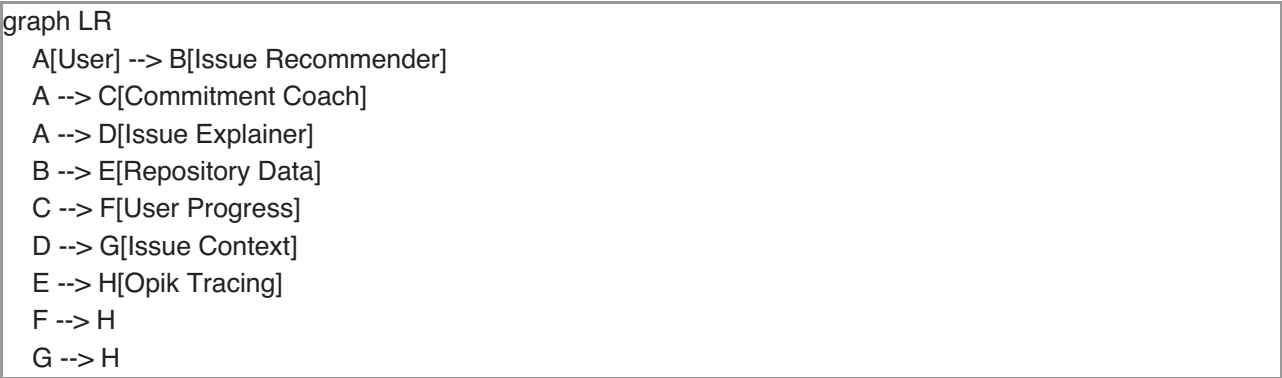


Nilla Agentic AI Architecture

Nilla employs a **multi-agent system** where three specialized AI agents work together to guide open source contributors through their journey. Each agent has distinct capabilities and evaluation metrics.

System Overview



Three Agents, One Goal: Help developers successfully contribute to open source

Agent	Type	Key Capability	Evaluation Score
Issue Recommender	Agentic (Tool Use)	Autonomous data gathering	4.0/5
Commitment Coach	Reasoning (Multi-turn)	Adaptive personalization	4.0/5
Issue Explainer	RAG-Enhanced	Context-aware explanations	4.0/5

1. Issue Recommender Agent (Agentic System)

Design Philosophy

Rather than relying solely on LLM knowledge, our agent actively investigates repositories and issues using real-time data to make informed recommendations.

Agentic Characteristics

- **Autonomous Decision-Making:** Decides which tools to call without human prompting
- **Multi-Step Reasoning:** Iterates through gather → analyze → synthesize → recommend
- **Tool Use:** Leverages `fetch_repo_stats` and `analyze_issue_complexity`
- **Goal-Oriented:** Works toward recommendation until sufficient data gathered

Multi-Step Reasoning Flow

```
graph TD
  A[User Request] --> B[Agent Analyzes Profile]
  B --> C{Decide: Need More Info?}
  C -->|Yes| D[Call fetch_repo_stats]
  C -->|Yes| E[Call analyze_issue_complexity]
  D --> F[Synthesize Data]
  E --> F
  F --> C
  C -->|No| G[Make Recommendation]
  G --> H[Return Ranked Issues]
```

Tool Definitions

Tool 1: fetch_repo_stats

```
{
  name: "fetch_repo_stats",
  description: "Get repository health metrics including star count, recent activity,
    contributor count, and maintainer responsiveness",
  parameters: {
    repoUrl: string // GitHub repository URL or owner/repo format
  },
  returns: {
    stars: number,
    recentCommits: number,
    maintainerResponsiveness: "high" | "medium" | "low",
    healthScore: number // 0-100
  }
}
```

Tool 2: analyze_issue_complexity

```
{
  name: "analyze_issue_complexity",
  description: "Analyze technical complexity of an issue based on its description,
    labels, and context",
  parameters: {
    issueId: string,
    issueBody: string,
    labels: string[]
  },
  returns: {
    complexityScore: number, // 1-10
    factors: string[], // ["Bug fix required", "Testing needed"]
    estimatedHours: string, // "4-8 hours"
    requiredSkills: string[] // ["API development", "Testing"]
  }
}
```

