

— WHY WE HAVEN'T SEEN TRUE PRIVATE AI BEFORE 2026

The Death of the Token Tax.

How four independent shifts just made Private AI not merely feasible, but economically superior — for most enterprise workloads. And how stacking optimizations flipped the economics of enterprise AI by summer 2026.

● DISTILLATION

● INFERENCE · LIFEBOAT

● HARDWARE ↓ 1,000×

● MODEL ISOLATION

Intelligence Unshared: The AI Sovereignty Papers — this is the fifth in a six-paper series on Private AI.
Find the series at iterate.ai/resources/white-papers

PREPARED BY

Iterate.ai · Engineering & Research
With the Lifeboat inference team

EDITION

Technical Brief
The economics of Private AI

SUBJECT

The Inflection
Why now, and not in 2023



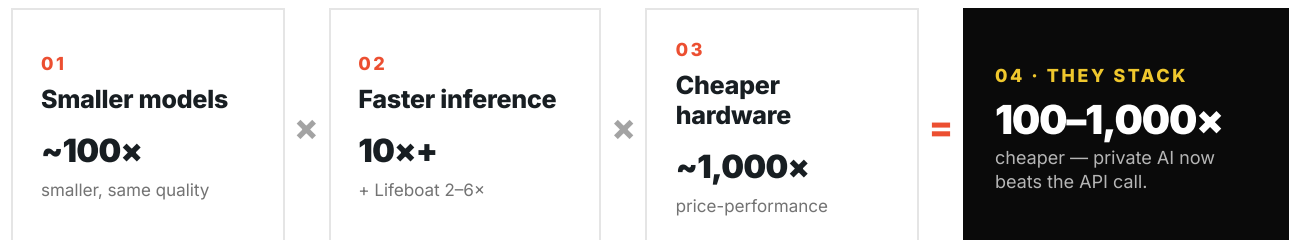
ABOUT THIS PAPER

For years, the honest answer to “why not run your own AI?” was cost and capability.

That constraint just broke. Between 2024 and 2026, four independent technological shifts converged — and the window they opened is still widening.

This brief is written for the engineering and platform leaders who have to decide whether to keep renting intelligence, or to build and run it themselves.

FOUR SHIFTS, ONE COMPOUNDING EFFECT



Each shift is real on its own. The fourth convergence is that they multiply — the compound effect, not any single advance, is the story. Full math on page 09.

HOW TO READ IT

Each of the four convergences is real on its own. The argument of this paper is that they *stack* — and that the compound effect, not any single advance, is what makes Private AI suddenly cheaper than the API call. The numbers on page 09 are the whole point.

WRITTEN FOR

CTOs, platform and infrastructure leaders, and the architects deciding where enterprise inference should run.

PUBLISHED BY

Iterate.ai — San Jose, CA & Denver, CO. Private AI infrastructure for the enterprise.



CONTENTS

What's inside.

Four convergences, stacked. The benchmark table on p.07 and the stacked-economics on p.09 are the two pages most readers return to.

THE SETUP · THE CONSTRAINT THAT JUST BROKE		04
PART ONE · THE FOUR CONVERGENCES		05 — 09
01	Smaller Models, Same Intelligence Distillation closed the gap between “affordable” and “frontier.”	05
02	The Inference Optimization Revolution Continuous batching, PagedAttention, RadixAttention, TurboQuant.	06
02	Lifeboat — 2–6× Beyond the Best Open Source Iterate’s KV-cache innovations, and what they do to the per-token cost.	07
03	Hardware Costs Collapsed A 1,000× price-performance improvement, and where it came from.	08
04	The Convergence — Stack the Gains None of these work alone. Multiplied, they reach 100–1,000×.	09
THE ROAD TO PRIVATE AI · 2017–2026 TIMELINE		10
WHERE INFERENCE RUNS · SHARED VS. PRIVATE GPU		11
PART TWO · WHY NOW, AND THE WINDOW AHEAD		12
05	Why Now? Why Not 2023? What changed in four years — and the 18–24 month first-mover window.	12
THE ACCELERATION · THE PAST & THE FUTURE OF THE CURVE		13 — 14
THE ECOSYSTEM ERA · DATA, FLYWHEELS & WHO OWNS THE INTELLIGENCE		15
BUILD. RUN. GOVERN. · COLOPHON & REFERENCES		16



THE SETUP

The economics used to be brutal: rent intelligence, or don't use frontier AI at all.

