# Guide to Using Large Multimodal Models v1.1

*Technical Supplement 4 — Governance & Deployment: Orchestration & Chaining of AI Workflows*

## About This Supplement

### Purpose & Scope
This supplement provides governance-level guidance for building controlled, multi-step AI workflows that integrate reasoning, tool use, and mandatory human oversight. It expands on the practical framework introduced in Section 2 of the Core Guide.

**Audience:** Managers, System Architects, Project Managers, AI Engineers
**Prerequisites:** Familiarity with the C.G.A.F.R. framework (Core Guide § 2) and verification workflow (Core Guide § 4).
**Outcome:** Ability to design, document, and execute auditable multi-step AI workflows using the "Decide → Act → Verify → Record" orchestration pattern.

**Key Objectives:**

- Systematically identify, assess, and mitigate ethical, operational, and reputational risks in LMM projects.
- Align LMM applications with organizational values and regulatory requirements from the initial design phase.
- Create a standardized risk assessment template for consistent evaluation and governance oversight of all AI use cases.

Developed by Russell Nida

# Contents

# 1 Concept: The LMM as a Workflow Controller

While Section 3 of the Guide introduced basic chaining (e.g., "generate, then ask for risks"), orchestration formalizes this into a controlled system. Here, the LMM acts as a reasoning engine and controller, making explicit decisions about the next step based on prior outputs. This enables workflows that dynamically route tasks, trigger tools, and escalate to humans based on predefined logic and confidence scores.

A well-orchestrated workflow explicitly manages four actions at each step:

| Action | Function | Example |
|---|---|---|
| **Decide** | Interpret the context and select the next action. | "If risk score > 8, trigger the 'High-Risk Alert' protocol; otherwise, summarize." |
| **Act** | Execute the chosen operation using an LMM or external tool. | "Run a Python script to analyze data or generate a report." |
| **Verify** | Evaluate the output's correctness or adherence to criteria. | "Check that the generated SQL query runs and summary covers all data points." |
| **Record** | Log inputs, outputs, reasoning, and confidence. | "Append Timestamp, Step ID, Prompt, Output, Confidence, Verification Result." |

**Key Principle:** Each step must define its **purpose**, **input**, **success criteria**, and **verification gate.**

---

# 2 The Core Orchestration Pattern: "Decide → Act → Verify → Record"

Once you understand the conceptual loop from Section 1, this section shows how to express it directly as a model instruction. This loop is the fundamental unit of a controlled workflow.

**Step Prompt Template**

Context: You are executing step [Step ID] of the "[Workflow Name]" workflow.
Previous Step Output: [Insert prior output].
Overall Goal: [Describe objective].

Your Instructions:
DECIDE: Based on the context, what is the single most logical next action?
ACT: Execute that action. Provide the core output as specified.

VERIFY: Self-evaluate your output.
  Confidence Score (1–10):
  Uncertainties or Missing Data:
RECORD: Combine ACT output + VERIFY summary.

**Critical Caution:** This pattern is a **risk amplifier**. An early hallucination can corrupt all downstream steps. Always implement **verification gates** after critical steps—either automated tests or human sign-offs. The "Record" output is your primary forensic tool and is non-negotiable for accountability.

---

# 3 Design Template – Documenting Each Step

| Field | Description | Example |
|---|---|---|
| Step ID | Unique label for tracking | WF-1.3 |
| Purpose | Why this step exists | Validate extracted numbers |
| Inputs | Files, data, or prior outputs | CSV rows 5–50 |
| Model / Tool | LMM or API used | GPT-5 + Python sandbox |
| Expected Output | Format & criteria | JSON report + confidence score |
| Verification Method | Pass/fail rule | Cross-check sum = total column |
| Reviewer / Owner | Human responsible | QA Lead – J. Doe |

