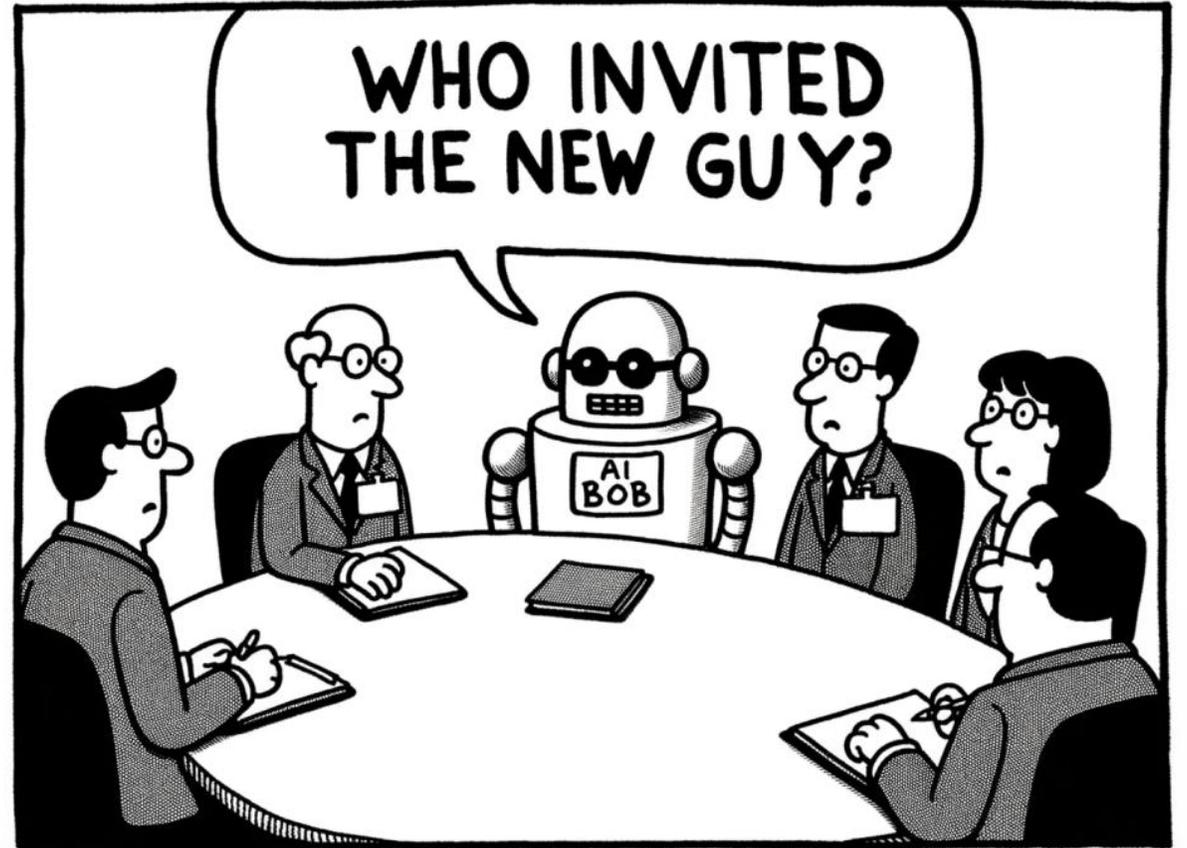# Operationalizing AI: Infrastructure to Agents
## Unlocking Business Value at Scale

**2026**

**Amit Nandi**

# AgentOps

# Investment Banking
## Large Quantitative Models (LQMs)

```
+--------------+
|    LQMs      |
| Quantitative |
| Models       |
+--------------+
```

# Investment Banking

Large Quantitative Models (LQMs) + Small Language Models (SLMs)

```
+--------------+        +----------------+
|    LQMs      |        |     SLMs       |
| Quantitative|        | Small Lang.    |
| Models       |        | Models         |
+--------------+        +----------------+
```

# Investment Banking
## Large Quantitative Models (LQMs) + Small Language Models (SLMs) +RAGs

```
+--------------+      +----------------+      +-----------------+      +------------------+
|    LQMs      |      |     SLMs       |      |   Tools/Apps    |      |    Data APIs     |
| Quantitative |      | Small Lang.    |      | (RAG, Search,   |      | Market, Risk, CRM|
| Models       |      | Models         |      |  Analysis)      |      | Pricing, KYC     |
+--------------+      +----------------+      +-----------------+      +------------------+
```

# Investment Banking
Large Quantitative Models (LQMs) + Small Language Models (SMLs) + RAGs

```
+-------------+      +-----------------+     +------------------+    +------------------+
| LQMs        |      | SLMs            |     | Tools/Apps       |    | Data APIs        |
| Quantitative|      | Small Lang.     |     | (RAG, Search,    |    | Market, Risk, CRM|
| Models      |      | Models          |     | Analysis)        |    | Pricing, KYC     |
+-------------+      +-----------------+     +------------------+    +------------------+

    ↳ Risk, Pricing,    ↳ Legal, Policy,       ↳ WebSearch,          ↳ Real-time + Batch
      Credit, Market      Compliance, Domain     RAG, BI Dash          Pipelines, Streaming

                        +-----------------------------+
                        | PROMPT LAYER                |
                        |-----------------------------|
                        | • Prompt Version Control    |
                        | • Prompt Templates & Macros |
                        | • Guardrails & Overrides    |
                        +-----------------------------+
```

