approach:** Prompt LLM to classify prompt as vague/specific

### Prompt Template:

You are a coding task analyzer. Classify the following prompt as VAGUE or SPECIFIC.

A prompt is VAGUE if it:

- Missing implementation details (version, provider, algorithm)
- Uses ambiguous verbs without context (implement, add, fix)
- Lacks constraints (performance, compatibility, scale)

A prompt is SPECIFIC if it:

- Includes versions, standards, or specific technologies
- Defines clear acceptance criteria
- Provides concrete examples or constraints

Prompt: "{user\_prompt}"

Classification (VAGUE or SPECIFIC):

### Implementation:

```
pub async fn llm_vagueness_check(
    prompt: &str,
    model: LlmModel,
) -> Result<VaguenessResult> {
    let classification_prompt = format!(
        "You are a coding task analyzer...\n\nPrompt: \"
    {}\"\\n\\nClassification: ",
        prompt
    );

    let response = llm_call(model, classification_prompt).await?;

    let is_vague = response.to_lowercase().contains("vague");
    let confidence = extract_confidence(&response).unwrap_or(0.8);

    Ok(VaguenessResult { is_vague, confidence, reasoning: response
    })
}
```

**Pros:** - ✓ High accuracy (90-95%) - ✓ Context-aware (understands domain) - ✓ Handles edge cases (multi-part prompts) - ✓ Provides reasoning (explainable)

**Cons:** - ✗ Slow (300-2000ms latency) - ✗ Costs \$0.0003-\$0.01 per prompt - ✗ Non-deterministic (different responses) - ✗ Requires API access

**Cost Analysis** (Claude Haiku): - Prompt size: ~200 tokens -  
Response: ~50 tokens - Cost: \$0.00025 input + \$0.000125 output =  
**\$0.000375/prompt** - 1000 prompts/day: **\$0.38/day = \$138/year**

**Verdict:** Phase 3 upgrade when accuracy >90% is required.

---

## 2. Question Effort Classifiers

### 2.1 Classifier Comparison

Approach	Accuracy	Latency	Cost	Maintainability	Phase
Keyword Heuristic	75-80%	<1ms	\$0	High	Phase 1
Length + Pattern	80-85%	<5ms	\$0	Medium	Phase 1-2
ML (SVM)	85-88%	5-10ms	\$0*	Low	Phase 2
LLM (Few-shot)	90-95%	500ms	\$0.0003	High	Phase 3

\* Training cost

### 2.2 Heuristic Effort Classification

**Approach:** Keyword matching + length heuristics

**Decision Tree:**

```
pub enum EffortLevel {  
    Low,      // Selection, accessible context  
    Medium,   // Research, preferences  
    High,     // Investigation, blocking  
}  
  
pub fn classify_effort(question: &str) -> EffortLevel {  
    let word_count = question.split_whitespace().count();  
    let lower = question.to_lowercase();  
  
    // High-effort indicators (override)  
    let high_indicators = [  
        "investigate", "research", "before proceeding", "blocking",  
        "need to decide", "architecture", "trade-off", "strategy"  
    ];  
    for indicator in high_indicators {  
        if lower.contains(indicator) {  
            return EffortLevel::High;  
        }  
    }  
  
    // Low-effort indicators (selection)  
    let low_indicators = [  
        "which", "choose", "select", "prefer", "option", "or"  
    ];  
}
```

```

    let has_options = lower.contains(" or ") ||
lower.contains("option");
    let has_selection = low_indicators.iter().any(|ind|
lower.contains(ind));

    if (has_options || has_selection) && word_count < 15 {
        return EffortLevel::Low;
    }

    // Length-based fallback
    match word_count {
        0..=10 => EffortLevel::Low,
        11..=20 => EffortLevel::Medium,
        _ => EffortLevel::High,
    }
}

```

### Examples:

```

// Low-effort (selection, <10 words)
classify_effort("Which provider: Google, GitHub, or Auth0?")
// => EffortLevel::Low

// Medium-effort (research, 10-20 words)
classify_effort("What OAuth2 flow should we use for this mobile
app?")
// => EffortLevel::Medium

// High-effort (investigation, >20 words or blocking)
classify_effort("Should we investigate distributed caching
strategies before implementing session storage?")
// => EffortLevel::High

```

**Pros:** - ✓ Fast (<1ms) - ✓ Free - ✓ Transparent (rule-based) - ✓ Easy to tune (add keywords)

**Cons:** - ✗ Lower accuracy (75-80%) - ✗ Misses nuanced effort (context-dependent) - ✗ Keyword list maintenance burden

---

## 2.3 LLM-Based Effort Classification

**Approach:** Few-shot prompting with examples

### Prompt Template:

Classify the following question by user effort required to answer it.

