



# Newsroom Infrastructure for AI Experimentation

Jonathan Soma  
Columbia University  
@dangerscarf · js4571@columbia.edu

<https://bit.ly/dh26-infra>

# What we' ll be looking at

- Python tools
  - Gradio
  - Streamlit
  - Ngrok/Cloudflare Tunnel/other tunnels (but not!!)
- Platforms
  - Prompt engineering and tests: Braintrust
  - *Doing it all themselves: n8n???*
- Staying flexible and model/vendor agnostic
  - Pydantic
  - LiteLLM
  - OpenRouter

# What we' ll be looking at

- ~~Python tools~~ Inviting feedback early
  - Gradio
  - Streamlit
  - Ngrok/Cloudflare Tunnel/other tunnels (but not!!)
- ~~Platforms~~ Embracing domain experts
  - Prompt engineering and tests: Braintrust
  - *Doing it all themselves: n8n???*
- Staying flexible and model/vendor agnostic
  - Pydantic
  - LiteLLM
  - OpenRouter

# Inviting feedback early

- You need something **shareable**
- Don't make other people run your code!
- **The best approaches have maybe changed since early 2026**

# Gradio

```
import gradio as gr

def greet(name):
    return f"Hello, {name}!"

with gr.Blocks() as demo:
    name = gr.Textbox(label="Name")
    output = gr.Textbox(label="Output")

    btn = gr.Button("Submit")
    btn.click(greet, inputs=name, outputs=output)

demo.launch()
```

Box Threshold 0.25 Text Threshold 0.25 Labels  
a person.  
a mountain.

Run

Image Annotated Image

a person a person a person a person a person a pe  
a mountain a person a mountain a mountain a moun  
a person a person

Labels  
a person. a mountain.

Box Threshold 0.25 Text T 0.

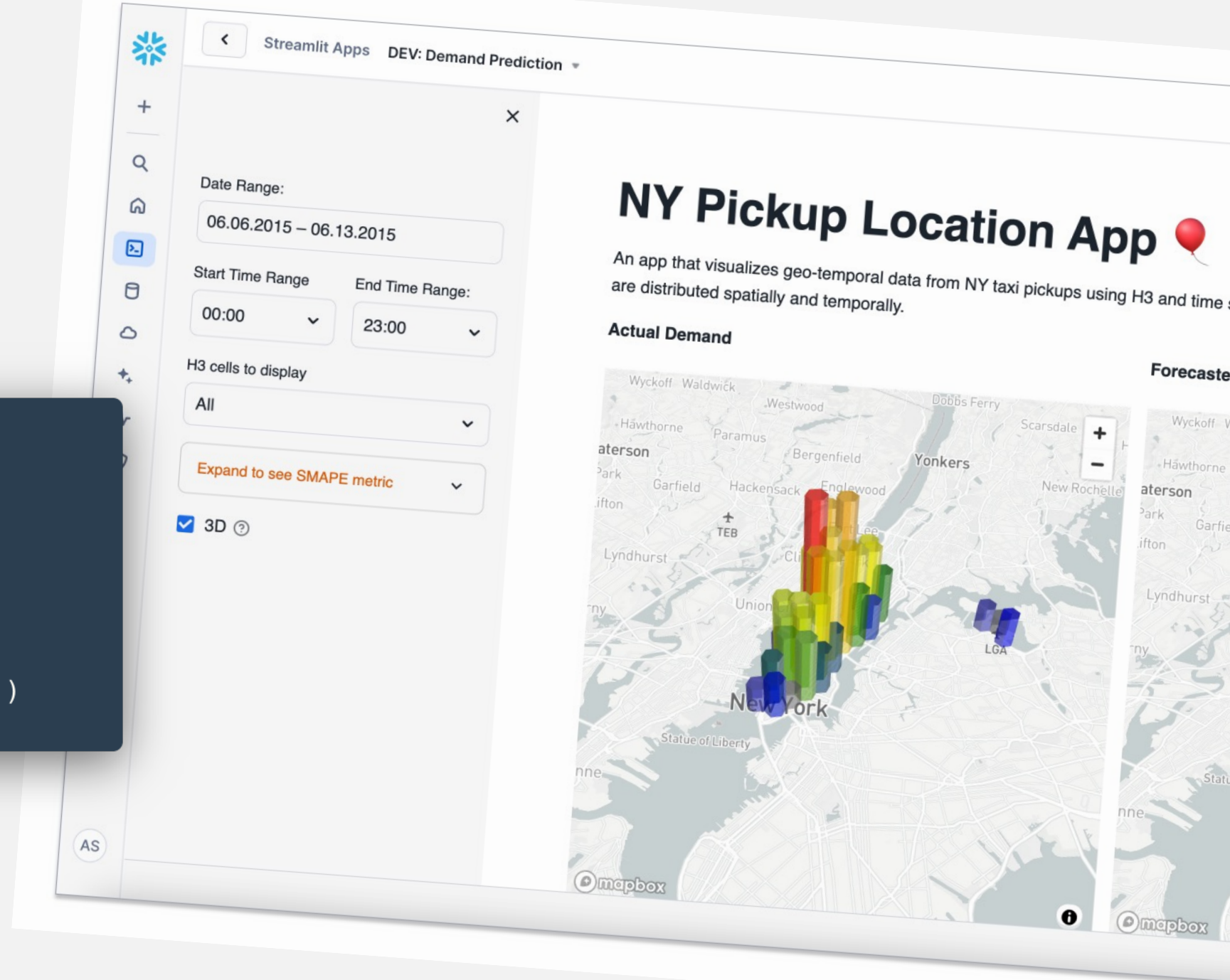
# Streamlit



```
import streamlit as st

name = st.text_input("Name")

if st.button("Submit"):
    st.write(f"Hello, {name}!")
```



Newsroom Infrastructure for

jsoma.github.io/workshop-newsroom-ai-infra/

Dark Light

DATAHARVEST 2026

# Newsroom Infrastructure for AI Experimentation

Jonathan Soma

Small local demos for OCR, text-to-speech, semantic PDF search, and Streamlit data browsing.

Open in Codespaces ↗ View si

- Open in Colab
- Colab (code-along)
- Codespaces

**OCR with Gradio: simple**  
A tiny PDF OCR interface using RapidOCR.  
↓ Download .ipynb

Try it ▾

Read

# 01-gradio-ocr-simple-ANSWERS.ipynb

File Edit View Insert Runtime Tools Help

Commands + Code + Text | **▶ Run all** Copy to Drive

[ ]



**# Install required packages**

```
!pip install --upgrade --quiet gradio ipyv
```

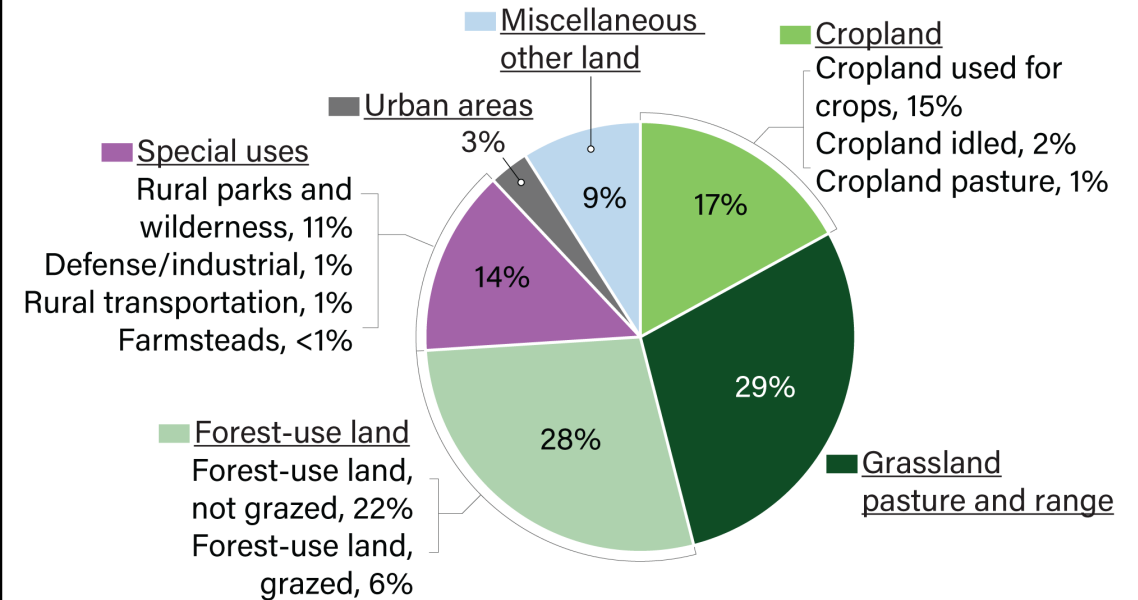
```
print('✓ Packages installed!')
```

# Custom everything

- Codex + Claude Code can now one-shot a *lot* of stuff
- Infrastructure and custom tooling is cheap
- “Home-cooked software”



Share of major U.S. land use, 2017



Note: Sub-components may not sum to the component totals because of rounding. Forest-use includes land that serves commercial forest uses, including grazing, as opposed to land that has forest cover but is used for other purposes. Miscellaneous includes uses such as land in cemeteries, golf courses, mining areas, quarry sites, marshes, swamps, sand dunes, bare rocks, deserts, tundra, and other unclassified land, as well as some—but not all—industrial, commercial, and residential sites in rural areas.

