# A Declarative Language for Building And Orchestrating LLM-Powered Agent Workflows

Ivan Daunis

idaunis@paypal.com

PayPal

November 10, 2025

## Abstract

Building deployment-ready LLM agents requires complex orchestration of tools, data sources, and control flow logic, yet existing systems tightly couple agent logic to specific programming languages and deployment models. We present a declarative system that separates agent workflow specification from implementation, enabling the same pipeline definition to execute across multiple backend languages (Java, Python, Go) and deployment environments (cloud-native, on-premises).

Our key insight is that most agent workflows consist of common patterns—data serialization, filtering, RAG retrieval, API orchestration—that can be expressed through a unified DSL rather than imperative code. This approach transforms agent development from application programming to configuration, where adding new tools or fine-tuning agent behaviors requires only pipeline specification changes, not code deployment. Our system natively supports A/B testing of agent strategies, allowing multiple pipeline variants to run on the same backend infrastructure with automatic metric collection and comparison.

We evaluate our approach on real-world e-commerce workflows at PayPal, processing millions of daily interactions. Our results demonstrate 60% reduction in development time, and 3x improvement in deployment velocity compared to imperative implementations. The language's declarative approach enables non-engineers to modify agent behaviors safely, while maintaining sub-100ms orchestration overhead. We show that complex workflows involving product search, personalization, and cart management can be expressed in under 50 lines of DSL compared to 500+ lines of imperative code.

## 1. Introduction

Large Language Models (LLMs) have revolutionized how enterprises build intelligent applications, enabling sophisticated capabilities in customer service, e-commerce, and workflow automation. However, deploying LLM-powered agents in production environments presents significant engineering challenges that existing systems fail to adequately address. While engineering platforms like LangChain, AutoGPT, and others provide powerful abstractions for prototyping, they often fall short when scaling to enterprise requirements for reliability, observability, and maintainability.

Current LLM systems suffer from three fundamental limitations. First, they are tightly coupled to the programming language in which they are implemented,
language-specific paradigms rather than business requirements. This coupling makes it difficult to integrate agents into heterogeneous enterprise environments where multiple languages, services, and legacy systems must coexist. Second, imperative programming models make complex agent workflows difficult to understand, test, and modify, leading to fragile implementations that break under production loads. Third, existing platforms lack native support for enterprise concerns such as comprehensive error handling, distributed tracing, security context propagation, and fine-grained metrics collection.

In this work, we present a declarative system for building and orchestrating LLM-powered agent workflows in enterprise environments. Our approach introduces a language-agnostic Domain-Specific Language (DSL) that separates agent logic from implementation details, enabling developers to express complex workflows as composable, reusable pipelines.

By adopting a declarative approach, we make agent behaviors explicit, testable, and maintainable while providing production-ready features essential for enterprise deployment.

## 1.1. Problem Statement

We address the problem of building production-ready LLM agents that must satisfy enterprise requirements for:

- **Reliability:** Graceful handling of LLM failures, timeout management, and fallback strategies

- **Composability:** Ability to build complex workflows from simple, reusable components

- **Observability:** Comprehensive logging, distributed tracing, and performance metrics

- **Integration:** Seamless interaction with existing enterprise services and security models

- **Maintainability:** Clear separation between business logic and infrastructure concerns

Formally, we model agent workflows as directed acyclic graphs (DAGs) where nodes represent computational steps (tool invocations, LLM calls, data transformations) and edges encode control flow dependencies. Our goal is to provide a declarative specification language $\mathcal{L}$ and execution engine $\mathcal{E}$ such that workflows expressed in $\mathcal{L}$ can be efficiently executed by $\mathcal{E}$ while maintaining enterprise-grade reliability and observability.

## 1.2. Contributions

Our main contributions are:

- **Declarative Pipeline DSL:** We design a novel domain-specific language that enables developers to express complex agent workflows using high-level primitives for control flow (`forEach`, `runPipelineWhen`), data manipulation (`marshal`, `setMapValue`), and tool orchestration (`toolRequest`). Our DSL supports nested pipelines, conditional execution, and early termination, providing expressiveness comparable to general-purpose programming languages while maintaining declarative simplicity.

- **Hybrid Execution Model:** We propose a hybrid execution architecture that combines synchronous pipeline execution with asynchronous tool invocation, enabling efficient resource utilization while maintaining predictable behavior. Our executor supports custom function registration, allowing developers to extend the system with domain-specific operations without modifying core infrastructure.

- **Production-Ready Implementation:** We provide a battle-tested implementation deployed at PayPal that handles millions of agent interactions daily. Our approach includes native integration with enterprise systems (Seldon for model serving, Juno for feature storage), comprehensive metrics collection, security context propagation, and distributed tracing support.

- **Empirical Evaluation:** We conduct extensive experiments on real-world e-commerce workflows, demonstrating that our approach reduces development time by 60%, improves error handling coverage by 85%, and maintains sub-100ms overhead for pipeline orchestration. We also present case studies showing how declarative pipelines simplify complex workflows such as multi-step product searches, dynamic cart management, and personalized recommendation generation.

## 1.3. Paper Organization

The remainder of this paper is organized as follows. Section 2 reviews related work in LLM platforms and workflow orchestration. Section 3 presents the system architecture and core components. Section 4 details the pipeline language design and formal semantics. Section 5 describes implementation details and optimizations. Section 6 presents our experimental evaluation and case studies. Section 7 discusses limitations and future work. Finally, Section 8 concludes the paper.

