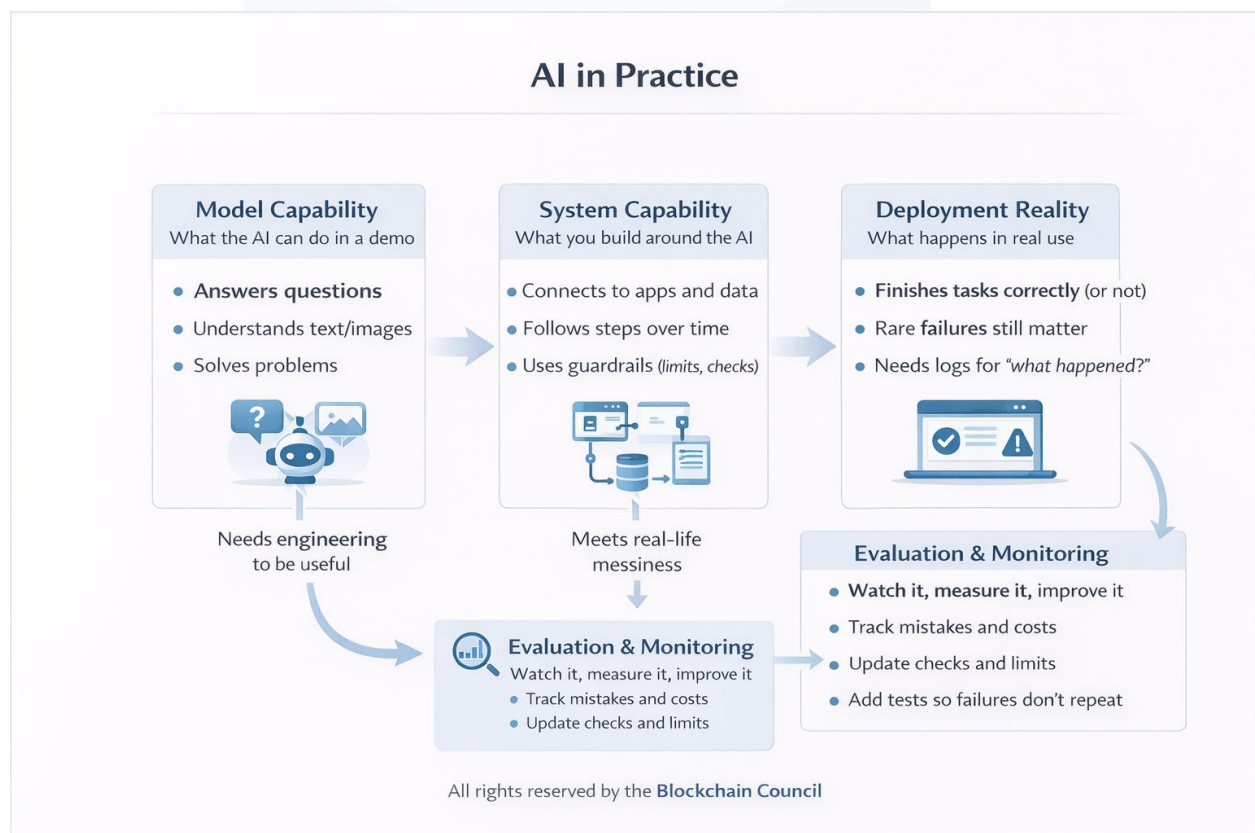


# The State of AI: 2025 in Review, 2026 in Focus

## Introduction

The public narrative around modern AI often compresses complex, system-level change into a sequence of model releases. This framing is convenient but increasingly incomplete. By the end of 2025, the practical impact of AI depended less on single benchmarks and more on whether model capability could be translated into system capability, systems that connect to tools and data, keep state over long horizons, resist adversarial inputs, and operate under cost and policy constraints. In parallel, the governance environment moved from anticipation to implementation, tightening expectations around documentation, monitoring, and accountability. These shifts motivate a paper that does two things at once: it records what changed in 2025 in a way that is operationally meaningful, and it sets out predictions for 2026 that can be evaluated against observed signals rather than rhetorical claims.



