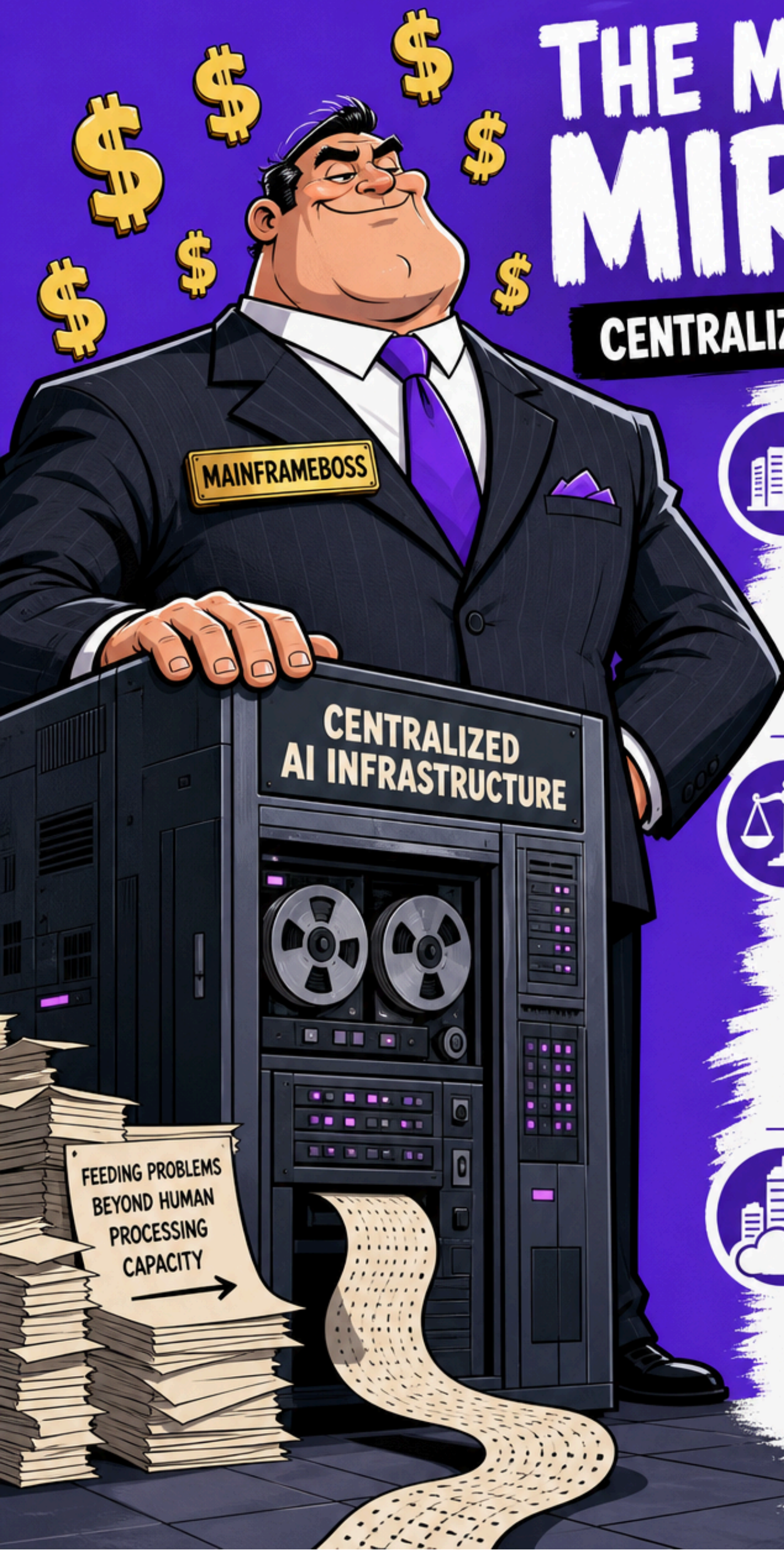


THE MAINFRAME MIRROR:

CENTRALIZED AI INFRASTRUCTURE



Today's AI infrastructure bears resemblance to the mainframe era of the **1960s and 1970s**. The scale is incomprehensible by 1968 standards, but the organizational logic remains: centralize computational resources, run them at maximum utilization, and pipe in problems that exceed human processing capacity.