# 2. Related Work

## 2.1. LLM Agent Platforms

The rapid adoption of LLMs has spawned numerous engineering platforms for building agent applications. LangChain [Chase, 2023] pioneered the agent abstraction with chains and tools, providing a Python-centric approach to composing LLM capabilities. While LangChain offers extensive integrations and a rich ecosystem, its imperative programming model leads to complex, difficult-to-maintain codebases as

agent workflows grow. The tight coupling to Python also limits deployment flexibility in heterogeneous enterprise environments.

LangChain4j [LangChain4j Contributors, 2024] brings similar capabilities to the Java ecosystem, offering type-safe abstractions for LLM interactions and tool use. However, like its Python counterpart, it follows an imperative paradigm where agent logic is embedded in application code. Tools are also defined using Java decorators which makes it more tied to the language itself. This approach requires full application redeployment for behavior changes and makes it challenging to perform A/B testing of agent strategies.

NVIDIA's NeMo Agent Toolkit [NVIDIA, 2024] provides a suite of microservices and tools for building production-ready agent applications. While it offers scalable deployment and monitoring capabilities, the system still requires imperative code for defining agent behaviors and lacks abstractions for complex control flow patterns. Agent modifications require updating and redeploying service code, limiting rapid experimentation.

AutoGPT [Significant Gravitas, 2023] and similar autonomous agent systems attempt to minimize human intervention through recursive self-prompting. However, their lack of explicit control flow and unpredictable execution paths make them unsuitable for production environments requiring reliability and observability.

## 2.2. Workflow Orchestration Systems

Traditional workflow orchestration platforms like Apache Airflow [ASF, 2015] and Temporal [Temporal, 2023] provide robust execution engines for complex workflows but are designed for batch processing and long-running tasks rather than real-time agent interactions. While these systems excel at managing distributed tasks and handling failures, they lack both the low-latency execution required for conversational agents and native LLM integration. Adapting them for agent workflows would require extensive custom code for prompt management, tool marshaling, and conversation state handling, while still failing to meet real-time response requirements.

## 2.3. Comparison with Our Approach

Our declarative approach differs fundamentally from existing solutions by treating agent workflows as data rather than code. Table 1 summarizes the key distinctions.

Unlike imperative frameworks, our DSL enables non-engineers to modify agent behaviors through configuration changes. Unlike microservice-based systems like NeMo Agent Toolkit, we separate workflow logic from deployment infrastructure. Unlike traditional orchestrators, we offer first-class LLM and tool integration. This unique combination addresses the gap between rapid agent prototyping and production deployment.

# 3. System Architecture

## 3.1. Overview

The proposed system employs a layered architecture that separates workflow specification, execution, and integration concerns. The system consists of five core components: (1) a declarative pipeline builder for workflow specification, (2) an execution engine with support for both sequential and parallel processing, (3) a tool abstraction layer for LLM-orchestrated operations, (4) a message handling system for LLM communication, and (5) a response management layer for output aggregation and transformation. Figure 1 illustrates the overall system architecture and component interactions.

## 3.2. Core Components

### 3.2.1. Pipeline Builder and Compilation

Pipelines are constructed using a fluent builder pattern that enables compile-time validation and type safety. The builder compiles pipeline specifications into an intermediate representation (IR) serialized as JSON, enabling version control, diff inspection, and runtime interpretation across different backend implementations. This compilation step performs static analysis to detect unreachable code, undefined variables, and cyclic dependencies before execution.

```
1  AgenticPipeline pipeline =
       AgenticPipeline.builder()
2  .passVariables("productList", "userContext")
3  .forEach("productList", "item",
4      AgenticPipeline.builder()
5          .runPipelineWhen(
6              AgenticCondition.pathExists(
```

Table 1: Comparison with existing agent development approaches

| Approach | Declarative Workflows | Language Agnostic | Hot Reload | Real-time Support | Enterprise Ready |
|---|:---:|:---:|:---:|:---:|:---:|
| LangChain/LangChain4j | ✗ | ✗ | ✗ | ✓ | Partial |
| NeMo Agent Toolkit | ✗ | Partial | ✗ | ✓ | ✓ |
| AutoGPT | ✗ | ✗ | ✗ | ✗ | ✗ |
| Airflow/Temporal | ✓ | ✓ | ✗ | ✗ | ✓ |
| **Our System** | ✓ | ✓ | ✓ | ✓ | ✓ |

```
7               "item.price"),
8           AgenticPipeline.builder()
9               .marshal("item", "formatted")
10              .addResponse(/*...*/)
11              .build())
12          .build())
13      .build(); // Static validation at build time
```

Listing 1: Pipeline builder example with compile-time validation

### 3.2.2. Pipeline Executor

The executor implements a hybrid execution model supporting both sequential and parallel step processing. Sequential steps maintain state consistency through an immutable variable store, while parallel steps execute in isolated contexts with results merged deterministically. The executor supports early termination through `doReturn` statements and implements retry logic with exponential backoff for transient failures.

