

Sombra[®]

What Works in 2026

Off-the-Shelf vs. Custom AI

Table of Contents

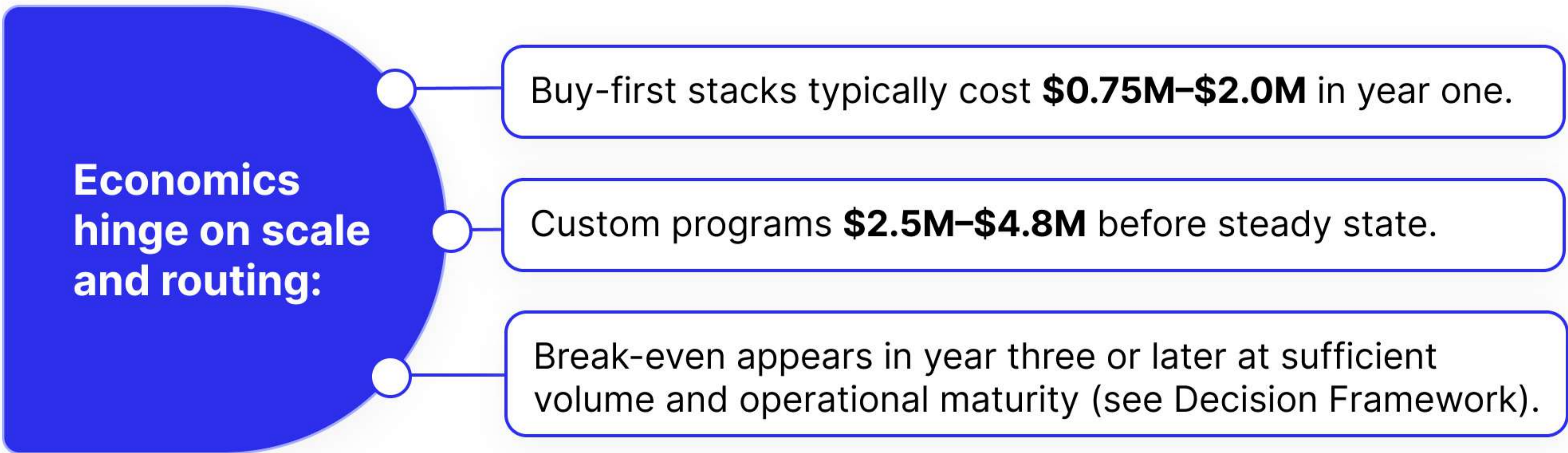
Executive Summary		3
01	Executive Market Context and Trends	5
02	Deciding How to Adopt AI in 2026	9
03	Buying Path: Speed with Guardrails	13
	Build Path: Control Through Ownership	15
	Hybrid Path: Owning the Core, Renting the Rest	17
04	Architecture Scope	21
	Regulatory & Operational Drivers	24
	Design Principles	26
	Data, MLOps, and Governance Readiness	27
	Integration	32
	KPIs and SLAs	34
05	Economics and TCO (12–24 months)	36
06	Proof-of-Concept Playbook	44
	Patterns and Case Snapshots	48
Appendixes Section		59
References		68

Executive Summary

Enterprise AI decisions in 2026 have shifted from a binary “buy or build” to a layered architecture strategy. Leaders combine packaged AI applications for speed, hosted model APIs for control points, and enabling infrastructure for retrieval, evaluation, and guardrails.

Hybrid adoption is now standard – buying where vendor capabilities meet requirements, building where proprietary control adds value. This mirrors prior multi-cloud adoption and reduces single-vendor dependency.

The time to value is measured in weeks. Off-the-shelf and platform options arrive with service commitments and governance features already in place, making them the fastest route to measurable results. Custom components are added when metrics demand tighter control, lower latency, or unique differentiation.



Regulation is now a design input. The EU AI Act, NIST AI Risk Management Framework, and ISO/IEC 42001 shape procurement and rollout plans.

Talent availability in MLOps, evaluation, and AI security still limits how far ownership can move down-stack – details in the Build Path section.

Bottom line

Start with the highest shelf that meets controls, own the middle services for portability and policy, and increase ownership only when metrics justify it.

Key Findings



Hybrid is the default.

Enterprises combine vendor-supplied capabilities with proprietary components to differentiate themselves, retaining strategic levers in-house.



Time-to-value has compressed.

Leaders expect measurable results within weeks, with “good-enough” quality and auditability as table stakes.



Economics depend on scale and latency.

Subscription models are suitable for low-volume workloads; custom builds pay off only when usage is sustained or when strict latency/compliance requirements are in place.



Governance is baked into design.

Regulatory timelines drive architecture and rollout plans, with traceability and deletion proof required from day one.



Vendor concentration remains high.

Approximately a dozen providers dominate the foundation-model market, making contract flexibility and integration control crucial.

Strategic Planning Assumptions

- **Hybrid will persist** through the planning horizon. Buyers will combine packaged workflows and hosted models with owned gateways, guardrails, retrieval, and evaluation services.
- **Procurement will demand portability** and runtime evidence. Expect SLAs/SLOs, exports for prompts/policies/embeddings, and signed deletion logs.
- **Time-to-value will stay short**, with baseline capability expected out of the box.
- **Talent scarcity will limit build velocity**, requiring staged increases in ownership only where teams can sustain evaluation harnesses, data contracts, and safety gates.
- **Regulatory cadence will drive sequencing** for capability rollout, audit artifacts, and incident-response readiness.

01

Market Context and Trends

Seven years ago, the enterprise AI market could be drawn on a single chart. **Today, it spans over 9,000 identifiable offerings**, including hyperscale cloud providers, model developers, MLOps platforms, and niche SaaS tools.

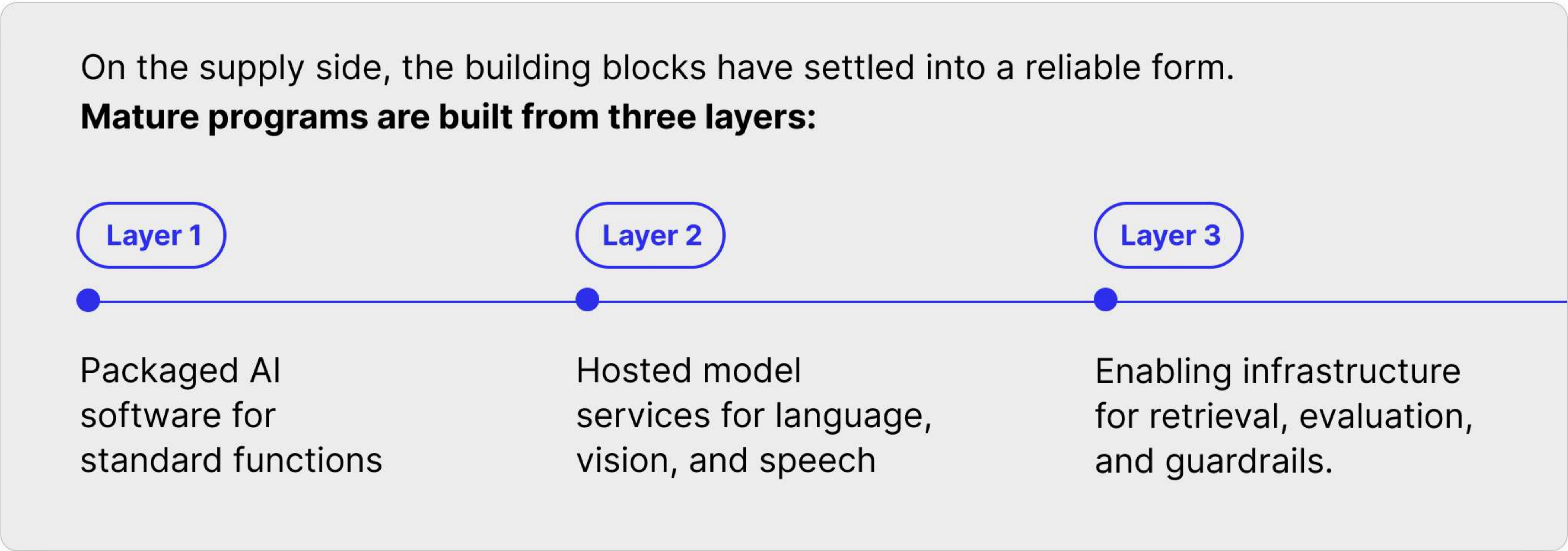
The terrain moves at **two distinct speeds**: a slow, stable core dominated by a handful of providers controlling computer and foundation models, and a fast, volatile edge where new applications appear, merge, or vanish in short cycles.

Choosing where to place spend is no longer just a procurement call; it determines how much control remains in-house, how much risk is absorbed, and how painful switching will be later.

Demand has shifted from tentative exploration to firm expectation. The pilots of 2024 and 2025 have matured into programs, and line-of-business leaders now assume time-to-value in weeks, not quarters.

*McKinsey's 2024 Global Survey found that **65% of organizations were using generative AI regularly**, nearly double the year before.*

The 2025 follow-up described a landscape of redesigned workflows, clearer governance roles, and direct, measurable impact on the bottom line.



This layered approach allows organizations to begin with available off-the-shelf capabilities and gradually bring more operations in-house as metrics justify the shift. The financial incentives favor quick starts but require more difficult decisions as operations scale up.

IDC projects that AI spending will grow by double digits through 2028, primarily driven by generative workloads.

Usage-based APIs and subscription tiers keep initial costs low; however, factors like high volumes, strict latency requirements, or regulatory compliance can shift the balance in favor of owned components and models.

Regulation now sits inside architecture diagrams rather than in footnotes.



The EU AI Act came into effect on **August 1, 2024**.



Prohibitions and AI literacy requirements began in **February 2025**, while obligations for general-purpose models will start in **August 2025**.



Most remaining requirements will become applicable in **August 2026**, with embedded high-risk systems being granted an extension until 2027.

Buyers are aligning their capability rollouts with this timeline and are investing in audit processes, model cards, and incident response systems well in advance. In addition to the EU framework, **the NIST AI Risk Management Framework and ISO/IEC 42001 have transitioned from theory to practice**, providing teams with a common language for governance and a standard for independent assurance.

People still set the limits. Scarce skills in sustained MLOps, evaluation engineering, and AI security define how far and how fast organizations can move toward deeper ownership. Buying reduces that operational load but locks more of the roadmap to vendor timelines; building removes that dependency but requires stable, well-resourced teams who can keep evaluation harnesses current, maintain data contracts, track cost telemetry, and enforce safety gates over time.

Commercial terms have matured in parallel. AI is now acquired like other strategic platforms, with requests for proposals that require evidence of quality, latency, and safety, aligned with regulatory frameworks.

Contracts specify how data is handled, how systems can be exited, and what portability guarantees exist for prompts, policies, and embeddings. Business cases are built on unit economics – the cost per assisted case, per resolved ticket, or drafted page – rather than abstract efficiency promises.

Vertical AI software has expanded rapidly, particularly in regulated sectors such as financial services, healthcare, and energy, where products often ship with compliance and workflow logic already embedded. Yet integration continues to slow deployment at scale.

Inconsistent formats, conflicting policies, and split ownership across teams keep connectors brittle and limit speed, no matter whether the approach begins with buying or building. In 2026, the decision is no longer framed as a one-time choice – it is an ongoing allocation of architectural control, adjusting as markets shift, regulations tighten, and capabilities mature.

The structure of the market reinforces one final reality: concentration at the foundation layer remains high, with **roughly a dozen providers controlling about 80% of market share**, and upgrade cycles that outpace most internal roadmaps.

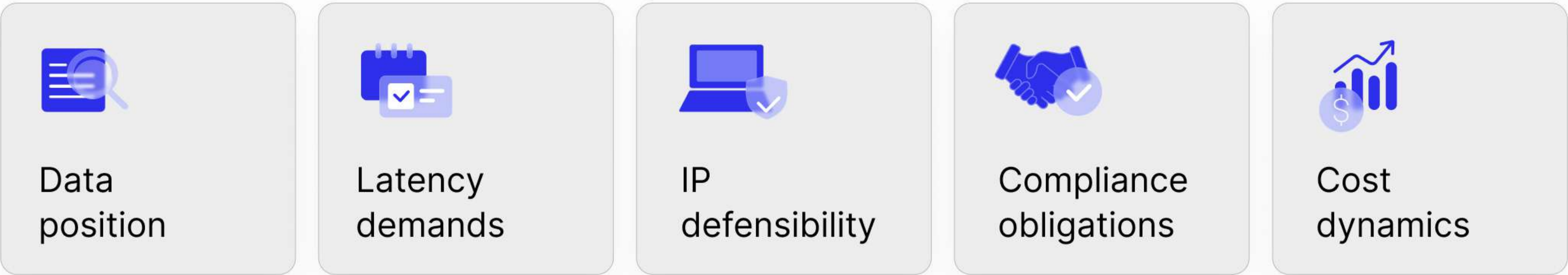
Map of AI Solutions 2026

Layer	Who's here	Market/Stats	Why it matters	Takeaway
Foundation Models & Compute	AWS Bedrock; Microsoft Azure AI; Google Cloud Vertex AI; OpenAI; Anthropic; Cohere; Mistral; Meta; xAI; Stability AI	~80% share across ~12 vendors	Core capability layer; high CapEx; upgrade cycles every 6–12 months	License at this layer; focus on contract terms and integration flexibility
MLOps, Governance & Data Infrastructure	Databricks; DataRobot; Hugging Face; Domino Data Lab; Snowflake; IBM Watsonx; OSS like MLflow, Kubeflow	~400 active products; >25% yearly churn	Maintains compliance, auditability; heavy feature overlap	Keep API-level flexibility so you can replace a tool without breaking workflows
Vertical AI SaaS	Healthcare: Tempus, PathAI, Olive AI; Finance: Zest AI, Darktrace, Ayasdi; Retail/Manufacturing: Covariant, Focal Systems, Uptake; Legal: Casetext, Harvey	5,000+ products; fastest adoption in healthcare, finance, manufacturing	Domain ROI; integration maturity varies	Keep API-level flexibility so you can replace a tool without breaking workflows
General Productivity AI	Jasper; Copy.ai; Otter.ai; Grammarly; Reclaim.ai; Superhuman	Thousands; high redundancy	Useful for broad tasks	Avoid overpaying for incremental features; consolidation is likely

02

Deciding How to Adopt AI in 2026

By 2026, the decision to buy or build AI will no longer be made on a single metric. It rests on a set of structural signals:



These signals work together rather than in isolation. Seeing these in context reveals the most sustainable and cost-effective path.