# Investment Banking
## Large Quantitative Models (LQMs) + Small Language Models (SMLs) + RAGs + Agents

```
                    +----------------+----------------+
                    |      AGENT LAYER                |
                    |  (Agent Ops, Orchestration)     |
                    |---------------------------------|
                    | • Agent Registry & Routing      |
                    | • Memory (Vector DB, Event Log) |
                    | • Tool Use + Reasoning Engine   |
                    | • Multi-Agent Collaboration     |
                    | • A2A (Agent-to-Agent) Protocol |
                    +----------------+----------------+
                                     |
        +----------------------------+--------+---------------+----------------------+
        |                            |                        |                      |
        ▼                            ▼                        ▼                      ▼
  +-------------+            +-------------+          +----------------+      +-----------------+
  |   LQMs      |            |    SLMs     |          |   Tools/Apps   |      |    Data APIs    |
  | Quantitative|            | Small Lang. |          | (RAG, Search,  |      | Market, Risk, CRM|
  | Models      |            | Models      |          |   Analysis)    |      | Pricing, KYC    |
  +-------------+            +-------------+          +----------------+      +-----------------+

  ↳ Risk, Pricing,          ↳ Legal, Policy,         ↳ WebSearch,           ↳ Real-time + Batch
    Credit, Market            Compliance, Domain        RAG, BI Dash          Pipelines, Streaming

                    +-----------------------------+
                    |       PROMPT LAYER          |
                    |-----------------------------|
                    | • Prompt Version Control    |
                    | • Prompt Templates & Macros |
                    | • Guardrails & Overrides    |
                    +-----------------------------+
```

# Investment Banking
## Large Quantitative Models (LQMs) + Small Language Models (SMLs) + RAGs + Agents

```
+------------------------------------------------+
|                USER INTERFACE                  |
|      Business, Ops, Analysts via AGUI/LLM Chat |
+-----------------------▲------------------------+
                        |
              +---------┴---------+
              |   LLM Interface   |     ← GPT, Claude, Gemini (UX Layer)
              |  Natural Language |
              +---------┬---------+
                        |
          +-------------┴----------------+
          |         AGENT LAYER          |
          |   (Agent Ops, Orchestration) |
          |------------------------------|
          | • Agent Registry & Routing   |
          | • Memory (Vector DB, Event Log)|
          | • Tool Use + Reasoning Engine |
          | • Multi-Agent Collaboration  |
          | • A2A (Agent-to-Agent) Protocol|
          +-------------┬----------------+
                        |
   +-------------+------┴------+-------------+------------------+
   |             |             |             |                  |
   ▼             ▼             ▼             ▼
+----------+  +----------+  +-------------+  +----------------+
|   LQMs   |  |   SLMs   |  |  Tools/Apps |  |    Data APIs   |
|Quantitative| | Small Lang.| | (RAG, Search,| | Market, Risk, CRM|
|Models    |  | Models   |  |  Analysis)  |  | Pricing, KYC   |
+----------+  +----------+  +-------------+  +----------------+

↳ Risk, Pricing,   ↳ Legal, Policy,   ↳ WebSearch,      ↳ Real-time + Batch
  Credit, Market     Compliance, Domain  RAG, BI Dash      Pipelines, Streaming

              +------------------------------+
              |         PROMPT LAYER         |
              |------------------------------|
              | • Prompt Version Control     |
              | • Prompt Templates & Macros  |
              | • Guardrails & Overrides     |
              +------------------------------+
```

# Investment Banking
## LQMs + SLMs + RAGs Orchestrated by Agents

```
+------------------------------------------------+
|                 USER INTERFACE                 |
|      Business, Ops, Analysts via AGUI/LLM Chat |
+------------------------------------------------+
                        ▲
                        |
              +---------+---------+
              |   LLM Interface   |        ← GPT, Claude, Gemini (UX Layer)
              |  Natural Language |
              +---------+---------+
                        |
              +---------+-------------------+
              |       AGENT LAYER           |
              |  (Agent Ops, Orchestration) |
              |-----------------------------|
              | • Agent Registry & Routing  |
              | • Memory (Vector DB, Event Log)|
              | • Tool Use + Reasoning Engine |
              | • Multi-Agent Collaboration |
              | • A2A (Agent-to-Agent) Protocol|
              +-------------+---------------+
                            |
    +----------------+------+--------+-------------------+
    |                |               |                   |
    ▼                ▼               ▼                   ▼
+-----------+  +-------------+  +-------------+  +----------------+
|   LQMs    |  |    SLMs     |  |  Tools/Apps |  |   Data APIs    |
| Quantitative| | Small Lang. |  | (RAG, Search,|  | Market, Risk, CRM|
| Models    |  | Models      |  |  Analysis)  |  | Pricing, KYC   |
+-----------+  +-------------+  +-------------+  +----------------+

↳ Risk, Pricing,  ↳ Legal, Policy,   ↳ WebSearch,   ↳ Real-time + Batch
  Credit, Market    Compliance, Domain  RAG, BI Dash    Pipelines, Streaming

              +-----------------------------+
              |        PROMPT LAYER         |
              |-----------------------------|
              | • Prompt Version Control    |
              | • Prompt Templates & Macros |
              | • Guardrails & Overrides    |
              +-----------------------------+

+------------------------------------------------+
|          INFRASTRUCTURE & OBSERVABILITY        |
|------------------------------------------------|
| • Compute: GPU/K8s/Serverless | • Data Mesh + Lakehouse |
| • CI/CD Pipelines & Telemetry | • Security & Governance |
| • LLM Gateway + Cost Controls | • Audit Logging + Monitoring|
+------------------------------------------------+
```

# Investment Banking
## Activating the Stack

```
+----------------- Business / Trader / Risk Officer --------------+
|             Natural Language Query / Instruction                |
+-------------------------- LLM Interface ------------------------+
| Broad Language Understanding | Context | Prompt Versioning      |
+------------------- Agent Orchestration Layer ------------------+
|   Memory | Workflow | Risk Guards | Governance | Tool Use       |
|              Routes requests intelligently                      |
+---------------- Models & Tools --------------------------------+
| LQMs (Pricing, VaR, Stress) | SLMs (Ops, AML, Enrichment)       |
+------------------------- Data Plane ---------------------------+
| RAG | Retrieval Fusion | Feature Store | Vector DB              |
+---------------------- Infrastructure Layer --------------------+
| GPUs | CI/CD | Data Lakes |  Observability | Security           |
+----------------------------------------------------------------+
```

