# AI & Engineering

A bit of history & usage.

sourcetoad

# Plan of Attack

- The beginning of the LLM

- A bit of History (GPT1-5)

- The Evolution of Tools

- AI in the Workplace

- Building with AI

- AI Supporting Software

sourcetoad

# Who?

- **Connor Tumbleson**

- @iBotPeaches

- Sourcetoad
  - Director of Engineering

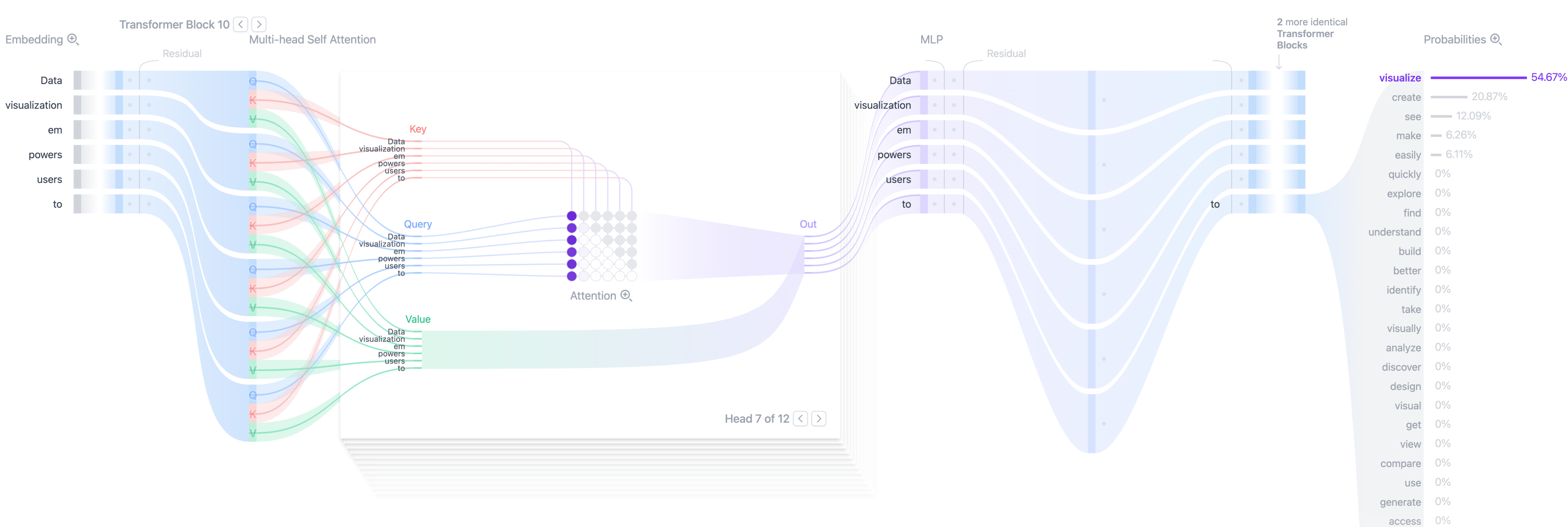sourcetoad

# Attention Is All You Need

- 2017 Research Paper
- Google Scientists
- 8 equally cited authors
- Cited over 200,000 times.

# Transformer Explainer

sourcetoad

ChatGPT 1-5

# GPT-**1**: The beginning (2018)

- Heavy research usage
- Trained from books
- Tokens in/out
- ~500 token limit

sourcetoad

# Hold up. **Tokens**?

**Tokens**
6

**Characters**
23

Hello my name is connor

Text    Token IDs

**Tokens**
6

**Characters**
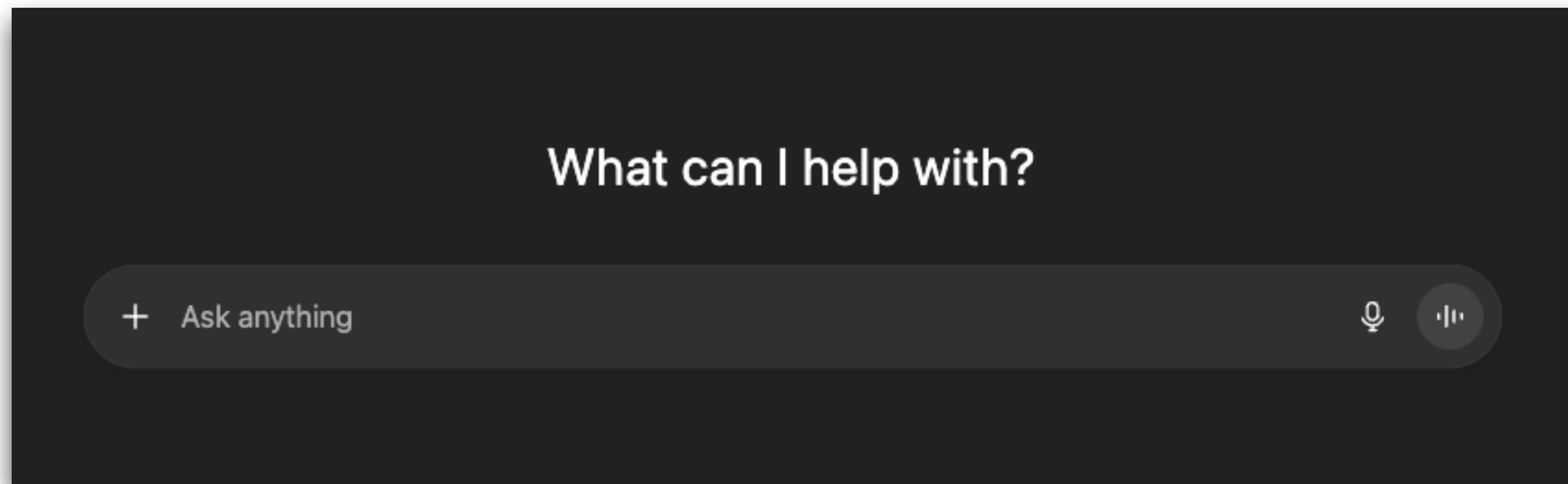23

[13225, 922, 1308, 382, 10850, 267]

