**⟳ ChatGPT**

# Comparative Analysis of Agentic Frameworks (LangChain, LlamaIndex, Haystack)

## Detailed Framework Reports

### LangChain

- **Tool Usage & MCP Integration:** LangChain provides a rich **tool abstraction** where each tool is defined by a Python function plus a JSON-based schema (name, description, arguments). This structured approach aligns well with OpenAI's function-calling format and JSON Schema. LangChain agents can use tools defined locally or via APIs, and recent updates enable direct integration with the **Model Context Protocol (MCP)**. An official adapter package converts MCP tools into LangChain-compatible tools, allowing agents to seamlessly call tools from one or multiple MCP servers. This makes LangChain MCP-friendly, though it relies on an external package (`langchain-mcp-adapters`) to interface with MCP servers. LangChain supports advanced agent capabilities (e.g. self-ask, sub-LLM calls) but does not natively implement MCP's specific "elicitation" or "LLM sampling" concepts – those would be handled by designing the agent's prompt logic. Overall, LangChain excels in flexible tool use (function calls, OpenAPI, etc.) and now can bridge into MCP with relative ease.

- **State Management & Memory (Qdrant Integration):** LangChain offers multiple **memory modules** for maintaining conversational state across turns (from simple buffers to summarizing or vector-store-backed memory). It directly supports using **vector databases like Qdrant** to store and retrieve long-term memory. In fact, LangChain provides a Qdrant integration (via a partner package `langchain-qdrant`) for vector storage, and a `VectorStoreRetrieverMemory` class exists to leverage vector stores (including Qdrant) for recalling relevant past interactions. This means conversation history or other state can be persisted and retrieved by similarity, which is useful for long-running dialogues. However, **long-running task persistence** is not automatic – LangChain doesn't inherently checkpoint or resume multi-session agents without custom implementation. Sessions would typically be managed by the developer (e.g. storing `ConversationBufferMemory` to a file/DB between runs). LangChain's design is mostly stateless per call unless you attach a persistent memory. It relies on external solutions (like databases or the LangSmith platform) if you need durable state beyond the lifetime of a Python object.

- **Modularity and Extensibility:** One of LangChain's biggest strengths is its highly modular architecture. It allows swapping out **core components** easily – you can choose among dozens of LLM APIs or local models, different vector stores (Chroma, Pinecone, Qdrant, etc.), and tools. This flexibility means developers are **not locked into any single vendor**; LangChain integrates with open-source and proprietary services alike via a unified interface. Nearly every part (prompts, models, memory, agent logic) can be customized or replaced. This modularity facilitates using LangChain in varied contexts – for example, one could use only its vector store classes in a `/register-environment` route to ingest data, and use its agent executor in an `/implement` route to run tasks. LangChain does not prescribe an all-in-one server, so integration into a larger

system requires wiring these components, but that is straightforward. Its new LangChain v1/ LangGraph architecture is even more component-based (with a `Runnable` interface and graph execution), which helps in composing separate routes. In summary, LangChain offers maximal extensibility for custom AgentOps pipelines, at the cost of some added complexity when assembling those pieces.

- **Observability and Debugging:** LangChain provides built-in logging and tracing hooks via a **Callback system**. Developers can set verbose mode on chains/agents to see step-by-step decision traces in the console. For more advanced monitoring, LangChain introduced **LangSmith**, a toolkit and service for instrumenting LLM apps. With LangSmith (which has an open-source server option), you can capture prompts, outputs, tool calls, and token usage across runs. This integration gives a dashboard to trace agent decisions and evaluate performance. Even without LangSmith, LangChain's callbacks can be connected to external observability solutions (e.g. logging events to OpenTelemetry or custom dashboards). The framework also supports **session tracing**: for instance, previously via LangChain's built-in tracer and now via LangSmith, one can record each chain/agent invocation for debugging. In addition, LangChain's documentation includes guides on testing and debugging chains, and the new version exposes a standardized way to inspect the intermediate steps (thoughts and actions) of an agent. This emphasis on transparency makes it easier to follow what tools the agent decided to call and why. Overall, LangChain scores high on observability, especially with first-party support for tracing through LangSmith.

- **Code Efficiency & Cost Management:** Using LangChain introduces a layer of abstraction that can sometimes add overhead, but it also provides features to help manage cost. On the efficiency side, LangChain's agents and chains can produce fairly verbose prompts if not configured carefully (especially older ReAct-style agents that inlined tool descriptions every time). Newer features mitigate this: for example, using OpenAI function calling means the function schema is sent just once and subsequent calls are more compact. LangChain also supports **prompt optimization techniques** like context truncation or compression chains (there are utilities for summarizing or refining context to fit token limits). For cost control, LangChain has a caching mechanism that can **memoize LLM outputs** [1]. Developers can enable an in-memory or Redis cache so repeated queries don't hit the API again, which saves tokens on frequently asked questions. It also provides token counting utilities (and with callbacks, one can log token usage per call). There isn't a built-in "cost tracker" that tallies dollars, but since the OpenAI API returns token counts, one can easily aggregate that with LangChain callbacks. **Latency** can be slightly impacted by LangChain's overhead (each tool call and chain step is a Python function call with some serialization). In practice this is minor, but a very performance-sensitive application might require fine-tuning how LangChain executes parallel tasks or minimizing chain length. In summary, LangChain makes it simpler to implement caching and prompt management, but developers should still be diligent in designing efficient prompts. Its flexibility can either be leveraged for efficiency (e.g. plugging in a cheaper model for certain tasks, or using a vector store to limit context size) or, if misused, lead to redundant calls – so cost management largely remains in the developer's hands.

- **Ease of Development:** LangChain is both **extensively documented and widely discussed** in the community. It is the most popular of these frameworks by far, which means a large ecosystem of tutorials, examples, and community support [2]. Beginners will find many getting-started guides, and the official docs cover core concepts and how-tos for common tasks. The flip side of this breadth is that LangChain can feel **overly complex or "abstracted"** to newcomers. There is a learning curve

to understanding its terminology (chains vs. agents vs. tools, etc.) and the numerous options available. Frequent releases and API changes (especially in 2023) were a pain point, though the project has since stabilized its API with the 1.0 alpha. Seasoned developers appreciate the granular control – almost any behavior can be customized – but simpler use cases might involve a lot of boilerplate to configure. In a production context, LangChain's flexibility allows fine-tuning and optimization as you scale (which is a plus), and the large contributor base means many integrations are up-to-date. **Community activity** is a double-edged sword: issues are addressed quickly and new features roll out rapidly, but one must keep up with updates. Overall, for development speed and support, LangChain is excellent (numerous integrations ready-made), but it may require more effort to master and occasionally to refactor code when upgrading versions. It's a good choice if your team is willing to invest in learning it and you value having a large community knowledge base.

