

Why Token Optimization Is a Gift to the Hyperscalers



UNCOVERALPHA

JUN 29, 2026 · PAID



99



4



14

Share

Hey everyone,

A few weeks ago I wrote a piece called [Most of the Economy Won't Run on the B Model](#), where I argued that the AI market will bifurcate: the frontier model goes small slice of work where intelligence is unbounded in economic value (drug discovery, novel math, the hardest agentic reasoning), and the middle of the economy — classification, extraction, summarization, routine code, support — runs on the cheapest model that clears the quality bar.

This article is the sequel to that one. Who captures the value when the world shifts from token maxing to token optimization?

The shift away from always-buy-the-best-model is, on its surface, bearish for the labs and looks like it should compress the whole stack. But it is quietly one of the bullish structural setups for the three hyperscalers — Microsoft, Amazon, and Google.

Let me explain why.

Think about a highway toll road. There are two businesses operating on it. The first is the company that manufactures the cars — they make the actual machine that does the work of getting you somewhere, and there's a real margin in a car. The second is the company that owns the tollbooth. The tollbooth owner doesn't care what car you drive, whether it's a Ferrari, Toyota, or a 12-year-old used Honda. Every one of them pays the same toll to cross the bridge. Now imagine a world where, suddenly, everyone realizes they're wasting money commuting to work in a Ferrari for no reason, and they all downgrade to the cheapest Honda to save money. The carmaker's revenue per vehicle collapses, but people don't drive less because they switched to a cheaper car. They drive more, because now

cheap enough to justify trips they'd never have taken before. And every one of the extra trips still crosses the bridge. The tollbooth's revenue goes up.

In the AI stack, the AI labs are the carmakers. The hyperscalers are the tollbooth. And in the next months, we are entering a period where most of the economy is trading down from the Ferrari to the Honda, while simultaneously driving 10x more miles.

From token maxing to token optimization

For the last 18 months, the dominant behavior in enterprise AI was token maxing: you found the single best model on the leaderboard, you pointed every workload at it, and you didn't think too hard about cost, because the whole thing was a pilot and the budget was small relative to the perceived upside of "does this even work?"

That era is ending fast. Companies are essentially blowing through their AI budgets a quarter. Altman said recently that the "my company spent my entire 2026 budget in Q1, can you make this more efficient?" complaint went from something that "never came up" to "all of a sudden a huge issue." And it's not just him; you can see it across the industry, from companies like Salesforce and Meta to many other smaller ones, all blowing their planned yearly budget in a matter of days.

The natural behavior change from this is that instead of having one model for everything, companies start routing: a small, cheap, often open-weight model handles the 80% of requests that are less complex, and only the genuinely hard requests escalate to the frontier. This is token optimization, and it doesn't reduce total token consumption — it accelerates it. The moment inference gets cheap enough, you stop rationing it. You run the agent in a loop. You let it read the whole codebase. You run it five times and vote on the answer. A single coding-agent session now chews through millions of tokens of context where a chatbot query used a few thousand.

You can see this in the hard numbers the hyperscalers themselves disclose. Microsoft said it processed over 100 trillion tokens in a single quarter in 2025, up 5x year-over-year, with a record 50 trillion in one month alone. By its fiscal Q3 2026 call, Microsoft said over 300 customers were on track to process more than a trillion tokens each

Foundry this year — and that this was accelerating 30% quarter-over-quarter. Google went from 480 trillion tokens per month at I/O in May 2025, to 980 trillion by July, 1.3 quadrillion by October — and in its Q1 2026 filing disclosed that its first-party models alone were processing more than 16 billion tokens per minute via direct API, up 60% in a single quarter.

At the same time, the price per unit of capability is falling substantially— the Standard HAI AI Index found inference cost for GPT-3.5-level performance fell more than 10-fold in two years, and a16z pegs the decline at roughly 10x per year for any fixed capability level. And yet, total tokens processed are growing several-fold per year.