A strong **proprietary data advantage** remains the single most decisive factor. Enterprises with exclusive, large-scale, and clean data streams can extract far greater value from custom development. In these cases, the AI system becomes more than a productivity tool – **it is a compounding asset**. Every use enhances the moat: data increases in value, models improve, and competitors find it harder to replicate. For such organizations, building is not just defensible; it is **strategically essential**.

When that advantage is absent – in sectors with commoditized datasets or heavy restrictions on data use – the equation flips. **Buying becomes the pragmatic choice**, especially when vendors already offer domain-specific models trained on relevant data and certified for compliance. In healthcare, for example, vendor readiness can outpace internal capability by years, making off-the-shelf adoption the fastest route to measurable outcomes.


Latency can tip the balance as well. In algorithmic trading, real-time fraud detection, or industrial automation, millisecond-level responsiveness may be the difference between profit and loss, or safety and failure. Although vendor platforms have improved, shared infrastructure still carries unavoidable latency. In these cases, **owning the serving stack is critical** to meeting operational guarantees.

Intellectual property is another lever. If AI capabilities are intended to differentiate in the market— **such as a proprietary risk model, an advanced recommendation engine, or specialized generative output** — dependence on an external provider risks both reliance and a loss of exclusivity. Here, internal builds or co-development with joint-IP agreements can secure long-term control, even at higher short-term cost.

In 2026, the costs for different deployment options are clearly defined.



For enterprise-scale, off-the-shelf solutions, which typically involve licensing, integration, and minimal customization, the first-year costs range from \$750,000 to \$2 million.



On the other hand, custom builds start at \$2.5 million and can go up to \$4.8 million in the first year. These custom projects require significant upfront investment in talent, infrastructure, and MLOps tools.

Break-even usually arrives in the third year or later, and only if usage is sustained and the architecture avoids major rework. For low-volume or exploratory workloads, subscription pricing remains more economical.

Compliance and governance can shift the decision in either direction. Vendor certifications can simplify adoption, but if jurisdiction-specific requirements can't be met externally, **in-house control becomes the safer route**. In regulated environments, the risk calculation often outweighs the speed advantage of buying.

The most resilient strategies treat these variables as **a weighted framework, not a binary switch**. Many start with vendor models for general capabilities, adding proprietary modules for high-impact areas. Others do the reverse – building the core while outsourcing peripheral needs. The strongest results come from aligning business priorities, technical realities, and projected economics into a coherent architecture that can flex over time.

The Quick Sneak Peek at Your Probable Paths

Term	Definition	Typical Examples	Primary Strength	Primary Weakness
Off-the-Shelf AI (Application)	Packaged AI apps delivered as SaaS; configurable, not deeply customizable	AI copilots in CRM/ITSM, contact-center AI, document AI suites	Fast time-to-value, low ops burden	Limited control, vendor roadmap risk
Off-the-Shelf AI (Model Service/API)	Hosted foundation or task models accessed via API; pay-per-use	General LLMs, vision/speech APIs, retrieval APIs	Breadth, scalability, predictable SLOs	Token costs, limited specialization
Custom AI (Fine-Tuned Model)	Adapting a base model on proprietary data for a domain task	Claims adjudication, underwriting notes, clinical coding	Higher task accuracy, moat via data	Ongoing training cost, eval burden
Custom AI (Domain-Specific/Small Model)	Purpose-built model with narrow scope and efficient serving	Low-latency edge, safety-critical text/vision	Latency, cost control, on-prem viability	Higher build complexity
Agentic System	Multi-tool workflow with planning/critique/act loops under policies	Complex L3 support, financial ops reconciliation	Autonomy on multi-step work	Harder assurance, safety gating needed
RAG (Retrieval-Augmented Generation)	Combine model generation with enterprise retrieval for grounded outputs	K-bases, procedures, product docs	Freshness, explainability, lower fine-tune need	Retrieval quality is a hard dependency
Evaluation Harness	Automated quality, safety, latency, and cost measurement pipeline	LLM-as-judge, golden sets, offline/online A/B	Objective decisioning	Requires ongoing curation
Orchestration Layer	Abstraction controlling model/tool selection, prompts, routing, guardrails	Gateways, routers, policy engines	Portability and safety controls	Added integration overhead

03

Buying Path, Build Path & Hybrid Path

Buying Path: Speed with Guardrails

*For many organizations, **buying remains the fastest and lowest-friction route to operational AI**. In 2026, that market spans three interconnected layers: AI SaaS applications, hosted model APIs, and cloud AI platforms. The right mix delivers speed while avoiding deep lock-in.*

AI SaaS tools embed capabilities directly into the systems employees already use:



Microsoft Copilot
in Office



Google Gemini
in Workspace



Salesforce
Einstein in CRM



ServiceNow Now
Assist in IT workflows

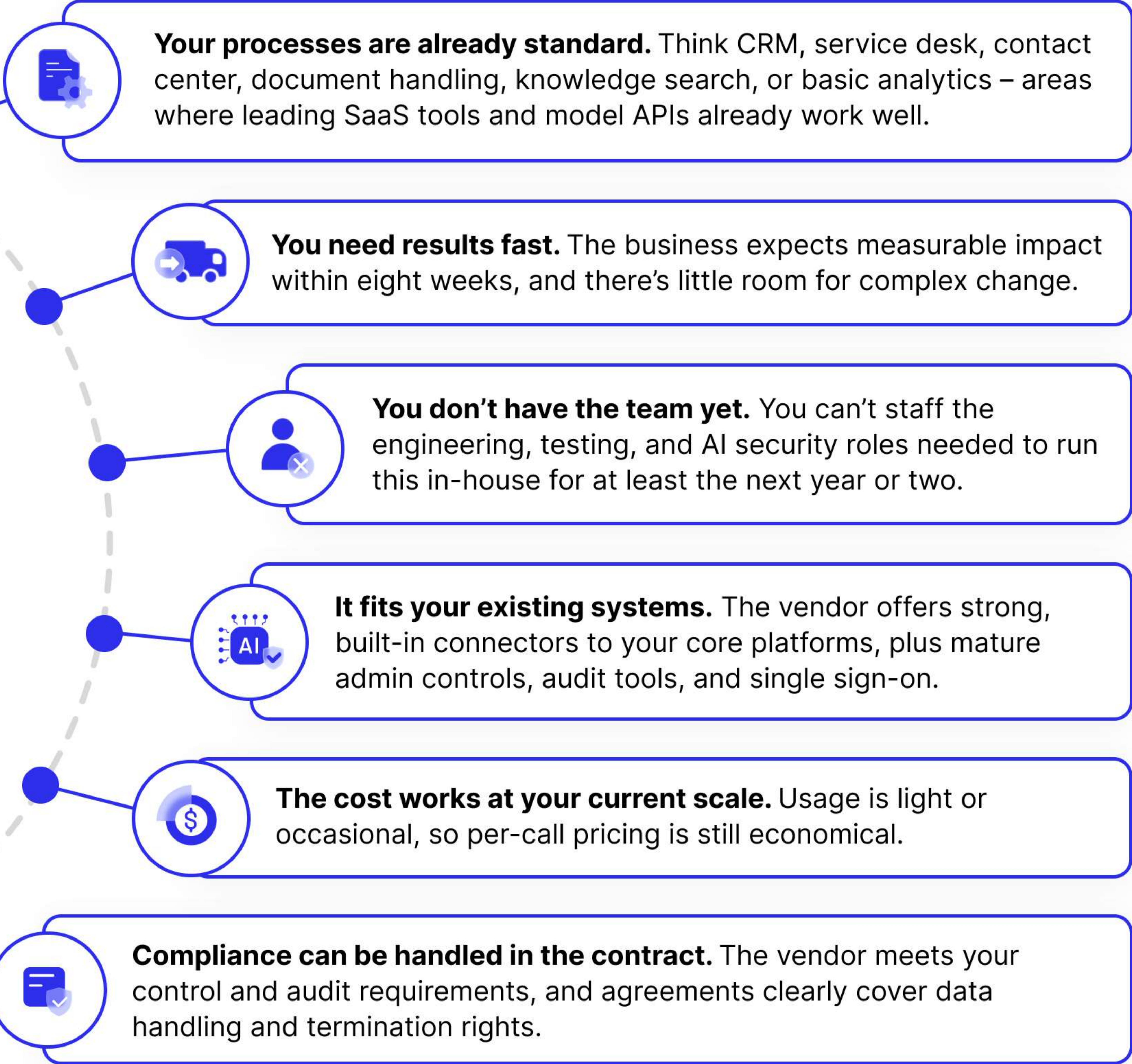
Because they sit on existing identity and permission frameworks, deployment can happen in weeks. The trade-off is strategic: **per-seat costs can escalate quickly, vendor roadmaps dictate your options, and customization depends entirely on extension points the vendor chooses to expose.**

Hosted model APIs offer a middle ground. They let enterprises retain control over their own interface, logic, and policies, while invoking external models via secure endpoints. Over the past year, major providers have aligned on a baseline privacy stance: **prompts and outputs are not used for model training by default**, and retention settings are clearly documented. This change removes one of the biggest historical blockers for legal and compliance teams.

Cloud AI platforms such as Azure AI Studio, Google Vertex AI, and AWS Bedrock offer curated model catalogs with **governance, networking, and audit controls already built in**. They now provide private connectivity as a standard feature, allowing inference to remain inside private network paths – a non-negotiable in regulated industries. For many, cloud familiarity drives platform choice because security, monitoring, and incident response practices transfer directly.

The buying path works best when leaders **know where each layer's boundaries lie** and integrate them into their own governance framework. SaaS delivers ready workflows, APIs enable controlled integration, and platforms centralize oversight. Each has weaknesses – SaaS limits deep customization, APIs require platform skill sets, and platforms demand internal product work to produce business outcomes. Price volatility remains a shared risk, making it critical to negotiate **renewal caps, usage thresholds, and exit clauses at the start.**

Triggers That Favor Buy



Signs you're here

You can hit your goals by configuring what's already available rather than building from scratch. You're comfortable following the vendor's roadmap, and your risk appetite leans toward external SLAs instead of internal ownership.

Build Path: Control Through Ownership

*Building AI means taking ownership of the full product lifecycle – from **design and delivery to compliance and continuous improvement**. Executed with discipline, it places quality, cost, and risk controls precisely where the business requires them. Done casually, it drains resources and duplicates effort.*

A sustainable build program starts with a dedicated platform team anchored by a **lead responsible for a shared gateway, guardrails, enterprise retrieval, and automated evaluation**. Supporting that lead are MLOps engineers, data engineers, and security/compliance specialists who ensure **controls such as SSO, private connectivity, and deletion proofs operate continuously**. A stable, skilled team – typically six to eight full-time roles – is the foundation for a viable build strategy.

The architecture is strict. All requests pass through a **gateway** for authentication, policy enforcement, and cost management. Retrieval runs before inference to pull enterprise knowledge under correct permissions. **Model routing sends easy cases to smaller, cheaper models while escalating complex ones**, and every action generates signed logs for export to monitoring systems. The evaluation harness runs automatically on any change, producing measurable quality, safety, latency, and cost reports for decision-makers.

Rollouts begin deliberately.



The first two weeks set KPIs, define data scope, and establish the gateway and logging.