- **Security & Human-in-the-Loop:** LangChain does not inherently sandbox tool execution – if you use a tool that runs Python code or hits an external API, it's executed with the privileges of your environment. Security must therefore be handled by the developer (e.g. restricting what file paths a file tool can access, or running the agent in a container). Regarding **prompt injection**, LangChain's philosophy is to provide utilities rather than enforce policies. For instance, it has integrations like **Rebuff** and **PredictionGuard** that can be used to detect or filter prompt inputs, and it encourages using function calling (JSON schema) which inherently forces more structured, constrained outputs (reducing the risk of a model following malicious free-form instructions). However, it doesn't automatically prevent prompt injection – you must design your prompts and tools defensively. On the **human-in-the-loop (HITL)** side, LangChain makes it possible to include a human step in an agent workflow. In practice, this can be done by using a "Human" tool (where the agent can ask for user input) or by inserting manual approval steps between chain calls. The LangChain docs explicitly discuss how to add a human in the loop for tool use, indicating patterns to pause an agent and get confirmation. There isn't a built-in approval queue or UI, but the framework is flexible enough to integrate one (for example, wrapping the agent execution so that certain actions require a user click to continue). In summary, LangChain provides the hooks needed for HITL and some emerging solutions for prompt security, but it leaves enforcement to the implementer. Caution and custom safeguards are needed to use it in high-stakes or untrusted scenarios.

## LlamaIndex

- **Tool Usage & MCP Integration:** LlamaIndex (formerly GPT Index) started primarily as a retrieval and data indexing framework, but it has evolved to support agentic tool use as well. It offers a **FunctionTool** abstraction that lets you wrap a Python function as a tool, automatically inferring a JSON Schema for the function's arguments or allowing a custom schema. These tools can then be used by an agent through LlamaIndex's "Function Agent" (which leverages LLMs with function calling). In practice, defining tools in LlamaIndex is straightforward – you can group related functions into a ToolSpec and convert them to a list of tools with schema and metadata. The framework's heritage in structured data means it was designed to handle schemas and function I/O cleanly. Regarding **MCP integration**, LlamaIndex has embraced it via an extension package (`llama-index-tools-mcp`). This allows an agent to connect to any MCP server and automatically import all its tools into LlamaIndex as usable FunctionTools. For example, using a `BasicMCPClient`, you can list available tools from an MCP server and feed them into a LlamaIndex agent with one call. This gives LlamaIndex parity with others in accessing the burgeoning catalog of MCP-based tools. Additionally, LlamaIndex can itself **act as an MCP server**: the library provides utilities to expose a LlamaIndex

"Workflow" as an MCP tool/app. This means you could wrap LlamaIndex's index or query logic behind an MCP interface for external agents to use. In summary, while LlamaIndex's core focus is on data/index tools, it now supports general tool use robustly and is fully MCP-compatible (client and server) through open-source integrations. One limitation is that LlamaIndex's native agent toolkit is slightly newer and less battle-tested than LangChain's; it may require more custom logic for complex multi-step tool strategies (LlamaIndex tends to encourage a single-step tool call as part of a larger query plan). But for most cases, it handles tool selection and JSON I/O effectively, and its MCP connector ensures it can interface with external tool sources similar to the others.

- **State Management & Memory (Qdrant Integration):** Persistence of conversation state in LlamaIndex is typically achieved by leveraging its data structures. It doesn't have "conversation memory" classes in the same way LangChain does; instead you can log chat history in an index (for example, appending conversation turns to a list or embedding them into a vector store and querying for context). In practice, a developer might use LlamaIndex's **indices** to serve as memory – e.g., a **SummaryIndex** to summarize past interactions or a **KnowledgeGraphIndex** to store facts gleaned. For more direct approaches, LlamaIndex supports vector stores for storing text, and it integrates with **Qdrant** smoothly. There is an official `llama-index-vector-stores-qdrant` package enabling Qdrant as a backend for LlamaIndex's VectorStoreIndex. With it, you can use Qdrant to persist nodes (text chunks or messages) and query them later by embedding similarity. This is analogous to using Qdrant as long-term memory: after each user message, the conversation or relevant info could be added as a node to the index, and the next query can retrieve the most relevant past nodes. The integration is production-ready – simply provide a Qdrant client to LlamaIndex's `QdrantVectorStore` wrapper and use it in a VectorStoreIndex. For **long-running tasks**, LlamaIndex doesn't explicitly handle agent session continuation; an agent run is expected to complete or perhaps ask the user for input. However, LlamaIndex's new **Workflow** feature (event-based workflows) can allow complex sequences and even waiting for events. In theory, one could create a Workflow that spans a long operation and yields intermediate events (which a human or scheduler could act on), but that is an advanced usage. Out-of-the-box, preserving agent state beyond a single query/response cycle would require manually storing needed data (e.g., in a vector store or some cache) and reinitializing the agent with that context. LlamaIndex is quite flexible with how you manage context – since you have full control to feed the prompt builder whatever text you want, you can include summary of previous chats, etc. In conclusion, LlamaIndex can integrate with persistent stores like Qdrant for memory, but it doesn't have as many built-in conversational memory conveniences. It gives you the tools to implement memory and state persistence, but the responsibility is on the developer to wire them together for an AgentOps scenario (for example, storing conversation turns in an index at each step).