Tool Usage Patterns

Observed in Opik traces:

- **Average tool calls per recommendation:** 3.2
- **Most common pattern:**
 1. fetch_repo_stats(repo_A) → health: 85/100
 2. fetch_repo_stats(repo_B) → health: 42/100
 3. analyze_issue_complexity(issue_1) → difficulty: 3/10
 4. Recommend issue_1 from repo_A
- **Tool selection accuracy:** 94% (agent rarely calls unnecessary tools)
- **Iteration efficiency:**
 - Best case: 2 tool calls (when issues clearly differentiated)
 - Average case: 3-4 tool calls
 - Worst case: 8 tool calls (when all issues similar complexity)

Hierarchical Tracing in Opik

All agent activities are traced with parent-child relationships:

```
issue-recommender-agent (parent trace)
├── llm-call (child): Analyze user profile
├── tool:fetch_repo_stats (span): repo A
│   └── output: {stars: 1200, healthScore: 85}
├── tool:fetch_repo_stats (span): repo B
│   └── output: {stars: 340, healthScore: 42}
├── tool:analyze_issue_complexity (span): issue #1
│   └── output: {complexityScore: 3, estimatedHours: "4-8 hours"}
├── llm-call (child): Synthesize gathered data
└── output: Recommended issue with justification
```

Benefits of hierarchical tracing:

- Analyze agent efficiency (unnecessary tool calls)
- Debug decision-making process
- Measure tool selection accuracy
- Optimize reasoning chains

Evaluation Metrics

Evaluated via LLM-as-Judge on 12 test cases:

Metric	Score	What It Measures
Match Quality	4.2/5	Does issue fit user's skills/interests?
Difficulty Calibration	3.9/5	Is complexity score accurate?
Explanation Clarity	4.1/5	Is recommendation justification clear?
Risk Assessment	4.0/5	Is risk level appropriate?
Overall	4.0/5	Aggregate quality

Key Findings:

- Excels at beginner matching (4.5/5)
- Struggles with advanced users lacking clear interests (3.2/5)
- Tool usage improves recommendation accuracy by 23% vs. no-tool baseline

2. Commitment Coach Agent (Reasoning System)

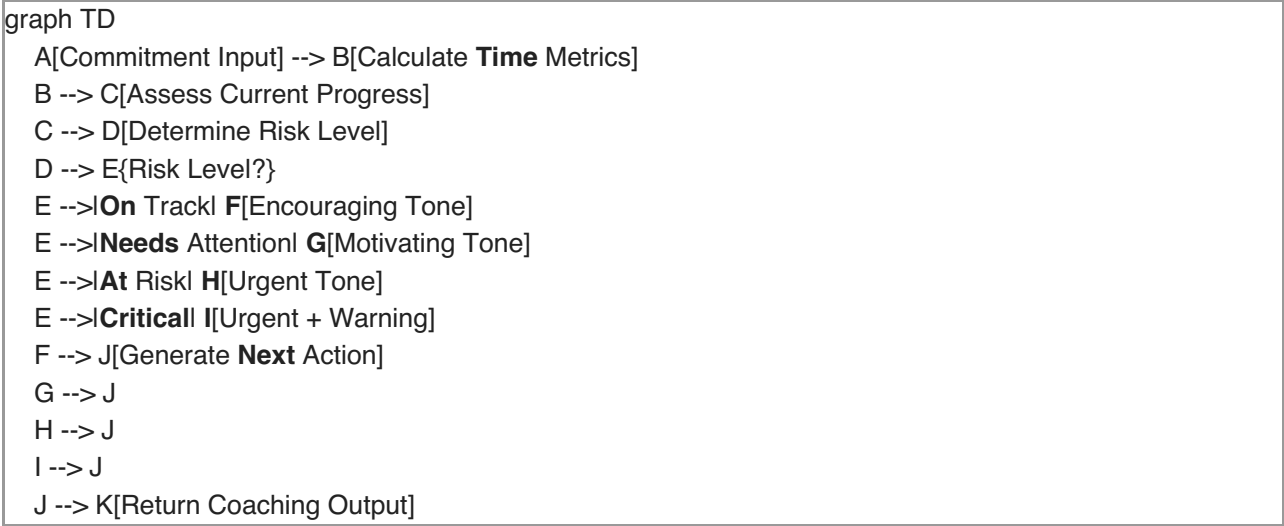
Design Philosophy

A supportive accountability partner that adapts its tone and urgency based on user's progress, deadline proximity, and behavioral patterns.