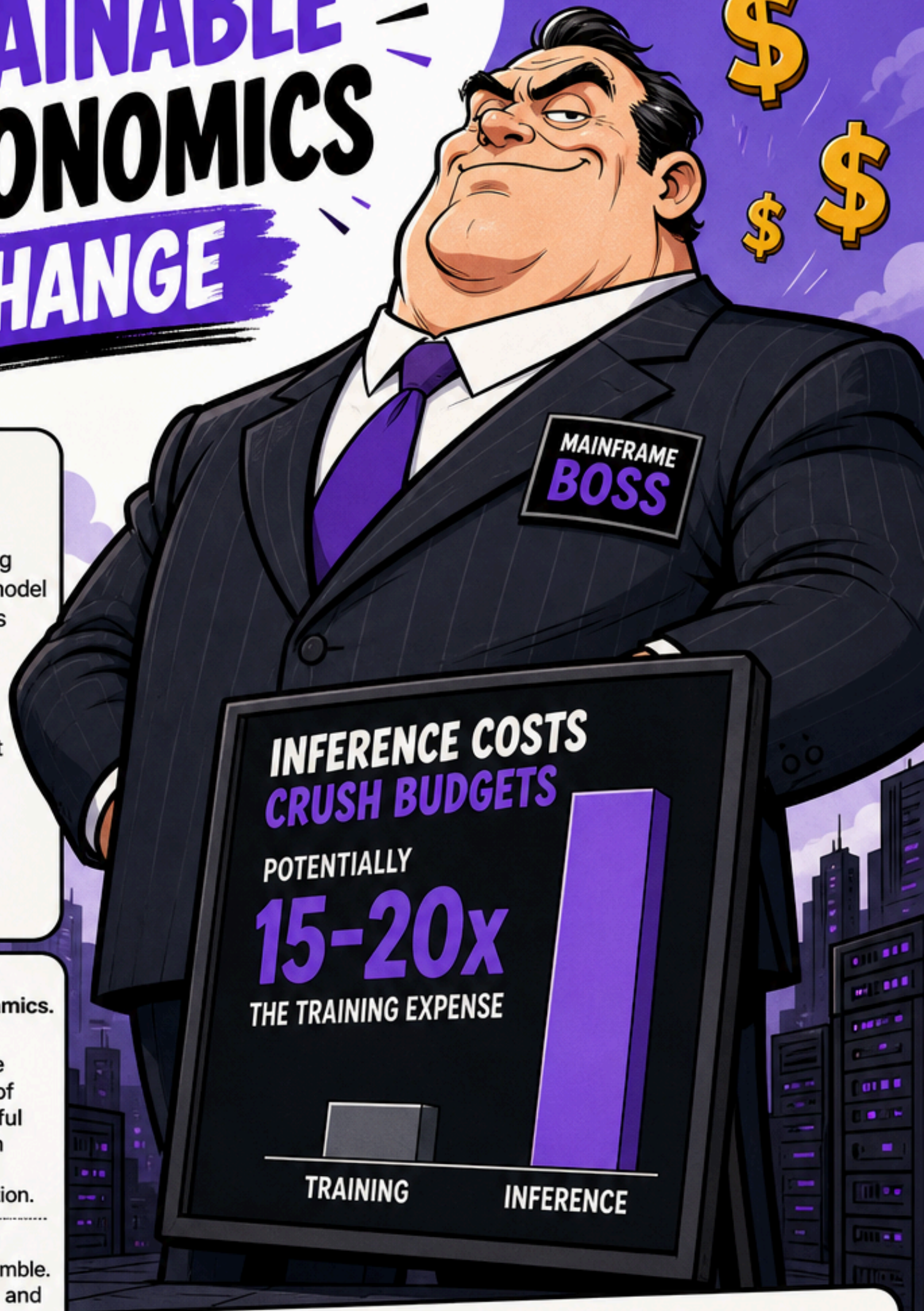


The economics show similarities. Organizations are hitting a tipping point where on-premises deployment may become more economical than cloud services for consistent, high-volume workloads. This may happen when cloud costs begin to exceed **60% to 70%** of the total cost of acquiring equivalent on-premises systems, making capital investment more attractive than operational expenses for predictable AI workloads.