Text    **Token IDs**

sourcetoad

# GPT-**2**: The Expanse (2019)
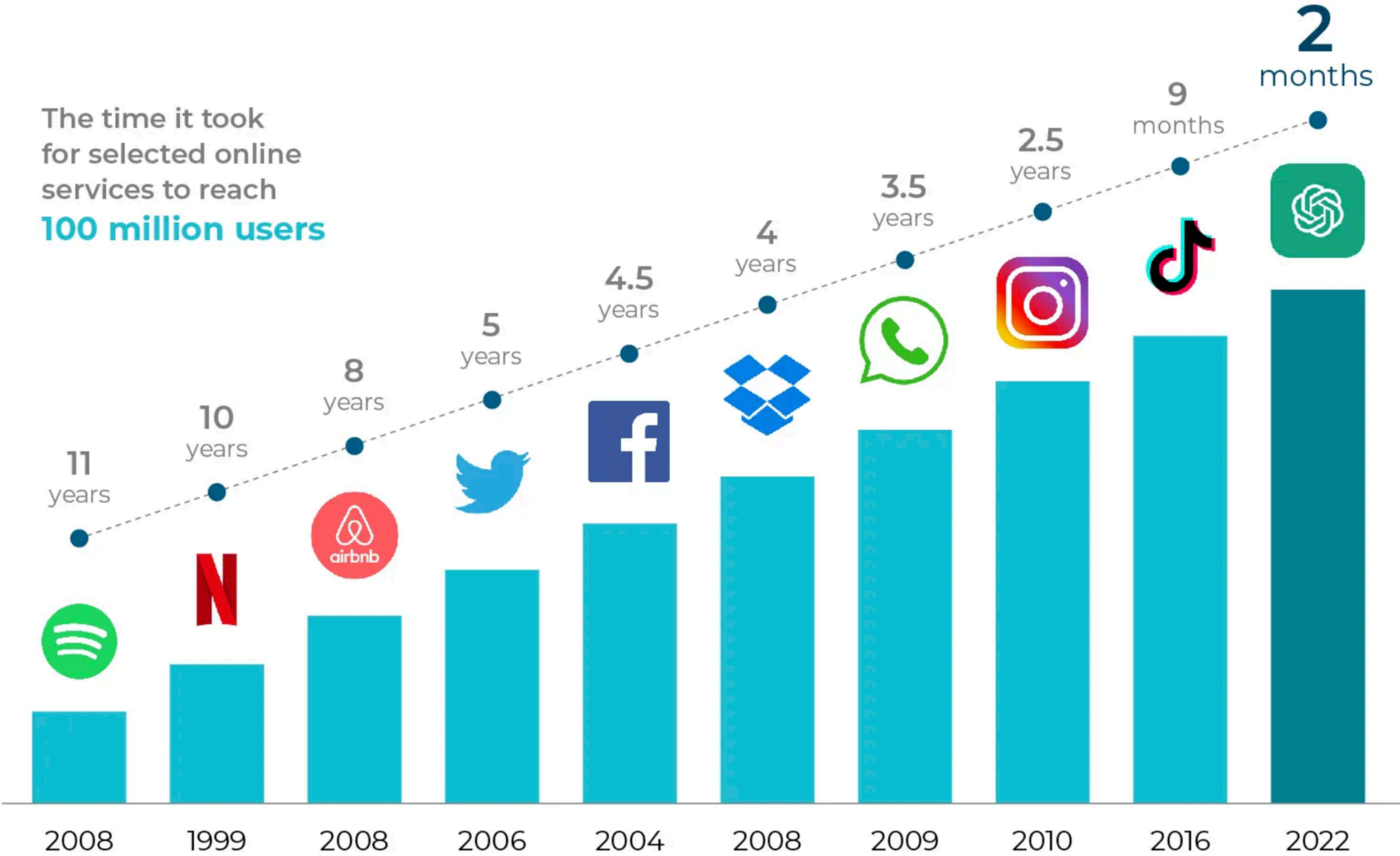
- Zero-shot training

- Trained from web/books

- ~1k token limit



sourcetoad

# GPT-**3**: The Explosion (2020)

- The commercialization begins

- Trained from everything?

- ~2k token limit

What can I help with?

+ Ask anything

sourcetoad

# GPT-**3**: The Explosion (2020)



**Chat-GPT sprints to 100 million users**

The time it took for selected online services to reach **100 million users**

- 11 years — 2008
- 10 years — 1999
- 8 years — 2008
- 5 years — 2006
- 4.5 years — 2004
- 4 years — 2008
- 3.5 years — 2009
- 2.5 years — 2010
- 9 months — 2016
- 2 months — 2022

**Source:** World of Statistics

# GPT-**3**: The Explosion (2020)

- ## /completions
  - `logit_bias`
  - `presence_penalty`
  - `seed`
  - `temperature`
  - `top_p`

## Create completion · Legacy

POST `https://api.openai.com/v1/completions`

Creates a completion for the provided prompt and parameters.

### Request body

**model** string Required
ID of the model to use. You can use the List models API to see all of your available models, or see our Model overview for descriptions of them.

**prompt** string or array Required
The prompt(s) to generate completions for, encoded as a string, array of strings, array of tokens, or array of token arrays.