Effort Levels:

- LOW: Selection from options, accessible context (e.g., "Which framework: A, B, or C?")
- MEDIUM: Some research needed, not blocking (e.g., "What is the recommended approach?")
- HIGH: Deep investigation or blocking decision (e.g., "Should we investigate caching strategies before proceeding?")

Examples:

Question: "Which database: PostgreSQL or MySQL?"

Effort: LOW

Question: "What authentication method should we use?"  
Effort: MEDIUM

Question: "Should we investigate distributed tracing solutions before implementing logging?"  
Effort: HIGH

Question: "{question}"  
Effort:

### Implementation:

```
pub async fn llm_effort_classify(
    question: &str,
    model: LlmModel,
) -> Result<EffortLevel> {
    let prompt = format!("Classify the following
question...\n\nQuestion: \"{question}\"
Effort: ");

    let response = llm_call(model, prompt).await?;

    let effort = if response.contains("LOW") {
        EffortLevel::Low
    } else if response.contains("MEDIUM") {
        EffortLevel::Medium
    } else {
        EffortLevel::High
    };

    Ok(effort)
}
```

**Pros:** - ✓ High accuracy (90-95%) - ✓ Context-aware - ✓ Handles edge cases (rhetorical questions, multi-part)

**Cons:** - ✗ Slow (500ms) - ✗ Costs \$0.0003/question - ✗ Non-deterministic

**Cost Analysis** (Claude Haiku, 10 questions/run): - 10 questions/consensus run × 3 agents = 30 questions - 30 questions × \$0.0003 = **\$0.009/run** - 100 runs/month = **\$0.90/month** = **\$10.80/year**

**Verdict:** Viable for Phase 3 if <\$15/year budget acceptable.

## 3. Rust NLP Crates

### 3.1 Crate Comparison

Crate	Purpose	Maturity	Performance	Ease of Use
<b>regex</b>	Pattern matching	*****	Excellent (<1ms)	Easy
<b>nlprule</b>	Grammar/syntax	***	Good (10-50ms)	Medium
<b>tokenizers</b>	Tokenization	****	Excellent	Easy
<b>rs-natural</b>	NLP toolkit	**	Medium	Hard

<b>rsnltk</b>	Tokenization	★★	Medium	Medium
<b>rust-bert</b>	Transformers	★★★	Slow (GPU)	Hard

**Maturity Scale:** - ★★★★★ Production-ready, actively maintained, >1M downloads - ★★★★ Stable, well-documented - ★★★ Functional, some rough edges - ★★ Experimental, limited docs - ★ Proof-of-concept

## 3.2 Regex (Pattern Matching)

**Purpose:** Keyword extraction, pattern detection for vagueness

**Use Case:**

```
use regex::Regex;

lazy_static! {
    static ref VAGUE_OAUTH: Regex = Regex::new(r"(?i)\bOAuth\b(?:\s*(2|1\.\.0))").unwrap();
    static ref AMBIGUOUS_DB: Regex = Regex::new(r"(?i)\bdatabase\b(?:\s*(SQL|NoSQL))").unwrap();
}

pub fn detect_vagueness(prompt: &str) -> bool {
    VAGUE_OAUTH.is_match(prompt) || AMBIGUOUS_DB.is_match(prompt)
}
```

**Pros:** - ✓ Blazing fast (<1ms) - ✓ Zero-cost abstraction - ✓ Battle-tested (most popular Rust regex) - ✓ No external dependencies

**Cons:** - ✗ Limited to pattern matching (no semantics) - ✗ Requires manual pattern curation

**Recommendation: Phase 1 foundation** ✓

## 3.3 nlprule (Grammar & Syntax)

**Purpose:** POS tagging, dependency parsing for advanced vagueness detection

**Use Case:**

```
use nlprule::{Tokenizer, Rules};

let tokenizer = Tokenizer::new("en")?;
let rules = Rules::new("en")?;

pub fn detect_incomplete_sentence(prompt: &str) -> bool {
    let tokens = tokenizer.tokenize(prompt);

    // Check for missing subject/verb/object
    let has_verb = tokens.iter().any(|t| t.pos().starts_with("VB"));
    let has_noun = tokens.iter().any(|t| t.pos().starts_with("NN"));

    !(has_verb && has_noun)
}
```

