# A Practical Guide for Designing, Developing, and Deploying Production-Grade Agentic AI Workflows

Eranga Bandara[a], Ross Gore[a], Peter Foytik[a], Sachin Shetty[a], Ravi Mukkamala[a], Abdul Rahman[b], Xueping Liang[c], Safdar H. Bouk[a], Amin Hass[g], Sachini Rajapakse[f], Ng Wee Keong[d], Kasun De Zoysa[e], Aruna Withanage[h], Nilaan Loganathan[h]

[a]*Old Dominion University, Norfolk, VA, USA*
[b]*Deloitte & Touche LLP, USA*
[c]*Florida International University, USA*
[d]*Nanyang Technological University, Singapore*
[e]*University of Colombo, Sri Lanka*
[f]*IcicleLabs.AI*
[g]*AnaletIQ, VA, USA*
[h]*Effectz.AI*

## Abstract

Agentic AI marks a major shift in how autonomous systems reason, plan, and execute multi-step tasks. Unlike traditional single model prompting, agentic workflows integrate multiple specialized agents with different Large Language Models(LLMs), tool-augmented capabilities, orchestration logic, and external system interactions to form dynamic pipelines capable of autonomous decision-making and action. As adoption accelerates across industry and research, organizations face a central challenge: how to design, engineer, and operate production-grade agentic AI workflows that are reliable, observable, maintainable, and aligned with safety and governance requirements. This paper provides a practical, end-to-end guide for design-

*Email addresses:* `cmedawer@odu.edu` (Eranga Bandara), `rgore@odu.edu` (Ross Gore), `pfoytik@odu.edu` (Peter Foytik), `sshetty@odu.edu` (Sachin Shetty), `mukka@odu.edu` (Ravi Mukkamala), `abdulrahman@deloitte.com` (Abdul Rahman), `xuliang@fiu.edu` (Xueping Liang), `sbouk@odu.edu` (Safdar H. Bouk), `aminhass@analetiq.com` (Amin Hass), `sachini.rajapakse@iciclelabs.ai` (Sachini Rajapakse), `awkng@ntu.edu.sg` (Ng Wee Keong), `kasun@ucsc.cmb.ac.lk` (Kasun De Zoysa), `aruna@effectz.ai` (Aruna Withanage), `nilaan@effectz.ai` (Nilaan Loganathan)

ing, developing, and deploying production-quality agentic AI systems. We introduce a structured engineering lifecycle encompassing workflow decomposition, multi-agent design patterns, Model Context Protocol(MCP), and tool integration, deterministic orchestration, Responsible-AI considerations, and environment-aware deployment strategies. We then present nine core best practices for engineering production-grade agentic AI workflows, including tool-first design over MCP, pure-function invocation, single-tool and single-responsibility agents, externalized prompt management, Responsible-AI-aligned model-consortium design, clean separation between workflow logic and MCP servers, containerized deployment for scalable operations, and adherence to the Keep it Simple, Stupid (KISS) principle to maintain simplicity and robustness. To demonstrate these principles in practice, we present a comprehensive case study: a multimodal news-analysis and media-generation workflow. By combining architectural guidance, operational patterns, and practical implementation insights, this paper offers a foundational reference to build robust, extensible, and production-ready agentic AI workflows.

*Keywords:* Agentic AI, Agentic AI Workflow, LLM, Model Context Protocol, Responsible AI

---

## 1. Introduction

The rapid advancement of Large Language Models (LLMs) [1, 2], Vision–Language Models (VLMs) [3, 4, 5], and tool-augmented reasoning has laid the foundation for a new paradigm in automation: agentic AI [6, 7]. Traditional LLM interactions follow a simple pattern in which a human provides a prompt and the model generates a response (as illustrated in the top half of Figure 1). In contrast, an AI agent can perform this same interaction autonomously: it can construct prompts, call models, interpret responses, and perform follow-up actions without direct human intervention (as illustrated in the bottom half of Figure 1). In essence, AI agents are software programs that use LLMs—together with tools, APIs, and external context—to execute tasks automatically. When multiple such agents collaborate—each with a specialized role such as searching, filtering, scraping, reasoning, validating, or publishing—they form agentic AI workflows [7]. These workflows enable systems that can reason, plan, take actions, monitor outcomes, retry intelligently, and iteratively refine their behavior. Modern agentic workflows integrate LLMs with external tools, structured memory,

search functions, databases, Model Context Protocol (MCP) servers, cloud services, and API-driven environments [8]. Rather than relying on a single monolithic prompt, the system delegates responsibilities across specialized agents to ensure modularity, determinism, and maintainability. This evolution has produced dynamic extensible pipelines capable of solving high-value real-world automation problems—from content generation and news analytics, to regulatory compliance, knowledge extraction, reasoning pipelines, and multimodal media synthesis [9]. By incorporating heterogeneous models (e.g., OpenAI, Gemini, Llama, Anthropic), deterministic tool calls, and environment-aware orchestration, agentic AI provides a flexible and powerful architectural approach for building scalable and production-grade AI systems.
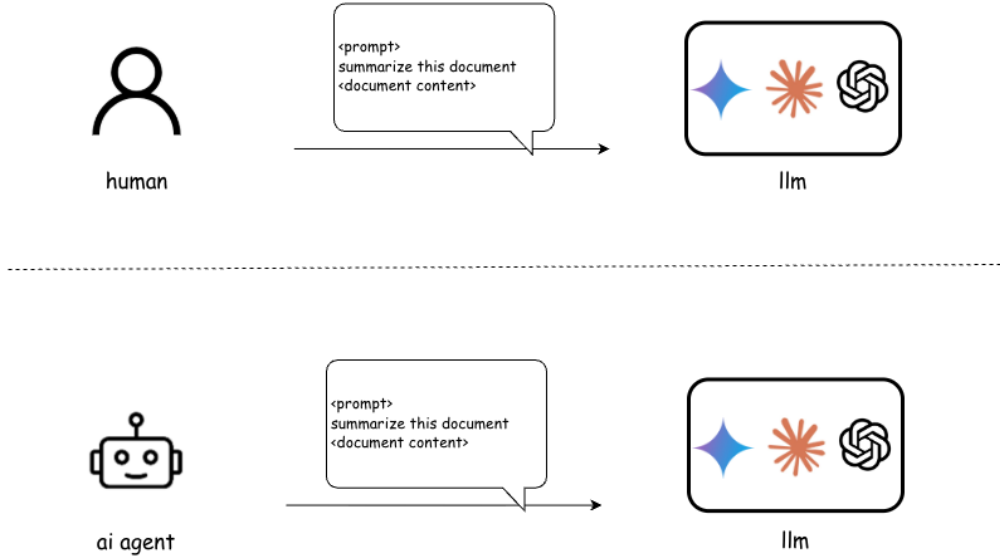


Figure 1: Human–LLM interaction versus autonomous AI agent–LLM interaction.

However, building production-grade agentic AI workflows remains challenging. While prototypes are easy to build with simple scripts or notebooks, scaling them into reliable, governed, and observable systems introduces multiple engineering complexities. In particular, design challenges include issues such as ways to decompose workflows into agents, making choices between tool calls and MCP actions [10, 11], ways to design deterministic orches-

tration, and methods to avoid implicit behaviors that lead to LLM drift or unpredictable execution paths [12]. The implementation challenges include managing multi-agent communication, handling tool schemas, maintaining prompt modularity, integrating heterogeneous model families, and enforcing responsible AI principles while ensuring output consistency. The operational challenges include running workflows reliably in production environments; managing concurrency, failures, retries, logging, and cost efficiency; securing tool access; monitoring agent traces; and ensuring reproducibility across model updates. Finally, the deployment challenges include containerizing complex multi-model systems, integrating with Kubernetes, managing versioning, exposing services through APIs and MCP servers, and supporting continuous delivery of agent updates and prompt modifications [13, 14].

As organizations move from experimentation to operational deployment, these challenges frequently result in brittle systems, inconsistent outputs, and high maintenance overhead. Without a disciplined engineering approach, agentic workflows can easily grow into opaque, unbounded, and error-prone pipelines that are difficult to debug, scale, or govern [12]. This paper addresses these gaps by presenting a practical, engineering-focused guide for designing, developing, and deploying production-grade agentic AI workflows. Drawing from real-world implementations and large-scale deployments, we outline a structured methodology for workflow decomposition, agent specialization, tool integration, safety mechanisms, orchestration strategies, and environment-aware deployment. We complement these principles with nine actionable best practices aimed at improving determinism, reliability, observability, and maintainability. Finally, we demonstrate these principles in action through a fully implemented case study—a multimodal news-analysis and media-generation workflow that autonomously scrapes web content, filters topics, generates podcast scripts through an agent consortium [15], consolidates reasoning output, and produces audio and video artifacts before publishing them to GitHub. This example showcases how heterogeneous agents can be composed into a unified pipeline and deployed reliably in production. The following are our main contributions of this research.

1. **A generalized engineering framework for production-grade agentic AI workflows.** We introduce a structured methodology for designing, developing, and deploying agentic systems using multi-agent orchestration, tool integration, and deterministic execution patterns suitable for real-world automation.

2. **A curated set of nine best practices for reliable and responsible–AI–enabled workflow design.** These practices encompass tool-over-MCP design, pure-function invocation, single-tool and single-responsibility agents, externalized prompt management, consortium-based reasoning for Responsible AI, separation of workflow logic and MCP servers, containerized deployment for scalability, maintainable workflow architecture, and adherence to the KISS principle for simplicity and robustness.

3. **A full implementation of a multimodal, multi-agent news-to-media workflow.** We present a complete case study demonstrating how the best practices are applied to a complex real-world workflow that performs feed discovery, topic filtering, content extraction, multi-LLM script generation, reasoning-based consolidation, audio/video synthesis, and automated GitHub publishing.

4. **An extensible blueprint for organizations adopting agentic AI in production.** By integrating containerization, Kubernetes orchestration, MCP accessibility, and Responsible-AI mechanisms (bias mitigation, reasoning audits, deterministic operations), the proposed system offers a robust template adaptable across domains such as compliance automation, media generation, analytics, and enterprise RPA.

The remainder of the paper is organized as follows. Section 2 presents the motivating use case, showcasing an end-to-end agentic AI workflow for multimodal news analysis, content synthesis, and media generation. Section 3 builds upon this use case by introducing a curated set of best practices for designing, developing, and deploying production-grade agentic AI workflows, with emphasis on architectural choices, orchestration patterns, tooling strategies, and Responsible-AI–aligned design principles. Section 4 details the implementation of the proposed workflow—including agent design, tool/function integration, reasoning-based consolidation, multimodal media generation, and deployment in a containerized production environment. Section 5 evaluates the system, assessing the performance and accuracy of individual agents and the effectiveness of the reasoning-based consolidation across models. Finally, Section 6 concludes the paper by summarizing key insights and outlining directions for future research, standardization, and broader adoption of robust and trustworthy agentic AI systems.