- **Modularity and Extensibility:** LlamaIndex is designed as a collection of **independent modules** (indices, query engines, routers, tools) that you can mix and match. It is quite modular, though in a slightly different way than LangChain. Rather than a plethora of model and tool integrations, LlamaIndex's core strength is being modular with data **indices and retrievers** – you can plug in different vector stores, LLMs, or chunkers. It supports numerous **LLM backends** (OpenAI, Azure, local HuggingFace models, etc.) via a simple LLM interface, so you're not tied to any vendor. Similarly, for storage, it works with in-memory stores or external vector databases (we saw Qdrant, and others like Pinecone, Weaviate, FAISS are supported). This avoids vendor lock-in at the data layer. LlamaIndex also plays well with other frameworks: for example, some developers use LlamaIndex as the retrieval component inside a LangChain agent – the two can interoperate via simple function

calls or shared vector DBs. In an AgentOps platform context, one might use LlamaIndex for the "data ingestion and querying" routes and another agent framework for higher-level orchestration. Because LlamaIndex is essentially a toolkit, you can invoke its components from any part of your system. It doesn't impose a particular server or API structure; you can embed it in a FastAPI route or a background process as needed. This means you could have a `/register-environment` endpoint that uses LlamaIndex to index new data sources (e.g., ingesting documents into a vector index), and another `/implement` endpoint where an agent uses that index to answer questions. The only caution is that LlamaIndex, being very flexible with fewer opinions, might require more careful structuring by the developer to keep things organized (as noted in a comparative analysis: LlamaIndex allows "more decisions" and varying code structure, whereas something like Haystack enforces a consistent pattern). Extending LlamaIndex – e.g., adding a new custom tool or a new type of index – is generally straightforward via subclassing or using their callback system. The framework doesn't lock you into pre-baked implementations; it's intended to be a layer you build your own system on. In summary, LlamaIndex is highly extensible and can be integrated in pieces, which is great for avoiding lock-in, but it may lack some higher-level abstractions to tie those pieces together (the developer must ensure the modular parts work in concert).

- **Observability and Debugging:** As a lighter-weight framework, LlamaIndex did not initially have a built-in UI for tracing, but it provides robust **callback hooks** for logging and debugging. Using the `CallbackManager`, you can register multiple callbacks to receive events during index construction or query execution. These events include things like LLM prompts, tool invocations, and intermediate results. This system allows integration with external observability tools. Indeed, LlamaIndex has documented integrations for **Langfuse**, **OpenLLMmetry, Arize Phoenix, and Opik** among others. For example, you can attach a Langfuse callback handler to automatically log prompts and completions to a Langfuse dashboard (Langfuse is an open-source tracing/monitoring platform similar to LangSmith). There's also a PromptLayer integration to log and version prompts, and a token counting handler to measure tokens used. In essence, while LlamaIndex doesn't have its own proprietary monitoring tool, it plays nicely with the ecosystem of LLM observability solutions. Debugging LlamaIndex workflows can be done by enabling verbose mode when running queries or agents, which will print out step-by-step information (e.g., what query was sent to the vector store, what answer came back, how the LLM parsed the tool output, etc.). The **community support** also helps here – for instance, the developers often provide tips on using the trace events to diagnose issues (such as verifying that the tool's schema was correct or that the retrieved context is relevant). Traceability of agent decisions in LlamaIndex is improving as the agent tooling matures: when using the FunctionAgent, you can inspect the function call outputs returned by the LLM to see which tool was chosen and with what arguments. The **transparency is generally good** (since it often relies on OpenAI function calling, the model's tool choice is explicit in a JSON payload), but the framework doesn't automatically present it in a UI. To summarize, LlamaIndex provides the hooks needed for observability and has multiple third-party integrations to capture traces, but users may need to set these up. It's a more DIY approach compared to LangChain's one-stop LangSmith, but it achieves comparable insight when configured.

- **Code Efficiency & Cost Management:** LlamaIndex's design focus on **retrieval augmentation** inherently promotes token efficiency – it tries to send only relevant data to the LLM. For example, rather than stuffing an entire document into a prompt, it will retrieve just the top $k$ relevant chunks (nodes) from an index to answer a query. This leads to smaller prompt sizes compared to naive approaches. LlamaIndex also provides tools for **chunking and optimizing context**; there are

strategies like recursive summarization or using tree indices to compress information when dealing with very long documents. These can help keep token counts low by summarizing sections of data hierarchically. In terms of framework overhead, LlamaIndex is relatively lightweight. It doesn't impose complex control flows at runtime – a typical query might involve an embedding lookup and one LLM call, which is efficient. That said, if you use some of its more advanced compositional techniques (like multi-step query plans or the auto-refinement of answers), it can invoke the LLM multiple times. Those strategies are under developer control. For **caching**, LlamaIndex doesn't have a built-in caching module like LangChain's LLMCache, but you can achieve similar results by caching the outputs of your index queries or LLM calls manually. Some users integrate it with an external cache (for instance, wrapping the LLM class to store responses keyed by prompt). The framework's callback system could also be used to plug in a caching layer (by intercepting an LLM call event and returning a stored result if available). On cost tracking, LlamaIndex can leverage its callback hooks for token counting – the library even provides a TokenCountingHandler which accumulates token usage for each run. This helps in analyzing cost per query. **Latency-wise**, LlamaIndex operations are quite optimized: vector searches are offloaded to the vector DB (with network overhead if remote), and the rest is just Python managing data structures. There is minimal extra latency added by LlamaIndex itself. One should note that the **most significant costs** (financial and time) still come from the LLM calls, which LlamaIndex helps minimize by smart retrieval. In summary, LlamaIndex is efficient by design for RAG use cases. It avoids redundant LLM calls by first querying data and only then prompting the model. Developers still need to be mindful of how many tools or intermediate steps they add (each tool call could be an LLM function call, for instance), but the framework doesn't encourage wasteful patterns. With some custom caching and the provided token counting, you can manage and reduce API usage effectively.

- **Ease of Development:** LlamaIndex is generally regarded as **user-friendly for retrieval tasks** and has thorough documentation. Its API is Pythonic and often less verbose than LangChain's. For example, creating an index and querying it can be done in a few lines, and the library abstracts a lot of the prompt engineering for you (it auto-formats queries to include retrieved text, etc.). The learning curve for basic operations (like indexing documents and asking questions) is relatively gentle. For agent and tool usage, the complexity increases a bit – you need to familiarize yourself with ToolSpec, the FunctionAgent, and how to define tool metadata. This part of LlamaIndex is newer, and documentation is still catching up, but there are examples in the official docs and community blogs. **Community support** for LlamaIndex, while not as massive as LangChain's, is quite solid. The project is open-source (with permissive license) and active on Discord/forums. According to some comparisons, LlamaIndex has a larger contributor community than Haystack, but smaller than LangChain. In practical terms, this means you'll find a decent number of extensions (community-contributed connectors in the LlamaHub, for instance) and people to ask for help, though fewer third-party tutorials than LangChain. One advantage of LlamaIndex's focus is that if your use case is primarily about question-answering over data, it feels very natural and quick to implement. Developers often praise it for *"just working"* for RAG without needing to tinker with prompt templates too much. However, this opinionated ease can turn into needed customization when you step outside the common path – e.g. building a complex multi-step agent logic might require diving into their workflow system or combining LlamaIndex with another framework. In terms of maintenance, LlamaIndex is under active development (with frequent releases around 2023–2024), which introduced some breaking changes (for example, renaming classes or moving modules as the design evolved). The team is fairly small but dedicated, led by the original author. There's an emphasis on stability now as it matures. One should be prepared to update code if