Note that <|endoftext|> is the document separator that the model sees during training, so if a prompt is not specified the model will generate as if from the beginning of a new document.

sourcetoad

# GPT-**4**: The Refinement (2023)

- Multimodal Input

- Trained from more and more.

- ~32k token limit

*GPT-4 can accept a prompt of text and images, which—parallel to the text-only setting—lets the user specify any vision or language task.*

# GPT-**4**: The Refinement (2023)

- New API `/chat/completions`

## Create chat completion

```
POST https://api.openai.com/v1/chat/completions
```

**Starting a new project?** We recommend trying Responses to take advantage of the latest OpenAI platform features. Compare Chat Completions with Responses.

sourcetoad

# GPT-**4**: The Refinement (2023)

- New API `/assistants`

## Create assistant <span>Beta</span>

```
POST https://api.openai.com/v1/assistants
```

Create an assistant with a model and instructions.

# GPT-**4**: The Refinement (2024)

- New API `/responses`

## Create a model response

`POST https://api.openai.com/v1/responses`

Creates a model response. Provide text or image inputs to generate text or JSON outputs. Have the model call your own custom code or use built-in tools like web search or file search to use your own data as input for the model's response.

sourcetoad

# GPT-5: The Reasoning (2025)

- Reasoning / Completion

- Trained from [redacted]

- ~512k token limit

ChatGPT 5 ⌄

GPT-5

**Auto**
Decides how long to think ✓

**Instant**
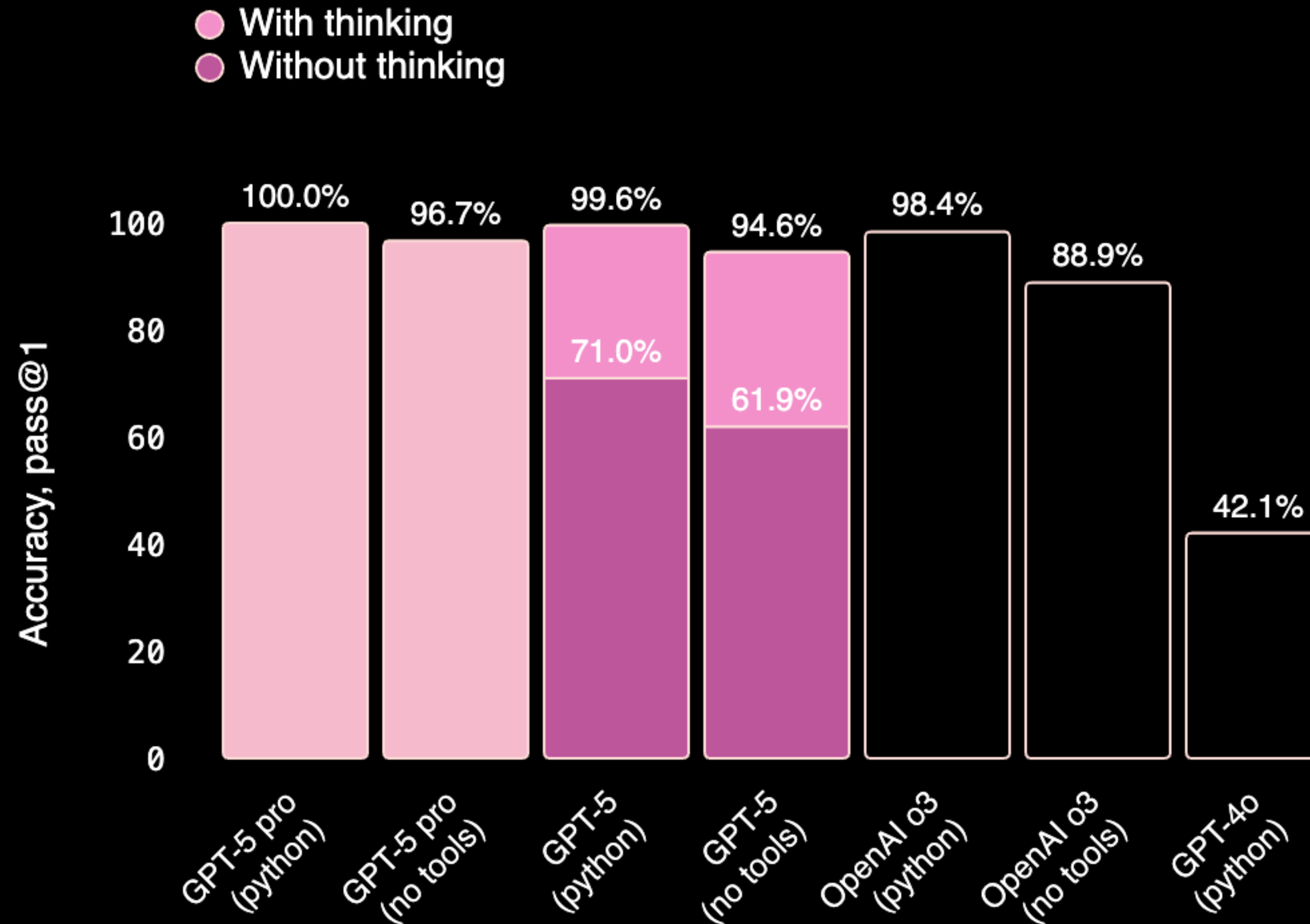Answers right away

**Thinking**
Thinks longer for better answers

**Pro**
Research-grade intelligence

Legacy models ›

# GPT-**5**: The Reasoning (2025)



AIME 2025
Competition math

● With thinking
● Without thinking

Accuracy, pass@1

- GPT-5 pro (python): 100.0%
- GPT-5 pro (no tools): 96.7%
- GPT-5 (python): 99.6% / 71.0%
- GPT-5 (no tools): 94.6% / 61.9%
- OpenAI o3 (python): 98.4%
- OpenAI o3 (no tools): 88.9%
- GPT-4o (python): 42.1%