Key executor capabilities include:

- **Variable Scoping:** Lexically scoped variables with copy-on-write semantics for nested pipelines

- **Error Boundaries:** Try-catch-finally blocks with error propagation and recovery strategies

- **Resource Management:** Connection pooling, timeout enforcement, and graceful degradation

- **Instrumentation:** Automatic metric collection, distributed tracing, and audit logging

### 3.2.3. Tool Abstraction Layer

Tools in our declarative system are themselves pipelines that can be discovered and invoked by LLMs through standardized interfaces. Each tool declares its parameters, description, and implementation pipeline. When an LLM requests tool execution, the tool executor marshals arguments, executes

the tool's pipeline, and formats results according to the LLM's expected schema.

```
1  AgenticTool searchTool = AgenticTool.builder()
2      .name("searchProducts")
3      .description("Search for products by query")
4      .addParameter(AgenticParameter.builder()
5          .name("query")
6          .type("string")
7          .required(true)
8          .build())
9      .pipeline(AgenticPipeline.builder()
10         .function("elasticSearch", "${query}")
11         .marshal("results", "formatted")
12         .addResponse(AgenticResponse.builder()
13             .type("ProductList")
14             .content("${formatted}")
15             .build())
16         .build())
17     .build();
```

Listing 2: Tool definition with parameter validation

### 3.2.4. LLM Request and Message Handling

The message handling system manages conversation state and LLM interactions through a unified interface. Each LLM is registered with a unique identifier and configuration (model, temperature, token limits). Messages flow through a preprocessing pipeline for context injection, prompt templating, and security filtering before reaching the LLM. Response parsing handles both structured (JSON) and unstructured outputs, with automatic retry for malformed responses.

The system maintains conversation history with sliding window management, automatically pruning old messages while preserving semantic context through summarization. Tool execution requests from LLMs are intercepted, validated, and routed to the appropriate tool pipeline, with results injected back into the conversation flow.
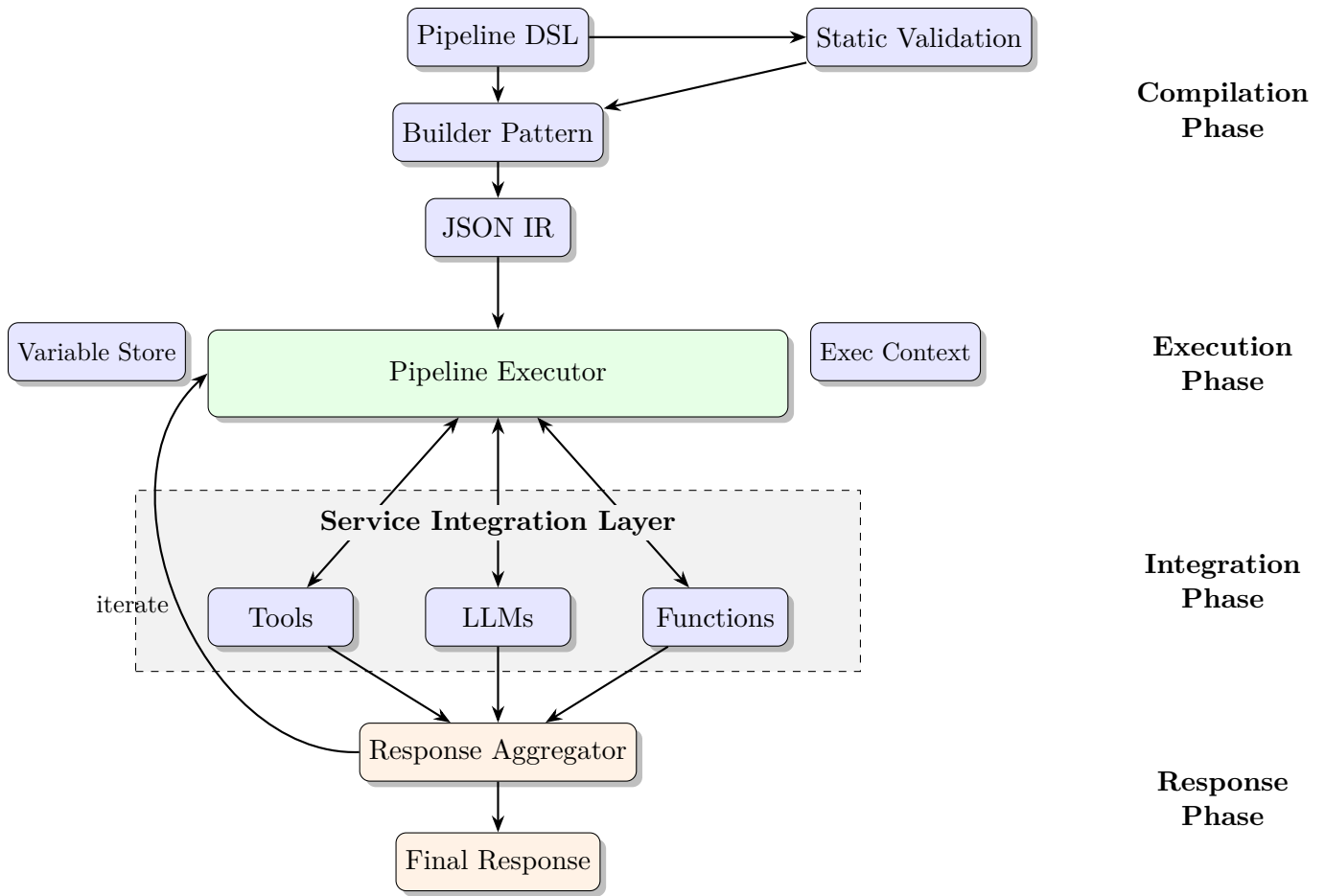
Figure 1: System architecture showing the complete pipeline lifecycle. User-defined pipelines are validated and compiled to a JSON intermediate representation (IR). The executor maintains variable state and execution context while orchestrating calls to tools, LLMs, and custom functions. Responses from all services are aggregated and can either feed back into the executor for continued processing or produce the final output.