Training GPT-4 took an estimated **\$100 million** in compute, more than 1,000 GPUs, water cooling, and a dedicated data center. Running it at scale demanded infrastructure only hyperscalers could afford.

For years, the answer to “why not run your own AI?” was simple: cost and capability. The intelligence you wanted lived inside models you could not afford to train and could not afford to serve. So you rented it — per token, from a handful of providers — or you did without.

None of the four is the story. The story is that they **stack** — and the compound effect is a 100–1,000× improvement in cost-performance.

That constraint just broke. Between 2024 and 2026, four independent technological shifts converged to make Private AI not just feasible, but economically superior for most enterprise workloads. **The window opened fast, and it is still widening.**

01

Smaller models, same intelligence

Distillation — 95–97% of frontier quality at 1/100th the size.

02

Inference software got 10× better

And Lifeboat goes 2–6× beyond the best open source.

03

Hardware costs collapsed

~1,000× price-performance gain since 2021.

04

The gains stack

Multiplied together: 100–1,000× cost-performance.



01 THE DISTILLATION BREAKTHROUGH

Smaller models, same intelligence.

In 2023, a 70B-parameter model was “small.” In 2026, production systems run on **3B-8B** models that deliver **95-97%** of frontier performance at 1/100th the size.

HOW THIS HAPPENED

Model distillation transfers knowledge from a massive “teacher” model — GPT-4, Claude Opus — into a compact “student.” The smaller student carefully copies the kinds of answers the giant teacher gives — think high-tech shrink-wrap: it squeezes the big model’s writing styles and step-by-step problem-solving into a far smaller size, giving you the same brainpower without a massive, expensive computer system to run it.

DeepSeek-R1’s 2026 result proved even complex reasoning can be distilled. Using **800,000** high-quality reasoning samples, it compressed a frontier reasoning model into an 8B student that outperformed directly-trained small models on mathematical benchmarks.

DISTILBERT

97% of BERT’s accuracy at 40% of the size and 60% faster inference.

LLAMA 3.2 3B

72% lower latency than Llama 3.1 405B, with comparable quality on most enterprise tasks.

WHAT A DISTILLED 7B MODEL ON A SINGLE ENTERPRISE GPU NOW HANDLES

- Document analysis and summarization
- Code generation and review
- Multi-step reasoning workflows
- Customer support automation
- Strategic research and synthesis
- No \$100M budget. No 1,000 GPUs.



The intelligence gap between “affordable” and “frontier” has collapsed.



02 THE OPTIMIZATION REVOLUTION

Inference software got 10x better.

Training a model once is expensive. Running it millions of times — *inference* — is where the real cost lives. In 2025–2026, inference software made a generational leap. For non-tech execs: **inference** = each time the AI answers; a **GPU** is the chip that does it; the **KV cache** is its short-term scratchpad.

Continuous batching

VLLM · SGLANG · TENSORRT-LLM

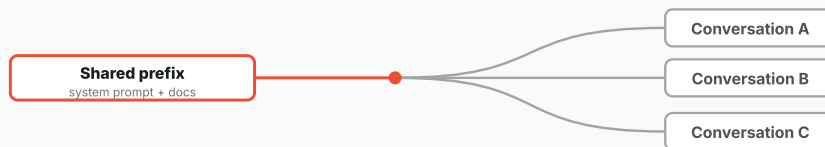
Traditional inference processed one request at a time, leaving GPUs ~80% idle. Continuous batching dynamically groups requests as they arrive, raising utilization from under 20% to over 70% — **3–10x higher throughput** on the same hardware. Picture it like a roller-coaster that fills every empty seat as riders arrive, instead of sending one person around alone.

PagedAttention & RadixAttention

KV CACHE OPTIMIZATION

PagedAttention treats the KV cache like virtual memory — small fixed-size blocks instead of huge contiguous allocations — cutting memory waste 60–80% and allowing **24x more concurrent requests**. SGLang’s RadixAttention stores caches in a radix tree, reusing shared prefixes (system prompts, document context) across requests and cutting time-to-first-token **30–50%**. Picture it like neat loose-leaf pages instead of one messy pile — and three students sharing one master copy of the same article.