# GPT-1 through GPT-5 Parameters

- Brains of the model

- Getting smarter

- Nearing AGI?

|  | Parameters | Year |
|---|---|---|
| GPT-1 | ~117 million | 2018 |
| GPT-2 | ~1.5 billion | 2019 |
| GPT-3 | ~175 billion | 2020 |
| GPT-4 | ~1.7 trillion | 2023 |
| GPT-5 | ?? | 2025 |

sourcetoad

LLM Tools

# LLM Tools

- Code Interpreter

- File Search

- Function Tool

- Web Search

- Computer Use

- Image Generation

- Local Shell

- Custom Tool Call

- Connectors

- MCP

sourcetoad

# Tool: Code Interpreter

- ## Write Code

- ## Run Code

```
curl https://api.openai.com/v1/responses \
  -H "Content-Type: application/json" \
  -H "Authorization: Bearer $OPENAI_API_KEY" \
  -d '{
    "model": "gpt-4.1",
    "tools": [{
      "type": "code_interpreter",
      "container": { "type": "auto" }
    }],
    "instructions": "You are a personal math tutor. When asked a math que
    "input": "I need to solve the equation 3x + 11 = 14. Can you help me?
  }'
```

sourcetoad

# Tool: Code Interpreter

- ## How Many R's

How many 'r' in 'strawberry'

**ChatGPT**

There are two 'r' characters in the word 'strawberry'.

**You**

Check again please

**ChatGPT**

Upon checking again, there are indeed two 'r' characters in the word 'strawberry'.

**You**

Are you 100% sure?

**ChatGPT**

Yes, I am 100% sure. The word 'strawberry' contains two 'r' characters.

**You**

Would you bet a million dollars on this?

**ChatGPT**

Yes, I would. The word "strawberry"

# Tool: File Search

- Create Vector Store

- Attach Files

  - images, text, docs.

- Search it

```json
{
  "output": [
    {
      "type": "file_search_call",
      "id": "fs_67c09ccea8c48191ade9367e3ba71515",
      "status": "completed",
      "queries": ["What is deep research?"],
      "search_results": null
    },
    {
      "id": "msg_67c09cd3091c819185af2be5d13d87de",
      "type": "message",
      "role": "assistant",
      "content": [
        {
          "type": "output_text",
          "text": "Deep research is a sophisticated capability
          "annotations": [
            {
              "type": "file_citation",
              "index": 992,
              "file_id": "file-2dtbBZdjtDKS8eqWxqbgDi",
              "filename": "deep_research_blog.pdf"
            },
```

sourcetoad

# Tool: Function Tool

- Create functions

- Expose to LLM

- LLM picks function

**Developer** ......................................... **Model**

**1** Tool Definitions + Messages
get_weather(location)
What's the weather in Paris?

**2** Tool Calls
get_weather("paris")

**3** Execute Function Code
get_weather("paris")
↓
{"temperature": 14}

**4** Results
All Prior Messages
{"temperature": 14}

**5** Final Response
It's currently 14°C in Paris.

# Tool: Web Search

- Search the web

- Cite the source

```
[
    {
        "type": "web_search_call",
        "id": "ws_67c9fa0502748190b7dd390736892e100be649c1a5ff9609",
        "status": "completed"
    },
    {

        "id": "msg_67c9fa077e288190af08fdffda2e34f20be649c1a5ff9609",
        "type": "message",
        "status": "completed",
        "role": "assistant",
        "content": [
            {
                "type": "output_text",
                "text": "On March 6, 2025, several news...",
                "annotations": [
                    {
                        "type": "url_citation",
                        "start_index": 2606,
                        "end_index": 2758,
                        "url": "https://...",
                        "title": "Title..."
                    }
                ]
            }
        ]
    }
]
```

sourcetoad

# Tool: Computer Use



**Screenshot**

**Computer environment**

**CUA Model**

**Computer tool call**

```
{
  type: "click",
  button: "left",
  x: 286,
  y: 102
}
```

**Application code**

**Action**

```
click (left, 286, 102)
```

sourcetoad

# Tool: Image Generation

- Describe

- Wait

- Image

- Iterate

"Generate an image of gray tabby cat hugging an otter with an orange scarf"



"Now make it look realistic"



sourcetoad

# Tool: Local Shell

```
>_ Welcome to Codex, OpenAI's command-line coding agent

> Sign in with ChatGPT to use Codex as part of your paid plan
  or connect an API key for usage-based billing

> 1. Sign in with ChatGPT
     Usage included with Plus, Pro, and Team plans
  2. Provide your own API key
     Pay for what you use

  Press Enter to continue
```

# Tool: Custom Tool Call

- Grammar (Lark)

- CFGs

- Conform

- Profit

```javascript
const response = await client.responses.create({
  model: "gpt-5",
  input: "Use the code_exec tool to print hello world to the console.",
  tools: [
    {
      type: "custom",
      name: "code_exec",
      description: "Executes arbitrary Python code.",
    },
  ],
});

console.log(response.output);
```

sourcetoad

# Tool: Connectors

- ## OpenAI Specific

- ## Company Connections

# Tool: MCP

**Chat interface**
Claude Desktop, LibreChat

**IDEs and code editors**
Claude Code, Goose

**Other AI applications**
5ire, Superinterface

**MCP**
Standardized protocol

**Data and file systems**
PostgreSQL, SQLite, GDrive

**Development tools**
Git, Sentry, etc.

**Productivity tools**
Slack, Google Maps, etc.

Bidirectional data flow

Bidirectional data flow

AI applications

Data sources and tools

sourcetoad

# MCP Basics

- Open Source

- Connect LLM with External Thing

- The "USB-C" of AI applications

# AI in the Workplace

# AI for **You - Summarize**

Re: Client Testing Account `External` → `Inbox ×` `Projects/Client ×`

### AI Overview

- Connor requested a new testing account after initial credentials failed, needing access for issue review.
- Client set up Connor's account and asked for confirmation of access.
- Connor's testing revealed a user failure and an inability to install a required SSL certificate for debugging.