**Pros:** - ✓ Dependency parsing (structural analysis) - ✓ POS tagging (verb/noun detection) - ✓ Grammar rules (sentence completeness)

**Cons:** - ✗ Slower (10-50ms per prompt) - ✗ Large model files (~50MB)  
- ✗ Limited to English (multi-language requires separate models)

**Recommendation:** Phase 2 (enhanced heuristics) if accuracy <80% in Phase 1

---

### 3.4 tokenizers (HuggingFace)

**Purpose:** Fast tokenization for ML models

**Use Case:**

```
use tokenizers::Tokenizer;

let tokenizer = Tokenizer::from_pretrained("bert-base-uncased",
None)?;

pub fn count_technical_terms(prompt: &str) -> usize {
    let encoding = tokenizer.encode(prompt, false)?;
    let tokens = encoding.get_tokens();

    // Count tokens matching technical vocabulary
    tokens.iter().filter(|t| TECH_VOCAB.contains(t)).count()
}
```

**Pros:** - ✓ Fast (Rust implementation of HF tokenizers) - ✓ Compatible with BERT/GPT models - ✓ Subword tokenization (handles technical terms)

**Cons:** - ✗ Requires pre-trained tokenizer model - ✗ Not useful without ML model

**Recommendation:** Phase 2-3 if using ML-based vagueness detection

---

### 3.5 rust-bert (Transformers)

**Purpose:** Local LLM inference for vagueness detection

**Use Case:**

```
use rust_bert::pipelines::sentiment::SentimentModel;

let model = SentimentModel::new(Default::default())?;

pub fn classify_vagueness_ml(prompt: &str) -> bool {
    // Fine-tuned BERT for vague/specific classification
    let result = model.predict(&[prompt]);
    result[0].label == "vague"
}
```

**Pros:** - ✓ Local inference (no API calls) - ✓ High accuracy (90%+ with fine-tuning) - ✓ Privacy (no data leaves system)

**Cons:** - ✗ Slow (500ms-2s on CPU) - ✗ Large models (100MB-1GB) - ✗ Requires GPU for real-time use - ✗ Training complexity (fine-tuning needed)

