

Leveraging LLMs for Agentic Workflow Orchestration with Apache Beam YAML Pipelines



About the Presenter



Charles Adetiloye is a Cofounder and Lead AI Engineer at MavenCode. He has well over 18 years of experience building large-scale distributed applications with extensive experience working and consulting with several companies implementing production grade GenAI / Agentic AI / ML platforms.



About MavenCode



MavenCode is an Artificial Intelligence Solutions Company with HQ in Dallas, Texas and a remote delivery workforce across multiple time zones. We do training, product development and consulting services with specializations in:

- Provisioning Scalable AI and ML Infrastructure - OnPrem and In the Cloud
- Development & Production Operationalization of ML platforms - OnPrem and In the Cloud
- Streaming Data Analytics and Edge IoT Model Deployment for Federated Learning
- Building out Agentic AI, LLM and ML pipelines at scale.





Today's Agenda



1. Introduction to agentic workflow orchestration
2. Overview of Apache Beam YAML pipelines
3. Leveraging LLMs for pipeline automation
4. Designing modular orchestration agents
5. Future directions, and Q&A

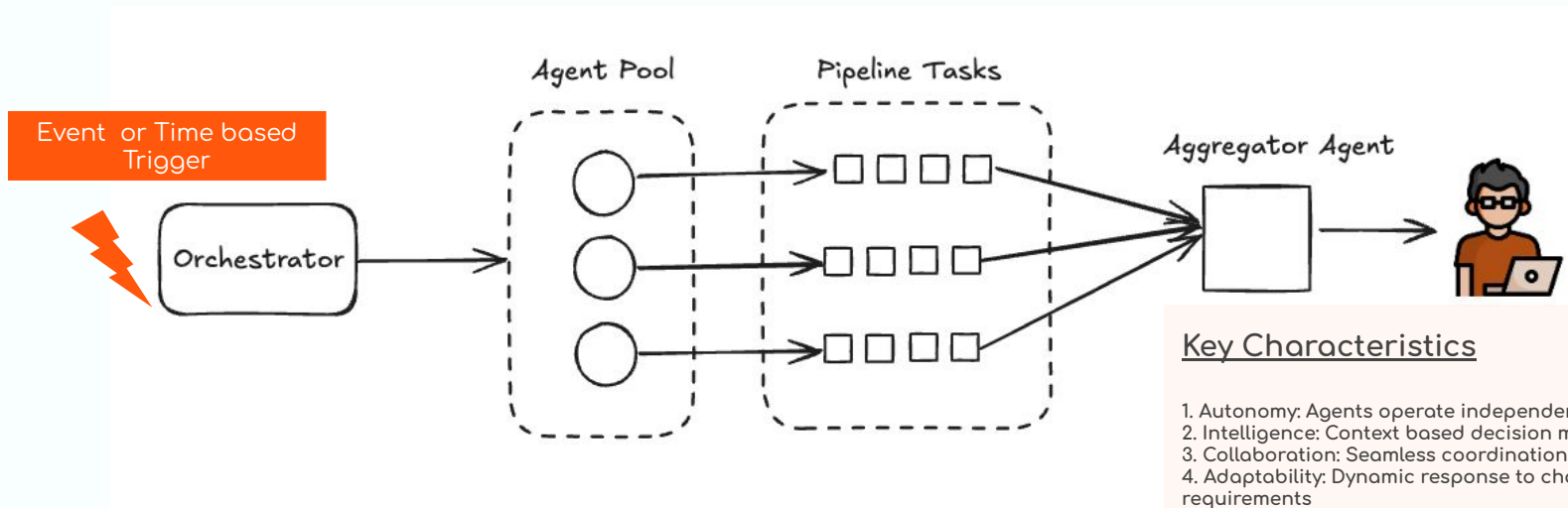
1. Introduction to Agentic Workflow Orchestration

What is Agentic Workflow Orchestration all About?

Agentic workflow orchestration leverages autonomous AI agents to dynamically coordinate, decide, and adapt across complex processes, enabling robust, scalable, intelligent, and resilient operations under varying conditions.

