**⊛ ChatGPT**

# Open-Source Alternatives for LLM Routing in *otel-ai* LLM Manager

The custom **LLM Manager** in `clayroach/otel-ai` can largely be replaced with robust open-source libraries. Below we evaluate each major aspect of its multi-model routing and suggest NPM packages that provide similar functionality – including prompt routing, tool/chain selection, provider failover, and orchestration across models. Each recommendation points to documentation for confirmation.

## LangChain.js – Multi-Model Prompt Routing and Tool Selection

LangChain.js is a popular framework that could replace much of the custom LLM manager logic. It provides out-of-the-box support for multiple LLM providers (OpenAI, Anthropic, local models, etc.) and advanced routing/agent capabilities:

- **Prompt/Chain Routing:** LangChain allows building *router chains* that direct an input to different prompts or sub-chains dynamically. For example, a *MultiRouteChain* can use one model to classify the query, then route it to an appropriate chain or prompt specialized for that query type [1]. This mirrors the *task-based model selection* in the custom manager, but using LangChain's built-in chain logic instead of custom code.
- **Agent Tool Selection:** LangChain's agent framework uses a language model to decide which tool or chain to use for a given step [2]. In an agent (ReAct) loop, the model can reason about the task and choose actions (e.g. call a database tool, use a code formatter) iteratively [2]. This means if the LLM Manager implemented tool or chain selection (for instance, deciding whether to run a "SQL query generation" chain vs. a "general answer" chain), a LangChain agent can handle that logic internally.
- **Multiple Model Orchestration:** LangChain supports using different models within one workflow. You can configure an agent with a *dynamic model* function that selects which LLM to use at runtime based on the context [3]. This enables sophisticated routing and cost optimization – e.g. using a fast local model for simple queries and a powerful GPT-4 for complex ones. In practice, your custom *model registry* and routing rules could be translated into LangChain's dynamic model selection or a RouterChain that picks the appropriate LLM. (LangChain v0.x also had MultiPromptChain and LLMRouterChain for similar purposes.)
- **Error Handling & Fallback:** While LangChain doesn't natively provide multi-provider failover like a gateway, it does offer utilities to retry calls or specify a default chain if one branch fails. For example, you could wrap LLM calls with `.withRetry()` or use a `MultiRouteChain` with a default fallback chain for errors. This would simplify any manual try/catch failover logic in the custom manager.

**Why use LangChain.js?** It's a well-supported library that covers prompt engineering, chaining of LLM calls, and integrating tools – all with minimal code. Instead of maintaining custom TypeScript classes for each provider and routing logic, you could use `ChatOpenAI`, `ChatAnthropic`, etc., and leverage LangChain's routing chains [1]. The agent API would handle tool selection and reasoning automatically [2]. In short, LangChain can serve as a high-level orchestrator to replace many bespoke parts of the LLM Manager.

## LlamaIndex TS – Agentic Workflows and Data Integration

[LlamaIndex.TS](#) (NPM package: `llamaindex` ) is another strong candidate for replacing custom chain logic, especially if your LLM manager interacts with data or external tools. LlamaIndex (the TypeScript port of the popular Python framework) focuses on building *agentic* applications connected to your own data [4] :

- **Agent Framework:** Like LangChain, LlamaIndex TS provides an agent abstraction that can use LLMs to plan and execute tasks. It supports multi-step reasoning where the agent can choose among multiple tools or actions. The library is designed to scale from simple prompts to complex multi-agent workflows [5] . This means your custom logic for deciding which action or model to use could be handled by a LlamaIndex agent in a standardized way.
- **Tool/Index Integration:** LlamaIndex's specialty is connecting LLMs with data sources. If the *otel-ai* LLM manager generates SQL or fetches telemetry data, LlamaIndex can integrate vector databases, indices, or APIs as tools. For example, you could create a VectorStore index of observability data and use LlamaIndex's `QueryEngineTool` so that an agent can decide when to query data versus when to call an LLM [6] [7] . This out-of-the-box toolset might replace custom implementations of data retrieval in the LLM manager.
- **Multiple Providers and Models:** LlamaIndex TS supports all the major LLM providers (OpenAI, Anthropic, Azure, local models via wrappers, etc.) [8] . You can specify which model to use for a given agent or tool. While LlamaIndex doesn't automatically route between models, it makes it easy to swap out the underlying LLM for different tasks. For instance, you might configure one agent to use a local Llama2 for cost savings and another to use GPT-4 for complex queries – all within the same framework.
- **Chain of Thought and Planning:** If the custom LLM manager used intermediate reasoning (as hinted by the *ReAct*-style approach in LangChain or the "model registry pattern"), LlamaIndex agents can similarly perform reasoning steps. The benefit is you get a tested planning loop instead of writing your own. This could improve the reliability of tool selection and reduce "quirks" that you mentioned had to be handled with adaptive routing.