HOW RADIXATTENTION REUSES MEMORY — THE “SHARED MASTER COPY”



The system prompt and document context are stored **once** and instantly reused across every conversation — instead of re-processing the same words for each one.

Google TurboQuant

DEEPMIND · ICLR 2026

A dual-mechanism approach — PolarQuant (rotation + Lloyd-Max quantization) plus QJL (1-bit residual correction) — that compresses the KV cache **6x** with no measurable accuracy loss and an **8x speedup** in attention on H100s. At 3-bit, the KV cache of Qwen2.5-3B at 8K tokens drops from 289 MB to ~58 MB. Picture it like a vacuum storage bag squeezing the wasted air out of the data, so it fits on smaller chips and moves faster.

THE COMPOUND EFFECT OF OPEN-SOURCE TOOLS

3–10x

THROUGHPUT
BATCHING

2–4x

MEMORY
PAGEDATTN

6x

KV CACHE
TURBOQUANT

2–5x

LATENCY
SPEC. DECODE

73%

ENERGY
SAVINGS



Lifeboat goes 2–6× beyond the best open-source engines.

vLLM, SGLang and TensorRT-LLM are the state of the art in open source. Building on Google's TurboQuant foundation, Iterate's **Lifeboat** engine delivers **2–6× better performance** — both faster and cheaper — through proprietary KV-cache innovations.

LIFEBOAT'S KV-CACHE INNOVATIONS

- 1 Adaptive KV compression**
 Dynamic per-layer ratios from attention-pattern analysis — **8–12×** compression on long-context workloads, accuracy intact.
- 2 Zero-copy KV reuse**
 Cross-request cache sharing without serialization overhead, cutting memory bandwidth **40–60%**.
- 3 Predictive eviction**
 ML-guided eviction that preserves reasoning-critical attention heads — **20–50%** cache reduction, near-lossless.
- 4 Hardware-aware scheduling**
 Custom CUDA kernels co-designed with the compression algorithms for H100/H200 tensor cores.

REAL-WORLD GAINS OVER VLLM / SGLANG

WORKLOAD	BASELINE	LIFEBOAT	GAIN
Short-context <4K tokens	12,500	25–30K	2–2.4×
Long-context 32K–128K tokens	3,200	12.8–19.2K	4–6×
Multi-turn prefix-heavy	8,500	25.5–34K	3–4×
Cost / 1M tok amortized	\$2.00	\$0.33–0.50	4–6×

Throughput in tokens/second unless noted. Baseline = best of vLLM / SGLang on identical hardware.

IN PLAIN ENGLISH

KV cache — the model's short-term memory of the current chat, kept handy so it needn't re-read every earlier word on each new one.

CUDA & kernels — NVIDIA's toolkit, plus the tiny tuned programs that make its GPUs do the math fast.

Serialization — repackaging data to copy it between processes; "zero-copy" skips that repackaging overhead.

H100 / H200 — NVIDIA's top data-center GPUs, the high-end chips most enterprise AI runs on.

A single H100 running Lifeboat now serves workloads that previously required **20–60 GPUs** on standard engines. The cost per query dropped from dollars, to cents, to **fractions of a cent** — making private inference routinely cheaper than the API call.



03 A 1,000× PRICE-PERFORMANCE GAIN

Hardware costs collapsed.

LLM inference costs have dropped ~10× annually since 2021. Equivalent model performance now costs roughly 1/1,000th of what it did four years ago.

GPU PRICING, THEN AND NOW

GPU	2022-23	2026	↓
A100 80GB	\$8-10/hr	\$1.50-1.80	5-6×
H100 SXM	\$7.50-11/hr	\$2.00-2.50	3-5×
Entry A10, L4	\$2-3/hr	\$0.50-1.20	2-5×
Spot fault-tolerant	—	\$0.20-0.60	batch

QUANTIZATION — SMALLER, SAME QUALITY

FP8: 50% less memory, 1.6× throughput on H100.