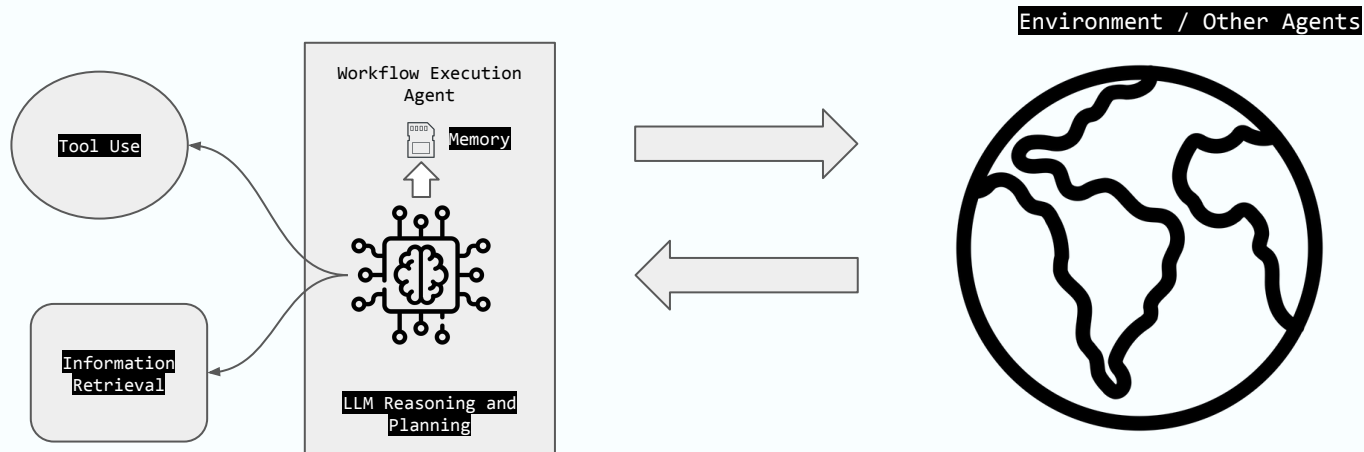
1. Introduction to Agentic Workflow Orchestration