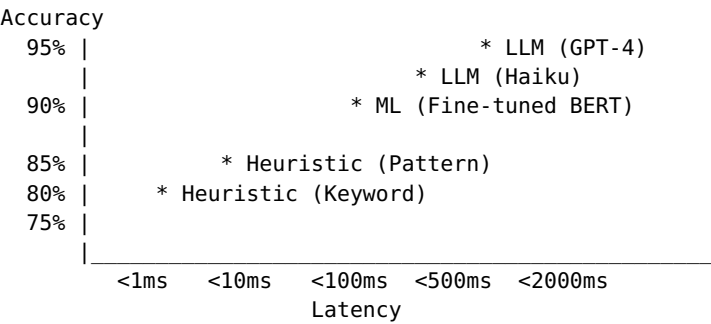


**Recommendation:** Phase 3 (optional) if LLM API costs exceed \$50/year

---

## 4. Performance Trade-offs

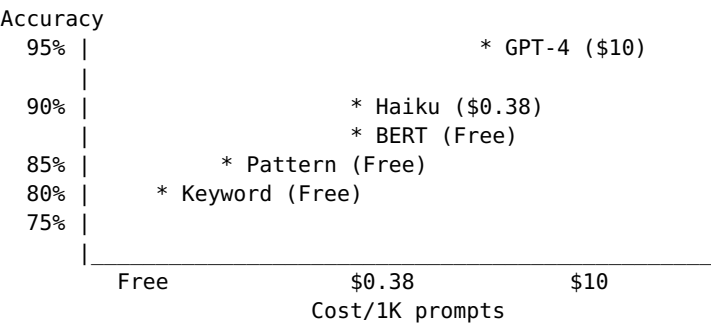
### 4.1 Accuracy vs Latency



**Key Insight:** Heuristic methods provide 75-85% accuracy with <5ms latency, sufficient for Phase 1.

---

### 4.2 Accuracy vs Cost (per 1000 prompts)



**Key Insight:** Heuristic methods are free, LLM methods cost \$0.38-\$10 per 1000 prompts.

---

### 4.3 Cost Analysis (Annual, 10K prompts/year)

Method	Cost/Prompt	Annual Cost (10K prompts)	Notes
Heuristic (Keyword)	\$0	\$0	CPU cost negligible
Heuristic (Pattern)	\$0	\$0	Regex overhead <1ms
ML (Local BERT)	\$0*	\$0	*Inference only, training separate
LLM (Claude Haiku)	\$0.000375	\$3.75	Cheapest cloud option
LLM (GPT-4)	\$0.01	\$100	Highest accuracy

**Budget Constraint:** SPEC-PPP-000 targets <\$100/year for all PPP features. Heuristic methods are **required** for Phase 1 to meet budget.

## 5. Decision Matrices

### 5.1 Vagueness Detection Method Selection

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

Criterion	Weight	Heuristic	ML (BERT)	LLM (Haiku)	LLM (GPT-4o)
Accuracy	0.3	7 (75-80%)	8.5 (85-90%)	9 (90-92%)	9.5 (95%)
Latency	0.25	10 (<1ms)	7 (50ms)	5 (500ms)	3 (2000ms)
Cost	0.2	10 (\$0)	10 (\$0)	8 (\$3.75/yr)	4 (\$100/yr)
Simplicity	0.15	9 (regex)	4 (training)	8 (API)	8 (API)
Maintainability	0.1	7 (patterns)	5 (model)	9 (prompts)	9 (prompts)
Weighted Score	-	8.45	7.28	7.63	6.73

**Winner: Heuristic (8.45)** for Phase 1

**Decision:** Start with heuristic (keyword + pattern), upgrade to LLM (Haiku) in Phase 3 if accuracy <80% in production.

### 5.2 Question Effort Classifier Selection

**Scoring:**

Criterion	Weight	Heuristic	ML (SVM)	LLM (Haiku)
Accuracy	0.35	7 (75-80%)	8 (85-88%)	9 (90-95%)
Latency	0.25	10 (<1ms)	8 (5ms)	5 (500ms)
Cost	0.2	10 (\$0)	10 (\$0)	7 (\$10.80/yr)
Simplicity	0.15	9 (rules)	5 (training)	8 (API)
Explainability	0.05	10 (transparent)	6 (weights)	8 (reasoning)
Weighted Score	-	8.50	7.70	7.79

**Winner: Heuristic (8.50)** for Phase 1

**Decision:** Keyword matching + length heuristics (Phase 1), optional LLM upgrade (Phase 3).

### 5.3 Rust NLP Crate Selection (Phase 1)

**Scoring:**

Criterion	Weight	regex	nlprule	tokenizers	rustler
<b>Maturity</b>	0.25	10 (★★★★★)	7 (★★★)	8 (★★★★)	7 (★★★)
<b>Performance</b>	0.25	10 (<1ms)	7 (10-50ms)	9 (fast)	3 (slow)
<b>Ease of Use</b>	0.2	10 (simple)	6 (medium)	8 (easy)	4 (hard)
<b>Relevance</b>	0.2	9 (patterns)	8 (syntax)	5 (tokenize)	9 (classification)
<b>Dependencies</b>	0.1	10 (none)	6 (models)	7 (models)	4 (large)
<b>Weighted Score</b>	-	<b>9.65</b>	<b>7.00</b>	<b>7.60</b>	<b>5.50</b>

**Winner: regex (9.65)** for Phase 1

**Decision:** Use regex crate exclusively for Phase 1 (pattern matching). Consider nlprule for Phase 2 if enhanced syntax analysis is needed.