Where the margin lives in a routed token economy

When a company calls the SOTA model directly through, say, an AI lab's first-party API, the lab captures the full economic rent of the token. There are many industry reports surfacing around the margins that providers like Anthropic have right now and many of them claim that the margins went substantially up this year to as high as 70% gross margins. That price embeds the lab's R&D, its brand, its leaderboard position — call it the “model-provider margin.”

And to be fair, when token optimization happens, there might not be a direct hit to margins in the short-term at these frontier labs (as there will always be use cases for the best AI model), but it does mean that growth acceleration becomes smaller than it would be if everyone stayed for every use case on the frontier path.

The company routes that same workload to an open-weight model — GLM 5.2, DeepSeek V3.2, Qwen3 Coder, Kimi K2.5, Llama, MiniMax M2.5 — running on a hyperscaler's managed inference. The “model-provider margin” essentially goes to zero, because the model is open-weight and nobody is charging a brand premium for it. But the token still has to run on somebody's GPUs, inside somebody's data center, behind somebody's managed API with its security, compliance, logging, and SLA. And that somebody are the hyperscalers. They still charge their full infrastructure margins on the token. AWS has historically run a roughly 35–38% operating margin on the token. Google Cloud, which lost money for years, posted operating margins above 33% in

2026 and is still climbing. That margin doesn't care whether the token came from a \$50/million frontier model or a \$1/million open-weight model. To return to the analogy that from before, the tollbooth charges the same toll regardless of the car.

So here is the structural beauty of it for the hyperscalers, and the structural danger for the labs:

- Per-token economics compress, but the compression lands almost entirely on the model layer, not the infrastructure layer. The AI lab's margins on simple workloads get squeezed. The hyperscaler's infrastructure margin is sticky.
- Total token volume explodes, and almost every single token crosses the hyperscaler's tollbooth. More usage, on a partly-depreciated, increasingly efficient installed base, means absolute infrastructure revenue and gross profit dollars go up even as the price of any individual token falls.

This is Jevons' paradox pointed directly at the cloud P&L. The hyperscalers squeeze more revenue out of the infrastructure they already own, and they don't need the model-provider margin to do it. Outside of Google, they were never in that business in the first place.

The orchestration layer brings real value

If the future is multi-model — and it clearly is — then someone has to own the orchestration: the layer that decides which model handles which request, holds the models, fine-tunes, runs the agent loop, manages memory and tool-calling, handles fallbacks when a model is down. This layer can bring immense value to those who own it, and all three hyperscalers are naturally positioned to become just that.



Amazon's Bedrock is the clearest example. It's no longer "a place to call Claude." In 2026, the catalog spans 18 providers and 110+ individually addressable model variants, and AWS bolted on Intelligent Prompt Routing, which automatically routes each request to the cheapest model in a family that can handle it — AWS claims up to 30% cost savings with no accuracy loss. On top of that sits Bedrock AgentCore, a production agent harness with built-in Runtime, Memory, Gateway, Browser, IDE, and Observability. Bedrock reached a multi-billion-dollar annualized run rate with customer spend growing 60% quarter-over-quarter across 100,000+ customers.

Microsoft's Azure AI Foundry and Google Vertex are the same play, although Google's preferred scenario is where the Gemini model family dominates the workloads.

The natural position for hyperscalers to be the orchestration layer is great because companies already have their data, security perimeter, billing, and compliance living in the cloud environments. They can pick from a menu of models." The lab moves from the front end — the thing the customer chooses and bonds with — to the back end, an interchangeable component sitting behind the hyperscaler's harness. Customers start building the fine-tunes, the RAG pipeline, the agent scaffolding, the evaluation suite in the cloud orchestration layer, not at the model layer; switching models becomes

something normal and not a migration task. One of my readers, who runs his own harness, put it perfectly in the comments of my last piece: the harness is worth as much as the model itself, and you're paying for it either way.

