# From Vibes to Scores

Build your own AI benchmark!

**Jonathan Soma**
js4571@columbia.edu

# What are we doing here?

Hi, I'm Soma!

- **Lede Program:** 10 weeks, remote-friendly, summer 2026
- **Data Journalism MS:** 12 months, deadline Jan 12

# Operationalize newsroom AI

- **Before:**
  - Little baby prompt
  - Custom GPT
  - Vibes, vibes, vibes
- **After:**
  - Multi-step agentic pipeline that you can put into prod!
  - Kidnapping knowledge from the tech world
  - Fame and fortune

Where this session is from

## We will learn how to…

- Compare effectiveness of different **AI prompts and models**

- Set up the **processes** to observe the "magic" going on inside of your AI tools

- **Structure eval test sets** to effectively cover your bases so you don't get blindsided

- **Build agentic AI** tools yourself (and what in the world agentic AI even is)

- Bring non-technical folks and **AI haters** to the table

- Why we even need to do this at all?? What about **vibes???**

# Stay in the Loop

Get exclusive updates, early access to new features, and weekly insights delivered straight to your inbox.

**Full Name**

John Doe

**Email Address**

john@example.com

Join 50,000+ Subscribers

No Spam

Weekly Updates

Premium Content

We respect your privacy. **Privacy Policy** · **Terms**

# Stay in the Loop

Get exclusive updates, early access to new features, and weekly insights delivered straight to your inbox.

**Full Name**

John Doe

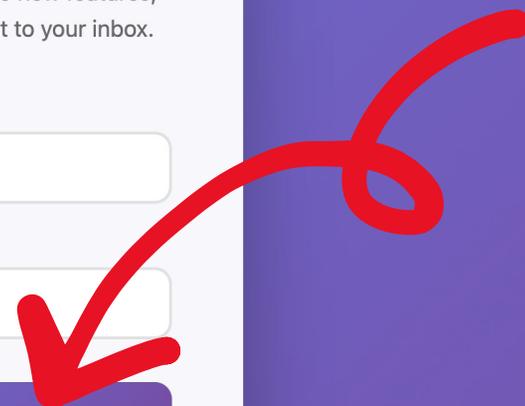**Email Address**

john@example.com

Join 50,000+ Subscribers

No Spam

Weekly Updates

Premium Content

We respect your privacy. **Privacy Policy** · **Terms**

AI is an unpredictable slot machine

(but we still love it)

what's the temperature of th...?

The surface of the Sun — what astronomers call the photosphere — has a...ping the photosphere (the visible...

average temperature of about **5,500 °C (5,800 K or ~9,940 °F).**

The **surface temperatu**...

"surface" we see) — is

Deeper inside, the ter...

Above it, the Sun's atmosphere actually gets much hotter:

- **Chromosphere:** ~10,000 °C
- **Corona:** over 1,000,000 °C

- **Core:** ~15 millio...
- **Corona** (the ou...

reaching 1–3...

Would you like me to explain why the outer atmosphere (corona) is hotter than the surface, even though it's farther away from the Sun's core?

Do you want me to expla...

though it's farther away from the Sun's c...

==inputs always change==

what's the temperature of the surface of the sun?

how many street cats live in athens?

what height should i mount wall speakers at?

What's your angle?

LET'S EMBRACE BOTH

It's unpredictable!

- Evaluation

- Evaluation

- *Evaluation!!!*

**Failures are good, actually?**



Business / Media

# 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

🌐 8:54 AM FRIDAY JANUARY 9, 2026

D.C.        BXL        LAGOS        RIYADH        BEIJING        SG

# SEMAFOR

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        SHARE ⬀

# Let's make it happen

- **Use case:** "generate opposing perspectives to be displayed on opinion stories"

- We will do BETTER than the LA Times!!

- Specification:
  - **Input:** …an article?
  - **Output:** …a response?

- **Braintrust!!!!**
- Arize Phoenix
- Comet Opik
- Langsmith
- Langfuse
- Promptfoo
- *A zillion more*

Playgrounds, not prompts!

Prompting meta-settings

models

API keys

# Create a dataset



**Left window (Datasets list):**

Little Columns

**Project**

default-otel-project

- Overview
- Logs
- Topics
- Monitor
- Review
- Playgrounds
- Experiments
- Datasets
- Prompts
- Scorers
- Parameters
- Tools
- SQL sandbox

Datasets

+ Dataset | ...ts view | Filter | Display

| | Name | Description |
|---|---|---|
| | Dataset 17 | – |
| | Stories | – |
| | Dataset 15 | – |
| | Dataset 14 | – |
| | Dataset 12 | – |
| | articles | – |
| | Dataset 11 | – |
| | Dataset 10 | – |
| | Dataset 9 2 | – |
| | Dataset 9 | – |

**Right window (Dataset 18):**

default-otel-project › Datasets

Dataset 18 ...

+ Row

## This dataset is empty

To populate this dataset, manually create a row, upload a CSV/JSON file, or input data programmatically

Drag and drop CSV or JSON file here, or click to select

Loop agent

Create a dataset

Changing options

Comparing prompts

My Project › Playgrounds

Upgrade

⚠ **Playground 1** ⋯

Diff ⬤

+ Experiments

● Base task

GPT-5 mini

Params

**User** ⌄

☐ Write a one-sentence response to the article titled "**{{input.headline}}**", content is:

**{{input.story}}**

+ Message part ⌄

+ Message  ◎ Tools  (×) Mustache ⌄  ⚡ MCP  T Text output ⌄

◯ Save prompt

Draft

● Comparison task

GPT-5 mini

Params

**System** ⌄

Write a GOOD!!!!!! one-sentence response to the article titled "**{{input.headline}}**", content is:

**{{input.story}}**

+ Message  ◎ Tools  (×) Mustache ⌄  ⚡ MCP  T Text output ⌄

◯ Save prompt

Draft

⊛ All rows view ⌄   ⚏ Filter   ⇅ Display   ⬇   + Row

+ Task 2  + Scorer  ⬡ Opinion Pieces

Overview
Logs
Topics
Monitor
Review
Playgrounds
Experiments
Datasets
Prompts
Scorers
Parameters
Tools
SQL sandbox
Settings

# Comparing prompts

100 years ago, Anaheim recalled its KKK city council. Why don't we remember?
Members of the KKK in Anaheim in 1924.
Members of the KKK in Anaheim in 1924. (Anaheim Public Library)
EL SEGUNDO CA DECEMBER 12, 2019 -- Gustavo Arellano, reporter for the Los Angeles Times.
By Gustavo Arellano
Columnist

Evaluating...

↗ Output

Gustavo Arellano's reminder that Anaheim once recalled KKK councilmembers — then largely erased that history — is a potent call to remember and teach how ordinary people can successfully resist organized racism before it returns.

↗ Output

Evaluating...   —

This article powerfully reminds us that forgetting Anaheim's KKK past isn't just an oversight, but a dangerous erasure of lessons crucial for confronting present-day hate.

---

2  headline: California's first partner pushes to regulate AI while Trump and tech bros thunder forward
story: |-
  California's first partner pushes to regulate AI while Trump and tech bros thunder forward
  California First Partner Jennifer Siebel Newsom, shown in 2023.
  California First Partner Jennifer Siebel Newsom, shown in 2023. (Christina House / Los Angeles Times)
  Anita Chabria. (Ricardo DeAratanha / Los Angeles Times)
  By Anita Chabria

↗ Output

Evaluating...