## 2. Use case: Podcast-Generation Agentic AI Workflow

To illustrate the practical value and architectural complexity of production-grade agentic AI systems, this paper uses a real-world workflow for automated podcast generation from live news sources. The workflow integrates heterogeneous LLM agents, web-scraping tools, multimodal generation models, and a GitHub publishing pipeline into a fully autonomous end-to-end system. Figure 2 shows the complete architecture.
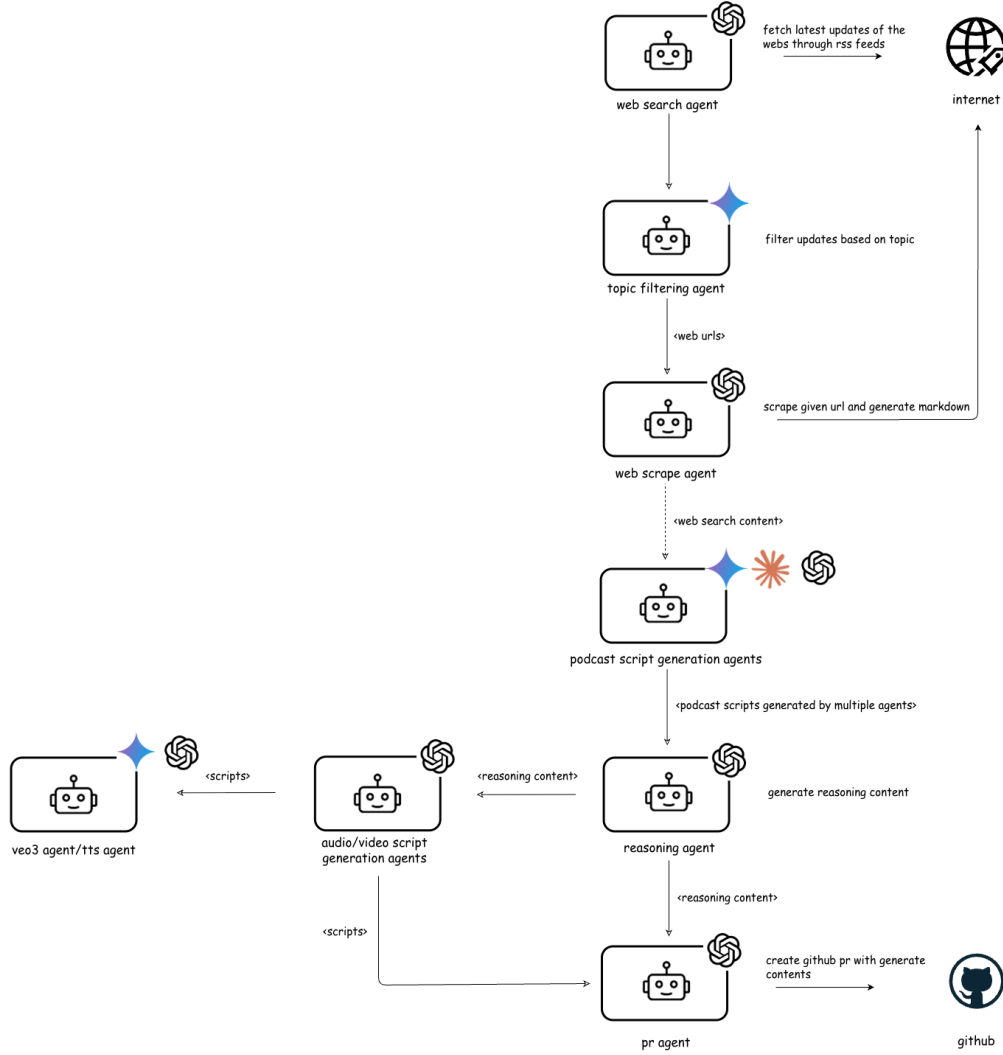
Figure 2: End-to-end agentic AI workflow for multimodal podcast generation.

At a high level, the podcast-generation workflow automates the entire life-cycle of transforming real-time news content into multimodal media assets. Given a topic and one or more source websites, the system autonomously discovers relevant news items, extracts, and summarizes their content, generates podcast scripts through a consortium of LLM agents, consolidates the output via a reasoning agent [16, 17, 18], and finally produces both audio and video podcast artifacts. The workflow not only synthesizes coherent narrative con-

tent from heterogeneous web sources, but also converts it into MP3 audio, MP4 video, and structured Veo-3 prompts—before packaging and publishing the results as a GitHub pull request [19]. This end-to-end automation showcases how an agentic AI pipeline can bridge web retrieval, content generation, multimodal synthesis, and software operations automation under a unified orchestration layer.

The workflow begins when the user provides a topic and one or more source URLs. These inputs trigger a pipeline of coordinated agents, each responsible for a well-defined function. The process starts with the Web Search Agent, which collects the latest updates from the internet by querying RSS feeds and MCP-based search endpoints [20]. Its output—a list of recent news articles—is then passed to the Topic Filtering Agent, which evaluates the relevance of each article to the user-specified topic and returns only the filtered URLs.

Next, the Web Scrape Agent extracts the full content of the web page for each selected URL. Rather than returning raw HTML, this agent converts each page into a clean, tool-generated Markdown representation, ensuring that downstream agents operate on consistent, structured text. The collected Markdown content serves as the primary knowledge base for the subsequent generation stack.

The system then invokes a consortium of Podcast Script Generation Agents [9], each powered by a different LLM provider (e.g., OpenAI, Gemini, Anthropic) [2, 21, 22]. These agents independently produce podcast scripts that describe the topic based on the scraped content. Because each model has different strengths, biases, and reasoning styles, their outputs form a diverse set of drafts.

To obtain a coherent and reliable final script, the system employs a Reasoning Agent, which consolidates the drafts by comparing them, resolving inconsistencies, removing speculative claims, and synthesizing a unified podcast narrative. This reasoning step not only improves accuracy, but also enforces responsible-AI constraints by grounding the final script strictly in the content extracted from the web [23].

With the consolidated script prepared, the workflow branches into multimodal output generation. Audio/Video Script Generation Agents transform the consolidated narrative into structured prompts suitable for text-to-speech (TTS) [24] and text-to-video (Veo-3) models [19]. A dedicated Veo Agent converts the video script into JSON-based VEO-3 instructions, while the TTS generator produces high-quality spoken audio from the podcast text.

Finally, all generated assets—including the consolidated script, individual agent outputs, audio files, video files, and Veo-3 instructions—are assembled by the PR Agent, which publishes them to GitHub via an MCP-integrated GitHub server [10]. The agent creates a new branch, commits the results, and automatically opens a pull request, completing the fully autonomous media-generation pipeline.

This use case demonstrates how agentic AI systems can orchestrate web-scale data acquisition, multi-agent LLM reasoning, multimodal content creation, and automated software operations within a single cohesive pipeline. It also illustrates why structured engineering practices are essential: each component of the workflow must be deterministic, auditable, observable, and robust enough to operate reliably in production environments. In the next section, we show how the workflow is enhanced by applying a curated set of best practices—transforming it into a maintainable, scalable, and responsible–AI–aligned agentic system suitable for real-world deployment [23].

## 3. Best Practices and Their Application

In this section, we discuss nine core best practices for designing, developing, and deploying production-grade agentic AI workflows. We illustrate how each practice strengthens the reliability, maintainability, scalability, and Responsible-AI characteristics of the system. Using the podcast-generation workflow as a running example, we show how applying these practices transforms the pipeline into a robust agentic architecture suitable for real-world deployment.

### 3.1. Tool Calls Over MCP

AI agents can integrate with external systems through MCP connections or through direct function calls (tool calls) [25]. MCP provides a standardized mechanism for structured communication between agents and external services, replacing many ad-hoc APIs with a unified interaction model. However, MCP integration also introduces additional layers of abstraction that can sometimes reduce determinism, complicate agent reasoning, or create ambiguous tool-selection behaviors. Further, it adds more complexity to configuring the MCP servers, etc, inside the workflow (see Figure 3).

In our workflow, the initial implementation relied on the GitHub MCP server to create pull requests for the generated podcast scripts. However, during evaluation, we observed recurring issues: the agent frequently made

ambiguous tool-selection decisions, inconsistently inferred invocation parameters, and occasionally failed with non-deterministic MCP responses (see Figure 3 and Figure 4). These challenges arose because the agent had to interpret multiple definitions of the MCP tool and reason through the metadata structure of the protocol, which increased cognitive load and introduced variability into the workflow. Although we repeatedly refined and adjusted the agent instructions to mitigate these issues, the behavior remained unstable and exhibited flickering, non-reproducible failures.

To address this, we replaced the GitHub MCP integration with a direct pull-request creation function that agents invoke explicitly (see Figure 5). This eliminated ambiguity, improved determinism, and ensured that the final step of the workflow—publishing the results to GitHub—was stable and predictable. The workflow became easier to debug, more auditable, and significantly more reliable in production environments. Further improvements are discussed in the following sections, including minimizing tool-set complexity and using pure function calls to reduce token overhead and inference variability.

```
 1  github_pr_agent = Agent(
 2      name="Pull Request Creation Agent",
 3      instructions=(
 4          "Create file with given file name and content(don't create file as base64 encoded) and push to given branch, if given branch
                does not exists create it from default branch first"
 5          "Then create a pull request, ensure the file is committed, pushed to the correct branch, and the pull request is opened
                targeting the repository's default branch. "
 6          "If pull request creation fails due to the error Error invoking MCP tool get_file_contents: Not Found: Resource not found: No
                commit found, recover by proceeding to create the branch and continue with pull request creation. Do not assume that the
                branch must have existing files or commits beforehand. "
 7      ),
 8  )
 9
10 ▾ async def main():
11      # configure google mcp server and github mcp server
12      async with MCPServerStdio(
13 ▾        params={
14              "command": "node",
15              "args": [GOOGLE_SEARCH_MCP_SERVER_PATH],
16 ▾            "env": {
17                  "GOOGLE_API_KEY": GOOGLE_API_KEY,
18                  "GOOGLE_SEARCH_ENGINE_ID": GOOGLE_CSE_ID  # make sure this matches your .env key
19              },
20          }
21      ) as google_search_server, MCPServerStdio(
22          cache_tools_list=True,
23 ▾        params={
24              "command": GITHUB_MCP_SERVER_PATH,
25              "args": ["stdio"],
26              "env": {"GITHUB_PERSONAL_ACCESS_TOKEN": GITHUB_TOKEN}
27          }
28 ▾    ) as github_server:
29          google_search_agent.mcp_servers = [google_search_server]
30          read_markdown_agent.mcp_servers = [google_search_server]
31          github_pr_agent.mcp_servers = [github_server]
32
33          # 1.web search throug mcp
34          # 2. scrape content through mcp
35          # 3. generate content
36
37          # 4. create PR with content with mcp
38          repo = GITHUB_REPO
39          branch = "podcast-v1"
40          file = "podcast-script.md"
41          input_text = (
42              f"Create a pull request in repository: '{repo}' "
43              f"branch: '{branch}', file name: '{file}', file content: \n{content}"
44          )
45          pr_result = await Runner.run(starting_agent=github_pr_agent, input=compliance_input_text)
```

Figure 3: Workflow integrated with an MCP server, illustrating the operational overhead of configuring and managing multiple MCP servers.

11

```
49  # observed error 1
50  Note: File does not exist in branch, will create new file Error invoking MCP tool create_or_update_file: Not Found: Resource not
        found: Branch podcast-v1 not found
51
52  # observed error 2
53  agents.exceptions.AgentsException: Error invoking MCP tool get_file_contents: Not Found: Resource not found: No commit found for the
        ref podcast-v1
54
55  # instrction of the agent which used to fix the error
56  instructions=(
57      "Create file with given file name and content(don't create file as base64 encoded) and push to given branch, if given branch does
            not exists create it from default branch first"
58      "Then create a pull request, ensure the file is committed, pushed to the correct branch, and the pull request is opened targeting
            the repository's default branch. "
59      "If pull request creation fails due to the error Error invoking MCP tool get_file_contents: Not Found: Resource not found: No
            commit found, recover by proceeding to create the branch and continue with pull request creation. Do not assume that the
            branch must have existing files or commits beforehand. "
60  ),
```

Figure 4: Failure cases observed in the MCP-integrated workflow.

```
1   # Tool: GitHub PR creation
2   @function_tool
3   def create_github_pr(repo: str, branch: str, file_name: str, file_content: str, token: str = None,
4       commit_message: str = "Automated PR from agentic workflow", base_branch: str = "master", pr_title: Optional[str] = None,
5 ▾ ):
6       """
7       Create a file in a new or existing branch and open a Pull Request in GitHub.
8       """
9       # implementation of pull request creation
10
11
12  # Agent: dedicated PR creation agent that must use the tool
13  github_pr_agent = Agent(
14      name="Pull Request Creation Agent",
15      instructions=(
16          "You are responsible for creating GitHub pull requests. Always use the `create_github_pr` tool to create or update \n"
17          " a file in a branch and open a pull request. Do not invent values; use exactly the repo, branch file_name, and \n"
18          "file_content provided in the user input."
19      ),
20      tools=[create_github_pr],
21  )
22
23
24  # Orchestration: invoke the agent with explicit parameters
25  repo = "erangaeb/agentic-workflows"
26  branch = "podcast-v1"
27  file_name = "podcast-script.md"
28
29  input_text = (
30      f"Create a pull request with the following parameters: - repo: {repo}, - branch: {branch} - file_name: {file_name}\n"
31      f"- file_content:\n{content}\n"
32  )
33
34  pr_result = await Runner.run(starting_agent=github_pr_agent, input=input_text)
```

Figure 5: Replacing MCP-integrated workflow with direct tool invocation.