The moat for the labs becomes shallow with most use cases but still remains important at the most complex ones (the ones where economic upside is not capped). But for hyperscalers, they benefit in all use cases that run the economy.

One might argue why companies couldn't move even outside of the cloud, but the problem is here that besides the normal reasons that people migrated a big part of workloads from on-prem to cloud (easier to scale, use, collaborate, manage), now on top of those reasons, you also have the problem of compute shortages where infrastructure outside of cloud environments is even harder to get and managing infrastructure becomes a much more difficult job. But an important aspect, I think that is forming is also the cybersecurity one. Given the recent developments here, it's clear that only a handful of companies will have first-row access to the most capable cyber models, to first shield their systems before the model is then available for everyone else. If you are an enterprise, you want safety, and it seems the only way you to get it is if you have your infrastructure at one of those preferred partners that get first-row access to cyber-capable models before everyone else does. The hyperscalers are those partners.

The hyperscalers' job, then, is almost simple to state: enable as many models as possible, make routing and fine-tuning and observability frictionless, and let customers optimize tokens to their heart's content. Every model they add makes the tollbooth more valuable.

The government approval window just made this worse for the labs

Now layer on a development from June that I don't think the market has connected to this thesis at all: the June 2, 2026 executive order: Promoting Advanced Artificial Intelligence Innovation and Security.

The headline version is that it sets up a voluntary framework where developers of "covered frontier models" — designation determined by the NSA via a classified

benchmarking process focused on cyber capabilities — give the federal government access to the model for up to 30 days before release to “other trusted partners.” (An earlier draft had a 90-day window; it was cut to 30 to avoid blunting U.S. competitiveness.) The order pointedly does not create mandatory licensing — but establishes a structured pre-release evaluation pathway, an AI cybersecurity clearinghouse run out of Treasury, and a government role in selecting which “trusted partners” get early access. The catalyst, per CFR’s reporting, was rising concern that models like Anthropic’s Claude Myths being able to autonomously find and exploit software vulnerabilities. Commerce ordered Anthropic to cut off non-U.S. access to Myths 5 and Fable 5 models on export-control grounds, and those models were removed from Bedrock days after launch. Now there is growing concern that, even when models are released to the public, U.S. citizens may get access before people from other countries. Hyperscalers have clients worldwide, so it becomes even more important for them to offer multiple models to companies and for companies to have an orchestration layer that isn’t just one AI lab, since they will have to juggle geopolitical compliance as well.

It strengthens the case

for entering AI through the hyperscaler. If you’re a regulated enterprise — a bank, hospital, a utility, exactly the critical-infrastructure operators the EO names — your safest path to frontier capability is increasingly through a cloud provider that handles the compliance, security review, data residency, and access-control machinery for you and gives you a menu of pre-vetted models. The EO effectively raises the compliance overhead of touching frontier AI, and compliance overhead is precisely the thing hyperscalers are built to absorb and resell.

There is another aftereffect of this government policy. It narrows the economic window the AI labs have to monetize a frontier model. A lab builds a genuinely state-of-the-art model. Historically, it gets to sell direct access to that frontier capability at a premium for as long as it takes competitors and open-weight models to catch up — that catch-up window is the lab’s profit window. But now insert a mandatory-in-practice review period and a “trusted partners first” distribution gate at the front

that window for the most capable (most cyber-capable, most monetizable) model frontier model's clock starts later and runs against open-weight competitors whose capabilities are compounding fast.

Also imagine a closed competitor releases a model that is less cyber-capable than the designated frontier model and therefore sails through review and reaches the open market faster, it might be days after the SOTA model, as the SOTA model was locked in the review process longer, meaning the duration of the economic value of the SOTA for the AI lab shrinks even more. The most capable model can paradoxically reach broad commercial availability slower than a slightly weaker one. The premium window — the period where a lab has the best model and can freely sell it to everyone. The slice of time where “best model” equals “uniquely monetizable model” is shrinking long as we have at least two companies competing in the SOTA race.