Core Architecture



1. Introduction to Agentic Workflow Orchestration

Agentic Reasoning + Task Execution

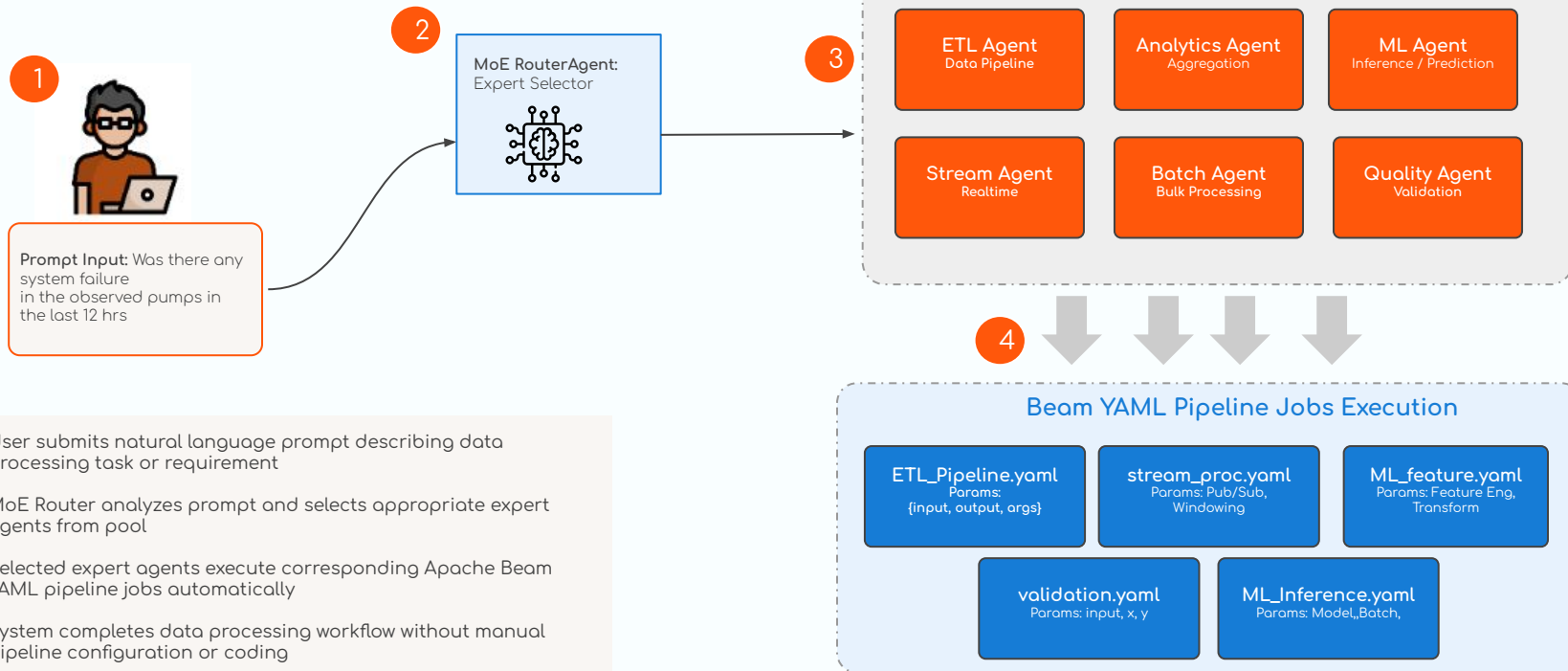


Key Objectives

1. Autonomous Decision-making aligned with Workflow goals
2. Coordinated multi-step task sequencing + resource allocation
3. Continuous monitoring with dynamic error recovery mechanism

1. Introduction to Agentic Workflow Orchestration

Agentic Execution of Apache Beam Workflow



2. Overview of Apache Beam YAML

Apache Beam YAML Pipelines

Apache Beam YAML is a declarative syntax for describing Apache Beam pipelines using YAML files. You can use Beam YAML to author and run a Beam pipeline without writing any code. This approach makes data processing pipelines more accessible by eliminating the need to write code in traditional Beam SDK languages.

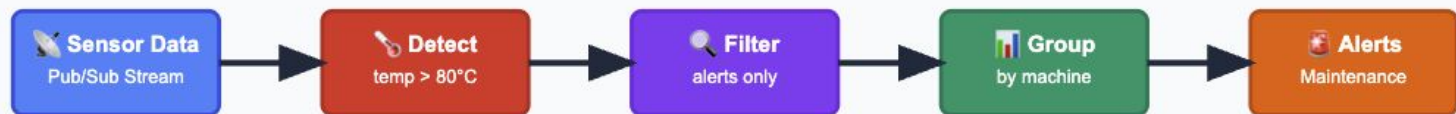
2. Overview of Apache Beam YAML

Key Benefits of Beam YAML pipeline

- **No-code Development:** Create sophisticated data pipelines using only YAML configuration without programming
- **Cloud-native ready:** Deploy instantly to Kubernetes / Google Dataflow with automatic scaling and management
- **Easy maintenance:** Update pipeline logic by editing text files instead of recompiling code
- **Self-documenting code:** YAML declarative nature makes AI-generated pipelines naturally explainable and auditable
- **Template-based generation:** GenAI can easily modify reusable YAML templates for different data scenarios

