# CogCache

# The most cost-effective way to access Azure OpenAI with no capacity limits.

UP TO
**50%**
CHEAPER

UP TO
**100x**
FASTER

COMPLETE
**AI control** AND **alignment**

CogCache is the most cost-effective way to access Azure OpenAI with no capacity limits.

By leveraging a global network of Azure supercomputing clusters, CogCache provides the highest performance and lowest cost on the market for Azure OpenAI tokens, enabling companies to unlock the full potential of generative AI.

## With CogCache, you can now:

- Benefit from the lowest cost on the market for Azure OpenAI tokens
- Achieve the highest performance with the fastest response times, lowest latency and full LLM grounding and alignment
- Access the most advanced LLMs with no capacity limits
- Pay only for what you need, with monthly and annual subscriptions
- Get started in minutes with a simple purchase process

## Reduce Costs and Carbon Footprint

Save up to 50% on LLM costs with CogCache for similar, reoccurring prompts and cut your carbon footprint, making your AI operations more sustainable and cost-effective.

## Asynchronous Alignment Scoring

CogCache uses Dynamic Constitutional AI to analyze and score cache content asynchronously, ensuring alignment and grounding of responses to uphold your brand integrity and comply with governance requirements.

## Full-stack LLM Observability

Gain real-time insights, track performance key metrics and view all the logged requests for easy debugging.

## Fast Cache Retrieval

100x your LLM response times with our two-tiered system uses high-speed in memory hashing and vector search to radically reduce response times and lower token rendering costs.
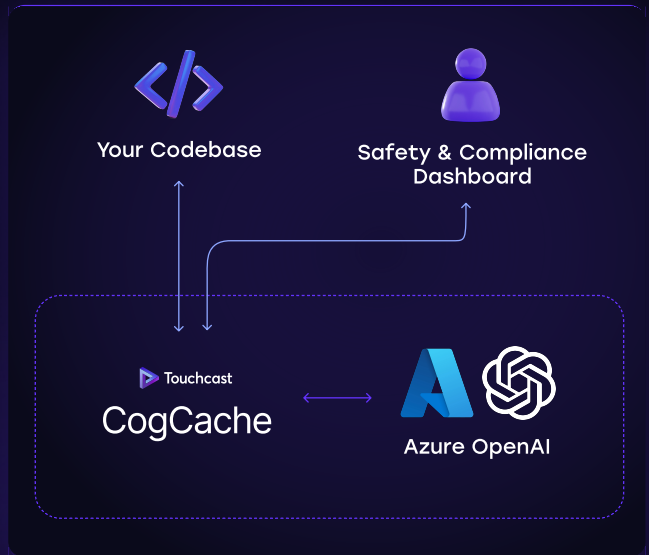
## Self-Healing Cache

CogCache's self-healing mechanism identifies and mitigates misaligned responses, permitting automatic and human-approved adjustments for cache integrity and system reliability.

## Temporal Relevance Tracking

The cognitive cache has the ability to discern if a cached result is temporally relevant and might no longer be applicable based on its content.

Azure OpenAI

# CogCache

## Your CogCache dashboard allows teams to monitor generative flow, set policy, and audit content with confidence.

Your Codebase

Safety & Compliance Dashboard

Touchcast
CogCache

Azure OpenAI

Get started in minutes: Simply point your Azure OpenAI calls to your CogCache instance. No code changes needed.

### Complete transparency and auditability of all generated content

This item was verified by a human.

⊙ 3
3 similar prompts

**Responsible AI › Cache**
All cache items (10)

⬇ Export   + Add cache item(s)

| | 2637 Cache items | 83% Aligned items |

Type to filter...   All applications   All alignment scores
Search through all cache items   Filter by application   Filter by alignment score

| | ITEM | ALIGNMENT | HITS | DATE |
|---|---|---|---|---|
| | Sales POC App — What is the metaverse? | 93 | 2 | 01/02/2023 |
| | Sales POC App — How can we integrate it into our sales processes? | 93 | 1 | 01/02/2023 |
| | RFP Copilot App — Why banks should be more oriented to Gen Z? | 55 | 1 | 01/02/2023 |
| | Sales POC App — What will be the impact of 5G on Fintech and what will be the economic impact on the overall market? | 80 | 1 | 01/02/2023 |
| | Sales POC App — How can we win more business using Gen AI? | 65 | 13 | 01/02/2023 |
| | HR Knowledge Base • RAG — Show me our policies on data retention | 40 | 20 | 01/02/2023 |
| | Sales POC App • RAG — Provide me a list of best practices within Microsoft for increasing brand awareness | 90 | 5 | 01/02/2023 |
| | Employee Support Chatbot — What have you done to avoid bias or discrimination in your responses? | 85 | 2 | 01/02/2023 |

This item was verified by a human.

3 similar prompts

Terms of Use • Privacy Policy • 2024 © Touchcast

### Ability to edit, align and correct cache entries

Use suggested answer

**Responsible AI › Cache › Manage cache item**
What will be the impact of 5G on Fintech and what will be the economic impact on the overall market?

Prompt & answer   Requests (3)   Analytics

**Application**
Customer Support Chatbot

**Main prompt**
What will be the impact of 5G on Fintech and what will be the economic impact on the overall market?

**Similar prompts**
How will 5G affect Fintech and its economic effect on the market?
What are the consequences of 5G on Fintech and its overall economic impact?
Add a similar prompt here...

**Answer**

The 5G technology has the potential to revolutionize the fintech industry, as it will enable faster data transmission and processing times. This will allow fintech companies to offer more efficient and secure services, such as faster payments, digital banking solutions, and investments services. 5G will also enable financial institutions to process large amounts of data in real-time, allowing them to make faster and more informed decisions.

5G will also enable the development of more secure and reliable online platforms for trading and investing, as well as for managing user data and financial transactions. We expect 5 trillion dollars of value to be created next year. For stock investors this will be a great opportunity that they should not miss.

Additionally, the introduction of 5G technology could have a disruptive impact on the traditional banking, but it is very unlikely to kill that way of doing business completely. The banks might need to dodge this bullet and adapt as fast as possible to the new reality and offer a mix of on-the-go and in-branch services.

Lastly, 5G will allow fintech companies to offer more innovative services, such as predictive analytics, AI, and ML. The CEO, Jon Doe said he is 100% sure this is the future and bets on it.

Reset changes   Save

Terms of Use • Privacy Policy • 2024 © Touchcast

## Always Available
Ensure throughput is always available at low latency with our reserved processing capacity (whether used or not).

## Flexible
Enjoy complete flexibility with no monthly commitments required.

## Cost-effective
Minimize costs by eliminating the need to use tokens on previously generated content.

## Smart
Empower your teams with the ability to monitor generated content, ensuring observability and control over AI processes.

## Easy
No implementation resources required - simply switch your endpoints and get started in minutes.

Touchcast
For more information: touchcast.com/cogcache
hello@touchcast.com   www.touchcast.com