The paper is organized around a simple premise: the most important developments in AI are increasingly those that change how systems behave in the wild. That includes technical

changes, tool-native reasoning, long-context workflows, multimodal interaction, but also institutional and infrastructural changes, such as procurement standards, regulatory milestones, security practices, and physical constraints on compute and energy. Because these factors interact, a “state of AI” account that focuses only on model capability trends risks missing the bottlenecks that determine deployment reality. Conversely, a forecast that focuses only on hype cycles risks missing the mechanisms that would make particular outcomes more or less likely.

To address this, the paper adopts a research-oriented synthesis approach with two commitments:

- Describe 2025 not as a list of releases, but as a set of converging system themes and constraints that changed what was feasible to ship.
- Frame 2026 predictions as hypotheses tied to measurable indicators, emphasizing uncertainty and falsifiability over confident narrative.

The combined 2025 and 2026 sections therefore serve different roles. The 2025 section provides a structured inventory of what became operationally normal, reasoning modes coupled with tool access, multimodal systems as default product surfaces, agents confronting reliability and security limits, and emerging integration protocols that reduce tool hookup friction. The 2026 section extends these themes forward and argues that the center of gravity shifts further toward evaluation, governance, and deployment discipline: the field will be judged by repeatable workflow outcomes, by regression testing and monitoring, by action permissioning and auditability, and by resilience to injection and tool abuse.

This framing is intended for a mass audience that still wants research-grade clarity. The goal is not to predict a single “breakthrough moment,” but to specify what observers should track to determine whether AI systems are becoming reliably useful and responsibly deployable. Accordingly, the paper emphasizes system-level metrics and operational artifacts:

- End-to-end task completion under explicit time and cost budgets
- Variance and tail-failure rates, not only average scores
- Deferral and human-in-the-loop design as part of correctness and safety
- Security controls (least privilege, sandboxing, telemetry) as core components of agent systems

- Governance evidence (documentation, logs, post-market monitoring) as a deployment prerequisite

The remainder of the paper proceeds in two main sections. The first, “State of AI , 2025 in Review,” documents the convergences and constraints that shaped shipped systems. The second, “Predictions for AI in 2026,” translates credible forecasts into a monitoring program centered on evaluation, controlled autonomy, security engineering, data provenance, infrastructure constraints, and compliance-driven operationalization. The paper closes by summarizing what would make 2025–2026 a coherent inflection period: not merely more capable models, but a change in how claims are tested and how responsibility is assigned when systems act.

## **Abstract**

This paper synthesizes two adjacent views of the artificial intelligence field: a structured account of what materially changed during 2025, and a set of testable predictions for 2026. Taken together, these sections describe a transition in how progress is produced and how it is judged. The 2025 record is marked by convergence, reasoning features becoming default behavior, multimodality widening into native system design, and agent-like tool use shifting from prototypes toward operational constraints such as reliability, security, governance, and integration standards. The 2026 outlook extends that trajectory but argues that the dominant story will be less about isolated model jumps and more about system discipline: evaluation moving from publicity to infrastructure, autonomy bounded by permissioning and verification, and deployment shaped by compliance timelines and physical constraints (compute, power, and data rights).

Rather than treating forecasts as statements of inevitability, the paper treats them as hypotheses and translates them into measurable indicators. It proposes a monitoring approach organized around system-level outcomes (end-to-end task completion under time and cost budgets), reliability as a distribution (variance and tail failures), security as an operational practice (threat models, sandboxing, incident response), and governance as evidence (documentation, audit trails, post-market monitoring). The central claim is that the combined 2025–2026 arc is best described as capability-rich yet autonomy-limited: models are increasingly effective inside well-instrumented domains, while open-world action remains constrained by compounding error, injection and tool abuse risk, and accountability demands. If the thesis holds, the most consequential progress markers will be comparatively mundane, logs, tests, permission boundaries, and repeatable workflow metrics, because those artifacts determine whether AI can move from impressive outputs to dependable outcomes.