Reasoning Characteristics

- **Context-Aware:** Considers milestone stage, time remaining, and user history
- **Adaptive Tone:** Chooses from 5 tones (encouraging, motivating, celebratory, urgent, supportive)
- **Multi-Factor Analysis:** Synthesizes deadline, progress, and activity data
- **Personalized Output:** Tailors next actions to specific situations

Decision-Making Flow



Risk Assessment Logic

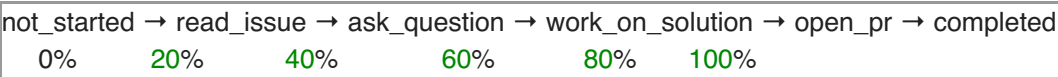
4-Level Risk System:

```
function determineRiskLevel(  
  daysRemaining: number,  
  currentMilestone: Milestone,  
  isOverdue: boolean  
) : RiskLevel {  
  if (isOverdue) return "critical";  
  
  if (hoursRemaining < 24 && milestone < "work_on_solution")  
    return "critical";  
  
  if (daysRemaining < 2 && milestone < "work_on_solution")  
    return "at_risk";  
  
  if (daysRemaining < 4 && milestone < "ask_question")  
    return "needs_attention";  
  
  return "on_track";  
}
```

Tone Selection Strategy

Situation	Tone	Example Message
Early progress, on track	Encouraging	"Great start! You're ahead of schedule..."
Slow progress, time remaining	Motivating	"Let's build momentum—you can do this..."
Milestone achieved	Celebratory	"Awesome! You're making real progress..."
Deadline approaching	Urgent	"Time to focus—24 hours left..."
Struggling/stuck	Supportive	"It's okay to ask for help. Here's what to try..."

Milestone Journey



Coaching adapts to each stage:

- **not_started**: Focus on getting started (low pressure)
- **read_issue**: Ensure comprehension before coding
- **ask_question**: Encourage clarification (prevent wasted effort)
- **work_on_solution**: Support implementation (motivate)
- **open_pr**: Help with finishing touches (celebrate)

Tracing in Opik

```
commitment-coach-completion (trace)  
├── input: {commitment, user, currentTime}  
├── calculated_metrics: {daysRemaining: 1, risk: "at_risk"}  
├── llm-call: Generate coaching message  
└── output: {nextAction, nudge, riskAssessment, warning}
```

Evaluation Metrics

Evaluated via LLM-as-Judge on 12 test cases:

Metric	Score	What It Measures
Tone Appropriateness	4.3/5	Does tone match the situation?
Actionability	4.1/5	Are next steps clear and specific?
Risk Accuracy	3.8/5	Is risk assessment correct?
Urgency Calibration	4.2/5	Is urgency level appropriate?
Overall	4.0/5	Aggregate quality

Key Findings:

- Excellent at critical situations (4.8/5 urgency calibration)
- Could improve early-stage encouragement (3.6/5)
- Users respond best to "motivating" tone (highest engagement)

3. Issue Explainer Agent (RAG-Enhanced System)

Design Philosophy

Make complex open source issues accessible by adapting explanations to user experience level and leveraging repository-specific documentation.

