

L2-M5: Introduction to AI Workflows - Breaking Down Complex Tasks into Multi-Step Sequences

1. Introduction to AI Workflows

The rapid evolution of Artificial Intelligence (AI), particularly in the domain of Large Language Models (LLMs) and Generative AI, has shifted the focus from isolated model performance to the **systemic integration** of AI components. A single, monolithic AI model is often insufficient for addressing complex, real-world problems that require multiple steps, diverse data sources, and dynamic decision-making [1]. This necessitates the adoption of **AI Workflows**, which are structured, multi-step processes designed to orchestrate various AI models, traditional software components, and human interactions to achieve a specific, complex objective.

An AI workflow is fundamentally a pipeline or a directed acyclic graph (DAG) where the output of one component serves as the input for the next. This architecture moves beyond the traditional "model-in-a-box" approach, enabling systems that are more reliable, interpretable, and capable of handling ambiguity and complexity [2]. Understanding the principles of workflow design is crucial for any intermediate to advanced practitioner in the AI workforce.

1.1. The Shift from Monolithic Models to Orchestrated Systems

Historically, AI development focused on creating single models optimized for a narrow task, such as image classification or simple text generation. However, tasks like "Develop a comprehensive market entry strategy for a new product" or "Automate the entire customer service resolution process" require a sequence of distinct cognitive steps: research, analysis, synthesis, drafting, and revision.

Feature	Monolithic AI Model	AI Workflow (Orchestrated System)
Task Complexity	Simple, narrow, single-step tasks	Complex, multi-faceted, multi-step tasks
Components	Single, large model (e.g., GPT-4)	Multiple specialized models, tools, and agents
Reliability	Single point of failure; opaque error handling	Distributed failure points; explicit error pathways
Cost Efficiency	High cost per inference for large models	Potential for cost optimization by using smaller, specialized models for subtasks [3]
Interpretability	Low; "black box" operation	Higher; each step's output is traceable and auditable

2. Task Decomposition: The Foundation of AI Workflows

Task Decomposition is the process of breaking down a large, complex goal into a series of smaller, manageable sub-tasks. This is the critical first step in designing any effective AI workflow. Just as a software engineer decomposes a large project into modules and functions, an AI system must decompose a complex prompt into actionable steps [4].

2.1. Methods for Task Decomposition

Several methodologies, often facilitated by LLMs themselves, are employed for effective task decomposition:

2.1.1. Chain-of-Thought (CoT) and Step-by-Step Planning

The simplest form of decomposition is prompting the LLM to think step-by-step. This technique, known as Chain-of-Thought (CoT) prompting, forces the model to articulate its reasoning process, which inherently involves task breakdown.

Example Prompt:

"Analyze the Q3 2025 earnings report for Company X, identify the three biggest risks, and draft a one-page executive summary."

CoT Decomposition: 1. **Retrieve and parse** the Q3 2025 earnings report for Company X. 2. **Analyze** the financial statements and management discussion to identify potential risks. 3. **Filter and rank** the identified risks to select the top three most significant ones. 4. **Synthesize** the findings into a concise, professional executive summary.

2.1.2. Hierarchical Task Planning (HTP)

HTP involves creating a tree-like structure where the main goal is the root, and sub-goals branch out, which can then be further decomposed into primitive actions. This is particularly useful for robotic or agentic systems that require grounding in physical or digital environments [5].

- **Goal:** Book a flight from London to New York.
 - **Sub-Goal 1: Search Flights**
 - Action 1.1: Query flight API for dates.
 - Action 1.2: Filter results by price and airline.
 - **Sub-Goal 2: Select and Book**
 - Action 2.1: Select the optimal flight.
 - Action 2.2: Initiate payment process.

2.1.3. ReAct (Reasoning and Acting)

The ReAct framework interweaves reasoning (CoT) and acting (tool use) [6]. The model reasons about the current state, decides on the next action (e.g., using a search engine, running code), observes the result, and then reasons again. This continuous loop of planning, execution, and observation is a form of dynamic, iterative task decomposition.

3. Core AI Workflow Design Patterns

Effective AI workflows are built upon established architectural patterns that govern the flow of data, control, and decision-making. These patterns ensure scalability, robustness, and maintainability.

3.1. Sequential Pipeline Pattern

The most straightforward pattern, where data flows linearly from one component to the next.

- **Structure:** Component A → Component B → Component C
- **Application:** Data preprocessing pipelines, simple document generation.
- **Example:**
 1. **Transcription Model:** Converts audio to text.
 2. **NER Model:** Extracts named entities (names, dates) from the text.
 3. **Summarization Model:** Generates a concise summary of the entity-rich text.