## 3.2. Direct Function Calls Over Tool Calls

Although tool calls provide a structured way for agents to interact with external systems, they introduce additional overhead and potential ambiguity. Every tool invocation requires the LLM to parse instructions, interpret parameter formats, and map natural language input to function

arguments—steps that increase token consumption and can lead to non-deterministic behavior [26]. In complex workflows, even carefully designed tools can produce unpredictable results if the agent misinterprets parameter names, defaults, or expected data structures.

For operations that do not require language reasoning, such as posting data to an API, committing a file to GitHub, performing database writes, or generating timestamps, tool calls are often unnecessary. Instead, these steps can be handled directly in the orchestration layer using pure function calls. Pure functions—functions executed directly by the workflow without involving the LLM—are deterministic, side-effect controlled, cheaper, faster, and fully testable [14].

In our workflow, the GitHub pull request step was originally based on a dedicated "PR Agent" that used the create tool _github _pr. Despite improving determinism relative to MCP, this still required the agent to reason about tool parameters and produce a structured call [27]. We eventually completely removed the agent for this operation and invoked create_github_pr directly from the workflow controller. This eliminated ambiguous tool formatting, removed unnecessary LLM reasoning, reduced token usage, and made the system significantly more stable (see Figure 6).

By shifting infrastructure-oriented tasks to pure functions and reserving tool calls for operations where language-driven reasoning is genuinely necessary, agentic workflows become simpler, more predictable, and better aligned with production software engineering principles.

```
 1  # Function: GitHub PR creation
 2  def create_github_pr(
 3      repo: str,
 4      branch: str,
 5      file_name: str,
 6      file_content: str,
 7      token: str = None,
 8      commit_message: str = "Automated PR from agentic workflow",
 9      base_branch: str = "master",
10      pr_title: Optional[str] = None,
11 ▾ ):
12      """
13      Create a file in a new or existing branch and open a Pull Request in GitHub. Requires a GitHub token with `repo` scope.
14      """
15      # implementation of pull request creation
16      ...
17
18
19  # Workflow orchestration (no agent involved in PR creation)
20 ▾ async def main():
21      # 1. Search feeds via direct function call
22      # 2. Scrape content via direct function call
23      # 3. Generate podcast content via agents
24
25      # 4. Create PR via direct function call with explicit parameters
26      timestamp = datetime.now().strftime("%Y%m%d-%H%M%S")
27      branch = f"podcast-{timestamp}"
28      file_name = f"podcast-{timestamp}.md"
29      commit_msg = f"Add podcast for topic: {topic}"
30      pr_title = f"Podcast: {topic} ({timestamp})"
31
32      pr_response = create_github_pr(repo=repo, branch=branch, file_name=file_name, file_content=podcast_pr_content,
33          token=GITHUB_TOKEN, commit_message=commit_msg, pr_title=pr_title,)
```

Figure 6: Direct function invocation replacing agent-based tool calls to reduce ambiguity and overhead.

## 3.3. Avoid Overloading Agents With Many Tools

Attaching multiple tools to a single agent increases prompt complexity and reduces reliability. When an agent is equipped with several tools, the model must first reason about which tool to invoke and how to structure the parameters—introducing unnecessary ambiguity and increasing the likelihood of incorrect or missing tool calls [26]. This cognitive overhead results in higher token usage, poorer accuracy, and inconsistent execution paths.

A more robust approach is to follow a "one agent, one tool" design whenever tool usage is required. The assignment of a single well-defined tool to each agent creates predictable roles, simplifies prompting, and eliminates tool-selection noise, allowing the agent to focus solely on parameter inference and execution [14]. This modular decomposition improves interpretability, eases debugging, and makes the workflow easier to scale and reuse across contexts.

In our workflow, we initially designed a single agent that used two tools: scrape_markdown and publish_markdown (Figure 7). The intention was to scrape the content of the webpage and then publish the extracted markdown

14

to third-party storage for audit purposes. However, during evaluation, we observed that the agent would often invoke only one tool, invoke them in the wrong order, or fail to call them entirely, especially when the prompt or input size increased [28, 29]. To address this, we decompose the design into two independent agents, each responsible for exactly one tool (Figure 8). This separation ensured deterministic behavior, eliminated missed tool calls, and significantly improved stability over repeated runs.

```
 1  @function_tool
 2  def extract_content_as_markdown(url: str):
 3      # Extract webpage content and convert it into markdown
 4      ...
 5
 6
 7  @function_tool
 8  def publish_content(content: str):
 9      # Publish extracted markdown to third-party storage for auditing
10      ...
11
12
13  aget = Agent(
14      name="Content Extraction and Publishing Agent",
15      instructions=(
16          "Your task has two steps:\n"
17          "1. Use the `extract_content_as_markdown` tool to extract markdown from the given URL.\n"
18          "2. Then publish the extracted markdown using the `publish_content` tool.\n\n"
19          "You must return the extracted markdown as your final output. Only use the tools provided. Do not invent values or call tools
                not listed here."
20      ),
21      model="gpt-5",
22      tools=[extract_content_as_markdown, publish_content],
23      output_type=str,
24  )
```

Figure 7: A single agent overloaded with multiple tool integrations.

15

```
 1  @function_tool
 2 ▾ def extract_content_as_markdown(url: str):
 3      # Extract webpage content and convert it into markdown
 4      ...
 5
 6
 7  @function_tool
 8 ▾ def publish_content(content: str):
 9      # Publish extracted markdown to third-party storage for auditing
10      ...
11
12
13  extract_agent = Agent(
14      name="Markdown Extraction Agent",
15      instructions=(
16          "You are responsible only for extracting webpage content. Use the `extract_content_as_markdown` tool to fetch the page\n"
17          "at the  given URL and return the extracted markdown as your final output.\n\n"
18          "Steps:\n"
19          "1. Call `extract_content_as_markdown` with the provided URL.\n"
20          "2. Do not call any other tools.\n"
21          "3. Return exactly the markdown you obtained, without adding extra commentary."
22      ),
23      model="gpt-5",
24      tools=[extract_content_as_markdown],
25      model_settings=ModelSettings(temperature=0.0),
26      output_type=str,
27  )
28
29
30  publish_agent = Agent(
31      name="Markdown Publishing Agent",
32      instructions=(
33          "You are responsible only for publishing markdown content. Use the `publish_content` tool to store the provided markdown\n"
34          "in third-party storage for auditing.\n\n"
35          "Steps:\n"
36          "1. Receive markdown text as input.\n"
37          "2. Call `publish_content` with that markdown.\n"
38          "3. Return a short status message indicating whether publishing succeeded."
39      ),
40      model="gpt-5",
41      tools=[publish_content],
42      model_settings=ModelSettings(temperature=0.0),
43      output_type=str,
44  )
```

Figure 8: Multiple agents, each assigned a single tool for clearer and more deterministic operation.

## 3.4. Single-Responsibility Agents

Closely related to avoiding multiple tools per agent is the principle of single-responsibility agents [30]. Just as good software design favors functions and classes that "do one thing well", production-grade agentic workflows benefit when each agent is responsible for a single, clearly defined task. When an agent is asked to handle multiple conceptual responsibilities—such as generation, validation, transformation, and side=effective actions—it becomes harder to prompt, harder to test, and more prone to subtle, non-deterministic failures [12].

In our workflow, an early design attempted to combine Veo-3 JSON prompt generation and video generation into a single agent. The agent received the consolidated podcast or video script as input and was instructed

to both (1) transform the script into a structured Veo-3 JSON specification and (2) "generate" the corresponding video [13, 19]. In practice, this blurred the boundary between planning (designing the video prompt) and execution (calling the Veo API and handling files), as shown in Figure 9. The LLM would sometimes produce malformed JSON, sometimes mix natural language with JSON, and sometimes "hallucinate" file paths or status messages about video generation that had not actually occurred.

```python
 1  @function_tool
 2  def script_to_video(veo_json_str: str, filename: str = "update.mp4") -> str:
 3      """
 4      Accepts a JSON string (Veo structured prompt), generates an 8s Veo 3 video, saves to filename.
 5      """
 6      ...
 7
 8
 9  agent = Agent(
10      name="Vio Video Generation Agent",
11      instructions=(
12          "You receive a finalized video script and must perform following tasks:\n"
13          "1. Convert the script into a valid Veo-3 JSON prompt.\n"
14          "2. Use the `script_to_video` tool to generate the video from that JSON.\n\n"
15          "Return BOTH of the following in your final answer:\n"
16          "- The Veo-3 JSON you used.\n"
17          "- The status or path returned by `script_to_video`.\n\n"
18          "You may think step by step, but the final answer must include the JSON and the video status."
19      ),
20      model="gpt-5",
21      tools=[script_to_video],
22      model_settings=ModelSettings(temperature=0.7),
23      output_type=str,
24  )
```

Figure 9: A single agent responsible for multiple internal functions, increasing complexity and ambiguity.

We resolved this by splitting the responsibilities: A Veo JSON Builder Agent whose only job is to take the final script and produce strictly valid Veo-3 JSON describing the video (scenes, timing, style, etc.). A separate, non-agent video generation function (script_to_video) that receives the JSON and interacts with the Veo-3 API to generate and store the MP4 file (see Figure 10). This single-responsibility design yields several benefits: prompts become simpler and more focused; the agent's output contract is clearer (always "valid Veo JSON" rather than "JSON plus narrative status"); and side effects (calling external APIs, handling retries, saving files) are handled deterministically in regular code. More generally, enforcing single-responsibility agents across the workflow makes the system easier to reason about, easier to debug, and safer to evolve as new capabilities are added [31].

17

```
1 ▾ def script_to_video(veo_json_str: str, filename: str = "update.mp4") -> str:
2       """
3       Accepts a JSON string (Veo structured prompt), generates an 8s Veo 3 video, saves to filename.
4       """
5       ...
6
7
8   veo_json_builder_agent = Agent(
9       name="Veo JSON Builder Agent",
10      instructions=load_prompt("veo_json_builder_agent.txt"),
11      model=gemini_model,
12      model_settings=ModelSettings(temperature=0.1),
13      output_type=str,
14  )
15
16
17 ▾ async def main():
18      # 1. generate veo3 json script
19      veo_prompt_input = (
20          f"TOPIC: {topic}\n\n"
21          "## Video Script\n"
22          f"{video_script_response.final_output}"
23      )
24      veo_json_response = await Runner.run(starting_agent=veo_json_builder_agent, input=veo_prompt_input,)
25
26      # 2. convert script to video
27      script_to_video(veo_json_response.final_output, filename=video_filename,)
```

Figure 10: Decomposition into multiple agents, each assigned a specific function based on the single-responsibility principle.

## 3.5. Store Prompts Externally and Load Them at Runtime

Embedding prompts directly within source code creates tight coupling, complicates version control, and restricts collaboration. A more flexible and maintainable approach is to store prompts as external artifacts—typically in Markdown or plain text—within a separate storage location, such as a GitHub repository, shared drive, or configuration service. Externalizing prompts enables non-technical stakeholders (e.g., policy teams, domain experts, content reviewers) to update and refine agent behavior without modifying application code [12].

In our workflow, all agent prompts are stored in a dedicated GitHub repository and dynamically loaded at runtime (see Figure 11). This decoupling allows collaborators to iterate on the language, constraints, and safety requirements of each agent independently of workflow deployment cycles. It also enables governance workflows such as review processes, version pinning, rollback, and controlled access.