# Scaling AI Operations: MLOps -> LLMOps -> AgentOps

Ops Framework



USER EXPERIENCE
HUMAN-AGENT INTERFACE

AGENT ORCHESTRATION

MODEL

CODE

DATA

INFRASTRUCTURE

# Scaling AI Operations: MLOps -> LLMOps -> AgentOps

Ops Framework

| | The 6 Dimensions |
|---|---|
| | **MLOps** |
| **USER EXPERIENCE HUMAN-AGENT INTERFACE** | - Analysts & Traders Dashboards<br>- Inferencing |
| **AGENT ORCHESTRATION** | - N/A<br>- Manual orchestration |
| **MODEL** | - ML Models<br>- Model Registry<br>- Model Decay/Training |
| **CODE** | - ETL Pipelines - CI/CD (Automation of ML Pipeline)<br>- Orchestration (Apache Airflow, Temporal.io, Kubernetes)<br>- API Services  - API Gateway (Model, Agent Integration) |
| **DATA** | - Data Lakehouse (Structured/Unstructured Data)<br>- Feature Store (Real-time & Batch Features)<br>- Data Pipelines (ETL, Preprocessing) |
| **INFRASTRUCTURE** | - CPU/GPU Resources (DEV/TEST/PRE-PROD/PROD)<br>- Cloud/On-prem Infrastructure (Elastic Compute, Kubernetes, GPU Provisioning) |

# Scaling AI Operations: MLOps -> LLMOps -> AgentOps
## Ops Framework

| The 6 Dimensions of AI/ML Operations | | |
| --- | --- | --- |
| | **MLOps** | **LLMOps** |
| **USER EXPERIENCE HUMAN-AGENT INTERFACE** | - Analysts & Traders Dashboards<br>- Inferencing | - LLM queries & drafting support<br>- Prompt Library<br>- RAGs |
| **AGENT ORCHESTRATION** | - N/A<br>- Manual orchestration | - N/A<br>- Basic workflows fallback logic |
| **MODEL** | - ML Models<br>- Model Registry<br>- Model Decay/Training | - LLMs, SLMs, Embeddings<br>- Foundational Model Gardens LLMs |
| **CODE** | - ETL Pipelines - CI/CD (Automation of ML Pipeline)<br>- Orchestration (Apache Airflow, Temporal.io, Kubernetes)<br>- API Services  - API Gateway (Model, Agent Integration) | - Prompt pipelines,<br>- APIs |
| **DATA** | - Data Lakehouse (Structured/Unstructured Data)<br>- Feature Store (Real-time & Batch Features)<br>- Data Pipelines (ETL, Preprocessing) | - Text corpora, vectors<br>- VectorDB (Embeddings, RAG) |
| **INFRASTRUCTURE** | - CPU/GPU Resources (DEV/TEST/PRE-PROD/PROD)<br>- Cloud/On-prem Infrastructure (Elastic Compute, Kubernetes, GPU Provisioning) | - High-memory GPU clusters<br>- Multi-GPU support for NLP & LLM |

# Scaling AI Operations: MLOps -> LLMOps -> AgentOps
## Ops Framework

| The 6 Dimensions of AI/ML Operations | | | |
|---|---|---|---|
| | **MLOps** | **LLMOps** | **AgentOps** |
| **USER EXPERIENCE HUMAN-AGENT INTERFACE** | - Analysts & Traders Dashboards<br>- Inferencing | - LLM queries & drafting support<br>- Prompt Library<br>- RAGs | - Full HITL integration, insights, overrides<br>- AGUI, HITL, Agent-to-Agent |
| **AGENT ORCHESTRATION** | - N/A<br>- Manual orchestration | - N/A<br>- Basic workflows fallback logic | - Orchestration<br>- Memory<br>- Reasoning<br>- Multi-agent coordination |
| **MODEL** | - ML Models<br>- Model Registry<br>- Model Decay/Training | - LLMs, SLMs, Embeddings<br>- Foundational Model Gardens LLMs | - LQMs, SLMs, LLMs coordinated |
| **CODE** | - ETL Pipelines - CI/CD (Automation of ML Pipeline)<br>- Orchestration (Apache Airflow, Temporal.io, Kubernetes)<br>- API Services  - API Gateway (Model, Agent Integration) | - Prompt pipelines,<br>- APIs | - Agent workflows (A2A), tool adapters (MCP) |
| **DATA** | - Data Lakehouse (Structured/Unstructured Data)<br>- Feature Store (Real-time & Batch Features)<br>- Data Pipelines (ETL, Preprocessing) | - Text corpora, vectors<br>- VectorDB (Embeddings, RAG) | - Live streams, knowledge stores, RAG |
| **INFRASTRUCTURE** | - CPU/GPU Resources (DEV/TEST/PRE-PROD/PROD)<br>- Cloud/On-prem Infrastructure (Elastic Compute, Kubernetes, GPU Provisioning) | - High-memory GPU clusters<br>- Multi-GPU support for NLP & LLM | - Multi-GPU, orchestration layer, persistent memory |

