# My AI Engineering Tech Stack

(How Open-Source Models *Actually* Run AI Coding at Scale)

Alex Ker
Engineer, Baseten

baseten

# Hi, I'm Alex





**And I work on AI infrastructure with these awesome + kind people at Baseten**

> Previously built ML pipelines at Neurable (BCI) and RL at LaunchDarkly
> Founded an AI incubator p-ai.org and contributed to Stanford HAI as an editor
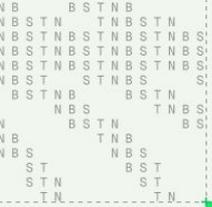> I like calisthenics and reading in my free time

Some of our customers <3, building at the frontier

We get to partner with them to implement and create the
patterns around agentic systems today

## Agenda

1. Limitations of closed-source models

2. Rise of open-source coding models

3. How to use OSS in your dev workflows

4. Performant inference for coding models

5. Takeways for AI native builders

# The problem with closed models

While closed source models bring a lot of advantages and out-of-the-box support and are an easy solution to get started, problems appear at scale.

- **Rate limits** - TPS
- **Cost** - multiples more expensive
- **Reliability** - Unmet SLAs
- **Speed** - Throughput
- **Control** - Fine-tuning, transparency

Why use closed models?
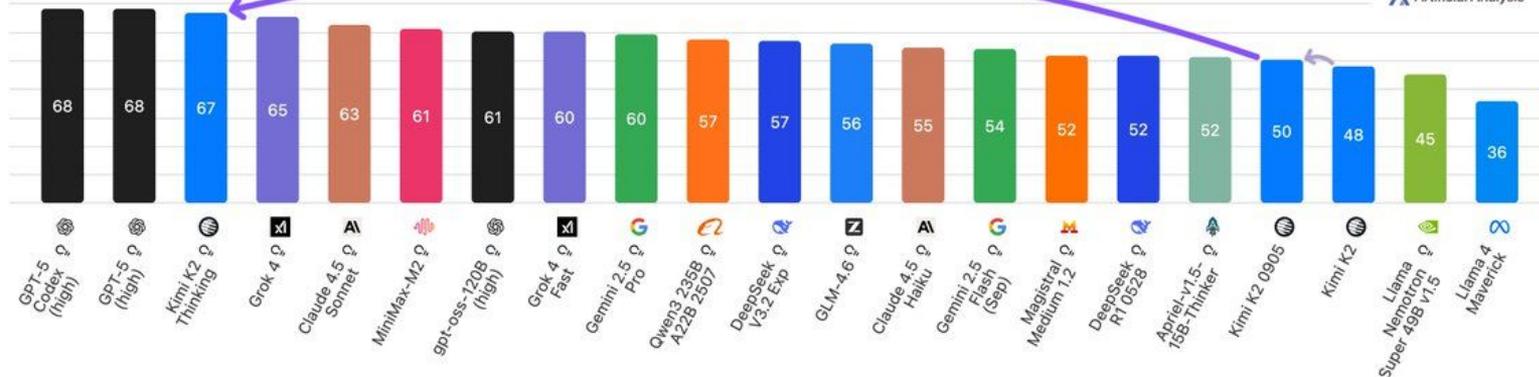
Because they're smart!

# But the quality gap is closing

As we recently saw the Kimi K2 Thinking drop, which is considered the smartest open-source model. It's benchmark scores are competitive with models like GPT-5 and Claude Sonnet 4.5.
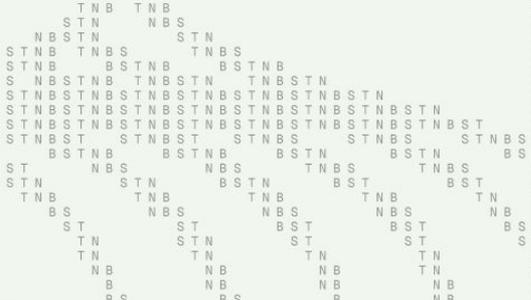