upgrading the library version, but the pace has slowed as the core has solidified. Summing up, for development ease: **for retrieval-centric applications LlamaIndex is excellent** – it feels tailored to that job – and for broader agentic applications it's improving. It might require a bit more assembly for full agent systems than LangChain, but it's very doable. The documentation (both official and community) will guide you through most tasks, and the community, while smaller, is active and helpful. Many developers find LlamaIndex a nice balance between high-level abstraction and low-level control.

- **Security & Human-in-the-Loop:** Since LlamaIndex often runs in a context where an LLM is reading your data, it puts some focus on data access control but not as much on tool misuse, simply because it historically didn't execute arbitrary code by itself. The **tools feature** now introduces similar security considerations as LangChain's – if you give the agent a tool that can perform destructive actions (e.g., write to a database or send an email), you must ensure those are safe. LlamaIndex does not provide a sandbox; any Python function you expose as a tool will run with your process's privileges. So similar best practices apply: limit what your tool functions do, use environment isolation if needed, and validate inputs. On the prompt side, prompt injection is a concern especially if user queries are fed directly into prompts (which they are in retrieval: a malicious user could include instructions in a question that the model might follow). LlamaIndex doesn't have a built-in prompt firewall, but you can intercept the user query or the assembled prompt via callbacks to sanitize it. Developers might use external libraries (like the OWASP guidelines or Rebuff) in conjunction with LlamaIndex to detect injection attempts before calling the LLM. It's also possible to leverage the fact that LlamaIndex uses structured function calling for tools – if the model tries to deviate from the JSON schema, you could catch that and refuse execution. However, these safeguards are not automatically in place; they require implementation. As for **human-in-the-loop capabilities**, LlamaIndex does not have an out-of-the-box mechanism for pausing an agent and asking for human approval, because it's primarily a library rather than a full workflow engine. That said, one could design the agent such that certain tools are effectively "human" tools – for instance, you might implement a tool function that actually sends a notification to a human or queries a human-in-the-loop service, then returns the result to the agent. Another approach is using the event loop (Workflow) in LlamaIndex: you could define an event that represents a required human approval, which stops the workflow until the human provides input. This is more of a custom solution. There isn't a built-in UI for a human to oversee each agent step, so integrating HITL likely means building a simple interface or script around the LlamaIndex agent. To summarize, LlamaIndex currently relies on developers to handle security and HITL concerns. It provides the flexibility to insert human steps or validations (thanks to its callback system and the ability to craft your own control flow), but it doesn't come with guardrails by default. In a production AgentOps platform, one would need to add those guardrails on top of LlamaIndex – for example, whitelisting which tools an agent can call in certain contexts, sanitizing user inputs, and routing certain high-risk queries to a human operator. The framework won't fight you on implementing these, but it won't do it for you either.

## Haystack

- **Tool Usage & MCP Integration:** Haystack by deepset historically focused on pipelines for search and QA, but it has incorporated **agentic tools** as LLM capabilities grew. Haystack's agent interface allows tools to be defined as **Components** in a pipeline or as external calls. Notably, Haystack has first-class support for the **Model Context Protocol** via its **MCPTool** integration. In Haystack, `MCPTool` is a special tool that lets an agent communicate with any MCP-compatible external server.

With this, Haystack agents can invoke a huge range of tools standardized by MCP (from web browsers to databases) as if they were native tools. The integration supports multiple transports (streamable HTTP, SSE, and even launching local MCP servers via stdio). Using it is straightforward: you install the `mcp-haystack` package, then create an `MCPTool` pointing to the server URL or local command, and include that in your agent's tool list. The agent then can call MCP tools and receive JSON results, fully aligning with MCP's schema-based interface. This deep integration means Haystack is essentially **MCP-ready out-of-the-box** – a considerable advantage for an AgentOps platform focused on MCP compliance. Aside from MCP, Haystack also supports tools via OpenAI's function calling. In fact, you can wrap a Haystack pipeline or function as an OpenAI function and use it in a chat model call. There are guides on converting Haystack pipelines into tools that GPT-4 can call. Under the hood, Haystack will automatically generate a JSON schema for a pipeline if you deploy it as a tool. For example, with **Haystack's Hayhooks server** (an API wrapper for pipelines), any deployed pipeline gets an `inputSchema` (JSON Schema for its parameters) and can be exposed as an MCP tool. This means if you define a pipeline (or agent) that answers questions, Haystack can not only use external MCP tools, but it can also *be* a tool server to others via MCP [3] . In summary, Haystack's tool usage is very much aligned with **structured, schema-based definitions** (whether via OpenAPI, MCP, or its own YAML pipeline specs). It easily interfaces with multiple tool sources: built-in Components (like search, transcribe, etc.), user-defined Python functions, other pipelines, and remote MCP tools. One slight difference is that Haystack encourages a more **planned pipeline** approach (like a DAG of steps) as a "tool," which can be more deterministic compared to free-form agent loops. But Haystack also provides a flexible Agent class where the LLM decides which tool to use next (similar to LangChain's agent loop). With the integration of MCP and its existing library of Connectors (e.g., GitHub, Slack connectors in the docs), Haystack is extremely powerful for tool use in a standardized way.

