# Advanced Design Patterns for Structuring CrewAI Agents with the Gemini API

This report details advanced design patterns for creating flexible and dynamic CrewAI frameworks. It focuses on using a primary agent to interpret natural language inputs and assemble specialized agent teams, integrating the Gemini API for complex reasoning, and structuring outputs for clarity and streaming.

## Core Architecture: The Coordinator and Specialist Model

A powerful and common pattern in CrewAI is to designate a primary agent as a coordinator or router that interprets a user's goal and delegates sub-tasks to a team of specialized agents [ref: 0-0, ref: 0-3].

### Intent-Based Delegation

A primary agent can be designed to classify a user's intent and orchestrate other agents accordingly. In a sample "ChatGPT clone" application, a primary `chat_agent` is configured to be the main user interface. Its goal is to understand the user's query, classify its intent, and delegate to the appropriate specialist [ref: 0-0].

- **Coordinator Agent (`chat_agent`):**
   - **Role:** Conversational Assistant [ref: 0-0].
   - **Goal:** Understand user intent and orchestrate specialized agents [ref: 0-0].
   - **Delegation:** Must have `allow_delegation: true` to pass tasks to other agents [ref: 0-0].
- **Specialist Agents:**
   - **`search_agent`:** Handles real-time factual lookups [ref: 0-0].
   - **`research_agent`:** Performs in-depth analysis on complex topics [ref: 0-0].
   - **`scraper_agent`:** Extracts specific data from websites [ref: 0-0].
   - **`image_agent`:** Generates images from descriptions [ref: 0-0].

This hierarchical structure allows the coordinator to manage the workflow, deciding which specialist is best suited for each part of the user's request [ref: 0-0].

### Conditional Routing with a Router Agent

A more advanced implementation of this pattern uses a dedicated `Router_Agent` to decide the execution path based on the query's content. In a Retrieval-

Augmented Generation (RAG) system, a `Router_Agent` determines whether to search a local vector store or the web [ref: 0-3].

- The `Router_Agent`'s goal is to "Route user question to a vectorstore or web search" [ref: 0-3].
- It uses a custom `router_tool` that returns a single keyword ('vectorstore' or 'web_search') to direct the flow [ref: 0-3].
- Subsequent tasks use the output of this router task as `context` to select the correct tool (`rag_tool` or `web_search_tool`) for information retrieval [ref: 0-3].

## Defining Flexible Agents and Tasks

CrewAI's project-based structure with YAML configuration is the recommended approach for creating maintainable and scalable agent systems [ref: 0-0, ref: 0-4].

### YAML Configuration

Defining agents and tasks in `agents.yaml` and `tasks.yaml` files separates configuration from code, which offers several benefits [ref: 0-0]:
- **Readability:** YAML is human-readable and easier to maintain [ref: 0-0].
- **Separation of Concerns:** Keeps agent/task definitions separate from business logic [ref: 0-0].
- **Collaboration:** Allows non-technical team members to modify agent personas and tasks without touching Python code [ref: 0-0].

Dynamic inputs are passed to agents and tasks at runtime using variable placeholders like `{user_input}` in the YAML files. These are populated when the crew is executed with the `kickoff()` method [ref: 0-0].

### Advanced Agent Configuration

CrewAI provides several parameters to fine-tune agent behavior, which can be set in the YAML file or in Python code [ref: 0-0].

| Option | Description |
|---|---|
| `llm` | The model name you want to use (e.g., `google/gemini-2.0-flash`) [ref: 0-0, ref: 0-2]. |
| `allow_delegation` | Controls whether the agent can request help from other agents [ref: 0-0]. |
| `verbose` | Enables detailed logging of agent actions for debugging [ref: 0-0]. |
| `max_rpm` | Limits requests per minute to external services [ref: 0-0]. |
| `max_iter` | Maximum iterations an agent can perform for a task [ref: 0-0]. |
| `max_execution_time` | Sets a time limit for task completion [ref: 0-0]. |

| `tools` | Array of tools the agent can use [ref: 0-0]. |
| `cache` | Determines if the agent should cache tool usage results [ref: 0-0]. |

### Best Practices for Agent and Task Design

- **Agents:** Design agents as specialists rather than generalists. They perform better with clearly defined, non-overlapping roles [ref: 0-0].
- **Tasks:** The majority of effort (80%) should focus on designing clear and specific tasks. Each task should have a single purpose, with explicit descriptions and a clear `expected_output` to guide the agent [ref: 0-0].

## Advanced Orchestration and Dynamic Assembly

For more complex and dynamic applications, CrewAI offers several advanced patterns beyond a simple sequential process.

### CrewAI Flows for Event-Driven Workflows

CrewAI Flows provide fine-grained, event-driven control over execution, combining direct LLM calls, crew-based processing, and regular Python code [ref: 0-2]. This is ideal for creating complex, multi-stage systems [ref: 0-2].