Jennifer Siebel Newsom's push for California-led AI guardrails is a timely and necessary counterbalance to a federal approach that risks ceding public safety and research to tech moguls more interested in profit than protecting children and society.

↗ Output

Evaluating...   —

While tech titans and Trump push for unchecked AI advancement, California's First Partner champions crucial regulations, prioritizing child safety over profit and demonstrating that responsible innovation is possible.

---

3  headline: Amid catastrophic loss, the unshakable allure of the San Gabriel Mountains
story: |-
  Amid catastrophic loss, the unshakable allure of the San Gabriel Mountains
  Meteorologist Edgar McGregor in Irwindale.
  Meteorologist Edgar McGregor knows the science behind Santa Ana winds and keeps a watchful eye out at the Santa Fe Dam Recreation Area in Irwindale.

↗ Output

Evaluating...

Steve Lopez's portrait of Edgar McGregor and the San Gabriel Mountains is a moving reminder that the region's beauty is inseparable from escalating risks — and that community stewardship, smarter planning and urgent climate action are all needed to protect lives and landscapes.

↗ Output

Evaluating...   —

Despite catastrophic losses, the San Gabriel Mountains' enduring allure lies in their dynamic beauty and the human spirit's resilience in the face of nature's formidable power.

Loop agent ⊘

# Evals, first step

- We need data to test this on.

# Where does the data come from?

- If you're lucky, you have it
- …but it's more exciting if you don't!
- The **Gulf of Comprehension** is bigger than you can imagine

# Synthetic data

- How do you make it?

- Ask an LLM?
  - Noooo
  - ………yessss?

- Let's do this in a structured way, with **dimensions**

# Dimensions of variation

- What are the ways each input could be different?
- Think about those that have the **most variance**, the biggest swings or the highest number of values
- **For submitted tips, how could each tip or each person submitting a tip be different?**
  - Knowledge of what's interesting to a journalist
  - Conspiracy theorist
  - Spam
  - Good at communicating in English
  - Length and amount of description

# Making inputs (breaking outputs)

- Generate based on **tuples**

- All combinations:
  - (A1, B1, C1), (A2, B1, C1), (A2, B2, C1)

- Use your intuition, filter what makes sense
  - A conspiracy theorist probably doesn't have a good idea

- Should only pick 4 or 5 dimensions, go with the ones that have the highest variation and most importance

- End up with **100 total examples**

What are our dimensions?



Business / Media

## 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

SIGN IN

8:54 AM  FRIDAY JANUARY 9, 2026

D.C.    BXL    LAGOS    RIYADH    BEIJING    SG

# SEMAFOR

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    SHARE ⬀

# How do we write our tests?

- Code-based evaluators

- LLM-as-judge

- Mandatory for **tracking performance over time** and **model drift**

# Code-based evaluators

- Fast and cheap! **Rely on these as much as possible.**
- Not "does this trace include email?"
  - **Search the JSON** to see that **send_email** tool was used
- Regular expressions and string matching are your best friends
- Don't ask an LLM "does this have 3 bullet points?"
  - Just check for **len(text.split(",")) === 3**
- Look for common "something broke" phrases:
  - As an LLM, Sorry I can't, etc

# LLM-as-judge

- Too much nuance or interpretation for code-based evals
- Only ever do Pass/Fail
- You need a secondary process to confirm it agrees with you, yet another gulf of specification to cross (you're running an eval on an eval)
- Each test should check for a **single specific failure mode**
- A good LLM as a judge should function as *scaling you*

Select scorers

Other projects

AutoEvals

Choose scorers for me

+ Create custom scorer

Draft

sk 2    + Scorer    Opinion Pieces 9

GPT    Lite    Comparison    Ta

Output    Evaluating... —

Gustav...
KKK co...
is a pot...
people can successfully resist organized racism before it
returns.

...fully reminds us that forgetting
Anaheim's KKK past isn't just an oversight, but a
dangerous erasure of lessons crucial for confronting
present-day hate.

Output    Evaluating...    Output

Loop agent

MCP ~    Text output ~          [[inputstory]]

                                 + Message part ~

+ Message    ⊙ Tools    (x) Mustache ~    ⚡ MCP    T Text output ~

                                 ⚲ Save prompt                    Draft ✎

Display    ⤓    + Row                                    Opinion Pieces 9 ~    ⚙

                          GPT-5                          Comparison        🏷 Ta

o, Anaheim recalled its KKK    ↗ Output                  ↗ Output              Evaluating...  —
t we remember?                 Gustavo A                 This article powerfully reminds us that forgetting
im recalled its KKK city       KKK counci                Anaheim's KKK past isn't just an oversight, but a
remember?                      is a potent call to remember and teach how ordinary    dangerous erasure of lessons crucial for confronting
n Anaheim in 1924.             people can successfully resist organized racism before it    present-day hate.
n Anaheim in 1924. (Anaheim    returns.

ER 12, 2019 -- Gustavo
  the Los Angeles Times.

                               **Loop agent** ⊕

first partner pushes to        ↗ Output                  Evaluating...    ↗ Output
  and tech bros thunder forward                                                    Loop agent ⊕

**Loop agent**

**Turn insights into improvements with ⊘ Loop**

Start a conversation to optimize prompts, generate dataset rows, and automate eval development.

- ⊙ Optimize prompts
- % Create code scorer
- ⊟ Generate dataset rows
- % Choose scorers

Build me a scorer that blah blah blah

Claude 4.6 Sonnet ⌄

# Time for the tech!

Pydantic, Pydantic AI, Braintrust SDK

Pydantic AI

🔍 Search

pydantic/pydantic-ai
v1.40.0  ⭐ 14.2k  🍴 1.5k

# Pydantic AI

*GenAI Agent Framework, the Pydantic way*

CI failing  Coverage  pypi v1.40.0  python 3.10 | 3.11 | 3.12 | 3.13  license MIT  Slack Join Slack

*Pydantic AI is a Python agent framework designed to help you quickly, confidently, and painlessly build production grade applications and workflows with Generative AI.*

FastAPI revolutionized web development by offering an innovative and ergonomic design, built on the foundation of Pydantic Validation and modern Python features like type hints.

Yet despite virtually every Python agent framework and LLM library using Pydantic Validation, when we began to use LLMs in Pydantic Logfire, we couldn't find anything that gave us the same feeling.

We built Pydantic AI with one simple aim: to bring that FastAPI feeling to GenAI app and agent development.