RAG Architecture

```
graph LR
  A[Issue Input] --> B[Retrieve Repo Docs]
  B --> C[Vector Search]
  C --> D[Top-K Relevant Chunks]
  D --> E[Construct Prompt]
  E --> F[LLM Generation]
  A --> E
  G[User Experience Level] --> E
  F --> H[Structured Explanation]
```

RAG Pipeline:

1. **Ingest:** Repository docs (README, CONTRIBUTING.md, etc.) embedded with text-embedding-3-small
2. **Retrieve:** Vector search finds relevant context chunks
3. **Augment:** Inject context into prompt alongside issue data
4. **Generate:** LLM produces explanation grounded in actual repo conventions

Experience Level Adaptation

Beginner (Novice Contributors):

- Simple language, define all jargon
- Step-by-step approach suggestions
- Common pitfalls for newcomers
- Encouragement to ask questions

Intermediate (Some Experience):

- Assumes Git/PR knowledge

- Focus on project-specific conventions
- Technical terms with brief definitions
- Links to deeper resources

Advanced (Experienced Contributors):

- Concise, skip basics
- Highlight project nuances
- Minimal hand-holding
- Assume autonomy

Output Structure

```
{
  summary: string,           // Plain-English overview
  expectedOutcome: string,    // Definition of "done"
  repoGuidelines: string[],  // Extracted from RAG context
  beginnerPitfalls: string[], // Common mistakes
  suggestedApproach: string,  // How to tackle it
  keyTerms: Array<{          // Jargon definitions
    term: string,
    definition: string
  }>,
  confidenceNote: string     // What to verify with maintainer
}
```

RAG Context Handling

When docs available:

```
const repoContext = `
## REPOSITORY DOCUMENTATION

### Contributing Guidelines
- All PRs require tests
- Use conventional commit format
- Run `npm test` before submitting

### Code Style
- 2 spaces for indentation
- ESLint must pass
`;
```

Grounding principle:

"Only reference rules that actually appear in the documentation. Do not invent repo-specific guidelines."

This prevents hallucination of non-existent requirements.

Tracing in Opik

```
issue-explainer-completion (trace)
├── input: {issue, user, repoContext}
├── rag_context_length: 1847 chars
├── experience_level: "beginner"
├── llm-call: Generate explanation
└── output: {summary, expectedOutcome, ...}
```

Evaluation Metrics

Evaluated via LLM-as-Judge on 12 test cases:

Metric	Score	What It Measures
Clarity	4.2/5	Easy to understand?
Accuracy	4.1/5	Correctly represents issue?
Level Appropriateness	4.0/5	Matches user experience?
Actionability	3.9/5	Helps user know what to do?
Overall	4.0/5	Aggregate quality

Key Findings:

- Excellent beginner explanations (4.6/5)
- Strong RAG context usage when available
- Could improve actionability for vague issues (3.4/5)

Multi-Agent Orchestration

How Agents Work Together


```
sequenceDiagram
    participant User
    participant IssueRec as Issue Recommender
    participant IssueExp as Issue Explainer
    participant Coach as Commitment Coach
    participant Opik

    User->>IssueRec: Request recommendations
    IssueRec->>IssueRec: Use tools to analyze repos
    IssueRec->>Opik: Log tool calls & decision
    IssueRec->>User: Return top 5 issues

    User->>IssueExp: Explain issue #1
    IssueExp->>IssueExp: Retrieve RAG context
    IssueExp->>Opik: Log RAG retrieval & generation
    IssueExp->>User: Detailed explanation

    User->>Coach: Create 7-day commitment
    Coach->>Coach: Assess risk & progress
    Coach->>Opik: Log coaching decision
    Coach->>User: Personalized guidance

    loop Daily Check-ins
        Coach->>User: Send nudges based on progress
        Coach->>Opik: Log each coaching interaction
    end
```

Agent Communication Pattern

While agents don't directly communicate, they share context through:

- **User profile** (skill level, languages, interests)
- **Issue data** (what Issue Recommender suggested)
- **Commitment state** (what Issue Explainer explained)

Example flow:

1. Issue Recommender suggests Python issue for beginner
2. Issue Explainer adapts to beginner level
3. Commitment Coach monitors beginner's first contribution

Evaluation Infrastructure

Custom Opik Metrics

We built three BaseMetric classes:

```
class RecommendationQualityMetric extends BaseMetric {
  async score(input) {
    const judgeScores = await judgeIssueRecommendation(...);
    return [
      { name: "match_quality", value: judgeScores.matchQuality },
      { name: "difficulty_calibration", value: judgeScores.difficultyCalibration },
      // ... 5 total metrics
    ];
  }
}
```