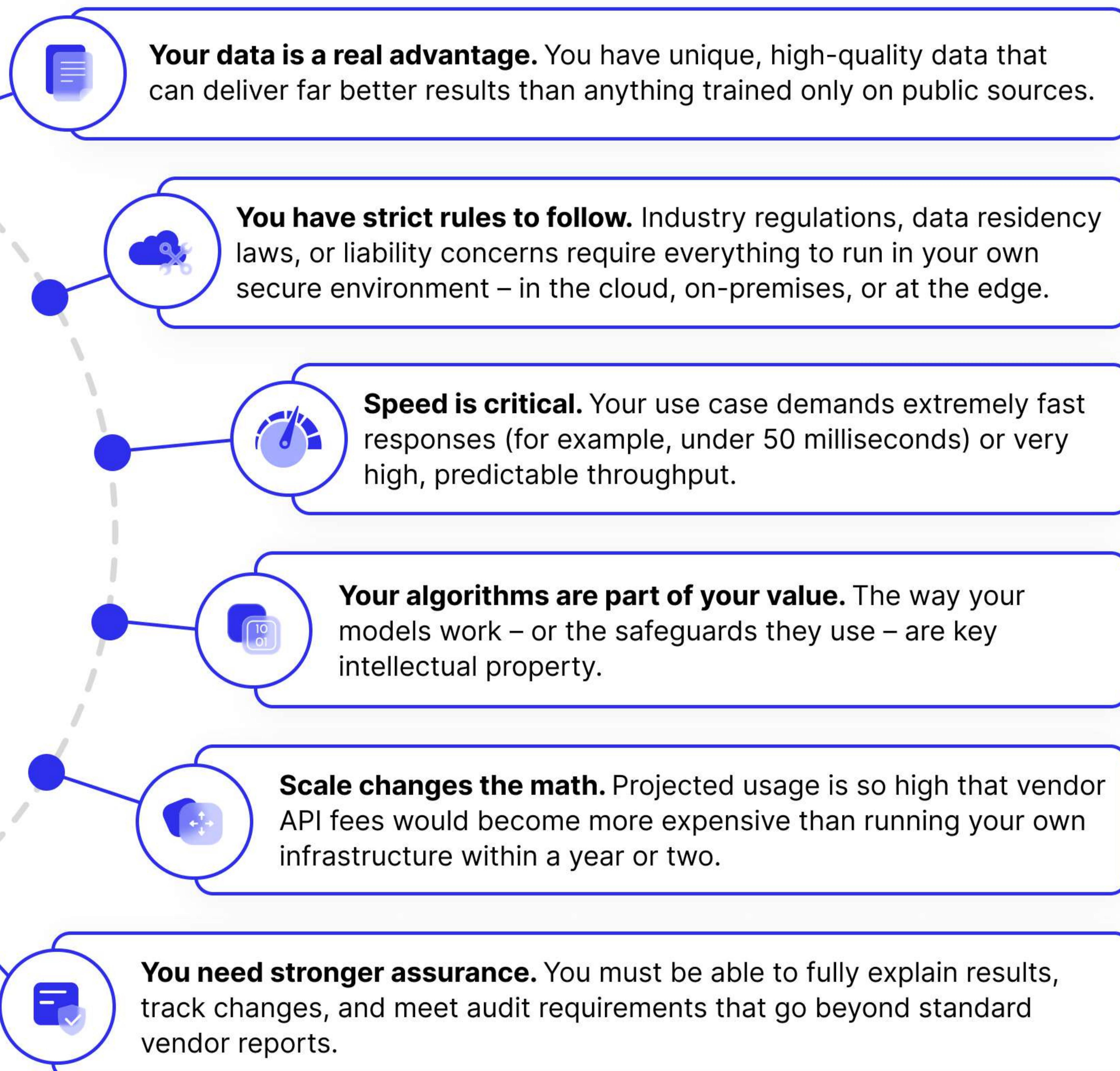
The next month tunes ingestion and retrieval quality, followed by complexity such as cost telemetry and live incident drills.



By day ninety, the system should be capable of producing auditable evidence of readiness for scale.

This approach gives leaders the levers to adjust cost vs. capability, shift workloads between models, and adapt to changing regulatory or operational demands. It is **the difference between simply running a model and running an accountable, adaptable AI service.**

Triggers That Favor Buy



Rule of thumb

If three or more of these points apply, lean toward building in-house (or a hybrid model with more ownership) – and budget for the teams and processes to keep it secure, well-tested, and reliable over time.

Hybrid Path: Owning the Core, Renting the Rest

By 2026, hybrid adoption has become the default architecture for enterprises that need both speed and control.

*The model is simple: **own the governance, retrieval, and logging layers** – rent the models and applications that can change without destabilizing the core.*

Several market shifts have locked in this pattern. Private connectivity to hosted models is now routine, removing the need to expose sensitive traffic to the public internet. Regulatory timelines push buyers to maintain **central evidence and control points**. Widely shared reference designs for retrieval-augmented generation and model routing make in-house quality improvements both practical and repeatable.

In a strong hybrid setup, the enterprise controls **identity and access at the gateway, runtime guardrails, retrieval over enterprise content with permission enforcement, a living evaluation suite, and portability contracts with vendors**. Everything else – foundation model endpoints, narrow SaaS workflows – remains externally sourced.

The advantage is flexibility. **You can start with vendor tools to prove value, then gradually shift critical or sensitive workloads to your own controlled stack**. Costs stay manageable by routing simpler cases to cheaper models, while keeping the capacity to swap vendors or reconfigure architectures when economics or compliance rules change.

The main pitfalls come from losing sight of where control lies:



Ceding too much to vendors

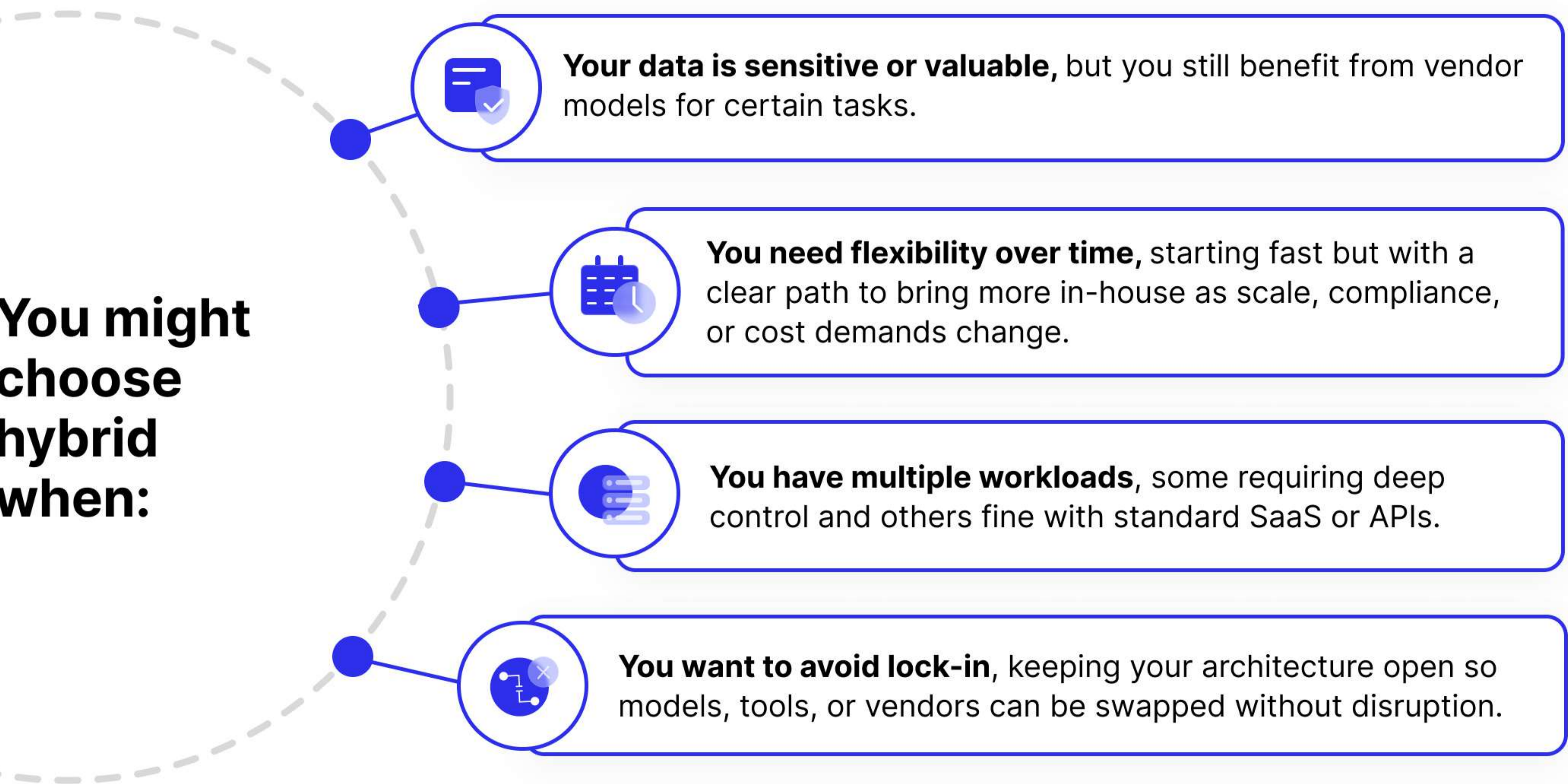


Overbuilding features that add no unique value, or allowing direct-to-model integrations that bypass governance.





When run with discipline, hybrid **keeps the strategic levers in-house while exploiting vendor progress at the edges**, ensuring the organization can adapt without starting from scratch.

When Hybrid Makes Sense

Hybrid works when you want the speed and maturity of vendor solutions but still need to keep the strategic levers – data, governance, and flexibility – under your own control. It’s the go-to model when neither a pure “buy” nor a full “build” fits the whole picture.

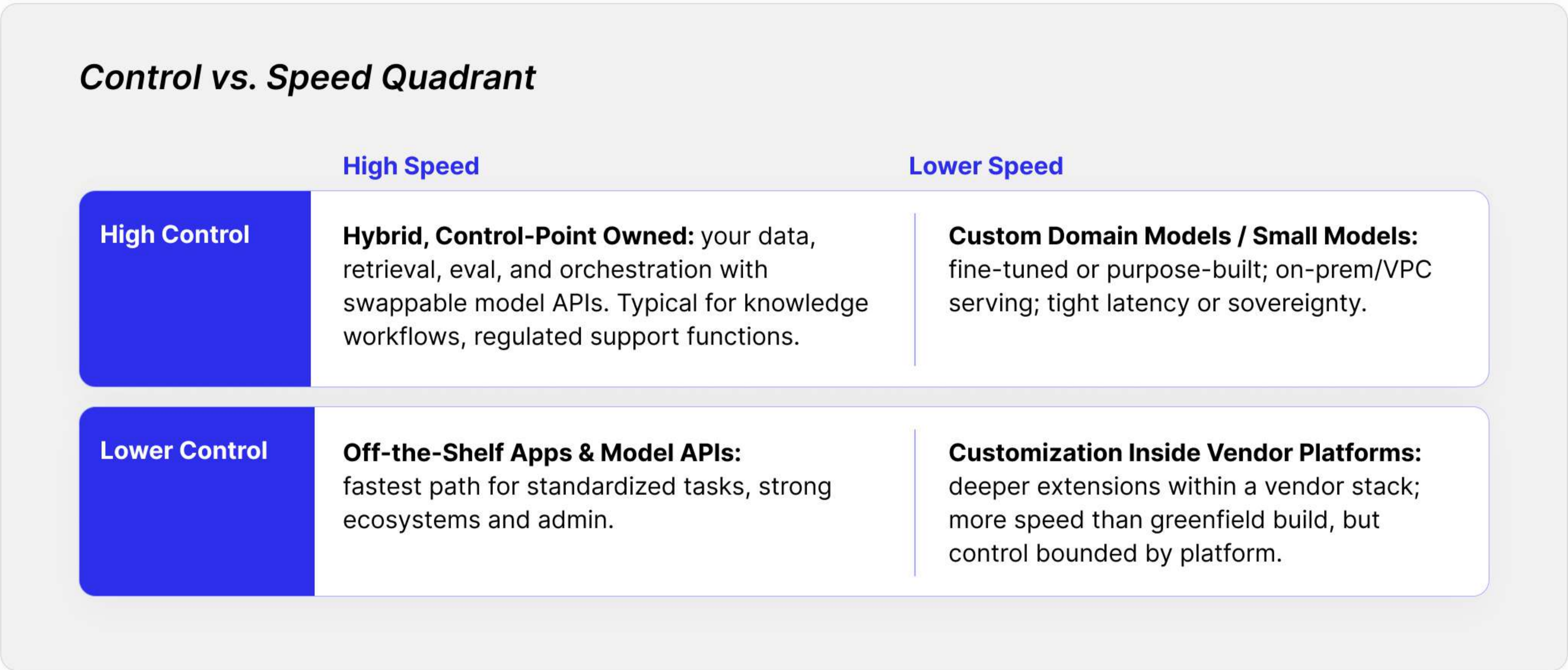


Typical hybrid patterns include:

 Data-Owned RAG: You manage enterprise retrieval and policies; vendors provide interchangeable model endpoints.	 Fine-Tuned Edge Models: Small, specialized models for core tasks, with overflow routed to vendor APIs.	 Orchestrated Agents: Your team controls the workflow engine, tools, and guardrails; vendors supply modular capabilities.	 Buy Now, Build Later: SaaS for quick wins, then gradually replace pieces with owned components as needed.
---	---	---	--

Key controls to keep in-house:

Data contracts, retrieval quality, evaluation metrics and golden sets, telemetry on performance and cost, policy guardrails, and the routing/orchestration layer.



How to use the quadrant

Place each use case by required control (assurance, sovereignty, latency, portability) and desired speed(time-to-value, change capacity). The quadrant rarely yields a single answer for all workloads; it maps a portfolio.

04

From Theory to Practice:

Building Blocks for
AI Implementation

Architecture Scope

This section defines how the buy-or-build decision translates into concrete components, ownership boundaries, and control points. It's where strategy becomes engineering reality.

“






An AI system isn't one decision, it's dozens – each defining what you own, what you rent, and what you trust.

Yurii Nakonechnyi,
Sombra CTO



Think of the system as a series of gates and highways: data flows in, is transformed, queried, enriched, passed through decision layers, and finally returned to the user in a controlled, auditable way.

At a high level, a **2026 reference architecture** for enterprise AI includes:

-  **Data ingress and governance layer** – connectors to structured, semi-structured, and unstructured sources; automated ingestion pipelines; real-time validation; governance hooks.
-  **Retrieval and orchestration core** – RAG pipelines tuned for latency and relevance, routing logic that selects models based on cost, capability, and compliance context.
-  **Evaluation harness** – continuous A/B testing of models, regression detection, compliance audit logs.
-  **Serving and integration layer** – APIs, SDKs, and direct integration into business applications, with latency controls and fallbacks.
-  **Control points** – admin dashboards, KPI dashboards, human-in-the-loop interfaces, and model override triggers.

The boundaries between owned and vendor-supplied components matter. In regulated industries, the architecture must allow for **quick swaps** – replacing a vendor's model or retrieval stack without breaking downstream workflows. That's not just risk management; it's operational survival.