## Why use Pydantic AI

1. **Built by the Pydantic Team**: Pydantic Validation is the validation layer of the OpenAI SDK, the Google ADK, the Anthropic SDK, LangChain, LlamaIndex, AutoGPT, Transformers, CrewAI, Instructor and many more. *Why use the derivative when you can go straight to the source?* 😀

Pydantic AI

3.10.13 (Python 3.10.13)

```python
@agent.tool_plain
async def perplexity_research(query: str):
    """Do research on a topic using Perplexity.

    Args:
        query: what to search for
    """
    search = client.search.create(query=query, max_results=10)
    return search.results


# Define the instructions for the agent

instructions = f"""
You are an investigative journalism assistant that triages news tips.

Your task is to:
1. Evaluate the tip for credibility and newsworthiness (score 1-10)
2. Determine if it is valid (score >= 5) or spam/invalid (score < 5)
3. Categorize it into the most appropriate beat from this list: local news, crime, politics, business, health, technology, environment, education, culture, sports, interr
4. Write a concise 2-3 sentence summary
5. Use Perplexity and/or Tavily search to conduct background research on the topic, entities, and location mentioned
6. Research the tipster's name and email to assess credibility (look for public profiles, previous journalism involvement, potential biases)
7. Generate 3-5 specific follow-up questions for further investigation
8. Compose a professional email response:
   - If valid: Thank them, indicate you're reviewing it, and ask the follow-up questions
   - If invalid: Politely decline and thank them for their interest

Be thorough in your research but efficient. Focus on verifiable facts and credible sources. Understand that people don't really know how to write tips, and you should fol
"""

# Define the output from the agent
class ProcessedTip(BaseModel):
    is_valid: bool
    credibility: int
    newsworthiness: int
    validity: int
    beat: str
    summary: str
    research: str
    follow_up_questions: list[str]
```

# Instrumentation vocabulary

- **Logging:** Tracking what happens
- **Tool/tool call:** Using a service - web search, mail, etc
- **Trace:** Everything you need to reproduce interaction
  - All of the user's inputs
  - All of the inputs/outputs to the LLM
- **Span:** A single step in the process
  - A tool call, the response, the time, the cost
- **Thread:** A group of traces, if multi-turn (ignore this??)
- **Observability:** Using instrumentation to track and analyze
- **Telemetry:** A more generic word for instrumentation

You are an investigative journalism assistant that triages news tips.

Your task is to:
1. Evaluate the tip for credibility and newsworthiness (score 1-10)
2. Determine if it is valid (score >= 5) or spam/invalid (score < 5)
3. Categorize it into the most appropriate beat from this list:
4. Write a concise 2-3 sentence summary
5. Use Perplexity search to conduct background research on the topic, entities, and location mentioned
6. Research the tipster's name and email to assess credibility (look for public profiles, previous journalism involvement, potential biases)
7. Generate 3-5 specific follow-up questions for further investigation
8. Compose a professional email response:
   - If valid: Thank them, indicate you're reviewing it, and ask the follow-up questions
   - If invalid: Politely decline and thank them for their interest

Be thorough in your research but efficient. Focus on verifiable facts and credible sources. Understand that people don't really know how to write tips, and you should follow up on most of the ones that are not obviously crazy people.

triage_tip
() 30.04s, 180484 tok

A\ agent run
30.03s, 180484 tok, Claude 4.5 Haiku

A\ chat claude-haiku-4-5
1.61s, 1677 tok, Claude 4.5 Haiku

running tools
1.44s

running tool: perplexity_research
0.41s

running tool: perplexity_research
0.55s

running tool: tavily_search
0.42s

running tool: tavily_search
0.46s

A\ chat claude-haiku-4-5
2.41s, 25606 tok, Claude 4.5 Haiku

running tools
0.85s

running tool: tavily_search
0.85s

running tool: perplexity_research
0.38s

running tool: tavily_search
0.36s

```python
@braintrust.traced
def create_powerpoint(path: str, source: str) -> str:
    prompt = f"""
    Create a powerpoint presentation named {path} using the content in {source}. Use the templates available in the MCP.
    For each slide, perform a tavily search for additional information needed to make the presentation complete.
    Do not make numeric claims like "improve productivity by 40%" without citations
    Do not include placeholder text like "Insert Image Here" or "Point A" or "Additional information"
    Include citations in the notes section of each slide.
    Maintain the order and structure of the slides in the source file.
    Slides should be as detailed as possible.
    """
    result = agent.run_sync(prompt,
                            instructions=instructions,)
    return result.output


path = "demo.pptx"
source = "vague-mcp-notes.txt"
result = create_powerpoint(path, source)
print(result)
```

# It's honestly not that complicated

- You have a normal LLM setup
- You give it tools
- You say use the tools until you're done
- Why is everyone so excited about
- *Because humans are so so so*

**ONE prompt**

# We get a tip. What's next?

- *IMMEDIATELY PUBLISH IT, IN FULL.*
- Judge whether it's newsworthy
- See whether it's already been covered
- Research the topic in general with preliminary research
- Research the person who submitted the tip (social, etc)
- Assign it to a beat/specific reporter
- Write a follow-up email, send it

# Multi-step process

A reader submits a tip

Is it good?

no → Ignore it

yes → Do some research on the topic

Do we need more information?

no → Send it to the reporter

yes → Generate follow-up questions → Email the reader

# AI-able steps

A reader submits a tip

Is it good?

no → Ignore it

yes → Do some research on the topic

Do we need more information?

no → Send it to the reporter

yes → Generate follow-up questions → Email the reader

I am a journalist in Durham, NC. I received the following tip. Is it worth pursuing?

From: Concerned Citizen (DO NOT PUBLISH)
Email: concerned.durhamite.north@protonmail.com
Tip: I am writing to you today regarding the absolute state of the traffic on Roxboro Street which has been a nightmare since they started talking about this two-way conversion nonsense. I was sitting at the light for twenty minutes yesterday, just watching the construction crews standing around, and I swear I saw the foreman handing a thick envelope of cash to a zoning board member who recused h due to a "conflict of interest." They w construction trailer near the library. I over the lens, but I know what I saw. It's just like when they built the highway through Hayti, money changing hands and nobody listening to the people. Also, the pothole on Geer St is getting bigger.

Short answer: yes

tGPT 5.2 ⌄

Short answer: **yes, it's worth a** *light* **pursuit**—but not as a corrupti
**public process and traffic impacts**, not an allegation of bribery un
facts.

Here's how I'd break it down.

---

## What in this tip has journalistic value?

### 1. Real, verifiable public issue

- Traffic disruption on Roxboro Street tied to a proposed
  documentable civic concern in Durham.
- Residents stuck at lights, construction delays, and poo
  *without* wrongdoing.

### 2. Named process + recent event

- The tip references:
  - A zoning board member
  - A recent recusal for "conflict of interest"
- That gives you a **paper trail**: agendas, minutes, disclosure forms, ethics filings.

### 3. Pattern framing (careful, but usable)

- The comparison to Hayti redevelopment (e.g., **Hayti**) signals a *community memo*
  rooted in historical displacement.
- That's useful context even if the cash claim goes nowhere.

↓

## What are the red flags?

Ask anything

---

✦

# What are the red flags?

## 🚩 The alleged bribery

- Anonymous source
- No usable photo
- "I know what I saw" language
- No corroboration
- Public accusation against a specific (even unnamed) official

This is **not publishable as an allegation** on its own and could
mishandled.

## 🚩 Mixed credibility signals

- The pothole on Geer Street being "getting bigger" is true bu
  venting rather than focusing.
- Emotional ton

---

## How I'd pursue it (low-risk, high-return)

### Step 1: Reframe the story

Instead of:

> "Possible bribery in Roxboro St project"

Start with:

St traffic delays amid two-way conversion debate"

### Step 3: On-the-ground reporting

If you have time:

- Visit the construction site near the library
- Observe crew activity and traffic flow
- Ask neutral questions like:

> "Residents say crews appear idle—what's causing the delays?"

You're gathering **observable facts**, not chasing a claim.

---

## Bottom line

- ✅ **Yes, pursue it** as a *civic infrastructure and transparency* story.
- ❌ **No, do not publish or imply bribery** based on this tip alone.
- 🕐 Let documents, timelines, and multiple sources lead you—*not* the tipster's conclusion.