INT4: 75% less memory — a 70B model fits on a single GPU.

GPTQ: post-training compression in under a minute, minimal accuracy loss.

OPEN-SOURCE MODELS

Llama 4 · Qwen3 · DeepSeek-V3/R1 · Mistral — production-ready, commercially licensed, frontier-class. No training from scratch, no licensing fees.



What used to require 8× A100 GPUs now runs on 1× H100 with quantization.

A100 vs. H100: the **H100** (Hopper) is NVIDIA's newer, far faster chip; the **A100** (Ampere) is the prior generation — both are U.S. export-restricted to China.

Total cost of ownership for private infrastructure — hardware, software, operations — is now **far below** the cumulative API spend for most mid-to-large enterprises, once you count volume vs. per-token pricing, governance overhead, intelligence-leakage risk, and vendor lock-in. **A 70B model on dedicated infrastructure costs less to run than querying GPT-4 at scale.**



04 THE CONVERGENCE

None of these work in isolation. Stack them.

The real breakthrough is multiplicative. Each advance is meaningful; together they compound into a **100–1,000×** cost-performance improvement.

Distilled 7B model vs. GPT-4 scale	100× smaller
×	
Quantization (FP8) weight precision	50% memory
×	
Lifeboat engine vs. vLLM / SGLang	2–6× throughput
×	
Cheaper hardware A100 spot, vs. 2023	5–6× cheaper
=	
Cost-performance	100–1,000×

A REAL-WORLD EXAMPLE — 1M TOKENS

	2023 · SHARED API GPT-4	2026 · PRIVATE + LIFEBOAT Llama 3.1 8B
Input	\$30.00	\$0.33
Output	\$60.00	\$0.50
Total / 1M	\$90.00	\$0.83
Your intelligence	trains competitors	stays private

AT 100M TOKENS / MONTH

\$9M/yr → **\$83K**/yr

And the Private AI version doesn't train your competitors.



WHY NOW · THE ROAD HERE

Private AI isn't a trend. It's a decade of breakthroughs stacking up.

Each milestone below is an independent engineering advance. Together, they broke the constraint that once made frontier AI affordable only to tech giants.

- 2017** ● **“Attention Is All You Need”**
The Transformer arrives — the “T” in ChatGPT stands for Transformer — but frontier power needs giant infrastructure only tech giants can afford.
- Nov 2022** ● **ChatGPT launches**
The world realizes software can think — yet running it means \$100M of GPUs or rented public APIs.
- 2023** ● **Frozen weights, models too big**
Meta’s Llama opens the open-weight era — but early open models stay too weak and 70B+ too bulky, so enterprises still ship data to public clouds.
- 2024** ● **Distillation & quantization**
Models shrink ~100× with no loss of quality — frontier intelligence now fits a standard office server.
- Late 2024** ● **Continuous batching · PagedAttention**
vLLM-class software ends one-at-a-time inference — dozens of requests at once, GPU waste collapses.
- Nov 2025** ● **TurboQuant & RadixAttention**
KV-cache compressed up to ~6× with no accuracy loss; shared prompts reused instantly across sessions.
- Jan 2026** ● **Hardware price drop**
The chip shortage ends — spot & cloud prices fall 5–6×. By late 2025 Iterate had already run a full LLM + RAG on a Qualcomm 6490 chip — no grid power, no internet, no cooling.
- May 2026** ● **Petabytes become queryable**
Iterate **vectorizes petabyte-scale mass storage** (e.g., NetApp boxes) on small, air-cooled GPUs — opening once-inaccessible healthcare and government archives to chat and agents.
- June 2026** ● **The convergence — Private AI is accessible to all companies**
All four shifts stack: private inference becomes **100–1,000× cheaper** than public API tokens.
- 2026–27** ○ **What’s next: on-the-fly learning**
Speculative. Models begin adapting their own weights at inference (test-time training). Powerful — but this is when even greater risks emerge, as models learn *inside themselves*, beyond outside audit. At this point, isolation (private AI) matters more than ever.

