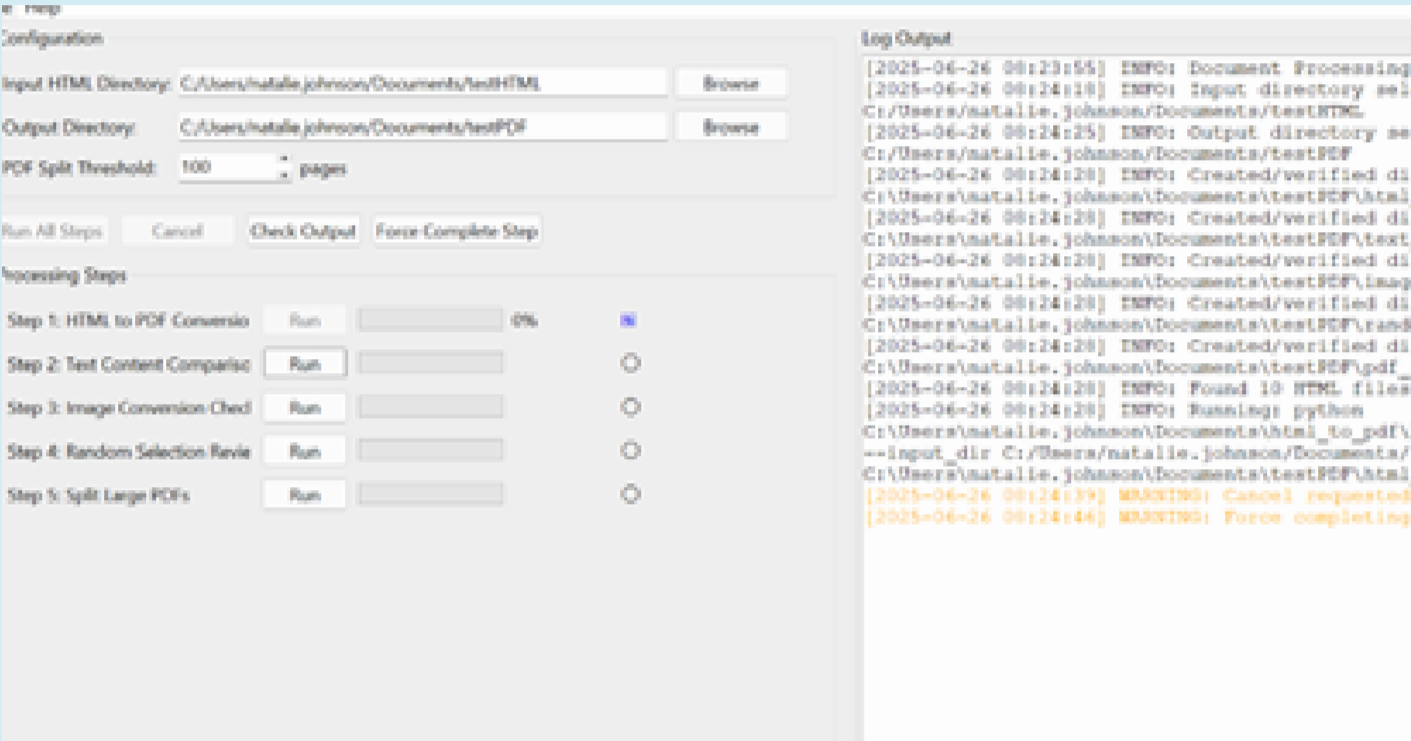