By Gemini; there may be mistakes. **Learn more**

# AI for **You - Meeting Notes**

## Notes from "Client Internal Sync"

These notes have been sent to invited guests in your organization.

**Open meeting notes**

The content was auto-generated on December 3, 2025, 4:20 PM EST and may contain errors.

### Summary

redacted

### Suggested next steps

→ Redacted Next Step 1

→ Redacted Next Step 2

Notes by **Gemini**

# AI for **You - Server Monitoring**



**Amazon GuardDuty**
A threat detection service that continuously monitors for compromised accounts, anomalous behavior, and malware

**Activate GuardDuty**
With a few steps in the console, monitor all your AWS accounts without additional software to deploy or manage

Amazon S3

Instance workloads

Databases

Accounts and users

Container workloads

Serverless

**Continuously analyze**
Automatically and continuously monitor AWS workloads and resources for potential threats at scale

**Intelligently detects threats**
GuardDuty uses machine learning anomaly detection, malware scanning, and integrated threat intelligence to identify and prioritize potential threats

**Take action**
Review detailed findings in the console, integrate into event management or workflow systems, and initiate AWS Lambda for automated remediation or prevention

sourcetoad

# AI for **You - GitHub Copilot PR**



Agents / Add playlist rotation change indicators and historical dates

`<> Code ▾`   `⇄ View pull request`

## Initial implementation ✓ 21m 51s ···

1. **Percentage Change Badges** - Added visual indicators showing +/- percentage changes for maps and gametypes when a previous rotation exists
2. **Added/Removed Indicators** - Clear badges showing when maps/gametypes are newly added (🟢 NEW) or removed (🔴 REMOVED)
3. **Date Display** - Shows current and previous rotation dates in the sidebar bottom left as requested

### 🔧 Technical Implementation:

1. **CompareRotations Action Class** - Calculates differences between current and previous playlist rotations
2. **Enhanced PlaylistPage Component** - Fetches previous rotation data and calculates changes
3. **Updated mode-breakdown Partial** - Displays change badges and handles different change types
4. **Model Relationship** - Added PlaylistChange relationship to Playlist model
5. **Comprehensive Testing** - Unit tests for comparison logic and integration tests for UI

### 🎨 Visual Improvements:

- **Green "+NEW" badges** for newly added items
- **Red "REMOVED" badges** for removed items (shows previous percentage)
- **Green "+X%" badges** for increased percentages
- **Yellow "-X%" badges** for decreased percentages
- **Date information** showing current and previous rotation dates

### 📊 User Experience:

The feature provides immediate visual feedback about how playlist compositions have changed, making it easy for users to:

What would you like to change?

Submitting will post a pull request comment

---

Overview | 📄 **Files changed** 8

○ 0 / 8 viewed  ⚙

> …ctions/Playlist/CompareRotations.php 📋 +63 ▪▪▪▪▪ ☐ Viewed 💬 ···

> app/Livewire/PlaylistPage.php 📋 ✛ +31 ▪▪▪▪▪ ☐ Viewed 💬 ···

∨ app/Models/Playlist.php 📋 ✛ +6 ▪▪▪▪▪ ☐ Viewed 💬 ···

```
          @@ -33,6 +33,7 @@
33    33     * @property-read string $image
34    34     * @property-read ?PlaylistStat $stat
35    35     * @property-read Collection<int, PlaylistAnalytic> $analytics
      36   +  * @property-read Collection<int, PlaylistChange> $changes
36    37     *
37    38     * @method static PlaylistFactory factory(...$parameters)
38    39     */

          @@ -137,4 +138,9 @@
137   138         {
138   139             return $this->hasMany(PlaylistAnalytic::class);
139   140         }
      141   +
      142   +    public function changes(): HasMany
      143   +    {
      144   +        return $this->hasMany(PlaylistChange::class);
      145   +    }
140   146     }
```

> …livewire/playlist-page.blade.php 📋 ✛ +22 -2 ▪▪▪▪ ☐ Viewed 💬 ···

∨ …laylist/mode-breakdown.blade.php 📋 ✛ +33 -1 ▪▪▪▪ ☐ Viewed 💬 ···

```
          @@ -3,14 +3,46 @@
```

sourcetoad

# AI for **You - GitHub Copilot Review**



**fastlane/Fastfile**  (Outdated)                    ⊹ Hide resolved

```
705  +        sh("bundle install")
706  +      end
707  +
708  +      sh("rake install")
```

**Copilot**  (AI)  2 weeks ago                              ···

The sh() call for "rake install" should use the array form for consistency with the PR's goal of hardening sh() calls. Line 95 properly uses the array form, and lines 702/705 at least have the excuse of being in exception handling, but this line has no such constraint.

| Suggested change |
| --- |
| −        sh("rake install") |
| +        sh("rake", "install") |

**Commit suggestion** ▾

☺  👍  👎   [Copilot](#) uses AI. Check for mistakes.

Reply...

**Unresolve conversation**    **iBotPeaches** marked this conversation as resolved.

sourcetoad

# AI for **You - GitHub Copilot Chat**

Ranked playlists and KDA formula details ⌄

How is the KDA formula produced?