2. Overview of Apache Beam YAML

Real-time Anomaly Detection and Preemptive Notification



- Read: Stream log data from Manufacturing Plant
- Detect Machines that are overheating
- Filter and keep the identified machines
- Send alert notifications by Machines

```
pipeline:
  type: chain

source:
  type: ReadFromPubSub
  config:
    topic: projects/[REDACTED]/topics/sensor-data
    format: JSON

transforms:
  # Detect sensor anomalies
  - type: MapToFields
    config:
      language: python
      fields:
        machine_id: machine_id
        temperature: temperature
        pressure: pressure
        vibration: vibration
      # Simple anomaly flags
      temp_alert: "temperature > 80"
      pressure_alert: "pressure > 150"
      vibration_alert: "vibration > 10"
      alert_level: |
        'CRITICAL' if temperature > 90 or pressure > 160
        else 'WARNING' if temperature > 80 or pressure > 140
        else 'NORMAL'

  # Only keep problematic readings
  - type: Filter
    config:
      language: python
      keep: "temp_alert or pressure_alert or vibration_alert"

  # Group by machine and count issues
  - type: Sql
    config:
      query: |
        SELECT
          machine_id,
          COUNT(*) as alert_count,
          MAX(temperature) as max_temp,
          MAX(pressure) as max_pressure,
          COUNT(CASE WHEN alert_level = 'CRITICAL' THEN 1 END) as critical_count
        FROM PCOLLECTION
        GROUP BY machine_id
```