3.2. Retrieval-Augmented Generation (RAG) Pattern

The RAG pattern is a critical advancement in LLM-based workflows, addressing the limitations of models' knowledge cutoffs and tendency to hallucinate. It augments the LLM's generation with relevant, external, and verifiable information [7].

- **Structure:** Query → **Retrieval System** (Vector DB) → Context Chunks → **LLM (Generation)** → Answer
- **Application:** Enterprise search, specialized Q&A, knowledge-grounded chatbots.
- **Step-by-Step RAG Workflow:**
 1. **User Query:** The user asks a question (e.g., "What is the policy on remote work?").
 2. **Embedding:** The query is converted into a vector representation.

3. **Retrieval:** The vector is used to search a Vector Database (containing embedded company documents) for the most semantically relevant text chunks.
4. **Prompt Construction:** The retrieved text chunks are combined with the original query and a system instruction into a single, comprehensive prompt (the "context").
5. **Generation:** The LLM generates an answer based *only* on the provided context, ensuring the response is factual and grounded in the source material.

3.3. Router/Orchestrator Pattern

In this pattern, a central **Router** or **Orchestrator** component (often an LLM itself) is responsible for analyzing the user's request and dynamically selecting the appropriate downstream tool or specialized model to handle the task [8].

Component	Function
Orchestrator (LLM)	Analyzes the input, determines the intent, and selects the next step.
Tool/Agent Pool	A collection of specialized models, APIs, or code functions (e.g., Calculator, Weather API, Database Query Agent).
Feedback Loop	The result from the selected tool is returned to the Orchestrator for evaluation or further action.

Real-World Application (Customer Service Bot): 1. **Input:** "I need to check my order status." 2. **Orchestrator Decision:** Intent is "Order Status." Route to the **Database Query Agent**. 3. **Input:** "What is the current stock price of AAPL?" 4. **Orchestrator Decision:** Intent is "Real-time Data." Route to the **Stock Market API Tool**.

3.4. Hierarchical Agentic Pattern

This pattern introduces a structure of specialized agents, often with a "Manager" or "Supervisor" agent overseeing "Worker" agents [9]. The Manager agent handles the high-level task decomposition and delegation, while the Worker agents execute specific sub-tasks, potentially using their own tools.

- **Manager Agent (Planner):** Breaks down the goal ("Write a research paper on quantum computing") into major sections (Introduction, Literature Review, Methodology, Results).
- **Worker Agent 1 (Researcher):** Executes the "Literature Review" task, using a search tool and a RAG pattern to gather sources.
- **Worker Agent 2 (Writer/Synthesizer):** Takes the output from the Researcher and drafts the section content.
- **Worker Agent 3 (Editor):** Reviews and refines the final draft for coherence and style.

4. Practical Applications and Implementation

Implementing AI workflows requires careful consideration of infrastructure, state management, and error handling. Modern workflow engines, such as Apache Airflow, Kubeflow, and specialized AI frameworks like LangChain or Microsoft's Semantic Kernel, provide the necessary scaffolding [10].

4.1. Case Study: Automated Investment Research

A financial institution seeks to automate the generation of daily investment summaries for a basket of 100 stocks.

Step	Component/Model	Function	Output
1. Data Ingestion	Custom Python Script (API Tool)	Fetches real-time stock prices, news headlines, and social media sentiment for 100 tickers.	Structured JSON data per ticker.
2. Sentiment Analysis	Specialized Fine-Tuned LLM (Sentiment Model)	Analyzes news and social media text for positive/negative sentiment scores.	Sentiment Score (0-1) and Key Phrases.
3. Risk Assessment	Traditional Machine Learning Model (XGBoost)	Uses historical data, price volatility, and sentiment scores to calculate a risk score.	Risk Score (Low, Medium, High).
4. Summary Generation	General-Purpose LLM (GPT-4)	Receives the JSON data, sentiment scores, and risk score. Instructed to generate a 250-word summary.	Final Investment Summary (Text).
5. Distribution	Email API Tool	Sends the generated summary to the distribution list.	Email Confirmation.

This workflow demonstrates the orchestration of traditional code, specialized ML models, and a powerful LLM for synthesis, highlighting the multi-step nature of complex AI tasks.

4.2. State Management and Persistence

A critical challenge in multi-step workflows is maintaining **state** and **context**. Since each step is often an independent call, the system must ensure that the necessary information from previous steps is correctly passed to the next.

- **Context Passing:** The output of Step N must be explicitly formatted and passed as input to Step N+1. In LLM workflows, this often means dynamically constructing a new, larger prompt that includes the results of prior actions.
- **Persistence:** For long-running workflows, the state must be persisted to a database (e.g., a key-value store or a relational database) to allow for recovery from failures or for human review at checkpoints.

5. Advanced Workflow Concepts: Self-Correction and Reflection

To achieve human-level performance, AI workflows must incorporate mechanisms for **self-correction** and **reflection**, moving beyond simple linear execution.