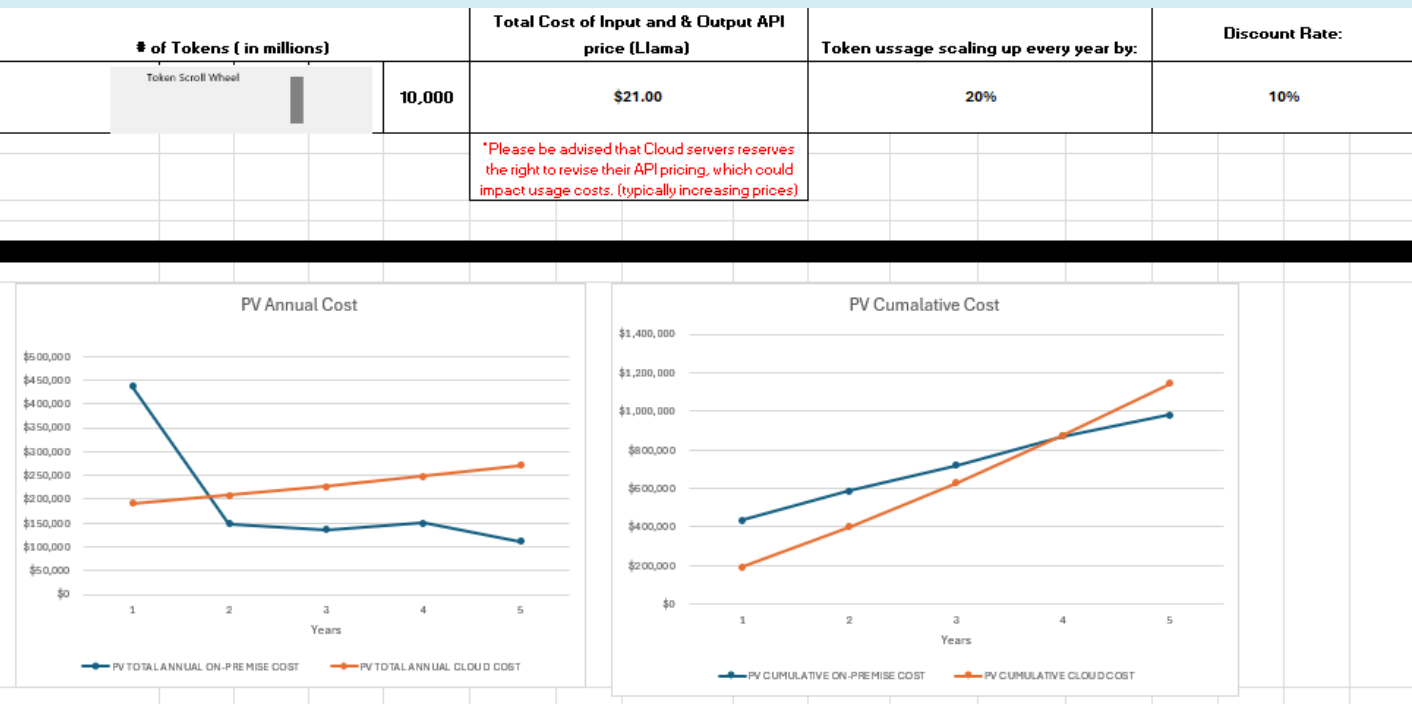# natalie lam johnson .com