2. Overview of Apache Beam YAML

2

```
# Add maintenance recommendations
- type: MapToFields
  config:
    language: python
    fields:
      machine_id: machine_id
      alert_count: alert_count
      max_temp: max_temp
      max_pressure: max_pressure
      critical_count: critical_count
    # Simple maintenance logic
    action_needed: |
      'SHUTDOWN' if critical_count > 2
      else 'MAINTENANCE' if alert_count > 5
      else 'MONITOR'

sink:
  type: WriteToPubSub
  config:
    topic: projects/[REDACTED]/topics/maintenance-alerts
    format: JSON

options:
  streaming: true
  runner: DataflowRunner
```

```
pipeline:
  type: chain

source:
  type: ReadFromPubSub
  config:
    topic: projects/[REDACTED]/topics/sensor-data
    format: JSON

transforms:
  # Detect sensor anomalies
  - type: MapToFields
    config:
      language: python
      fields:
        machine_id: machine_id
        temperature: temperature
        pressure: pressure
        vibration: vibration
      # Simple anomaly flags
      temp_alert: "temperature"
      pressure_alert: "pressure"
      vibration_alert: "vibration"
      alert_level: |
        'CRITICAL' if temperature > 100
        else 'WARNING' if temperature > 80
        else 'NORMAL'

  # Only keep problematic readings
  - type: Filter
    config:
      language: python
      keep: "temp_alert or pressure_alert or vibration_alert"

  # Group by machine and count issues
  - type: Sql
    config:
      query: |
        SELECT
          machine_id,
          COUNT(*) as alert_count,
          MAX(temperature) as max_temp,
          MAX(pressure) as max_pressure,
          COUNT(CASE WHEN alert_level = 'CRITICAL' THEN 1 END) as critical_count
        FROM PCOLLECTION
        GROUP BY machine_id
```

2. Overview of Apache Beam YAML

2

```
# Add maintenance recommendations
- type: MapToFields
```

Local Dev:

```
python -m apache_beam.yaml.main \
  --yaml_pipeline_file=digital_stream_pipeline.yaml \
  --runner=DirectRunner
```

Dataflow Runner:

```
gcloud dataflow yaml run digital-stream-process-monitoring \
  --yaml-pipeline-file=digital_stream_pipeline.yaml \
  --region=us-central1 \
  --max-workers=10 \
  --enable-streaming-engine
```

```
config:
  topic: projects/[REDACTED]/topics/maintenance-alerts
  format: JSON

options:
  streaming: true
  runner: DataflowRunner
```

3. Leverage LLM for Beam Pipeline Orchestration

Prompt Driven Pipeline Generation

Define transforms or trigger pipelines in natural language on the fly.

RAG Enhanced Context Retrieval

Fetch relevant pipeline templates, docs, or metrics to ground LLM decisions

Adaptive Branching Logic

Real-time rerouting of elements when anomalies or new conditions are detected

Schema & Code Synthesis

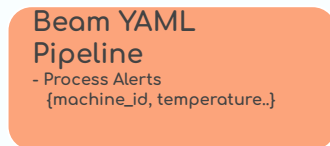
Auto-generate YAML/SDK snippets for connectors, transforms, and I/O.

Human-in-the-Loop Verification

Insert review checkpoints, leveraging RAG to surface past best practices

3. Leverage LLM for Beam Pipeline Orchestration

LLM Enriched Pipeline

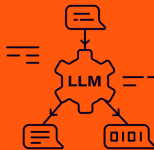


Sample Alert Data

```
{
  "machine_id": "001",
  "temp": "420",
  "pressure": 18.7,
  "process_unit": "Distillation",
  "alert_level": "critical"
}
```