## Artificial Analysis Intelligence Index

Artificial Analysis Intelligence Index v3.0 incorporates 10 evaluations: MMLU-Pro, GPQA Diamond, Humanity's Last Exam, LiveCodeBench, SciCode, AIME 2025, IFBench, AA-LCR, Terminal-Bench Hard, τ²-Bench Telecom
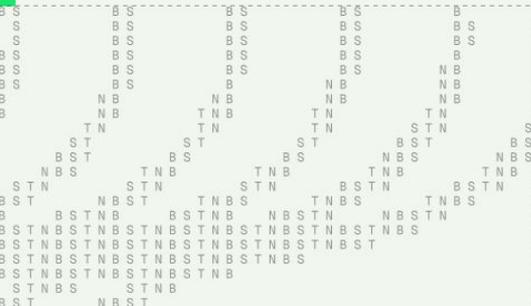
+ Add model from specific provider

Artificial Analysis

| Model | Score |
|---|---|
| GPT-5 Codex (high) | 68 |
| GPT-5 (high) | 68 |
| Kimi K2 Thinking | 67 |
| Grok 4 | 65 |
| Claude 4.5 Sonnet | 63 |
| MiniMax-M2 | 61 |
| gpt-oss-120B (high) | 61 |
| Grok 4 Fast | 60 |
| Gemini 2.5 Pro | 60 |
| Qwen3 235B A22B 2507 | 57 |
| DeepSeek V3.2 Exp | 57 |
| GLM-4.6 | 56 |
| Claude 4.5 Haiku | 55 |
| Gemini 2.5 Flash (Sep) | 54 |
| Magistral Medium 1.2 | 52 |
| DeepSeek R1 0528 | 52 |
| Apriel-v1.5-15B-Thinker | 52 |
| Kimi K2 0905 | 50 |
| Kimi K2 | 48 |
| Llama Nemotron Super 49B v1.5 | 45 |
| Llama 4 Maverick | 36 |

# Advantages of Open Source

**Latency** (developer experience)

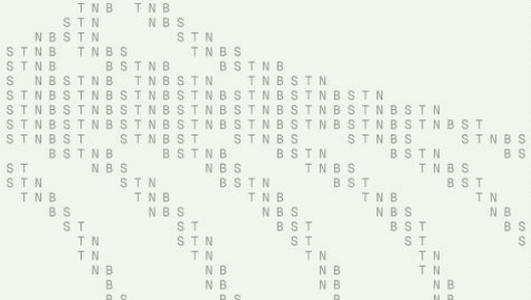**Reliability** (uptime at peaks)

**Cost** (per token)
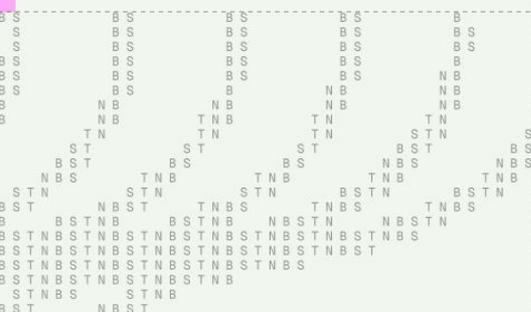
# What does AI native product engineers care about?

- It turns out…the same things.
  - People care about these principles in their devtools AND product.
- Central to building production-ready apps with standout UI/UX within a competitive market.

In 2025, the top 3 things we learned at Baseten:

1. **Reliability** is the new moat

2. **Speed** is the product

3. **Developers** want the wheel

The best open source

models today

# OpenRouter

- Unified access to 500+ models

- Automatic fallback for reliability

- Transparent metrics for providers

## The Unified Interface For LLMs

Better prices, better uptime, no subscription.

