PROVISIONAL PATENT APPLICATION
COVER SHEET
Title of Invention: UltraLLMOrchestrator System (UltrAI)
Inventor(s): Joshua Field 1625 SW Alder St. #510 Portland, OR 97205 Citizenship: United States of America
Correspondence Address: Joshua Field 1625 SW Alder St. #510 Portland, OR 97205
Application Filed By: Inventor
_____


## SPECIFICATION

Title: UltraLLMOrchestrator System (UltrAI)
Cross-Reference to Related Applications
[Not Applicable]


### Field of the Invention

This invention relates to advanced artificial intelligence systems, specifically to orchestrating and integrating multiple Large Language Models (LLMs) in a unified, extensible, and dynamically configurable analysis platform. The invention addresses the need for robust, multi-model reasoning, dynamic model management, and developer-friendly extensibility in AI-driven analytical workflows.

UltraLLM Orchestrator


### Background of the Invention

Large Language Models (LLMs) have revolutionized natural language processing, but the constraints of single-model architectures often limit their deployment in real-world analytical tasks. Individual LLMs, while powerful, are subject to inherent biases, knowledge gaps, and inconsistent reasoning. Current solutions typically rely on static, single-model pipelines or ad hoc multi-model integrations that lack structure, extensibility, and systematic quality control.

Existing LLM orchestration approaches are often rigid, requiring significant codebase changes to add new models or modify analytical workflows. They lack a unified framework for structured, multi-stage analysis and do not provide mechanisms for runtime extensibility, automated quality evaluation, or developer onboarding.

As a result, organizations and developers face high barriers to integrating new models, experimenting with orchestration strategies, and ensuring the reliability and depth of AI-generated insights.

There is a critical need for a system that: effectively orchestrates multiple LLMs via configurable, pattern-driven analysis pipelines; enables runtime registration and deregistration of models and orchestration logic; provides automated, multi-dimensional quality assessment of model outputs; facilitates rapid developer onboarding and integration of new models, patterns, and business logic; and delivers robust, contextually aware analytical outputs that exceed the capabilities of any individual LLM.

#### Technical Problems

1. Single-Model Limitations: Individual LLMs exhibit biases, knowledge gaps, and inconsistent reasoning that limit their effectiveness for complex analytical tasks.
2. Static Integration Challenges: Traditional LLM integration requires code modifications to add new models or change analysis workflows, creating development overhead and deployment delays.
3. Quality Evaluation: Assessing LLM output quality typically requires manual review, which doesn't scale and introduces subjectivity.
4. Developer Experience: Building with multiple LLMs typically requires mastering different APIs, authentication methods, and response formats.
5. Context Management: LLMs struggle with limited context windows and integrating document knowledge effectively.


### Summary of the Invention

The present invention, UltraLLMOrchestrSystem (UltrAI), is a next-generation AI analysis platform that orchestrates multiple LLMs through a dynamically configurable, multi-stage pipeline. UltrAI enables users and developers to:

- Register, configure, and orchestrate any number of LLMs at runtime, without codebase modification or downtime.

- Define and extend analysis patterns—structured, multi-stage workflows that guide LLMs through critique, fact-checking, perspective-taking, synthesis, and more.
- Leverage a plugin-like architecture for both models and orchestration logic, supporting rapid experimentation and integration of new AI capabilities.
- Automatically evaluate the quality of each model's output across multiple dimensions (coherence, technical depth, strategic value, uniqueness), with metrics tracked for continuous improvement.
- Onboard new developers and models quickly, supported by comprehensive documentation, onboarding guides, and a robust test suite.

Key innovations and advantages include:

- Dynamic model registry: Models can be added, removed, or reconfigured at runtime via a unified API, enabling rapid adaptation to new LLMs and providers.
- Configurable orchestration pipeline: The sequence, logic, and participants of each analysis stage are defined by extensible patterns, supporting arbitrary workflows and model combinations.
- Extensible plugin system: Both models and orchestration strategies are implemented as plugins, allowing for easy extension, replacement, or customization without core code changes.
- Automated, multi-dimensional quality evaluation: Each model's output is scored and tracked, enabling data-driven model selection, performance tuning, and transparent reporting.
- Developer-centric design: The system is built for rapid onboarding, with clear extension points, thorough documentation, and a comprehensive test suite to ensure reliability and ease of integration.
- Seamless document/context integration: Uploaded documents are processed, chunked, and incorporated as context in analysis tasks, maximizing the relevance and depth of LLM outputs.
- Integrated business logic and pricing: The platform supports flexible pricing models, usage tracking, and business logic, making it suitable for commercial deployment.

Technical Solutions

1. Single-Model Limitation Solution: The system orchestrates multiple models through configured patterns, leveraging their complementary strengths while mitigating individual weaknesses. Models critique and enhance each other's work, while the synthesis stage combines insights into a comprehensive analysis.
2. Static Integration Solution: The dynamic model registry and plugin architecture enable runtime registration of new models and orchestration patterns without code changes. Runtime registration dramatically reduces integration time and allows rapid adaptation to new LLM releases.
3. Quality evaluation Solution: The automated quality evaluation system scores outputs across multiple dimensions, providing objective metrics for comparing model performance and tracking improvements over time.
4. Developer Experience Solution: The unified orchestration API abstracts provider differences, offering consistent interfaces, comprehensive documentation, and a robust testing framework that significantly improves developer productivity.
5. Context Management Solution: The document processing and context management system intelligently chunks and integrates relevant information, maximizing the effective use of context windows across models.