Data Enrichment

- create incident reports
Add equipment context



```
import apache_beam as beam
import requests
import json
import logging
import os

# --- Configuration for your self-hosted vLLM endpoint ---
# The internal URL for your vLLM server.
VLLM_ENDPOINT = "http://your-vllm-host.internal:8000/v1/chat/completions"

def enrich_with_data(element: beam.Row) -> beam.Row:
    """
    Calls a self-hosted, unsecured vLLM endpoint to enrich data.
    Uses the OpenAI-compatible API format without an API key.
    """
    # 1. Format the prompt for the model.
    prompt = (
        f"Analyze the following industrial machine data and provide a one-sentence summary "
        f"of its operational status. Data: "
        f"Machine ID: {element.machine_id}, "
        f"Temperature: {element.temperature}°C, "
        f"Pressure: {element.pressure} psi, "
        f"Vibration: {element.vibration_freq} Hz."
    )

    # 2. Prepare the request headers. No Authorization is needed.
    headers = { 'Content-Type': 'application/json' }

    # 3. Prepare the OpenAI-compatible payload.
    # The model name must match the model you loaded into the vLLM server.
    payload = {
        "model": "mistralai/Mistral-7B-Instruct-v0.2", # <- IMPORTANT: Use your loaded model's name
        "messages": [
            {
                "role": "system", "content": "You are an expert AI for industrial machine monitoring.",
            },
            {
                "role": "user", "content": prompt
            }
        ],
        "max_tokens": 150,
        "temperature": 0.5
    }

    try:
        # 4. Make the unauthenticated API call to your local vLLM server.
        response = requests.post(VLLM_ENDPOINT, headers=headers, json=payload, timeout=30)
        response.raise_for_status()

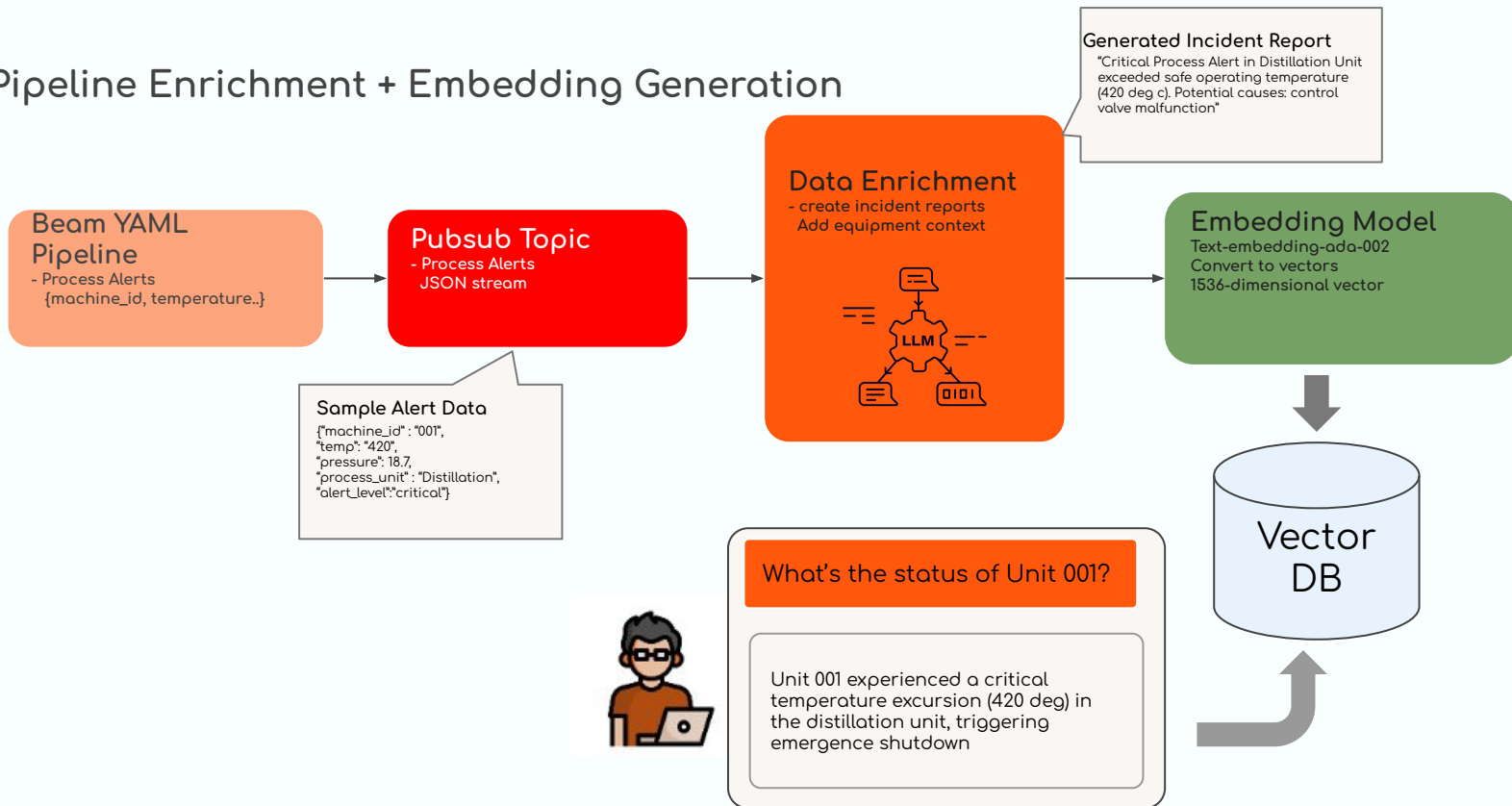
        # 5. Parse the OpenAI-compatible response.
        api_response = response.json()
        enriched_message = api_response['choices'][0]['message']['content'].strip()

    except requests.exceptions.RequestException as e:
        logging.error(f"vLLM API call failed for machine {element.machine_id}: {e}")
        enriched_message = "Error: Could not reach self-hosted SLM."
    except (KeyError, IndexError) as e:
        logging.error(f"Failed to parse vLLM API response for machine {element.machine_id}: {e}")
        enriched_message = "Error: Invalid API response format from self-hosted SLM."

    # 6. Return the enriched Beam Row.
    return beam.Row(["element_asdict()", enriched_message=enriched_message])
```