WHY THE GAINS COMPOUND

Distilled models 100× smaller × Quantization 50% memory × Lifeboat 2–6× throughput ×

Cheaper chips 5–6× = **100–1,000× better cost-performance**

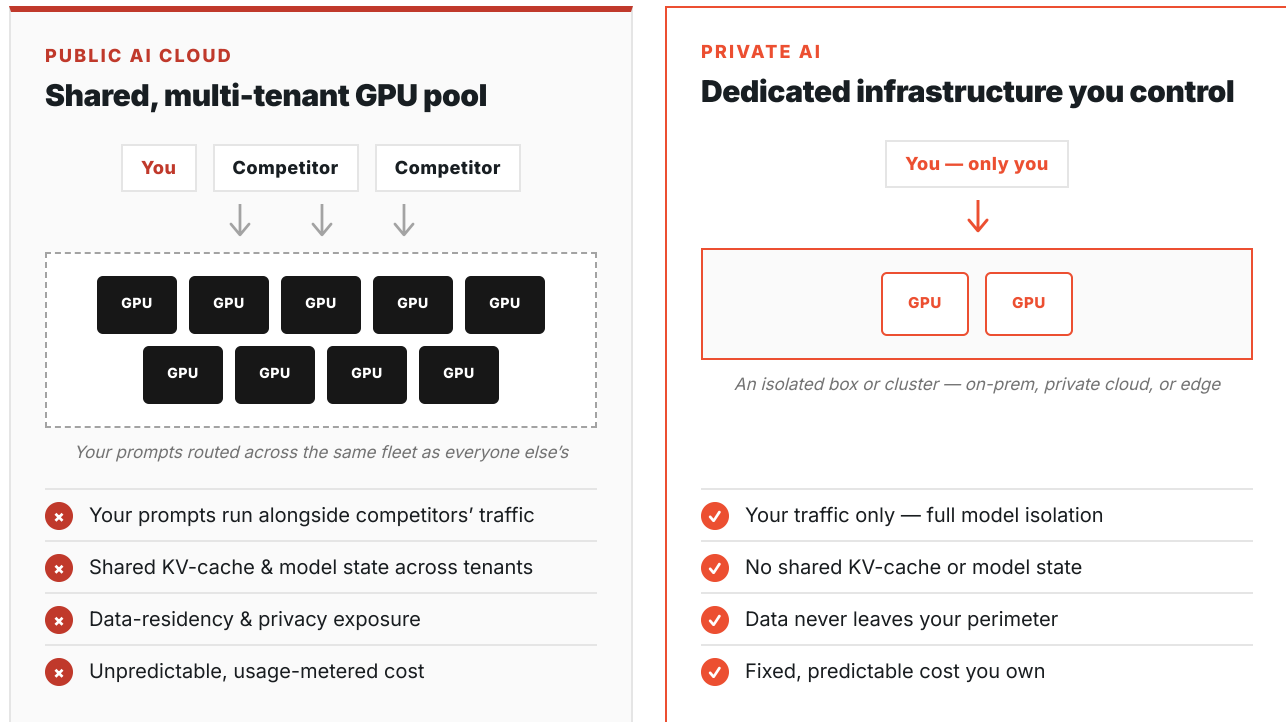
Iterate is on the invention train too — aimed squarely at Private AI: **8 granted AI patents, 13 pending, 10+ in draft** as of June 2026.



WHERE YOUR INFERENCE RUNS

Shared GPUs, or infrastructure you own?

The convergences make private AI *affordable*. But affordability is only half the story — **where** your tokens are processed decides who else can see them.



The architectural difference translates directly into corporate risk. Routing enterprise tokens through a multi-tenant fleet means your execution logic, context-retrieval patterns, and dynamic prompts are mixed into a shared pool — even if data-privacy contracts hold, the **security blast radius stays collective**. Private AI collapses this by trapping the entire execution lifecycle inside a dedicated container — **your operational wisdom compounds solely for you.**



05 WHY NOW? WHY NOT 2023?

The convergence happened fast — and most organizations haven't noticed yet.

	2023	2026
Models	Small models were too weak.	Distillation closed the capability gap.
Inference	Inefficient; GPUs mostly idle.	10x more efficient — and Lifeboat pushed it 2-6x further.
Hardware	Expensive; hyperscaler-only.	Costs dropped 5-6x.
Tooling	Immature.	Open-source tooling reached production grade.