# Detailed Description of the Invention

1. System Overview: UltrAI is a comprehensive AI analysis platform orchestrating multiple LLMs through configurable, pattern-based pipelines. The system provides enhanced analytical quality through multi-model synthesis and critique, rapid integration of new models and orchestration strategies, and automated, transparent quality evaluation and reporting.

2. System Architecture The UltrAI system employs a modular, client-server web architecture as illustrated in the accompanying figures.

   2.1. Frontend
       2.1.1. Built as a modern React SPA (TypeScript, Vite, Tailwind CSS).
       2.1.2. Provides intuitive interfaces for prompt input, document upload, model/pattern selection, and result visualization.
       2.1.3. Supports responsive design and rich user experience features (history, sharing, real-time feedback).
   2.2. Backend
       2.2.1. Python API using FastAPI, serving as the central orchestrator.
       2.2.2. Modularized with routers for analysis, documents, models, pricing, users, and health.
       2.2.3. Manages application state, business logic, and external integrations.
   2.3. Core Orchestration Engine (MultiLLMOrchestrator)

        2.3.1.     Receives analysis requests and orchestrates LLMs through configurable, multi-stage pipelines.
        2.3.2.     Supports runtime registration/deregistration of LLMs and orchestration patterns.
        2.3.3.     Enables both parallel and sequential execution of LLM calls.
        2.3.4.     Tracks token usage, quality metrics, and intermediate outputs for each stage
        2.3.5.     Extensible via subclassing or plugin registration, allowing custom orchestration logic and new analysis patterns.
    2.4.    Document Processing Module (UltraDocumentsOptimized)
        2.4.1.     Handles file uploads, text extraction, and intelligent chunking.
        2.4.2.     Integrates document context into analysis workflows, maximizing LLM relevance and depth
    2.5.    Pricing & Business Logic Module
        2.5.1.     Calculates costs based on token usage, model selection, and feature tiers.
        2.5.2.     Implements markups, discounts, and token efficiency factors for flexible business models.
        2.5.3.     Manages user authorization, usage tracking, and billing
    2.6.    Database
        2.6.1.     Stores user accounts, document metadata, analysis history, pricing configs, and cached results.
        2.6.2.     Supports collaborative features via shared analysis links.
    2.7.    External LLM Services Integration
        2.7.1.     Securely connects to third-party LLM providers (OpenAI, Anthropic, Google, Mistral, etc.).
        2.7.2.     Abstracts API differences and manages credentials.
    2.8.    Caching Layer
        2.8.1.     In-memory TTL cache (cachetools) for fast retrieval of repeated analysis requests.
        2.8.2.     Reduces latency and operational costs
    2.9.    Error Monitoring
        2.9.1.     Integrates with real-time error tracking (e.g., Sentry) for robust monitoring and rapid issue resolution.

3.    Detailed Features

    3.1.    Core Analysis & Orchestration
        3.1.1.     Multi-LLM Orchestration
            3.1.1.1.     Executes analysis tasks across user-selected LLMs, leveraging their complementary strengths.
            3.1.1.2.     Supports dynamic, runtime model registration and deregistration.
            3.1.1.3.     Orchestration pipelines are fully configurable, supporting arbitrary stages (e.g., initial, meta, synthesis) and model sets.
            3.1.1.4.     Enables parallel and sequential LLM execution, with intermediate outputs tracked and organized.
            3.1.1.5.     Token usage and quality metrics are collected for each model and stage, supporting transparent reporting and optimization.
        3.1.2.     Pattern-Based Analysis
            3.1.2.1.     Guides LLM interaction using extensible, structured patterns (e.g., critique, fact-check, perspective, confidence analysis).
            3.1.2.2.     Each pattern defines a unique, multi-stage workflow with custom prompt templates and instructions.
            3.1.2.3.     Patterns can be added, modified, or replaced without core code changes, supporting rapid experimentation.
        3.1.3.     "Ultra" Model Synthesis
            3.1.3.1.     Designates a user-selected "Ultra" LLM to synthesize intermediate results into a final, enhanced output.
            3.1.3.2.     Synthesis logic is pattern-driven and extensible, supporting new synthesis strategies as plugins.
    3.2.    Model Registration & Extensibility
        3.2.1.     New LLMs can be registered at runtime via a simple API, without codebase modification or redeployment.
        3.2.2.     Orchestration strategies and analysis patterns can be extended or replaced via subclassing or plugin registration.
        3.2.3.     Comprehensive documentation and a robust test suite support rapid developer onboarding and integration.
    3.3.    Document Handling
        3.3.1.     Uploaded documents are processed, chunked, and stored for use as context in analysis tasks.
        3.3.2.     Intelligent context integration maximizes the relevance and depth of LLM outputs.
    3.4.    Pricing, Billing & Business Model