3. Leverage LLM for Beam Pipeline Orchestration

Pipeline Enrichment + Embedding Generation



3. Leverage LLM for Beam Pipeline Orchestration

2

```

import apache_beam as beam
import requests
import json
import logging
import os

# --- Configuration for your self-hosted vLLM endpoint ---
# The internal URL for your vLLM server.
VLLM_ENDPOINT = "http://your-vllm-host.internal:8000/v1/chat/completions"

def enrich_with_data(element: beam.Row) -> beam.Row:
    """
    Calls a self-hosted, unsecured vLLM endpoint to enrich data.
    Uses the OpenAI-compatible API format without an API key.
    """

    # 1. Format the prompt for the model.
    prompt = (
        f"Analyze the following industrial machine data and provide a one-sentence "
        f"of its operational status. Data: "
        f"Machine ID: {element.machine_id}, "
        f"Temperature: {element.temperature}°C, "
        f"Pressure: {element.pressure} psi, "
        f"Vibration: {element.vibration_freq} Hz."
    )

    # 2. Prepare the request headers. No 'Authorization' is needed.
    headers = { "Content-Type": "application/json" }

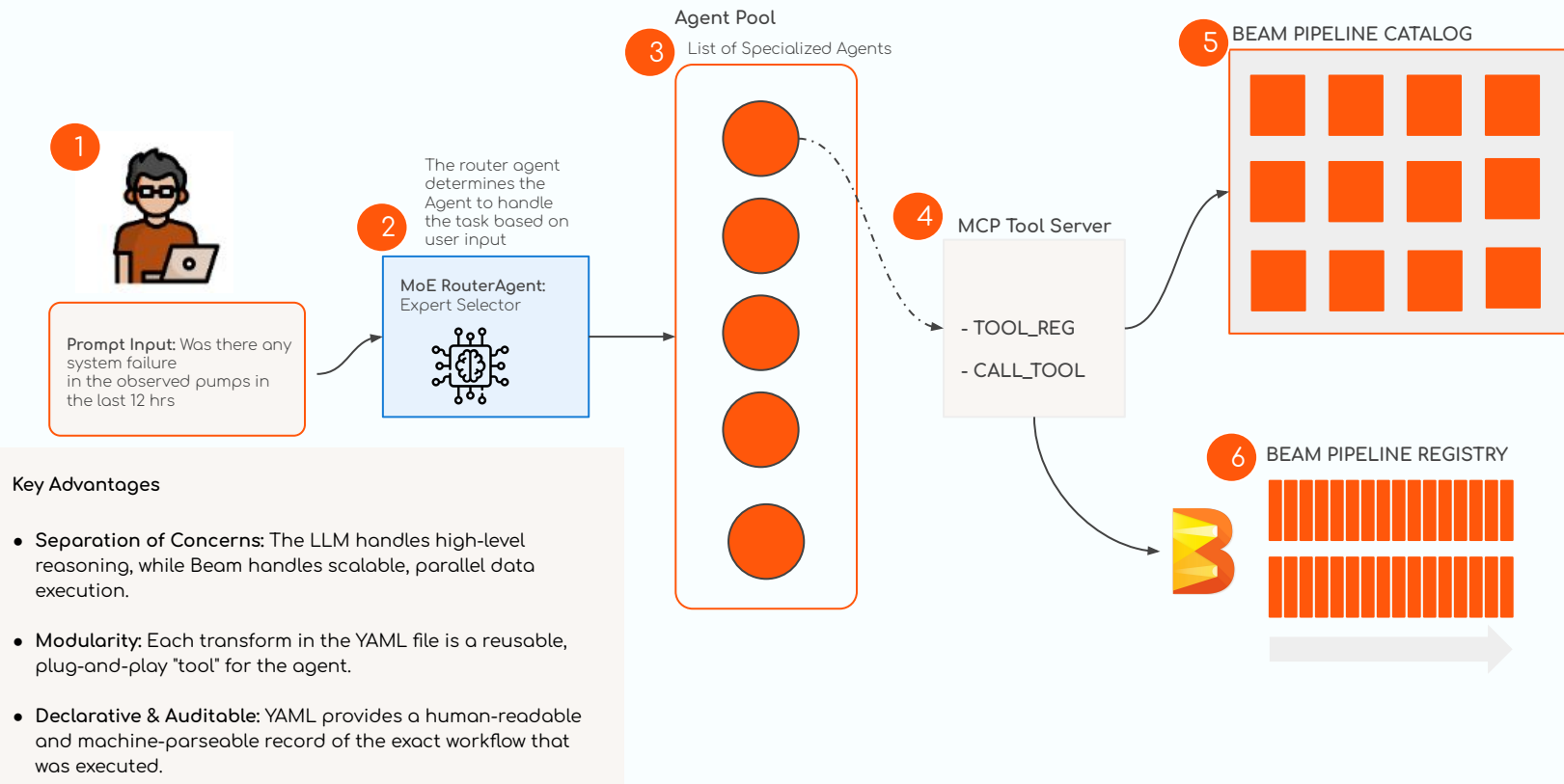
    # 3. Prepare the OpenAI-compatible payload.
    # The 'model' name must match the model you loaded into the vLLM server.
    payload = {
        "model": "mistralai/Mistral-7B-Instruct-v0.2", # <-- IMPORTANT: Use your local
        "messages": [
            { "role": "system", "content": "You are an expert AI for industrial mach
            { "role": "user", "content": prompt }
        ],
        "max_tokens": 150,
        "temperature": 0.5
    }

    try:
        # 4. Make the unauthenticated API call to your local vLLM server.
        response = requests.post(VLLM_ENDPOINT, headers=headers, json=payload, timeout=
        response.raise_for_status()

        # 5. Parse the OpenAI-compatible response.
        api_response = response.json()