Just as mainframes were accessible only to large institutions with substantial capital, today's AI capabilities remain concentrated among hyperscalers. In the public markets, hyperscaler companies are eating into their cash hoards and seeking alternative forms of financing to fund their AI **CAPEX (capital expenditures)** plans. Industry estimates suggest the top hyperscalers are collectively planning massive infrastructure investments in the coming years.

UNSUSTAINABLE UNIT ECONOMICS DRIVE CHANGE



Training gets the headlines, yet inference—the continuous, recurring cost of serving that model in production—crushes

budgets at potentially **15-20x the training expense.**



Industry analyses suggest that major language models could see inference costs **far exceed their initial training investments** over time.

far exceed their initial training investments over time.



This economic reality is creating challenging dynamics.

When metered billing is applied to an infrastructure stack that sits **idle 95%** of the time, the cost per useful token becomes a line-item emergency the moment a project moves into production.



5%



Enterprises locked in GPU capacity during the AI scramble. Now utilization sits at **5%** and the bill is due.



Industry estimates suggest **potential savings** for organizations willing to invest in dedicated infrastructure.

For sustained, high-utilization workloads, on-premises may deliver **substantially lower cost per million tokens** compared to cloud IaaS and commercial GenAI APIs.

ON-PREMISES BREAKEVEN POINT AGAINST CLOUD

ESTIMATED AT

**<4
MONTHS**

FOR HIGH-UTILIZATION
SCENARIOS

VS.

**12-18
MONTHS**

PREVIOUS GENERATION
CYCLES

THE EDGE REVOLUTION:

SPECIALIZED CHIPS ENABLE DISTRIBUTED AI



CENTRALIZED AI INFRASTRUCTURE



NPUs and specialized AI accelerators are now standard in smartphones, laptops, and industrial equipment.

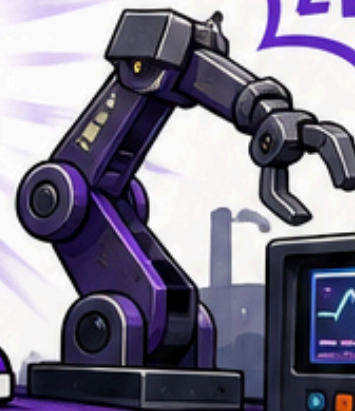
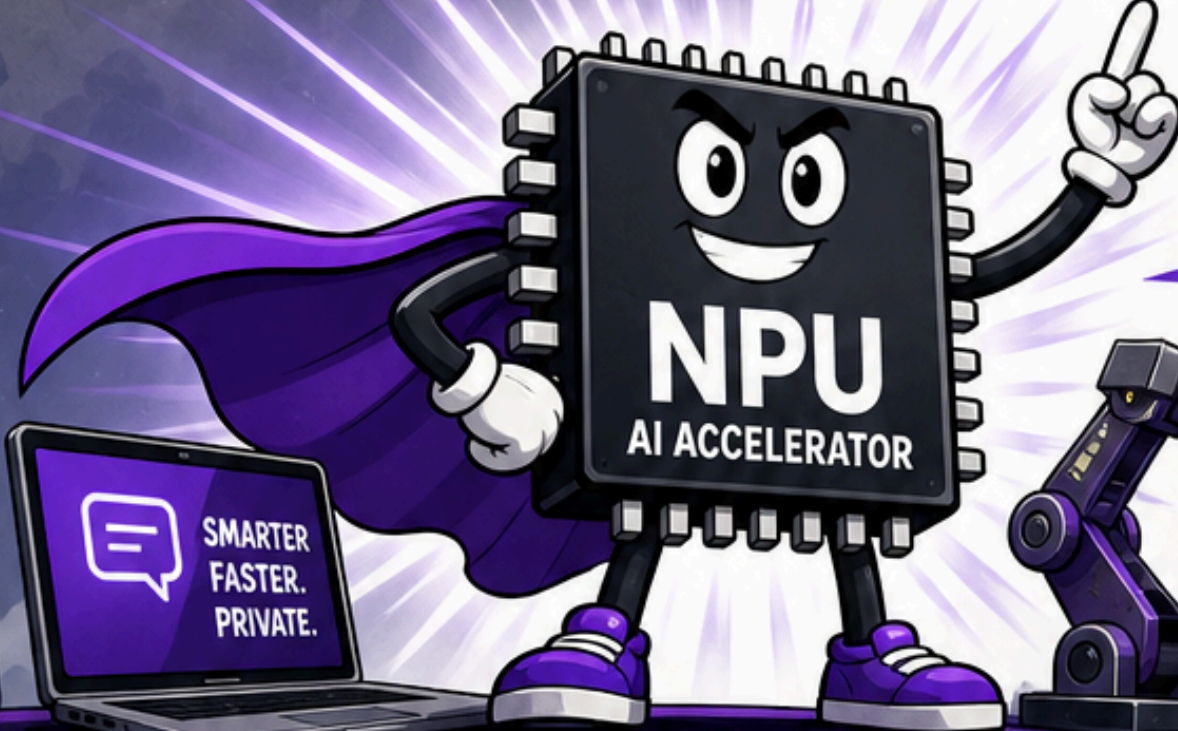