Key Architectural Building Blocks

The **core building blocks** are not static features; they're living systems that will change with model performance, regulatory pressure, and business priorities. In 2026, the most resilient AI architectures are modular, observable, and replaceable at the component level.

Retrieval-Augmented Generation

RAG has moved from proof-of-concept novelty to **mandatory infrastructure** for most enterprise AI use cases. In 2023–2024, the focus was on connecting a model to a knowledge base; in 2025–2026, the challenge is **governance and adaptability**.

Modern RAG stacks:



Use **domain-specific embeddings** rather than generic ones, improving semantic precision in finance, healthcare, and manufacturing contexts.



Run **multi-hop retrieval**, where queries are refined iteratively for deeper context coverage.



Integrate **sensitivity filters** that redact PII or regulated data before reaching the model, to satisfy compliance requirements like the EU AI Act's Article 10 (data governance).



We see RAG as the guardrail that keeps LLMs in the lane of truth. Without it, hallucination risk doesn't just go up – it becomes unmanageable at scale.

Yurii Nakonechnyi,
Sombra CTO

Model Routing

No single model will dominate every workload. Enterprises are already **routing** between multiple models – foundation models, fine-tuned vertical models, and open-weight local deployments.



Cost-performance trade-offs are now monitored in real time; expensive calls are reserved for high-value transactions.



Compliance-aware routing ensures models trained on non-EU data aren't used in EU-regulated workflows, reducing audit exposure.



Latency-aware routing keeps conversational systems responsive by sending simpler requests to faster models.

The winning architectures are **multi-model by design**, not by accident. Your files reflect the economics behind this shift.

In a typical scenario at **50M tokens/month**, a buy-first path totals roughly **\$2.1M over two years**, while a custom platform with owned infrastructure and team lands **around \$5.9M**.

Build becomes cheaper only at much higher volumes (your model shows a break-even near **~180 M tokens/month**), or when you need fine-tuning that the market doesn't offer.

Evaluation Harness

The days of once-a-quarter model testing are over. Continuous evaluation frameworks now:

1

Benchmark model performance **daily** against curated test sets.

2

Detect **concept drift** – subtle shifts in model output caused by upstream changes in vendor APIs.

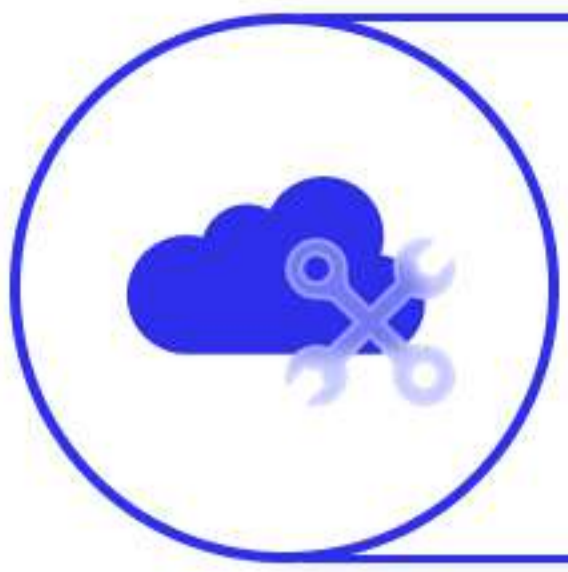
3

Feed error signals directly into retraining or prompt-adjustment pipelines.

Operational AI without evaluation is a black box – and black boxes don't survive in regulated or mission-critical environments.

Private Endpoints

Private endpoints are no longer a “nice to have.” With data residency laws tightening, they're becoming standard in enterprise AI contracts. They ensure:



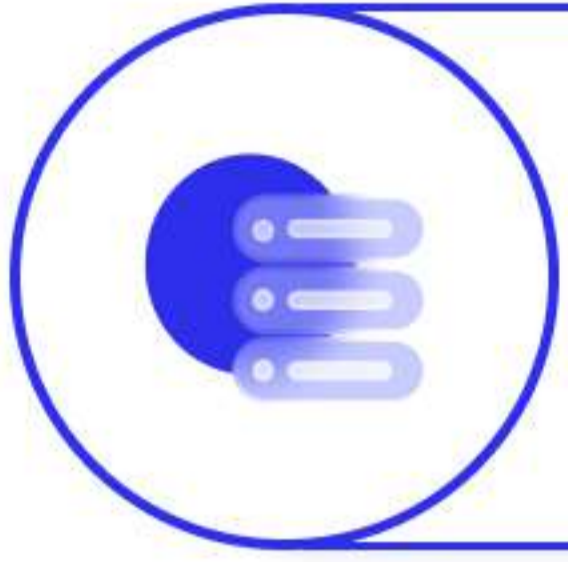
Traffic isolation –

no data leaves the corporate network or approved sovereign cloud.



Consistent performance –

avoiding shared public endpoints with variable latency.



Auditability –

logging requests and responses for internal review.

By 2026, we expect most serious deployments in finance, healthcare, and defense to treat **public API calls as a temporary phase** – the goal is private, controlled, and monitored execution.

Regulatory & Operational Drivers

Architectural choices in 2026 would be about **survivability under scrutiny**. Regulations, procurement cycles, and operational risks now shape the AI stack as much as model accuracy.

EU AI Act – Enforcement Calendar

The EU AI Act has moved from theory to clockwork. By mid-2025, high-risk AI systems (Annex III) must demonstrate compliance with requirements for data governance, transparency, human oversight, and robustness.

By 2026, enforcement widens to all regulated AI categories.
For architecture, this means:



Full data lineage tracking – knowing exactly what datasets contributed to a model's output.



Human-in-the-loop checkpoints in workflows that affect safety, rights, or financial decisions.



Risk classification tags on each model, determining how and where it can be used.

Failing these isn't just a fine – it can mean a **forced shutdown of production AI systems** in the EU.

NIST AI Risk Management Framework (RMF)

In the U.S., the NIST AI RMF has become the de facto playbook for federal agencies and regulated industries. Procurement now prioritizes vendors that demonstrate alignment with its four core functions: Govern, Map, Measure, and Manage, integrated into their architecture.

Architectural impacts include:

1

Pre-deployment “red-team” environments for adversarial testing.

2

Metric-driven bias detection pipelines connected to evaluation harnesses.

3

Escalation playbooks triggered automatically when KPIs breach tolerance.

Operational Realities – Vendor Churn and Economic Pressure

2024–2025 proved that the AI market is volatile. Vendors folded, pivoted, or restricted API usage without notice. The average enterprise AI program has seen two major vendor changes in the last 18 months.

Architecture must assume:

-  **Hot-swap capability** – the ability to replace a model or service in days, not quarters.
-  **Dual-vendor redundancy** in mission-critical paths.
-  **Contract clauses for exit and migration** as part of procurement.

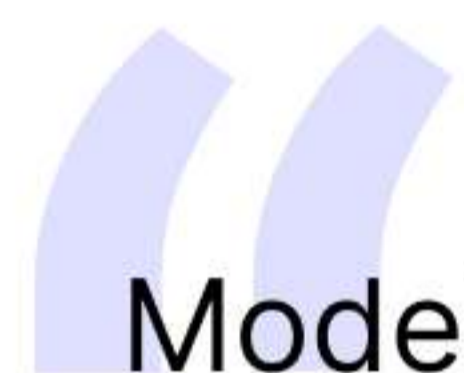
Economic pressure is also reshaping AI spending. Budgets are shifting from experimental pilots to **production systems with ROI in under 12 months**. That drives demand for architectures that are **composable, KPI-driven, and vendor-portable**.

Design Principles

The most resilient AI architectures share one trait – **they're built for change**. At Sombra, we see this as three core principles.

Own the Middle

Control the orchestration layer between models, data, and delivery endpoints. This is where **routing, context injection, evaluation, and governance** live. Owning the middle means you can switch providers, update evaluation logic, and adapt to new compliance rules without ripping apart the stack.



Model choice is temporary; orchestration is permanent. If you don't own the decision layer, you don't own the system.

Yurii Nakonechnyi,
Sombra CTO

Composability by Default

In practice, this means adopting open standards, modular APIs, and containerized microservices. It's not just a developer preference; it's survival in a market where half the vendors from 2024 won't exist in 2026.

Executives care about how fast the solution **feels**. Two practices dominate: **streaming** (show the first output quickly) and work **splitting** (fast, short responses for interactive steps; batch the heavy work).

Retrieval adds some delay but improves correctness; users accept that trade when they see citations.

The result is stable service-desk-like targets: sub-second for small actions, a couple of seconds for longer generations, and asynchronous jobs for heavy back-office runs.

Every component – from RAG modules to vector stores to private endpoints – should be **swappable without downtime**.

KPI-Driven Escalation

Architects should know when they're failing. **Define leading KPIs** (latency, relevance, safety metrics) and bind them to **automated escalation paths**. If a KPI exceeds the threshold, the system routes to a fallback model, initiates a human review, or disables the affected features until the issue is resolved.

This turns governance from a paper exercise into **live operational control**.

What this means for your roadmap

In regulated, high-stakes domains, compliance deadlines and vendor churn are not rare events – they're constants. Owning the middle, building for composability, and embedding KPI-triggered governance make the difference between a system that adapts in days and one that collapses under its own dependencies.

Data, MLOps, and Governance Readiness

Readiness determines the path. Weak data management, fragmented evaluation processes, and controls that exist only in documentation make a custom build slow, expensive, and risky. In this case, the most practical route is to adopt a working solution, deliver measurable results, and use the time saved to strengthen internal capabilities. Where these capabilities are already in place, a hybrid or fully custom build allows greater control over quality, risk, and cost.

In 2026, the **buy-versus-build choice depends on the strength of current foundations.**

Data Foundation

The priority is ensuring the system can access the correct sources, provide citations, and respect access boundaries. This requires a clear inventory of wikis, policies, tickets, specifications, and knowledge bases in scope. Update frequency must be known. Sensitive or restricted records must be tagged and enforced at retrieval. Many accuracy failures originate in outdated or poorly governed sources.

The solution is structured data cleaning, consistent preprocessing, and permission checks at retrieval. In cost-benefit models, improving retrieval and ingestion often produces a larger accuracy gain than replacing the model vendor, while requiring less investment.

MLOps and Product Hygiene

Maintain a registry of all models, prompts, and policies deployed in production, with named approvers. Keep a small but representative golden set of test cases reflecting actual workloads. Every change to content, prompts, or models is tested against this set before release. Post-launch monitoring compares live traffic with the baseline to detect quality drops early.

This approach supports continuous operation and allows vendor performance to be evaluated against internal benchmarks.

Operational Governance

Governance functions as an active layer in architecture. Sensitive inputs are masked before reaching a model. High-risk outputs are blocked or sent for human review. Major actions create a detailed audit trail that includes the source documents, model version, prompt, and policy permitting the action. Model cards and run books document intended use, limitations, and known failure modes. Including these controls early reduces the risk of compliance issues and prevents delays in deployment.

Your Data Trust/Readiness checklist

The readiness checklist measures whether your data, MLOps, and governance capabilities can support a production AI system. Each item covers a stage in the lifecycle, from data source discovery to model retirement. The focus is on reliability, compliance, and cost control.

One gap can be managed if mitigated. Multiple gaps mean the environment is not stable enough for scaling. In that case, fix the foundations before adding model complexity. Utilizing a checklist for each project establishes a common quality standard and provides a clear, evidence-based decision for launch.

Use this to decide if a use case is “go,” “fix,” or “defer.” If two or more items are “no,” you do not have launch-ready governance—fix those before adding model complexity

-  **Sources & lineage:** We can list every source, owner, refresh cadence, and last ingest; we can show lineage for any answer.
-  **Access at retrieval:** Permissions from source systems are enforced at query time; citations never expose restricted content.
-  **Golden set:** We have 50–200 representative questions with expected answers/citations; it's refreshed quarterly by business owners.
-  **Registry:** Model, prompt, and policy versions are tracked with approvals; roll-back takes minutes, not days.
-  **Evaluation harness:** Every change runs offline tests; results show quality, safety, latency, and cost.
-  **Online guardrails:** Redaction, safety filters, and approval gates are active and logged.
-  **Monitoring & drift:** Live dashboards show accuracy, data and prediction drift, costs per task; alerts page the owner on thresholds. (If you use cloud platforms, your minimum bar is their built-in drift monitors.)
-  **Audit pack:** Model cards, run books, deletion/retention procedures, and export formats are up to date; logs stream to SIEM.
-  **Human oversight:** Clear points where a person reviews or approves high-impact actions; fallbacks defined.
-  **Exit plan:** We've practiced exporting prompts, policies, embeddings, and logs; deletion and crypto-erase are verifiable.

The same capabilities, split three ways

The capability table shows how responsibility shifts across buy, hybrid, and build approaches. Each row covers a function that shapes quality, risk, and cost. In a buy scenario, the vendor handles most operations, and your role is to configure, verify, and monitor.

A hybrid approach enables you to oversee critical elements such as ingestion, retrieval, and governance, while vendors handle the other components. In contrast, a build approach places full responsibility in your hands. This mapping aligns your readiness with the operational requirements of each option, making your decision more defensible.

Capability (kept short; only what changes the decision)	Buy (vendor-provided, you verify)	Hybrid (you own the middle)	Build (you own end-to-end)
Data quality & access at answer time	Configure source connectors + permissions; demand citations and per-user access checks	You run ingestion/retrieval with your access rules; vendors plug into it	Same as Hybrid, plus your own indexing/storage SLAs
Evaluation harness	Vendor dashboards; you bring a small golden set and require exports	You run offline tests + online A/B for all vendors/models via your gateway	Same as Hybrid, with your own test data mgmt and regression suites
Gateway & guardrails	Vendor's	Yours (routing, retries, cost telemetry, redaction/policy filters)	Yours
Registry (models/prompts/policies)	Vendor's change log; you require version exports	Yours (single source of truth; rollback in minutes)	Yours
Drift & incidents	Vendor monitoring; you require alerts & runbooks	Your monitors page owners; vendor signals feed your SIEM	Your monitors only
Compliance artifacts	Vendor model cards, data handling, deletion proofs	You standardize artifacts; vendors must match them	You produce all artifacts yourself

Security, Compliance, and Regulatory Mapping

Security and compliance requirements influence architectural choices from the first design session. The safe rule is to own controls that auditors and regulators will expect to see running in real time. These include identity, data handling, logging, and deletion. Outsource only where a vendor can meet those controls faster than internal teams can.