Key components of a Flow include:
- **State Management:** A Pydantic `BaseModel` class is used to define and maintain the application's state across different steps [ref: 0-2].
- **`@start()` Decorator:** Marks the initial function of the flow, which often handles user input [ref: 0-2].
- **`@listen()` Decorator:** Creates an event-driven relationship, triggering a function upon the completion of another [ref: 0-2].
- **Direct LLM Calls:** Use the `LLM` class for simple, structured AI interactions without the overhead of a full agent and crew [ref: 0-2].
- **Crew Integration:** Flows can seamlessly `kickoff` a crew as one of the steps in the larger process [ref: 0-2].

### Hierarchical Process and Planning

CrewAI supports different execution processes that can be set when creating a crew [ref: 0-4].
- **Sequential Process:** The default process where tasks are executed one after another [ref: 0-4].
- **Hierarchical Process:** A manager agent coordinates the crew, delegating tasks and validating outcomes. This requires specifying a `manager_llm` in the crew definition [ref: 0-3, ref: 0-4].
- **Planning Feature:** By setting `planning=True` and providing a

`planning_llm`, a planner agent will generate a step-by-step plan before each iteration and add it to the task descriptions, enabling the crew to adapt its approach dynamically [ref: 0-1].

### Dynamic Generation of Agents and Tasks

For maximum flexibility, agents and tasks can be generated programmatically at runtime. One demonstrated approach involves a script that:
1. Loads available local LLM models (e.g., by running `ollama list`) [ref: 0-1].
2. Dynamically creates an `Agent` and a `Task` for each model [ref: 0-1].
3. The task prompts the model to describe its own ideal agent persona (role, goal, backstory) [ref: 0-1].
4. The results can be used to dynamically generate the `agents.yaml` and `tasks.yaml` files for future use [ref: 0-1].

This pattern allows the system to build a crew from a dynamic list of available resources [ref: 0-1].

## Integrating the Gemini API

Integrating Gemini or any other LLM is straightforward. The LLM can be specified at the agent level or used for direct calls.

- **Agent-Level Configuration:** In `agents.yaml` or when creating an `Agent` instance in Python, set the `llm` parameter. For Gemini, the format would be `provider/model-id`, for example: `google/gemini-2.0-flash` [ref: 0-2].
- **Direct LLM Calls in Flows:** When you need a simple, structured response, you can initialize the `LLM` class directly with a Gemini model: `llm = LLM(model="google/gemini-2.0-flash")` [ref: 0-2].
- **Custom LLM Setup:** For models not natively supported or requiring specific configurations (like Groq), you can use `langchain-groq` or a similar library to instantiate the chat model and pass it to the agent's `llm` parameter. This pattern is adaptable for the Gemini API client [ref: 0-3].

## Structuring and Streaming Outputs

Ensuring consistent, structured output is crucial for reliable agentic systems.

### Structured Outputs with Pydantic and JSON

CrewAI tasks can be configured to produce validated, structured outputs [ref: 0-4].
- **`output_pydantic`:** By assigning a Pydantic `BaseModel` to this task attribute, you instruct the agent to return an output that conforms to the model's

schema. This ensures the output is structured and validated [ref: 0-4].
- **`output_json`:** Similarly, you can use this attribute to ensure the task's output is a valid JSON object, optionally conforming to the structure of a Pydantic model [ref: 0-4].

The final output of a crew is the output of its final task. If that task uses `output_pydantic` or `output_json`, the `CrewOutput` object can be accessed like a dictionary or Pydantic object, making it easy to handle structured data [ref: 0-4].

### Output Validation with Guardrails and Graders

To ensure the quality and accuracy of outputs, especially in complex reasoning tasks, you can implement validation layers.
- **Task Guardrails:** A `guardrail` is a Python function passed to a `Task` that validates its output. If validation fails, the function returns an error message to the agent, which will attempt to fix its output. The `guardrail_max_retries` parameter limits the number of attempts [ref: 0-4].
- **Grader Agents:** A more sophisticated pattern involves adding dedicated agents to the crew to grade and validate information. A RAG pipeline example uses a `Grader_agent` to check for relevance, a `hallucination_grader` to check if the answer is supported by facts, and an `answer_grader` to assess overall usefulness before generating the final response [ref: 0-3].

### Streamable Updates for User Interfaces

While deep, token-level streaming from the agent process is not explicitly detailed, a practical approach for creating a responsive UI is demonstrated using Streamlit [ref: 0-0].
1. The user's query is sent to the crew using `crew.kickoff()` [ref: 0-0].
2. After the crew completes its work and returns the final result object, the application extracts the raw text response [ref: 0-0].
3. A "typewriter effect" is implemented in the front-end to display the response character-by-character, simulating a real-time stream for the user [ref: 0-0].

This provides a smooth user experience by rendering the complete, validated response in a visually appealing, stream-like manner [ref: 0-0].