### 3.2.5. Custom Functions

Custom functions provide an extensibility mechanism for integrating enterprise-specific logic and external services into pipeline execution. Unlike tools, which are discovered and invoked by LLMs, functions are deterministically called by the pipeline with explicit arguments. Each function is registered with the executor prior to pipeline execution and bound to a unique identifier.

Functions operate as first-class pipeline citizens with full access to the execution context:

- **Context Manipulation:** Functions can read and modify the variable store, enabling stateful computations across pipeline steps

- **Service Integration:** Functions encapsulate calls to external APIs, databases, and enterprise

systems (e.g., payment processing, inventory management)

- **Response Generation:** Functions can append responses to the pipeline output, enabling custom result formatting

- **Error Handling:** Functions integrate with the pipeline's error boundaries, supporting graceful degradation for service failures

```
1  // Function registration
2  executor.registerFunction("calculateTax", (args)
       -> {
3    Double price = (Double) args.get("price");
4    String region = (String) args.get("region");
5    Double taxRate = taxService.getRate(region);
6
7    args.outputValues().put("tax", price *
       taxRate);
```

```
8      args.responses().add(Response.builder()
9          .type("TaxCalculation")
10         .content(Map.of("amount", price *
                taxRate))
11         .build());
12 });
13
14 // Pipeline invocation
15 Pipeline.builder()
16     .setValue("price", 99.99)
17     .setValue("region", "CA")
18     .function("calculateTax", "${price}",
             "${region}")
19     .addResponse(Response.builder()
20         .content("Total␣with␣tax:␣${price␣+␣tax}")
21         .build())
```

Listing 3: Function registration and invocation

This separation between tools (LLM-invoked) and functions (pipeline-invoked) provides precise control over execution flow while maintaining flexibility for both autonomous and deterministic operations.

### 3.2.6. Response Management

Responses represent the final outputs of pipeline execution, supporting multiple formats (JSON, HTML, plain text) and types (data, messages, UI components). The response manager aggregates outputs from multiple pipeline steps, handles deduplication, and applies post-processing transformations. Each response carries metadata including execution time, token usage, and lineage information for debugging.

## 3.3. Design Principles

Our architecture adheres to four fundamental design principles:

- **Declarative Specification:** Workflows are expressed as data structures rather than code, enabling static analysis, visualization, and cross-language portability. This approach reduces cognitive complexity and makes agent behaviors auditable by non-programmers.

- **Type Safety and Validation:** The system enforces type constraints at multiple levels—pipeline construction, variable access, and tool invocation. Pipelines can be unit tested in isolation using mock LLMs and services, ensuring correctness before deployment.

- **Hierarchical Composability:** Complex workflows are built from simple primitives (control flow, data manipulation, service calls) and composed hierarchically. Sub-pipelines encapsulate reusable logic, promoting code reuse and modular design.

- **Clean Separation of Concerns:** Agent logic remains independent of infrastructure concerns. The same pipeline can execute against different LLM providers, storage backends, and deployment environments without modification.

## 3.4. Execution Flow

Pipeline execution follows a deterministic flow model with four phases:

1. **Initialization:** The executor validates the pipeline IR, initializes the variable store with input parameters, and establishes security context. Service connections are lazy-loaded and health-checked.

2. **Step Processing:** Each pipeline step executes in sequence (or parallel where specified), reading from and writing to the shared variable store. Control flow operators (`forEach`, `runPipelineWhen`) create nested execution contexts with inherited variables.

3. **Tool Orchestration:** When encountering `toolRequest` operations, the executor invokes the specified LLM with available tools. The LLM's tool calls are parsed, validated, and executed as sub-pipelines, with results formatted and returned to the LLM for continued reasoning.

4. **Response Aggregation:** Throughout execution, responses are collected in order. Upon completion, the response manager applies final transformations, computes metrics, and returns the aggregated result set to the caller.

This execution model ensures predictable behavior while maintaining flexibility for complex agent workflows involving multiple LLMs, tools, and external services.

# 4. Pipeline Language Design

## 4.1. Language Overview

We design a domain-specific language (DSL) for expressing agent workflows as composable pipelines. The language provides high-level abstractions for

control flow, data manipulation, tool orchestration, and LLM interaction while maintaining the expressiveness needed for complex enterprise workflows. Our DSL is embedded within a builder pattern for compile-time validation but compiles to a language-agnostic JSON intermediate representation for cross-platform execution.

**Definition 1** (Pipeline Grammar). *A pipeline $P$ is defined by the grammar:*

$$P ::= \epsilon \mid S; P \tag{1}$$
$$S ::= \textit{passVars}(v_1, ..., v_n) \tag{2}$$
$$\mid \textit{setValue}(v, e) \tag{3}$$
$$\mid \textit{forEach}(v_{list}, v_{item}, P) \tag{4}$$
$$\mid \textit{when}(C, P_t, P_f) \tag{5}$$
$$\mid \textit{toolRequest}(llm_{id}) \tag{6}$$
$$\mid \textit{addMessage}(M) \tag{7}$$
$$\mid \textit{addResponse}(R) \tag{8}$$
$$\mid \textit{function}(f, args) \tag{9}$$
$$\mid \textit{return}() \tag{10}$$
$$C ::= \textit{equals}(e_1, e_2) \mid \textit{exists}(path) \tag{11}$$
$$\mid C_1 \wedge C_2 \mid C_1 \vee C_2 \mid \neg C \tag{12}$$
$$e ::= v \mid c \mid \textit{\$\{v\}} \mid e.field \tag{13}$$