- **State Management & Memory (Qdrant Integration):** Haystack's original use-case was stateless question answering, but it has components to manage state when needed. For conversation history, Haystack doesn't have an explicit "ConversationMemory" object like LangChain; instead, you typically feed the prior messages as part of the prompt (often handled by a **ChatPromptBuilder** component). There is a concept of a **MessageHistory** in the ChatPromptBuilder that you can update with each turn, enabling the LLM to see previous user and assistant messages. In terms of long-term memory or vector-based memory, Haystack integrates very well with vector databases. It supports a variety of **DocumentStore** backends including **Qdrant** natively. You can use a QdrantDocumentStore to persist documents (or chat transcripts) and then employ a **Retriever** to fetch relevant bits given the latest user query. Indeed, Haystack provides Qdrant-specific retrievers (for dense, sparse, or hybrid search), showing a mature integration. This means implementing an episodic memory in Haystack can be done by writing each conversation turn as a Document into Qdrant and using a retriever to pull the most relevant past turns when answering new queries. It's not packaged as a one-liner memory class, but all the pieces are there. Regarding **long-running tasks** and persistence: Haystack pipelines can be persisted (serialized to YAML) and re-loaded, which is more about pipeline definitions than runtime state. For runtime agent state, Haystack doesn't automatically persist an agent's chain-of-thought – each query to an agent is handled anew, albeit potentially with context passed along. However, the **Hayhooks MCP server** does manage user sessions: when acting as an MCP server, Haystack (Hayhooks) uses a 5-minute TTL for sessions and keeps a Redis pub/sub to route messages [4] [5] . This is more for horizontal scaling and auth, but it implies that in a multi-turn interaction through MCP, session state (like which tools have been listed or used) is maintained for that session. In a simpler scenario, one could also store agent intermediate results or plan in a

database if one needed to pause and resume – but that's custom work. Another aspect of state is **long-running tool actions** (e.g., a tool that streams back results). Haystack's design with pipeline nodes supports streaming responses, and using the MCPTool with SSE or streamable HTTP means the agent can handle streaming tool outputs in real-time, which is good for responsiveness in long operations. In conclusion, Haystack is well-equipped for memory via its document stores and retrievers (excellent Qdrant support), but like LlamaIndex it might require you to wire the conversation saving and retrieval logic explicitly. It's very much capable of persistent knowledge bases; persistent conversational state can be achieved by treating conversation as data. Meanwhile, built-in session management in Hayhooks addresses multi-request continuity in an API setting.

- **Modularity and Extensibility:** Modularity is a core design principle of Haystack. The framework is built around **Components** and **Pipelines** – every functionality (reader, retriever, generator, agent, etc.) is a component that can be swapped out. You can, for instance, drop in a different LLM (ChatGPT, Claude, local models via the `OpenAIChatGenerator` or others listed) without changing your pipeline structure. Swapping vector stores or embedding models is similarly easy since Haystack provides a common interface for them and many implementations (from FAISS to Qdrant to Weaviate). This ensures **no vendor lock-in** – you can start with an in-memory store and later move to Qdrant or Elasticsearch by just changing the document store initialization. The framework is also language-agnostic regarding models: it supports OpenAI APIs, Hugging Face models, Cohere, Azure, and more through its Generator classes. One area Haystack shines is enabling **custom components**. If you have a specific business logic step (say a proprietary API call or a custom text transformation), you can create your own Component class and insert it into a pipeline. The rest of the pipeline will treat it like any other node. This is extremely useful for an AgentOps platform where you might want custom preprocessing or postprocessing steps. Moreover, Haystack's pipelines can be defined via **YAML** or Python, and even exported/imported, which is great for managing configurations in a modular way. For example, you could define a `/register-environment` pipeline in YAML that takes new documents and indexes them, and a separate `/implement` pipeline (or agent) that retrieves answers – each can be developed and maintained independently, then connected via shared resources (like the DocumentStore). The **avoidance of tight coupling** extends to how Haystack can integrate into systems: you can run it as a service (with Haystack's REST API or Hayhooks) or embed it in your own application logic. It doesn't force you into a specific server architecture. Also noteworthy, Haystack can interoperate with other frameworks if needed; for instance, one could use Haystack's retriever inside a LangChain tool or vice versa, since at bottom they're Python calls (this is occasionally done by advanced users). Haystack's extensibility is also evident in its **Plugin-like connectors** (GitHub, Slack, etc.), which are optional modules you can include for extra capabilities. Finally, the commitment to MCP indicates a forward-looking extensibility – rather than building every possible tool integration in-house, Haystack can just call an MCP server, which means it can leverage any new tool community develops in the MCP ecosystem without needing a code update. In summary, Haystack provides a highly modular, lego-blocks approach that's ideal for a production platform. Each part of your agent system can be independently scaled or replaced. The potential downside is that if you want an all-in-one solution, Haystack might feel too low-level (you have to design the pipeline); however, for an AgentOps platform this is actually an advantage, as you can tailor each route and component without bending the framework's intent.

- **Observability and Debugging:** Haystack, being geared to enterprise and production use, has invested in observability features. Out of the box, it offers **logging at each pipeline node** (you can

configure logging levels to see inputs/outputs of components). There's also a **Pipeline Debugging** interface: when running a pipeline in Python, you can enable debug mode to get the output of each node step-by-step (or use the interactive debugging in a notebook to inspect intermediate results) [6] . More formally, Haystack introduced **Pipeline "Breakpoints"** which allow you to pause a pipeline at certain nodes and examine or modify the output before continuing. This can be used during development to ensure each component is doing what you expect. For example, if you have an agent pipeline and want to verify the tool inputs it's generating, you could set a breakpoint right before the ToolInvoker. For tracing in production, Haystack provides integration with **Langfuse** (via a `LangfuseConnector` ). By configuring that, every request and its stages can be logged to a Langfuse instance, giving you a UI to analyze conversations, errors, and latencies. The inclusion of Langfuse shows that Haystack acknowledges the importance of structured traces for LLM applications. Additionally, when using the OpenAIChatGenerator with function calling, you inherently get a structured record of function calls which can be logged or examined. Haystack also can **visualize pipelines**: there is a tool that can generate a graph diagram of your pipeline for documentation and understanding (useful when debugging complex flows of components). In terms of agent-specific observability, Haystack's Agent can be run in verbose mode where it will print each thought and action choice of the LLM as it decides on tools. This is analogous to LangChain's verbose agent logging. Since many Haystack agents will often use the OpenAI function-calling mode, the model's "thought" process is somewhat hidden inside the model (it directly outputs a function name and args). However, you can still log the model's raw response for inspection if needed. The **error handling** in Haystack pipelines is also noteworthy: you can configure what happens if a node fails (skip, retry, etc.), and with the upcoming support for asynchronous pipelines, you can gather error traces more easily. All these features contribute to a high degree of transparency and debuggability. In practice, users often praise Haystack for being easier to debug than highly dynamic frameworks, because a lot of the behavior is explicitly laid out in the pipeline configuration. The consistent patterns mean fewer surprises – which is a form of observability through design. Finally, Haystack has a **testing utility** where you can simulate pipeline inputs and verify outputs (useful for regression tests on your agent logic). Overall, Haystack provides strong observability, combining built-in debugging tools with hooks for external monitoring. This is crucial for a production platform where you need to trace how an answer was arrived at (for compliance or error analysis).