Source: USDA Economic Research Service estimates based on data from USDA, National Agricultural Statistics Service; USDA, Farm Service Agency; USDA, Natural Resources Conservation Service; USDA, Forest Service; U.S. Department of Commerce, Bureau of the Census; U.S. Department of Defense; U.S. Department of the Interior, Bureau of Land Management, National Park Service; U.S. Geological Survey; and Utah State University.

SORT BY

Urban growth (highest first) ▾

CHART TYPE

Stacked Area

Line

Stream

Slope

Stacked Bar

Horizon

TREEMAP VIEW

Year for treemap:

Off ▾

YEAR RANGE

1945 - 2017

Start:



End:



COLORS

Grassland

Forest

Cropland

Urban

Other

FEATURED COUNT

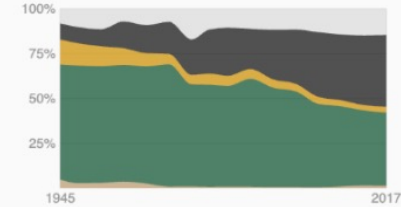
8 states ▾

# Concrete Ate the Eastern Seaboard

How urban sprawl reshaped land use across 48 states, 1945-2017

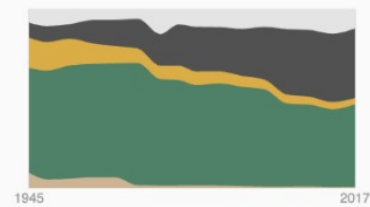
Grassland Forest Cropland Urban Other

### Massachusetts



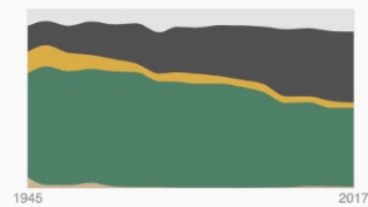
Route 128 turned woodlands into tech parks. Forest 64%→41%; parkland surged 0.4%→12%.

### Connecticut



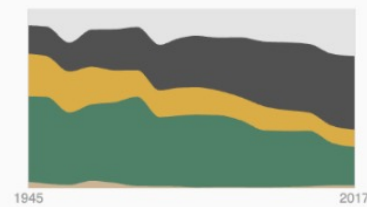
NYC commuters drove sprawl. Urban quintupled; cropland 17%→4%.

### Rhode Island



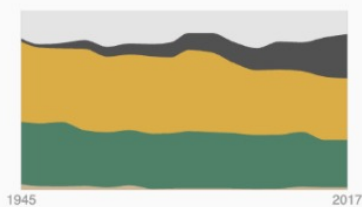
Already 14% urban in 1945 — hit 40% by 2017. Smallest state ran out of room.

### New Jersey



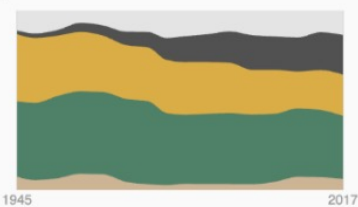
Lost more farmland to development than any state. Parks 0.5%→13% as preservation fought back.

### Delaware



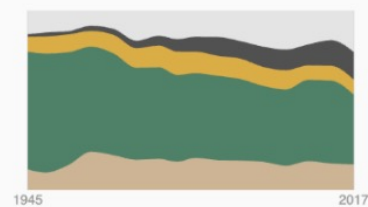
Urban 2%→25%. Wilmington suburbs spread south through farmland.

### Maryland



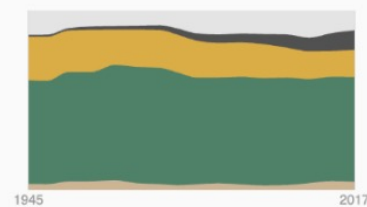
DC suburbs ate Montgomery County. Half its farmland gone by the 1960s.

### Florida



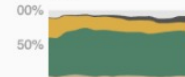
A/C made it possible. Forest halved 66%→39%; parkland exploded 0.4%→17%.

### North Carolina



Research Triangle + Charlotte: 105 acres/day consumed. 2nd in projected farmland loss by 2040.

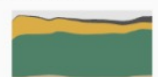
### Georgia



### Ohio



### South Carolina



### Pennsylvania



### Tennessee



### New Hampshire



### Virginia



### Indiana



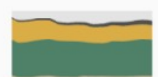
### Illinois



### New York



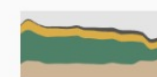
### Michigan



### Alabama



### California



### Louisiana



### Texas



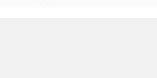
### Washington



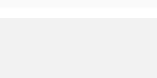
### Kentucky



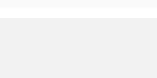
### West Virginia



### Arkansas



### Oklahoma

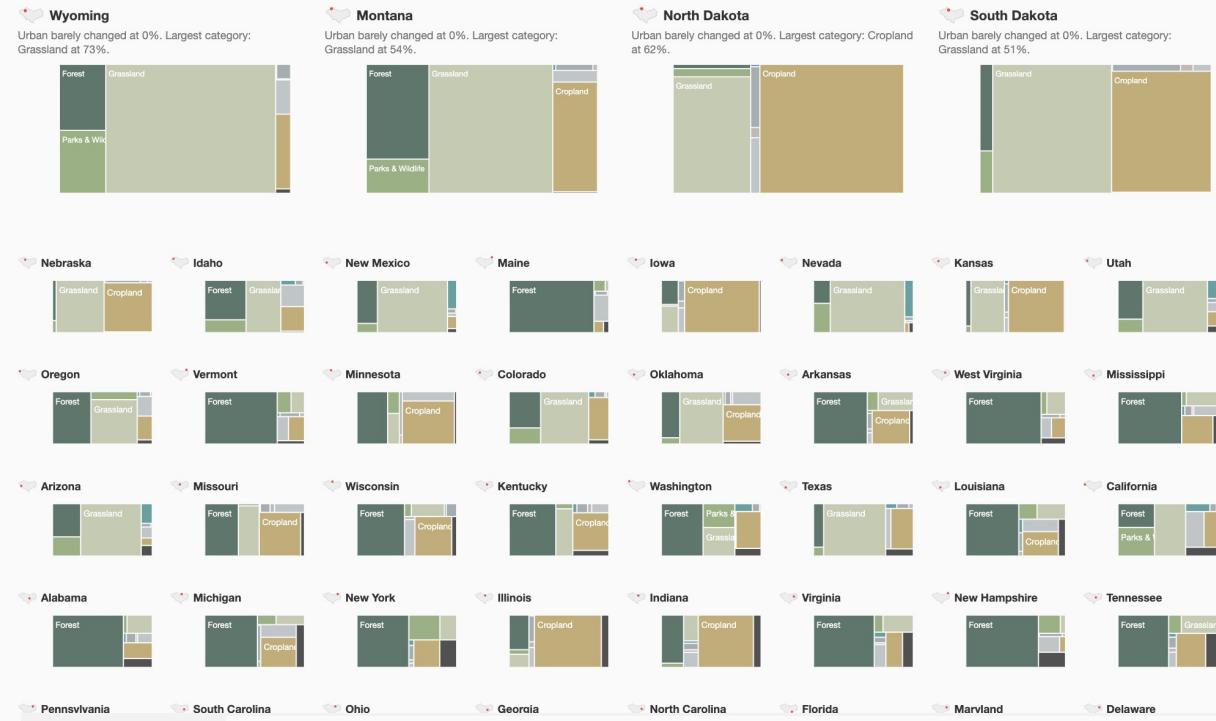


<https://wongpeiting.github.io/land-use/>

# The States Still Untouched by Sprawl in 2017

Land-use composition by state, ranked by least urban growth

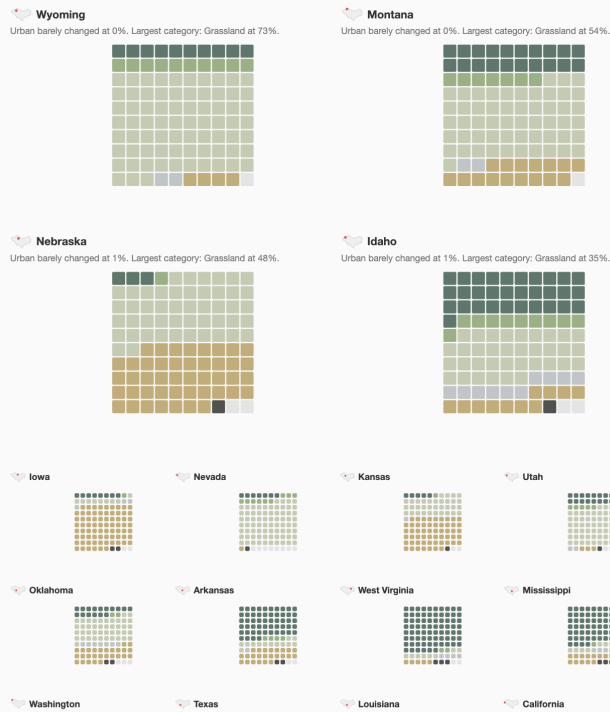
Forest Parks & Wildlife Grassland Defense Transport Farmsteads Misc. Cropland Urban Other