Maintaining prompts externally further supports continuous improvement practices, including A/B testing, prompt red-teaming, and evolving

Responsible-AI rules, all without requiring code redeployments. This separation of concerns significantly enhances maintainability, transparency, and organizational agility when operating agentic AI systems in production [29].

```python
1   topic_filtering_agent = Agent(
2       name="TopicFilteringAgent",
3       instructions=load_prompt("topic_filtering_agent.txt"),
4       model=gemini_model,
5       model_settings=ModelSettings(temperature=0.1),
6       output_type=List[str],
7   )
8
9   read_markdown_agent = Agent(
10      name="ReadWebpageMarkdownAgent",
11      tools=[extract_content_as_markdown],
12      instructions=load_prompt("read_markdown_agent.txt"),
13      model=OPENAI_MODEL_NAME,
14      model_settings=ModelSettings(temperature=1),
15      output_type=str,
16  )
17
18  gemini_podcast_generation_agent = Agent(
19      name="PodcastGenerationAgentGemini",
20      instructions=load_prompt("podcast_generation_agent.txt"),
21      model=gemini_model,
22      model_settings=ModelSettings(temperature=0.1),
23      output_type=str,
24  )
```

Figure 11: Loading externalized prompts dynamically at runtime.

## 3.6. Responsible AI Agents

Single-model outputs often suffer from well-known limitations such as hallucinations, reasoning inconsistencies, and subtle or overt biases. To address these challenges, production-grade agentic workflows benefit from a multi-model consortium architecture, where several specialized LLMs (e.g., Gemini, GPT, Claude, Llama, Pixtral, Qwen) [32, 4, 5] independently generate outputs that are later synthesized by a dedicated reasoning agent [23, 33]. This design strengthens the workflow in several key ways: Higher accuracy through cross-model agreement; reduced bias by incorporating diverse model behaviors and training distributions; greater robustness to model updates or drift; and better alignment with Responsible AI principles, especially accountability and verifiability.

19

The reasoning agent functions as the final auditor in the pipeline. Rather than creating new content from scratch, it performs structured consolidation tasks such as conflict resolution, logical consistency checking, factual alignment, deduplication, and relevance filtering. Its role is to critically evaluate the drafts produced by individual model agents and produce a harmonized, trustworthy final output [9, 17].

In our workflow, Responsible-AI behavior is realized by combining multiple LLMs in parallel and routing their output through a dedicated reasoning LLM (e.g., OpenAI GPT-oss) [33, 17, 18], as illustrated in Figure 12. Each agent—whether tasked with planning, content generation, or validation—can interface with this model consortium to obtain richer, multi-perspective responses. For example, during podcast generation, the workflow collects draft scripts from agents based on Anthropic Claude, OpenAI GPT-5, and Gemini [2, 21, 22], Figure 13. The reasoning agent then synthesizes these drafts to produce a responsible and well-structured final script that reflects consensus rather than the idiosyncrasies of any single model.

This ensemble-based reasoning mechanism improves transparency, mitigates risk, and increases the reliability of agentic outputs—providing a solid foundation for building Responsible-AI-aligned workflows suitable for production deployment.
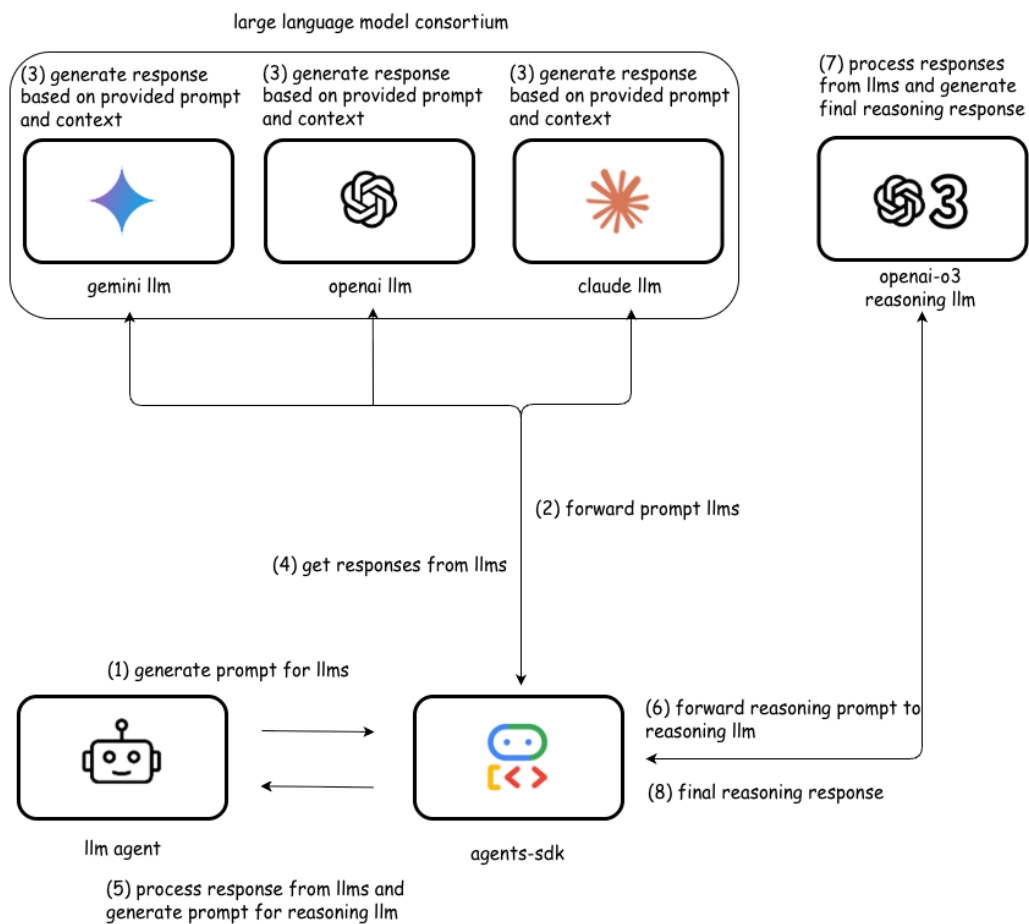
large language model consortium

(3) generate response
based on provided prompt
and context

(3) generate response
based on provided prompt
and context

(3) generate response
based on provided prompt
and context

(7) process responses
from llms and generate
final reasoning response

gemini llm

openai llm

claude llm

openai-o3
reasoning llm

(2) forward prompt llms

(4) get responses from llms

(1) generate prompt for llms

(6) forward reasoning prompt to
reasoning llm

(8) final reasoning response

llm agent

agents-sdk

(5) process response from llms and
generate prompt for reasoning llm

Figure 12: Integration flow of the LLM consortium with the reasoning LLM

```
 1  gemini_podcast_generation_agent = Agent(
 2      name="PodcastGenerationAgentGemini",
 3      instructions=load_prompt("podcast_generation_agent.txt"),
 4      model=gemini_model,
 5      model_settings=ModelSettings(temperature=0.1),
 6      output_type=str,
 7  )
 8
 9  openai_podcast_generation_agent = Agent(
10      name="PodcastGenerationAgentOpenAI",
11      instructions=load_prompt("podcast_generation_agent.txt"),
12      model=OPENAI_MODEL_NAME,
13      model_settings=ModelSettings(temperature=1),
14      output_type=str,
15  )
16
17  anthropic_podcast_generation_agent = Agent(
18      name="PodcastGenerationAgentAnthropic",
19      instructions=load_prompt("podcast_generation_agent.txt"),
20      model=llama_model,
21      model_settings=ModelSettings(temperature=1),
22      output_type=str,
23  )
24
25  reasoning_agent = Agent(
26      name="PodcastGenerationReasoningAgent",
27      instructions=load_prompt("podcast_generation_reasoning_agent.txt"),
28      model=OPENAI_MODEL_NAME,
29      model_settings=ModelSettings(temperature=1),
30      output_type=str,
31  )
32
33 ▾ async def run_workflow(topic, days=7, urls=["https://www.cbsnews.com/"]):
34      # 1. generate podcast scripts through gemini, openai, anthropic agents
35
36      # 2. generate reasoning podcast script
37      reasoning_intro = (
38          "You are given draft podcast scripts produced by multiple models. "
39          "Compare them, resolve conflicts, and craft one consolidated, on-topic script "
40          "that keeps the strongest shared facts and drops speculative or conflicting bits. "
41          "Stay grounded in the drafts—no new facts."
42      )
43      reasoning_input = (
44          f"TOPIC: {topic}\n\n{reasoning_intro}\n\n" + "\n\n".join(reasoning_input_parts)
45      )
46      reasoning_response = await Runner.run(
47          starting_agent=reasoning_agent,
48          input=reasoning_input,
49      )
```

Figure 13: Coordinated operation of the agent consortium and the reasoning agent.

## 3.7. Separation of Agentic AI Workflow and MCP Server

It is a common practice to expose the Agentic AI workflows function through an MCP server, allowing any MCP-enabled client—such as Claude

Desktop, VS Code extensions, or LM Studio—to seamlessly integrate with the workflow [34]. To achieve this, the workflow itself should be served via a REST API, while the MCP server acts as a thin orchestration layer that simply forwards MCP tool calls to the underlying API [11]. A key design principle is decoupling the agentic workflow engine from the MCP server. Instead of embedding workflow logic inside the MCP server, the architecture cleanly separates the backend workflow logic (the agentic pipeline, multi-agent orchestration, and tool integrations), the MCP server (a lightweight adapter that exposes workflow endpoints as MCP tools [25]), and user-facing MCP clients (external tools that can invoke workflow capabilities without direct code integration). This separation improves maintainability, supports independent scaling of components, and ensures long-term adaptability as LLMs, tools, and API specifications evolve. It also keeps the MCP server simple, stable, and safe—while allowing the workflow backend to iterate rapidly. The resulting modular architecture is shown in Figure 14.
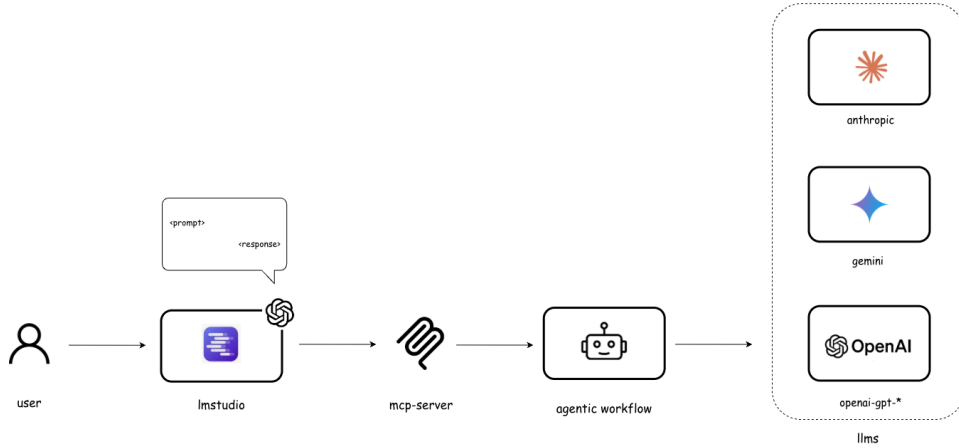


Figure 14: Exposing the agentic workflow functions through an MCP server.

### 3.8. Containerized Deployment

In production environments, agentic AI workflows and their accompanying MCP servers should be deployed using containerization technologies such as Docker and orchestrated with platforms like Kubernetes. Containerization provides a consistent and reproducible runtime environment that eliminates configuration drift and ensures that workflows behave identically in development, staging, and production environments [35].