Start a message...  →

Featured Models                    View Trending ⧉

**Gemini 2.5 Pro** ✦
by google
196.0B          2.7s          -13.54%
Tokens/wk       Latency       Weekly growth

**GPT-5**
by openai
72.0B           7.2s          +9.98%
Tokens/wk       Latency       Weekly growth

**Claude Sonnet 4.5**
by anthropic
659.8B          2.0s          -11.97%
Tokens/wk       Latency       Weekly growth

**20T**
Monthly Tokens

**4.2M+**
Global Users

**60+**
Active Providers

**500+**
Models

baseten

# GLM 4.6 on OpenRouter

## Baseten

ⓘ  US  fp4  📄  ⏸  🔑

| | | |
|---|---|---|
| Latency | Throughput | Uptime |
| 0.19s | 141.0tps | ▮▮▮ |

| Total Context | Max Output | Input Price | Output Price | Cache Read | Cache Write | Input Audio | Input Audio Cache |
|---|---|---|---|---|---|---|---|
| 262K | 163.8K | $0.60 | $2.50 | -- | -- | -- | -- |

# OpenRouter Stats 1)



### Top Models

Weekly usage of models across OpenRouter

**February 16, 2026**

| Model | Value |
|---|---|
| Others | 5.64T |
| MiniMax M2.5 | 2.57T |
| Kimi K2.5 | 1.04T |
| Gemini 3 Flash Preview | 859B |
| GLM 5 | 803B |
| DeepSeek V3.2 | 745B |
| Grok 4.1 Fast | 669B |
| Claude Opus 4.6 | 643B |
| Claude Sonnet 4.5 | 534B |
| Trinity Large Preview (free) | 449B |
| Total | 14T |

# OpenRouter Stats 2)



🏆 LLM Leaderboard                    This Week ⌄

1.  🔴  **MiniMax M2.5**           3,24T tokens
        by minimax              ↑5 405%

2.  ⚫  **Kimi K2.5**              1,24T tokens
        by moonshotai           ↓18%

3.  Ⓩ  **GLM 5**                  1,03T tokens
        by z-ai                 ↑1 448%

4.  ✦  **Gemini 3 Flash Previ...** 816B tokens
        by google               ↑10%

5.  🔵  **DeepSeek V3.2**          738B tokens
        by deepseek             ↓0%

6.  Ⓐ  **Claude Opus 4.6**        619B tokens
        by anthropic            ↑40%

# Let's break down the top 3

**MiniMax M2.5** - small, fast, smart; good default choice

**GLM 5** - very smart, but slow from high active parameter count

**Kimi Thinking (and 2.5)** - reasoning thus slow, best multimodal model

# MiniMax M2.5
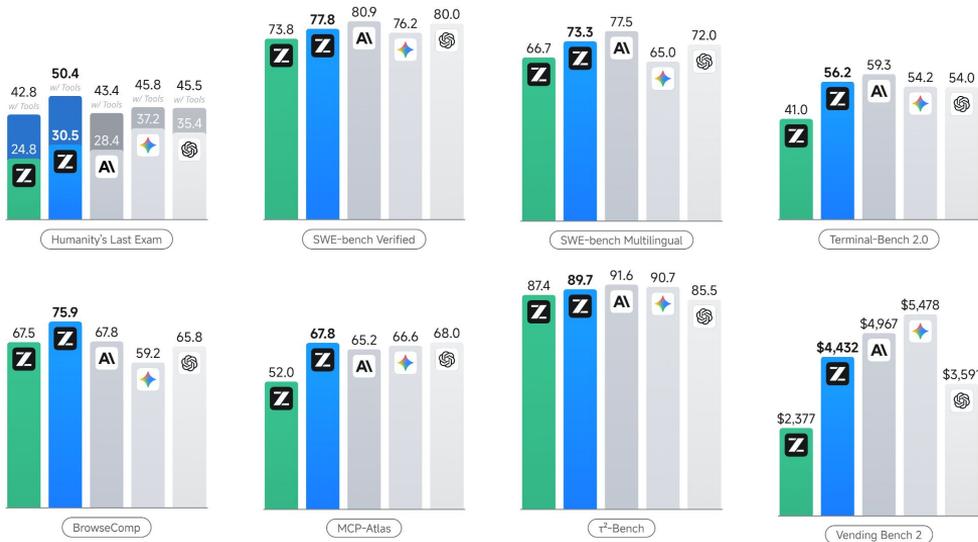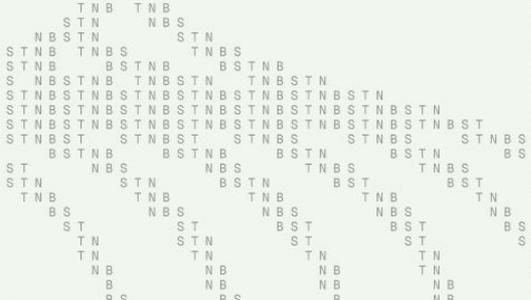
small, fast, smart; good default choice