## The States Still Untouched by Sprawl in 2017

Land-use composition by state, ranked by least urban growth

Forest Parks & Wildlife Grassland Misc. Cropland Urban Other



**SORT BY**  
Urban growth (lowest)

Sorted category position:  
Top

**TIME-SERIES CHARTS**  
Stacked Area Line  
Stream Stacked Bar  
% Change Sparkline

**NON-TIME-SERIES CHARTS**  
Treemap Waffle

Year:  
2017

**COLOR SCHEME**  
Land (default)

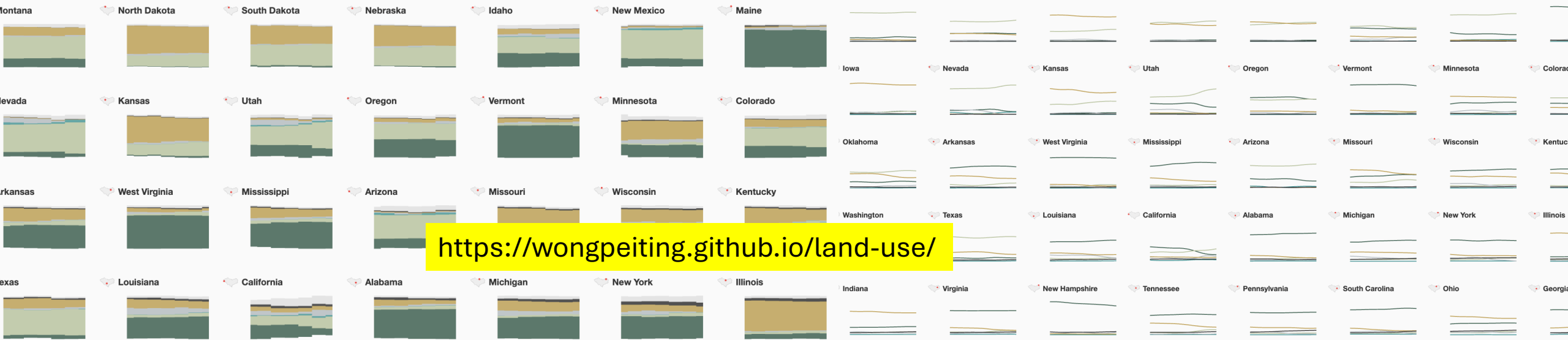
**TWEAKS**  
Desaturate non-sorted categories: 55%

Curve:  
Basis (soft)

Show state mini-map

Compact columns:  
8

**YEAR RANGE**  
1987 - 2017  
Start: [Slider]  
End: [Slider]



<https://wongpeiting.github.io/land-use/>

# Leveraging domain expertise

- Writing/tweaking prompts
- Developing test cases
- Scoring responses
- Fine-tuning LLM judges
- This all falls under **evals**

# Review Classifications

Confirmed: **600** / 600



**Confirm** →

Skip →

Unreviewed

All

Human only

AI only

**Export CSV**

Reset Cache

## GLOSSARY

### SUBJECT — WHAT IS IT ABOUT?

#### Celebrating Trump

The video's entire purpose is hero worship. GOAT framing, dramatic montages of Trump looking powerful, victory laps. Trump appears in almost every video — only pick this if the point IS Trump, not a policy he's in.

#### Attacking opponents

Mocking Democrats (Schumer, Pelosi, Biden, Jeffries), attacking media/CNN, ridiculing journalists. The mockery is the point. If it blames Democrats for an immigration issue, it's this — not enforcement.

#### Law enforcement

Immigration, ICE raids, deportation, border wall, "illegals", CBP, crime stats, DHS recruitment. The enforcement apparatus. If Trump appears alongside ICE footage, it's still this.

#### War & military

Iran strikes, Operation Epic Fury, military operations, combat footage, troops deploying. If it shows actual military hardware or strike footage, it's war — even if Trump appears heroically.

#### Culture & sports

Sports events, holidays, entertainment, Olympics, lifestyle moments. The Cristiano Ronaldo meeting, March Madness, St Patrick's Day. The soft stuff with no policy edge.

#### Governing & diplomacy

Executive orders, economic numbers ("\$10.5 trillion"), energy policy, legislation, foreign dignitary meetings, trade deals, summit handshakes. If Trump is shown with policy achievements (gas prices, executive orders), it's this — not "trump".

#### Promotion / meta

WH app launch, internships, roundups, "subscribe", DHS recruitment videos credited to DHS. Content about the account or government recruiting itself.

#### Unclear

Genuinely ambiguous. Cryptic posts with no context. If you can't tell what the video is about after watching it, pick this.

### PACKAGING — INSTITUTIONAL TO TEETERING ON FICTIONAL

#### ① Official

Press conference, signing ceremony, official statement. Clean editing. Talking heads. What government video normally looks like. No style, no flair.

#### ② Direct address



2025-11-13 · #273/600 · 1.7M views

**HUMAN VERIFIED**

The most iconic duo ❤️🇺🇸

[Open on TikTok](#) ↗️

**Notes:** Shows Trump and Melania holding hands and walking down a hallway at the White House. Music: Take My Breath Away (Love Theme from "Top Gun"), by Berlin.

## SUBJECT

Celebrating Trump

Attacking opponents

Law enforcement

War & military

Culture & sports

Governing & diplomacy

**Promotion / meta**

Unclear

## PACKAGING — ① INSTITUTIONAL → ⑦ TEETERING ON FICTIONAL

① Official

② Direct address

**③ Produced/cinematic**

④ TikTok-native

⑤ Pop culture mashup

⑥ Internet meme

⑦ Game interface

## TAGS

### AUDIO

**Trending song**

Meme audio clip

Bass drop / bass-boosted

**Dramatic / epic score**

Original audio

Hype language

### VISUAL EFFECTS

Deep-fried / distortion

Speed ramp / slow-mo

Greyscale to color

Stock hero / B-roll reuse

### NARRATIVE STRUCTURE

Comparison / before-after

Highlight reel / montage

Gamification (scores, achievements)

Fictional overlay on real footage

Punchline / reveal

Call & response to critics

### CONTENT FLAGS

Expletive / swearing

Names a specific person to mock

Real military footage used

AI-generated or CGI visuals

Trump dancing

Troll / shitpost

Provocative / edgy

Recruitment

**Aura farming**

Hyper-masculine

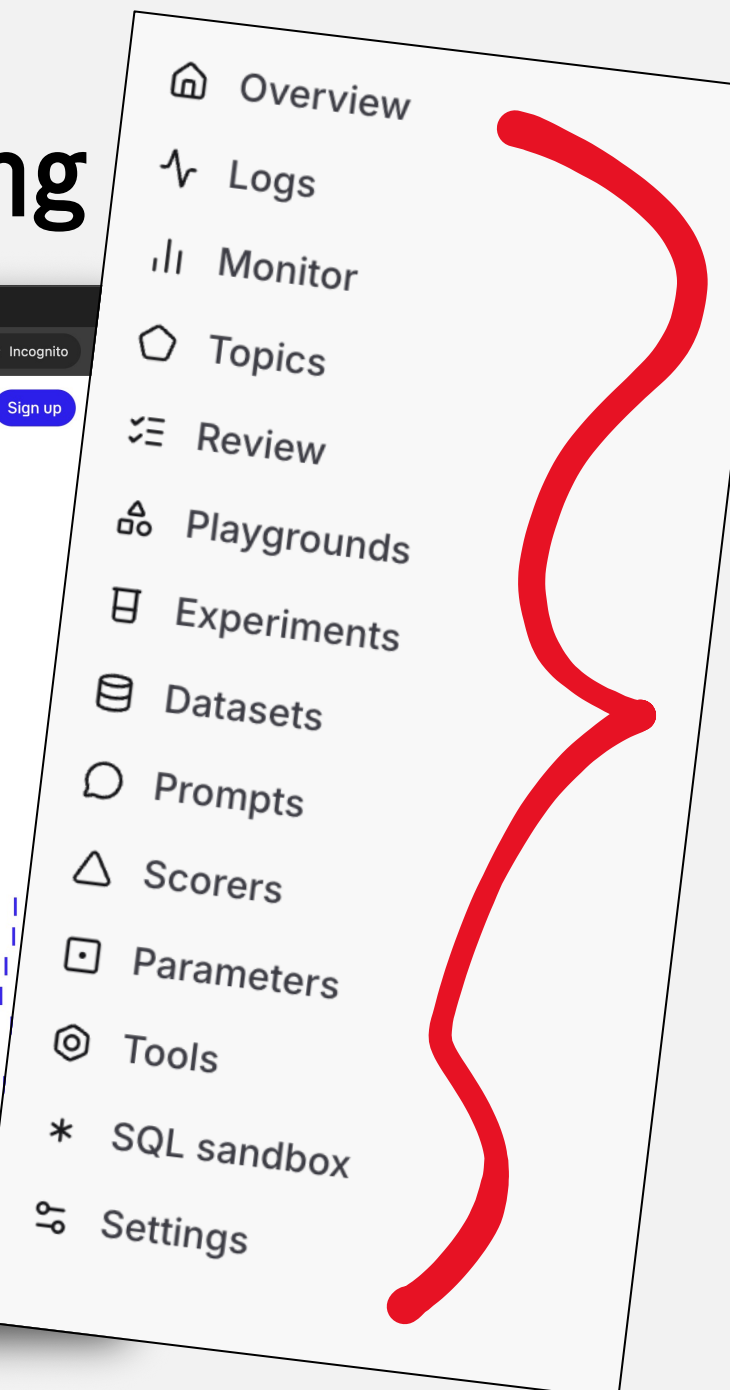
MAGA Minute

## NOTES

Shows Trump and Melania holding hands and walking down a hallway at the White House. Music: Take My Breath Away (Love Theme from "Top Gun"), by Berlin.