If you want, I can help you:

- Draft a records request
- Map out who to interview first
- Pressure-test whether the recusal itself is newsworthy

cs angle

om last week
y
rds
ted to the project
I show up **on paper before it ever shows up in an envelope.**

- ❌ **No, do not publish or imply bribery** based on this tip alone.
- 🧭 Let documents, timelines, and multiple sources lead you—*not* the tipster's conclusion.

If you want, I can help you:

- Draft a records request
- Map out who to interview first
- Pressure-test whether the recusal itself is newsworthy

That's usually where the real story hides.

What are some follow up questions I could ask the submitt[...]
brief research has come up with this research:

Recent news about Roxboro Street in Durham, NC includes s[...]
crime and transportation-planning updates, but there do not [...]
to be major new incidents reported there in the last few days [...]
specifically.

Recent crime incidents
In late November 2025, Durham police investigated an early-mo[...]
shooting at a gas station on North Roxboro Street around 3:15 a.[...]
where officers found shell casings on the ground; authorities were [...]
still working to learn whether anyone was hurt and if a suspect had
been arrested.

Also in late November 2025, a man was critically injured in what

## 1. Lock down *time, place, and context* (most important)

These questions test whether the sighting is anchored in reality.

- What date and approximate time did you observe this? (day of week + time window)
- Was this on North or South Roxboro Street?
- Which library was nearby? (Main Library vs. branch)
- Where exactly were you stopped in traffic when you saw this?
- How long did you observe the interaction—seconds or minutes?

If they can't answer these cleanly, the tip weakens significantly.

## 2. Clarify *who* they believe they saw (without naming)

Avoid asking for accusations; ask for descriptions.

- You mentioned the person looked like a zoning board member—what makes you say that? (build, hair, clothing, prior encounters, media photos)
- Had you ever seen this person in real life before, or only in photos?
- Was anyone else present or nearby who might also have seen this?
- Did the individual arrive or leave in a vehicle? If so, can you describe it?

This helps distinguish *recognition* from *assumption*.

Ask anything

Write me the email

Here's a **neutral, legally safe follow-up email** that keeps the source engaged, signals professionalism, and avoids validating any allegation. You can adjust tone slightly depending on your outlet's style.

Email

**Subject** Follow-up on Roxboro Street tip

Hello,

Thank you for taking the time to write and for sharing your concerns about traffic conditions and construction activity on Roxboro Street. I'm doing some preliminary reporting to better understand what residents are experiencing and to verify timelines and processes around the project.

To help me assess what you observed, could you clarify a few details?

- About when did this occur (date and approximate time)?
- Where exactly were you located, and where were the individuals you mentioned standing?
- Was this on North or South Roxboro Street, and which nearby library or landmark was closest?
- How long did you observe the interaction?
- You mentioned the person resembled a zoning board member—had you seen this individual in person before, or was your recognition based on photos or coverage?

You also noted that you took a photo, though it may be obstructed. Even if it doesn't clearly show the interaction, it could still help confirm the location and timing. If you're comfortable sharing it confidentially, I'd be happy to take a look.

Follow-up on Roxboro Street tip

To                                                                                                 Cc Bcc

Follow-up on Roxboro Street tip

Hello,

Thank you for taking the time to write and for sharing your concerns about traffic conditions and construction activity on Roxboro Street. I'm doing some preliminary reporting to better understand what residents are experiencing and to verify timelines and processes around the project.

To help me assess what you observed, could you clarify a few details?

- About when did this occur (date and approximate time)?
- Where exactly were you located, and where were the individuals you mentioned standing?
- Was this on North or South Roxboro Street, and which nearby library or landmark was closest?
- How long did you observe the interaction?
- You mentioned the person resembled a zoning board member—had you seen this individual in person before, or was your recognition based on photos or coverage?

You also noted that you took a photo, though it may be obstructed. Even if it doesn't clearly show the interaction, it could still help confirm the location and timing. If you're comfortable sharing it confidentially, I'd be happy to take a look.

Separately, I'm also interested in hearing more about how the construction and traffic conditions have affected you or others who travel this corridor regularly.

I appreciate you reaching out and understand your request not to be identified. Any information you share will be handled carefully as I continue to verify the facts.

Best,
[Your Name]

Gmail
Outlook
Default email app

# LLM Pipelines/Automated workflows

Let the humans take a break, use the AI to fill those gaps

chatgpt.com/?temporary-chat=true

ChatGPT 5.2 ⌄

◌ Temporary Chat   •••

what are some dog shelters in Arizona?

web_search("Arizona dog shelters")

Here are **some reputable dog shelters and rescue organizations in Arizona** where you can adopt, foster, volunteer, or support dogs in need:

🐾 **...o Area**

- **Arizona Humane Society Papago Park Campus** – Large shelter with adoptions, medical care, and supportive programs.  `azhumane.org`
- **Arizona Humane Society South Mountain Campus** – Another major campus of the Arizona Humane Society in Phoenix.  `azhumane.org`
- **AAWL Adoption Center (Arizona Animal Welfare League Main Shelter)** – No-kill shelter offering adoption, training resources, and community support.  `Arizon...`
- **Maricopa County Animal Care and Control - West Valley Shelter** / **Maricopa County Animal Care and Control** – Public shelters in Phoenix & Mesa with many dogs available for adoption.  `Maricopa C...`
- **Pima Animal Care Center** – Large animal shelter in Tucson where dogs are available for adoption.
  `Pima C...`

# Tools in LLMs

- Allow an LLM to do *stuff* that isn't text generation
- web_search("Whatever you want to look for")
- generate_image("whatever you want an image about")
- other_things("I don't know, it's infinite??")

**Step 1: Reframe the story**

Instead of:

> "Possible bribery in Roxboro St project"

Start with:

> "Residents frustrated by Roxboro St traffic delays amid two-way conversion deba

This gives you a reason to:

- Talk to DOT / city transportation staff
- Interview commuters and nearby businesses
- Ask about timelines, staffing, and delays

**"ok, do that"**

**Step 2: Quietly check the ethics angle**

Without publishing or accusing:

- Pull zoning board minutes from last week
- Identify who recused and why
- Check:
  - Campaign finance records
  - Business affiliations
  - Any contracts connected to the project

If there *is* something real, it will show up **on paper before it ever shows up in an envelope.**

# Agentic AI workflows

Let's REALLY give humans take a break and hand over control to AI

# One-step vs multi-step pipeline

- **Find me dog shelters in Arizona**
  - web_search("dog shelters Arizona")
  - OK sounds good here are the results
- **Find me contact info for every dog shelter in Arizona**
  - web_search("dog shelters Arizona")
  - Visit website 1
  - web_search("Happy Paws Rescue email address")
  - Visit website 2
  - …etc…etc…etc…etc…………
- Agentic AI: **uncertain when it's done**

The user talks to the LLM through the agent, which makes decisions (using the LLM) about what to do next.

You are an investigative journalism assistant that triages news tips.