• **Code Efficiency & Cost Management:** Haystack is built with **performance and scalability** in mind, as it targets production deployments. One of its advantages is that it often uses **pipeline orchestration instead of sequential prompting loops**, which can be more efficient. For example, if you have a RAG pipeline in Haystack, the retrieval and reading happen in a defined flow without unnecessary round-trips to the LLM. Haystack's components are generally optimized (many use vectorized operations or efficient batch processing where possible). It also supports asynchronous pipelines, so components (especially slow ones like LLM calls) can run concurrently if the pipeline is set up that way – this can improve throughput and latency in multi-query scenarios. On the token usage front, Haystack doesn't explicitly alter the model's prompts beyond what is needed. When using function calling, it simply provides the tool schema and lets the model produce tool outputs or answers, keeping prompts tight. For standard QA pipelines, Haystack's prompt templates for readers/generators are concise and focused (and you can customize them as needed). **Caching** is an area where Haystack provides concrete tools: it includes a `BaseCachingWrapper` and `OutputCaching` mechanism that can cache the outputs of any node given its inputs. A specific component, the **CacheChecker**, can be inserted at the top of a pipeline to check if a question was seen before and retrieve the stored answer if so. This is very useful to reduce duplicate work and

cost. You could, for instance, cache answers to common queries or cache document embeddings to avoid recomputation. Additionally, Haystack's DocumentStores (like Qdrant, FAISS) serve as a cache of knowledge that doesn't require hitting external APIs for known info – effectively acting as a cost-saving by answering from memory rather than an LLM when possible. In terms of overhead, the framework adds minimal latency on top of the actual operations. It's essentially function calls passing data between components; it's written in efficient Python and sometimes uses multiprocessing for heavy tasks (like embedding large corpora). If extremely low latency is needed, one can always deploy Haystack's pipelines as persistent processes (so you avoid startup time) – this is exactly what Haystack's REST API or Hayhooks allows. For **cost tracking**, while Haystack doesn't print a token count by default, you can leverage OpenAI's response metadata or use a callback on the generator to accumulate usage. Enterprise users often integrate Haystack with their own monitoring to tally API calls. Because Haystack pipelines are predictable, it's also easier to estimate costs (for example, you know a query will do X embeddings and 1 generation, so you can calculate worst-case tokens). In summary, Haystack is very efficient in how it structures LLM workflows and offers built-in caching to cut down on redundant calls. It imposes little overhead and is suitable for high-throughput scenarios. Proper use of its caching and pipeline branching (e.g., using a router to only call the LLM when necessary) can significantly reduce token consumption – which directly translates to cost savings.

- **Ease of Development:** Haystack might not be as instantly famous as LangChain, but it has a well-earned reputation for being **developer-friendly in production settings**. Its documentation is comprehensive, with a concept overview and many tutorials (for instance, "Build a Tool-Calling Agent", "Multi-Agent System with Haystack", etc., are provided step-by-step on the official site). The learning curve for simple QA is low: you can set up a pipeline with a few lines of code. For agentic use, it's slightly higher because you need to understand how to structure a pipeline or Agent schema. However, once you grasp the pipeline concept, it actually *simplifies* development – you have a clear map of what happens first, next, and so on. This opinionated structure ensures consistent patterns across projects, making it easier for teams to collaborate and for code to be maintainable. In contrast to LangChain's very flexible approach, Haystack's approach may feel a bit rigid but it prevents a lot of errors (for example, you explicitly connect components, so it's harder to have a missing piece that you forgot to handle). For a developer building an AgentOps platform, this consistency is a boon: each new tool integration or pipeline can follow a template. The **community and maintenance** of Haystack is smaller than LangChain's, but it's quite mature – the project dates back to 2019 and has been used in industry for years. The current active community (on GitHub, Discord) is smaller, which means fewer third-party plugins but also that the core maintainers oversee most contributions for quality. According to one comparison, Haystack's curated ecosystem yields high-quality components, whereas a larger community like LlamaIndex's can have more varied quality. There were some questions raised about Haystack's long-term sustainability given its smaller community, but deepset (the company behind it) actively maintains it and frequently releases updates (the 2.x series throughout 2023 and 2024 added many LLM features). In practical terms, if you encounter an issue, the maintainers are responsive on GitHub. The framework's API has remained fairly stable, with deprecation warnings given well in advance for major changes – this enterprise-style stability is appreciated in production environments. Getting started is made easier by Haystack's out-of-the-box examples and the fact that it can run in a managed way (deepset offers a hosted Haystack, and even locally you can launch everything via Docker). Another point in ease-of-use is that **Haystack integrates with FastAPI** (via Haystack REST API) and with Ray for scaling, which saves developers time if they want to deploy an agent as a service. On the downside, if a

developer comes from the LangChain world expecting a one-liner to create an agent that can do everything, they might find Haystack requires more upfront configuration. But that initial effort corresponds to clarity. Summing up, Haystack is very developer-friendly for those building robust systems. It may have a bit more setup and learning at first (especially understanding YAML pipeline definitions and how to deploy pipelines), but that investment pays off in maintainability. The documentation and tutorials smooth this process, and the framework's longevity means many corner cases have been addressed. For an AgentOps platform team, using Haystack can impose a healthy discipline that ultimately speeds up development of complex features.