THE WINDOW IS OPEN — BUT NOT FOREVER

First-movers in Private AI get **18-24 months** of compounding advantage:

- Lower costs while competitors burn cash on API calls.
- Faster iteration — no rate limits, no API downtime.
- Proprietary intelligence that compounds with every query.
- No training leakage to competitors still on shared models.

By the time the industry wakes up to model isolation, the leaders will be unreachable.

THE QUESTION

The question isn't "Can we afford Private AI?"

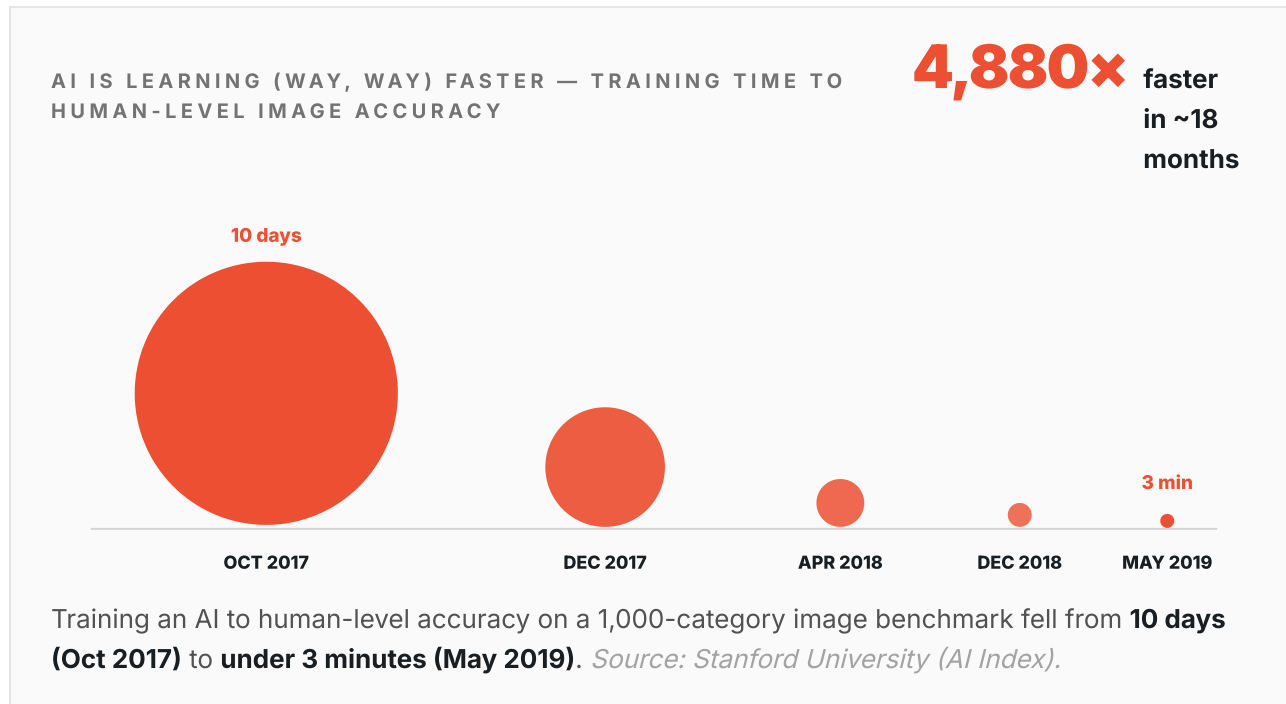
It's "Can we afford to keep renting intelligence from our competitors?"



THE ACCELERATION · THE PAST

We've watched this exact curve before.

The four convergences didn't come out of nowhere. The same exponential curve that crushed training time and chip cost years ago is the one still compounding today.



AND THE SILICON GOT TINY & CHEAP

The **AI (NPU) portion** of a flagship smartphone chip by 2020 cost just **\$1.42–\$5.10** — Huawei \$1.42, Samsung \$3.50, Apple \$5.10. The chip is so small that **roughly 50 of them fit on a U.S. penny** — AI cheap enough to live in your pocket, at the edge, off the grid.

Source: Deloitte, 2020. Historical figures.