```

4. Designing Modular Orchestration Agents



4. Designing Modular Orchestration Agents

Pipeline Tools Catalog Agent Tools for Apache Beam Pipelines



GET /tools Get Tools



POST /stream-pipeline-tool Stream Agent Tool



POST /etl-pipeline-tool Etl Agent Tool



POST /analytics-pipeline-tool Analytics Agent Tool



POST /batch-pipeline-tool Batch Agent Tool



POST /ml-pipeline-tool ML Agent Tool



POST /quality-pipeline-tool Quality Agent Tool



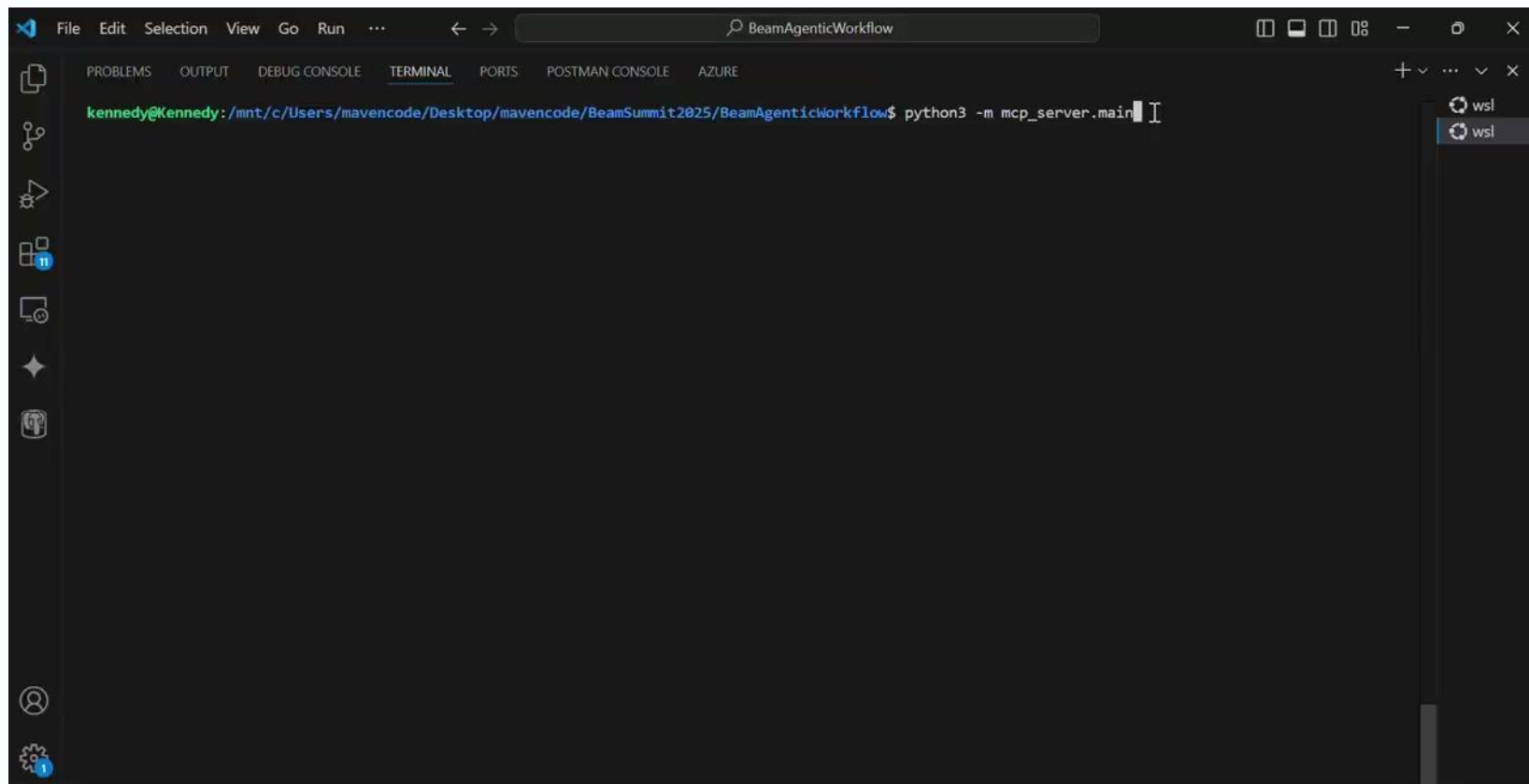
POST /pipeline-run-tool Pipeline Run Tool



4. Designing Modular Orchestration Agents

GET	/tools	Get Tools	▼
POST	/stream-pipeline-tool	Stream Agent Tool	▼
POST	/etl-pipeline-tool	Etl Agent Tool	▼
POST	/analytics-pipeline-tool	Analytics Agent Tool	▼
POST	/batch-pipeline-tool	Batch Agent Tool	▼
POST	/ml-pipeline-tool	ML Agent Tool	▼
POST	/quality-pipeline-tool	Quality Agent Tool	▼
POST	/pipeline-run-tool	Pipeline Run Tool	▼
Schemas			
AnalyticsToolInput > Expand all object			

4. Designing Modular Orchestration Agents



The image shows a screenshot of a Visual Studio Code (VS Code) terminal window. The terminal is open to the 'TERMINAL' tab, showing a command prompt for a user named 'kennedy' on a machine named 'Kennedy'. The current directory is '/mnt/c/Users/mavencode/Desktop/mavencode/BeamSummit2025/BeamAgenticWorkflow'. The command entered is 'python3 -m mcp_server.main'. The terminal output is empty. The VS Code interface includes a menu bar at the top with 'File', 'Edit', 'Selection', 'View', 'Go', and 'Run'. The left sidebar shows various icons for file explorer, search, and other tools. The right sidebar shows a list of open files, including 'wsl'.

```
kennedy@Kennedy: /mnt/c/Users/mavencode/Desktop/mavencode/BeamSummit2025/BeamAgenticWorkflow$ python3 -m mcp_server.main
```

7. Future Directions

- **Dynamic YAML Synthesis from Natural Language:** Agent translates natural language goals into executable `beam.yml` files by selecting transforms from a catalog.
- **Adaptive Workflow Repair via YAML Modification:** Agent analyzes pipeline errors, programmatically edits the `beam.yml` to fix the issue, and automatically relaunches the job.
- **Declarative A/B Testing of Pipeline Logic:** Agent generates multiple `beam.yml` versions for A/B testing transforms, then launches another YAML pipeline to compare results.
- **Automated Discovery and Integration of New Transforms:** Agent automatically scans documentation for new transforms, adding them to its catalog for use in future YAML pipelines.

Charles Adetiloye

QUESTIONS?

<https://www.linkedin.com/in/charlesadetiloye/>

github.com/MavenCode