These chips deliver up to **10 trillion** operations per second while consuming just **2.5 watts** of power.



That's at least **6x** more efficient than traditional CPUs and mainstream GPUs for neural network tasks.

THE CHIPS ARE FINALLY INSIDE EVERYTHING.



INFLECTION POINT

NPUs are now standard components in consumer and industrial hardware. From premium to standard. That creates an enormous installed base for on-device AI.



MODELS SHRANK, NOT THE CAPABILITY

Compression techniques — quantization, pruning, knowledge distillation — have gotten much better.



SMALL MODELS, BIG IMPACT

1-7B parameter models now handle real tasks competently on normal consumer devices.

DISTRIBUTED AI. REAL-WORLD IMPACT. RIGHT AT THE EDGE.

MARKET FORCES DRIVING DEMOCRATIZATION OF AI

The same market forces that drove the mainframe-to-PC transition are building momentum in AI. This decentralization of AI mirrors the shift from mainframe to personal computing or the rise of cloud computing, each democratizing access to computational power in different ways.



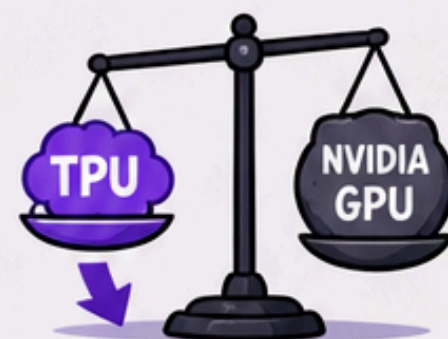
It democratizes access to advanced AI capabilities, moving them from the exclusive domain of hyperscale data centers to billions of everyday devices.

This transition is akin to the personal computing revolution, where computational power became accessible to individuals, or the cloud computing era, which provided scalable infrastructure on demand.