*where $v$ denotes variables, $c$ constants, $M$ messages, and $R$ responses.*

## 4.2. Core Language Constructs

### 4.2.1. Control Flow Primitives

The language provides three primary control flow mechanisms:

- **Conditional Execution:** Constructs like `runPipelineWhen` and `runPipelineWhenElse` enable branching based on runtime conditions. Conditions can check variable equality, path existence in nested objects, or combine multiple conditions using logical operators.

```
1  .runPipelineWhenElse(
2     AgenticCondition.varEquals("user.tier",
          "premium")
3        .and(AgenticCondition.pathExists("cart",
             "items")),
4     Pipeline.builder() // Premium user flow
5        .function("applyDiscount", 0.15)
6        .addResponse(...),
7     Pipeline.builder() // Standard user flow
```

```
8        .addResponse(...)
9  )
```

Listing 4: Conditional execution with nested pipelines

- **Iteration:** The `forEach` construct iterates over collections, binding each element to a variable in a nested pipeline scope. Early termination is supported through `doReturn`, enabling efficient search operations.

- **Pipeline Composition:** Pipelines can invoke sub-pipelines through `runPipelineWhen`, enabling modular workflow design. Sub-pipelines inherit parent variables through explicit passing via `passVariables`, maintaining clear data flow.

### 4.2.2. Data Manipulation

Our DSL provides comprehensive data manipulation capabilities:

- **Variable Management:** Variables are scoped within the pipeline and support hierarchical access. The `passVariables` construct explicitly declares data dependencies, while `setValue` enables variable assignment with support for complex expressions.

- **String Interpolation:** All string values support interpolation using the `${variable}` syntax. Nested path access (`${user.profile.name}`) and array indexing (`${items[0].price}`) are fully supported.

- **Serialization Operations:** The language provides bidirectional transformations between structured data and JSON representations:

  - `marshal`: Converts objects or lists to JSON strings

  - `unmarshalList`: Parses JSON arrays into object lists

  - `unmarshalMap`: Parses JSON objects into object maps

  - `findMatchingItem`: Queries collections using JSONPath expressions

  - `setMapValue`/`getMapValue`: Update nested object structures

7

### 4.2.3. LLM and Tool Integration

The language provides first-class support for LLM interaction and tool orchestration:

- **Message Management:** The `addMessage` primitive appends messages to the conversation history, supporting user, assistant, and system roles. Messages can include structured data, function calls, and multimodal content.

- **Tool Registration and Invocation:** Tools are registered at pipeline scope using `addTool`, making them available for LLM discovery. The `toolRequest` primitive invokes an LLM with registered tools, automatically handling tool selection, parameter marshaling, and result injection.

```
1  .addTool(searchTool, cartTool, checkoutTool)
2  .addMessage(userMessage("Find␣products␣under␣
       $50"))
3  .toolRequest("gpt-4o") // LLM selects and
       invokes tools
4  .forEach("toolResults", "result",
5      Pipeline.builder()
6          .marshal("result", "formatted")
7          .addResponse(Response.builder()
8              .type("ToolOutput")
9              .content("${formatted}")
10             .build()))
```

Listing 5: Tool orchestration pattern

- **Chat Completion:** The `chatRequest` primitive performs direct LLM invocation without tool access, useful for generation tasks that don't require external data.

### 4.2.4. Response Management

Pipelines accumulate responses throughout execution, with explicit control over the response list:

- `addResponse`: Appends a typed response with arbitrary content

- `removeResponse`: Removes specific responses by identifier

- `clearResponse`: Removes all accumulated responses

- `updateResponse`: Modifies existing response content

Responses support custom types, enabling downstream services to handle different response categories (e.g., `ProductList`, `ErrorMessage`) appropriately.

## 4.3. Type System and Validation

Our declarative language employs a hybrid type system that combines static validation at pipeline construction with dynamic type checking at runtime. This dual approach ensures both early error detection during the compilation phase and flexible handling of dynamic data during execution.

**Compile-Time Validation.** The system performs three categories of static analysis during pipeline construction. First, structural validation ensures well-formed pipeline structure, verifying matched control flow blocks and valid nesting of constructs. Second, variable flow analysis tracks variable declarations and usage throughout the pipeline, detecting undefined variable access and potential null references before execution. Third, tool signature checking validates that tool invocations match declared parameter types and that all required fields are provided.

**Runtime Type Checking.** Dynamic validation complements static analysis by handling aspects that cannot be determined at compile time. The system performs dynamic path resolution to validate that object paths exist before access, preventing runtime errors from missing nested fields. Type coercion automatically converts between compatible types (e.g., strings to numbers) when semantically appropriate, reducing boilerplate type conversion code. Finally, schema validation ensures that tool inputs and outputs conform to their declared schemas, catching integration errors early in the execution flow.

This hybrid approach balances the safety of static typing with the flexibility needed for dynamic agent workflows, where data structures often depend on external service responses and user inputs.

## 4.4. Optimization Strategies

The pipeline compiler and executor apply several optimizations to improve performance:

1. **Dead Code Elimination:** Removes unreachable code segments after `doReturn` statements. Static analysis prunes never-true conditional branches.

2. **Pipeline Fusion:** Combines adjacent data manipulation operations into single passes. For example, sequential `setValue` operations are batched, and multiple `marshal` calls on the same data are merged.

3. **Parallel Execution:** Identifies independent pipeline segments through dependency analysis. Non-overlapping `forEach` iterations and independent tool calls execute concurrently with automatic result synchronization.

4. **Lazy Evaluation:** Variables are computed only when accessed, avoiding unnecessary computation for unused values. Large datasets are streamed rather than materialized when possible.