The threat model is familiar. Data leakage is minimized when all requests are routed through private network paths, and sensitive information is masked before being processed by the model.

Unauthorized actions are prevented through role checks and approvals that occur in real time. Accountability is ensured by maintaining signed logs that link each output to its source material, the model version, the prompt used, and the governing policy.

Most enterprise AI risk falls into three categories: data leakage, unauthorized actions, and weak accountability.

If a vendor cannot demonstrate that these controls are in place and functioning with your data, achieving compliance will be impossible.

Private connectivity closes a long-standing gap between control and capability. AWS Bedrock offers PrivateLink from customer VPCs. Azure AI and Azure OpenAI provide private endpoints inside virtual networks. Google Vertex AI supports Private Service Connect. These features allow regulated workloads to run in vendor environments without exposing traffic to the public internet.

The buy, hybrid, and build options each map cleanly to these requirements. Buy works when a product already demonstrates SSO, private networking, filtering, signed logs, and deletion. Hybrid works when you keep identity, logging, and guardrails at your gateway, connect vendors or model APIs behind it, and pull their logs into your system. Build is justified when policy or scale requires fully controlled infrastructure.

Three deliverables make compliance tangible:



A control catalog listing each safeguard and its runtime enforcement method.



A compliance matrix linking each control to EU AI Act provisions, NIST RMF, and ISO/IEC 42001 clauses.



A risk register naming owners, mitigations, and test evidence for priority risks such as leakage, unsafe actions, and missing logs.

Integration and Change Overhead

The ability to absorb change determines how quickly AI moves from prototype to daily use. Successful programs keep integration simple. Most standardize on single sign-on, a few stable connectors, and a single gateway for telemetry, rate limits, and retries. Complex or custom integrations slow every future update.

When considering a product for purchase, the primary concern is whether it can seamlessly integrate into existing systems without requiring additional effort. It should be able to connect with work environments such as CRM, ITSM, content repositories, and email platforms, as well as send events to the logging and ticketing tools.

For building, achieving these same integration goals necessitates engineering work to create gateways and establish data contracts. The hybrid approach uses a central gateway to handle identity and policy checks while allowing for the integration of products or APIs with minimal customization.

Integration affects long-term cost. Each week saved on setup is a week spent improving quality or adoption.

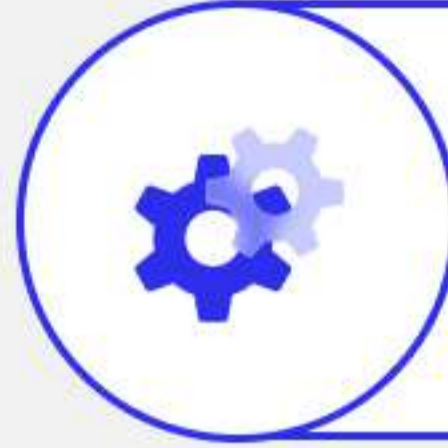
Four operational loops must be ready:



Identity – SSO with role-based access.



Observability – Unified monitoring for quality, speed, and unit cost, feeding alerts and BI.



Incident – Named owners for triage, engineering escalation timelines, and rollback methods.



FinOps – Usage tracked to cost units like tokens, GPU hours, and storage for accurate forecasting.

These loops apply regardless of the adoption model. The difference is whether you adopt the vendor's version or create your own. Buy lets you use vendor loops immediately. Hybrid and Build require internal loops from the start.

People and training are as important as integration. Rollouts should start in one unit, measure outcomes, and expand when results are proven. Training should cover when to use the assistant, how to identify errors, and how to escalate. Support channels must be active on launch day with clear ownership and response times.

A simple estimator helps decide readiness.

Count the systems to connect, the identity domains, and the workflows in scope for the first release.

High counts with a small platform team favor Buy or Hybrid. Low counts with existing gateways and monitoring make Build viable.

Three planning documents ensure delivery alignment: a rollout plan identifying the first business unit, metrics, and expansion criteria; a training plan with designated owners and timelines; and a living integration estimate that is updated following a brief technical test.

The same gateway, guardrails, retrieval, evaluation, and logging layers make future integration easier and reduce the cost of change. They also provide end-to-end evidence for regulators, aligned with the EU timeline and NIST/ISO expectations.

KPIs and SLAs: Operationalizing AI in 2026

In 2026, enterprises no longer treat AI metrics as background dashboards. They are contractual control surfaces that determine whether a system can run at scale, survive audit, and justify its economics. KPIs define performance; SLAs bind it to enforcement. Together they align technical operations with business outcomes, finance discipline, and regulatory assurance.

From Scorecards to Runtime Levers

Earlier AI programs measured quality and latency as reporting functions. That model collapsed once workloads moved into regulated, high-stakes domains. Today KPIs and SLAs are embedded into architecture as triggers. A breached latency threshold routes requests to a smaller backup model.

A safety KPI breach invokes human-in-the-loop review. An SLA violation against uptime or governance artifacts escalates to vendor penalties or contract renegotiation.

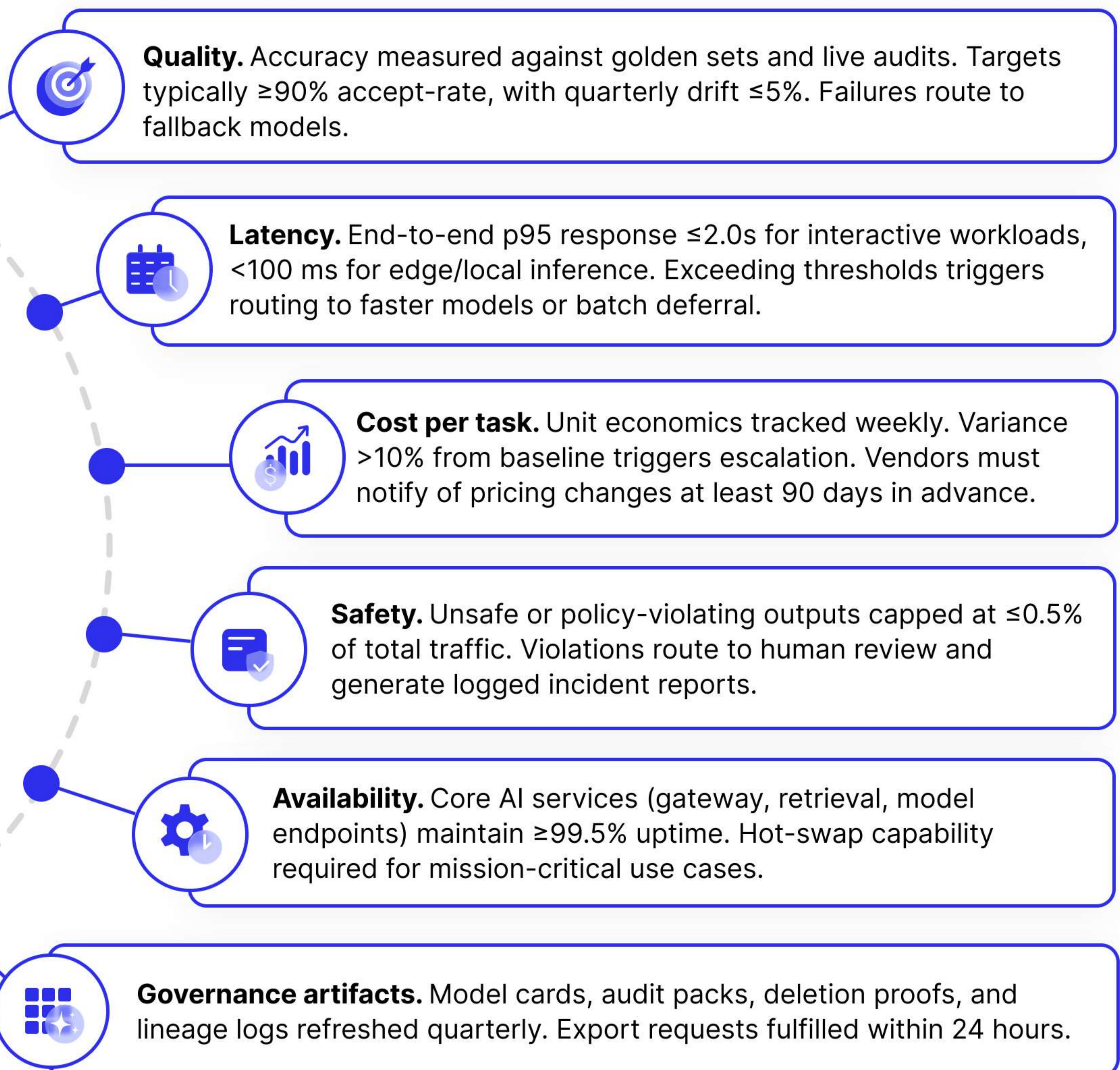
This change has made KPIs/SLA design as important as model selection. Finance teams require \$/task reporting to validate ROI; compliance requires deletion proofs and lineage exports; operations teams require p95 latency guarantees to maintain service desk or fraud-detection SLAs.

Without this evidence, systems are blocked at procurement review.

The Core Domains

Despite sector variation, most enterprises converge on the same six domains for KPIs and SLAs. These cover the lifecycle from model performance to governance evidence, each mapped to operational enforcement.

Core KPI and SLA domains in enterprise AI (2026):



Economic and Operational Implications

When KPIs and SLAs are treated as runtime contracts, they reshape economics. Latency KPIs encourage routing architectures that send simpler traffic to cheaper models, reserving premium endpoints for high-value cases.

Cost SLAs cap renewal surprises by binding per-token or per-seat pricing. Availability SLAs encourage teams to implement hybrid redundancy — using one vendor endpoint alongside an internal fallback — to avoid breach penalties.

The operational impact is equally strong. Teams design monitoring to page owners the moment thresholds are crossed. Evaluation harnesses re-run nightly to detect drift before it accumulates. Finance, compliance, and engineering share the same dashboards, so trade-offs are visible across functions.

KPI and SLA Domains for 2026

Domain	Typical KPI Target	SLA Commitment	Enforcement Mechanism
Quality	≥90% accept-rate, drift <5%/quarter	Vendor provides exports; platform enforces	Golden set regression + live audits; fallback routing
Latency	≤2.0s p95 interactive; <100ms edge	Vendors publish p95 reports; penalties apply	Routing layer redirects to smaller/faster models
Cost per task	Stable within ±10% of baseline	Price change notice ≥90 days; renewal caps	Weekly \$/task dashboards; FinOps alerting
Safety	≤0.5% unsafe outputs	Vendor artifacts + internal approval gates	Runtime guardrails + escalation to human review
Availability	≥99.5% uptime for endpoints/gateway	Penalties for breach; redundancy required	Hot-swap failover; SIEM alerts
Governance artifacts	Quarterly refresh; 24h export fulfillment	On-demand delivery of logs/cards/deletions	Export APIs; SIEM integration; compliance checks

Closing Note

By 2026, KPIs and SLAs will act as the operating system of enterprise AI. They decide when systems escalate, when vendors are liable, and when regulators are satisfied. Programs that design them as runtime levers — rather than static reports — retain control under pressure and scale without losing predictability.

05

Economics and TCO

(12–24 months)

This section defines how the buy-or-build decision translates into concrete financial outcomes. The comparison uses modeled scenarios, current market pricing, and operational cost patterns we at Sombra see in live deployments.

The goal is to give decision-makers a clear view of baseline costs, scaling behavior, and where the economic crossover points occur.

Baseline Scenario

Current benchmarks show a substantial gap in startup costs.



Packaged AI services can deliver production-ready capabilities with total first-year spend between **\$0.75 million** and \$2.0 million, depending on scope, integration depth, and compliance requirements.



Building in-house with current-generation GPU infrastructure, engineering staff, and governance processes typically **requires \$2.5 million to \$4.8 million in year one.**



In a **buy** scenario, enterprise-grade AI SaaS products have matured to the point where most include identity integration, role controls, pre-built connectors, and telemetry dashboards. Vendors package these capabilities into a subscription that can be rolled out quickly. This model keeps capital expenditure low, but it means you are inside the vendor's guardrails for how the system operates.

A **build** approach reverses the equation. You take control over the architecture and infrastructure, but you also take on the capital outlay. Cloud inference using, for example, NVIDIA H100 GPUs in early 2025 **averages \$1.80 to \$2.50 per GPU hour** for on-demand pricing, with committed-use discounts bringing it closer to **\$1.20–\$1.50**. Training runs, even at smaller scales, can consume thousands of GPU hours, and sustained inference loads can drive monthly costs into six figures.



Beyond raw compute, storage, data engineering, MLOps pipelines, and security hardening, add recurring expenses in both engineering time and cloud bills.

Cost Structure Differences

The two approaches differ in how costs behave over time. Buying packaged AI is dominated by variable costs tied to usage – API calls, storage, and service tiers. It scales linearly with adoption, which keeps early expenditure predictable but can lead to high unit costs at very large volumes.

In most buy-first stacks, two items dominate the bill: per-user licensing and usage charges. Price increases of 10–30% at renewal are common and often accompanied by premium charges for high-touch support. Unless contracts set caps, these increases should be assumed in planning.