**Why use LlamaIndex TS?** It provides a *higher-level workflow* for building LLM applications with your data and tools, which aligns with the LLM Manager's goal of multi-model orchestration. The framework is designed for *agentic RAG (Retrieval-Augmented Generation)* scenarios, so it naturally handles combining LLM calls with data lookups or calculations. By using LlamaIndex, you could remove custom glue code and let its agents decide when to call Claude vs. GPT, or when to use a database tool – based on your prompts and provided tools. The TypeScript SDK is open-source and geared towards Next.js and Node apps, making integration straightforward [4] .

## OpenRouter – Unified Multi-Provider API with Fallbacks

If the current LLM Manager's primary job is to route requests to different providers (OpenAI, Anthropic, local) and handle failover, **OpenRouter** offers a plug-and-play solution. OpenRouter is an API gateway (by Together AI) that allows you to access **all major models through one API endpoint** [9] . You can replace low-level provider SDK calls and custom fallback logic with OpenRouter's unified interface:

- **Single API for Multiple Models:** OpenRouter supports providers like OpenAI, Anthropic, Google PaLM, etc., under a single API key. Instead of instantiating separate API clients (one for GPT-4, one for Claude, etc.), you send requests to OpenRouter and specify the model ID (e.g.

`"openai/gpt-4"` or `"anthropic/claude-2"` ). The request format is OpenAI-compatible, so you can even use OpenAI's own Node SDK pointed at OpenRouter's base URL [10] . This eliminates the need for a custom `LLMClient` interface – the API itself abstracts the providers.

- **Automatic Model Selection:** OpenRouter provides an *Auto Router* feature – you can use the special model ID `"openrouter/auto"` to let their system choose the best model for a given prompt [11] . This is effectively **prompt-based routing** out of the box. For example, if your LLM Manager tries to pick a model based on the query content (as was hinted by performance insights), `openrouter/auto` can do that classification and selection for you. It uses an AI classifier (NotDiamond) to route to a high-quality model suited to the prompt [11] .
- **Provider Failover:** OpenRouter natively supports failover ordering. You can supply a list of model endpoints in a `models` array, and it will **try the next model if the primary one fails** due to downtime, rate limit, or content filtering [12] . For instance, you might request with `"models": ["openai/gpt-4", "anthropic/claude-2"]` . If GPT-4's API is unavailable or errors out, OpenRouter will automatically call Claude as a fallback [12] . This kind of robust retry logic could replace your manual failover code and increase reliability (no single point of failure).
- **Unified Monitoring and Cost Control:** While not strictly part of routing, using OpenRouter means all model usage funnels through one service. OpenRouter will return which model was ultimately used in the response, along with usage tokens, etc. It also handles choosing *the least expensive available provider* and can enforce rate limits across providers [13] [14] . This could complement your observability goals – you'd have a central place to monitor prompt latencies and costs, rather than gathering metrics from each provider separately.

**Why use OpenRouter?** It significantly simplifies multi-LLM orchestration by handling provider selection and failovers for you. In code, this means you could remove the complex branching logic (which likely checks which model to call or catches errors to retry on another). Instead, a single API call to OpenRouter can achieve intelligent routing and high availability. OpenRouter isn't a client library per se but a service; however, it's free to use (bring your own API keys for each provider) and it aligns well with the *"intelligent model routing"* goal of your project. The fact that it **automatically falls back** on errors [12] and even has an auto-select mode for best model [11] addresses both the *provider failover* and *prompt routing* aspects of the LLM Manager.

## Portkey – Open-Source LLM Gateway for Conditional Routing & Resilience

[Portkey](#) is an open-source AI gateway that can be self-hosted. It provides enterprise-grade routing, orchestration, and observability for LLMs – effectively an in-house alternative to the custom manager, with many features ready-made:

- **Configurable Conditional Routing:** Portkey allows you to define routing rules that decide which model or provider to use based on custom conditions (without changing your application code) [15] . For example, you could route queries containing the word "SQL" to a specialized SQL model, or send large-context requests to Claude instead of GPT-4. This *conditional routing* logic is declared in a config and enforced by Portkey's gateway [16] , replacing hard-coded branching in your TypeScript.
- **Automatic Retries and Fallbacks:** Portkey has built-in strategies for **retrying failed requests and failing over** to secondary models/providers [17] . It will handle transient errors or unresponsive primary models by calling another model per your configured order. This is very similar to OpenRouter's failover, but Portkey gives you more control and is fully open source. The gateway can

also act as a circuit breaker – e.g. temporarily stop using a provider that starts erroring frequently [18].