Your task is to:
1. Evaluate the tip for credibility and newsworthiness (score 1-10)
2. Determine if it is valid (score >= 5) or spam/invalid (score < 5)
3. Categorize it into the most appropriate beat from this list:
4. Write a concise 2-3 sentence summary
5. Use Perplexity search to conduct background research on the topic, entities, and location mentioned
6. Research the tipster's name and email to assess credibility (look for public profiles, previous journalism involvement, potential biases)
7. Generate 3-5 specific follow-up questions for further investigation
8. Compose a professional email response:
   - If valid: Thank them, indicate you're reviewing it, and ask the follow-up questions
   - If invalid: Politely decline and thank them for their interest

Be thorough in your research but efficient. Focus on verifiable facts and credible sources. Understand that people don't really know how to write tips, and you should follow up on most of the ones that are not obviously crazy people.

triage_tip
() 30.04s, 180484 tok

A\ agent run
30.03s, 180484 tok, Claude 4.5 Haiku

A\ chat claude-haiku-4-5
1.61s, 1677 tok, Claude 4.5 Haiku

running tools
1.44s

running tool: perplexity_research
0.41s

running tool: perplexity_research
0.55s

running tool: tavily_search
0.42s

running tool: tavily_search
0.46s

A\ chat claude-haiku-4-5
2.41s, 25606 tok, Claude 4.5 Haiku

running tools
0.85s

running tool: tavily_search
0.85s

running tool: perplexity_research
0.38s

running tool: tavily_search
0.36s

```python
@braintrust.traced
def create_powerpoint(path: str, source: str) -> str:
    prompt = f"""
    Create a powerpoint presentation named {path} using the content in {source}. Use the templates available in the MCP.
    For each slide, perform a tavily search for additional information needed to make the presentation complete.
    Do not make numeric claims like "improve productivity by 40%" without citations
    Do not include placeholder text like "Insert Image Here" or "Point A" or "Additional information"
    Include citations in the notes section of each slide.
    Maintain the order and structure of the slides in the source file.
    Slides should be as detailed as possible.
    """
    result = agent.run_sync(prompt,
                            instructions=instructions,)
    return result.output


path = "demo.pptx"
source = "vague-mcp-notes.txt"
result = create_powerpoint(path, source)
print(result)
```

# It's honestly not that complicated

- You have a normal LLM setup

- You give it tools

- You say use the tools until you're done

- Why is everyone so excited abou~~~

- *Because humans are* <u>*so so so*</u>

**ONE prompt**

MwO0R2kySK1sl

Star 167,629

n8n AI Beta  Ask Build

v 5  + Add tag

Publish  Save

Editor  Executions  Evaluations

Make me a workflow that blah blah blah blah

50/50 onthly credits left ⓘ  Get more

Daily AI news digest

Lead qualification and call scheduling

Multi-agent research workflow  RAG knowledge assistant

YouTube video chapters  Invoice processing pipeline

Summarize emails with AI  Daily weather report

Pydantic AI

△ Pydantic AI          ⚙ 🔍 Search          pydantic/pydantic-ai
                                            ◇ v1.40.0  ☆ 14.2k  ⑂ 1.5k

# Pydantic AI

*GenAI Agent Framework, the Pydantic way*

CI failing  Coverage  pypi v1.40.0  python 3.10 | 3.11 | 3.12 | 3.13  license MIT  Slack Join Slack

*Pydantic AI is a Python agent framework designed to help you quickly, confidently, and painlessly build production grade applications and workflows with Generative AI.*

FastAPI revolutionized web development by offering an innovative and ergonomic design, built on the foundation of Pydantic Validation and modern Python features like type hints.

Yet despite virtually every Python agent framework and LLM library using Pydantic Validation, when we began to use LLMs in Pydantic Logfire, we couldn't find anything that gave us the same feeling.

We built Pydantic AI with one simple aim: to bring that FastAPI feeling to GenAI app and agent development.

## Why use Pydantic AI

1. **Built by the Pydantic Team**: Pydantic Validation is the validation layer of the OpenAI SDK, the Google ADK, the Anthropic SDK, LangChain, LlamaIndex, AutoGPT, Transformers, CrewAI, Instructor and many more. *Why use the derivative when you can go straight to the source?* 😀

```python
@agent.tool_plain
async def perplexity_research(query: str):
    """Do research on a topic using Perplexity.

    Args:
        query: what to search for
    """
    search = client.search.create(query=query, max_results=10)
    return search.results


# Define the instructions for the agent

instructions = f"""
You are an investigative journalism assistant that triages news tips.

Your task is to:
1. Evaluate the tip for credibility and newsworthiness (score 1-10)
2. Determine if it is valid (score >= 5) or spam/invalid (score < 5)
3. Categorize it into the most appropriate beat from this list: local news, crime, politics, business, health, technology, environment, education, culture, sports, intern
4. Write a concise 2-3 sentence summary
5. Use Perplexity and/or Tavily search to conduct background research on the topic, entities, and location mentioned
6. Research the tipster's name and email to assess credibility (look for public profiles, previous journalism involvement, potential biases)
7. Generate 3-5 specific follow-up questions for further investigation
8. Compose a professional email response:
   - If valid: Thank them, indicate you're reviewing it, and ask the follow-up questions
   - If invalid: Politely decline and thank them for their interest

Be thorough in your research but efficient. Focus on verifiable facts and credible sources. Understand that people don't really know how to write tips, and you should fol
"""

# Define the output from the agent
class ProcessedTip(BaseModel):
    is_valid: bool
    credibility: int
    newsworthiness: int
    validity: int
    beat: str
    summary: str
    research: str
    follow_up_questions: list[str]
```

# But sometimes things go wrong

The further humans are, the more nervous you should be getting

CNN Business

Watch    Subscribe

Markets →

| | | |
|---|---|---|
| DOW | 49,266.11 | 0.55% ▲ |
| S&P 500 | 6,921.46 | 0.01% ▲ |
| NASDAQ | 23,480.02 | 0.44% ▼ |

Fear & Greed Index →

47

Business / Media

# 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

Events   Email Briefings                                                          SIGN IN

8:54 AM  FRIDAY JANUARY 9, 2026

| D.C. | BXL | LAGOS | RIYADH | BEIJING | SG |
|---|---|---|---|---|---|

# SEMAFOR

*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                    SHARE ⬏

# Evaluation

- Every project is magic and special and personal
- "Does this AI model hallucinate?" is not good enough
- You want to trace **everything that's happening**
- You want to judge whether the outputs are good enough

Processed tips (filled in)

File   Edit   View   Insert   Format   Data   Tools   Extensions   Help   Gemini

Share

D50 | I am writing to you today regarding the absolute state of the traffic on Roxboro Street which has been a nightmare since they started talking about this two-way conversion nonsense. I was sitting at the light for twenty minutes yesterday, just

Summarize this data

| D | | | E | F | G | H | I | J | L |
|---|---|---|---|---|---|---|---|---|---|
| tip | | | itySt | newsworthy | isValid | beat | followUpQuestions | summary | tipResearch |

**D (tip):** I am writing to you today regarding the absolute state of the traffic on Roxboro Street which has been a nightmare since they started talking about this two-way conversion nonsense. I was sitting at the light for twenty minutes yesterday, just watching the construction crews standing around, and I swear I saw the foreman handing a thick envelope of cash to a man who looks exactly like that zoning board member who recused himself from the vote last week due to a "conflict of interest." They were standing right in front of the construction trailer near the library. I took a picture but my thumb was over the lens, but I know what I saw. It's just like when they built the highway through Hayti, money changing hands and nobody listening to the people. Also, the pothole on Geer St is getting bigger.