## What is

**R**etrieval
**A**ugmented
**G**eneration?

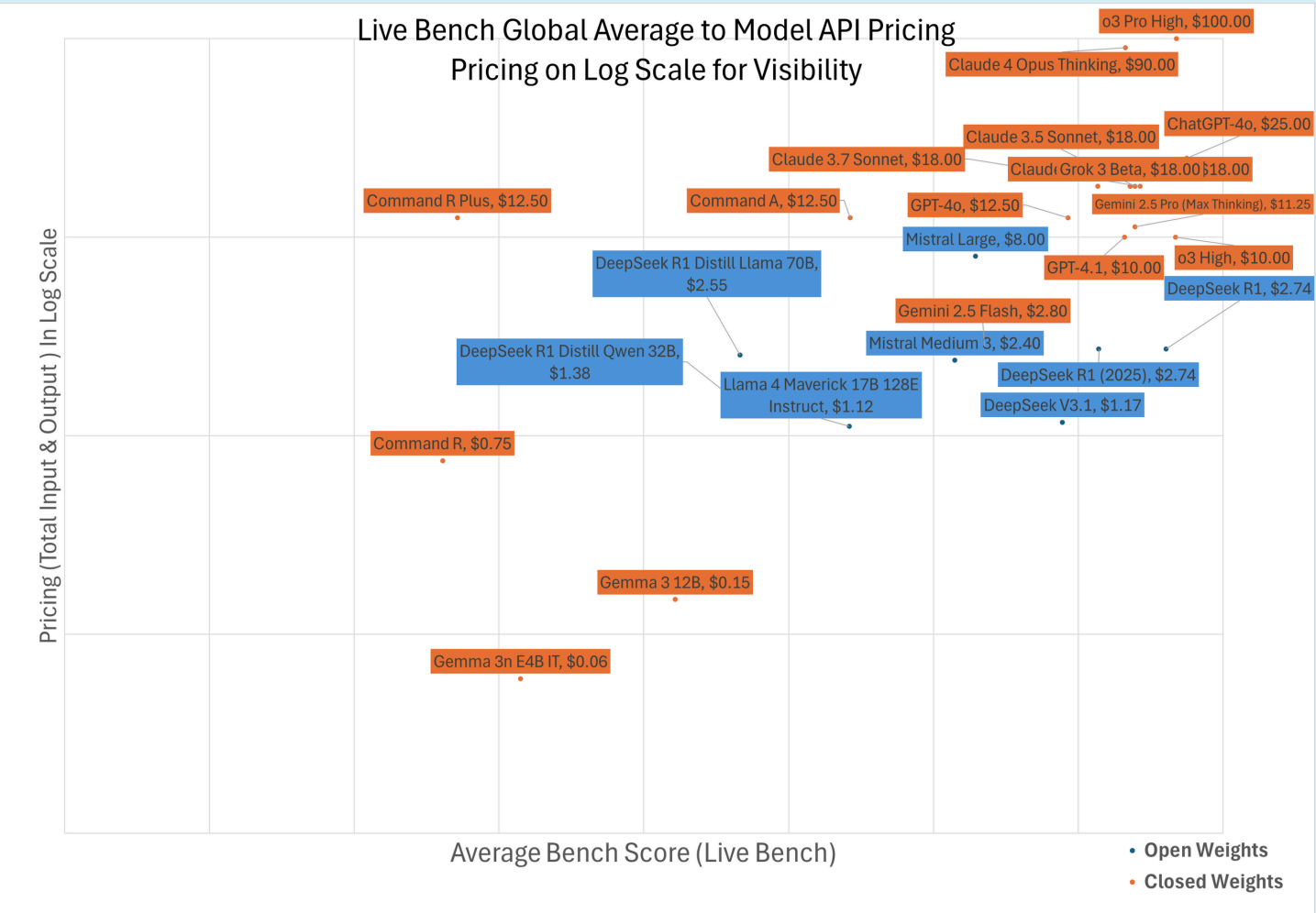Combines LLM with an external knowledge base (ex. TDC) to provide more accurate answers

RAG HTML → PDF

gui app

*These are clickable!*

*\*Filtered Out Compressed & Chinese Models*

Live Bench Global Average to Model API Pricing
Pricing on Log Scale for Visibility

## Economics of Gen AI

*These are clickable!*

### Large Language Model Cost Efficiency Analysis: Hardware and Software Trends (2023-2025)
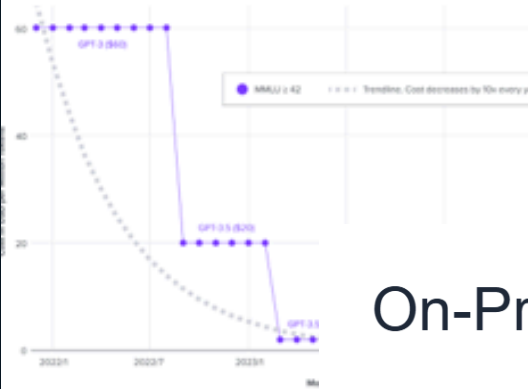
*Parts of this report are generated. I have taken measures to ask people around the office the accuracy of particular sections to evaluate where I have limited expertise.*