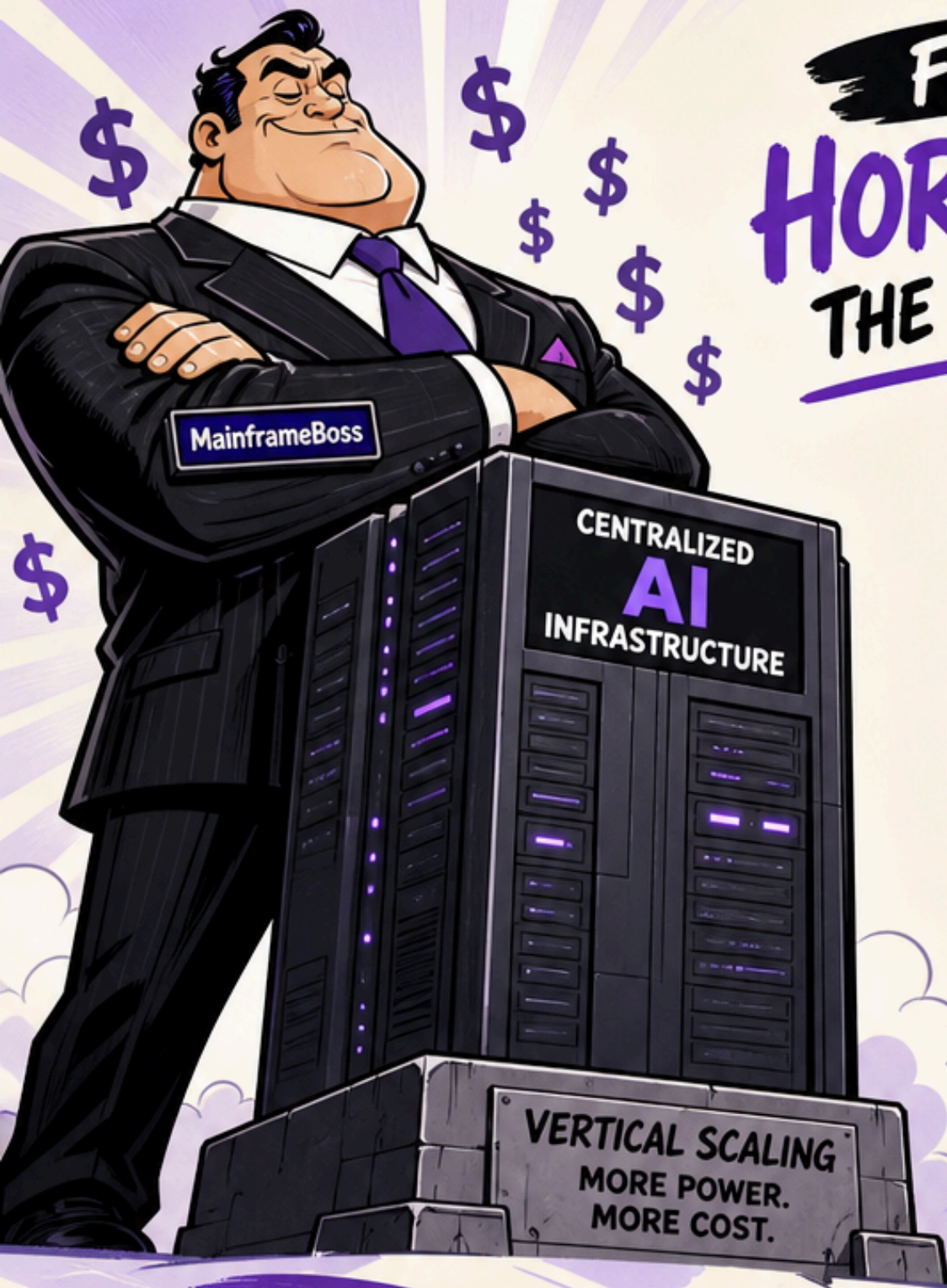


COST PRESSURES ARE ACCELERATING THIS DEMOCRATIZATION.

SemiAnalysis estimates that Google could cut its cost per computation relative to Nvidia by 62% with its own TPU for internal workloads. This mirrors how PC manufacturers eventually offered better price-performance than mainframes for many workloads.



FROM VERTICAL TO HORIZONTAL SCALING: THE AGENT NETWORK PARADIGM



AI agents are autonomous software entities designed to achieve specific goals, execute tasks independently, and make real-time decisions. They can operate on their own or collaborate within an “agentic mesh,” a network where agents coordinate seamlessly with other agents, tools, and transactional systems.



By combining multiple AI technologies, these agents can be orchestrated in a modular way, making it possible to handle even highly complex workflows from end to end.



NARROW FIRST. PROVE, THEN EXPAND.

- ✓ Narrow single-function agents scale more reliably than broad multi-function ones.
- ✓ Start with a single, well-defined task and measurable output.
e.g., document classifier, data enrichment pipeline, routing agent.
- ✓ Broad, open-ended agents fail at scale due to compounding quality variance and untestable edge cases.
- ✓ Expand scope only after the narrow version proves stable for 90+ days.

SHIFTING TO
HORIZONTAL
SCALING.