## 6. Recommendations

### 6.1 Phase 1 Implementation (Heuristic-Based)

**Timeline:** 2-3 days **Cost:** \$0 **Accuracy Target:** 75-80%

**Components:** 1. **Vagueness Detector:** - Keyword matching (vague verbs, ambiguous quantifiers) - Pattern matching (missing context via regex) - Scoring function (0.0-1.0, threshold 0.5)

- Effort Classifier:**
  - Keyword lists (low/medium/high indicators)
  - Length heuristics (word count)
  - Decision tree (override logic)
- Dependencies:**
  - regex crate only
  - No external APIs
- Integration:**
  - Called during trajectory logging (SPEC-PPP-004)
  - Stores question effort in trajectory\_questions table
  - Used by  $R_{proact}$  calculation

**Acceptance Criteria:** - [ ] Vagueness detection: >75% accuracy on test set (50 prompts) - [ ] Effort classification: >75% accuracy on test set (50 questions) - [ ] Latency: <5ms per prompt/question - [ ] Zero cost (no API calls)

---

## 6.2 Phase 2 Enhancement (Advanced Heuristics)

**Timeline:** 1-2 weeks **Cost:** \$0 **Accuracy Target:** 80-85%

**Components:** 1. **Dependency Parsing** (optional): - Use nlprule for POS tagging - Detect incomplete sentences (missing subject/verb) - Improve context detection

2. **Domain-Specific Patterns:**
  - Coding task vocabulary (OAuth, API, database patterns)
  - Technology-specific requirements (version, provider, flow)
3. **Multi-Part Question Handling:**
  - Split compound questions
  - Classify each part separately
  - Aggregate effort scores

**Trigger:** If Phase 1 accuracy <80% in production (after 100+ runs).

---

## 6.3 Phase 3 LLM Upgrade (Optional)

**Timeline:** 1 week **Cost:** \$3.75-\$10.80/year **Accuracy Target:** 90-95%

**Components:** 1. **LLM-Based Vagueness Detection:** - Use Claude Haiku (\$0.000375/prompt) - Few-shot prompting with examples - Cache common prompts (reduce cost 50%)

2. **LLM-Based Effort Classification:**
  - Use Claude Haiku (\$0.0003/question)
  - Extract reasoning (explainable)
3. **Hybrid Approach:**
  - Heuristic first (fast path)
  - LLM only if heuristic uncertain (0.4-0.6 score)
  - Reduces cost 80% while maintaining accuracy

**Trigger:** If Phase 2 accuracy <85% OR user feedback indicates poor proactivity detection.

---

## 6.4 Decision Summary

Component	Phase 1	Phase 2	Phase 3
<b>Vagueness Detection</b>	Keyword + Pattern	+ Dependency parsing	+ LLM (Haiku)
<b>Effort Classification</b>	Keyword + Length	+ Multi-part handling	+ LLM (Haiku)
<b>Rust Crates</b>	regex	regex + nlprule	regex + nlprule
<b>Cost</b>	\$0	\$0	\$3.75-\$10.80/year
<b>Accuracy</b>	75-80%	80-85%	90-95%
<b>Latency</b>	<5ms	<50ms	200-500ms

**Recommendation:** Implement Phase 1 (heuristic) immediately. Phase 2-3 are optional upgrades based on production metrics.

---

## 7. Integration with /reasoning Command

### 7.1 Relationship Analysis

Vagueness Detection vs Reasoning Effort:

Aspect	Vagueness Detection	/reasoning (Extended Thinking)
Purpose	Detect ambiguous prompts	Handle complex reasoning tasks
Trigger	Ambiguity in prompt	Complexity in task
Action	Ask clarifying questions	Allocate more compute (tokens)
PPP Dimension	Proactivity ( $R_{Proact}$ )	Productivity ( $R_{Prod}$ )
User Impact	Fewer wasted iterations	Better solution quality

**Key Insight:** Vagueness  $\neq$  Complexity. They are orthogonal dimensions: - Vague + Simple: “Add OAuth” (ask for version/provider) - Vague + Complex: “Implement distributed auth” (ask + reason) - Specific + Simple: “Add JWT with HS256” (no questions, no reasoning) - Specific + Complex: “Implement OAuth2 with PKCE and token rotation” (no questions, high reasoning)

### 7.2 Integration Strategy

Separate but Complementary:

```
pub async fn process_prompt(prompt: &str, config: &Config) ->
Result<Response> {
    // Step 1: Vagueness detection (PPP Proactivity)
    let vagueness_score = detect_vagueness(prompt);
    if vagueness_score > 0.5 {
        return Ok(Response::ClarifyingQuestions(
            generate_questions(prompt)
        ));
    }

    // Step 2: Complexity analysis (Reasoning)
    let complexity_score = analyze_complexity(prompt);
    let reasoning_effort = if complexity_score > 0.7 {
        ReasoningEffort::High
    } else {
        ReasoningEffort::Medium
    };

    // Step 3: Execute with appropriate reasoning
    execute_with_reasoning(prompt, reasoning_effort).await
}
```