**SWE-bench Verified Score Evolution**
**Anthropic vs OpenAI vs Google vs MiniMax**

Legend:
- Anthropic (Sonnet / Opus)
- OpenAI (GPT)
- MiniMax (M Series)
- Google (Gemini)

Y-axis: SWE-bench Verified Score (%)

Data points:
- Sonnet 3.7 — 62.3%
- o3 — 69.1%
- Sonnet 4 / Opus 4 — 72.7%
- Opus 4.1 — 74.5%
- Sonnet 4.5 — 77.2%
- Opus 4.5 — 80.9%
- Opus 4.6 — 80.8%
- GPT-5 — 72.8%
- GPT-5.1 — 76.3%
- GPT-5.2 — 80.0%
- Gemini 2.5 Pro — 63.8%
- Gemini 3.0 Pro — 76.2%
- M1 — 56.0%
- M2 — 69.4%
- M2.1 — 74.0%
- M2.5 — 80.2%

X-axis: Feb '25, Apr '25, Jun '25, Aug '25, Oct '25, Dec '25, Feb '26

*Note: Scores are from official company announcements. Testing scaffolds may vary.*

# GLM 5

very smart, but slow from high active parameter count



LLM Performance Evaluation: Agentic, Reasoning and Coding

8 benchmarks: Humanity's Last Exam, SWE-bench Verified, SWE-bench Multilingual, Terminal-Bench 2.0, BrowseComp, MCP-Atlas, τ²-Bench, Vending Bench 2

Legend: GLM-4.7, GLM-5, Claude Opus 4.5, Gemini 3 Pro, GPT-5.2 (xhigh)

Humanity's Last Exam: 42.8 w/ Tools (24.8), 50.4 w/ Tools (30.5), 43.4 w/ Tools (28.4), 45.8 w/ Tools (37.2), 45.5 w/ Tools (35.4)

SWE-bench Verified: 73.8, 77.8, 80.9, 76.2, 80.0

SWE-bench Multilingual: 66.7, 73.3, 77.5, 65.0, 72.0

Terminal-Bench 2.0: 41.0, 56.2, 59.3, 54.2, 54.0

BrowseComp: 67.5, 75.9, 67.8, 59.2, 65.8

MCP-Atlas: 52.0, 67.8, 65.2, 66.6, 68.0

τ²-Bench: 87.4, 89.7, 91.6, 90.7, 85.5

Vending Bench 2: $2,377, $4,432, $4,967, $5,478, $3,591

# Kimi K2 Thinking Breakdown

## Interesting Stats

- 1T parameter model

- Leader on humanity's' last exam

- Tau-squared bench for agentic tool use

- 200-300 continuous tool calls

## Architecture

- Sparse MoE architecture

- Bigger vocab compared with DeepSeek R1

- More experts, fewer attention heads per expert

# Kimi K2 Thinking: Interleaved Reasoning

Model alternates between thinking and tool use

Compared to chain of thought, where all the reasoning is upfront, interleaved:

- Mimics humans – Reflect after each action, adjust approach dynamically
- Maintains coherence
- Solves complexity – Breaks ambiguous problems into adaptive subtasks through continuous feedback

**Example:**
Solved PhD-level geometry problem using 23 interleaved reasoning + tool call cycles

## Kimi K2 Thinking: Tool Call Training

# 5-Step Training Pipeline:

- Gathered 3K+ real tools from GitHub

- Generated 20K+ synthetic tools via clustering

- Created synthetic agents with varied configurations

- Simulated diverse multi-turn usage scenarios

- Evaluated and filtered all trajectories for quality

# Kimi K2.5: coding with vision



Original Video (Cr: Jesko Jet)