INFRASTRUCTURE BUILT FOR DISTRIBUTION

Industry adoption of standardized protocols like MCP is accelerating, with growing support from major AI providers:



Emerging agent-to-agent communication protocols are being developed to handle horizontal coordination across organizational and platform boundaries.



AGENTIC MESH

Agents coordinate seamlessly with other agents, tools, and transactional systems.



MORE AGENTS. MORE AGILITY.
BETTER OUTCOMES.

THE MID-TIER IMPLEMENTATION GAP



The transition from mainframes to PCs wasn't direct. It required a new ecosystem of software vendors, system integrators, and support providers. Similarly, AI's potential shift to distributed deployment is creating opportunities for mid-tier implementers to **bridge the gap** between platform giants and widespread adoption.



Industry analysts predict that agentic AI demand may accelerate faster than commercial software tooling can reliably support, creating opportunities where **services providers – especially global system integrators – could thrive**, even as traditional software-as-a-service vendors feel pressure.



Current data suggests implementation challenges. Recent surveys indicate that while AI agent pilots are becoming common among enterprises, successful production deployment remains challenging, with many organizations **struggling to scale beyond pilot phases**.



This gap creates opportunities for specialized implementers who can navigate the complexities of production deployment. Organizations with production-scale deployments were not spending more on AI overall — their total AI budgets were comparable to stalled organizations. The difference was allocation: successful scalers spent proportionally **more on evaluation infrastructure, monitoring tooling, and operational staffing**, and proportionally **less on model selection and prompt engineering**. The data suggests that scaling failure is a **build-vs-operate imbalance, not an underspending problem**.



PLATFORM GIANTS

MID-TIER IMPLEMENTERS
(SYSTEM INTEGRATORS, PARTNERS, SERVICES PROVIDERS)
BRIDGING THE GAP

WIDESPREAD ADOPTION

THE MID-TIER ECOSYSTEM



SOFTWARE VENDORS



SYSTEM INTEGRATORS



SUPPORT PROVIDERS

THE COMING TRANSFORMATION

The historical parallel suggests AI's potential "PC moment" could arrive faster than the mainframe transition.

WHAT TOOK
15 YEARS
COULD TAKE
5-7 YEARS
WITH AI



MAINFRAME TRANSITION



15 YEARS
TO MAINSTREAM ADOPTION

MODERN TRANSFORMATION



5-7 YEARS
TO MAINSTREAM ADOPTION

SEVERAL FACTORS SUGGEST THIS SHIFT MAY BE UNDERWAY:

1 ECONOMIC PRESSURE



Extended interconnection timelines, constrained data center capacity, and rising infrastructure costs are now intersecting directly with approved AI budgets.

Entering 2026, the question is no longer whether AI will require substantial infrastructure support, but how those requirements reshape the financial pathways available.

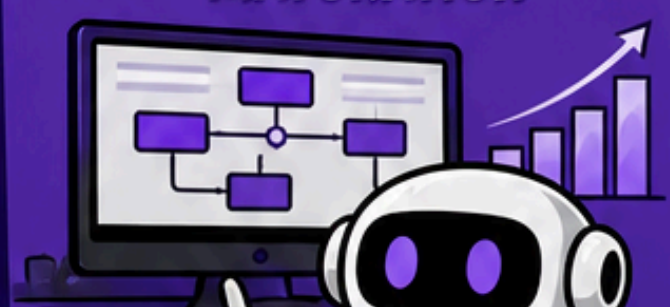
2 TECHNOLOGY READINESS



Edge computing is evolving rapidly with the rise of Edge AI. Smaller and more efficient models, often called Small Language Models or Micro LLMs, are being designed to run directly on devices.

This could allow laptops, vehicles, and smart home systems to understand language, detect patterns, and make decisions without cloud dependency.

3 MARKET MATURATION



- ✓ FRAMEWORKS EMERGING
- ✓ PRODUCTION DEPLOYMENTS GROWING
- ✓ SELECTION DECISIONS OPERATIONAL

The agent systems landscape is evolving rapidly. While early implementations required custom orchestration layers, frameworks are emerging with growing production deployments, and selection decisions are becoming more operational rather than experimental.

Supramono



Discover. Build. Sell. One AI Venture Engine.

<https://supramono.com>