**Enter** confirm + next · **→** skip · **←** prev · **S** export

# Observability/evals/prompting




braintrust.dev (and 2,000 others)

The LA Times' new AI tool sympathized with the KKK. Its owner wasn't aware until hours later

By Liam Reilly, CNN  
8 min read · Published 3:30 PM EST, Wed March 5, 2025

22 comments



Washington Post's AI-generated podcasts rife with errors, fictional quotes

Intelligence for the New World Economy

Exclusive / Washington Post's AI-generated podcasts rife with errors, fictional quotes

Max Tani  
Media Editor, Semafor

Dec 11, 2025, 5:08pm EST MEDIA



Playground 2 - Playgrounds

braintrust.dev/app/Little%20Columns/p/default-otel-project/playgrounds/Playground%202

Little Columns / default-otel-project / Playgrounds

Playground 2 ...

Diff + Experiments Run

Base task

GPT-5 mini

User

You're a journalist. Is the following tip worth researching? `{{tip}}`

+ Message part

+ Message @ Tools Mustache MCP

Structured output worth\_rese...

Save prompt Draft

Comparison task

Claude 4.5 Haiku

User

You're a journalist. Is the following tip worth researching? `{{tip}}`

+ Message part

+ Message @ Tools Mustache MCP

Structured output worth\_rese...

Save prompt Draft

All rows view Filter Display Row

+ Task 2 + Scorer 1 Dataset 3 50

% Worth Researching Clear

Input

GPT-5 mini Base Claude 4.5 Haiku

Tradeoff Base has better latency. Comparison has

1 email: carter.linda.durham@gmail.com name: Linda Carter tip: I lost my left earring somewhere between the DPAC parking deck and the entrance to the theater last night. It is a gold hoop. If anyone at the station found it, or if a viewer turned it in, please let me know. It has great sentimental value. worth\_researching: false

Output Dec 10, 2025 0s 279 0.00 Output Dec 10, 2025 0.1s 36

Loop agent



```
email: carter.linda.durham@gmail.com
name: Linda Carter
tip: I lost my left earring somewhere between the DPAC parking deck and the entrance to the theater last night. It is a gold hoop. If anyone at the station found it, or if a viewer turned it in, please let me know. It has great sentimental value.
```

```
email: durham.watchdog.11@protonmail.com
name: Anonymous
tip: I saw a drone flying over the American Tobacco Campus near the water tower. It was hovering very still. I think it might be spying on the people eating at the restaurants. Is this legal? You should warn people not to eat outside until we know who is piloting it.
```

```
email: k.oconnor55@verizon.net
name: Kevin O'Connor
tip: My grandson, Timmy, just learned how to tie his shoes all by himself. He is five years old and
```

Input

GPT-5 mini

Base

Claude 4.5 Haiku

Tradeoff Base has better latency. Comparison has be

email: carter.linda.durham@gmail.com  
 name: Linda Carter  
 tip: I lost my left earring somewhere between the DPAC parking deck and the entrance to the theater last night. It is a gold hoop. If anyone at the station found it, or if a viewer turned it in, please let me know. It has great sentimental value.

Output Dec 10, 2025 0s 279 0.00  
 worth\_researching: false

Output Dec 10, 2025 0.1s 36 0.0  
 worth\_researching: false

% Worth Researching 100%

% Worth Researching 100%

email: durham.watchdog.11@protonmail.com  
 name: Anonymous  
 tip: I saw a drone flying over the American Tobacco Campus near the water tower. It was hovering very still. I think it might be spying on the people eating at the restaurants. Is this legal? You should warn people not to eat outside until we know who is piloting it.

Output Dec 10, 2025 0s 151 0.00  
 worth\_researching: true

Output Dec 10, 2025 0.1s 36 0.0  
 worth\_researching: false

% Worth Researching 0%

% Worth Researching 100%

email: k.oconnor55@verizon.net  
 name: Kevin O'Connor  
 tip: My grandson, Timmy, just learned how to tie his shoes all by himself. He is five years old and

Output Dec 10, 2025 0s 279 0.00  
 worth\_researching: true

Output Dec 10, 2025 0.1s 36 0.0  
 worth\_researching: false

Lo

# Scorers: code-based vs LLM-as-judge

```
from typing import Any
import re

def handler(input: Any, output: Any, expected: Any, metadata: dict[str, Any]):
    """
    A heuristic scorer to detect overly personal language.
    Fast and explainable, but imperfect.
    """
    output_text = str(output) if output else ""
    text = output_text.strip().lower()

    # Patterns that indicate personal/informal language
    PERSONAL_PATTERNS = [
        r"\bthank(s| you)\b",
        r"\bi (really )?(appreciate|feel|think|believe|agree|disagree)\b",
        r"\bi'm\b",
        r"\bmy (view|opinion|take|heart)\b",
        r"\bdear\b",
        r"\bsincerely\b",
        r"\bhere'?s\b",
        r"\bgreeting\b",
    ]

    matched_patterns = []
    for pattern in PERSONAL_PATTERNS:
        if re.search(pattern, text, flags=re.IGNORECASE):
            matched_patterns.append(pattern)

    return {
        "score": 0.0 if matched_patterns else 1.0,
        "name": "Not Too Personal",
        "metadata": {
            "matched_patterns": matched_patterns,
            "text": text,
            "input": input,
            "output": output,
            "expected": expected,
            "metadata": metadata,
        }
    }
```

You are grading whether the text uses overly personal, informal, or therapeutic language.

Evaluate the TEXT for:

- Greetings (e.g., "Dear", "Hi")
- Thanks or sign-offs (e.g., "Thank you", "Sincerely", "Best regards")
- First-person emotional framing (e.g., "I appreciate", "I feel", "I believe")
- Conversational or therapeutic tone