~50 CHIPS ON A PENNY

Consumer products are getting brains and memories — courtesy of edge-AI chipmakers like **Syantiant** and AI software makers like **Iterate.ai**. Each newly intelligent device also becomes a new stream of data — fuel for the services built on top of it.



THE ACCELERATION · THE FUTURE

And it isn't slowing down.

The curve that made Private AI feasible in 2026 doesn't flatten — it steepens. Every quarter of waiting is more expensive than the last.

NOW: AND IT KEEPS GOING

"20 years from now, the rate of change will be 4x what it is now — and in 40 years, 16x."

Said another way: a child who is 10 today will, by 60, experience a year's worth of change in about 11 days.

Michael Simmons, "How Fast Will the World Change in Ten Years?" (drawing on Ray Kurzweil). Illustrative projection.

"If the rate of change outside an organization is faster than the rate of change inside, the end is near."

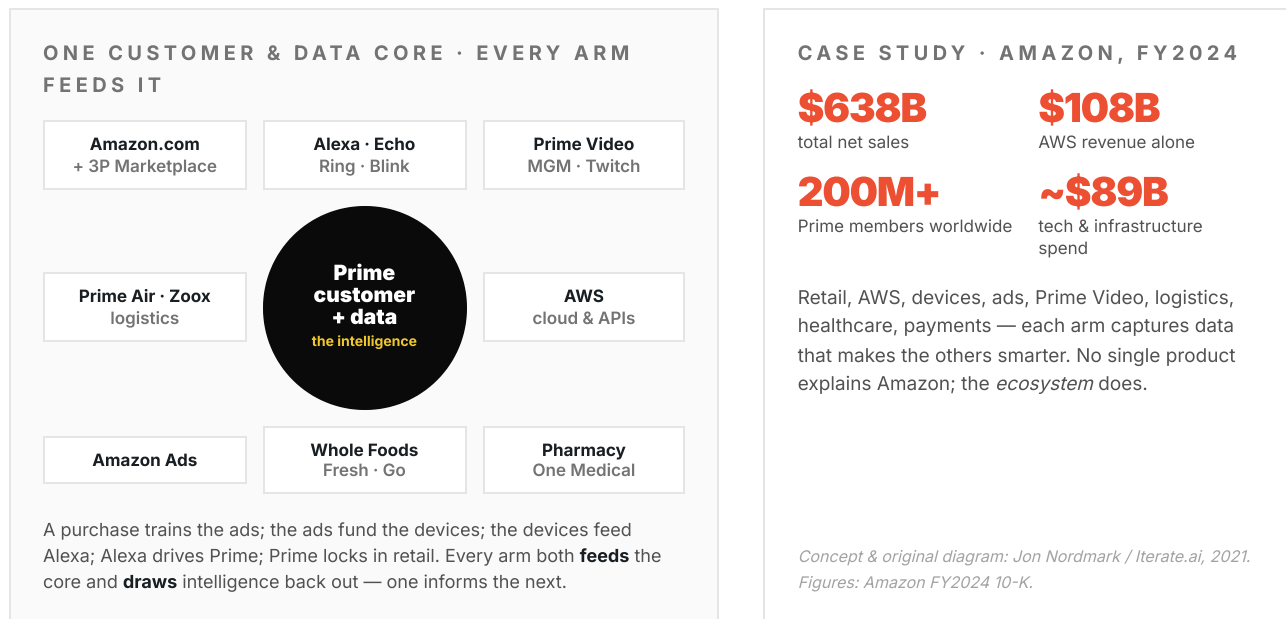
Former GE CEO **Jack Welch** championed relentless adaptability — famously advising leaders to *"change before you have to"* to avoid becoming obsolete in a shifting market. It has rarely been more true than now.



THE ECOSYSTEM ERA

Cheap intelligence everywhere builds ecosystems — not products.

When a chip that costs a few dollars gives any product a brain and a memory, every product becomes a sensor. Data gathered by one feeds the next — and the company becomes a self-reinforcing **flywheel**. Amazon was the early template: not a company, but an ecosystem.



THE STAKES

The most powerful companies of the next decade won't look like the companies we know — they'll be **ecosystems**, compounding intelligence across every product they touch.