3.4.1.    Flexible pricing models based on token usage, model selection, and feature
                  tiers.
        3.4.2.    Supports markups, discounts, token efficiency factors, and business model
                  simulation.
        3.4.3.    Usage tracking, reporting, and billing integration for commercial deployment.
    3.5.    User Interface & Experience
        3.5.1.    Intuitive, responsive UI for prompt input, document management, model/pattern
                  selection, and result visualization.
        3.5.2.    Features include history management, sharing, real-time feedback, and dynamic
                  cost estimation.
    3.6.    System Operations & Development
        3.6.1.    Modular, testable codebase with clear extension points for models, patterns,
                  and business logic.
        3.6.2.    Comprehensive test suite, health/metrics endpoints, and performance monitoring.
        3.6.3.     Supports configuration via environment variables/files and scalable deployment
                  (Docker, Vercel, etc.).

4.  Enablement Details

    4.1.    Dynamic Model Registry Implementation: The dynamic model registry is implemented
            through an abstraction layer that:
        4.1.1.    Defines a provider-agnostic interface for LLM interactions
        4.1.2.    Employs the adapter pattern to normalize different LLM APIs
        4.1.3.    Implements connection pooling for efficient resource management
        4.1.4.    Supports runtime configuration updates without system restart
        4.1.5.    Maintains capability metadata for intelligent model selection
        4.1.6.    Provides health monitoring and circuit breaking for reliability
        4.1.7.    Implementation involves:
            4.1.7.1.    A central registry class that manages model registrations
            4.1.7.2.    Provider-specific adapter classes implementing a familiar interface
            4.1.7.3.    Dynamic loading of adapter modules through dependency injection
            4.1.7.4.    Capability discovery through provider API introspection
            4.1.7.5.    Configuration validation using JSON schema

    4.2.    Plugin Architecture Enablement: The plugin architecture for extending orchestration
            logic:
        4.2.1.    Uses a modular design with well-defined interfaces
        4.2.2.    Implements a registry pattern for pattern discovery and loading
        4.2.3.    Supports hot-swapping of components without service interruption
        4.2.4.    Allows for the composition of patterns and sub-patterns
        4.2.5.    Provides a templating system for prompt generation
        4.2.6.    Implementation details:
            4.2.6.1.    Interface definitions for orchestration patterns
            4.2.6.2.    Factory methods for pattern instantiation
            4.2.6.3.    Configuration schema for pattern parameterization
            4.2.6.4.    Service locator pattern for discovering available patterns
            4.2.6.5.    Hook points for extending existing patterns

    4.3.    Quality Evaluation System: The automated quality evaluation system:
        4.3.1.    Defines multi-dimensional scoring criteria (coherence, accuracy, etc.)
        4.3.2.    Implements both rule-based and model-based evaluation strategies
        4.3.3.    Supports customizable scoring algorithms for different use cases
        4.3.4.    Maintains historical performance data for trend analysis
        4.3.5.    Provides visualization and reporting capabilities
        4.3.6.    Implementation details:
            4.3.6.1.    Evaluation dimension definitions and algorithms
            4.3.6.2.    Model response parsing and feature extraction
            4.3.6.3.    Statistical normalization of scores across dimensions
            4.3.6.4.    Persistent storage of quality metmaintain arting and visualization
                        components

5.  **Alternative Embodiments**

    5.1.    Edge Computing LLM Orchestration: The UltrAI system can be adapted for edge computing
            environments where bandwidth or privacy concerns limit cloud LLM usage. In this
            embodiment:
        5.1.1.    Lightweight local LLMs run on edge devices
        5.1.2.    Orchestration logic prioritizes local models, falling back to cloud models only
                  when necessary
        5.1.3.    Smart caching reduces repeated cloud queries
        5.1.4.    Privacy-sensitive data remains on local devices