5.1. Reflection-Based Workflows

Reflection involves an agent or a dedicated LLM step reviewing the output of a previous step against a set of criteria or the original goal [11].

- **Process:**
 1. **Execution:** Agent generates a draft.
 2. **Reflection:** A second, critical LLM (the Reflector) receives the draft and the original prompt.
 3. **Critique:** The Reflector assesses the draft for tone, accuracy, completeness, and adherence to constraints.
 4. **Revision:** The critique is passed back to the original Agent (or a new Revision Agent) along with the draft for an iterative correction loop.

This pattern is essential for high-stakes tasks like legal drafting or complex coding, where a single-pass generation is prone to errors.

5.2. Tool Use and Dynamic Function Calling

Modern workflows leverage **dynamic function calling** (or "tool use") where the LLM is given a description of available external functions (e.g., `get_weather(city)`, `run_sql_query(query)`). The LLM's role is to decide *when* and *how* to call these functions by outputting a structured function call object, which is then executed by the workflow engine.

Tool Use Benefit	Description
Grounding	Access to real-time, external, or proprietary data, preventing hallucinations.
Computation	Offloading deterministic, complex calculations (e.g., math, code execution) to reliable tools.
Action	Interacting with the external world (e.g., sending an email, updating a database).

6. Conclusion and Key Takeaways

The future of applied AI lies not in isolated models but in sophisticated, orchestrated workflows. The ability to decompose a complex task into a sequence of smaller, manageable sub-tasks is the foundational skill for designing robust AI systems. By applying established design patterns—such as Sequential Pipelines, RAG, and the Orchestrator/Agentic Hierarchy—practitioners can build scalable, cost-effective, and highly reliable AI solutions. The integration of reflection and dynamic tool use further elevates these systems, allowing them to self-correct and interact dynamically with the real world.

Key Takeaways

1. **AI Workflows are Essential:** Complex problems require orchestrated systems, not monolithic models.
 2. **Decomposition is Key:** Utilize methods like CoT, HTP, and ReAct to break down goals into actionable steps.
 3. **Master Design Patterns:** The RAG pattern is crucial for knowledge-grounding; the Orchestrator pattern is vital for dynamic decision-making and tool selection.
 4. **Prioritize Robustness:** Implement strong state management and persistence for long-running processes.
 5. **Embrace Iteration:** Integrate Reflection and self-correction loops to improve the quality and reliability of outputs.
-

References

- [1] Malhotra, L. (2025). *The Architecture of AI Workflows: Designing Beyond the Model Layer*. Dev.to. https://dev.to/leena_malhotra/the-architecture-of-ai-workflows-designing-beyond-the-model-layer-45ld
- [2] Microsoft Azure. (2025). *AI Agent Orchestration Patterns*. Azure Architecture Center. <https://learn.microsoft.com/en-us/azure/architecture/guide/ai-agent-design-patterns>
- [3] Amazon Science. (n.d.). *How task decomposition and smaller LLMs can make AI more affordable*. <https://www.amazon.science/blog/how-task-decomposition-and-smaller-langs-can-make-ai-more-affordable>
- [4] Zhai, W. et al. (2025). *A Survey of Task Planning with Large Language Models*. IComputing. <https://spj.science.org/doi/10.34133/icomputing.0124>
- [5] Ma, F. et al. (2024). *Task Navigator: Decomposing Complex Tasks for Multimodal Large Language Models*. CVPR Workshop. https://openaccess.thecvf.com/content/CVPR2024W/MAR/papers/Ma_Task_Navigator_Decomposing_Complex_Tasks_for_Multi...
- [6] Yao, S. et al. (2022). *ReAct: Synergizing Reasoning and Acting in Language Models*. arXiv preprint arXiv:2210.03629. (Conceptual reference for ReAct framework)
- [7] Lewis, P. et al. (2020). *Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks*. Advances in Neural Information Processing Systems (NeurIPS). (Conceptual reference for RAG)
- [8] Mikulski, B. (2025). *Comprehensive Guide to AI Workflow Design Patterns with LLMs*. <https://mikulskibartosz.name/ai-workflow-design-patterns>
- [9] Vellum AI. (2025). *Agentic Workflows in 2025: The ultimate guide*. <https://www.vellum.ai/blog/agentic-workflows-emerging-architectures-and-design-patterns>
- [10] Microsoft Azure. (2025). *AI Architecture Design*. Azure Architecture Center. <https://learn.microsoft.com/en-us/azure/architecture/ai-ml/>

- [11] Shinn, N. et al. (2023). *Reflexion: an autonomous agent with dynamic memory and self-reflection*. arXiv preprint arXiv:2303.11366. (Conceptual reference for Reflection)