By packaging the workflow engine, supporting tools, and the MCP server into isolated containers, organizations gain several operational advantages: 1) Portability: Containers encapsulate dependencies, model configurations, tool clients, and runtime libraries, enabling seamless deployment across any cloud or on-premise environment. 2) Scalability: Kubernetes can automatically scale agentic workflows based on load—spinning up additional replicas during high-traffic periods or downscaling to save resources. 3) Resilience: With built-in health checks, container restarts, and self-healing mechanisms, Kubernetes ensures that transient failures do not disrupt production workloads. 4) Security and Governance: Container isolation, network policies, secret management, and role-based access control (RBAC) enable robust security boundaries around AI workflows [27, 30]. 5) Observability: Kubernetes integrates well with logging, metrics, and tracing systems (e.g., Prometheus, Grafana, OpenTelemetry), providing visibility into multi-agent pipelines. 6) Continuous Delivery: Containers plug naturally into CI/CD pipelines, allowing predictable deployments and automated rollbacks [36].

In our implementation, both the agentic AI workflow and the MCP server were fully Dockerized and deployed to a Kubernetes cluster, Figure 15. This architecture allows for independent scaling of the workflow compute, reasoning components, scraper functions, and MCP interface nodes. It also enables blue–green deployments, canary releases, and safe iteration on individual components without impacting the entire system. Containerized deployment, therefore, provides the operational foundation required for stable, scalable, and production-grade agentic AI systems.
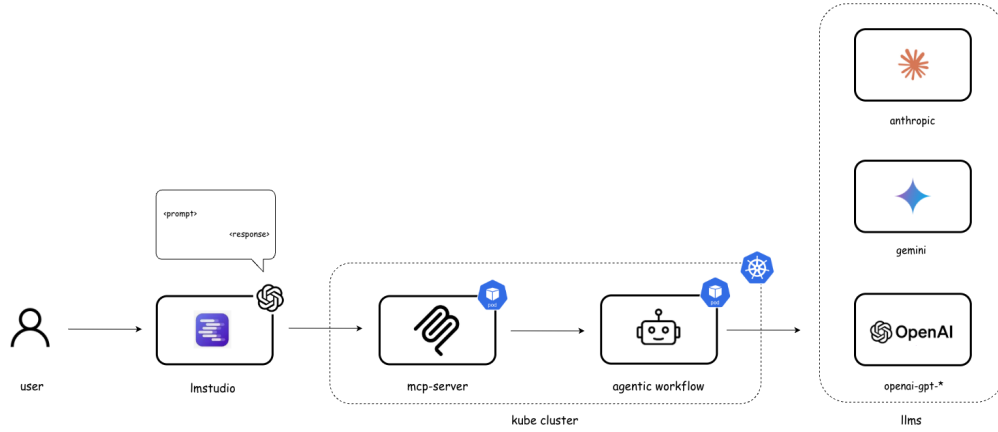
Figure 15: Deployment of the agentic workflow and MCP server on a Kubernetes cluster.

## 3.9. Keep it Simple, Stupid

Complexity is one of the biggest threats to the reliability and maintainability of agentic AI workflows. Unlike traditional enterprise software systems—where layered architecture, elaborate design patterns, and deeply nested abstractions are common—agentic workflows operate under a different paradigm. Their primary purpose is not to implement intricate logic internally, but rather to delegate reasoning, generation, and decision-making to LLMs and specialized agents. Therefore, applying the Keep It Simple and Stupid (KISS) principle is essential [14].

First, workflow implementations should avoid unnecessary structural complexity, over-engineering, or traditional architectural patterns that add little value in agentic systems. Introducing multiple indirections, deep inheritance, or microservice-like decomposition often leads to brittleness rather than clarity. Agentic workflows benefit far more from flat, readable, function-driven designs, where each component is responsible for a single task and the orchestration logic remains transparent and lightweight [14].

Second, simplicity improves reliability. Each layer of additional complexity introduces more opportunities for ambiguity in agent behavior, tool invocation mismatches, or unintended side effects. A simple workflow makes agent decision pathways clearer, reduces unexpected LLM behaviors, and improves the predictability of tool calls.

Third, KISS-aligned workflows integrate more naturally with modern AI-assisted development tools such as Claude Code, GitHub Copilot, or OpenAI

25

Codex [37]. These tools excel when the project structure is intuitive and the codebase is concise. A clean and simple workflow makes it significantly easier for AI-assisted coding tools to generate correct patches, refactor logic, or suggest improvements, accelerating iteration and reducing engineering overhead.

Finally, maintaining simplicity also supports long-term extensibility. By keeping workflow logic minimal and delegating most cognitive tasks to LLMs and agents, the system remains adaptable to new models, new tools, and new runtime environments without requiring architectural rewrites [26].

In summary, adhering to the KISS principle ensures that agentic AI workflows remain easy to understand, easy to debug, LLM-friendly, and scalable as the surrounding tooling ecosystem evolves.

## 4. Implementation

The proposed agentic AI workflow use case was fully implemented using the OpenAI Agents SDK [38], providing a structured and extensible foundation for multi-agent orchestration, tool integration, and deterministic function execution. The development process incorporated all optimization principles introduced in Section 3, including tool-first design, pure-function invocation, single-responsibility agents, externalized prompts, and the use of a multi-model reasoning consortium to ensure accuracy, reliability, and alignment of Responsible-AI [39].

To support seamless external integration, the workflow backend was exposed through a dedicated REST API, and a corresponding MCP server was developed to bridge MCP-enabled clients with the workflow [10]. The MCP server acts as a lightweight adapter layer, forwarding tool calls directly to the workflow API, thereby maintaining a clean separation between orchestration logic and communication interfaces. This separation also ensures long-term adaptability as models and tools evolve independently.

The complete implementation—including the workflow logic, agent definitions, tool functions, prompt templates, and deployment artifacts—has been released as open-source in dedicated GitLab repositories, enabling full reproducibility and community-driven refinement. The primary workflow implementation is available in the Podcast Workflow repository [40], while the corresponding MCP server implementation is published separately in the Podcast Workflow MCP Server repository [41]. Both projects include Dockerfiles and Kubernetes deployment manifests, supporting cloud-native,

production-grade deployment [42]. The workflow and MCP server were containerized using Docker and deployed via Kubernetes, providing scalability, workload isolation, high availability, and operational resilience suitable for enterprise environments.

To validate interoperability, the workflow's MCP server was integrated into LM Studio, enabling full end-to-end testing with local LLMs acting as MCP clients [34]. This integration confirmed that the workflow can be invoked, parameterized, and executed directly through an MCP-compatible interface. Figure 16 demonstrates how users provide natural-language requests that LM Studio routes to the MCP server. Figure 17 shows how LM Studio interprets these requests and automatically maps them to the workflow's input parameters. Finally, Figure 18 illustrates the successful invocation of the workflow function through the MCP server and the corresponding system response. Together, these interactions validate seamless end-to-end connectivity between LM Studio and the deployed agentic workflow.
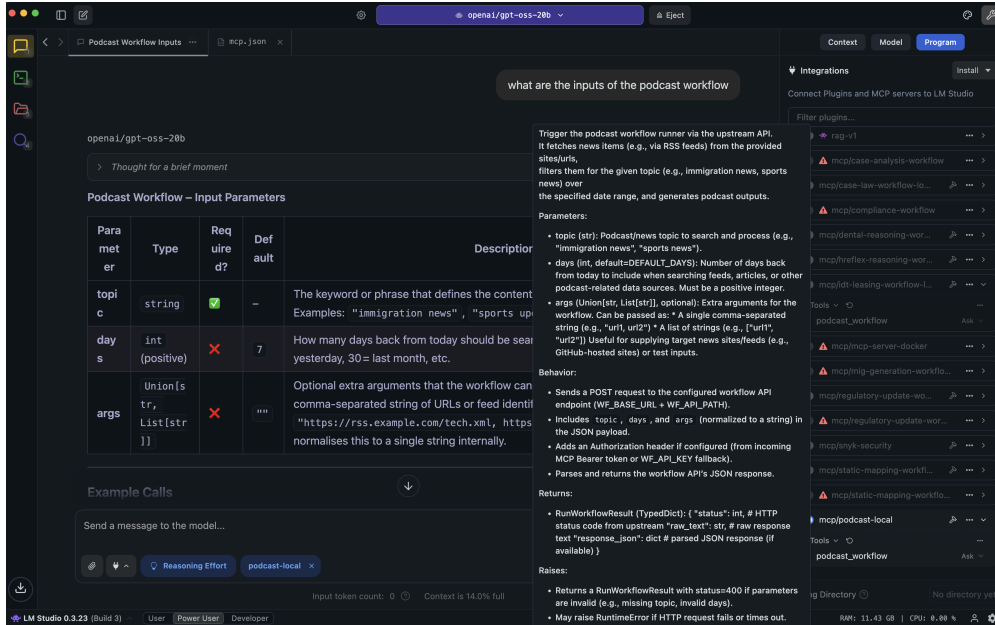


Figure 16: Integration of the workflow with LM Studio via the MCP server, showing how workflow inputs are provided.
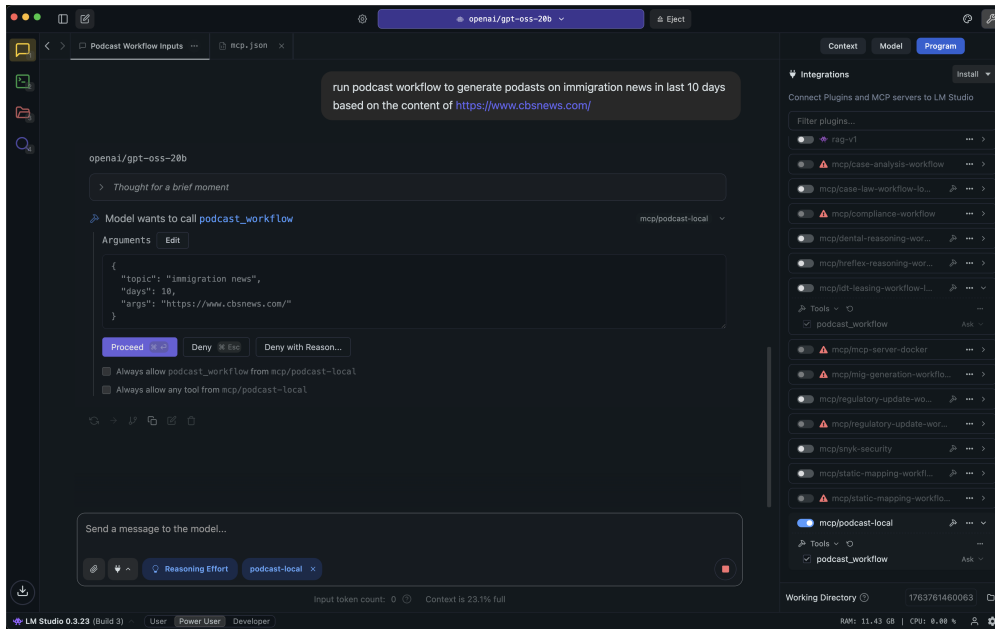
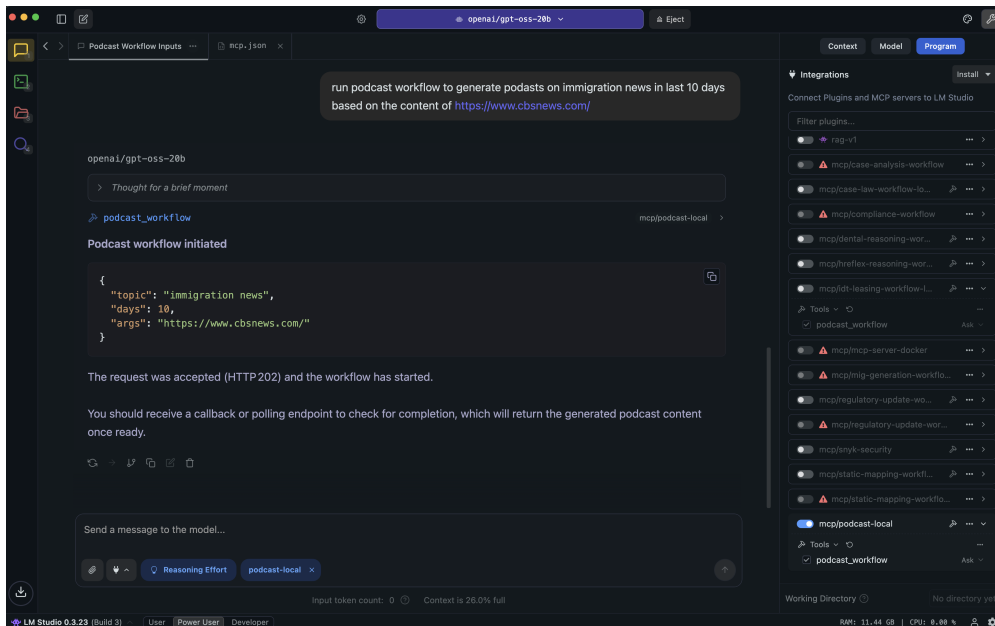Figure 17: LM Studio interpreting the user input and mapping it to the workflow MCP server's input parameters.