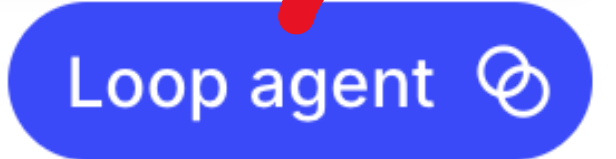
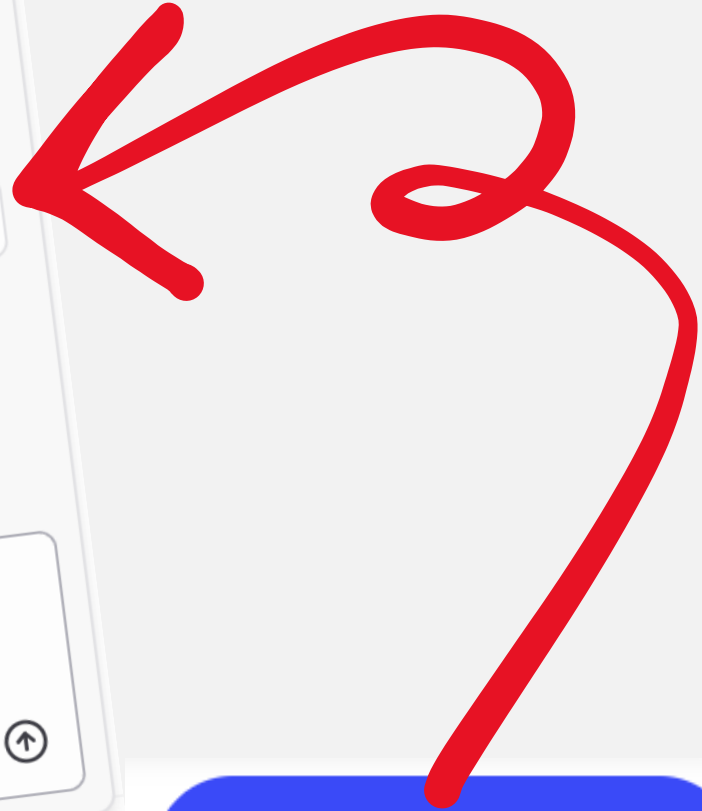
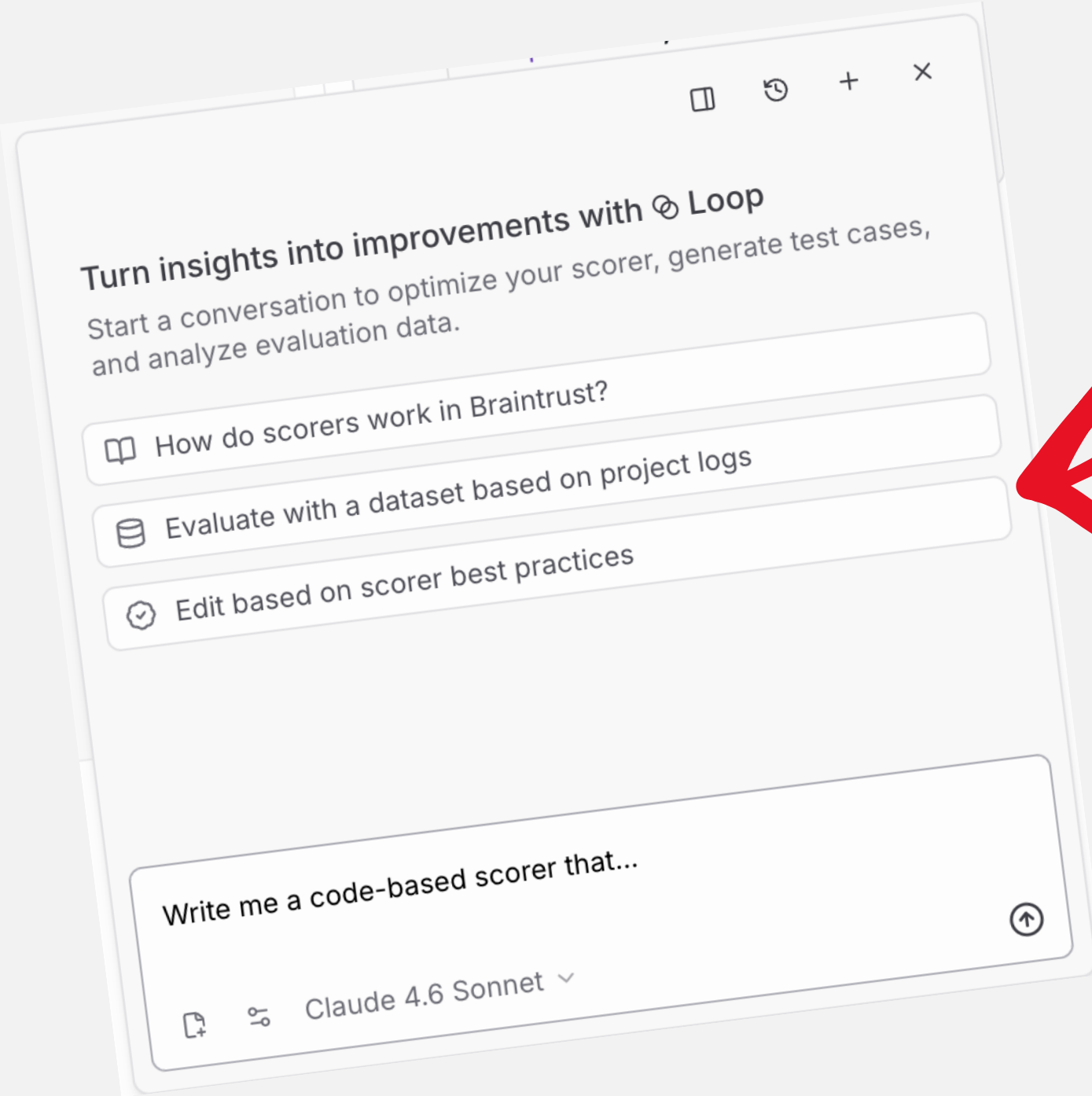
Grade as:

- **\*\*a\*\*** if the text is professional, analytical, and avoids personal language entirely
- **\*\*b\*\*** if the text has mild first-person framing once (like "I think") but is otherwise professional
- **\*\*c\*\*** if the text contains greetings, thanks, sign-offs, or strongly personal/therapeutic tone

TEXT:

{{output}}

Return ONLY the letter **\*\*a\*\***, **\*\*b\*\***, or **\*\*c\*\***.

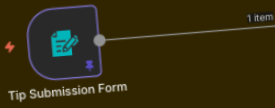


# Your evals todo-list cheatsheet

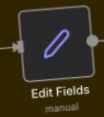
1. Build a dataset of possible **inputs**
  - Real-life are best, but there are some techniques for synthetic ones
  - Have a grouchy teammate who says “it’ll *NEVER* work in *THIS* case” try to write the hard ones.
  - “If it can’t do these successfully, it can’t go to production”
2. Make your **prompt** alone or with input
3. Throw the **inputs** and **prompt** together to get **outputs**
4. Have the good-spirited review a spreadsheet of **inputs + outputs** and score them.
  - Yes/No is much easier to deal with than 1–5 scores
5. Use their responses to teach an **LLM judge** about edge cases
6. **Repeat!!!** Although if you get to 100% then your evals probably aren’t difficult enough

# What about n8n?

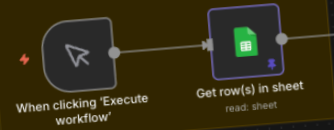
Get test or live input from form



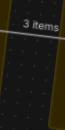
Standardize input about tips



Get test input from Google Sheet



Get test input from Google Sheet

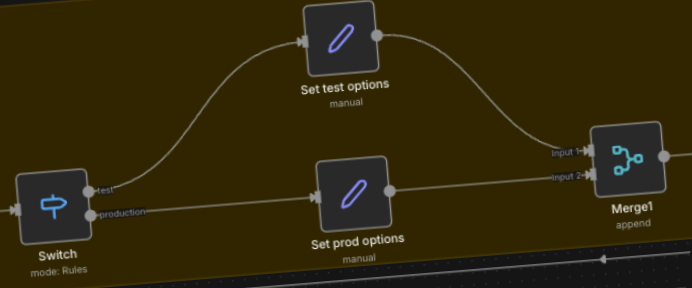


Standardize input about tips



Is it test mode? Is it production mode?

Update settings accordingly (if you have any)  
One big improvement you could make is set the document or sheet ID that you're updating based on whether it's test or prod! We're always updating the same sheet, which is probably not the best approach (it's easy to do, though)



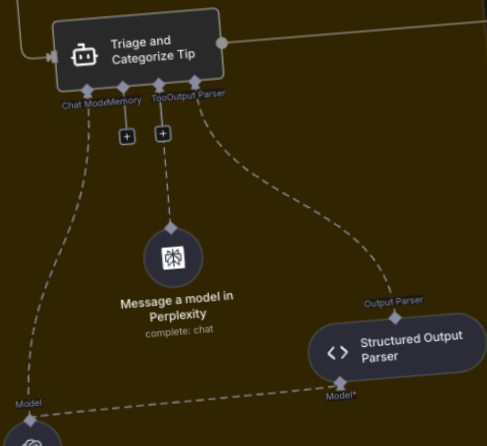
Is it test mode? Is it production mode?

Update settings accordingly (if you have any)  
One big improvement you could make is set the document or sheet ID that you're updating based on whether it's test or prod! We're always updating the same sheet, which is probably not the best approach (it's easy to do, though)



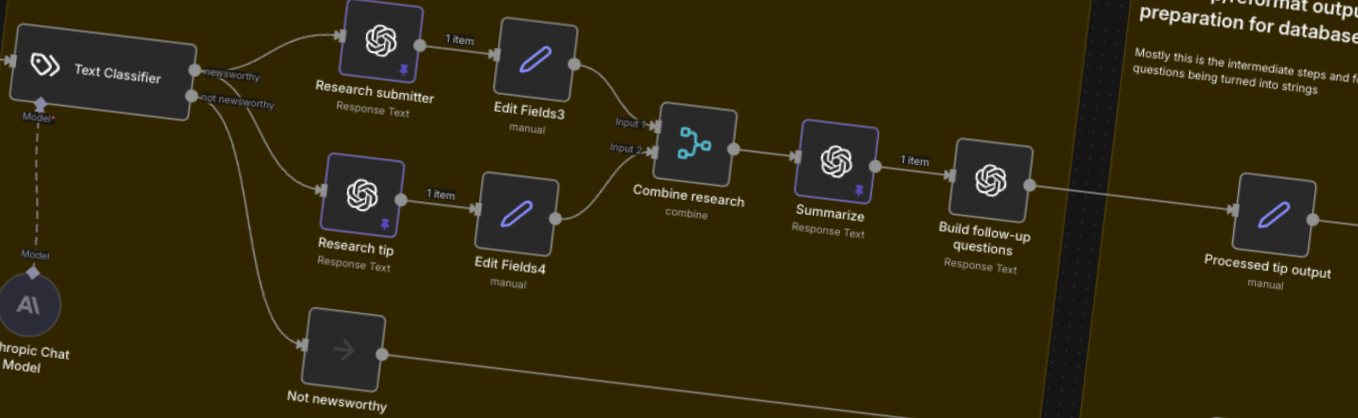
Investigate and triage the tip.

In order to track the agent, intermediate steps must be turned on.  
Auto-fix format is also turned on in the Structured Output tool.  
LLM model is determined from 'Set universal options'



Research + summarization

In this example, we're using many different specific queries to the LLM to do our research.