That is the prize — and the warning. Whoever owns the intelligence layer owns the flywheel: your customers, your operations, and the collective brain of your workforce. **The question every enterprise must answer is whether that brain lives inside its own walls — or someone else's.**



PRIVATE AI: BUILD. RUN. GOVERN.

Built, run, and governed on infrastructure you control.

Iterate.ai builds, runs, and governs Private AI infrastructure for banks, hospitals, insurers, retailers, big tech, and data centers. Deployments are designed for the strictest requirements — HIPAA, SOC 2, ISO 27001, GDPR, and government standards (FedRAMP-ready). Models run on your infrastructure with zero external connectivity, so data never leaves your environment.

Powered by **Lifeboat** — our proprietary inference engine, 2–6× faster and cheaper than the best open-source stacks.

REFERENCES & SOURCES

- 01** Sanh et al., *DistilBERT* — 97% of BERT at 40% size.
- 02** Meta, *Llama 3.1 / 3.2* — 8B–405B, commercially licensed.
- 03** DeepSeek, *R1* reasoning distillation (800K samples) & *V3*.
- 04** vLLM — *PagedAttention*, continuous batching.
- 05** SGLang — *RadixAttention* prefix reuse.
- 06** Google DeepMind, *TurboQuant* (PolarQuant + QJL) — ICLR 2026.
- 07** Alibaba, *Qwen 2.5 / 3*; Mistral AI, *Mistral*.
- 08** *GPTQ* — post-training quantization.
- 09** Iterate.ai, *Lifeboat* inference-engine benchmarks (internal).
- 10** Industry GPU spot/on-demand pricing, 2022–2026.

A note on scope. This brief was prepared by a combination of Iterate.ai's engineering team, Claude (Anthropic), Perplexity, Gemini (Google), and Iterate's private *Generate* AI platform. Performance figures reflect internal Lifeboat benchmarks and publicly available research as of June 2026; results vary by workload, model, and hardware. It is intended for technical and strategic evaluation.

● Data privacy

Data never leaves your environment. Zero external connectivity by design.

● Model isolation

Full model isolation — your intelligence compounds for you, not for a shared model.

● Hardware control

Your hardware, your models, your data — on infrastructure you own and operate.




ABOUT ITERATE.AI

Private AI infrastructure for the enterprise.

Iterate.ai builds, runs, and governs private AI infrastructure for banks, hospitals, insurance companies, retailers, big tech, and datacenters. Founded in Silicon Valley and Colorado in 2015 by one team that helped invent the iPhone and another that sold \$1.65 billion worth of travel bags before its exit to Samsonite.

The four convergences in this brief — distillation, efficient inference, collapsing hardware cost, and the gains that stack — are why we serve enterprises that have decided their AI workloads, and the intelligence those workloads accumulate, belong inside their own walls.

SOME OF ITERATE'S PRODUCT FAMILY

 Lifeboat Inference acceleration. Same model, 2–6× faster and dramatically lower cost-per-token.	 Generate A no-code platform for building and running AI agents. Privately.	 AgentWatch AI governance and observability of employee LLM usage and agent behavior.
---	--	--

RECOGNITION

01 20 Hottest AI Software Companies 2025 · 2026 Channel Reseller News	05 AI 100 2023 · 2024 · 2025 KM World
02 Best Innovation in AI for Healthcare 2026 AI Tech Awards · <i>Generate for Healthcare</i>	06 Best Workplaces for Innovators 2024 Fast Company · <i>AI + Robotics</i>
03 Best AI Edge Deployment 2026 Pinnacle Awards	07 Technology of the Year 2024 InfoWorld · <i>AI/ML Models</i>
04 Best Use of AI in Healthcare 2026 Pinnacle Awards	08 Best in Business 2023 · 2024 Inc. Magazine · <i>AI + Data</i>

SERIES

**Intelligence Unshared:
The AI Sovereignty**

Paper: 5 of 6

June 2026

© 2026 Iterate.ai

HEADQUARTERS

**San Jose, CA
Denver, CO**

WEBSITE

Iterate.ai
hello@iterate.ai

PUBLICATIONS

iterate.ai/resources/white-papers
iterate.ai/resources/books