**Executive Summary**

Industry analysis suggests that the cost of large language model (LLM) inference has declined substantially over the past two years, with estimates indicating **5-10× annual cost reductions** for equivalent performance benchmarks. This trend appears driven by convergent improvements in GPU architectures, memory systems, software optimizations, and competitive market dynamics. However, the sustainability of these dramatic cost reductions faces potential constraints from supply chain limitations, manufacturing yields, and physical scaling challenges that warrant careful consideration in strategic planning.
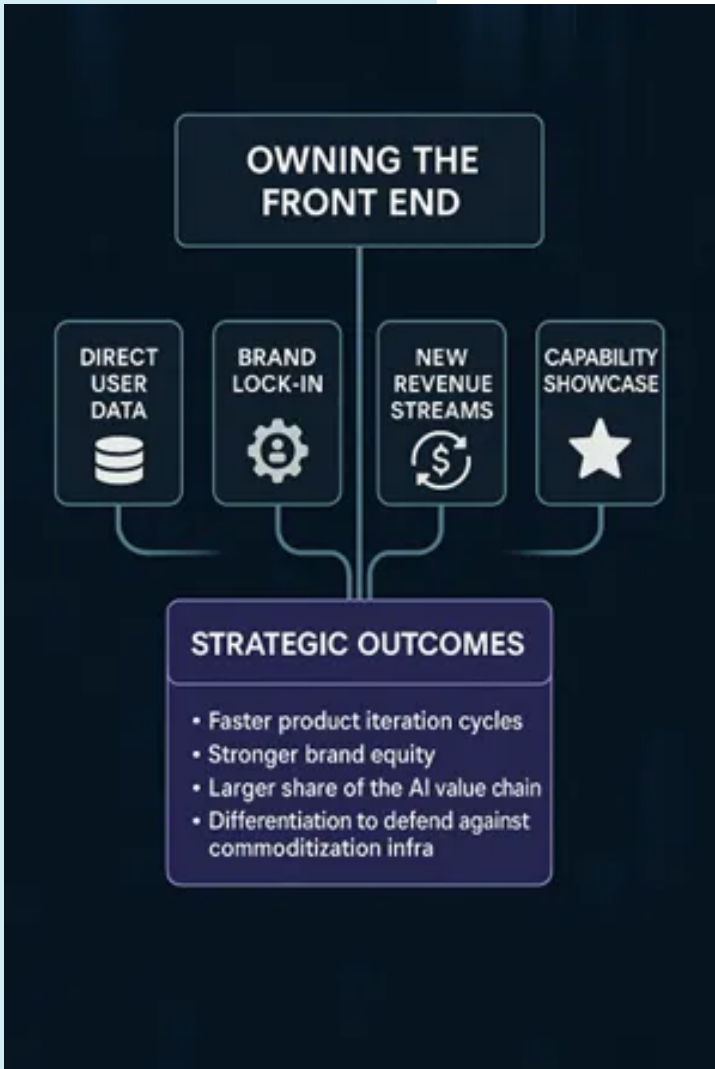
Early indicators show that models achieving GPT-3-level performance (MMLU score ~42) may have seen pricing fall from approximately $60 per million tokens in late 2021 to under $0.10 per million tokens by [...] cost reduction of 500-1000× over this period. While this trend has significant [...] adoption, decision-makers should consider quality variations, use case specificity, and [...] evaluating these metrics.

Cost of the Cheapest LLM with a Minimum MMLU Score of 42

https://a6z.com/llmflati[...]

### On-Premise vs. Cloud Decision Helper

Welcome! This tool will help you decide whether an on-premise self-hosting solution or a cloud service is better suited for your project's generative AI needs.

For each statement, rate how much you agree with it for your specific needs on a scale of 1 to 10.

Start Assessment

**OWNING THE FRONT END**

DIRECT USER DATA · BRAND LOCK-IN · NEW REVENUE STREAMS · CAPABILITY SHOWCASE

**STRATEGIC OUTCOMES**

- Faster product iteration cycles
- Stronger brand equity
- Larger share of the AI value chain
- Differentiation to defend against commoditization infra

## WeShare Pages

- [Costs of Generative AI](#)
- [On-Premise (Open Weight) vs Cloud Service (Closed Weight)](#)
- [5 Year CapEx for on-Prem vs. Cloud Server Hosting](#)
- [Case for Newer (more Expensive) Models](#)

*Using Labled Data*

*Reward-Based, graded output*

## Vocab I've Learned

- **SFT**: Supervised Fine Tuning
- **RFT**: Reinforcement Fine Tuning
- **Mixed Precision**: Combining number precisions for faster computation
- **Hill Climbing**: [read here](#)
- **NRE Assessment**: non-recurring engineering cost
- **POC**: Proof of Concept