           5.1.5.    Synchronization protocols maintaina consistent model registry across distributed systems
- 5.2.    Real-time Collaborative Analysis: Users can extend the system to support real-time collaborative analysis, where multiple users:
  - 5.2.1.    Simultaneously contribute to prompt refinement
  - 5.2.2.    Select different model combinations for comparative analysis
  - 5.2.3.    Add supplementary documents and context in real-time
  - 5.2.4.    Comment on and annotate intermediate results
  - 5.2.5.    Vote on preferred analysis patterns and synthesis approaches
  - 5.2.6.    Track revisions and contributions through a version control mechanism
- 5.3.    Autonomous Adaptive Orchestration:  An advanced embodiment leverages reinforcement learning to autonomously optimize orchestration:
  - 5.3.1.    The system learns optimal model selection for different query types
  - 5.3.2.    Analysis patterns evolve based on quality metrics and user feedback\
  - 5.3.3.    Prompt refinement strategies improve automatically through usage
  - 5.3.4.    Cost-optimization balances quality against token usage
  - 5.3.5.    Anomaly detection identifies potential issues in model outputs
  - 5.3.6.    Recovery strategies mitigate LLM hallucinations and errors
- 5.4.    Domain-Specific Orchestration: Users can leverage the system for vertical domains through:
  - 5.4.1.    Custom domain-specific analysis patterns (e.g., medical diagnosis, legal analysis)
  - 5.4.2.    Industry-specific document processing and knowledge integration
  - 5.4.3.    Specialized evaluation metrics for domain-specific quality assessment
  - 5.4.4.    Integration with domain knowledge bases and ontologies
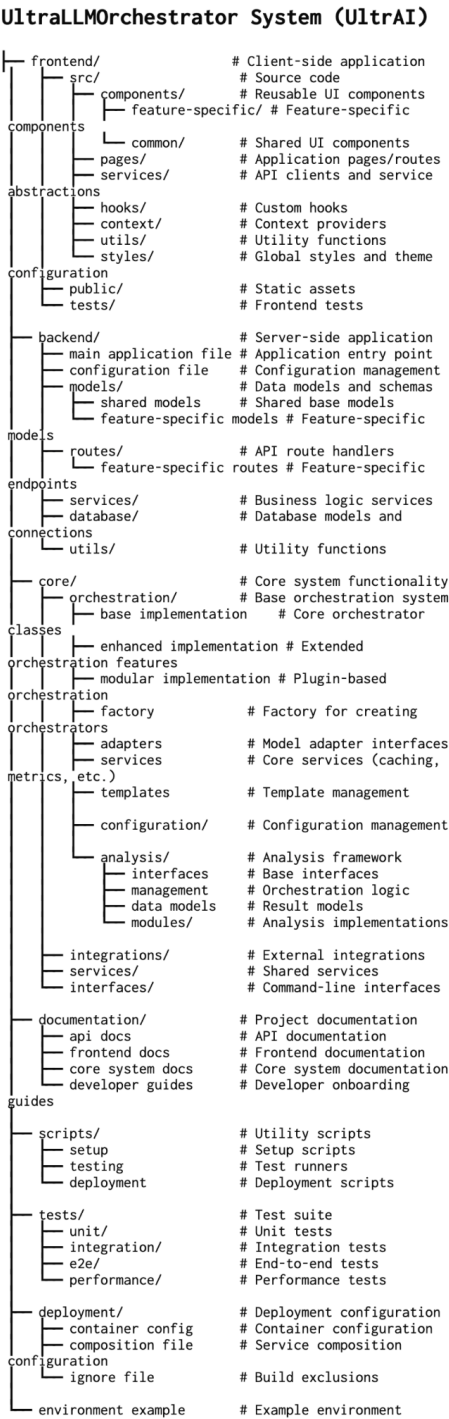  - 5.4.5.    Compliance checks for domain-specific regulations and standards

## Figures/Diagrams

Brief Description of the Drawings

- Figure 1: Generalized Directory Structure: The UltraLLMOrchestrator (UltrAI) follows a comprehensive directory structure emphasizing modularity, concerns separation, and clear functionality organization.
- Figure 2: Dynamic Model Registration Flowchart: Sequence of steps for adding a new LLM at runtime, including parameter validation, provider connectivity check, registry update, and success/error handling.
- Figure 3: Analysis Stage Pipeline Flowchart: Internal orchestration loop that executes a multi-stage analysis pattern—preparing prompts, querying models, validating responses, tracking metrics, and performing final synthesis.
- Figure 4: Model Quality Evaluation Flowchart: Process for scoring model responses across multiple quality dimensions, normalizing and aggregating the scores, updating historical metrics, and returning a comprehensive assessment.
- Figure 5: End-to-End Analysis Workflow Flowchart: High-level view of the complete user journey: from submitting an analysis request and (optionally) processing documents, through staged orchestration and synthesis, to result delivery, caching, and usage-metric recording.

# FIGURE 1

Generalized Directory Structure: The UltraLLMOrchestrator (UltrAI) follows a comprehensive directory structure that emphasizes modularity, separation of concerns, and clear organization of functionality:

```
UltraLLMOrchestrator System (UltrAI)

├── frontend/              # Client-side application
│   ├── src/               # Source code
│   │   ├── components/     # Reusable UI components
│   │   │   ├── feature-specific/ # Feature-specific
components
│   │   │   └── common/      # Shared UI components
│   │   ├── pages/         # Application pages/routes
│   │   ├── services/      # API clients and service
abstractions
│   │   ├── hooks/         # Custom hooks
│   │   ├── context/       # Context providers
│   │   ├── utils/         # Utility functions
│   │   └── styles/        # Global styles and theme
configuration
│   ├── public/            # Static assets
│   └── tests/             # Frontend tests
│
├── backend/               # Server-side application
│   ├── main application file # Application entry point
│   ├── configuration file    # Configuration management
│   ├── models/            # Data models and schemas
│   │   ├── shared models     # Shared base models
│   │   └── feature-specific models # Feature-specific
models
│   ├── routes/            # API route handlers
│   │   └── feature-specific routes # Feature-specific
endpoints
│   ├── services/          # Business logic services
│   ├── database/          # Database models and
connections
│   └── utils/             # Utility functions
│
├── core/                  # Core system functionality
│   ├── orchestration/     # Base orchestration system
│   │   ├── base implementation    # Core orchestrator
classes
│   │   ├── enhanced implementation # Extended
orchestration features
│   │   ├── modular implementation # Plugin-based
orchestration
│   │   ├── factory           # Factory for creating
orchestrators
│   │   ├── adapters          # Model adapter interfaces
│   │   ├── services          # Core services (caching,
metrics, etc.)
│   │   ├── templates         # Template management
│   │   │
│   │   ├── configuration/    # Configuration management
│   │   │
│   │   └── analysis/         # Analysis framework
│   │       ├── interfaces    # Base interfaces
│   │       ├── management    # Orchestration logic
│   │       ├── data models   # Result models
│   │       └── modules/      # Analysis implementations
│   │
│   ├── integrations/      # External integrations
│   ├── services/          # Shared services
│   └── interfaces/        # Command-line interfaces
│
├── documentation/         # Project documentation
│   ├── api docs           # API documentation
│   ├── frontend docs      # Frontend documentation
│   ├── core system docs   # Core system documentation
│   └── developer guides   # Developer onboarding
guides
│
├── scripts/               # Utility scripts
│   ├── setup              # Setup scripts
│   ├── testing            # Test runners
│   └── deployment         # Deployment scripts
│
├── tests/                 # Test suite
│   ├── unit/              # Unit tests
│   ├── integration/       # Integration tests
│   ├── e2e/               # End-to-end tests
│   └── performance/       # Performance tests
│
├── deployment/            # Deployment configuration
│   ├── container config   # Container configuration
│   ├── composition file   # Service composition
configuration
│   └── ignore file        # Build exclusions
│
└── environment example    # Example environment

/
```