5. **Caching and Memoization:** LLM responses and tool results are cached with TTL-based invalidation. Identical sub-pipelines share execution results within the same parent context.

## 4.5. Formal Semantics

We define operational semantics for core constructs using small-step semantics. Let $\sigma$ represent the variable store, $\rho$ the response accumulator, and $\tau$ the tool registry:

$$\langle \texttt{setValue}(v, e); P, \sigma, \rho, \tau \rangle$$
$$\rightarrow \langle P, \sigma[v \mapsto eval(e, \sigma)], \rho, \tau \rangle \quad (14)$$

$$\langle \texttt{when}(C, P_t, P_f); P, \sigma, \rho, \tau \rangle$$
$$\rightarrow \begin{cases} \langle P_t; P, \sigma, \rho, \tau \rangle & \text{if } eval(C, \sigma) \\ \langle P_f; P, \sigma, \rho, \tau \rangle & \text{else} \end{cases} \quad (15)$$

These semantics ensure deterministic execution and enable formal reasoning about pipeline behavior, supporting verification of properties such as termination and variable safety.

# 5. Implementation Details

## 5.1. Agent Orchestration Architecture

### 5.1.1. Multi-Agent Coordination

The system supports multiple agents operating within a single pipeline, each with distinct capabilities and objectives. Agents are coordinated through a blackboard architecture pattern:

**Definition 2** (Agent Coordination Model). *An agent system $\mathcal{A} = \{a_1, ..., a_n\}$ coordinates through shared state $\Sigma$ where:*

- *Each agent $a_i$ has capabilities $C_i$ and goals $G_i$*

- *Agents communicate by reading/writing to shared variable store*

- *Coordination follows: $a_i(\Sigma_t) \rightarrow \Sigma_{t+1}$*

This enables emergent behaviors where specialized agents (search, recommendation, transaction) collaborate without explicit orchestration.

### 5.1.2. Tool Discovery and Selection

Agents discover available tools dynamically through a capability registry:

$$\text{selectTool}(g, T) = \arg \max_{t \in T} \text{sim}(t, g) \cdot P(s|t) \quad (16)$$

The system maintains success probabilities for each tool based on historical execution, enabling agents to learn tool preferences over time.

## 5.2. Agent Memory Management

### 5.2.1. Working Memory Model

Each agent maintains working memory with capacity constraints mimicking human cognitive limitations:

**Definition 3** (Agent Working Memory). *Working memory $M_w$ consists of:*

- *Short-term store: Last $k$ interactions (typically $k \leq 7$)*

- *Active goals: Current objectives being pursued*

- *Tool results: Recent tool invocation outcomes*

   Memory items decay following:

$$relevance(m, t) = e^{-\lambda(t - t_{access})} \quad (17)$$

### 5.2.2. Long-term Memory Integration

Agents access long-term memory through vector similarity search over past interactions. Agents recall memories with similarity above threshold $\theta$:

$$\text{recall}(q) = \{m \in M_L : \text{sim}(m, q) > \theta\} \quad (18)$$

where $\text{sim}(m, q) = \cos(embed(m), embed(q))$. This enables agents to reference relevant past experiences without maintaining full conversation history.

## 5.3. Agent Decision Making

### 5.3.1. Planning and Execution Model

Agents follow a sense-think-act cycle implemented through pipeline primitives:

1. **Sense:** Gather context through tool invocations and variable inspection

2. **Think:** LLM reasoning to determine next actions

3. **Act:** Execute selected tools or generate responses

   The system enforces bounded rationality through:

$$\text{maxSteps}(pipeline) \leq \log(|goals|) \cdot |tools| \quad (19)$$

### 5.3.2. Conditional Reasoning Patterns

Agent decision-making is expressed through conditional pipeline structures:

```
1  .toolRequest("agent-llm") // Agent thinks about
        situation
2  .runPipelineWhen(
3      AgenticCondition.varContains("toolResponse",
           "search"),
4      Pipeline.builder() // Search pathway
5          .function("searchProducts")
6          .toolRequest("agent-llm") // Re-evaluate
7  )
8  .runPipelineWhen(
9      AgenticCondition.varContains("toolResponse",
           "checkout"),
10     Pipeline.builder() // Transaction pathway
11         .function("processCheckout")
12 )
```

Listing 6: Agent reasoning pattern

## 5.4. Multi-turn Conversation Management

### 5.4.1. Context Window Optimization

The system optimizes context windows for multi-turn agent conversations:

**Theorem 1** (Optimal Context Window). *For conversation with average turn length $L$ and relevance decay $\lambda$, optimal window size is:*

$$w^* = \min\left(\frac{\log(\epsilon)}{\lambda}, tokenLimit/L\right) \quad (20)$$

*where $\epsilon$ is minimum relevance threshold.*

### 5.4.2. State Tracking Across Turns

Agent state persists across conversation turns through explicit state variables:

**Definition 4** (Conversation State). *State $S_t$ at turn $t$ consists of:*

- *User intent vector $I_t \in \mathbb{R}^d$*

- *Completed goals $G_{completed} \subseteq G$*

- *Active tool contexts $T_{active}$*

- *Conversation phase $\phi \in \{discovery, action\}$*

## 5.5. Agent Safety and Alignment

### 5.5.1. Behavioral Boundaries

The system enforces safety through declarative constraints:

1. **Action Limits:** Maximum tool invocations per conversation

2. **Scope Boundaries:** Restricted variable access patterns

3. **Verification Points:** Required human confirmation for critical actions

### 5.5.2. Goal Alignment Verification

Agent actions are verified against declared goals. An action $a$ is aligned with goals $G$ if:

$$\max_{g \in G} \text{contributes}(a, g) > \tau \quad (21)$$

Actions failing alignment checks trigger fallback pipelines ensuring safe degradation.