Clean up/reformat output preparation for database

Mostly this is the intermediate steps and follow-up questions being turned into strings



NEW: know what to automate before you build with CrewAI Discovery

Sign In Sign Up



Enterprise Open Source Docs Resources Pricing

The screenshot displays the CrewAI Discovery interface within a browser window. The main workspace, titled 'Untitled Project', features a 'Canvas' view with a 'Run' button. A workflow is being built, starting with a 'Triggers' block (currently set to 'No triggers configured') which connects to a 'Monitor News Coverage' task. This task is configured to search for and analyze all news articles mentioning a specific company name and their executive leadership. The workflow then branches into three parallel tasks, each powered by 'gpt-5.4': 'News Reputation Researcher' (monitoring and gathering recent news articles), 'Social Media Reputation Researcher' (searching and analyzing social media mentions), and 'Reputation Analyst' (synthesizing findings from news and social media research). A large play button is centered over the workflow. On the left, a chat interface shows a conversation where the AI assistant is creating a reputation monitoring crew with three agents: a news researcher, a social media researcher, and a reputation analyst. At the bottom of the browser window, a video player controls are visible, showing a progress bar at 0:23 and a total duration of 3:06. Below the browser window, two buttons labeled 'No Code' and 'CLI' are displayed.

# Staying model/vendor agnostic

- You can't experiment if you're married to a vendor!
- Vendor lock-in is the *enemy*. Some examples:
  - **OpenAI**: Custom GPTs, Assistants API, specific parts of Agents SDK
  - **Anthropic**: Claude managed agents, Claude Code settings
  - **Google**: Gems, Workspace integrations, specific parts of ADK
- **Although** DIYing might get you a worse product at a higher cost
- Ability and cost are in constant flux
  - ...although except Gemini is always the cheapest

```
from openai import OpenAI
client = OpenAI()

response = client.responses.create(
    model="gpt-5.5",
    input="Write a one-sentence bedtime story about a unicorn."
)

print(response.output_text)
```

```
import anthropic

client = anthropic.Anthropic()

message = client.messages.create(
    model="claude-opus-4-7",
    max_tokens=1024,
    messages=[{"role": "user", "content": "Hello, Claude"}],
)

print(message.content)
```

```
from google import genai

client = genai.Client(api_key='GEMINI_API_KEY')
response = client.models.generate_content(
    model='gemini-2.5-flash',
    contents={'text': 'Why is the sky blue?'},
)

print(response.text)
```

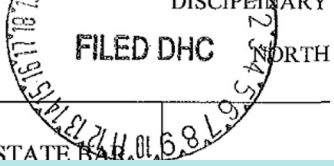
# LiteLLM

- Call any provider using the same `completion()` interface — no re-learning the API for each one
- Consistent output format regardless of which provider or model you use
- Built-in retry / fallback logic across multiple deployments via the `Router`
- Self-hosted `LLM Gateway (Proxy)` with virtual keys, cost tracking, and an admin UI



```
from litellm import completion
import os

response = completion(
    model="gemini/gemini-pro",
    messages=[
        {"role": "user", "content": "Hello, how are you?"}
    ]
)
```



7. Defendant filed a Motion to Show Cause September 2014.

8. Defendant failed to serve the opposi September 2014.

10. A.J. traveled from another county to appe

11. A.J. asked Defendant several times to

Panel enters the following:

CONCLUSIONS OF LAW

1. All parties are properly before the Hearing Panel, and the Panel has jurisdiction

discipline pursuant to N.C. Gen. Stat. Professional Conduct in effect at the time

- (a) By failing to notify A.J. that the hearing had been continued from the 30 September 2014 calendar and by otherwise failing to maintain communication with A.J., Defendant failed to keep his client reasonably informed about the status of the matter in violation of Rule 1.4(a)(3);
- (b) By failing to respond to A.J.'s inquiries about the status of the matter, Defendant failed to comply promptly with a reasonable request for information in violation of Rule 1.4(a)(4); and
- (c) By failing to provide a response to the fee dispute in good faith in the fee dispute resolution process.

# Structured outputs with Pydantic

NICHOLAS S. ACKERMAN, Attorney,

The composed pursuant to Plaintiff, t S. Ackern parties stip order and freely and any way th

Bas consent of the follow

1. North Car Chapter 84 Carolina S

2. admitted to rules, regulations, and Rules of Professional Conduct of the laws of the State of North Carolina.

3. Defendant was properly served with process, Hearing Panel with due notice to all parties.

4. During the relevant period referred to herein practice of law in Greensboro, Guilford County, North Carolina

5. On 25 September 2014, A.J. retained Defenc contempt order and custody modification.

20. On 12 F required Defendant to p

21. Defendar

22. On 9 Ma dispute resolution proces 15G0195.

23. Defendant represent A.J. resolution process and appear in court on her behalf.

```

class FilingDetails(BaseModel):
    plaintiff: str
    defendant: str
    charges: str
    city_of_practice: str
    penalty: str
  
```

```

{
  'Plaintiff': 'The North Carolina State Bar',
  'Defendant': 'Nicholas S. Ackerman',
  'Charges': 'Failure to keep client reasonably inform the status of the matter; Failure to comply promptly w reasonable request for information; Failure to partici good faith in the fee dispute resolution process',
  'City_of_practice': 'Greensboro',
  'Penalty': 'Suspension of license for one (1) year, for two (2) years as long as Defendant complies with conditions.'
}
  
```

CONCLUSIONS REGARDING DISCIPLINE

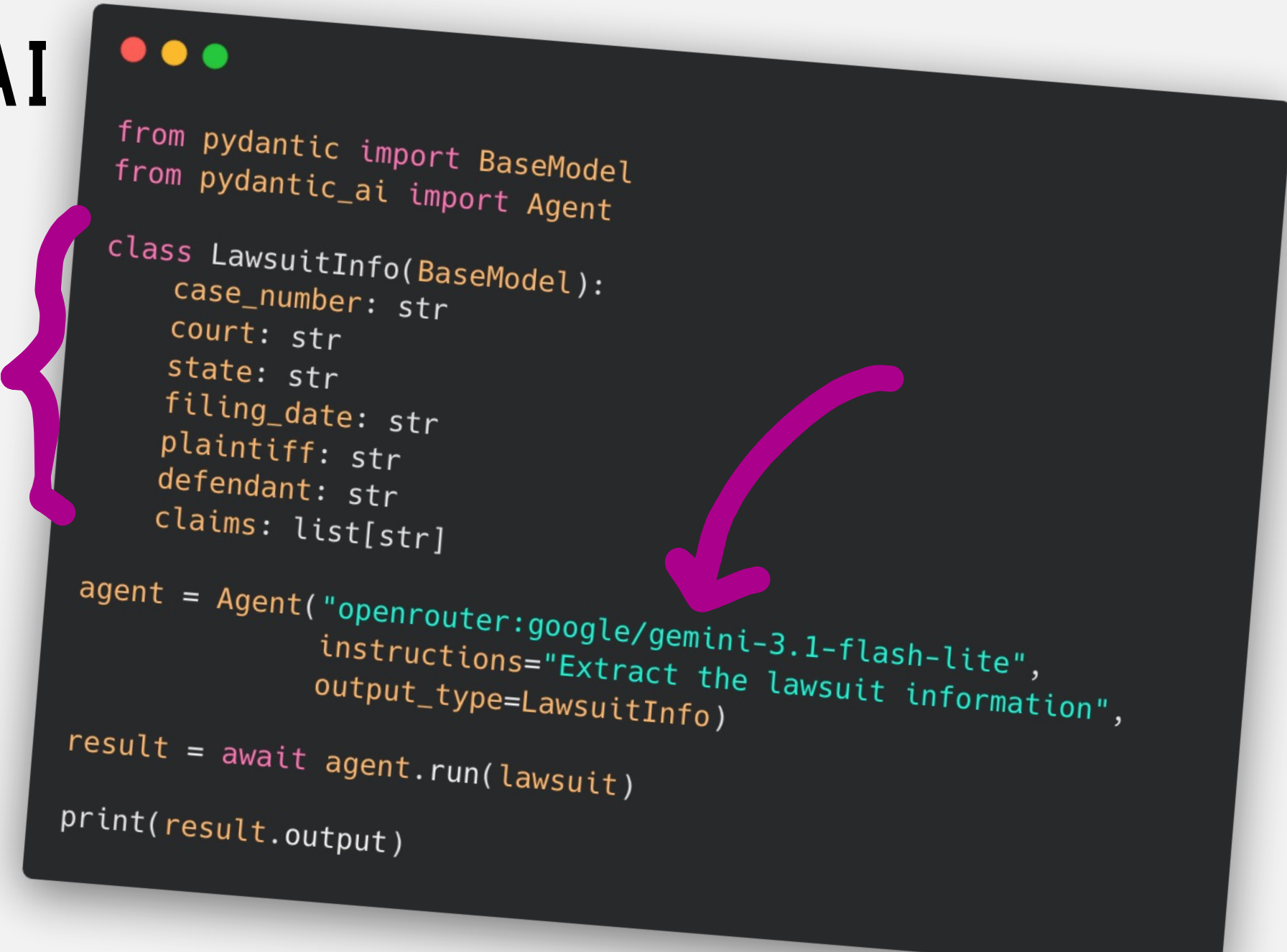
# Pydantic AI

```
from pydantic import BaseModel
from pydantic_ai import Agent

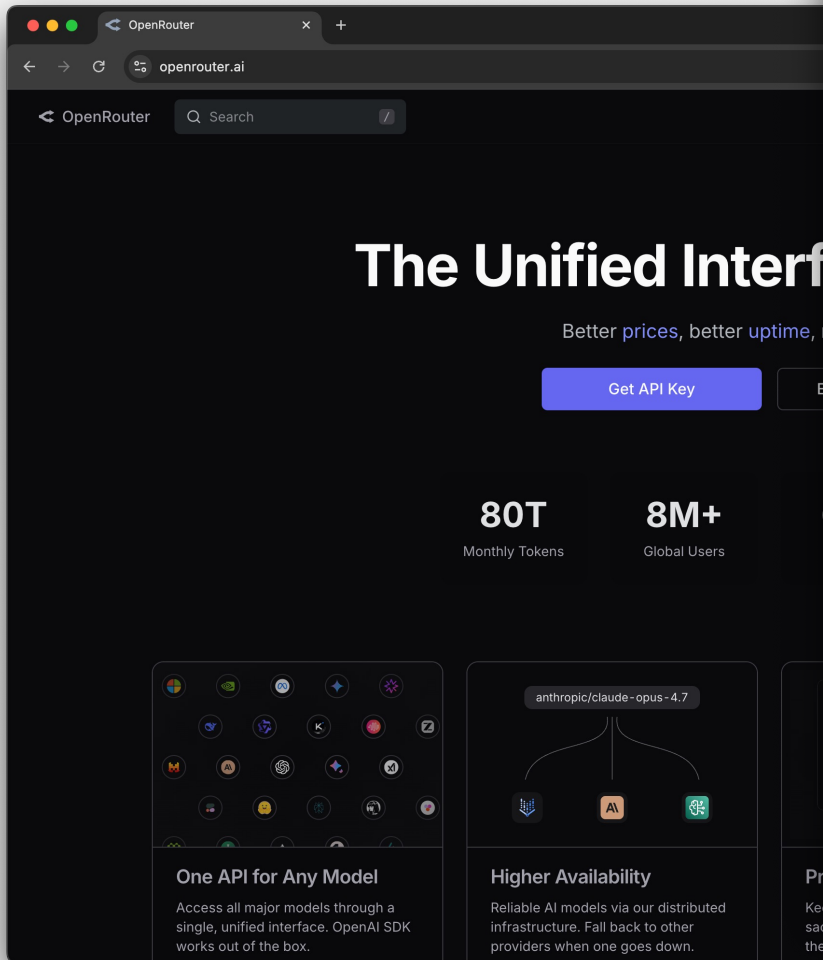
class LawsuitInfo(BaseModel):
    case_number: str
    court: str
    state: str
    filing_date: str
    plaintiff: str
    defendant: str
    claims: list[str]

agent = Agent("openrouter:google/gemini-3.1-flash-lite",
             instructions="Extract the lawsuit information",
             output_type=LawsuitInfo)

result = await agent.run(lawsuit)
print(result.output)
```