6

FIGURE 2

Dynamic Model Registration Flowchart: Sequence of steps for adding a new LLM at runtime, including parameter validation, provider connectivity check, registry update, and success/error handling
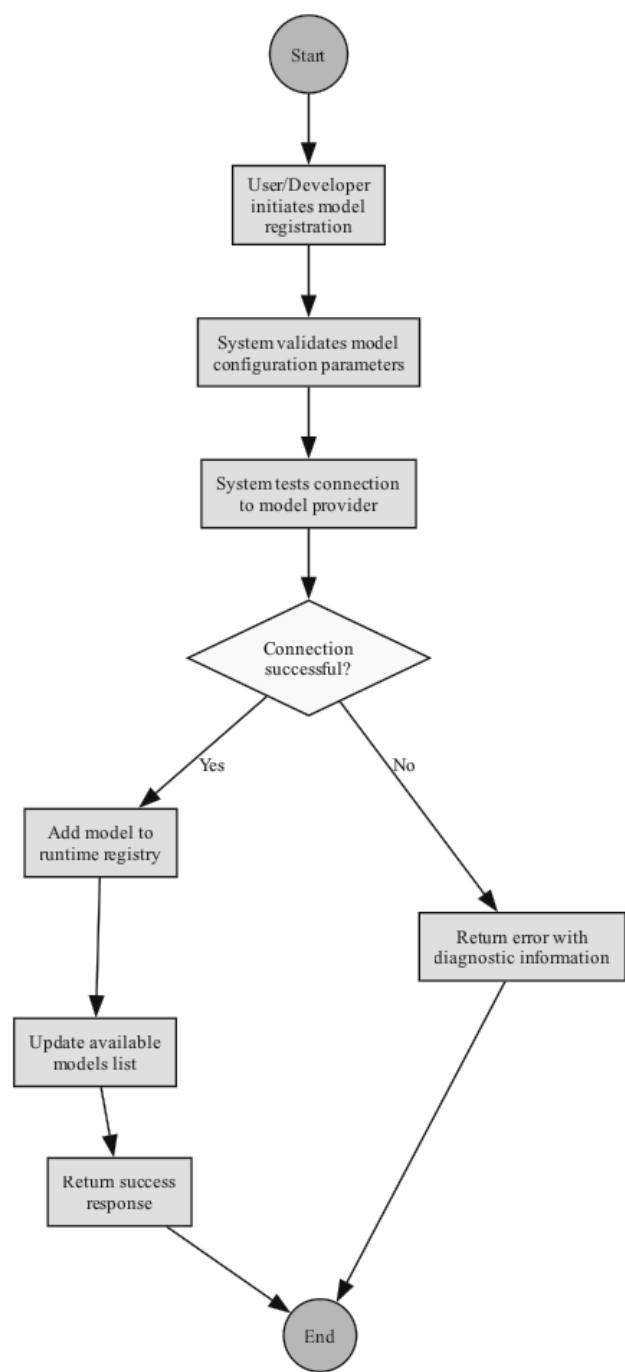
"