# Scaling AI Operations: MLOps -> LLMOps -> AgentOps
## Ops Framework

### The 6 Dimensions of AI/ML Operations

| | KEY FEATURES | OBSERVABILITY | GOVERNANCE |
|---|---|---|---|
| **USER EXPERIENCE HUMAN-AGENT INTERFACE** | - UX Dashboards (Real-time Metrics, Performance)<br>- Human-in-the-loop (HITL) Interfaces<br>- Agent-to-Agent (A2A) Communication<br>- Agent-to-Human (A2H) Interaction (Chatbots, Assistants) | - User Behavior Analytics<br>- Model Feedback Tracking<br>- UI/UX Monitoring | - Ethical Usage Guidelines<br>- Privacy Management for Users<br>- Human-in-the-loop Audit Trails |
| **AGENT ORCHESTRATION** | - Orchestrates LQMs + SLMs workflows<br>- Multi-Agent Coordination (Task Delegation, Temporal Workflows)<br>- Memory management, reasoning, multi-agent collaboration<br>- Fallback logic, Implements policies<br>- Risk checks, and escalation thresholds<br>- MCP / A2A / AGUI protocols | - Agent Behavior Analytics<br>- Agent Reliabilty Tracking<br>- A2A - Agent 2 Human Monitoring<br>- Tool Calling MCP<br>- Token consumption | - Agent Risk and Trust Guardrails<br>- Explainability Audit Trails<br>- Accountability Audit Trails<br>- Kill Switch Activation |
| **MODEL** | - Model Registry (Versioning, Metadata)<br>- Foundation Model Garden (Pre-trained models)<br>- A/B Testing (Model Evaluation)<br>- Continuous Model Training<br>- Model Inferencing | - Model Performance Metrics<br>- Latency Monitoring<br>- Model Stability | - Responsible AI (Ethics, Transparency<br>- Explainability (Black-box Models)<br>- Accountability for Model Behavior |
| **CODE** | - Infrastructure as Code (IaC)<br>- ETL Pipelines - CI/CD (Automation of ML Pipeline)<br>- Orchestration (Apache Airflow, Temporal.io, Kubernetes)<br><br>- API Services - API Gateway (Model, Agent Integration) | - Code Quality Monitoring<br>- CI/CD Pipeline Monitoring<br>- Error Tracking | - Ethical Usage Guidelines<br>- Privacy Management for Users<br>- Human-in-the-loop Audit Trails |
| **DATA** | - Data Lakehouse (Structured/Unstructured Data)<br>- Feature Store (Real-time & Batch Features)<br>- VectorDB (Embeddings, RAG)<br>- Data Pipelines (ETL, Preprocessing) | - Data Quality Monitoring<br>- Drift Detection<br>- Anomaly Detection | - Data Privacy and Security<br>- Model Fairness & Bias Checking<br>- Access Control (Data Permissions) |
| **INFRASTRUCTURE** | - CPU/GPU Resources (DEV/TEST/PRE-PROD/PROD)<br>- Cloud/On-prem Infrastructure (Elastic Compute, Kubernetes, GPU Provisioning) | - Performance Monitoring<br>- Cost Tracking<br>- Security (Access, Compliance)<br>- SLA Management | - Data Lineage<br>- Model Lineage<br>- Responsible AI (Fairness, Bias)<br>- Guardrails (Model and Data Validation) |

# ML/AI People
## Operations

**Platform Team**
**Secure Cloud/Data/ML Platform**

**Business**
**Viz Dashboards, ML Adoption, & ROI**

**Advance Analytics Team**
**Data Lake**

**Data Science Team**
**Experimentation & MLOps**

### Data Engineer
Prepare & Ingest data building ETL pipelines

### Data Scientist
Create the best ML models to solve business problems

### MLOps Engineer/Admin
Standardize CI/CD, user/service role, model consumption, testing and deployment methodology

### Business Stakeholder Product Owners
Define business problem, business KPIs, and make business decisions

### Data Owners
Manage data sharing and provide access

### ML Engineer
Collaborate with DS to productionize ML

### Security
Assess data, user, and service access creating policies and guardrails

### Business Stakeholder Data & ML Consumers
Consumers of ML results from other BUs, driving business decision making

**Risk & Compliance**
**Approve & Review Models**

### Architects/ SysOps Engineer
Standardize account infrastructure, connectivity, user roles implementation

### Auditors/Risk & Compliance
Review models, data sources, code artifacts

# ML/AI + LLMs + Agentic AI People
## Operations

**Advance Analytics Team**
**Data Lake**

### Data Engineer
Prepare & Ingest data building ETL pipelines

### Data Owners
Manage data sharing and provide access

**Labeler Team**
**Data Preparation at Scale**

### Data Labelers/Editors
Label or edit billions of Data for FM models and hundreds of data for fine tuning interacting with data lake using a dedicated website

---

**Data Science Team**
**Experimentation & MLOps**

### Data Scientist
Create the best ML models to solve business problems

### ML Engineer
Collaborate with DS to productionize ML

**Data Science Team Extension**
**Context Adaptation**

### Fine Tuners
Select the corresponding FM, evaluate the model & design the deployment method/infrastructure

---

**Platform Team**
**Secure Cloud/Data/ML Platform**

### MLOps Engineer/Admin
Standardize CI/CD, user/service role, model consumption, testing and deployment methodology

### Security & Architects
Assess data, user, and service access creating policies and infrastructure

**Risk & Compliance**
**Approve & Review Models**

### Auditors/Risk & Compliance
Review models, data sources, code artifacts

---

**Business**
**Viz Dashboards, ML Adoption, & ROI**

### Business Stakeholder Product Owners
Define business problem, business KPIs, and make business decisions

### Business Stakeholder Data & ML Consumers
Consumers of ML results from other BUs, driving business decision making

**End-Users**
**Consume Generative AI applications**

### Generative AI End-users
Consume Generative AI solutions as black box, share data and rate the quality of output

---

**GenAI Application Team**
**Integrate GenAI models in applications**

### Generative AI Developers
Select, test, evaluate the FM, filter inputs/outputs, and develop the generative AI application back-end (e.g. LangChain Experts)