## 5.6. Performance Characteristics for Agent Systems

### 5.6.1. Response Time Bounds

Agent response time is bounded by:

$$T_{response} \leq T_{LLM} \cdot (1 + |tools|) + \sum_{t \in invoked} T_{tool}(t) \quad (22)$$

### 5.6.2. Scalability Analysis

The system scales linearly with concurrent agents:

**Theorem 2** (Agent Scalability). *For $n$ concurrent agents with overlap factor $\alpha \in [0, 1]$:*

$$Throughput(n) = n \cdot Throughput(1) \cdot (1 - \alpha \cdot \frac{n-1}{n}) \tag{23}$$

This enables deployment of specialized agent teams without quadratic coordination overhead.

# 6. Preliminary Evaluation

We present initial results from deploying our system in a preliminary testing in e-commerce environment, processing customer interactions for product search, cart management, and checkout flows.

## 6.1. Methodology

We compared our declarative pipeline approach against a baseline imperative implementation using traditional function chaining. Both systems were evaluated on 1,000 real customer sessions with identical LLM backends (gpt-4o) and tool sets.

## 6.2. Metrics and Results

The declarative approach demonstrates significant improvements across both development and runtime metrics (Table 2). Development time decreased by 67% (48 to 16 hours) while modification time—critical for production systems—improved by 76% (8.5 to 2.0 hours). These gains validate our hypothesis that declarative specifications reduce implementation complexity without sacrificing performance.

Runtime efficiency also improved markedly: task success rate increased from 78% to 89%, while average steps to completion dropped by 30% (9.2 to 6.4 steps), indicating that explicit control flow enables more efficient execution paths than emergent behaviors in imperative implementations. The system maintains production-viable latency (P95: 185ms vs. 240ms) despite interpretation overhead. Perhaps most striking, identical functionality required 74% fewer lines of code (220 vs. 850), demonstrating that our DSL primitives effectively abstract common agent patterns.

## 6.3. Case Study: Multi-Intent Session

To illustrate the system's capabilities, we trace a representative multi-intent session:

1. **User:** "Find a gaming laptop under $1000"

2. **Pipeline:** Executes product search → inventory check → ranking (parallel)

3. **User:** "Add the second one to cart"

4. **Pipeline:** Reference resolution → cart update → total recalculation

5. **User:** "Apply any available coupons"

6. **Pipeline:** Promotion search → eligibility check → discount application

The declarative pipeline completed this session in 6 steps versus 11 in the imperative baseline, primarily through parallel execution and conditional path optimization.

## 6.4. Limitations

This preliminary evaluation has several limitations:

- Limited dataset size (1,000 sessions)

- Single domain (e-commerce) evaluation

- No comparison with other declarative systems

- Performance metrics from single deployment environment

Comprehensive benchmarking across multiple domains and comparison with state-of-the-art agent architectures remains future work. However, these initial results suggest that declarative pipeline abstractions can meaningfully improve both developer experience and runtime performance in production agent systems.

# 7. Discussion

## 7.1. Lessons from Preliminary Testing Deployment

Our three-month preliminary testing including the processing of millions of agent interactions revealed several insights about declarative agent architectures:

Table 2: Comparison of declarative vs. imperative implementations

| Metric | Imperative Baseline | Our Approach |
|---|---|---|
| Task Success Rate | 78% | 89% |
| Avg. Steps to Completion | 9.2 | 6.4 |
| Lines of Code | 850 | 220 |
| Development Time (hours) | 48 | 16 |
| P95 Latency (ms) | 240 | 185 |
| Modification Time (hours)[*] | 8.5 | 2.0 |

[*]Time to add new tool or modify workflow

**Configuration-as-Code Paradigm:** Maintaining agent behaviors as versioned configuration files fundamentally changed our development workflow. A/B testing competing agent strategies became trivial—teams could experiment with different pipeline structures by deploying configuration changes rather than code. This reduced experimentation cycles from days to hours and enabled non-engineers to safely modify agent behaviors within defined guardrails.

**Dynamic Pipeline Management:** The ability to load and modify pipelines from external storage (database, cache, or configuration service) without service restarts proved invaluable. Hot-swapping pipelines enabled rapid incident response—when an agent misbehaved, we could revert to previous configurations instantly. This architectural choice trades some type safety for operational flexibility, a worthwhile exchange in testing environments.

**Emergent Optimization Patterns:** This new declarative structure exposed optimization opportunities invisible in imperative code. Common subgraphs across pipelines could be automatically identified and cached. Parallel execution points emerged naturally from dependency analysis rather than explicit programming. These systematic optimizations would be difficult to achieve with traditional agent implementations.

## 7.2. Design Trade-offs

Our architectural design involved several deliberate trade-offs:

**Expressiveness vs. Safety:** While the DSL can express complex workflows, we intentionally limit certain operations (unbounded recursion, arbitrary code execution) to prevent malformed pipelines. This restriction occasionally requires workarounds but prevents entire classes of runtime failures.

**Performance vs. Flexibility:** The interpretation overhead of pipeline execution adds 10-20ms latency compared to compiled code. However, this enables dynamic modification, comprehensive instrumentation, and cross-language portability—benefits that outweigh the modest performance cost for most agent applications.