# FIGURE 3

Analysis Stage Pipeline Flowchart: Internal orchestration loop that executes a multi-stage analysis pattern—preparing prompts, querying models, validating responses, tracking metrics, and performing final synthesis.
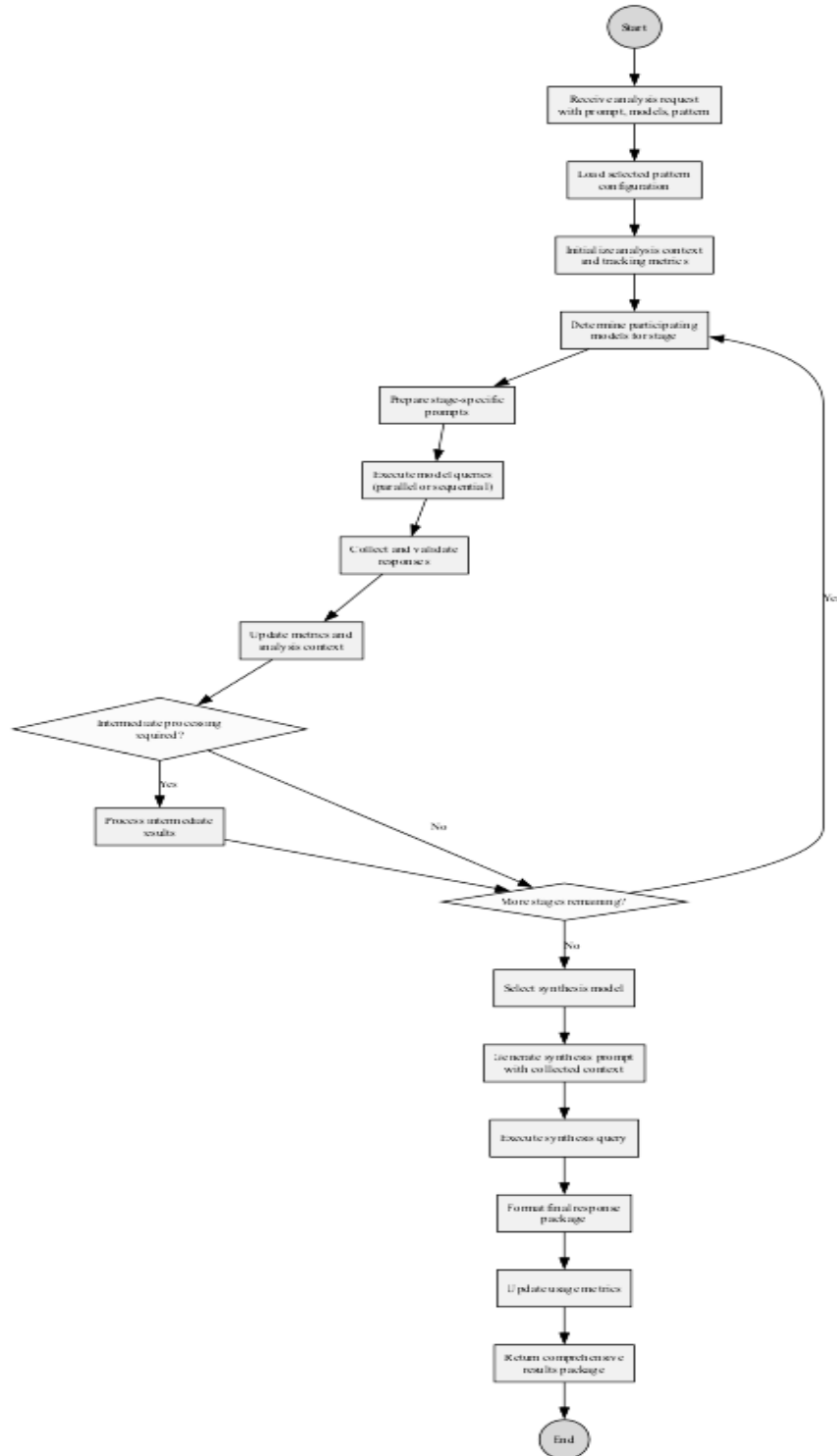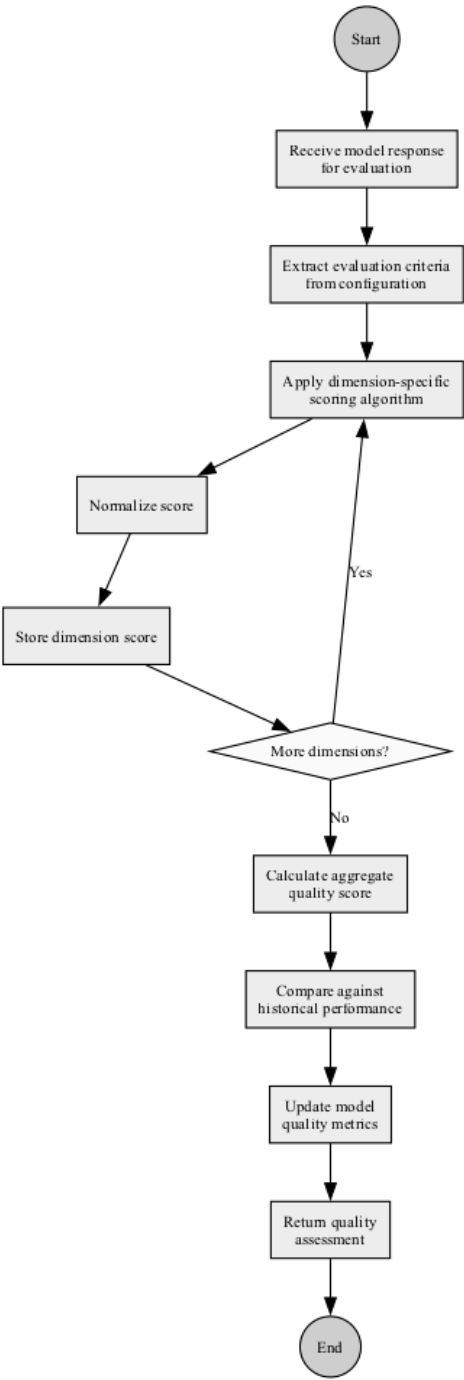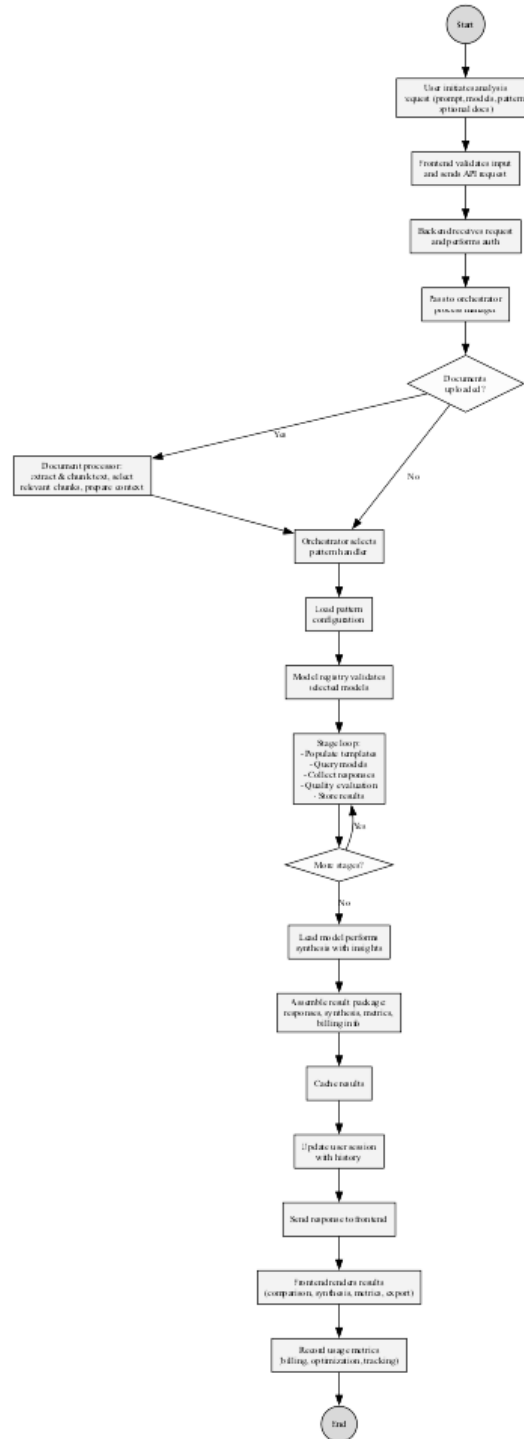
FIGURE 4