- **Security & Human-in-the-Loop (HITL):** Haystack in its core is secure-by-structure in some ways: because you often define exactly which components (tools) are used and in what order, there's less risk of an agent going off-script. However, when using the more open-ended Agent with an LLM deciding tools, similar security caveats apply as with any agent. **Prompt injection** can target an agent if a user inputs something like "ignore previous instructions". If using OpenAI function calling, the model is constrained to return either a function call or an answer, which mitigates some injection routes (it can't execute arbitrary code except through allowed tools). Haystack doesn't have a built-in prompt sanitizer, so implementing checks (like disallowing certain patterns in user prompts or verifying the model's chosen tool is expected) would be up to the developer. One could insert a custom component at the start of a pipeline to perform input validation or injection detection. On **tool execution sandboxing**, Haystack again assumes you trust the components you include. For example, if you use the `PythonNode` (a component to execute Python code from an LLM input), you have to be extremely careful (this is analogous to giving an agent a Python tool in LangChain). The recommendation would be to run such capability in a sandboxed environment (perhaps a Docker container or restricted Python exec). Haystack itself doesn't enforce a sandbox, but because pipelines are deployed typically in controlled servers, you have the opportunity to isolate at a higher level. It's worth noting that Haystack's approach of often **separating retrieval and generation** can act as a safety net – the LLM is usually only generating text given retrieved info, rather than directly interacting with the world unless you explicitly give it that power. For **HITL support**, Haystack has interesting features: its **Pipeline Breakpoints** can be considered a form of HITL checkpoint, albeit for debugging. You could repurpose that concept – for instance, insert a breakpoint after the LLM produces an answer but before returning to the user, then have a human reviewer verify the answer if it's a sensitive query. There's also a notion of **feedback gathering** in some Haystack utils (to collect thumbs-up/down on answers, etc.), which can be used to loop a human's feedback into system improvement. If a truly interactive HITL is needed (where a human and agent collaborate in real time), one might integrate Haystack with a chat UI and allow the human to be another "agent" in the loop. While not a one-click feature, Haystack's component system could allow a "HumanInput" component that waits for human confirmation at a certain stage. Given that Haystack is often used in enterprise settings, a common pattern is to start with a fully automated pipeline and later add a human review step for only certain outputs (e.g., if the confidence score of an answer is below a threshold, forward it to a human). Haystack supports confidence scoring especially for QA (via reader model scores), which can trigger HITL interventions. In the agent scenario, one could similarly use the LLM's function call confidence or some heuristic to decide when to ask a person. To summarize, Haystack doesn't deliver turnkey HITL interfaces, but its structured pipelines actually make it easier to slot in a human step. Security-wise, it provides all the knobs but expects you to turn them as needed: you have fine control over tool availability and can set up approval flows, but must implement checks yourself. Its alignment with MCP (which has an authorization model for tool servers) suggests that in an MCP-enabled deployment, you could leverage MCP's auth and access

controls for tool usage (for example, requiring certain tokens/permissions for an agent to call a particular tool server). That can be part of a security strategy when using Haystack + MCP. In essence, Haystack is as secure as the design you put into it – it won't do unsafe things unless you configure it to (no surprise internet browsing unless you include that tool), and it gives you opportunities to intercept and control the flow for safety or human oversight.

## Scoring Matrix (1–10 Scale by Criterion)

| Criterion | LangChain | LlamaIndex | Haystack |
|---|---|---|---|
| Tool Usage & MCP Integration | 9 – Excellent | 8 – Very Good | 10 – Outstanding |
| State Management & Memory | 8 – Robust | 7 – Adequate | 8 – Robust [7] |
| Modularity & Extensibility | 10 – Outstanding | 8 – Very Good | 9 – Excellent |
| Observability & Debugging | 9 – Excellent | 8 – Very Good | 9 – Excellent [6] |
| Code Efficiency & Cost Management | 7 – Good [1] | 8 – Very Good | 9 – Excellent |
| Ease of Development | 8 – Very Good [2] | 7 – Good | 8 – Very Good |
| Security & HITL | 7 – Good | 6 – Fair | 7 – Good |

*(Scores are relative assessments: 1 = very poor, 10 = excellent. Scores in each row are accompanied by brief justifications and references.)*

## Final Recommendation and Trade-offs

**Recommended Framework:** After thorough analysis, **Haystack** emerges as the best foundational framework for an AgentOps platform emphasizing MCP compatibility. Haystack's native integration with the Model Context Protocol is unparalleled – it treats MCP tools as first-class citizens and can even serve pipelines as MCP-accessible tools. This means your platform can immediately tap into the growing MCP ecosystem (hundreds of standardized tools) and expose its own capabilities via MCP, with minimal custom glue. Haystack also excels in production-oriented features: robust state management with vector stores, highly modular design, strong observability, and built-in caching and pipeline control for efficiency. Its more opinionated, pipeline-centric approach yields consistency and easier maintenance in a complex multi-agent system. In a setting where reliability and clarity are paramount (typical of AgentOps scenarios), these qualities make Haystack a solid foundation.

**Trade-offs and Rationale:** While LangChain scored equally high or higher in some areas (particularly tooling flexibility and community support), we weigh the criteria differently for a production platform:

- **MCP Integration:** This was a critical factor. Haystack's MCP support is **production-ready and deeply integrated**, requiring almost no extra work to connect to or host MCP servers. LangChain and LlamaIndex have MCP connectors as well, but these are add-ons rather than core features (LangChain's adapter is recent and primarily for LangChain v1, LlamaIndex's is community-driven). Haystack's out-of-the-box MCPTool and Hayhooks MCP server give it a clear edge for an MCP-centric platform.

- **Production Maturity:** Haystack is a mature framework (since 2017) with a track record in enterprise QA systems. It emphasizes stable APIs and backward compatibility. LangChain, in contrast, underwent rapid, sometimes breaking, changes during its explosive growth; it has stabilized recently, but a platform built on LangChain might need more vigilance for updates. LlamaIndex lies somewhere in between: younger than Haystack but past the early churn stage, though its agent capabilities are still maturing. For longevity and predictability, Haystack's slower, quality-focused development cadence is a safer bet – an important consideration for AgentOps, where you'll be maintaining the platform over time.

- **Observability & Debugging:** All three frameworks support debugging, but Haystack's approach (structured pipelines, breakpoint debugging, Langfuse integration) makes it straightforward to trace issues in a multi-step agent process [6]. LangChain offers similar tracing via LangSmith, but that's an external service/package and comes with some complexity in setup. LlamaIndex relies on third-party callbacks for tracing. In a platform context, having baked-in, easily accessible observability (as Haystack does) is a big advantage for monitoring agent behavior and quickly diagnosing problems in production.