**E:** 6  **F:** TRUE  **G:** TRUE  **H:** Local News

**I (followUpQuestions):** - Please provide the exact date, time, and precise location of the observed handoff (e.g., construction trailer area near the Roxboro/Mangum project by the library, including any nearby cross streets or landmarks).
- Do you recall the name or official title of the zoning-board member you mentioned? If possible, share any public identifiers (such as committee name, district, or seat) where the person...
- Do you have any additional evidence beyond memory (photo, video, notes, or social posts)? If you have a photo, can you share it or describe any visible details (even if the face is unclear) and the timestamp.
- Were there any other witnesses present who can corroborate what you saw? If yes, please provide their contact information or ask them to contact us...
- Have you seen any official recusal or conflict-of-interest disclosures or minutes mentioning this board member around the time you observed? If you have dates, links, or document names, please share.
- Can you confirm that the tip refers specifically to the Roxboro Street and Mangum Street two-way conversion project, and not another project? If you have related documents or dates in mind, share them.
- Regarding the Geer Street pothole note, can you specify the exact block or cross-streets and describe the pothole's size or condition, plus any city ticket or service-request numbers if you recall?

**J (summary):** Tip alleges a foreman handed a "thick envelope of cash" to a zoning-board member who recused himself for a supposed conflict of interest, observed near the Roxboro/Mangum street-conversion construction site. Our review finds the project is publicly documented and in a multi-year design/public-engagement process with milestones (design began 2025; 25% design anticipated late 2025/early 2026; construction window 2027–2028). There is no public corroboration in minutes or credible reports of bribery or a recent crime; the tip is therefore unverified as presented; pursue public-records requests for meeting minutes, conflict disclosures, and related procurement/payments, plus interviews with city staff/consultants, and verify the Geer Street pothole status. Report with clear distinctions between verified facts and unverified allegations.

**L (tipResearch):** Here's a practical, journalist-facing breakdown of what you can pursue the... without getting tangled in unveri... separated what... confirmed in public records from what... and I've in... concrete, date-stamped leads you ca...

What's publicly verifiable about Roxboro Street / Mangu... conversion topic)
- The project exists... is being actively studied/desig... Durham in partne... with NCDOT.
- The City says... sign for the two-way conversion beg... review of 25%... cept plans was anticipated in fall 202... Roxboro... Mangum St. as part of a plan to conve... two-way... ([durhamnc.gov]...https://www.durhamnc.gov/5203/Roxb... utm_source=openai))
- The City's design pa... Roxboro St. and Mangu... notes a multi-phase p... with Phase I running Apri... 2026 and Phase II... gagement around fall 2025, targe... design mile... a construction window later (pot... depending on... g and approvals). ([durhamnc.gov]...ps://www.durhamnc.gov/5405/Roxb... Design?utm_source=openai))
- Public engagement has taken place and continues.
- The City posted notices and hosted public events (th... referenced in coverage of October 2025) to show prelim... feedback. ([durhamnc.gov](https://www.durhamnc.gov/m/NewsFla... m_source=openai))
- An independent outlet (IndyWeek)... ered the Octo... and noted general timelines (const... on expected 202... the near term to move the design... ward). ([indyweek.com](https://indy... com/news/roxboro-m... hape-up-in-durham/?utm_source=openai))
- The latest public-facing... dule, funding, and desig... multi-year process rath... an a completed project.
- A 2025 Durham p... summarizing the project repea... anticipated... window and notes the project i... City's CIP, wit... sign team (Stantec + Just Cities ... ([durhamnc.gov](... ps://www.durhamnc.gov/5203/Roxb... utm_source=openai))
- A separate design page (Phase II) reiterates the eng... the 25% design milestone... with comments accepted th...

+ | live | Most recent test run | golden set | Most recent test run (agentic self-hosted) | Most recent test run (agentic)

triage_tip 30.04s, 180484 tok

agent run 30.03s, 180484 tok, Claude ...

chat cla... 1.61s, 1677 tok, Claude 4.5 H...

**Tool span**

running tool: perplexity research

+ Dataset ...

chat cla... 1.61s, 1677 tok, Claude 4.5 H...

running tools 1.44s

running tool: perplexity_research 0.41s

running tool: perplexity_research 0.55s

running tool: tavily_search 0.42s

running tool: tavily_search 0.46s

chat cla... 2.41s, 25606 tok, Claude 4.5...

running tools 0.85s

running tool: tavily_search 0.85s

running tool: perplexity_research 0.38s

running tool: tavily_search 0.36s

chat cla... 2.12s, 48605 tok, Claude 4.5...

running tools 0.37s

running tool: tavily_search 0.37s

chat cla... 3.72s, 50325 tok, Claude 4.5...

running tools 0.8s

running tool: tavily_search 0.8s

chat cla... 16.69s, 54271 tok, Claude 4....

---

**triage_tip** 30.04s, 180484 tok

**agent run** 30.03s, 180484 tok, Claude ...

**Tool span**

running tool: tavily_search

**Tool span**

running tool: tavily_search

+ Dataset ... Flag for review

Metrics

Start          Duration        Offset
23m ago        0.7972s         12.5396s

↗ Input

Pretty

query          Mitch Colvin embezzlement criminal investigation audit 2025

topic          news

↓ Output

Pretty

0 {4}

content        Image 23: Protest erupts in Baltimore following U.S. strikes on Venezuela ### Protest erupts in Baltimore following U.S. strikes on Venezuela 04:43 10h agoImage 24: Tips for paying off that holiday debt and managing your finances in 2026 ### Tips for paying off that holiday debt and managing your finances in 2026 00:54 12h agoImage 25: Fatal shooting sparks first homicide investigation of 2026 in Baltimore City ### Fatal shooting sparks first homicide investigation of 2026 in Baltimore City 01:54 12h agoImage 26: Ravens flock to Pittsburgh for high stakes season finale ### Ravens flock to Pittsburgh for high stakes season finale 01:28 14h agoImage 27 Registration open for "Miles for Maryland" challenge ### Registration open for "Miles for Maryland" challenge 00...

Loop

---

Offset
2.5831s

e North Carolina mayor embezzlement

Sheriff Faces Ouster Bid
Alleged Threats To
maker. By Hayley Fowler
uary 6, 2026, 6:40 PM
To view the full article,
ster now. Try a seven
FREE Trial. ### Recent
les By Hayley, Agita

Loop

braintrust.dev/app/Little%20Columns/p/default-otel-project/playgrounds/LA%20Times%20prompts

Little Columns

**Project**

default-otel-project

- Overview
- Logs
- Monitor
- Review
- **Playgrounds**
- Experiments
- Datasets
- Prompts
- Scorers
- Settings

**Recent projects**

- Classifying News Articles ...
- My Project
- mixed-otel-braintrust-pyth...

default-otel-project › Playgrounds

Upgrade

△ **LA Times prompts** ···

Diff ⬤   + Experiments   ▷ Run

⬤ Base task

⬤ Comparison task

A\ Claude 4.5 Haiku ▾    ⇄ Params ▾

A\ Claude 4.5 Haiku ▾    ⇄ Params ▾

System ▾

System ▾

You are a helpful assistant

You are an editor on the opinion desk of a major US paper.

User ▾

User ▾

▢ Write a response to the piece {{ input.story }}

▢ You are writing a dissenting response to an opinion article for a newsroom feature.

+ Message part ▾

ARTICLE HEADLINE: {{input.headline}}
ARTICLE TEXT: {{input.story}}

+ Message   ⊙ Tools   (») Mustache ▾   ⬚ MCP   T Text output ▾

Write ONE response with these rules:
- Output plain text only (no title, no header, no bullets, no numbered lists).
- 150–200 words.
- No greetings, no thanks, no sign-off, no "here's the response".
No em dashes (—

▢ latimes-response   75447eba ▾            Update...