👥 20 references ⌄

The KDA (Kills, Deaths, Assists) formula in the LeafApp_Infinite project is displayed in UI tooltips and code as:

**KDA = (Kills + (Assists × 0.3)) ÷ Deaths**

This is shown in:

- UI abbreviations:
  - "KDA" column/heading:
    `(Kills + (Assists * .3)) / Deaths` (example: <u>overview-stats.blade.php</u> and <u>stats.blade.php</u> )
- The value itself is stored pre-computed in the database and set via the backend.

So in Halo Infinite as implemented here:

> **KDA = (Kills + 0.3 × Assists) / Deaths**

The app does not appear to compute KDA on the fly, but this is the accepted calculation presented to the user and stored as the `kda` value.

👍 👎 ⧉ ↻ | ⌄

# Junie

# OpenCode

# ChatGPT

What can I help with?

Ask anything

Company knowledge

# AI for **You - Chatting**

AI missing the mark

# Atlassian Intelligence - Rovo

# Atlassian Intelligence - Rovo Part 2

**Activity**

| All | **Comments** | History | Work log | Approvals |

Summarize 6 comments

**Comments summary**    🔒 Only visible to you    ⌄    ✕

- Text does not comply with Atlassian's "Acceptable use policy".

ⓘ Uses AI. Verify results.    Rovo

**Add internal note**    /    **Reply to customer**

**Pro tip:** press [M] to comment

sourcetoad

# Gmail - Suggested Response



**Suggested reply**

Hi ████████████

Okay, I'll check the Jira ticket for context and see what needs updating.

Thanks,
Connor

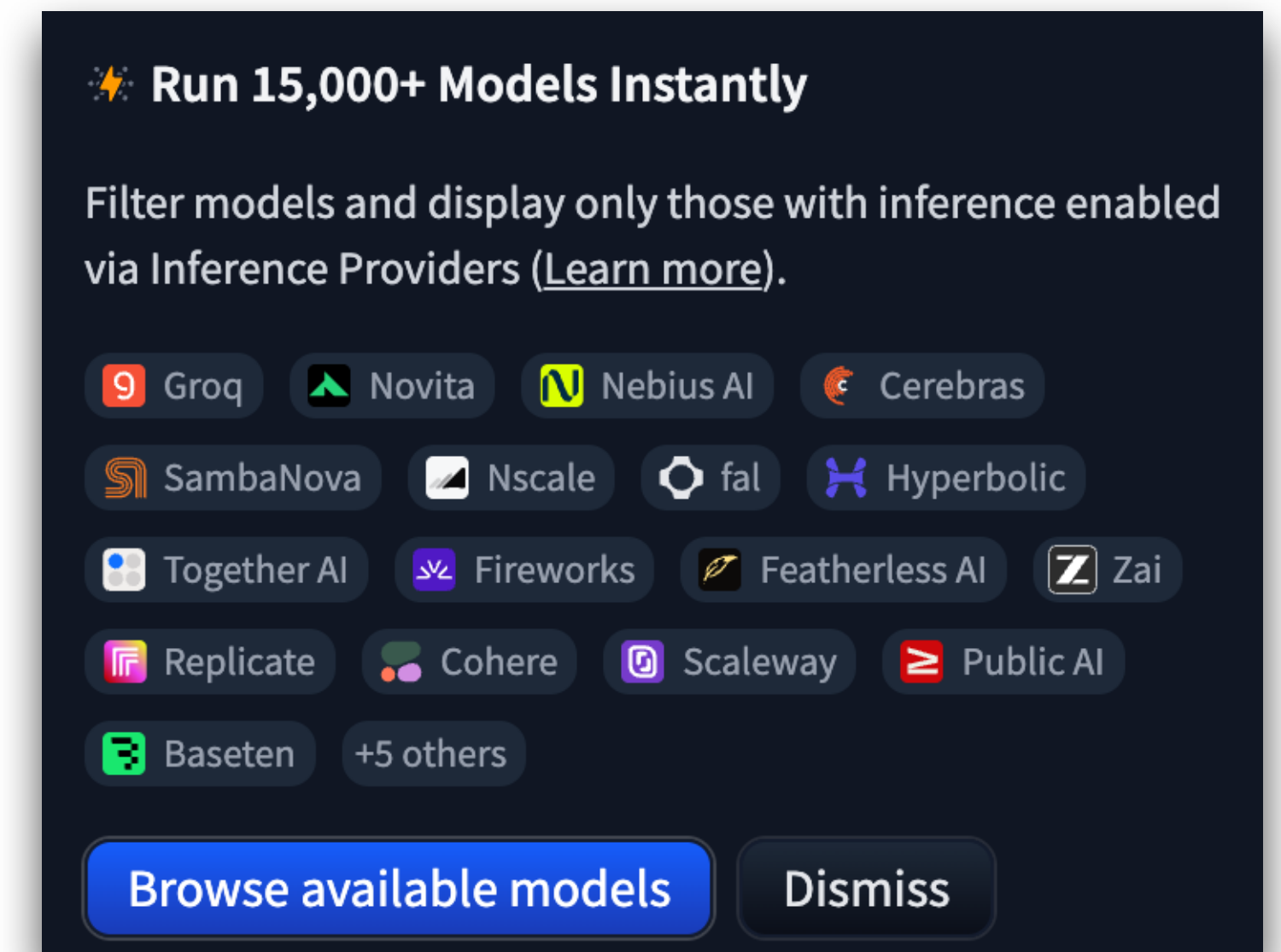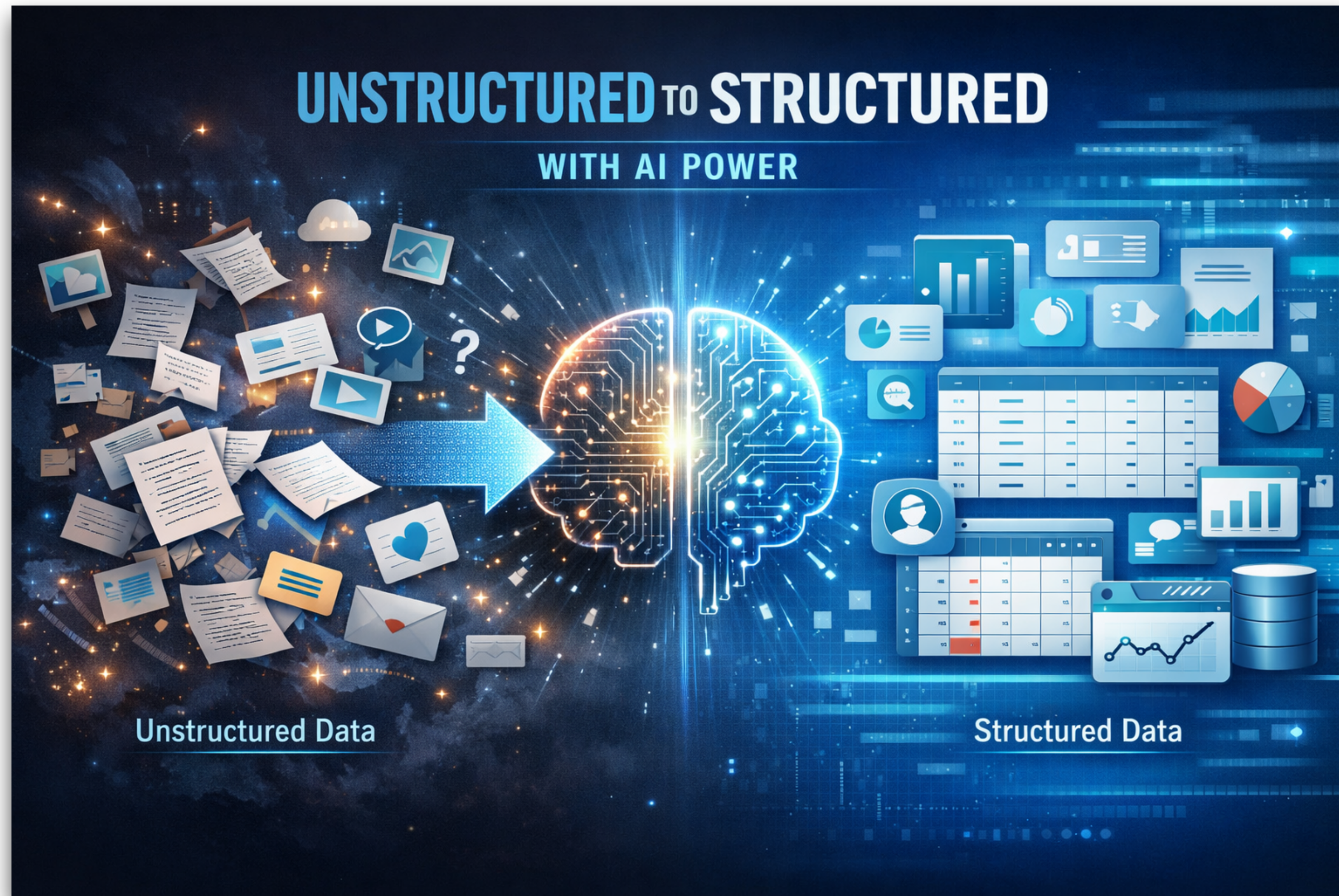[Reply] [Reply all] [Forward] [Share in chat] **New**