Kimi's Output

# My favorite dev

## workflows and harnesses

# Claude Code CLI

## (cheaper, faster, and just as smart)



**Ian Nuttall** ✔
@iannuttall

Claude Code pro tip:

Use Kimi K2 as your agent by changing the base URL and setting the API key before running claude.

all prompts and code requests will now go through Kimi, in the CC terminal you're used to!

```
~/dev/ianslist   main  1 • +8 -450
export ANTHROPIC_BASE_URL=https://api.moonshot.ai/anthropic
export ANTHROPIC_AUTH_TOKEN=sk-YOURKEY
claude
```

Shell    auto (claude 4 sonnet)

6:28 AM · Jul 14, 2025 · **204.3K** Views

💬 91       ⟳ 144       ♡ 1.5K       🔖 2K

```
19    ### Speed (Throughput) (Higher is better)
20    | Model | Tokens per Second |
21    |--------|---------------------|
22    | Claude Sonnet 4.5 | ~60 TPS |
23    | GLM-4.6 on Baseten | ~100+ TPS |
24
25    ## Quick Start
26
```

Problems   Output   Debug Console   **Terminal**   Ports

```
alexker@Alexs-MacBook-Pro baseten-claude-code % ./setup_baseten_claude_code
.sh
🚀 Starting LiteLLM proxy...
✅ Starting LiteLLM proxy on http://localhost:4000
🍎 Use Ctrl+C to stop
```

```
alexker@Alexs-MacBook-Pro baseten-claude-code % ./launch_claude_code.sh
```

bash

zsh

cursor agent to run Agent · ⌘K to generate command

cursor agent to run Agent · ⌘K to generate command

main    0    0

Instead of hacking these CLIs, what if there was a way to build in open source native IDE?

# Vercel AI SDK V5

- Instantiate provider objects that could be used directly in your Next.js application

- Compatibility across streaming, async, batch etc

- Patterns to help you build AI apps in your language of choice

## Model APIs

You can select Baseten models ↗ using a provider instance. The first argument is the model id, e.g. `'moonshotai/Kimi-K2-Instruct-0905'` : The complete supported models under Model APIs can be found here ↗.

```
1  const model = baseten('moonshotai/Kimi-K2-Instruct-0905');
```

## Example

You can use Baseten language models to generate text with the `generateText` function:

```
1  import { baseten } from '@ai-sdk/baseten';
2  import { generateText } from 'ai';
3
4  const { text } = await generateText({
5    model: baseten('moonshotai/Kimi-K2-Instruct-0905'),
6    prompt: 'What is the meaning of life? Answer in one sentence.',
7  });
```

Baseten language models can also be used in the `streamText` function (see AI SDK Core).

baseten

# LangChain

- Agent Builder

- For knowledge workers, non technical

- Create an agent based on guided questions

# Cline

- Bring your own key

- Plan/Act modes

- Prewritten system prompts

# OpenClaw

- Your own Jarvis, get started by choosing a model, and list of apps to be connected

- people have been able have tested crazy workflows
  - Booking restaurant reservations through voice call
  - Building apps while texting it feedback on whatsapp

- Kimi K2.5 is 8x cheaper and 4.5x faster than Claude Opus*



Baseten vs. Claude Opus pricing (lower is better)

# Cursor

- my personal go-to setup
  - **Composer** (for faster, non-complex queries)
  - **Claude Code** (for planning projects, harder tasks that I need to be done correctly, not worried about time)
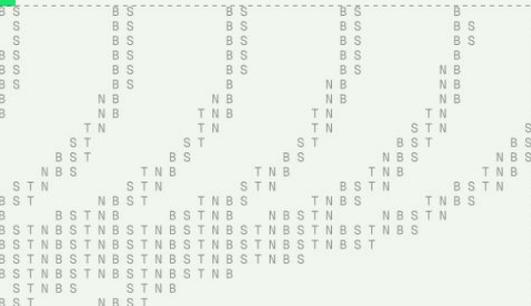
# Performant inference