### AppDev/DevOps
Develop the front-end of the GenAI application

### Prompt Engineers
Design the input/output prompts to adapt the solution to the context and test the initial version
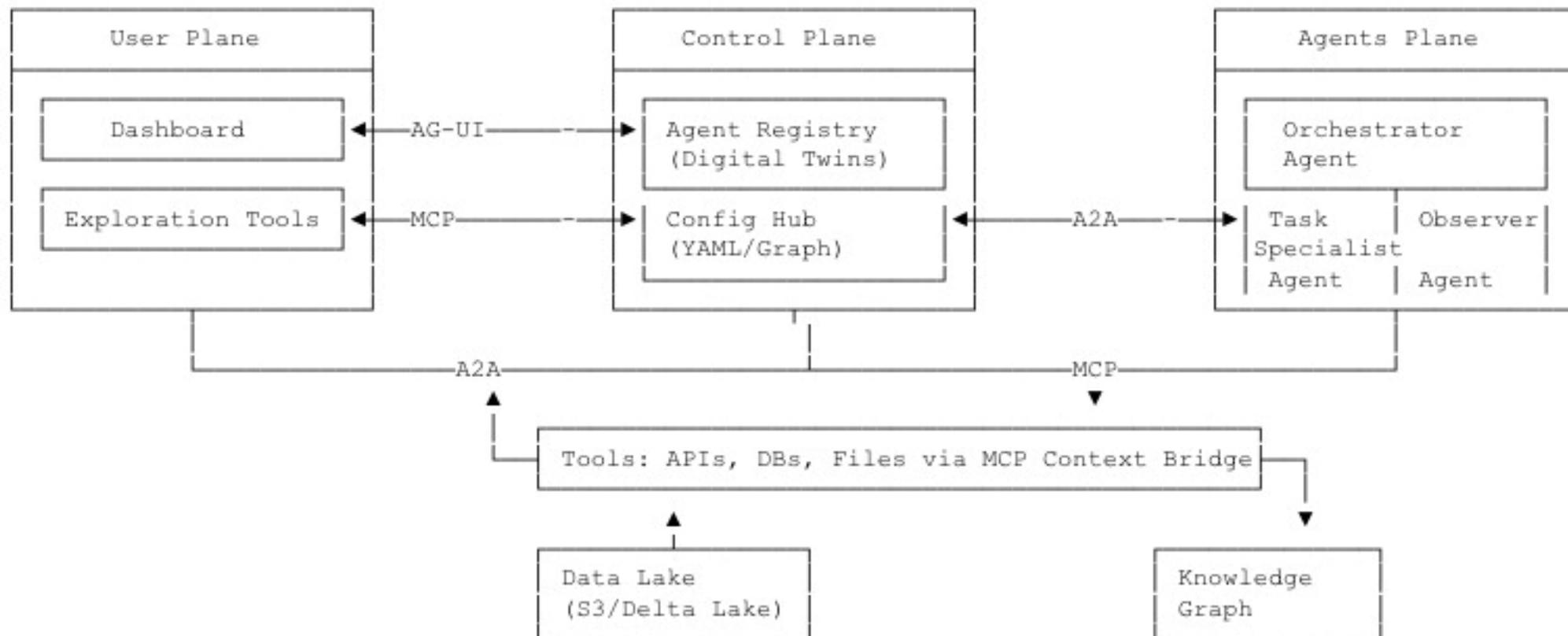
### Prompt Testers
Test at scale the generative. AI solution (back-end/ front-end) and feed their results to the prompt catalog
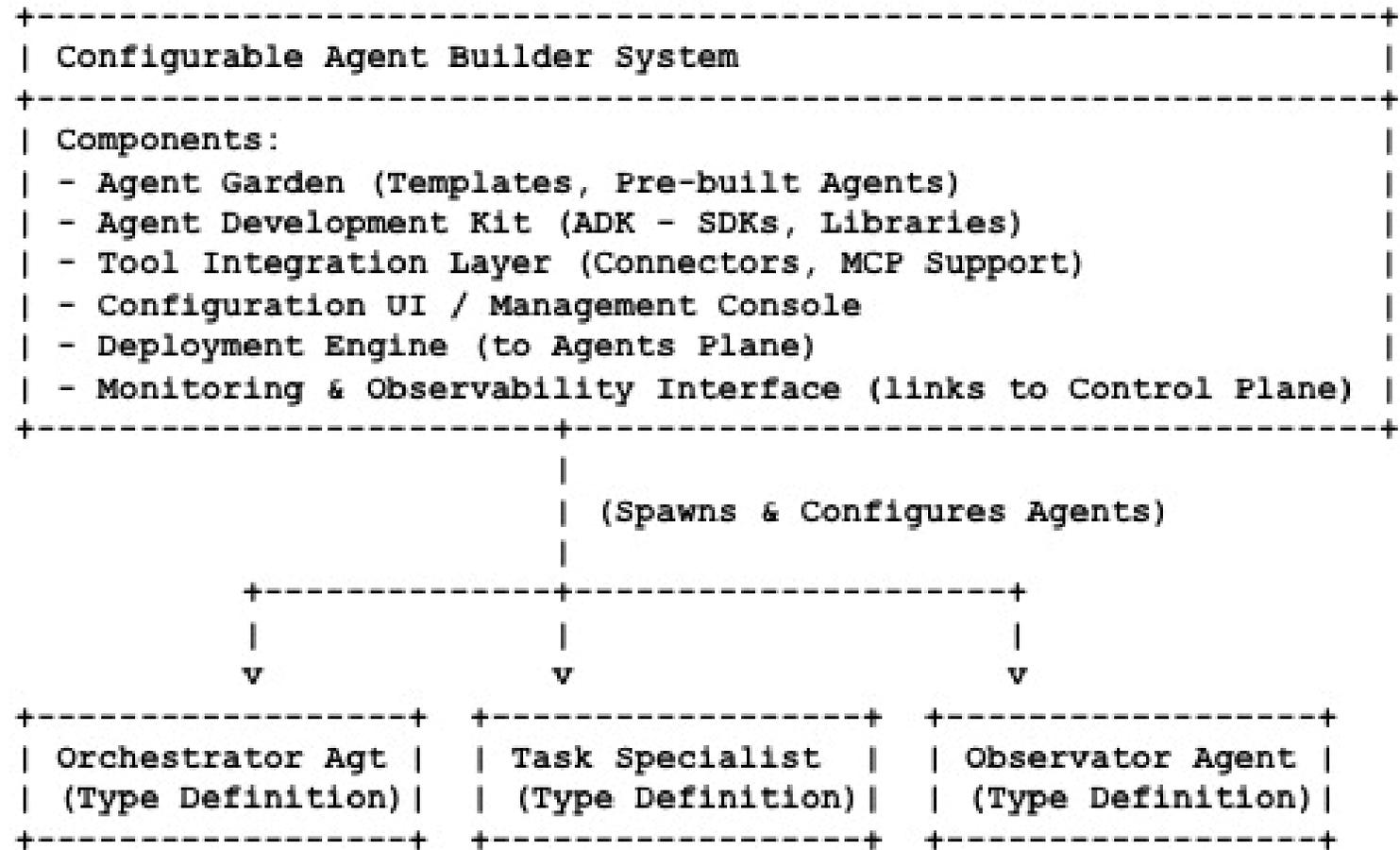
# Agentic AI
## User Plane - Control Plane - Agents Plane