▢ latimes-response-better   3ea00a5a (latest) ▾

⊘ Default view ▾   ⊜ Filter   ⇄ Display   ⬇   + Row

+ Task 2   ⊹ Scorer 5   ⊟ Dataset 11 9 ▾   ⇄

% LAT: Word Count ... ✕   % LAT: Not Too Per... ✕   % LAT: No Em Dashes ✕   % LAT: LLM: Person... ✕   % LAT: LLM: Respo... ✕   Clear

☐   ⇲ Input                              ⬤ latimes-response                Base   ⬤ latimes-response-better       Comparison   ⊙ Tags

1   headline: Amid national sorrow, Trump stops    ↗ Output      ⊟ Dec 18, 2025  ⏱ 6.6s  ⊞ 398  $0.004   ↗ Output    ⊟ Dec 18, 2025  ⏱ 5s  ⊞ 252  $0.003   —
    pretending it's about the 'worst of the worst'
    story: |-                                      **A Response to Chabria's**                           Chabria conflates rhetoric with policy in ways that obscure
      Amid national sorrow, Trump pretending it's  **Column on Trump's**                                 rather than clarify the actual stakes. Her invocation of
      about the 'worst of the worst'               **Immigration Rhetoric**                              "Great Replacement Theory" and Nazi imagery rests on
      Staff Sgt. Andrew Wolfe and Spec. Sarah Beckstrom                                                  interpreting government language through an ideological
      of the West Virginia National Guard.         Anita Chabria's column raises important concerns about lens rather than examining what deportation policy
      Air Force Staff Sgt. Andrew Wolfe and Army Spec. the rhetoric surrounding immigration policy in the current actually targets. The article does not provide evidence that
      Sarah Beckstrom of the West Virginia National Guard.                                               enforcement actions have focused on citizens or legal
      Beckstrom was killed and Wolfe wounded by a gunman in                                              immigrants, despite her suggestion that "disruptive
      Washington, D.C. (Associated Press)                                                                populations" could encompass them. Existing law already
      Anita Chabria. (Ricardo DeAratanha / Los Angeles
      Times)
      By Anita Chabria                             % LAT: LLM: Person...  100%                           % LAT: LLM: Person...  100%

**Free plan usage** ↗

Logs          0.0024 of 1 GB ▭
Scores        8 of 10,000 ▭

✳ SQL sandbox

Little Columns

default-otel-project › Playgrounds

Upgrade

**Project**

default-otel-project

- Overview
- Logs
- Monitor
- Review
- **Playgrounds**
- Experiments
- Datasets
- Prompts
- Scorers
- Settings

**Recent projects**

- Classifying News Articles ...
- My Project
- mixed-otel-braintrust-pyth...

LA Times prompts ···

Diff ◯ + Experiments ▷ Run

Default view ▽ | ≡ Filter | ≡ Display | ⬇ | + Row

+ Task 2 | + Scorer 5 | ⊟ Dataset 11 9 ▽

% LAT: Word Count ... ✕ | % LAT: Not Too Per... ✕ | % LAT: No Em Dashes ✕ | % LAT: LLM: Person... ✕ | % LAT: LLM: Respo... ✕ | Clear

| | ↘ Input | ● latimes-response | Base | ● latimes-response-better | Comparison | ⬚ Tags |
|---|---|---|---|---|---|---|

**Free plan usage** ↗
Logs 0.0024 of 1 GB
Scores 8 of 10,000

✳ SQL sandbox

---

1 | headline: Amid national sorrow, Trump stops pretending it's about the 'worst of the worst' story: |-
    Amid national sorrow, Trump stops pretending it's about the 'worst of the worst'
    Staff Sgt. Andrew Wolfe and Spec. Sarah Beckstrom of the West Virginia National Guard.
    Air Force Staff Sgt. Andrew Wolfe and Army Spec. Sarah Beckstrom of the West Virginia National Guard. Beckstrom was killed and Wolfe wounded by a gunman in Washington, D.C. (Associated Press)
    Anita Chabria. (Ricardo DeAratanha / Los Angeles Times)
    By Anita Chabria
    Columnist |
    Follow
    Dec. 2, 2025 3 AM PT

1

↗ Output | 📅 Dec 18, 2025 ⏱ 6.6s ⊞ 398 $0.004

**A Response to Chabria's Column on Trump's Immigration Rhetoric**

Anita Chabria's column raises important concerns about the rhetoric surrounding immigration policy in the current

% LAT: LLM: Person... 100%
% LAT: LLM: Respon... 100%
% No Em Dashes 0%
% Not Too Personal 100%
% Word Count 150... 0%

↗ Output | 📅 Dec 18, 2025 ⏱ 5s ⊞ 252 $0.004 | —

Chabria conflates rhetoric with policy in ways that obscure rather than clarify the actual stakes. Her invocation of "Great Replacement Theory" and Nazi imagery rests on interpreting government language through an ideological lens rather than examining what deportation policy actually targets. The article does not provide evidence that enforcement actions have focused on citizens or legal immigrants, despite her suggestion that "disruptive populations" could encompass them. Existing law already

% LAT: LLM: Person... 100%
% LAT: LLM: Respon... 100%
% No Em Dashes 100%
% Not Too Personal 100%
% Word Count 150... 100%

---

2 | headline: These federal judges are building a legal wall against Trump's assault on transgender rights story: |-
    These federal judges are building a legal wall against Trump's assault on transgender rights
    Attorney General Pam Bondi
    Attorney General Pam Bondi, seen here testifying before the Senate Judiciary Committee, is the flag carrier of President Trump's assault on transgender rights, but her campaign has been questioned by federal judges. (Win McNamee / Getty Images)
    Los Angeles Times columnist Michael Hiltzik
    By Michael Hiltzik
    Business Columnist |
    Follow
    Dec. 2, 2025 3 AM PT

1

↗ Output | 📅 Dec 18, 2025 ⏱ 8.9s ⊞ 405 $0.005

**A Response to Hiltzik's "Federal Judges Building a Legal Wall Against Trump's Assault on Transgender Rights"**

% LAT: LLM: Person... 100%
% LAT: LLM: Respon... —
% No Em Dashes 0%
% Not Too Personal 100%
% Word Count 150... 0%

↗ Output | 📅 Dec 18, 2025 ⏱ 8.2s ⊞ 220 $0.004 | —

Hiltzik frames judicial skepticism of DOJ subpoenas as evidence of improper government overreach, but the rulings themselves may reflect narrower concerns about procedure rather than the substance of federal authority. Judge Kearney's decision to quash requests for patient identities rests partly on privacy grounds, a legitimate judicial interest that does not necessarily vindicate the broader claim that Congress forbade federal involvement in medical regulation. Federal agencies routinely

% LAT: LLM: Person... 100%
% LAT: LLM: Respon... 100%
% No Em Dashes 100%
% Not Too Personal 100%
% Word Count 150... 100%

---

3 | headline: In Trump's regime, Catholics are among the

↗ Output | 📅 Dec 18, 2025 ⏱ 11.4s ⊞ 610 $0.005 | ↗ Output | 📅 Dec 18, 2025 ⏱ 4.5s ⊞ 231 $0.003 | —

Little Columns

default-otel-project › Playgrounds

Upgrade

**Project**

default-otel-project

🏠 Overview
📈 Logs
📊 Monitor
☑ Review
⚙ Playgrounds
🗄 Experiments
🗃 Datasets
💬 Prompts
% Scorers
⚙ Settings

🔺 LA Times prompts ···

Diff ⚪ + Experiments ▷ Run

System
You are a helpful assistant

User ⌄
☐ Write a response to the piece {{ input.story }}

+ Message part ⌄

+ Message   ⚙ Tools   (x) Mustache ⌄   ⚘ MCP   T Text output ⌄

💬 latimes-response   75447eba ⌄                    Update...