# OpenRouter



The Unified Interface

Better prices, better uptime, more models

[Get API Key](#)

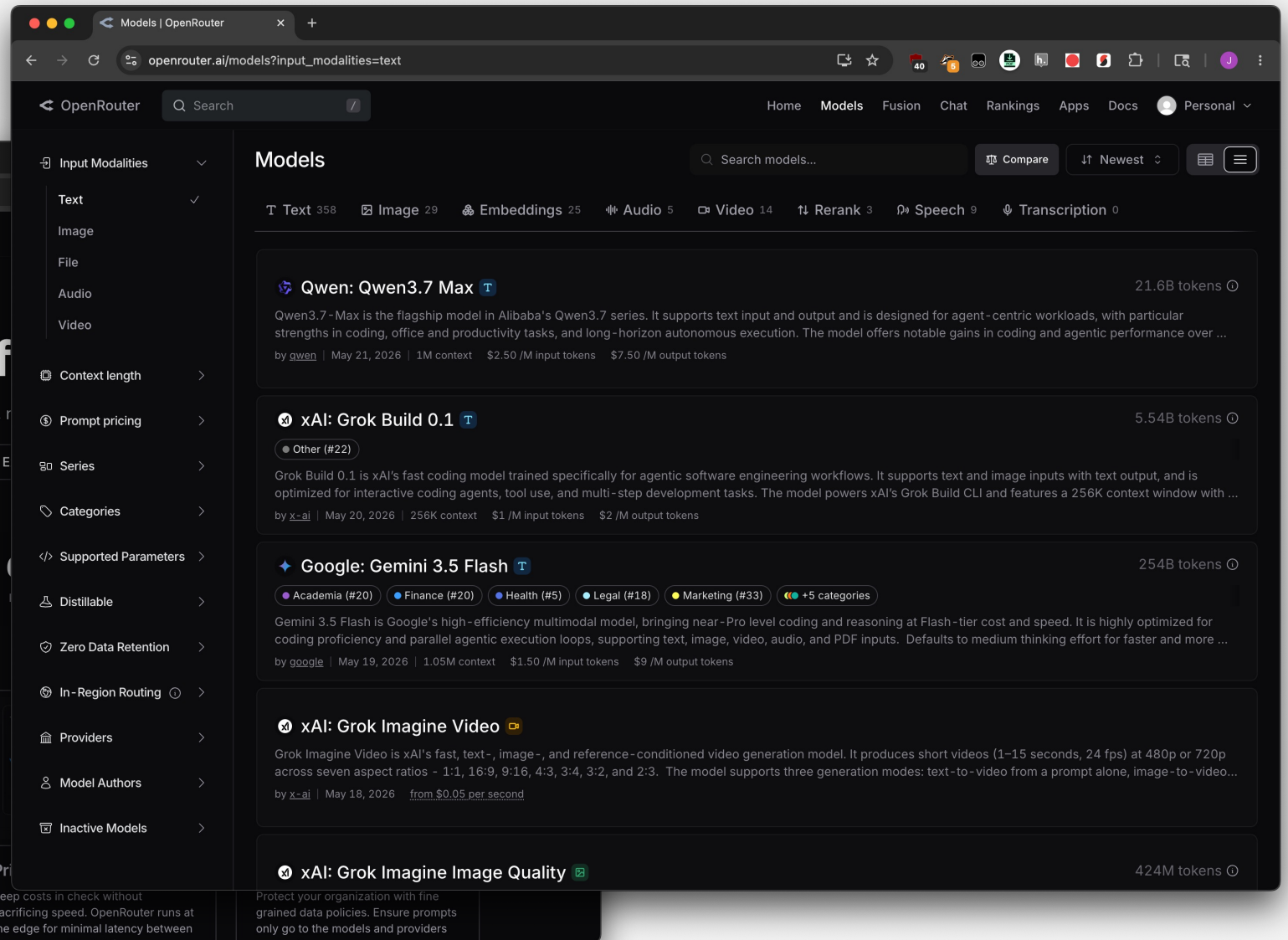
**80T** Monthly Tokens

**8M+** Global Users

**One API for Any Model**  
Access all major models through a single, unified interface. OpenAI SDK works out of the box.

**Higher Availability**  
Reliable AI models via our distributed infrastructure. Fall back to other providers when one goes down.

**Pr...**  
Keep costs in check without sacrificing speed. OpenRouter runs at the edge for minimal latency between...



Models | OpenRouter

openrouter.ai/models?input\_modalities=text

OpenRouter Search

Home Models Fusion Chat Rankings Apps Docs Personal

## Models

Search models...