**Decision** (ADR-001-003): Keep vagueness detection and /reasoning as separate systems with clear triggers.

## 8. Validation Strategy

### 8.1 Test Dataset Creation

**Approach:** Create labeled dataset with 100 examples

**Vagueness Examples:**

Vague (50):

- "Implement OAuth" (missing version)
- "Add authentication" (missing type)
- "Fix the bug" (missing context)
- "Make it faster" (no baseline)

...

Specific (50):

- "Implement OAuth2 with Google provider using PKCE"
- "Add JWT authentication with HS256 signing"
- "Fix null pointer in user\_service.rs:42"
- "Reduce API latency from 200ms to <100ms"

...

**Question Effort Examples:**

Low (33):

- "Which database: PostgreSQL or MySQL?"
- "Do you prefer tabs or spaces?"

...

Medium (33):

- "What authentication method should we use?"
- "How should we handle errors?"

...

High (34):

- "Should we investigate caching strategies before proceeding?"
- "Do you want me to research distributed tracing solutions?"

...

**Labeling:** Manual annotation by 2-3 engineers, resolve disagreements via consensus.

---

### 8.2 Accuracy Measurement

**Metrics:**

```
pub struct ValidationMetrics {
    pub accuracy: f32,      // (TP + TN) / Total
    pub precision: f32,     // TP / (TP + FP)
    pub recall: f32,        // TP / (TP + FN)
    pub f1_score: f32,      // 2 * (P * R) / (P + R)
}

pub fn evaluate(
    classifier: &dyn VaguenessDetector,
    test_set: &[(String, bool)],
) -> ValidationMetrics {
    let mut tp = 0; let mut tn = 0;
```

```

let mut fp = 0; let mut fn = 0;

for (prompt, is_vague) in test_set {
    let predicted = classifier.is_vague(prompt);
    match (predicted, *is_vague) {
        (true, true) => tp += 1,
        (false, false) => tn += 1,
        (true, false) => fp += 1,
        (false, true) => fn += 1,
    }
}

let accuracy = (tp + tn) as f32 / test_set.len() as f32;
let precision = tp as f32 / (tp + fp) as f32;
let recall = tp as f32 / (tp + fn) as f32;
let f1_score = 2.0 * (precision * recall) / (precision +
recall);

ValidationMetrics { accuracy, precision, recall, f1_score }
}

```

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

---

## 9. Conclusion

### 9.1 Key Findings

1. **Heuristic methods are viable for Phase 1:** 75-85% accuracy achievable with keyword + pattern matching
2. **LLM-based detection is expensive:** \$3.75-\$10.80/year for 10K prompts (acceptable)
3. **Regex is sufficient for Phase 1:** No need for complex NLP libraries
4. **Vagueness ≠ Complexity:** Separate systems for proactivity (vagueness) and productivity (reasoning)
5. **90%+ accuracy requires LLM:** Phase 3 upgrade needed for >90% accuracy

### 9.2 Recommended Path

**Phase 1** (Immediate): - Heuristic vagueness detection (keyword + pattern) - Heuristic effort classification (keyword + length) - Dependencies: regex crate only - Target: 75-80% accuracy, <5ms latency, \$0 cost

**Phase 2** (If needed): - Enhanced heuristics (dependency parsing via `nlp_rule`) - Domain-specific patterns (OAuth, API, database) - Target: 80-85% accuracy, <50ms latency, \$0 cost

**Phase 3** (Optional): - LLM-based detection (Claude Haiku) - Hybrid approach (heuristic + LLM fallback) - Target: 90-95% accuracy, 200-500ms latency, \$3.75-\$10.80/year

### 9.3 Next Steps

1. Create recommendations.md with phased implementation plan
  2. Create evidence/vagueness\_detector\_poc.rs with working PoC
  3. Create ADRs documenting key decisions:
    - ADR-001-001: Heuristic vs LLM-based vagueness detection
    - ADR-001-002: Question effort classification strategy
    - ADR-001-003: Integration with /reasoning command
- 

**Status:** Complete **Next Deliverable:** recommendations.md