# Agentic AI
## Conceive - Craft - Customize - Control

```
+------------------------------------------------------------------------+
| Configurable Agent Builder System                                      |
+------------------------------------------------------------------------+
| Components:                                                            |
| - Agent Garden (Templates, Pre-built Agents)                           |
| - Agent Development Kit (ADK - SDKs, Libraries)                         |
| - Tool Integration Layer (Connectors, MCP Support)                     |
| - Configuration UI / Management Console                                |
| - Deployment Engine (to Agents Plane)                                  |
| - Monitoring & Observability Interface (links to Control Plane)        |
+-------------------------------+----------------------------------------+
                                |
                                | (Spawns & Configures Agents)
                                |
            +-------------------+--------------------+
            |                   |                    |
            v                   v                    v
+---------------------+ +--------------------+ +--------------------+
| Orchestrator Agt    | | Task Specialist    | | Observator Agent   |
| (Type Definition)   | | (Type Definition)  | | (Type Definition)  |
+---------------------+ +--------------------+ +--------------------+
```

# Agentic AI
## Agent Factory



Agent Factory

▶ Blueprint Catalog
    - Orchestrator Template
    - Task Agent Template
▶ Skill Composer
    - MCP Tool Binding
    - A2A Workflow Designer

Validation
▶ Security Scanning
▶ Protocol Compliance

Runtime Controller

▶ Container Orchestration
▶ MCP Context Injection
▶ A2A Network Policies

Agent Instance

| MCP Bridge | A2A Comms | Memory Mgmt | Tools |
|------------|-----------|-------------|-------|

# Agentic AI
## Orchestrator Agents + Task Specialist Agents

| Orchestrator Agent | | |
|---|---|---|
| Input | Memory | Tools/Actions |
| ▶ User Prompts<br>▶ A2A Messages<br>▶ MCP Context<br>▶ File Uploads | STM:<br>  – 8k Token<br>LTM:<br>  – Vector DB<br>Context:<br>  – Session | ▶ A2A Broker<br>▶ MCP Gateway<br>▶ API Gateway |
| Reasoning/Planning | | |
| ▶ LLM (Mixtral 8x22B)<br>▶ Chain-of-Thought<br>▶ Tree-of-Thought<br>▶ Constraint Solver | ▲<br>└─A2A Feedback Loop─<br>▼ | |

```
Task 1
Task 2
```

| Task Specialist Agent | | |
|---|---|---|
| Input | Memory | Tools/Actions |
| ▶ A2A Tasks<br>▶ MCP Context<br>▶ Sensor Data | STM:<br>  – 4k Token<br>LTM:<br>  – Graph DB<br>Context:<br>  – Workflow | ▶ Domain APIs<br>▶ MCP Tools<br>▶ A2A Proxies |
| Reasoning/Planning | | |
| ▶ Fine-Tuned LLM<br>▶ Rule Engine<br>▶ Pattern Matching<br>▶ Optimization | ▲<br>└─MCP Context Sync─<br>▼ | |

```
Action 1
Action 2
```

# Agentic AI
## Observer Agents



```
┌─────────────────────────────────────────────────────────┐
│                    Observer Agent                        │
├───────────────────┬─────────────────┬───────────────────┤
│ Input             │ Memory          │ Tools/Actions     │
├───────────────────┼─────────────────┼───────────────────┤
│ ▶ Log Streams     │ STM:            │ ▶ A2A Override    │
│ ▶ A2A Monitors    │   - 2k Token    │ ▶ MCP Hooks       │
│ ▶ MCP Probes      │ LTM:            │ ▶ Killswitch      │
│                   │   - Audit DB    │                   │
│                   │ Context:        │                   │
│                   │   - Policies    │                   │
├───────────────────┴─────────────────┴───────────────────┤
│         Reasoning/Planning                               │
├──────────────────────────────────────────────────────────┤
│ ▶ Policy LLM                ▲                            │
│ ▶ Anomaly Detection          └─Real-time Alerting───     │
│ ▶ Multi-LLM Consensus       ▼                           │
│ ▶ Compliance Checker                                     │
└────────────────────┬────────────────┬──────────────────┘
                     │ Alert 1        │
                     │ Action 2       │
                     └────────────────┘
```

# Agentic AI - From Infrastructure to Agents
## Scaling for Business Value

```
+------------------------------- UX / Human Interface -----------------------------+
| Users interact via LLM interface; insights, explanations, alerts presented |
| Supports human-in-the-loop (HITL) and oversight on agent-driven actions |
+------------------------------ Agent Orchestration Layer --------------------------+
| Orchestrates LQMs + SLMs workflows |
| Memory management, reasoning, multi-agent collaboration, fallback logic |
| Implements policies, risk checks, and escalation thresholds |
+------------------------------------ Model Layer ---------------------------------+
| LQMs: Quantitative models for pricing, portfolio, risk calculations |
| SLMs: Task-specific language models (summarization, reporting) |
| LLMs: Interface layer mediating user interaction and natural language querying |
+--------------------------------- Code / Orchestration ---------------------------+
| Agent workflows: scheduling, retries, branching, tool adapters |
| CI/CD pipelines for LQMs, SLMs, prompts, and agent rules |
+------------------------------------ Data Layer ----------------------------------+
| Structured market data, portfolio data, risk factors, document corpora |
| Embeddings, RAG pipelines, feature stores |
| Real-time streams integrated with agent memory and orchestration |
+------------------------------- Infrastructure Layer ------------------------------+
| Multi-GPU clusters for LQMs & LLM inference |
| High-availability, low-latency compute for trading operations |
| Observability, logging, and secure environment |
+----------------------------------------------------------------------------------+
```

AI

# Thank You

AI powered Business

# AI/ML Ops
## Framework

### USER EXPERIENCE / HUMAN-AGENT INTERFACE

| Key Features | Observability | Governance |
|---|---|---|
| - UX Dashboards (Real-time Metrics, Performance)<br>- Human-in-the-loop (HITL) Interfaces<br>- Agent-to-Agent (A2A) Communication<br>- Agent-to-Human (A2H) Interaction (Chatbots, Assistants) | - User Behavior Analytics<br>- Model Feedback Tracking<br>- UI/UX Monitoring | - Ethical Usage Guidelines<br>- Privacy Management for Users<br>- Human-in-the-loop Audit Trails |