Compare Newest

Text 358 Image 29 Embeddings 25 Audio 5 Video 14 Rerank 3 Speech 9 Transcription 0

- Qwen: Qwen3.7 Max** (Text) 21.6B tokens  
Qwen3.7-Max is the flagship model in Alibaba's Qwen3.7 series. It supports text input and output and is designed for agent-centric workloads, with particular strengths in coding, office and productivity tasks, and long-horizon autonomous execution. The model offers notable gains in coding and agentic performance over ...  
by qwen | May 21, 2026 | 1M context | \$2.50 /M input tokens | \$7.50 /M output tokens
- xAI: Grok Build 0.1** (Text) 5.54B tokens  
Other (#22)  
Grok Build 0.1 is xAI's fast coding model trained specifically for agentic software engineering workflows. It supports text and image inputs with text output, and is optimized for interactive coding agents, tool use, and multi-step development tasks. The model powers xAI's Grok Build CLI and features a 256K context window with ...  
by x-ai | May 20, 2026 | 256K context | \$1 /M input tokens | \$2 /M output tokens
- Google: Gemini 3.5 Flash** (Text) 254B tokens  
Academia (#20) Finance (#20) Health (#5) Legal (#18) Marketing (#33) +5 categories  
Gemini 3.5 Flash is Google's high-efficiency multimodal model, bringing near-Pro level coding and reasoning at Flash-tier cost and speed. It is highly optimized for coding proficiency and parallel agentic execution loops, supporting text, image, video, audio, and PDF inputs. Defaults to medium thinking effort for faster and more ...  
by google | May 19, 2026 | 1.05M context | \$1.50 /M input tokens | \$9 /M output tokens
- xAI: Grok Imagine Video** (Image) 424M tokens  
Grok Imagine Video is xAI's fast, text-, image-, and reference-conditioned video generation model. It produces short videos (1-15 seconds, 24 fps) at 480p or 720p across seven aspect ratios - 1:1, 16:9, 9:16, 4:3, 3:4, 3:2, and 2:3. The model supports three generation modes: text-to-video from a prompt alone, image-to-video...  
by x-ai | May 18, 2026 | from \$0.05 per second
- xAI: Grok Imagine Image Quality** (Image) 424M tokens

Models - Top Weekly | OpenR

openrouter.ai/models?input\_modalities=text&output\_modalities=text&order=top-weekly

OpenRouter Search

Home Models Fusion Chat Rankings Apps Docs Personal

Models Search models... Compare Top Weekly

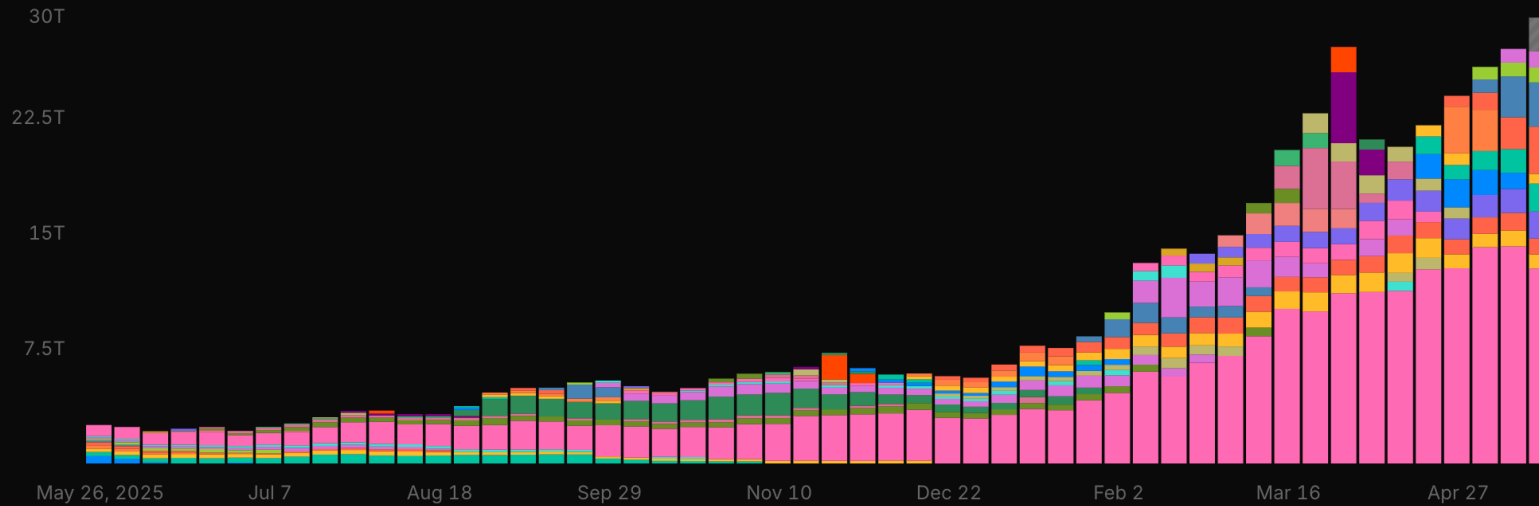
Text 358
  Image 29
  Embeddings 25
  Audio 5
  Video 14
  Rerank 3
  Speech 9
  Transcription 0

Model Name	Weekly Tokens	Input	Output	Context	Released
<a href="#">DeepSeek: DeepSeek V4 Flash</a>	3.46T	\$0.10	\$0.20	1,048,576	Apr 23, 2026
<a href="#">Tencent: Hy3_preview</a>	3.29T	\$0.066	\$0.26	262,144	Apr 22, 2026
<a href="#">Anthropic: Claude Opus 4.7</a>	1.95T	\$5	\$25	1,000,000	Apr 16, 2026
<a href="#">Anthropic: Claude Sonnet 4.6</a>	1.91T	\$3	\$15	1,000,000	Feb 17, 2026
<a href="#">Owl Alpha</a>	1.22T	\$0	\$0	1,048,756	Apr 28, 2026
<a href="#">Google: Gemini 3 Flash Preview</a>	1.18T	\$0.50	\$3	1,048,576	Dec 17, 2025
<a href="#">DeepSeek: DeepSeek V4 Pro</a>	1.1T	\$0.435	\$0.87	1,048,576	Apr 23, 2026
<a href="#">DeepSeek: DeepSeek V3.2</a>	1.07T	\$0.252	\$0.378	131,072	Dec 1, 2025
<a href="#">StepFun: Step 3.5 Flash</a>	755B	\$0.09	\$0.30	262,144	Jan 29, 2026
<a href="#">MoonshotAI: Kimi K2.6</a>	700B	\$0.73	\$3.49	262,144	Apr 20, 2026
<a href="#">Google: Gemini 2.5 Flash Lite</a>	655B	\$0.10	\$0.40	1,048,576	Jul 22, 2025
<a href="#">NVIDIA: Nemotron 3 Super (free)</a>	642B	\$0	\$0	1,000,000	Mar 11, 2026
<a href="#">Google: Gemini 2.5 Flash</a>	607B	\$0.30	\$2.50	1,048,576	Jun 17, 2025
<a href="#">MiniMax: MiniMax M2.7</a>	598B	\$0.279	\$1.20	204,800	Mar 18, 2026
<a href="#">OpenAI: GPT-5</a>	518B	\$5	\$20	1,050,000	Apr 24, 2026

- Input Modalities
  - Text
  - Image
  - File
  - Audio
  - Video
- Context length
- Prompt pricing
- Series
- Categories
- Supported Parameters
- Distillable
- Zero Data Retention
- In-Region Routing
- Providers
- Model Authors
- Inactive Models

## 📊 Top Models

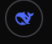
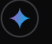

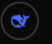

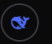




Weekly usage of models across OpenRouter



## 🏆 LLM Leaderboard

Compare the most popular models on OpenRouter ⓘ

This Week ↕

1.	 <b>DeepSeek V4 Flash</b> by <a href="#">deepseek</a>	3.21T tokens ↑73%	6.	 <b>Gemini 3 Flash Preview</b> by <a href="#">google</a>	1.13T tokens ↓0%
2.	 <b>Hy3 preview</b> by <a href="#">tencent</a>	3.08T tokens ↑16%	7.	 <b>DeepSeek V4 Pro</b> by <a href="#">deepseek</a>	1.03T tokens ↑17%
3.	 <b>Claude Opus 4.7</b> by <a href="#">anthropic</a>	1.84T tokens ↑18%	8.	 <b>DeepSeek V3.2</b> by <a href="#">deepseek</a>	1.03T tokens ↑3%
4.	 <b>Claude Sonnet 4.6</b> by <a href="#">anthropic</a>	1.82T tokens ↑18%	9.	 <b>Step 3.5 Flash</b> by <a href="#">stepfun</a>	728B tokens ↑10%
5.	 <b>Owl Alpha</b> by <a href="#">openrouter</a>	1.15T tokens ↑46%	10.	 <b>Kimi K2.6</b> by <a href="#">moonshotai</a>	675B tokens ↓41%

Name ⓘ
✕

look at this key!

Credit limit (optional) ⓘ

5

Reset limit every... ⓘ

Daily
⌵

Expiration ⓘ

30 days
⌵

Key will expire on: Jun 23, 2026, 7:14 AM EDT

Create

Expires	Last Used	Usage	Limit	
May 30, 2026, 3:50 PM EDT	May 23, 2026, 3:51 PM EDT	\$ 0.042	\$2	WEEK
Jun 21, 2026, 12:09 PM EDT	May 23, 2026, 8:39 PM EDT	\$ 4.83	\$10	TODAY
May 27, 2026, 1:25 PM EDT	May 21, 2026, 1:55 AM EDT	\$ 0.128	unlimited	TOTAL
Never	May 17, 2026, 2:00 AM EDT	\$ 3.91	unlimited	TOTAL
Expired	May 13, 2026, 4:24 PM EDT	\$ 34.50	unlimited	TOTAL
Never	May 8, 2026, 8:05 PM EDT	< \$0.001	\$20	WEEK
Never	May 14, 2026, 2:05 PM EDT	\$ 3.23	\$5	TODAY
Never	Apr 16, 2026, 9:35 PM EDT	\$ 2.60	\$20	TODAY

# But you need evals to do it!!!

- Otherwise you won't know if the change is worth it.



# Newsroom Infrastructure for AI Experimentation

Jonathan Soma  
Columbia University  
@dangerscarf · js4571@columbia.edu

<https://bit.ly/dh26-infra>