ChatGPT Enterprise at \$60–\$75 per user per month is a common reference point.

At 1,000 users, this equates to \$0.78M–\$0.90M annually before accounting for usage.

Building in-house frontloads costs in infrastructure and talent. The fixed investment is high, but the marginal cost per inference or workflow step decreases as usage grows.

Once deployed at scale, in-house systems can operate at significantly lower per-unit costs than API-based models, provided utilization remains high.

Usage fees depend on the model mix and workload.

A premium model priced at \$0.03 per thousand tokens keeps costs modest at moderate request rates, while a smaller model at \$0.002 per thousand tokens is more efficient for routine queries.

At high query-per-second loads, token-based billing can become a significant portion of spend. This billing model remains the standard across platforms such as Azure OpenAI and AWS Bedrock.

Market Pricing Context

Public cloud providers have stabilized AI model pricing in 2025, with API-based GPT-class services ranging from **\$0.002 to \$0.015 per 1,000 tokens for input** and **\$0.006 to \$0.04 for output**. High-context, multimodal, or domain-specific APIs may have prices that are up to three times higher than standard rates. Storing proprietary embeddings and vector indexes incurs ongoing costs, particularly at an enterprise scale.



GPU rental rates in major regions average **\$2.50 to \$3.50 per H100 GPU-hour** under committed usage. Using spot pricing and preemptible instances can reduce costs by 30 to 50%, although this comes with a risk to service levels.



Owning hardware shifts the equation further: a single H100 PCIe card lists around **\$30,000 to \$35,000** with a three-year service life. This brings the effective **per-hour cost under \$1.50**, but requires upfront capital and ongoing management.

Hybrid Positioning

Hybrid strategies combine vendor APIs for rapid delivery with in-house model hosting for high-volume, sensitive, or latency-critical workloads. The economics here depend on accurate workload segmentation. For example, routing 80% of daily inference volume to an in-house RAG stack and leaving 20% on premium APIs can cut total compute spend by 40–60% after the second year.

However, hybrid requires duplicated capabilities – monitoring, scaling, and update pipelines on both sides – which increases operational complexity and demands more specialized staff. The initial investment is lower than full in-house but higher than pure buy.

Economic Tilt Factors

Several drivers can shift the cost curve in favor of one model over another.

Volume predictability matters

Steady, high-volume workloads amortize the fixed cost of in-house infrastructure faster, while unpredictable demand makes pay-per-use APIs safer.

Model stability is another driver

Stable workloads tied to mature models benefit from long-term hardware ownership, while rapidly evolving needs may justify subscription-based vendor models to avoid stranded investment.

Data sensitivity can also tip the scales, since keeping regulated datasets inside controlled infrastructure avoids vendor compliance premiums and transfer risk.

Every architecture has **hidden costs** that surface over time.

For buy scenarios, these often include premium API usage beyond the base subscription, per-seat overages, or compliance add-ons for specific geographies.

For build scenarios, the hidden costs come from sustaining the engineering team, handling security audits, and refreshing infrastructure as model architectures and hardware evolve.

Scaling dynamics matter as much as starting costs. AI systems tend to grow in usage once they prove useful, and this growth can double or triple operating costs in a short period. Inference traffic that begins as a pilot in one department can quickly expand across the company, consuming more GPU hours or API calls than budgeted.

Without careful monitoring, this growth can erode ROI, even if the initial implementation looked efficient.

This is where **CFO and CTO alignment** becomes decisive. The finance leader sees recurring costs and risk exposure, while the technology leader focuses on performance, reliability, and capability.

When these perspectives are aligned early, scaling can be planned with guardrails—such as automated cost caps, load balancing across cheaper endpoints, or scheduled retraining windows that avoid peak pricing periods.



We see architecture as an economic control surface. The way you split ownership between vendors and your own infrastructure determines not just what the system can do, but how its costs behave over time. We design for predictability as much as for performance.

Yurii Nakonechnyi,
Sombra CTO

When architecture, economics, and governance are designed together, the AI platform's cost profile stays predictable. The enterprise gains the ability to grow usage without losing financial discipline, which is critical as AI moves from pilot to core infrastructure.

Regulatory Influence on Cost

Forthcoming compliance deadlines, especially the **EU AI Act in 2026**, will influence economic decisions as well.

Systems handling high-risk use cases will face stricter testing, documentation, and monitoring requirements. Vendor APIs may absorb part of this burden, but will likely increase prices to reflect added compliance overhead.

Self-hosted solutions will need additional budget for audit tooling, explainability modules, and periodic external certification.



These costs are often underestimated in early planning but can **reach 10–15% of total AI operating expenses** in regulated industries.

Price comparison

Cost driver	Buy (SaaS / API)	Build (your infra & team)	Hybrid (“own the middle, rent the rest”)
Seat licenses (per-user / mo)	Microsoft 365 Copilot: \$30/user/mo (EA) Google Workspace Gemini: \$20(Business) / \$30(Enterprise) GitHub Copilot Business: \$19/user/mo_	N/A	Often mixed: pay seats where work happens (e.g., Copilot/Gemini for office suites), but keep core logic in your platform to avoid feature-creep seats. Your docs model Year-1 \$0.75M–\$2.0M for Buy at enterprise scale.
LLM API usage (tokens)	OpenAI (illustrative): GPT-4.1 \$... per 1M tokens (tiered); GPT-4o mini \$0.15/M input & \$0.60/M output (very low)	If self-hosting, token fees go to GPU & ops instead; your files benchmark \$0.002/1K (small) to \$0.03/1K(premium) for planning.	Hybrid mixes both: cheap models for “first pass,” escalate hard prompts to premium models; API plus small-model routing is the main cost lever (your docs emphasize this).
GPU compute (inference/ training)	Included in vendor fee (opaque)	Cloud H100 on-demand (8× GPUs): AWS p5.48xlarge starting ~\$55.04/h (~ \$6.9/ GPU-h) Azure ND96isr H100 v5 ~\$98.32/h (~ \$12.3/GPU-h) Your files use a fully-loaded enterprise assumption of ~\$38/ GPU-hto include networking, support, capacity insurance.	Hybrid mixes both: cheap models for “first pass,” escalate hard prompts to premium models; API plus small-model routing is the main cost lever (your docs emphasize this).
Staffing	1–2 FTEs to integrate & vendor-manage: ~\$320k/yr in your model.	6 FTEs (MLE/MLOps) ~\$960k/yrbaseline; Year-2 maintenance 20–40%of Year-1 dev cost.	Platform team (2–3 FTE) to own the gateway/guardrails/eval; keep vendors behind it. Your docs assume \$320k–\$480k/yr as typical.
Governance & audits	Typical ~\$80k/yr (evidence packs, audits).	~\$120k/yr (tooling + internal effort).	Same as Build for the runtime you own; vendor must feed your logs.
One-time setup / services	Integration/pro services: \$100k–\$400k common; can rival Year-1 license if scope is broad.	\$1.5M setup is a realistic placeholder (data, pipelines, hardening) in your model; Year-1 \$2.5M–\$4.8M all-in is typical at enterprise scope.	Usually light: build the “middle” once (gateway, eval, retrieval); vendors plug in. Budget \$150k–\$300k initial platform lift depending on starting point (from your ranges).
2-year TCO (illustrative, 50M tokens/ mo)	≈ \$2.1M	≈ \$5.9M	≈ \$2.6–\$3.5M (Buy seats + API + your small platform).
Break-even signal	–	Your files show Build overtakes Buy only around ~180M tokens/mo (with those GPU & staffing assumptions).	Hybrid moves the crossover earlier if you keep GPUs busy and route easy prompts to small models.

06

Proof-of- Concept Playbook

Patterns and
Case Snapshots

Proof-of-Concept Playbook

A modern PoC is the first controlled run of your production AI stack in miniature. It encompasses identity, data, model calls, guardrails, logging, and governance checks applied at scale, all focused on a narrow scope to allow the team to learn quickly and minimize risk.

The goal is to surface integration gaps, data quality issues, and policy conflicts early, while the system is still easy to change.

Stakeholders now expect this early stage to perform. A PoC that feels slow, brittle, or hard to audit will lose support quickly.

The sequence below distills the PoC process into a set of steps that a cross-functional team can execute with clear success criteria. Each step explains what to do, why it matters now, and which controls or platform features can make it work.

AI adoption is mainstream – 78% of firms report using it in at least one function – and most executives have already experienced commercial AI tools.

Step 1

Fix the business objective and the guardrails

Pick **one business outcome** and write it like an OKR: e.g., “Reduce average handle time (AHT) in Tier-1 support by **10%** without lowering CSAT.” Name the **unit of work** (“resolved ticket,” “approved claim”), because every metric later depends on this denominator. Then set **non-negotiables** up front: what data may/ may not be used; what must be redacted; which actions (send email, file a claim, push to CRM) require human approval.

In 2025, you can bake these into runtime policy rather than PowerPoint: all three clouds support **private endpoints** (Azure private endpoints, AWS **PrivateLink**, Google **Private Service Connect**) so data doesn’t have to traverse the public internet; each publishes data-handling policies you can cite in your internal review.

A PoC that ships fast but ignores guardrails will stall at security review. A PoC that nails guardrails early **unblocks** Security and Compliance when you ask to pilot with real users. Your docs also show that governance decisions (SSO, logs, deletion proofs) are what separate a PoC that scales from one that dies in procurement.

Step 2

Gate the data: access, quality, lineage

In a PoC worth doing, the model answers **from your content**. That means you must confront data contracts and permissions at the start, not the end. Implement **retrieval with access checks at query time** (users only see citations they're allowed to see) and **redact-then-retrieve** for sensitive fields. Azure, Google, and AWS now publish first-party RAG guides you can mirror; they all treat retrieval as a **control surface** (not just an accuracy trick).

Two practical gates before you proceed:



Access gate: demonstrate that a user lacking permission cannot retrieve restricted sources (prove with a log export).



Quality gate: run a “**golden set**” of 50–200 real questions with expected answers and sources; if top-k retrieval is off by more than, say, **10–15 pp** precision@k versus baseline search, fix data and chunking before touching models. Your files emphasize that this single discipline prevents wasted model work later.

Step 3

Establish a baseline with an off-the-shelf reference

Create a reference experience using a reputable SaaS or straightforward API so stakeholders can understand what “good enough” looks like today in terms of latency, quality, and cost per task. This is not your end state; it's the control group you will beat (or decide you won't).

Since adoption is already widespread, your business partners likely use Copilot, Gemini, or Einstein on a daily basis; let this set expectations and keep you accountable regarding time-to-value.

Why this matters

The TCO sections show the first-12-months cost of “Buy” is often \$0.75M–\$2.0M including integration, versus \$2.5M–\$4.8M to “Build.” A baseline shows if the fast path is already enough – and gives Finance real \$/task for comparison.

Step 4

Stand up the evaluation harness (quality, latency, \$/task, safety)

Treat evaluation like CI/CD. Every change – data, prompt, model, policy – re-runs the golden set offline and feeds a small online A/B when traffic arrives.

Your dashboard should track:

1

Task quality (accept-rate vs. human label or judge-model with spot audit)

2

Latency (p50/p95 end-to-end)

3

Cost per task (tokens, seats, or GPU time)

4

Safety (blocked/approved outputs, policy infractions)

This matches NIST's Generative AI Profile for continuous risk measurement and ISO/IEC 42001's AIMS obligations. Weekly monitoring of cost/token and hallucination rate avoids renewal shocks. Remember that teams that watch **cost/token and hallucination rate weekly** learn faster and avoid nasty surprises at renewal.

Step 5

Iterate features against the baseline (earn complexity)

Now you earn every ounce of complexity by beating the baseline **with evidence**:



RAG first: once your retrieval gates clear, expect a **step-change** in grounded accuracy without touching training. All three clouds publish 2025 RAG design guides–use their playbooks (chunking, metadata, re-rankers) rather than inventing your own.



Routing next: send easy queries to small/cheap models, escalate hard ones to bigger models. This is where unit economics bend. (Azure/Vertex docs increasingly show routing as a first-class pattern.)



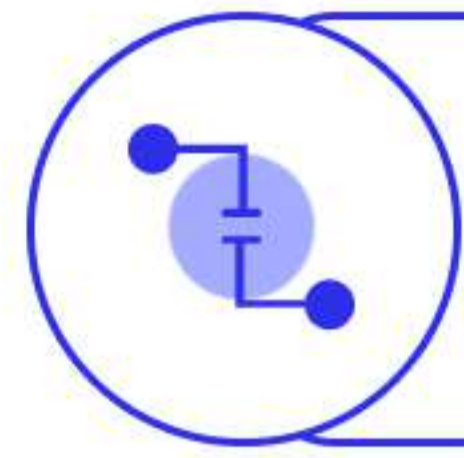
Fine-tuning later: only when retrieval and routing plateau, and only if your proprietary data is dense enough to move the needle.

Each change should get a “delta” vs. baseline: +X pp accuracy, –Y% cost, –Z ms latency. No delta? Roll it back.

Step 6

Security and Compliance Checks

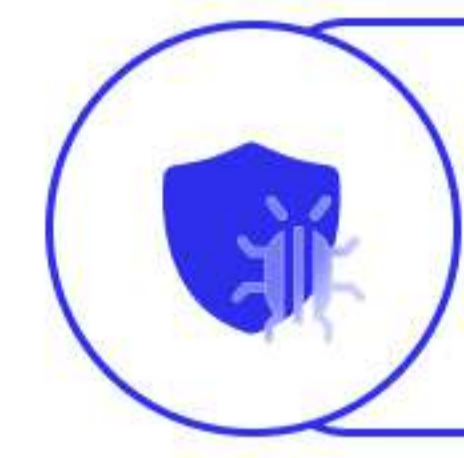
Before you even whisper “pilot,” **prove controls work in runtime:**



Network isolation: show private connectivity to model endpoints (screenshot the **PrivateLink/PSC/Private Endpoint** configuration and a packet-trace with no public egress).