## Keywords

- Artificial intelligence
- Frontier models
- Reasoning and inference-time compute
- Multimodality
- Agents and tool use
- Multiagent systems
- Evaluation and benchmarking
- Workflow completion metrics
- Reliability and uncertainty calibration
- Prompt injection and tool abuse
- AI security engineering
- Audit trails and post-market monitoring
- Data provenance and licensing
- Compute, energy, and infrastructure constraints
- EU AI Act compliance

## State of AI - 2025 in Review

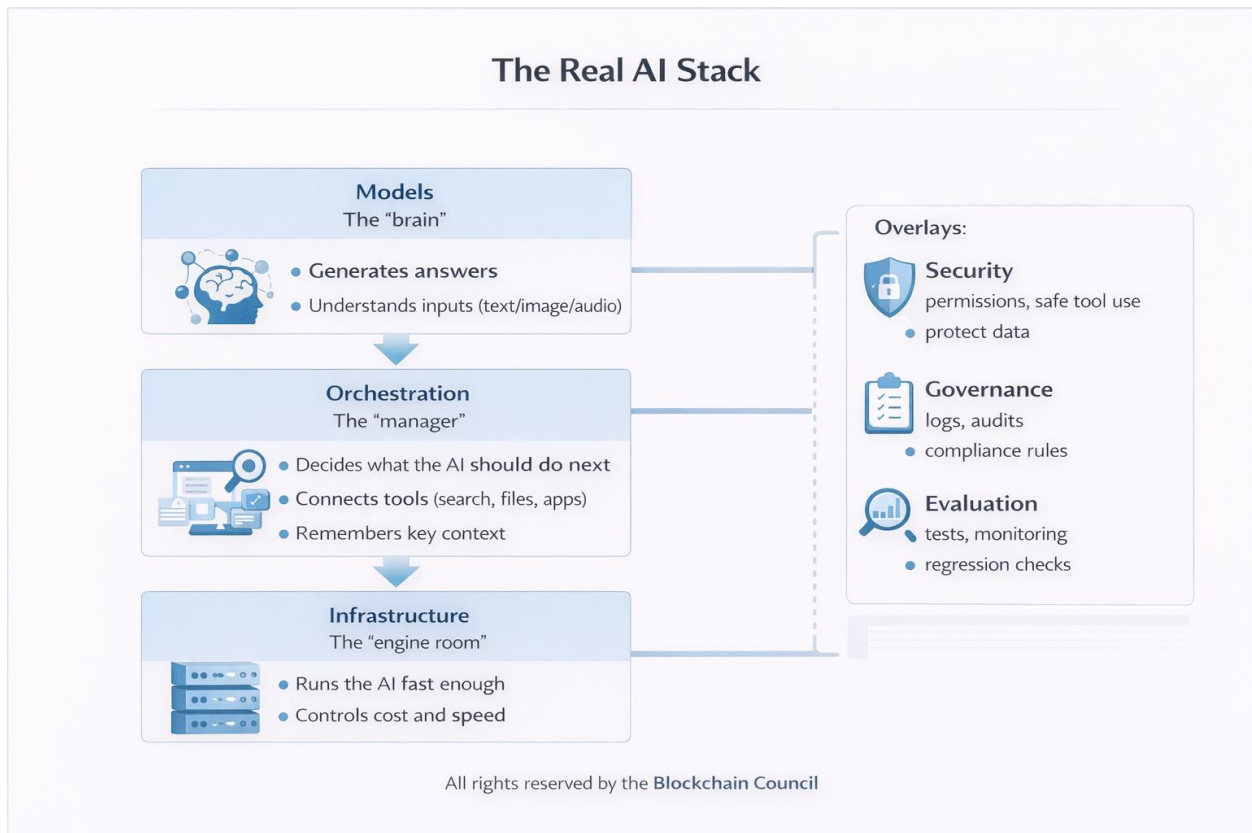
### Executive synthesis

The 2025 AI scene did not turn on a single headline jump. Instead, the year is best read as a set of converging threads that had been building for several cycles and finally met in shipped systems.

First, “reasoning” moved from a specialist feature into default behavior. Across major providers, the story shifted from “here is a chat model” to “here is a model family that can allocate more compute per request, run tools, and keep track of longer working state.” Second, multimodality widened from bolt-on image handling into end-to-end systems that treat text, images, and audio as part of one interaction loop. Third, “agents” moved from demo-first prototypes toward the hard edges of production: error rates, privilege boundaries, governance, monitoring, and the unglamorous work of connecting models to data and tools.