Building with AI

# Choosing a Model

- **S**LM vs **L**LM (Small vs Large)

- Reasoning vs Completion

- Budget

- Modality

- Availability

- Limits (Time, Context)



⚡ **Run 15,000+ Models Instantly**

Filter models and display only those with inference enabled via Inference Providers (Learn more).

Groq    Novita    Nebius AI    Cerebras
SambaNova    Nscale    fal    Hyperbolic
Together AI    Fireworks    Featherless AI    Zai
Replicate    Cohere    Scaleway    Public AI
Baseten    +5 others

**Browse available models**    Dismiss

sourcetoad

# Adding Structure to the Unstructured

# Translating the World

Can you translate "A cat is in the house" into Spanish, French, Korean, Japanese?

Sure! Here are the translations:

- **Spanish:** *Un gato está en la casa.*
- **French:** *Un chat est dans la maison.*
- **Korean:** *고양이가 집 안에 있어요.*
- **Japanese:** *猫が家の中にいます。*

If you want more natural or alternative phrasing for any of them (casual vs. formal), just let me know 🙂

Do you like this personality?

AI Supporting Software

# LangFuse

Python SDK **JS/TS SDK**

```typescript
import { LangfuseClient } from "@langfuse/client";
import OpenAI from "openai";

const langfuse = new LangfuseClient();
const openai = new OpenAI();

async function handleRequest(userInput: string) {
  // Fetches the latest 'production' version if no 'label' is provided
  // Caches client-side, revalidates in background, instant subsequent calls
  const prompt = await langfuse.prompt.get("support-reply", { type: "chat" });

  // Prompt can contain variables and placeholders
  // Filled at runtime via compile()
  const compiledPrompt = prompt.compile({ tone: "friendly", input: userInput });

  // Works with any model, supports text and chat message formats
  const res = await openai.chat.completions.create({
    model: "gpt-5",
    messages: compiledPrompt,
  });

  return res.choices[0].message.content;
}
```

</> Quickstart

📄 Prompt **langfuse-docs-assistant-text** ⓘ 🏷️

**Versions** Metrics

🔍 Search versic  + New

**# 10**  ● production

9/10/2025, 11:42:34 PM by Jannik Maierhoefer

**# 9**

9/4/2025, 5:04:45 PM by Jannik Maierhoefer

**# 8**

9/4/2025, 4:47:34 PM by Jannik Maierhoefer

**# 7**

**Prompt** Config Linked Generations

Text Prompt

You are a helpful Langfuse support a inks and minimal, runnable examples.