- **Multi-Model Orchestration:** With Portkey, you can access over 1,600+ models via a unified REST API [19]. It supports all major cloud providers and open-source models. Like OpenRouter, you specify which model to use per request, but Portkey also supports *load balancing across API keys*, canary testing new models, and even multi-modal (image, audio) requests through the same interface [20]. This means your *"model registry"* from the custom implementation could be externalized to Portkey's config – you register API keys and model endpoints once, and then ask Portkey to route or distribute requests among them.
- **Caching and Optimization:** A standout feature is Portkey's **semantic cache** [20]. The gateway can cache LLM responses (exact or semantically similar prompts) to save cost and improve latency. If your observability app often gets repeated or similar questions, Portkey could return cached answers instantly, reducing the need to hit an LLM every time. This goes beyond what a basic LLM manager typically does and can drastically cut down API usage costs.
- **Observability and OTel Integration:** Given *otel-ai*'s focus on OpenTelemetry, it's worth noting Portkey also emphasizes observability. It can emit traces and metrics for all LLM calls passing through it, and it's designed to plug into tools like OpenTelemetry for monitoring and debugging [21]. By adopting Portkey, you gain a pre-built way to trace LLM requests, monitor latency per model, and even set *budget limits or rate limits* on usage [22].

**Why use Portkey?** It's essentially a drop-in, open-source **LLM Ops platform** that covers routing, fallback, caching, and monitoring. For a production-grade system, Portkey could replace the entire LLM Manager module: you would route requests to the Portkey gateway (which you can run locally or as a service) and configure in YAML/JSON how to handle them. The flexibility is immense – e.g., you can say *"use Claude for queries labeled analytics, GPT-4 for general queries; if both fail, use local LLama"* – without writing new code each time. All the features align with the needs identified in *otel-ai*: intelligent routing, model selection, cost-awareness, and reliability through retries [20]. Since Portkey is open-source, you can integrate it tightly and even contribute if needed. It's backed by a community focusing on LLM orchestration best practices.

## Conclusion

In summary, the custom multi-LLM routing in *otel-ai* can be mapped to high-level libraries that are actively maintained and feature-rich. **LangChain.js** provides a developer-friendly way to implement prompt routing, chains, and tool usage inside your code [1] [2]. **LlamaIndex TS** offers an agent-based approach especially useful if your LLMs need to work with external data or tools in a conversation [5]. **OpenRouter** and **Portkey** approach the problem at the API/gateway level – they externalize routing logic so your code remains simple while still achieving model selection and failover (OpenRouter via a hosted service with unified API [12], Portkey via a self-hosted gateway with rich policies [20]).

Each of these can replace parts of the LLM Manager:

- *LLM provider clients*: Use LangChain's integrated model classes or OpenRouter/Portkey's unified API instead of bespoke client implementations.
- *Prompt/Query routing*: Use LangChain RouterChains or OpenRouter's auto-model routing to decide which model/prompt to use [11].
- *Tool/Chain selection*: Use LangChain or LlamaIndex agents to decide actions and tools at runtime, rather than hard-coding sequences [2].

- *Provider failover*: Rely on OpenRouter or Portkey's fallback mechanisms to handle errors gracefully [12] [20] .
- *Multi-model orchestration*: Utilize Portkey or LangChain's multi-model support to orchestrate complex workflows with monitoring and caching built-in [20] .

By leveraging these open-source projects, you can accelerate development and gain proven capabilities, all while focusing on the unique features of *otel-ai* (AI-native observability) rather than reinventing LLM orchestration from scratch. Each library comes with comprehensive documentation and a community, as linked above, to support the transition.

---

[1] MultiRouteChain | LangChain.js

https://v02.api.js.langchain.com/classes/langchain.chains.MultiRouteChain.html

[2] [3] Agents - Docs by LangChain

https://docs.langchain.com/oss/javascript/langchain-agents

[4] [5] [8] LlamaIndex.TS - Build LLM-powered document agents and workflows

https://ts.llamaindex.ai/

[6] [7] How to build LLM Agents in TypeScript with LlamaIndex.TS | by Emanuel Ferreira | LlamaIndex Blog | Medium

https://medium.com/llamaindex-blog/how-to-build-llm-agents-in-typescript-with-llamaindex-ts-a88ed364a7aa

[9] OpenRouter

https://openrouter.ai/

[10] [11] [12] [14] Model Routing | Dynamic AI Model Selection and Fallback | OpenRouter | Documentation

https://openrouter.ai/docs/features/model-routing

[13] Principles | OpenRouter's Core Values and Mission

https://openrouter.ai/docs/principles

[15] [16] [17] [21] Portkey: An open-source AI gateway for easy LLM orchestration | InfoWorld

https://www.infoworld.com/article/3835182/portkey-an-open-source-ai-gateway-for-easy-llm-orchestration.html

[18] [20] [22] AI Gateway - Portkey Docs

https://portkey.ai/docs/product/ai-gateway

[19] Portkey: Production Stack for Gen AI Builders

https://portkey.ai/