---

# 4 Applied Example: From Sensor Data to Alert

| Field | Step 1: Anomaly Detection | Step 2: Risk Assessment | Step 3: Alert Drafting |
|---|---|---|---|
| Step ID | WF-1.1 | WF-1.2 | WF-1.3 |
| Purpose | Identify anomalies in sensor log. | Classify the severity of anomalies. | Draft notification for high severity. |
| Inputs | sensor_data.csv | Output from WF-1.1 | Output from WF-1.2 |
| Model / Tool | LMM + Code Interpreter | LMM | LMM |
| Expected Output | JSON list of anomalies. | Risk classification: Low/Medium/High. | Draft email body. |
| Verific | Code runs, valid JSON. | Cross-check baseline. | HUMAN REVIEW |

| Field | Step 1: Anomaly Detection | Step 2: Risk Assessment | Step 3: Alert Drafting |
|---|---|---|---|
| ation Method | | | REQUIRED before send. |
| Reviewer / Owner | Data Eng – A. Smith | Data Eng – A. Smith | Incident Mgr – B. Jones |

# 5 Scaling and Governance Guidelines

- **Start Simple:** Begin with linear 3–4 step chains before branching.
- **Isolate Models by Function:** Use separate system prompts for generation, verification, and summarization to reduce bias.
- **Version Control Everything:** Track workflow definition, prompts, and model versions with validation results.
- **Mandate Human Checkpoints:** Any output affecting policy, safety, or people requires human sign-off. Design checkpoints into the workflow, not as afterthoughts.
- **Integrate with Audit Systems:** The "Record" outputs should feed directly into logs or compliance dashboards for review.

# 6 Cost Management for Orchestrated Workflows

Orchestration transforms prompting into process. Treat each chain as a **mini-workflow** with explicit decision points, risk controls, and documentation. The goal is not autonomy—it is **accountable automation.**

When moving from individual prompts to multi-step orchestrated workflows, cost control evolves from a personal concern to an operational requirement. This section provides a framework for monitoring and optimizing token usage across deployed AI systems.

## 6.1 Understanding Cost Drivers in Orchestration

The primary cost drivers shift in an orchestrated environment:

| Cost Driver | Description | Management Strategy |
|---|---|---|
| **Iteration Amplification** | A single poorly-designed step, repeated across hundreds of workflow runs, wastes significant budget. | Invest in robust prompt design and testing for each step in the workflow template. |

| Cost Driver | Description | Management Strategy |
|---|---|---|
| **Context Accumulation** | Passing large intermediate outputs between steps bloats the context window, increasing the cost of every subsequent step. | Design workflows to pass concise summaries or structured data (e.g., JSON) between steps, not verbose narratives. |
| **Tool Execution Overhead** | Unnecessary or inefficient use of tools (e.g., web search, code execution) adds cost without value. | Gate tool use behind clear decision logic. Use the cheapest effective tool for the task. |

## 6.2 Optimization Strategies for System Architects

1. **Implement Usage Logging and Alerts:** Ensure your orchestration platform logs token consumption per step. Set up alerts for abnormal usage patterns that suggest a step is failing or misconfigured.

2. **Benchmark and Right-Size Models:** Not every step requires the most powerful (and expensive) model. Use cheaper, faster models for simple steps like formatting or routing, and reserve high-cost models for complex reasoning tasks.

3. **Design for Conciseness:** Enforce output constraints in every step's prompt. A step should be instructed to "Output a JSON object with keys X, Y, Z" rather than generating a long paragraph that the next step must painfully process.

4. **Build Cost into the "Verify" Step:** In the core "Decide → Act → **Verify** → Record" loop, include a cost-awareness check. For example, a verification step could flag outputs that are abnormally long for further review.

## 6.3 The Cost-Reliability Trade-Off

Optimization must not compromise the verification and safety gates that make the workflow reliable. The goal is to eliminate *waste*, not to skip necessary review. A failed workflow that costs $0.50 but produces a flawed output used in a client report is far more expensive than a reliable one that costs $2.00.

**Key Takeaway:** Treat token cost as a key performance indicator (KPI) for your AI workflows. Managing it effectively is a sign of a mature, well-governed deployment.

---

# 7 Key Takeaway

Orchestration transforms prompting into process. Treat each chain as a **mini-workflow** with explicit decision points, risk controls, and documentation. The goal is not autonomy—it is **accountable automation.**

# 8 Cross-References

- For verification standards and troubleshooting, see **Technical Supplement 2 – Verification, Auditing & Troubleshooting.**

- For data customization and RAG/fine-tuning frameworks, see **Technical Supplement 5 – Customization & Fine-Tuning.**

- For security controls, see **Technical Supplement 6 – Security & Governance.**