Figure 18: Integration of the workflow with LM Studio via the MCP server, illustrating the invocation of the workflow function.

## 5. Evaluation

The evaluation of the proposed agentic workflow focuses on the performance and accuracy of four core components: the podcast script generation agents, the reasoning agent, the video-script generation agent, and the Veo-3 JSON builder agent. The first stage of evaluation examines the podcast generation agents, which operate as a multi-model consortium composed of Llama, OpenAI, and Gemini. The prompt template used to instruct these agents is shown in Figure 19, while Figures 20, 21, and 22 present representative podcast scripts produced by each model.

These outputs demonstrate the natural diversity that emerges from heterogeneous LLMs. Llama tends to generate concise, structured summaries; OpenAI produces more detailed, narrative-driven content; and Gemini emphasizes stylistic flow and contextual framing. This diversity is advantageous, as it captures different semantic and stylistic dimensions of the underlying news content. However, it also introduces inconsistencies, emphasis drift, and occasional factual variations—highlighting the need for a downstream consolidation mechanism.

```
 1  ROLE
 2  You write concise, engaging two-person podcast scripts based on scraped source content for a requested
    TOPIC. Keep it conversational and fact-grounded; no citations in the dialog.
 3
 4  TOPIC HANDLING
 5  • If the first non-empty line starts with 'TOPIC:' use that exact topic string to scope the script.
 6  • Otherwise infer the topic from the sources' titles/content and stick to it.
 7  • Stay on-topic; ignore unrelated material even if present in sources.
 8
 9  INPUT FORMAT (PLAIN TEXT)
10  You receive one or more repeated blocks:
11    ## Source
12    <URL or identifier>
13
14    ## Content
15    <page content in markdown or text>
16
17  Blocks may repeat. Noise (menus/footers/share bars) may be present—ignore it.
18
19  PROCESSING RULES
20  1) Read all sources, extract the 4–8 most relevant facts/points tied to the TOPIC.
21  2) Do NOT invent facts. If details are missing, acknowledge them briefly in-dialog instead of
    hallucinating.
22  3) Avoid marketing fluff; focus on substance.
23  4) Keep language natural for text-to-speech (short sentences, minimal jargon unless explained).
24  5) Use only information from the provided sources; no outside knowledge.
25
26  OUTPUT FORMAT (MARKDOWN)
27  ## Podcast Script
28  - Title: <short, catchy, on-topic>
29  - Hosts: <Host A> and <Host B>
30  - Episode Summary: <1–2 sentences anchoring to sources, no citations>
31
32  ### Script (alternating turns)
33  - Host A: <hook/intro referencing TOPIC>
34  - Host B: <response and setup>
35  - <Alternate Host A / Host B for 10–14 turns, weaving in extracted facts. Include light context/
    explanations so a listener can follow without reading the sources.>
36  - Host A: <wrap-up with key takeaway>
37  - Host B: <closing/call-to-action, e.g., "thanks for listening">
38
39  STYLE
40  - Conversational, friendly, and concise; avoid monologues longer than 2 sentences per turn.
41  - Attribute facts to "reports," "updates," or "the article" instead of specific refs; keep proper nouns
    from sources.
42  - If sources conflict or are thin, surface that briefly in-dialog.
43  - No bullet lists inside the script turns; keep them as plain sentences.
```

Figure 19: Prompt template used by the Podcast Script Generation Agents.

```
1  # Podcast script (Gemini)
2
3  ## Podcast Script
4  - Title: Immigration Crackdown Expands to Charlotte
5  - Hosts: Alex and Ben
6  - Episode Summary: We're discussing recent immigration news, specifically a new federal operation in
   Charlotte, North Carolina, and its planned expansion to other cities.
7
8  ### Script (alternating turns)
9  - Alex: Hey Ben, there's been a lot in the news lately about immigration, and it seems federal
   operations are really ramping up in certain areas.
10 - Ben: That's right, Alex. The latest reports indicate that federal immigration agents have officially
   launched a new operation in Charlotte, North Carolina.
11 - Alex: So, what's this operation called, and who exactly are they targeting?
12 - Ben: It's been dubbed "Charlotte's Web," and the Department of Homeland Security says it's
   specifically targeting immigrants living in the area illegally, particularly those identified as
   "criminal illegal aliens" and "public safety threats."
13 - Alex: Interesting. And who's leading this effort? I heard a name mentioned that's been involved in
   similar actions before.
14 - Ben: You're thinking of Commander Gregory Bovino. He's reportedly involved, and he previously led a
   controversial crackdown in the Chicago area. This operation in Charlotte is expected to involve scores
   of Border Patrol agents.
15 - Alex: Wow, "scores of agents." The reports even suggest the potential use of armored vehicles and
   special operations teams, which sounds pretty intense for a city operation.
16 - Ben: It does. And as you might expect, this has drawn strong criticism. Local Democratic officials,
   including Charlotte Mayor Vi Lyles, have voiced concerns, saying the operation is causing "unnecessary
   fear and uncertainty" in the community.
17 - Alex: That's a significant point, highlighting the local impact. Is Charlotte the only city seeing
   this kind of increased activity, or are there plans for other areas?
18 - Ben: According to the updates, Charlotte is just one step. After this, Border Patrol's focus is
   expected to shift to New Orleans, where they're planning an operation called "Catahoula Crunch,"
   potentially deploying up to 200 agents there.
19 - Alex: So, it sounds like a broader, widening crackdown on illegal immigration, with Charlotte and New
   Orleans being key targets in this current phase.
20 - Ben: Exactly. It's a significant development in federal immigration enforcement, impacting multiple
   cities and communities.
21 - Alex: Thanks for breaking that down, Ben. It's clear these operations are having a big impact on the
   communities involved.
22 - Ben: Absolutely. Thanks for listening, everyone. We'll keep an eye on these developments.
```

Figure 20: Podcast script generated by the Gemini model.

```
1  # Podcast script (OpenAI)
2
3  ## Podcast Script
4  - Title: Charlotte's Web: DHS Crackdown Hits Charlotte, Eyes New Orleans
5  - Hosts: Host A and Host B
6  - Episode Summary: Reports say DHS launched "Charlotte's Web" in Charlotte as part of a broader
   immigration crackdown, with Border Patrol teams making arrests and local leaders pushing back. Internal
   documents point to a large-scale operation now and a planned shift to New Orleans next.
7
8  ### Script (alternating turns)
9  - Host A: Big immigration update today. DHS says it just launched "Charlotte's Web" in Charlotte, part
   of a widening crackdown on people living in the U.S. illegally.
10 - Host B: The agency frames it as a public safety push, saying they're surging officers to remove
   threats. A spokesperson argued there have been too many victims of criminal illegal aliens.
11
12 - Host A: Reports say scores of Border Patrol agents are involved, and internal documents suggest
   armored vehicles and special operations teams could be used.
13 - Host B: Videos posted on social media show green-uniformed Border Patrol agents making arrests around
   the region, which lines up with the timing of the launch.
14
15 - Host A: A familiar name is tied to this. Commander Gregory Bovino, who led a weeks-long, controversial
   crackdown in the Chicago area, is expected to be part of the Charlotte effort.
16 - Host B: That Chicago context is why critics in North Carolina are already vocal. Democratic officials,
   including Charlotte's Mayor Vi Lyles, say the operation is causing unnecessary fear and uncertainty.
17
18 - Host A: DHS counters that this is about public safety, but the reports don't detail how agents are
   prioritizing targets beyond the "criminal illegal aliens" label.
19 - Host B: Scale matters here. It's not just Charlotte. After this push, the focus is expected to shift
   to New Orleans.
20
21 - Host A: And that New Orleans phase even has a code name—"Catahoula Crunch"—with plans for as many as
   200 agents, according to the reporting.
22 - Host B: The operation in Charlotte began Saturday, and the on-the-ground videos appeared the same day,
   so this is moving fast.
23
24 - Host A: One more note: the deployment to Charlotte had been flagged earlier, and now it's official and
   visible.
25 - Host B: We'll watch for updates on who's being arrested, how tactics evolve, and when the New Orleans
   shift starts.
26
27 - Host A: Key takeaway: DHS has kicked off a large, highly visible operation in Charlotte, framed as
   targeting public safety threats, but it's drawing sharp local criticism and may soon expand to New
   Orleans.
28 - Host B: Thanks for listening. We'll keep tracking the developments and bring you the latest on
   immigration news.
```

Figure 21: Podcast script generated by the OpenAI model.

```
 1  # Podcast script (Llama)
 2
 3  ## Podcast Script
 4
 5  - Title: "Immigration Crackdown in Charlotte"
 6  - Hosts: Rachel and Mike
 7  - Episode Summary: Federal Immigration agents have begun operating in Charlotte, North Carolina as part
    of the Trump administration's crackdown on illegal immigration. We'll explore the details and reactions
    to this bold move.
 8
 9  ### Script
10
11  Rachel: Welcome back to our show, folks! Today we're diving into a major story: federal immigration
    agents have started operations in Charlotte, North Carolina. Mike, what do we know about this
    development?
12
13  Mike: That's right, Rachel. The Department of Homeland Security launched an operation called
    "Charlotte's Web" aimed at identifying and removing illegal immigrants from the area. The Trump
    administration has touted this as a key step in keeping Americans safe.
14
15  Rachel: According to reports, Border Patrol agents are involved in the operation, which may involve
    armored vehicles and special operations teams. Hundreds of agents could be dispatched to Charlotte.
16
17  Mike: Yeah, that's right. And it's not just Charlotte, folks. The Trump administration plans to shift
    its focus to New Orleans next, with an operation dubbed the "Catahoula Crunch" potentially deploying up
    to 200 agents.
18
19  Rachel: Now, I've got some clips from local leaders in North Carolina expressing concerns about this
    operation. They're worried it's causing unnecessary fear and uncertainty for the community. Mike?
20
21  Mike: Yeah, that's definitely a concern. Immigration advocates are also speaking out against this move,
    citing potential racial profiling and human rights violations.
22
23  Rachel: The administration is pushing back on these criticisms, though, saying they're committed to
    public safety. What do you make of all this, folks? Share your thoughts with us!
24
25  Mike: And don't forget to tune in next time for more immigration news and analysis! Thanks for listening.
26
27  (Note: I've condensed the original text into a conversational script format.)
```

Figure 22: Podcast script generated by the Llama model.

To resolve these differences and produce a final, authoritative output, the workflow employs a dedicated reasoning agent responsible for synthesizing the consortium's drafts into a unified script. The reasoning agent prompt—shown in Figure 23—explicitly instructs the model to compare, cross-validate, and reconcile the outputs of the three podcast agents. It retains only information consistently supported across drafts, while removing speculation, correcting emphasis drift, and resolving contradictory statements.