### AGENT ORCHESTRATION LAYER

| Key Features | Observability | Governance |
|---|---|---|
| - Orchestrates LQMs + SLMs workflows<br>- Memory management, reasoning, multi-agent collaboration<br>- Fallback logic, Implements policies<br>- Risk checks, and escalation thresholds<br>- MCP / A2A / AGUI protocols | - Agent Behavior Analytics<br>- Agent Reliabilty Tracking<br>- A2A - Agent 2 Human Monitoring<br>- Tool Calling MCP<br>- Token consumption | - Agent Risk and Trust Guardrails<br>- Explainability Audit Trails<br>- Accountability Audit Trails<br>- Kill Switch Activation |

### MODEL LAYER

| Key Features | Observability | Governance |
|---|---|---|
| - Model Registry (Versioning, Metadata)<br>- Foundation Model Garden (Pre-trained models)<br>- A/B Testing (Model Evaluation)<br>- Multi-Agent Coordination (Task Delegation, Temporal Workflows) | - Model Performance Metrics<br>- Latency Monitoring<br>- Model Stability | - Responsible AI (Ethics, Transparency<br>- Explainability (Black-box Models)<br>- Accountability for Model Behavior |

### CODE LAYER

| Key Features | Observability | Governance |
|---|---|---|
| - Infrastructure as Code (IaC)<br>- ETL Pipelines - CI/CD (Automation of ML Pipeline)<br>- Orchestration (Apache Airflow, Temporal.io, Kubernetes)<br>- API Services - API Gateway (Model, Agent Integration) | - Code Quality Monitoring<br>- CI/CD Pipeline Monitoring<br>- Error Tracking | - Ethical Usage Guidelines<br>- Privacy Management for Users<br>- Human-in-the-loop Audit Trails |

### DATA LAYER

| Key Features | Observability | Governance |
|---|---|---|
| - Data Lakehouse (Structured/Unstructured Data)<br>- Feature Store (Real-time & Batch Features)<br>- VectorDB (Embeddings, RAG)<br>- Data Pipelines (ETL, Preprocessing) | - Data Quality Monitoring<br>- Drift Detection<br>- Anomaly Detection | - Data Privacy and Security<br>- Model Fairness & Bias Checking<br>- Access Control (Data Permissions) |

### INFRASTRUCTURE LAYER

| Key Features | Observability | Governance |
|---|---|---|
| - CPU/GPU Resources (DEV/TEST/PRE-PROD/PROD)<br>- Cloud/On-prem Infrastructure (Elastic Compute, Kubernetes, GPU Provisioning) | - Performance Monitoring<br>- Cost Tracking<br>- Security (Access, Compliance)<br>- SLA Management | - Data Lineage<br>- Model Lineage<br>- Responsible AI (Fairness, Bias)<br>- Guardrails (Model and Data Validation) |

# MLOps - Road to Production

**ML Use Case**
Business Impact
Data Sources
Transformation
Lineage
ML Algorithms
Scenarios
….

**Artefacts**
Business Case
Solutions Architecture
Data Lineage and Transformation
ML Algos
Test Cases

## Model Governance

**mlflow**

| Tracking | Projects | Models |
|---|---|---|
| Experiments run logs metadata | packaging | registry |

## DEVOPS - CI/CD

| Git | Pipelines | Endpoint |
|---|---|---|
| Code Repos | CI/CD | Deploy |

**Artefacts**
Configs
Logging
Script
Code
Feature extraction ML pipeline
Package dependencies
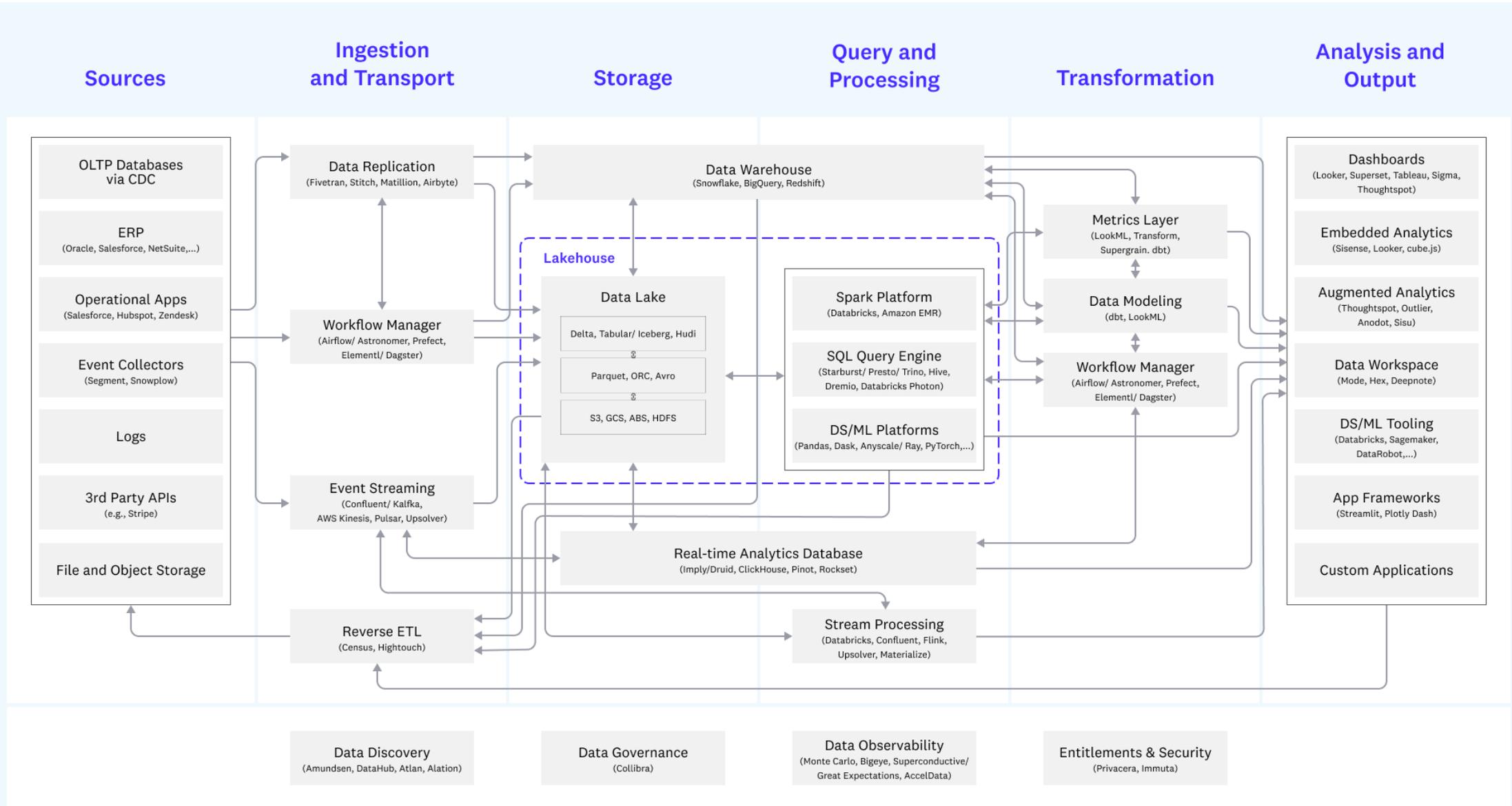Model version serialization metadata