Model Quality Evaluation Flowchart: Process for scoring model responses across multiple quality dimensions, normalizing and aggregating the scores, updating historical metrics, and returning a comprehensive assessment."

FIGURE 5

End-to-End Analysis Workflow Flowchart: High-level view of the complete user journey: from submitting an analysis request and (optionally) processing documents, through staged orchestration and synthesis, to result delivery, caching, and usage-metric recording.

# Abstract

The UltraLLMOrchestrator System (UltrAI) is a robust, extensible AI analysis platform that orchestrates multiple Large Language Models (LLMs) through dynamically configurable, pattern-based pipelines. UltrAI enables runtime registration and orchestration of any number of LLMs, supports extensible analysis patterns, and provides automated, multi-dimensional quality evaluation. The system's modular architecture, illustrated in the diagrams, supports rapid integration of new models, patterns, and business logic, making it suitable for research and commercial deployment. UltrAI delivers enhanced analytical quality, developer-friendly extensibility, and seamless adaptation to advances in AI, setting a new standard for multi-model orchestration and AI-driven analysis.

# Claims

## Core System Claims

1. A system for orchestrating multiple Large Language Models (LLMs), comprising:
   - 1.1. a dynamic model registry configured to register and deregister LLM models at runtime;
   - 1.2. a configurable orchestration pipeline supporting multi-stage analysis workflows;
   - 1.3. an extensible pattern framework defining structured analytical processes;
   - 1.4. a quality evaluation module for automated assessment of model outputs; and
   - 1.5. a synthesis engine for combining results from multiple models.
2. The system of claim 1, wherein the dynamic model registry supports adding new LLM providers without requiring code modifications.
3. The system of claim 1, wherein the configurable orchestration pipeline supports both parallel and sequential execution of LLM queries.
4. The system of claim 1, wherein the extensible pattern framework includes predefined patterns for critique, fact-checking, perspective analysis, and confidence assessment.
5. A method for orchestrating multiple Large Language Models (LLMs) for enhanced analysis, comprising:
   - 5.1. registering multiple LLM models in a dynamic registry;
   - 5.2. receiving an analysis request with a prompt and selected models;
   - 5.3. executing the prompt across models chosen according to a configurable pattern;
   - 5.4. evaluating the quality of responses across multiple dimensions;
   - 5.5. synthesizing a combined response based on individual model outputs; and
   - 5.6. returning the synthesized response and individual model responses.