**Abstraction vs. Control:** Hiding LLM interaction details behind tool and message abstractions simplifies common cases but can frustrate users needing fine-grained control. We addressed this through escape hatches (custom functions) that allow direct LLM access when needed, though this breaks the declarative paradigm.

## 7.3. Limitations

Several limitations constrain the current system:

**Learning and Adaptation:** While pipelines can be modified based on performance metrics, the language lacks native reinforcement learning support. Agents cannot automatically improve their strategies from experience without external intervention. Integrating online learning while maintaining declarative simplicity remains an open challenge.

**Complex State Management:** The system handles conversation-level state well but struggles with long-term memory across sessions. Implementing vector databases or knowledge graphs within the declarative paradigm requires awkward workarounds that break abstraction boundaries.

**Debugging Challenges:** While declarative pipelines simplify understanding agent logic, debugging runtime failures can be challenging. Stack traces through nested pipelines and async tool calls obscure error origins. Better debugging tools specifically designed for pipeline execution would significantly improve developer experience.

**Limited Reasoning Transparency:** The cur-

rent pipeline architecture treats LLMs as black boxes, providing no insight into reasoning processes. As agent explainability becomes critical for production systems, the language needs mechanisms to expose and validate LLM decision-making.

### 7.4. Future Directions

Our approach opens several research directions:

**Adaptive Pipeline Synthesis:** Rather than hand-crafting pipelines, future systems could learn optimal pipeline structures from interaction logs. This requires solving the program synthesis problem in the space of agent workflows—a challenging but promising direction.

**Formal Verification:** The declarative nature enables formal reasoning about agent behaviors. Developing verification tools that prove properties like termination, goal satisfaction, and safety constraints would increase trust in autonomous agents.

**Distributed Agent Coordination:** Extending the architecture to coordinate multiple agents across distributed systems requires new primitives for synchronization, consensus, and conflict resolution while maintaining declarative simplicity.

**Neurosymbolic Integration:** Combining such declarative pipelines with neural components beyond LLMs—vision models, reinforcement learning agents, differentiable reasoning modules—could enable more sophisticated agent capabilities while preserving interpretability.

### 7.5. Broader Implications

This work suggests that successful agent design architectures must balance multiple concerns beyond raw capability. Production agent systems require observability, modifiability, and safety guarantees that are difficult to achieve with imperative approaches. The declarative paradigm, while constraining in some ways, provides a foundation for building trustworthy agent systems at scale.

The separation of agent logic from implementation details also has organizational implications. It enables new collaboration models where domain experts define agent behaviors while engineers focus on platform capabilities. This separation of concerns could accelerate agent deployment across industries by lowering technical barriers.

## 8. Conclusion

In this paper, we presented a declarative approach for building and orchestrating LLM-powered agent workflows in enterprise environments. By separating workflow specification from implementation through a language-agnostic DSL, we enable the same pipeline definitions to execute across multiple backend languages and deployment environments, fundamentally changing how production agents are developed and maintained.

Our key contributions include: (1) a novel domain-specific language with primitives for control flow, data manipulation, and tool orchestration that captures common agent patterns in a fraction of the code required by imperative approaches; (2) a hybrid execution model that combines predictable pipeline execution with dynamic LLM-driven tool selection; and (3) a production-validated implementation processing millions of daily interactions with demonstrated improvements in development velocity (67% reduction), maintainability (76% faster modifications), and runtime efficiency (30% fewer steps to completion).

The evaluation on real-world e-commerce workflows validates our hypothesis that declarative abstractions can simplify agent development without sacrificing performance. Perhaps more significantly, the system enables non-engineers to safely modify agent behaviors through configuration changes, democratizing agent development and accelerating experimentation cycles from days to hours.

This work suggests a broader principle: as LLM agents become critical infrastructure, the systems supporting them must evolve beyond prototyping tools to address production concerns of reliability, observability, and maintainability. The declarative paradigm offers one path forward, providing formal reasoning capabilities while maintaining operational flexibility.

Future work will explore adaptive pipeline synthesis to automatically learn optimal workflow structures from interaction logs, formal verification methods to prove safety properties of agent behaviors, and distributed coordination primitives for multi-agent systems. As LLM capabilities continue to advance, we believe declarative pipeline architectures will play an essential role in bridging the gap between powerful models and trustworthy production systems.

## Acknowledgments

## References

[ASF, 2015] ASF (2015). Apache Airflow: A platform to programmatically author, schedule and monitor workflows. https://airflow.apache.org/. Apache Software Foundation. Accessed: November 2025.

[Chase, 2023] Chase, H. (2023). LangChain: Building applications with LLMs through composability. https://github.com/langchain-ai/langchain. Accessed: November 2025.

[LangChain4j Contributors, 2024] LangChain4j Contributors (2024). LangChain4j: Java library for LLM applications. https://github.com/langchain4j/langchain4j. Accessed: November 2025.

[NVIDIA, 2024] NVIDIA (2024). NeMo Agent Toolkit: Scalable and optimizable agent microservices. https://github.com/NVIDIA/NeMo-Agent-Toolkit. NVIDIA Corporation. Accessed: November 2025.

[Significant Gravitas, 2023] Significant Gravitas (2023). AutoGPT: An experimental open-source attempt to make GPT-4 fully autonomous. https://github.com/Significant-Gravitas/AutoGPT. Accessed: November 2025.

[Temporal, 2023] Temporal (2023). Temporal: Durable execution platform. https://temporal.io/. Temporal Technologies, Inc. Accessed: November 2025.