## Observability - Monitoring & Audit

| Run / Logs | Model Data Drift | Conf / Logs |
|---|---|---|
| **FinOps** Cost | **SecOps** Data/EndPt | **Audit** RespAI |

### DEV/Sandbox
Exploration

**ML Feature Engineering ML Model**
Feature Engineering
Train Model
Test Model
Hyper Parameters

**Test**
Unit Test
Integration Test

**Dev Data**
Sample Data
Schemas + Metadata
Data Profiling + Joins
Partitions
Train / Test Datasets
**Add. Exploratory Data**

### DEV/VALIDATION
Scaling & Tuning Model Training

**ML Feature Engineering ML Model**
Feature Engineering
Train Model (Valid Data)
Test Model (Valid Data
Hyper Parameters

**Valid Model to train**
**Serialized Valid Model**

**Feature Store**
**Metadata – Lineage - Feature**

**Test**
Integration Test
Stress Test
Performance Tuning

**Validation/Test Data**
Historical Data
Data partition
Checkpointing
Failover/Rollback

### PRE-PROD/STAGING
Continuous Integration

**ML Feature Engineering ML Model**
Feature Engineering
Train Model (PreProd)
Test Model (PreProd Data
Hyper Parameters

**Prod Model Train**
**Prod Model Predict**

**Feature Store**
**Metadata – Lineage - Feature**

**Test**
Integration Test
Stress Test
Performance Tuning

**Pre-Prod Data**
Stress Test Data
Train / Test Datasets

### PROD
Continuous Training & Deployment

**ML Feature Engineering Model Continuous Training**
Feature Engineering
Train Model (Prod Data)
Test Model (Prod Data)
Hyper Parameters

**ML Model Deploy Serve**
Predict API
Model Performance Decay

**Feature Store**
**Metadata – Lineage - Feature**

**Prod Data**
Historical / Live Data

**A/B TESTING**

API

Batch

Streaming

**Gateway**

**App**

data Extract Refresh

data Extract Refresh
Sampling Anonymize

data Duplicate Refresh

F/Wall

F/Wall

F/Wall

Data Science Led Activities

Engineering Led Activities

# Data Platform
## Orchestration - Catalogue - Observability - Security



https://a16z.com/emerging-architectures-for-modern-data-infrastructure/

# ML/AI Platform
## Orchestration - Model & Data Drift - Audit-ability

**Data Transformation**  **Model Training and Development**  **Model Inference**  **Integration**



Data Labeling
(Scale, Labelbox, Snorkel, Sagemaker)

Model Diagnostics
(Labelbox, Scale, Nucleus, Aquarium)

Data Sources

Workflow Manager
(Airflow, Prefect, Pachyderm, Elementl/ Dagster, Tecton, Kubeflow)

Query Engines
(Presto, Hive)

Feature Store
(Tecton, Feast, Databricks)

Feature Server
(Tecton, Feast, Databricks)

Data Science Libraries
(Spark, Pandas, Numpy, Dask...)

Pre-trained Models
(Hugging Face, ModelZoo, PyTorch/TensorFlow)

ML APIs
(OpenAI, Cohere, AWS, GCP, Azure)

Data Science / Machine Learning Platform
(Jupyter, Databricks, Domino, Sagemaker, H2O, Colab, Deepnote, Noteable)

Model Registry
(MLflow, Sagemaker, Algorithmia, Hugging Face)

Batch Predictor
(Spark, etc.)

App Framework
(Flask, Streamlit, Rasa, etc.)

Clients

ML Framework
(Scikit-learn, XGBoost, MLlib)

Compiler
(OctoML/ TVM)

Online Model Server
(TensorFlow Serving, Ray Serve, Seldon)

Vector Database
(Faiss, Milvus, Pinecone)

DL Framework
(TensorFlow, Keras, PyTorch, H2O)

Validation
(Robust Intelligence, Calypso)

RL Libraries
(Gym, Dopamine, RLlib, Coach)

Auditing
(Credo, Armilla)

Distributed Processing
(Spark, Ray, Dask, Kubeflow, PyTorch, Tensorflow)

Experiment Tracking
(Weights & Biases, MLflow, Comet, ClearML)

Monitoring
(Arize, Fiddler, Arthur, Truera, WhyLabs, Gantry)

Low Code ML
(DataRobot, H2O, Databricks AutoML, Google AutoML, Continual, Mage, MindsDB, Obviously AI, Roboflow, Akkio)

# LLM Stack

## Foundation Models - Prompt - Fine Tune - RAG



https://a16z.com/emerging-architectures-for-llm-applications/