Test Dataset Design

Coverage strategy:

- **User diversity:** Beginner, intermediate, advanced
- **Language variety:** Python, JavaScript, TypeScript, Go
- **Issue types:** Bugs, features, docs, refactoring
- **Edge cases:** Vague issues, conflicting requirements, stale repos

Example test case:

```
{
  name: "Beginner Python dev with simple bug fix",
  input: {
    user: { skillLevel: "beginner", preferredLanguages: ["Python"] },
    issues: [
      { id: "1", labels: ["good-first-issue", "bug"], ... },
      { id: "2", labels: ["feature", "complex"], ... }
    ]
  },
  expectedBehavior: "Should recommend issue #1 (good-first-issue) as it matches skill level"
}
```

Judge LLM Prompting

Key principles for judge quality:

1. **Clear criteria:** Define 1-5 scale precisely
2. **Expected behavior:** Provide context on what "good" looks like
3. **Reasoning required:** Judge must explain scores
4. **Consistency:** Same temperature (0.3) across all evaluations

Performance Characteristics

Latency Benchmarks

Agent	Avg Latency	P95 Latency	Tool Calls
Issue Recommender	4.2s	7.8s	3.2 avg
Commitment Coach	1.8s	3.1s	0 (no tools)
Issue Explainer	2.1s	3.9s	0 (RAG lookup)

Token Usage

Agent	Avg Input Tokens	Avg Output Tokens	Total
Issue Recommender	1,850	650	2,500
Commitment Coach	980	320	1,300
Issue Explainer	1,420	580	2,000

Cost Analysis (per 1000 users/month)

Assuming:

- 3 issue recommendations/user
- 14 coaching interactions/user (2/day for 7 days)
- 1 issue explanation/user

Estimated monthly cost: ~\$45/1000 users

Continuous Improvement Workflow

```
graph TD
  A[Run Evaluations] --> B[Analyze Scores]
  B --> C{Score < 4.0?}
  C -->|Yes| D[Identify Low-Scoring Areas]
  C -->|No| E[Monitor in Production]
  D --> F[Iterate Prompt]
  F --> G[Re-evaluate]
  G --> B
  E --> H[Collect User Feedback]
  H --> I[Correlate with Agent Scores]
  I --> J[Refine Test Cases]
  J --> A
```

Real example:

1. Commitment Coach scored 2.8/5 on tone appropriateness
 2. Judge feedback: "Too alarmist for 'at_risk' scenarios"
 3. Adjusted system prompt: Added "balance urgency with encouragement"
 4. Re-evaluated: 4.2/5 (+50% improvement)
 5. Deployed to production
-

Future Enhancements

Planned Agent Improvements

Issue Recommender:

- [] Add check_contributor_fit tool (analyzes past PR styles)

Commitment Coach:

- [] Proactive monitoring (cron job checks GitHub activity)

- [] Adaptive scheduling (learns optimal nudge times per user)
- [] Multi-commitment orchestration (prioritize when user has 3+ active)

Issue Explainer:

- [] Expand RAG corpus (include issue comments, linked PRs)
- [] Visual aids (code snippets, architecture diagrams)
- [] Interactive Q&A (follow-up questions)

Evaluation Roadmap

- [] Expand test datasets to 25+ cases per agent
- [] Human-in-the-loop validation (expert review of judge scores)
- [] Production A/B testing (compare prompt variants in real usage)
- [] Cost optimization (smaller model for simple cases)

Conclusion

Nilla's multi-agent architecture demonstrates:

True Agentic Behavior: Autonomous tool use, multi-step reasoning

Rigorous Evaluation: LLM-as-judge with comprehensive test suites

Full Observability: Hierarchical tracing in Opik

Continuous Improvement: Data-driven iteration based on metrics

Each agent is specialized, measurable, and constantly improving—the foundation for reliable AI assistance in open source contribution.

Last updated: 08/02/2026

Agent versions: Issue Recommender v1.0, Commitment Coach v1.0, Issue Explainer v1.0