Across the year, the most visible progress came from systems work rather than a single benchmark jump. Tool use became a built-in expectation across multiple model families. Long-context pipelines matured and became normal in developer interfaces. Vendor ecosystems began to settle on common ways for models to connect to external tools, data sources, and other agents. The practical result was a more operational AI stack, with clearer lines between:

- Models as capability sources (language, vision, and other modalities)
- Orchestration as the control layer (agent frameworks, tool routing, memory policies)
- Infrastructure as the cost and latency layer (inference serving, training clusters, scheduling, monitoring)



Two strategic tensions shaped the year.

One tension was the widening gap between capability at the frontier and dependable autonomy. Consumer-facing “computer use” agents and coding agents improved a great deal, but “hands off” deployment remained limited by compounding error rates, prompt-injection risk, and brittle interactions with graphical interfaces and real-world constraints (time, payments, identity, policy boundaries). Product stories that framed 2025 as the breakout year for general agents met stronger skepticism in year-end coverage, including commentary in *The New Yorker*.

The second tension was the hardening of geopolitics around compute and regulation.

- In Europe, the EU AI Act’s staged timeline moved from paper to practice. Prohibited practices and AI literacy duties began applying on February 2, 2025. Obligations for general-purpose AI models began applying on August 2, 2025. Those dates forced product teams to treat compliance as a near-term delivery constraint rather than a future policy item.
- In the United States, the policy posture tilted toward accelerating domestic AI leadership. Executive Order 14179 landed in January 2025, and the White House

released an AI Action Plan in July 2025, while export controls and enforcement signals continued to change.

- In China, late-2025 proposals targeted emotionally interactive AI services, and industrial policy stressed domestic shares of semiconductor equipment for new capacity.

## The 2 Big AI Tensions in 2025

### Ability vs Trust

**Models got smarter** — agents still needed limits

- Capability jumped (reasoning, tools, multimodal)
- Real-world autonomy lagged (multi-step errors, safety risks)



Most “agents” worked best with *confirm steps*.

### Rules vs Resources

**AI raced forward while constraints tightened**

- Regulation & compliance moved from “coming soon” to “now”
- **Compute, cost,** and power became hard limits (chips, energy, grid delays)



Shipping AI meant balancing *policy deadlines* with *infrastructure reality*.

All rights reserved by the Blockchain Council

In parallel, infrastructure constraints became more visible than in prior years. NVIDIA’s roadmap advanced from Blackwell into Blackwell Ultra for the second half of 2025, and later introduced Rubin CPX as a GPU class aimed at massive-context inference. The underlying signal was that inference, not only training, was becoming a dominant cost center for reasoning-heavy and agent-style workloads.

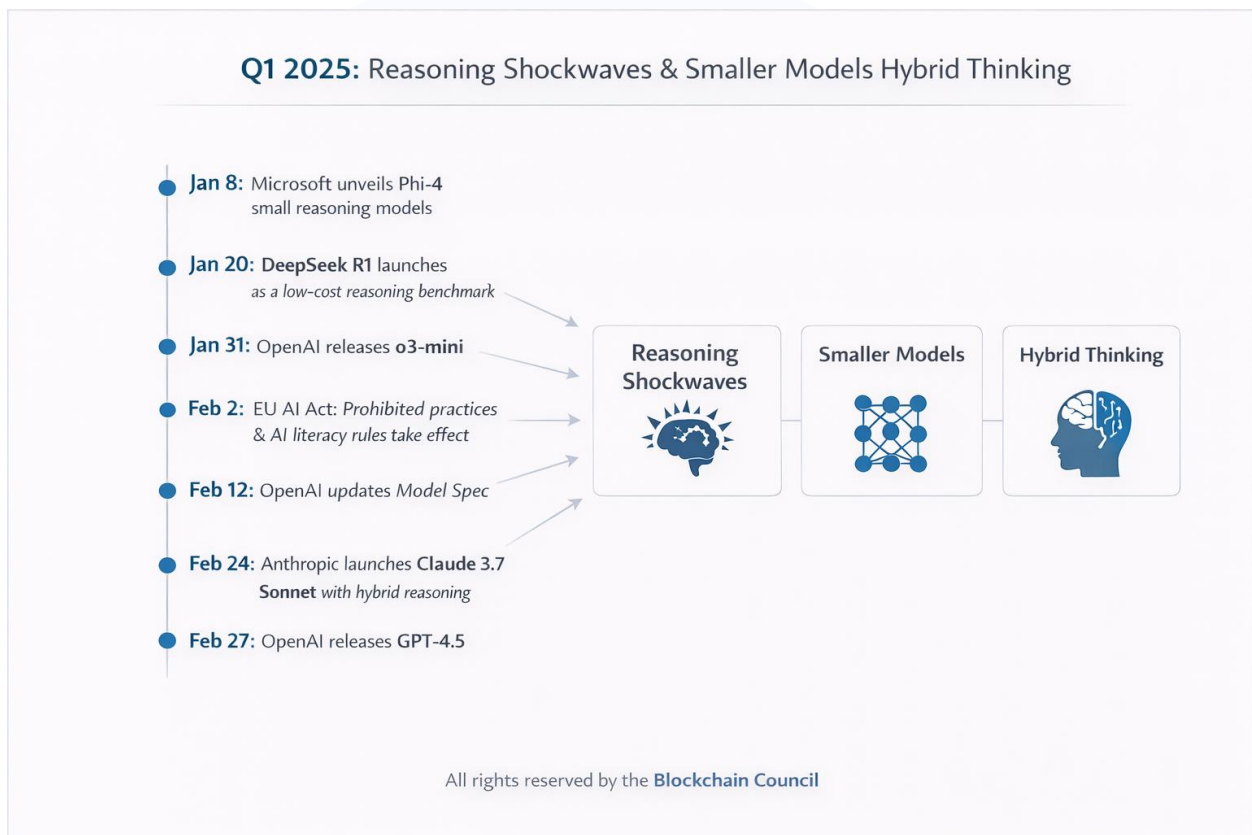
Energy and the grid moved from back-office detail into board-level planning. The International Energy Agency estimated data centers consumed about 415 TWh of electricity in 2024 (about 1.5% of global electricity use) and projected rapid growth toward 2030. Late-2025 reporting described developers leaning on on-site generation to bypass multi-year grid delays.