## System Architecture Claims

6. A multi-model orchestration system comprising:
   - 6.1. a model abstraction layer providing a unified interface to multiple LLM providers;
   - 6.2. a pattern registry storing configurable analysis workflows;
   - 6.3. a context management system maintaining state across orchestration stages; and
   - 6.4. a response synthesis engine for integrating multiple model outputs.
7. The claim 6 system comprises a caching subsystem that stores and retrieves responses based on input similarity and configurable time-to-live parameters.
8. The system of claim 6, wherein the context management system supports document chunking and relevance-based retrieval for incorporating document knowledge into LLM prompts.

## Method Claims

9. A method for pattern-based orchestration of language models, comprising:
   - 9.1. receiving a pattern selection, springing an analytical workflow;
   - 9.2. loading the pattern configuration, including stage definitions and participant rules;
   - 9.3. executing each stage according to the pattern configuration;
   - 9.4. collecting intermediate results between stages; and
   - 9.5. synthesizing a final response based on the collected results.
10. The method of claim 9 further comprises dynamically selecting participant models for each stage based on capability requirements and performance metrics.

11.    The method of claim 9, wherein executing stages includes both parallel and sequential model execution based on stage configuration.


Quality Evaluation Claims

12.    A method for multi-dimensional quality evaluation of language model responses, comprising:
    12.1.    receiving a response from a language model;
    12.2.    evaluating the response across multiple quality dimensions;
    12.3.    calculating dimension-specific scores using configurable evaluation algorithms;
    12.4.    computing an aggregate quality score;
    12.5.    storing the quality metrics in a performance database; and
    12.6.    using the quality metrics to inform future model selection.
13.    The method of claim 12, wherein the quality dimensions include at least coherence, technical depth, strategic value, and uniqueness.
14.    The method of claim 12, wherein evaluating includes using another language model to perform meta-evaluation of responses.


Dynamic Registry Claims

15.    A system for runtime management of language model providers, comprising:
    15.1.    a model registry interface for adding and removing model providers;
    15.2.    a capability discovery mechanism for determining model capabilities;
    15.3.    an adapter factory creating provider-specific adapters implementing a unified interface;
    15.4.    a configuration validator ensuring model configurations meet system requirements; and
    15.5.    a registry event system notifying dependent components of registry changes.
16.    The system of claim 15 is further comprised of a health monitoring subsystem that tracks model availability and performance.


Document Processing Claims

17.    A method for integrating document context in multi-model language processing, comprising:
    17.1.    receiving document uploads;
    17.2.    extracting text content from various document formats;
    17.3.    segmenting content into semantically meaningful chunks;
    17.4.    generating vector embeddings for each chunk;
    17.5.    selecting relevant chunks based on query similarity;
    17.6.    integrating selected chunks into LLM context windows; and
    17.7.    tracking chunk usage for attribution and explanation.
18.    The method of claim 17 further comprises a token budget management system that optimizes chunk selection based on model-specific context window limitations.


Business Logic Claims

19.    A system for usage-based pricing of multi-model language processing, comprising:
    19.1.    a token tracking mechanism monitoring consumption across models;
    19.2.    a tiered pricing engine with configurable pricing rules;
    19.3.    a cost estimation service providing pre-execution estimates; and
    19.4.    a billing integration system for usage reporting.
20.    The system of claim 19, wherein the pricing engine supports model-specific markups, volume discounts, and specialized pattern pricing.


User Interface Claims

21.    A method for visualizing comparative language model analysis, comprising:
    21.1.    displaying multiple model responses in a configurable side-by-side view;
    21.2.    highlighting similarities and differences between responses;

21.3.  visualizing confidence and quality metrics using interactive graphical elements;
21.4.  providing an integrated synthesis view that references contributing models; and
21.5.  enabling user feedback on individual responses and synthesis quality.

22.  The method of claim 21, further comprising an explanation view that illustrates how the system arrived at its synthesis.


## Error Handling and Reliability Claims

23.  A method for ensuring reliable operation in a multi-model orchestration system, comprising:
23.1.  implementing model-specific circuit breakers that prevent cascading failures;
23.2.  providing configurable fallback strategies when models are unavailable;
23.3.  implementing timeout management tailored to model-specific response patterns;
23.4.  logging detailed diagnostics for failure analysis; and
23.5.  automatically retrying failed operations according to configurable policies.

24.  The method of claim 23 further comprises an operation monitoring system that detects performance degradation and initiates mitigating actions.


## Pattern Extension Claims

25.  A method for extending analysis capabilities in a language model orchestration system, comprising:
25.1.  defining a pattern specification schema for analytical workflows;
25.2.  implementing a pattern registration mechanism;
25.3.  providing a template system for stage-specific prompt generation;
25.4.  enabling pattern composition through an inheritance mechanism; and
25.5.  supporting runtime pattern parameter configuration.

26.  The method of claim 25, wherein patterns can define custom quality evaluation criteria specific to the analytical workflow.