

Agentic Engineering

Engineering practices for agentic AI systems

- [Configuring Custom Models: Vision & Parallel Tool Calls](#)

Configuring Custom Models: Vision & Parallel Tool Calls

Overview

This page documents configuration patterns for using custom models (e.g. Qwen, custom fine-tunes) with **LiteLLM** as a proxy and **OpenCode** as the agentic frontend.

When using non-standard models, both LiteLLM and OpenCode require explicit configuration to advertise and support advanced capabilities like vision (image input) and parallel tool calls. Without proper setup, these features silently fail.

Vision / Multimodal Support for Custom Models

Configuring image support for custom models in a llama.cpp + LiteLLM + opencode setup.

The Problem

LiteLLM does **not** infer vision support for arbitrary `custom_openai` models. OpenCode also performs its own preflight modality check and does **not** read LiteLLM's `supports_vision` metadata for custom providers. Both layers must be configured independently.

LiteLLM Side

- Add `supports_vision: true` under the LiteLLM model's `model_info` so LiteLLM and its model metadata advertise image capability
- LiteLLM's OpenAI-compatible/custom OpenAI chat path does not normally strip OpenAI-style image blocks — it preserves `image_url` content and normalizes string image URLs into `{ "url": ... }`

OpenCode Side

- For opencode custom providers, image support must be declared per model with `modalities.input` including `image`
- If `modalities` is set, include both `input` and `output`
- Without `modalities.input: ["text", "image"]`, attachments get replaced with an error: `ERROR: Cannot read image (this model does not support image input). Inform the user.`

Working OpenCode Model Shape

```
"Qwen3.6-35B-A3B": {
  "name": "Qwen3.6-35B-A3B",
  "limit": {
    "context": 262144,
    "output": 65536
  },
  "modalities": {
    "input": ["text", "image"],
    "output": ["text"]
  }
}
```

Parallel Tool Calls

Enabling parallel tool calls for a custom model served through LiteLLM, used by OpenCode.

The Problem

The `parallel_tool_calls: true` in `model_info` is **metadata only** — it does not auto-forward to the API. Without explicit pass-through, the downstream server never receives the signal.

OpenCode Side

- **No config change needed** in `opencode.jsonc`
- OpenCode uses the AI SDK which auto-detects parallel tool call capability from the provider
- The `tool_call` boolean in model config is for declaring tool support, not parallel behavior
- Your setup (`jamesravey_litellm` provider with `@ai-sdk/openai-compatible`) handles this at the SDK level

LiteLLM Side

Fix: Add `parallel_tool_calls` inside `litellm_params`:

```
"litellm_params": {  
  ...  
  "parallel_tool_calls": true  
}
```

Prerequisites

1. LiteLLM must be **v1.61.0+** (parallel tool calls pass-through added then)
 2. The downstream inference server (vLLM, TGI, etc.) must support parallel tool calls
-

Files Reference

- `~/.config/opencode/opencode.jsonc` — OpenCode model config (no changes needed for parallel tool calls)
- LiteLLM model config — add `parallel_tool_calls` to `litellm_params` and `supports_vision: true` to `model_info`