## Guidelines
- Make extensive use of the tools pr n.
- Getting an overview of the docs is
- Answer the question using only thi
- Before using any tools, reason ab
- After tool responses, reason about n; if not, use the tools again or gi
- If a user asks where the traces/lc ease refer to the Langfuse Demo Proj kpwwm0m000gmm094odg11gi/traces, US D 00qmad0772bw0tol/traces)

## Prompt Management

Version-control prompts collaboratively, deploy/roll-back instantly to different environments, support for templates, variables, and A/B testing. Cached client-side for 0 latency/availability impact.

📖 Documentation ↗

▷ Watch Demo ↗

# OpenRouter

# The Unified Interface For LLMs

Better prices, better uptime, no subscriptions.

**Get API Key**    **Explore Models**

**25T**
Monthly Tokens

**5M+**
Global Users

**60+**
Active Providers

**500+**
Models

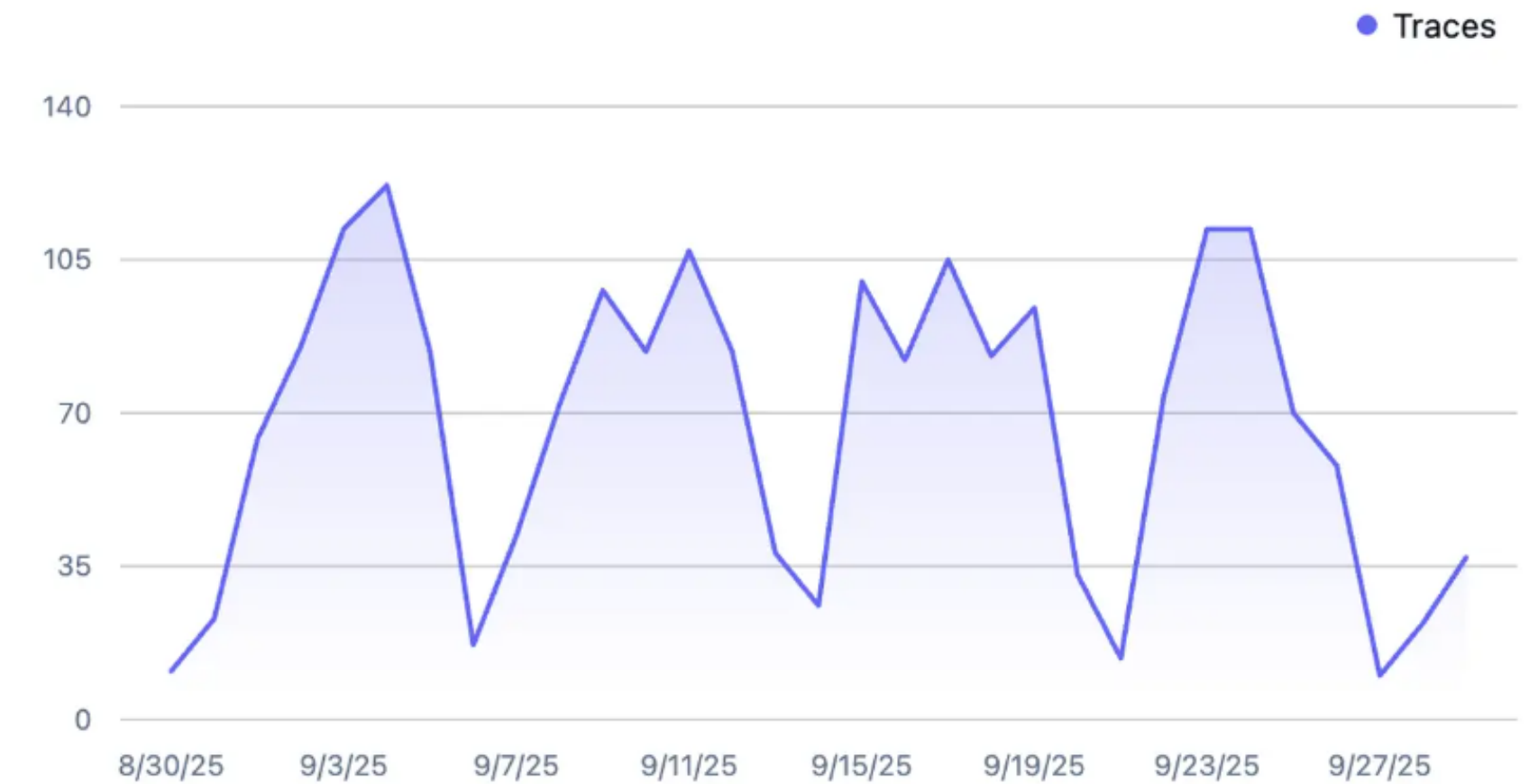# In Closing

# In Closing

- Research a model
- Tools are powerful
- Structure always
- Introspection is key

**Traces by time**

Traces     Observations by Level

**2.08K** Traces tracked

- Traces

140

105

70

35

0

8/30/25   9/3/25   9/7/25   9/11/25   9/15/25   9/19/25   9/23/25   9/27/25

**User consumption**

Token cost     Count of Traces

**$58.93** Total cost

sourcetoad

# Thanks!

connortumbleson.com

@iBotPeaches