System
You are an editor on the opinion desk of a major US paper.

User ⌄
☐ You are writing a dissenting response to an opinion article for a newsroom feature.

ARTICLE HEADLINE: {{input.headline}}
ARTICLE TEXT: {{input.story}}

Write ONE response with these rules:
- Output plain text only (no title, no header, no bullets, no numbered lists).
- 150–200 words.
- No greetings, no thanks, no sign-off, no "here's the response".
- No em dashes (—)

💬 latimes-response-better   3ea00a5a (latest) ⌄

**Recent projects**

🕐 Classifying News Articles ...
📄 My Project
🕐 mixed-otel-braintrust-pyth...

⊘ Default view ⌄   Reset   Save as   ▤ Filter   ⚙ Display   ⬇   + Row

+ Task 2   + Scorer 5   🗄 Dataset 11 9 ⌄

% LAT: Word Count ... ✕   % LAT: Not Too Per... ✕   % LAT: No Em Dashes ✕   % LAT: LLM: Person... ✕   % LAT: LLM: Respo... ✕   Clear

% LAT: LLM: Personal Language Check
**80%** AVG                              ↘ 3
**100%** -20%                           ↘ 3

% LAT: LLM: Response Quality
**90%** AVG                        ↗1 ↘2
**95.56%** -6%                     ↗1 ↘2

% No Em Dashes
**0%** AVG                              ↘ 9
**100%** -100%                         ↘ 9

% Not Too Personal
**100%** AVG
**100%**

% Word Count 150-200
**0%** AVG                              ↘ 9
**100%** -100%                         ↘ 9

⏱ Duration
**88.75s** SUM                          ↘ 9
**49.07s** +39.68s (+0.45)             ↘ 9

⏱ LLM duration
**88.59s** SUM                          ↘ 9
**48.87s** +39.71s (+0.45)             ↘ 9

Prompt tokens
—
—

Completion tokens
—
—

**Free plan usage**
Logs      0.0024 of 1 GB
Scores    8 of 10,000

✳ SQL sandbox

braintrust.dev/app/Little%20Columns/p/default-otel-project/scorers/78f40d8b-df97-4788-85f0-0daf00d552ec

Little Columns

default-otel-project › Scorers

Upgrade

**Project**

default-otel-project

- Overview
- Logs
- Monitor
- Review
- Playgrounds
- Experiments
- Datasets
- Prompts
- **Scorers**
- Settings

**Recent projects**

- Classifying News Articles ...
- My Project
- mixed-otel-braintrust-pyth...

**Free plan usage**

| Logs | 0.0024 of 1 GB | |
| Scores | 8 of 10,000 | |

✳ SQL sandbox

% LAT: Word Count 150-200 ⋯

🏷 ⎘ 🗑 📻 Online scoring rule | Evaluate in ⌄ | Save version

**Name**

LAT: Word Count 150-200

**Slug**

word-count-150-200-536b

🐍 Python

```python
1  from typing import Any
2
3  def handler(input: Any, output: Any, expected: Any, metadata: dict[str,
   Any]):
4      """
5      Scores based on word count targeting 150-200 words.
6      - Under 100: 0.0 (too short)
7      - 100-149: 0.6 (a little short)
8      - 150-200: 1.0 (perfect)
9      - 201-250: 0.6 (a little long)
10     - Over 250: 0.0 (too long)
11     """
12     output_text = str(output) if output else ""
13     word_count = len(output_text.split())
14
15     if word_count < 100:
16         score = 0.0
17         category = "too_short"
18     elif word_count < 150:
19         score = 0.6
20         category = "slightly_short"
21     elif word_count <= 200:
22         score = 1.0
23         category = "perfect"
24     elif word_count <= 250:
25         score = 0.6
```

▷ Run | 🔲 Activity

**Variables source**

| Editor | Dataset | Logs |

Enter `input` object

The assumed input schema for scorers is {`input`, `expected`, `metadata`, `output`}. Learn more

```json
JSON ⌄
1 ⌄ {
2     "input": "",
3     "output": "",
4     "expected": "",
5     "metadata": {}
6 }
```

⊘ Generate

Test

Loop ⊘

braintrust.dev/app/Little%20Columns/p/default-otel-project/scorers/8cab965e-e543-4f9a-afb6-4d03b76a9089

125% — +  Reset

Little Columns ⌄

Project

**default-otel-project** ⌄

🏠 Overview
〰️ Logs
📊 Monitor
✅ Review
🎛 Playgrounds
📦 Experiments
🗄 Datasets
💬 Prompts
**% Scorers**
⚙️ Settings

Recent projects

🕐 Classifying News Articles …
🕐 My Project
🕐 mixed-otel-braintrust-pyth…

Free plan usage ↗
Logs          0.0024 of 1 GB
Scores        8 of 10,000

✳️ SQL sandbox

default-otel-project › Scorers

**% LAT: LLM: Response Quality**  ···

🏷  ⧉  🗑  ((•)) Online scoring rule   **Evaluate in** ⌄   **Save version**

Upgrade  📖  ❓  🔍  ⚙️  👤

**Name**
LAT: LLM: Response Quality

**Slug**
llm-response-quality-9cf0

**Prompt**
GPT-5 mini                                        ⌄   | ⇄ Params ⌄

System ⌄                                              ⧉  +  ⌄

You are grading whether RESPONSE is a good dissenting counterargument to ARTICLE.

Evaluate for:
- Clearly engages 1-2 concrete claims from the article
- Offers a plausible counterargument
- Analytical, fact-focused tone (not melodramatic)
- Civilly written (no rudeness or harassment)
- Grounded in the article (doesn't invent facts not in the text)
- Hedges appropriately when article lacks evidence

ARTICLE:
**{{input}}**

RESPONSE:
**{{output}}**

Grade as:
- **a** if the response clearly engages concrete claims, offers plausible counterarguments, maintains civil analytical tone, and doesn't invent facts
- **b** if the response is on-topic but generic, has weak counterarguments, or makes some vague claims
- **c** if the response is mostly unrelated, rude, hallucinates facts not in article, or is inappropriate

▷ Run      🕐 Activity

**Variables source**

| Editor | Dataset | Logs |

**Enter input object**

The assumed input schema for scorers is {input, expected, metadata, output}. Learn more

JSON ⌄

```
1  {
2    "input": "",
3    "output": "",
4    "expected": "",
5    "metadata": {}
6  }
```

⟲ Generate

Test

Loop ⊗

## Learn more

- **Workflow tools:** n8n, Pydantic AI *(not Pydantic!!)*
- **Observability:** Braintrust, Arize Phoenix, Comet Opik
- **Evaluation/theory:** Hamel Husain, Shreya Shankar
- **Product development:** News Product Alliance
- **Newsroom AI class:** https://littlecolumns.com/courses/ai-newsroom/
- **Conferences:** NICAR! (oh hey we're here???)
- **Email me!** js4571@columbia.edu

# From Vibes to Scores

Build your own AI benchmark!

**Jonathan Soma**
js4571@columbia.edu