# for coding models

# Case study: **Autocomplete**

- Care about TTFT

- Long prefill, short decode
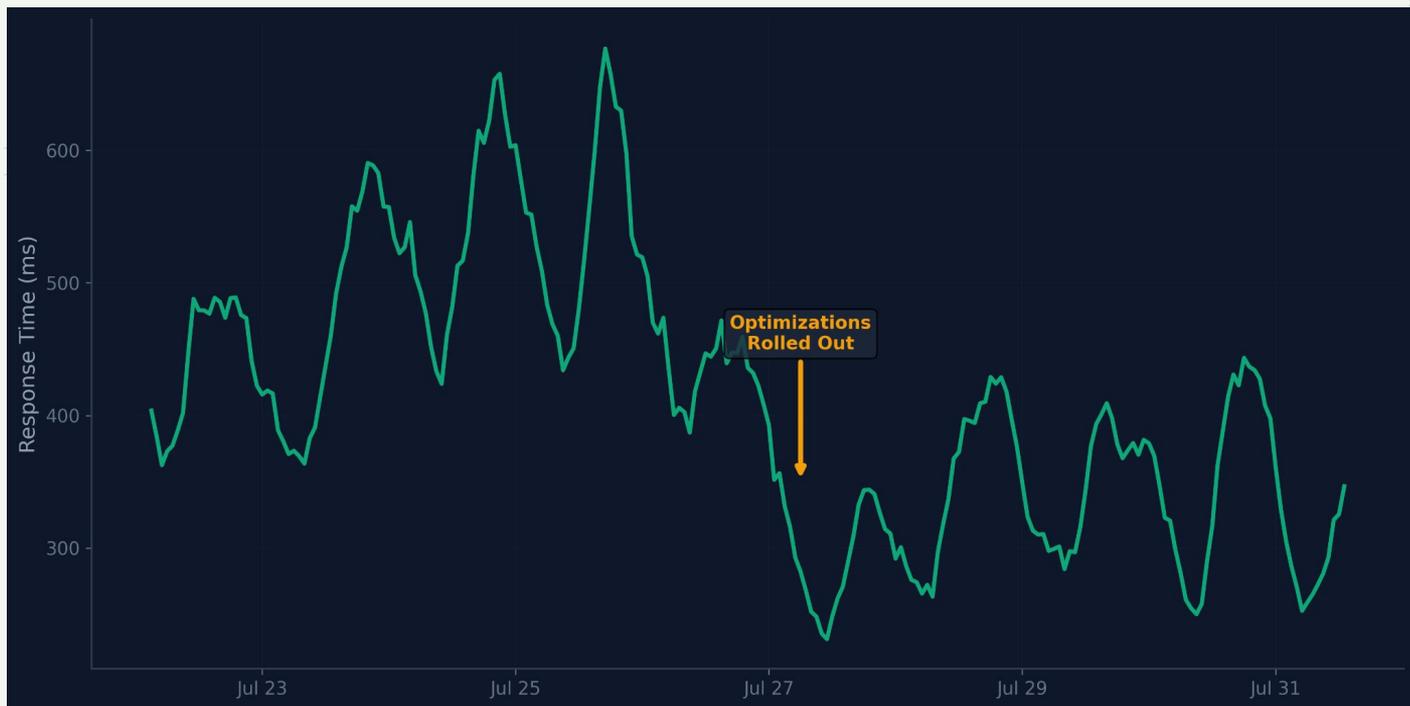
- Keeps up with user workflow

# How?

| KV cache re-use | KV-aware routing | N-gram speculation |
|-----------------|------------------|--------------------|

# Results!

# General Key Takeaways

1. If you're just using Claude Code or Codex, you're missing out

2. Open models have closed the gap

3. Developers have built a rich ecosystem of tools around them

4. The best thing is to **experiment** and not limit yourself

5. If you care about build delightful UX with AI, you should care about: **performance, reliability & control**

baseten

# Actionable Steps

- **Delegate** and **review** like a staff engineer
- Develop your **own taste** of models & harnesses combinations. Spend time and let them do things for you!
- The models are very capable given the *right information*
    - The golden question: "Is there sufficient context for the agent to figure out what the answer is?"
    - This may solve most of your problems

baseten

# Thank you!

instagram.com/thealexker

x.com/thealexker

linkedin.com/in/alex-ker

baseten