The resulting consolidated script, illustrated in Figure 24, shows marked improvements in clarity, factual stability, and narrative coherence. The reasoning agent effectively reduces hallucination risk by grounding its synthesis

in multi-model agreement. This consensus-driven approach not only improves overall accuracy but also aligns the workflow with Responsible AI principles—balancing diverse model perspectives and mitigating single-model bias.

```
1  You are the final podcast reasoning agent. Given a TOPIC and three draft podcast scripts (from OpenAI,
   Gemini, Anthropic), produce one consolidated, best-of script.
2
3  Inputs:
4  - **TOPIC**
5  - **Draft scripts** from the three models. Each draft follows the podcast format (title, hosts, summary,
   alternating turns).
6
7  Objectives:
8  1) Compare the drafts to identify shared facts/themes and conflicts. Prefer details mentioned by
   multiple drafts; drop speculative or conflicting items.
9  2) Stay grounded in the drafts—no new facts. If something is unclear or missing, acknowledge it briefly
   in-dialog instead of inventing.
10 3) Produce a single polished podcast script that is concise, engaging, and ready for text-to-speech.
11
12 Output (Markdown, strict):
13 ## Podcast Script
14 - Title: <short, on-topic, best option or blended>
15 - Hosts: <Host A> and <Host B> (reuse names from drafts when possible; otherwise pick simple names)
16 - Episode Summary: <1-2 sentences grounded in draft content, no citations>
17
18 ### Script (alternating turns)
19 - Host A: <hook/intro referencing TOPIC and main angle>
20 - Host B: <response/setup>
21 - <Alternate Host A / Host B for 10-14 turns, weaving in the strongest facts/insights from the drafts.
   If drafts conflict or lack detail, surface that lightly (e.g., "reports vary on timing").>
22 - Host A: <wrap-up/key takeaway>
23 - Host B: <closing/call-to-action>
24
25 Style:
26 - Conversational, friendly, concise; max 2 sentences per turn.
27 - Keep jargon minimal or quickly explained. Attribute to "reports/updates/source" rather than model
   names.
28 - Do not include inline citations or reference lists; rely only on draft content.
29 - If a useful fact appears in only one draft, use it only if it does not conflict with the others and
   fits the TOPIC.
```

Figure 23: Prompt template used by the Reasoning Agent for cross-model consolidation.

```
 1  # Podcast script (Consolidated)
 2
 3  ## Podcast Script
 4  - Title: Charlotte's Web: DHS Operation in Charlotte, Eyes New Orleans
 5  - Hosts: Alex and Ben
 6  - Episode Summary: Reports say DHS launched "Charlotte's Web" in Charlotte, targeting people labeled as
    criminal illegal aliens and public safety threats. The effort is drawing local criticism and is expected
    to shift next to New Orleans under "Catahoula Crunch."
 7
 8  ### Script (alternating turns)
 9  - Alex: Big immigration update today—reports say DHS has launched an operation in Charlotte called
    "Charlotte's Web."
10  - Ben: The agency frames it as a public safety push, focusing on people identified as criminal illegal
    aliens and public safety threats.
11
12  - Alex: Scale-wise, reports point to scores of Border Patrol agents, with internal documents suggesting
    armored vehicles and special operations teams could be used.
13  - Ben: Videos posted on social media show green-uniformed Border Patrol agents making arrests around the
    region, lining up with the launch.
14
15  - Alex: A familiar name is tied to this—reports link Commander Gregory Bovino to the effort; he
    previously led a controversial crackdown in the Chicago area.
16  - Ben: Local Democratic officials, including Charlotte Mayor Vi Lyles, say the operation is causing
    unnecessary fear and uncertainty.
17
18  - Alex: DHS counters that it's about keeping communities safe, but beyond those labels, details on how
    targets are prioritized aren't clear in the reporting.
19  - Ben: Timing-wise, some updates say the Charlotte phase began Saturday, though specifics could evolve
    as more information comes in.
20
21  - Alex: It's not just Charlotte; after this push, the focus is expected to shift to New Orleans.
22  - Ben: That phase is reportedly code-named "Catahoula Crunch," with plans for up to 200 agents.
23
24  - Alex: Worth noting, some of the heavier tactics mentioned—like armored vehicles—come from documents
    and reports, so we'll watch to see what's actually used on the ground.
25  - Ben: Key takeaway: a large, highly visible DHS operation is underway in Charlotte amid sharp local
    criticism, with New Orleans likely next.
26
27  - Alex: We'll keep tracking who's being targeted, any changes in tactics, and when the New Orleans shift
    starts.
28  - Ben: Thanks for listening—subscribe for updates as this story develops.
```

Figure 24: Final consolidated podcast script generated by the Reasoning Agent.

The video-script generation agent and the Veo-3 JSON builder agent were also evaluated to verify seamless multimodal content generation. The video-script agent reliably transforms the consolidated podcast script into scene-based descriptions that preserve narrative fidelity and temporal coherence. Figure 25 presents the prompt template used by the video-script generation agent, while Figure 26 shows a representative output illustrating how the agent structures the script into coherent, visually aligned scenes.

```
1 ROLE
2 Create a tight, visual-ready micro video script from a podcast script about the given TOPIC. Keep it
  grounded in the provided content—no new facts.
3
4 INPUT
5 - TOPIC (first line starts with 'TOPIC:' or is provided in the system header)
6 - Podcast script text
7
8 OUTPUT (STRICT)
9 Title: <concise, on-topic>
10 Duration: 8-10 seconds
11 Beat 1 (0-2s): Narration - <≤14 words> | Visuals - <concise action/setting, no on-screen text>
12 Beat 2 (2-4s): Narration - <≤14 words> | Visuals - <...>
13 Beat 3 (4-6s): Narration - <≤14 words> | Visuals - <...>
14 Beat 4 (6-8/10s): Narration - <≤14 words> | Visuals - <...>
15
16 RULES
17 - Use only facts/themes present in the podcast script; if timing/details are unclear, keep narration
   non-specific.
18 - Conversational tone, plain language; avoid legalese and filler.
19 - Visuals: describe imagery, movement, and setting; avoid text overlays, logos, or calls-to-action.
20 - No links, citations, or markdown fences—just the lines above.
```

Figure 25: Prompt template used by the Video Script Generation Agent.

```
1 ## Video Script
2
3 Title: DHS Operation "Charlotte's Web"
4 Duration: 8-10 seconds
5 Beat 1 (0-2s): Narration - DHS launched "Charlotte's Web" in Charlotte, targeting public safety
  threats. | Visuals - Quick shot of a city street, then blurred image of agents in uniform.
6 Beat 2 (2-4s): Narration - Local officials criticize the operation, citing unnecessary fear in the
  community. | Visuals - Close-up of a concerned face in a crowd, then a shot of a city hall building.
7 Beat 3 (4-6s): Narration - Social media shows Border Patrol agents making arrests in the region. |
  Visuals - Handheld phone footage style of green-uniformed agents, quick cut to an arrest.
8 Beat 4 (6-8/10s): Narration - The operation will shift to New Orleans next, code-named "Catahoula Crunch.
  " | Visuals - Map animation showing a line from Charlotte to New Orleans, then a skyline shot of New
  Orleans.
```

Figure 26: Video script output generated by the Video Script Generation Agent.

The Veo-3 JSON builder agent was evaluated for structural correctness, schema alignment, and integration fidelity with Google Veo-3. As shown in Figure 27, the agent receives a structured prompt describing the podcast narrative and corresponding visual plan. Figure 28 presents a sample generated Veo-3 JSON prompt, which is syntactically valid, interpretable, and executable by Veo-3's video-generation API. Across multiple test runs, the agent consistently produced well-formed JSON specifications capable of driving automated video generation without requiring manual correction.

```
 1  ROLE
 2  Convert a 4-beat micro video script into a concise Veo 3 JSON prompt for rendering.
 3
 4  INPUT
 5  - TOPIC
 6  - Video script with Beats 1-4 (narration + visuals)
 7
 8  OUTPUT (STRICT JSON ONLY)
 9  {
10    "scene_description": "<timeline 0-8/10s with beat-by-beat visuals and narration cues>",
11    "visual_style": "<overall look and feel in a short phrase>",
12    "camera_movement": "<simple camera moves per beat>",
13    "main_subject": "<what is primarily shown>",
14    "background_setting": "<where it happens>",
15    "lighting_mood": "<lighting & mood>",
16    "audio_cue": "<SFX/music cues aligned to beats>",
17    "color_palette": ["#hex1", "#hex2"],
18    "dialog": "<the 4 narration sentences joined in order>",
19    "subtitles": "OFF"
20  }
21
22  RULES
23  - Map beats to exact time ranges: 0-2s, 2-4s, 4-6s, 6-8/10s.
24  - Keep every field brief and production-ready; avoid long prose or repetition.
25  - Reflect any geography/setting cues from the script subtly in visuals; no on-screen text.
26  - Return ONLY the JSON object (no markdown fences or commentary).
```

Figure 27: Prompt template used by the Veo-3 JSON Builder Agent.

```
1  ## Veo 3 JSON Prompt
2
3  ```
4  {
5    "scene_description": "0-2s: Quick shot of a bustling city street, then a blurred, indistinct image of
      agents in uniform. Narration: DHS launched \"Charlotte's Web\" in Charlotte, targeting public safety
      threats. 2-4s: Close-up of a concerned face in a diverse urban crowd, followed by a clear shot of a
      city hall building. Narration: Local officials criticize the operation, citing unnecessary fear in the
      community. 4-6s: Handheld phone footage style showing green-uniformed agents moving, with a quick cut
      to an agent making an arrest. Narration: Social media shows Border Patrol agents making arrests in the
      region. 6-8/10s: An animated map highlights Charlotte, then draws a line to New Orleans, concluding
      with a vibrant skyline shot of New Orleans at dusk. Narration: The operation will shift to New Orleans
      next, code-named \"Catahoula Crunch.\"",
6    "visual_style": "Gritty, realistic news report aesthetic with a mix of professional and user-generated
      content styles.",
7    "camera_movement": "0-2s: Quick pan, then static blur. 2-4s: Dolly in, then static wide. 4-6s: Shaky
      handheld, quick cut. 6-8/10s: Animated map zoom, then slow reveal.",
8    "main_subject": "DHS agents, urban citizens, cityscapes of Charlotte and New Orleans, animated map.",
9    "background_setting": "Urban streets, city hall, general city environments of Charlotte and New
      Orleans.",
10   "lighting_mood": "Neutral, realistic daylight for most scenes, transitioning to a slightly somber or
      dramatic dusk for the New Orleans skyline.",
11   "audio_cue": "Subtle urban ambient sounds, news report underscore music, brief sound of a crowd murmur.
      ",
12   "color_palette": ["#4CAF50", "#607D8B", "#2196F3", "#FFC107", "#B0BEC5"],
13   "dialog": "DHS launched \"Charlotte's Web\" in Charlotte, targeting public safety threats. Local
      officials criticize the operation, citing unnecessary fear in the community. Social media shows Border
      Patrol agents making arrests in the region. The operation will shift to New Orleans next, code-named
      \"Catahoula Crunch.\"",
14   "subtitles": "OFF"
15  }
16  ```
```

Figure 28: Veo-3 JSON prompt generated by the Veo-3 Builder Agent.

Together, these evaluations confirm that each agent performs its specialized role effectively, with the workflow producing stable, coherent, and high-quality multimodal outputs. The results also demonstrate the value of agent specialization—podcast script generation, reasoning consolidation, video-script generation, and Veo-3 JSON construction—working in coordination to deliver a complete, production-aligned agentic workflow.

## 6. Conclusions and Future Work

Agentic AI represents a fundamental shift in how autonomous systems are engineered, enabling AI-driven pipelines that can reason, act, observe, and iteratively refine their outputs across complex, multi-step tasks. However, achieving production-grade reliability requires far more than chaining LLM calls together—it demands disciplined engineering, modular design, and operational rigor. In this paper, we present a practical, end-to-end guide for designing, developing, and deploying agentic AI workflows

38

that are maintainable, deterministic, auditable, and aligned with Responsible AI principles. Through a multimodal news-to-podcast generation workflow, we demonstrated how agent orchestration, tool/function separation, model-consortium reasoning, containerized deployment, and MCP-based accessibility can be combined to form a scalable and extensible automation architecture. The nine best practices distilled from this case study provide actionable guidance for practitioners seeking to balance flexibility with robustness—ensuring that workflows remain interpretable, secure, and stable even as underlying LLMs evolve. By unifying architectural patterns, operational best practices, multi-agent reasoning, and deployment strategies, this work contributes a reusable blueprint for building production-grade agentic systems. As organizations increasingly rely on AI to automate sophisticated workflows, these principles will be essential to advance reliability, safety, and long-term maintainability. Future work will explore adaptive evaluation pipelines, workflow self-monitoring, and tighter safety/guardrail integrations to further enhance the trustworthiness of agentic AI deployments. For future work, we plan to extend the proposed best practices to a broader set of agentic AI workflow–automation use cases, validating their effectiveness across diverse domains and increasingly complex multi-agent pipelines.