Data handling: export logs tying each answer to sources, model/prompt/policy versions, and user; run a deletion test and keep evidence. This aligns to **EU AI Act** record-keeping and to **NIST/ISO** expectations for traceability.



Safety: demonstrate that unsafe outputs are blocked or routed for approval (show both a blocked case and an approved case in the logs).

Why this matters

PoCs without runtime evidence will not scale in EU operations.

Step 7

Pilot plan, success criteria, and go/no-go

A PoC without a disciplined pilot transition is just a demo. The pilot phase is where you move from lab traffic to real users and measure how the system performs in a live, but limited, environment.

The point isn't to prove the concept again – it's to validate that the entire delivery chain works under production conditions: integrations hold, metrics stay within agreed thresholds, support teams can respond, and governance controls remain effective.

Define scope and containment

Choose one business unit, queue, or site. Limit the scope so you can attribute outcomes directly to the AI system without contamination from unrelated process changes.

Define:

- Number of users or agents
- Workflows or request types in scope
- Timeframe (start/stop dates, typical 4–8 weeks)
- Any exclusions (e.g., certain customer segments, high-risk transactions)

Set success and failure criteria up front





KPIs should map directly to the business objective from Step 1.

- For example:**
- *AHT reduction*: $\geq 10\%$ improvement over baseline with no CSAT drop
 - *Cost per task*: \leq baseline cost established in Step 3
 - *p95 latency*: ≤ 2.0 seconds end-to-end
 - *Safety*: \leq agreed threshold for blocked outputs or policy violations

Also define stopping conditions – e.g., budget cap reached, drift in retrieval precision, safety incidents over limit. These are not negotiable mid-pilot unless formally re-approved.

Prepare operational readiness

Before the first real request flows, confirm:

-  **Incident response:** named owners, escalation paths, rollback procedures tested in a tabletop exercise
-  **Monitoring and alerting:** dashboards active for latency, cost, safety, accuracy; alerts configured to page the right owner
-  **Support channel:** open and staffed on day one (with response times defined)
-  **Change freeze:** no unrelated system changes during the pilot window

Run the pilot as a controlled experiment

All traffic in scope passes through the same gateway, guardrails, and logging stack used in PoC. Apply online A/B if volume allows, keeping a control group to measure net effect. Track weekly trends and investigate anomalies immediately – don't wait until the end of the pilot.

Make the go/no-go decision based on evidence

At the end of the pilot, evaluate against the criteria you set:

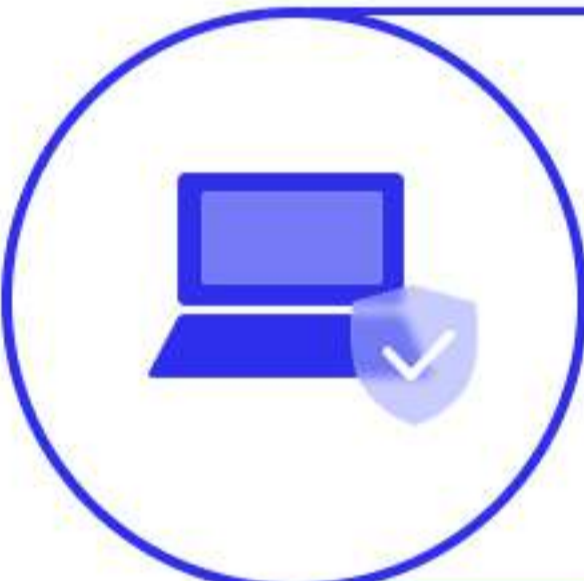
- 1 Did KPIs improve to target without degrading quality or safety?
- 2 Were costs per task stable or improving?
- 3 Did all governance controls function as designed?
- 4 Were incidents resolved within expected timelines?


If yes – scale to the next business unit or workflow, with the same controls. If no – pause and fix the weakest link, often retrieval quality or evaluation coverage, before adding complexity.


Sombra's Take

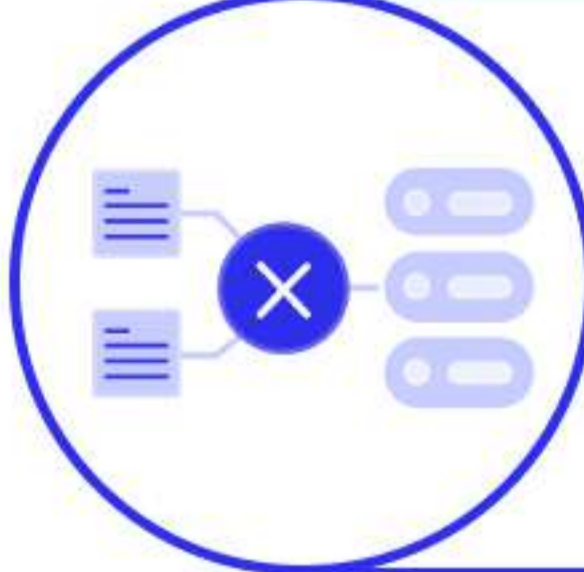
This step is where most AI initiatives stumble, because teams either skip formal success criteria or treat them as soft guidelines. In 2026, with compliance deadlines fixed and stakeholder expectations high, a pilot that lacks measurable thresholds is a political liability. You want the decision to be boring and data-driven – not a subjective debate. The strongest pilots I've seen have a written charter, an incident rehearsal before launch, and a clean "pass/fail" table at the end. If you can't make the go/no-go call in one meeting with the evidence on the screen, the pilot wasn't designed right.

Observations and 2025 trends to keep you out of trouble

- 

Private by default is normal now. You can satisfy Network and Data Protection without building everything yourself; Azure, AWS, and Google document **private endpoints** for AI services. Use them from day one—even in PoC—so there's no re-work later.
- 

RAG is a governance tool as much as an accuracy tool. It keeps answers **grounded in your content** and gives you auditable citations—an expectation in NIST/ISO language and a practical requirement under the EU timeline.
- 

Measure \$/task early. Your materials highlight weekly **cost/token** and hallucination monitoring at leaders like Walmart—PoCs that track this avoid “surprise” invoices and know when routing to smaller models pays off.
- 

Lock-in is architectural, not just contractual. If any app or team can hit a model **without** your gateway, you've lost policy and cost control. Your docs call this out: put an abstraction layer in front of vendors so models/apps are **swappable**.

Follow this sequence and you'll avoid the classic trap your files warn about – an impressive demo that can't survive production – and you'll be able to defend your choices to Security, Compliance, and Finance with artifacts they recognize. The point of a 2025 PoC isn't a wow moment; it's to prove you can ship, measure, and pass audit on a small canvas, then scale with confidence.

Patterns and Case Snapshots To Help With The Decisions

In 2026, the right AI adoption path is rarely decided in a planning deck. It emerges under pressure – when a launch date is fixed, compliance is closing in, or cost curves start to bend. The following cases come from programs that made it through PoC, went live, and kept operating under real-world constraints.

Each snapshot shows the conditions, the choice, the architecture, the outcome, and the lesson learned. The goal: give you field-tested reference points you can adapt, not just theory.

Case 1

When Buy wins:

The service desk that needed a 90-day result

Context

A global manufacturer had rising ticket volumes and an AHT target for Q2. IT had no bandwidth to build.

Decision

Buy a service-desk assistant embedded in the existing ITSM suite.

Architecture

Keep identity and logging in the company's gateway; use the vendor's native connectors for tickets and knowledge; stream usage and quality to the BI stack.

Result

Value appeared in weeks: **AHT down ~12%**, first-contact resolution **+6 points**. Security signed off because traffic ran through **private cloud endpoints** (Azure/AWS/GCP now provide these by default), so no "public internet" debate.

Pitfall surfaced

Seat creep at renewal. They fixed it by moving to pooled licenses and by routing simple Q&As into a cheaper, vendor-provided small model.

Case 2

When Buy wins:

Sales content where the workflow is already standardized



Context

A B2B SaaS firm wanted faster proposal cycles across 14 countries.

Decision

Buy the suite assistant (Office/Docs/Slides) + CRM-specific AI for snippets.

Architecture

Nothing exotic: SSO/SCIM, DLP rules, and a simple review workflow for external copy.

Result

Content throughput up ~30%, but the real win was adoption—people used it because it lived inside their daily apps. Vendors are normalizing this as a suite feature (Microsoft/Gemini/Einstein/Now Assist), so time-to-value is unusually short.

HYPERLINK →

Pitfall

Governance drift. They negotiated “no training on our data by default” into the MSA—even though public policies already say this—so legal never had to re-litigate it.

Case 3

When Buy wins:

Edge extraction with latency guarantees.



Context

A logistics company needed invoice and bill-of-lading extraction in depots with unstable links.

Decision

Build small, fine-tuned models and run inference locally; push summaries to the cloud when connected.

Architecture

Router → SLM (millisecond answers) → queue for sync; cloud model only on ambiguous cases.

Result

Sub-100 ms local responses and predictable throughput; cloud costs dropped because 80–90% of traffic never left the site.

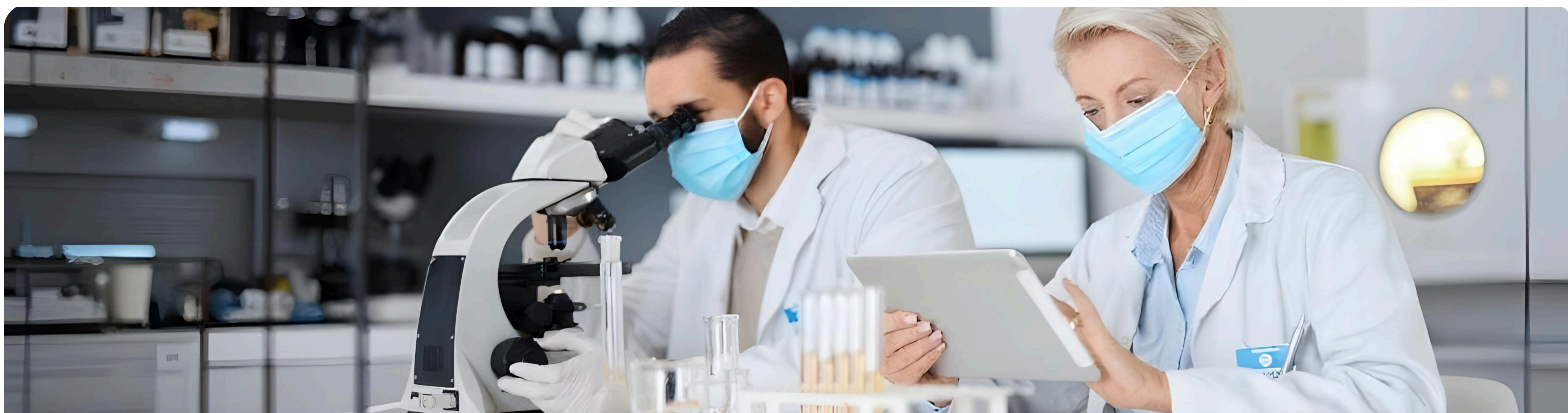
Pitfall

Confidence calibration. They added human spot-checks until the SLM's "uncertain" flag was trustworthy.

Case 4

When Buy wins:

Regulated research with strict audit



Context

A pharma R&D group needed literature triage with perfect traceability.

Decision

Build the retrieval layer and evaluation harness in-house; use cloud models behind **private endpoints**.

Architecture

Their assets were the ingestion/metadata rules and a ruthless evaluation harness; models were swappable (Azure/Vertex/Bedrock) over private links.

Result

Higher researcher throughput and painless audits—the logs could reconstruct any answer with sources and versions.

Pitfall

Ingestion quality. They staffed one data engineer full-time to maintain chunking and metadata—cheap compared with fine-tuning too early.

Case 5

When Buy wins:

The enterprise knowledge assistant



Context

A bank wanted a cross-department “answer engine” with permissions and citations.

Decision

Own the middle (gateway, guardrails, retrieval, evaluation), rent model endpoints and a few narrow SaaS apps.

Architecture

Requests hit the company gateway → redaction → **RAG** over sources with **access checks at retrieval** → model router (small model first, big model when needed) → guardrails → logs to SIEM.

Result

Accuracy rose steadily with better retrieval; unit cost fell as routing matured. When the EU AI Act obligations were staged in (2025–2026), nothing changed in code, the evidence (logs, model cards, deletion proofs) already existed.

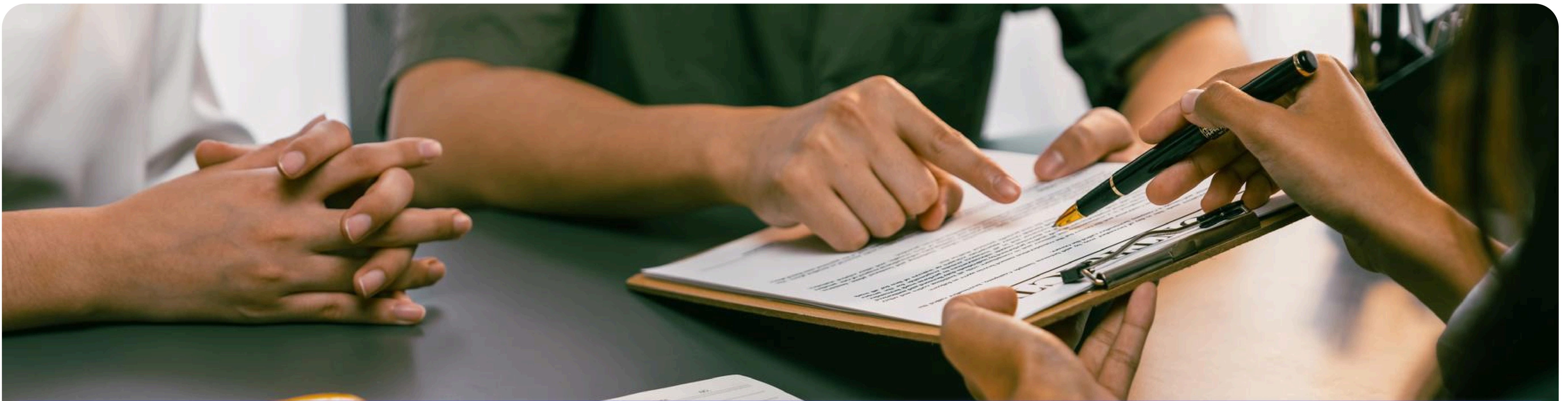
Pitfall

One team bypassed the gateway for a “quick demo,” hiding spend and policy. IT shut it down and made the gateway mandatory.