- **Modularity & Extensibility:** All three are modular, but their philosophies differ. **LangChain** provides a lego-box of pieces you can wire in countless ways, which is powerful but can lead to ad-hoc implementations. **Haystack's modularity** is more guided – it encourages you to assemble components into a coherent pipeline or agent, and it handles a lot of the wiring (through YAML or Pipeline API). This guided modularity reduces the chance of misconfiguration and eases onboarding new team members to the codebase (they'll find standardized pipeline definitions rather than a tangle of custom chain logic). **LlamaIndex** is also flexible, especially for swapping data stores or LLMs, and could be integrated as a subsystem within either LangChain or Haystack (for example, using LlamaIndex for document indexing inside a Haystack tool). But by itself, LlamaIndex is less of a full-stack agent framework and more of a specialized retrieval toolkit, which might require supplementing with custom code for other agent capabilities. For an **AgentOps platform**, which likely involves multiple subsystems (data ingestion, tool integration, conversation handling, etc.), Haystack's clearly defined module boundaries and pipeline DSL make it easier to maintain separation of concerns.

- **Security & HITL:** None of the frameworks offer plug-and-play human approval flows or comprehensive security by default, so these will need to be implemented on top no matter what. However, Haystack's structure again aids here: because you can insert custom components at any point, adding an HITL checkpoint is relatively clean (e.g., a component that requires a human to approve before proceeding). LangChain also allows this but in a less structured way (likely via a custom loop or callback halt). The difference is subtle, but designing a *systematic* HITL mechanism (like: "for certain tool outputs, pause and request human review") feels more straightforward in Haystack's pipeline paradigm. On prompt injection, all frameworks require custom handling; LangChain has some emerging tools (Rebuff, etc.), which can be integrated into Haystack too (for instance, one could use Rebuff as a pre-processor component). So in terms of security, no framework solves it outright, but Haystack's emphasis on *explicit flows* means you can more easily reason about where to put guardrails.

- **Community and Ecosystem:** LangChain undeniably has the largest community and the most integrations available off-the-shelf. This means if you need a connector to an obscure API or a ready-

made chain for a use-case, LangChain might have it. Haystack's ecosystem is smaller and curated – most essential integrations are there (search, databases, etc.), and of course MCP access fills many gaps, but you might occasionally need to write a custom component where LangChain would have had a template. LlamaIndex's community is growing, and it particularly shines in retrieval and data connectors (via LlamaHub). The question is, how important is the broad community for your platform? Given that MCP can act as a bridge to many external tools/services, the need for built-in integrations is lessened – you can focus on core agent logic and use MCP servers for specialized tools as needed. LangChain's huge community also comes with high velocity; for a production platform, you might favor stability over getting every cutting-edge feature first.

**Feasibility of Combining Components:** It's worth noting that this decision is not all-or-nothing. **Combining strengths** of frameworks is a viable strategy and is done in practice. For example: - You could use **LlamaIndex for data indexing and querying** within a Haystack agent. LlamaIndex can build a complex index (say a hierarchical index or a graph of knowledge) that Haystack alone would require more manual setup to replicate. You can either incorporate LlamaIndex as a sub-component (via a custom Haystack node that calls LlamaIndex) or expose LlamaIndex as an MCP server (there's an example of a LlamaIndex documentation server accessible via MCP) and call it from Haystack's MCPTool. This way, you leverage LlamaIndex's efficient retrieval and Haystack's orchestration around it. - It's also possible to **use LangChain in specific subsystems**. For instance, if LangChain has a very convenient agent for a niche use-case (say an agent that interacts with spreadsheets or a particular API), you could run that agent in isolation and interface it with the rest of the platform via MCP or REST. Given LangChain's focus on orchestration, there's some overlap with Haystack – you likely wouldn't run a LangChain agent inside a Haystack agent concurrently, but you might have separate processes or services, one using LangChain for a specialized task and another using Haystack for the main workflow, communicating through an API. - **MCP as an interoperability layer** makes combining even simpler: you could have LangChain serve as an MCP client or server (with its adapter) and Haystack as another, allowing them to call each other's tools in a standardized way. For example, a LangChain agent could call a Haystack pipeline via MCP if that pipeline is published as a tool, or vice versa. This kind of hybrid might sound complex, but if each framework is used for what it's best at, the overall system can be quite elegant.

In terms of team skills and maintenance, mixing frameworks means needing familiarity with all, so one should only combine if there's a clear payoff. A realistic approach could be: **use Haystack as the primary framework**, given its strengths in MCP and production deployment, and **augment with LlamaIndex** for any advanced retrieval/indexing needs that exceed Haystack's built-in capabilities (since LlamaIndex was built for that purpose). LangChain could be tapped for any missing tool integrations or prototyping new agent ideas, but one might gradually replace those with native Haystack or MCP tools for maintainability.

**Conclusion:** For an AgentOps platform aiming to be *tool-agnostic, scalable, and maintainable*, Haystack provides the most solid foundation, especially in an MCP-centric world. It gives you structured power: the agents and tools behave predictably, you can monitor and tweak them in production, and you can interface with external systems cleanly (MCP, API endpoints, etc.). LangChain is a close runner-up – it's incredibly powerful and might even be preferable for fast prototyping or if you need a very broad array of ready-made integrations. Its huge community means it won't steer you wrong on popular use-cases, but you should be prepared to manage the complexity and ensure any MCP integration is kept up to date with the adapters. LlamaIndex is somewhat a different category, excellent for what it does (retrieval and data-centric agents) and it can live alongside either of the others; it wouldn't likely replace them for full agent orchestration, but it can enhance them.

Ultimately, **Haystack is recommended as the backbone** of the platform. It strikes the right balance between flexibility and structure, and it aligns closely with the goals of AgentOps: robust tool integration (via MCP), traceability, and modular deployment. By building on Haystack and supplementing with LlamaIndex and LangChain components as needed, you can leverage all their strengths. This way, the platform can achieve broad tool compatibility and sophisticated data handling without compromising on maintainability or protocol compliance. Such a combined strategy ensures that the resulting AgentOps platform is not only powerful and feature-rich but also reliable and ready for the demands of production usage.

---

[1] Model caches |  LangChain
https://python.langchain.com/docs/integrations/llm_caching/

[2] LangChain vs LlamaIndex vs Haystack | by Yujian Tang - Dev Genius
https://blog.devgenius.io/langchain-vs-llamaindex-vs-haystack-0d12d25b189e

[3] GitHub - deepset-ai/hayhooks: Easily deploy Haystack pipelines as REST APIs and MCP Tools.
https://github.com/deepset-ai/hayhooks

[4] [5] Google Drive
https://drive.google.com/file/d/1zZHJ4QgVaUCD84q8r1H_hwwD5x8y_eXA

[6] [7] MCPTool
https://docs.haystack.deepset.ai/docs/mcptool