# References

[1] S. Arora, B. Yang, S. Eyuboglu, A. Narayan, A. Hojel, I. Trummer, C. Ré, Language models enable simple systems for generating structured views of heterogeneous data lakes, arXiv preprint arXiv:2304.09433 (2023).

[2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., Language models are few-shot learners, Advances in neural information processing systems 33 (2020) 1877–1901.

[3] J. Zhang, J. Huang, S. Jin, S. Lu, Vision-language models for vision tasks: A survey, IEEE Transactions on Pattern Analysis and Machine Intelligence (2024).

[4] P. Agrawal, S. Antoniak, E. B. Hanna, B. Bout, D. Chaplot, J. Chudnovsky, D. Costa, B. De Monicault, S. Garg, T. Gervet, et al., Pixtral 12b, arXiv preprint arXiv:2410.07073 (2024).

[5] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge, et al., Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution, arXiv preprint arXiv:2409.12191 (2024).

[6] D. B. Acharya, K. Kuppan, B. Divya, Agentic ai: Autonomous intelligence for complex goals–a comprehensive survey, IEEE Access (2025).

[7] E. Bandara, R. Gore, X. Liang, S. Rajapakse, I. Kularathne, P. Karunarathna, P. Foytik, S. Shetty, R. Mukkamala, A. Rahman, et al., Agentsway–software development methodology for ai agents-based teams, arXiv preprint arXiv:2510.23664 (2025).

[8] E. Bandara, A. Hass, S. Shetty, R. Mukkamala, R. Gore, A. Rahman, S. H. Bouk, Deep-stride: Automated security threat modeling with vision-language models, in: 2025 International Conference on Software, Telecommunications and Computer Networks (SoftCOM), 2025, pp. 1–7.

[9] R. Gore, E. Bandara, S. Shetty, A. E. Musto, P. Rana, A. Valencia-Romero, C. Rhea, L. Tayebi, H. Richter, A. Yarlagadda, et al., Proof-of-tbi–fine-tuned vision language model consortium and openai-o3 reasoning llm-based medical diagnosis support system for mild traumatic brain injury (tbi) prediction, arXiv preprint arXiv:2504.18671 (2025).

[10] X. Hou, Y. Zhao, S. Wang, H. Wang, Model context protocol (mcp): Landscape, security threats, and future research directions, arXiv preprint arXiv:2503.23278 (2025).

[11] E. Bandara, S. Shetty, R. Mukkamala, R. Gore, P. Foytik, S. H. Bouk, A. Rahman, X. Liang, N. W. Keong, K. De Zoysa, et al., Model context contracts-mcp-enabled framework to integrate llms with blockchain smart contracts, arXiv preprint arXiv:2510.19856 (2025).

[12] A. Singh, A. Ehtesham, S. Kumar, T. T. Khoei, Enhancing ai systems with agentic workflows patterns in large language model, in: 2024 IEEE World AI IoT Congress (AIIoT), IEEE, 2024, pp. 527–532.

[13] G. Marvin, N. Hellen, D. Jjingo, J. Nakatumba-Nabende, Prompt engineering in large language models, in: International Conference on Data Intelligence and Cognitive Informatics, Springer, 2023, pp. 387–402.

[14] E. Bandara, S. Shetty, R. Mukkamala, R. Gore, P. Foytik, S. H. Bouk, A. Rahman, X. Liang, N. W. Keong, K. De Zoysa, et al., Model context contracts-mcp-enabled framework to integrate llms with blockchain smart contracts, arXiv preprint arXiv:2510.19856 (2025).

[15] E. Bandara, R. Gore, S. Shetty, R. Mukkamala, C. Rhea, A. Yarlagadda, S. Kaushik, L. De Silva, A. Maznychenko, I. Sokolowska, et al., Standardization of neuromuscular reflex analysis–role of fine-tuned vision-language model consortium and openai gpt-oss reasoning llm enabled decision support system, arXiv preprint arXiv:2508.12473 (2025).

[16] J. Wang, A tutorial on llm reasoning: Relevant methods behind chatgpt o1, arXiv preprint arXiv:2502.10867 (2025).

[17] S. Agarwal, L. Ahmad, J. Ai, S. Altman, A. Applebaum, E. Arbus, R. K. Arora, Y. Bai, B. Baker, H. Bao, et al., gpt-oss-120b & gpt-oss-20b model card, arXiv preprint arXiv:2508.10925 (2025).

[18] G. Mondillo, M. Masino, S. Colosimo, A. Perrotta, V. Frattolillo, Evaluating ai reasoning models in pediatric medicine: A comparative analysis of o3-mini and o3-mini-high, medRxiv (2025) 2025–02.

[19] T. Wiedemer, Y. Li, P. Vicol, S. S. Gu, N. Matarese, K. Swersky, B. Kim, P. Jaini, R. Geirhos, Video models are zero-shot learners and reasoners, arXiv preprint arXiv:2509.20328 (2025).

[20] N. I. Ali, A. H. Soomro, S. Soomro, A. G. Memon, et al., Rss feeds filtering from multiple sources using automated techniques of natural language processing, The Asian Bulletin of Big Data Management 4 (1) (2024) Science–4.

[21] M. Imran, N. Almusharraf, Google gemini as a next generation ai educational tool: a review of emerging educational technology, Smart Learning Environments 11 (1) (2024) 22.

[22] N. I. Ali, A. H. Soomro, S. Soomro, A. G. Memon, et al., Rss feeds filtering from multiple sources using automated techniques of natural language processing, The Asian Bulletin of Big Data Management 4 (1) (2024) Science–4.

[23] E. Bandara, R. Gore, A. Yarlagadda, A. H. Clayton, P. Samuel, C. K. Rhea, S. Shetty, Standardization of psychiatric diagnoses–role of fine-tuned llm consortium and openai-gpt-oss reasoning llm enabled decision support system, arXiv preprint arXiv:2510.25588 (2025).

[24] R. R. Manku, Y. Tang, X. Shi, M. Li, A. Smola, Emergenttts-eval: Evaluating tts models on complex prosodic, expressiveness, and linguistic challenges using model-as-a-judge, arXiv preprint arXiv:2505.23009 (2025).

[25] M. M. Hasan, H. Li, E. Fallahzadeh, G. K. Rajbahadur, B. Adams, A. E. Hassan, Model context protocol (mcp) at first glance: Studying the security and maintainability of mcp servers, arXiv preprint arXiv:2506.13538 (2025).

[26] T. Masterman, S. Besen, M. Sawtell, A. Chao, The landscape of emerging ai agent architectures for reasoning, planning, and tool calling: A survey, arXiv preprint arXiv:2404.11584 (2024).

[27] G. Ntousakis, J. J. Stephen, M. V. Le, S. S. L. Chukkapalli, T. Taylor, I. M. Molloy, F. Araujo, Securing mcp-based agent workflows, in: Proceedings of the 4th Workshop on Practical Adoption Challenges of ML for Systems, 2025, pp. 50–55.

[28] E. Bandara, S. H. Bouk, S. Shetty, R. Gore, S. Kompella, R. Mukkamala, A. Rahman, P. Foytik, X. Liang, N. W. Keong, K. De Zoysa, Vindsec-llama — fine-tuned meta's llama-3 llm, federated learning, blockchain and pbom-enabled data security architecture for wind energy data platforms, in: 2025 International Wireless Communications and Mobile Computing (IWCMC), 2025, pp. 120–126. doi:10.1109/IWCMC65282.2025.11059510.

[29] B. Perak, S. Beliga, A. Meštrović, Incorporating dialect understanding into llm using rag and prompt engineering techniques for causal commonsense reasoning, in: Proceedings of the Eleventh Workshop on NLP for Similar Languages, Varieties, and Dialects (VarDial 2024), 2024, pp. 220–229.

[30] A. Yehudai, L. Eden, A. Li, G. Uziel, Y. Zhao, R. Bar-Haim, A. Cohan, M. Shmueli-Scheuer, Survey on evaluation of llm-based agents, arXiv preprint arXiv:2503.16416 (2025).

[31] E. Bandara, S. H. Bouk, S. Shetty, S. Roy, R. Mukkamala, A. Rahman, P. Foytik, X. Liang, N. W. Keong, K. De Zoysa, Llama-recipe — fine-tuned meta's llama llm, pbom and nft enabled 5g network-slice orchestration and end-to-end supply-chain verification platform, in: 2025 IEEE 22nd Consumer Communications & Networking Conference (CCNC), 2025, pp. 1–6. doi:10.1109/CCNC54725.2025.10976116.

[32] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al., The llama 3 herd of models, arXiv preprint arXiv:2407.21783 (2024).

[33] Y. Zhang, S. Mao, T. Ge, X. Wang, A. de Wynter, Y. Xia, W. Wu, T. Song, M. Lan, F. Wei, Llm as a mastermind: A survey of strategic reasoning with large language models, arXiv preprint arXiv:2404.01230 (2024).

[34] Y. Yan, S. Wang, J. Du, Y. Yang, Y. Shan, Q. Qiu, X. Jia, X. Wang, X. Yuan, X. Han, et al., Mcpworld: A unified benchmarking testbed for api, gui, and hybrid computer use agents, arXiv preprint arXiv:2506.07672 (2025).

[35] E. Bandara, S. Shetty, R. Mukkamala, A. Rahman, P. Foytik, X. Liang, K. De Zoysa, N. W. Keong, Devsec-gpt — generative-ai (with custom-trained meta's llama2 llm), blockchain, nft and pbom enabled cloud native container vulnerability management and pipeline verification platform, in: 2024 IEEE Cloud Summit, 2024, pp. 28–35. doi:10.1109/Cloud-Summit61220.2024.00012.

[36] E. Bandara, S. H. Bouk, S. Shetty, R. Mukkamala, A. Rahman, P. Foytik, R. Gore, X. Liang, N. W. Keong, K. De Zoysa, Sre-llama–fine-tuned meta's llama llm, federated learning, blockchain and nft enabled site reliability engineering (sre) platform for communication and networking software services, arXiv preprint arXiv:2511.08282 (2025).

[37] X. Yu, L. Liu, X. Hu, J. Liu, X. Xia, Where are large language models for code generation on github?, arXiv preprint arXiv:2406.19544 (2024).

[38] E. Chen, C. Lin, X. Tang, A. Xi, C. Wang, J. Lin, K. R. Koedinger, Vtutor: An open-source sdk for generative ai-powered animated pedagogical agents with multi-media output, arXiv preprint arXiv:2502.04103 (2025).

[39] I. H. Sarker, Llm potentiality and awareness: a position paper from the perspective of trustworthy and responsible ai modeling, Discover Artificial Intelligence 4 (1) (2024) 40.

[40] E. Bandara, Multimodel podcast generation agentic ai workflow implementation, GitLab Repository, available at: `https://gitlab.com/rahasak-labs/podcast-workflow` (Oct. 2025).

[41] E. Bandara, Mcp server for the multimodel podcast generation agentic ai workflow, GitLab Repository, available at: `https://gitlab.com/rahasak-labs/podcast-workflow-mcp-server` (Oct. 2025).

[42] A. Sojan, R. Rajan, P. Kuvaja, Monitoring solution for cloud-native devsecops, in: 2021 IEEE 6th International Conference on Smart Cloud (SmartCloud), IEEE, 2021, pp. 125–131.