Case 6

When Buy wins:

Claims triage with agentic steps



Context

An insurer needed triage plus a few follow-up actions across three systems.

Decision

Keep orchestration and approvals in-house; rent the model API.

Architecture

The agent could plan steps but had **guardrails**: idempotent tools, pre-conditions, and human sign-off for actions that moved money.

Result

Cycle time down ~18%;

zero policy breaches in six-week pilot because the guardrails lived in the company's runtime, not in vendor code.

Pitfall

Unbounded loops in early runs; the team added timeouts and step budgets.

Case 7

When Buy wins:

Global rollout under a hard compliance clock



Context

A retailer had to expand AI assistants to the EU and U.K. while the **EU AI Act** deadlines loomed.

Decision

Deploy vendor apps and API endpoints **only via private connectivity**; keep logs/evidence centrally with their own evaluation harness.

Architecture

Azure/Bedrock/Vertex endpoints over PrivateLink/Private Endpoint/PSC; model choice by region; central SIEM and deletion tests.

HYPERLINK →

Result

Procurement sailed through because the run-time evidence matched **NIST/ISO** expectations and the EU timeline (GPAI obligations in **Aug 2025**, broader duties by **Aug 2, 2026**).

HYPERLINK →

Pitfall

None major; the lock-in risk was contained by contractually guaranteed exports (prompts, embeddings, logs) and a tested exit path.

Case 8

When Buy Fails Despite Speed



Context

A regional bank bought a domain-specific AI chatbot for customer service to meet a 6-month launch goal.

Decision

Off-the-shelf vendor platform with pre-built banking intents.

Architecture

Direct integration into the call-center platform, minimal customization.

Result

Launched on schedule, but accuracy plateaued at ~70% because the vendor's model couldn't ingest the bank's internal policy documents. Customer satisfaction dropped and the bot was scaled back.

Pitfall

Assumed vendor training data would cover proprietary processes; never validated retrieval on internal content before go-live.

Case 9

When Build Fails Despite Technical Success



Context

A healthcare provider built an in-house clinical summarization tool to meet strict HIPAA and latency requirements.

Decision

Fully owned ingestion, retrieval, model hosting, and evaluation.

Architecture

On-prem GPUs, private endpoints, domain-specific model trained on de-identified data.

Result

Technically sound – sub-second latency, high accuracy, passed audits – but cost per note stayed 4× higher than SaaS alternatives. The program was mothballed after a year because no one could justify the economics.

Pitfall

Ignored volume economics; built for peak security needs without segmenting workloads that could have been safely offloaded to cheaper vendor models.

Appendix Section

Appendix 1

RFP Question Bank

Tight, evidence-seeking questions to drop directly into procurement.

Data & Privacy

- “Point me to your public policy confirming that our prompts/outputs are not used to train your models by default. Confirm the same in the contract.” (Azure/Google/AWS/OpenAI all publish this—match it.)
- “What’s your retention window for prompts/outputs and logs? Can we configure zero data retention? Provide the doc.”
- “Do you support customer-managed keys? Describe key rotation and crypto-erase on revocation.”

Connectivity & Security

- “Can all inference run through private endpoints (PrivateLink/Private Service Connect/Azure Private Endpoint)? Show the diagram.”
- “Demonstrate SSO (SAML/OIDC), RBAC, and SCIM provisioning with our IdP. Provide de-provision SLA.”

SLA/SLO & Support

- “Commit to $\geq 99.9\%$ availability and name your latency SLO for the target regions. Provide your status page history.”
- “Define Sev-1 response times for L2 and L3. Provide named escalation contacts.”
- “Explain your credit schedule and how credits scale with incident duration and scope.”

Compliance & Audit

- “Map your controls to EU AI Act, NIST AI RMF/GenAI Profile, and ISO/IEC 42001. Provide model cards and an audit evidence pack.”

Portability & Exit

- “Provide export formats for data, embeddings, prompts, and policies. Commit to deletion within 30 days (with proof), and include X assistance hours for termination.”

Extensibility

- “Show webhooks/SDKs and how we can route requests via our gateway and our retrieval layer.”

Search (Additional)

- “How is search/retrieval evaluated for accuracy (precision/recall/coverage)? Provide benchmarks and methodology.”
- “Do you support multi-hop or re-ranking in retrieval? Provide performance data.”
- “Can retrieval enforce per-user permissions at query time? Demonstrate with exportable logs.”

Appendix 2

Red-Flag List

If you see any of these, stop the deal.

- Training on your data by default, or vague opt-out language that contradicts the provider’s public privacy page.
- No option for private connectivity; inference must traverse the public internet.
- No export of prompts, policies, or embeddings; “proprietary format only.”
- SLA credits that cap at trivial amounts or exclude the incidents that matter (e.g., latency breaches, extended downtime).
- Logs that cannot be exported to your SIEM; no ability to reconstruct “who ran what, using which sources and model version.”

Appendix 3

The Quick Decision Tree

1 Is the use case standardized across your industry?

- Yes → **Buy baseline.**
- No → go to 2.

2 Will proprietary data deliver ≥ 15 -point quality lift or unique capability?

- Yes → go to 3.
- No → **Buy or Hybrid with owned retrieval.**

3 Do sovereignty, safety, or latency requirements demand private/on-prem inference (e.g., P99 < 50 ms; regulated workloads)?

- Yes → **Build (or Hybrid with owned models).**
- No → go to 4.

4 Do projected volumes make per-call/API pricing uneconomical within 12–24 months?

- Yes → **Build or Hybrid with portability plan.**
- No → go to 5.

5 Can you sustain MLOps, evaluation, and AI security teams for 12–24 months?

- Yes → **Build or Hybrid based on ROI.**
- No → **Buy/Hybrid now; revisit when capacity exists.**

6 Vendor posture check (roadmap fit, exit terms, data handling, SLAs):

- No → **Bias toward Hybrid/Build.**
- Yes → **Proceed with chosen path and PoC.**

Appendix 4

Vendor Scorecard

How to use

1. Assign a weight (1–5) to each section based on your priorities.
2. Score each vendor 1–5 using the anchors.
3. Multiply $(\text{Score} \div 5) \times \text{Weight}$.
4. Sum across sections.
Keep the comments/evidence block—this is what leadership will actually read.

1 Business Fit & Use-Case Coverage

What: Fit for current and next 12–18 months.

Why: Misfit means delays and missed KPIs.

Evidence: References, demo on your data, roadmap slide.

PoC: Top 5 journeys, measure time saved and exception handling.

Anchors:

- 1 – Generic demo, core flows missing
- 3 – Most flows supported, gaps with roadmap
- 5 – Full coverage now, strong roadmap, industry references

2 Data & Model Fit

What: How well the solution works with your data (RAG, adapters, refresh).

Why: Biggest accuracy gains come from data fit, not model brand.

Evidence: Lift vs. baseline, connectors, indexing rules.

PoC: A/B test vs. current process.

Anchors:

- 1 – Works only on public/examples
- 3 – Basic adapters, partial lift
- 5 – Strong lift across corpus, transparent retrieval quality

3 Security (Access, Network, Data Protection)

What: SSO, RBAC, SCIM, private networking, encryption.

Why: Reduces breach risk and audit pain.

Evidence: Network diagrams, incident runbooks.

PoC: Test deprovision, review logs, run incident tabletop.

Anchors:

1 – Weak roles, public endpoints

3 – SSO + roles, VPC option, IR plan

5 – Full SSO/SCIM, least privilege, private links, tested IR

4 Compliance & Governance

What: EU AI Act, NIST RMF, ISO/IEC 42001 alignment.

Why: Avoids fines and blocked launches.

Evidence: Control matrix, model cards, data handling docs.

PoC: Run DPIA/impact assessment, test deletion/export.

Anchors:

1 – “We’re compliant” with no proof

3 – Partial evidence and mapping

5 – Evidence-backed, audit-ready, clear responsibilities

5 Reliability & Performance (SLOs)

What: Availability, latency, throughput.

Why: Slow or unstable AI disrupts adoption.

Evidence: Published SLOs, status page history.

PoC: Load test at peak.

Anchors:

1 – No SLOs, unstable under load

3 – Basic SLOs, some spikes

5 – Strong SLOs, predictable under stress

6 Observability & Quality Management

What: Dashboards, eval harness, drift alerts.

Why: No visibility = no improvement.

Evidence: Logs, test reports, BI exports.

PoC: Review eval run, simulate drift.

Anchors:

- 1 – Minimal logging, no harness
- 3 – Basic dashboards, some A/B
- 5 – Full evaluation pipeline, role-based dashboards

7 Integration & Extensibility

What: Connectors, APIs/SDKs, webhooks.

Why: Faster integration, lower change costs.

Evidence: Connector list, API docs, sandbox.

PoC: Connect top 3 systems, test extension.

Anchors:

- 1 – Custom everywhere, brittle
- 3 – Standard connectors, stable APIs
- 5 – Rich connectors, clean SDKs, documented extensions

8 Operability & Support

What: Day-2 ops, L2/L3, TAM.

Why: Incidents are inevitable.

Evidence: Severity definitions, escalation ladder.

PoC: Open ticket, run QBR.

Anchors:

- 1 – Email-only, slow escalation
- 3 – Clear Sev targets, quarterly reviews
- 5 – Fast L2/L3 SLAs, proactive TAM, proven path

9 Commercials & 12–24-Month TCO

What: Pricing model, unit economics, cost controls.

Why: Entry is cheap, scale is not.

Evidence: Pricing sheet, TCO model, routing/caching levers.

PoC: Run week of traffic, compare invoice.

Anchors:

1 – Opaque pricing, surprise fees

3 – Forecastable, some levers

5 – Transparent, cost-reduction playbook, good discounts

10 Lock-in & Exit

What: Exports, portability, termination assistance.

Why: Protects options, improves vendor behavior.

Evidence: Export docs, deletion proofs, exit SLAs.

PoC: Test export/import, verify deletion.

Anchors:

1 – No exports, proprietary formats

3 – Partial exports, some portability

5 – Clean exports, orchestration-friendly, vendor assistance

11 Vendor Viability & Roadmap Confidence

What: Financial stability, track record, references.

Why: Avoid tying workflows to shaky suppliers.

Evidence: Financials, roadmap delivery, references.

PoC: Review past releases, talk to customers.

Anchors:

1 – Weak finances, missed releases

3 – Decent track record, references

5 – Strong balance sheet, enterprise-grade references

Score Summary Table

Section	Weight (1-5)	Vendor Score (1-5)	Weighted Score	Comments/ Evidence
1. Business Fit & Coverage				
2. Data & Model Fit				
3. Security				
4. Compliance & Governance				
5. Reliability & Performance				
6. Observability & Quality Mgmt				
7. Integration & Extensibility				
8. Operability & Support				
9. Commercials & TCO				
10. Lock-in & Exit				
11. Vendor Viability				
Total				

References

1	Sombra. <i>Deterministic vs. Generative AI.</i>	Link
2	EY. <i>Should Organisations Buy AI Systems or Build Them?</i>	Link
3	IBM. <i>AI Adoption Challenges.</i>	Link
4	IBM. <i>AI Orchestration.</i>	Link
5	OpenAI Blog. <i>In-House AI: The Machine Learning Lifecycle.</i>	Link
6	Thematic. <i>Cost to Build an AI Feedback Analytics Platform.</i>	Link
7	Hugging Face. <i>AI TCO Calculator.</i>	Link
8	Plain English. <i>Understanding the Costs of Large Language Models.</i>	Link
9	CIO Influence. <i>Enterprise Observability Framework Considerations.</i>	Link
10	Sombra. <i>Conversational AI Agent Case Study.</i>	Link
11	Sombra. <i>Operational Cost of Low-Trust Data.</i>	Link
12	LinkedIn. KubeVisor GmbH. <i>Mitigating Generative AI Hallucination.</i>	Link
13	KPMG. <i>The Evolution of Build vs. Buy.</i>	Link
14	AWS. <i>MLOps Checklist.</i>	PDF Link
15	Microsoft. <i>Machine Learning Model Checklist.</i>	Link
16	Sombra. <i>Operational Cost of Low-Trust Data.</i>	[Duplicate; already listed at #34]
17	Sombra. <i>Data Wealth in Asset Management.</i>	Link
18	European Union. <i>Artificial Intelligence Act (EU AI Act).</i>	Link
19	NIST. <i>AI Risk Management Framework (AI RMF).</i>	Link
20	ISO. <i>ISO/IEC 42001 – Artificial Intelligence Management System Standard.</i>	Link
21	McKinsey & Company. <i>The State of AI in 2024 Global Survey.</i>	Link
22	Gartner. <i>Generative AI for Business.</i>	Link
23	AWS. <i>AWS PrivateLink Overview.</i>	Link
24	Databricks. <i>AI & Machine Learning Platform.</i>	Link

Sombra[®]