And who is indifferent to all of this? The tollbooth. The hyperscalers host the frontier models and the open-weight models and the not-quite-frontier closed models. Whichever one wins the customer's routing decision in any given quarter, the toll still crosses their bridge. They are structurally hedged across the entire model layer which is exactly where you want to be standing when the model layer is getting commoditized and regulated at the same time.

Summary

The intuitive read on “token optimization” is bearish: prices are collapsing, the models are getting commoditized, surely this compresses the whole AI trade. And it is bearish — for the labs' pricing power on commodity workloads (but not for the most complex tasks), and for anyone whose thesis depends on a single model maintaining a duration premium.

But the same forces that compress the model layer expand the infrastructure layer's take. Here's the asymmetry as I see it:

The hyperscalers win on volume, mix, and stickiness simultaneously. Volume, because cheaper tokens get consumed in vastly greater quantity and everyone pays the infrastructure toll — you can watch it in the token-processing disclosures growi

several-fold a year against per-token prices falling 10x a year. Mix, because as workloads shift to cheap and open-weight models, the lab margin in each token evaporates but the hyperscaler margin is untouched, so the hyperscalers capture larger share of a token's total economics than before. And stickiness, because the orchestration harness — Bedrock AgentCore, Foundry, Vertex — is where the customer's real switching costs accrue, not the model.

On top of that, this lands on a cost base that is increasingly paid-for. The 2026 capex is enormous — AWS guiding to ~\$200B, Microsoft over \$100B, Google \$175–185B — the bear case is that this capex never earns a return. But if you can squeeze an order of magnitude more useful tokens out of installed hardware and keep filling new hardware with exploding demand, you get operating leverage going vertical: incremental revenue landing on partly-depreciated, dramatically-more-efficient-per-token infrastructure. That's the scenario where cloud converts from a capex incinerator to a cash machine.

The labs, by contrast, are being squeezed into a narrower and narrower strip of the value map: the genuinely frontier, genuinely unbounded-value workloads where someone will still pay the Ferrari price — and even that strip now has a regulatory speed bump bolted to the front of it. It's a real business. Anthropic's trajectory shows there's enormous value there. But it's a narrower business than the “every token runs on the best model” dream, and it's one where the customer relationship increasingly belongs to the cloud provider sitting in front of them.

The market right now is mostly priced for a world where capex keeps compounding and the hyperscalers are making a losing bet on it. I think that has it close to backwards. The token-optimization shift everyone reads as a deflationary threat to the owners of the tollbooth, a Jevons-paradox tailwind with a regulatory moat that is in for free.

As always, I hope you found this article valuable. I would appreciate it if you could share it with people you know who might find it interesting.

Thank you!

Disclaimer:

I own hyperscaler Meta (META), Amazon (AMZN), Microsoft (MSFT), Google (GOOGL) stock.

Nothing contained in this website and newsletter should be understood as investment or financial advice. All investment strategies and investments involve the risk of loss. Past performance does not guarantee future results. Everything written and expressed in this newsletter is only the writer's opinion and should not be considered investment advice. Before investing in anything, know your risk profile and if needed, consult a professional. Nothing on this site should ever be considered investment advice, research, or invitation to buy or sell any securities.



99 Likes · 14 Restacks

Discussion about this post

Comments Restacks



Write a comment...



William Dwyer William's Substack 29 de jun.

Liked by UncoverAlpha

I enjoy your articles and your views.

Visit my substack: gentleinspiration.substack.com

Friendly Christian Advice for a Happy Life

LIKE (1) REPLY



Murali Murali 8d Edited

"On top of that, this lands on a cost base that is increasingly paid-for. The 2026 capex is enormous. AWS guiding to ~\$200B, Microsoft over \$100B, Google \$175–185B — and the bear case is that the capex never earns a return. But if you can squeeze an order of magnitude more useful tokens