Within this context, 2025 looks like the year the field tried to turn “model capability” into “system capability,” and in the attempt exposed the bottlenecks: security, realistic evaluation, integration standards, and infrastructure at the level demanded by broad deployment.

## Frontier model releases and capability direction

The 2025 release cadence from leading providers shared a common structure: a “model family” was shipped together with an operating envelope, tool access, context length, safety and policy artifacts, and developer primitives, rather than as a stand-alone blob of weights.

That packaging matters for research and deployment. It changes what it means to compare models. It also changes what engineering teams must build around them. When tool access and long context become first-class, failure modes shift from single-turn mistakes toward multi-step drift, tool misuse, and subtle instruction conflict inside long working state.



## OpenAI: from scaled pretraining to tool-native reasoning families

OpenAI's 2025 path illustrates the year's broader move from "a new model" to "a new family plus a working environment."

Early 2025 brought GPT-4.5 as a research preview (February 27, 2025). It was positioned as a forward step in pretraining scale and post-training, stressing improved pattern recognition and generation without framing it as a pure reasoning specialist.

In parallel, OpenAI expanded its reasoning line:

- o3-mini (January 31, 2025) targeted cost-conscious STEM reasoning in both ChatGPT and the API.
- o3 and o4-mini (April 16, 2025) were described as reasoning-focused models with full tool capability, including web browsing and Python, and with emphasis on image input for multimodal reasoning.

In April 2025, OpenAI introduced the GPT-4.1 family in the API (April 14, 2025), pointing to long-context handling up to 1 million tokens and improved coding and instruction following. Long context became central to the 2025 framing of agent systems: persistence across multi-step plans, larger tool schemas, and repository-scale coding all become more plausible when the model can carry more state. At the same time, long context pushes teams into new problems:

- Retrieval drift, where relevant details get buried or overwritten
- Compaction errors, where summarization loses constraints that later matter
- Instruction conflict, where old goals collide with new ones and the model resolves them unpredictably

The mid-year inflection was GPT-5 (August 7, 2025), presented as the flagship with improvements across professional tasks. A notable part of the release story was evaluation in health through HealthBench.

HealthBench, introduced in May 2025, was a rubric-based benchmark built with physician input. Its design choice matters: instead of treating health tasks as trivia, it asks whether the model behaves in clinically sensible ways, including how it expresses uncertainty and how it handles safety-sensitive advice.

Late 2025 continued the "family plus artifacts" pattern:

- GPT-5.2 (December 11, 2025) was positioned for professional work and long-running agents.

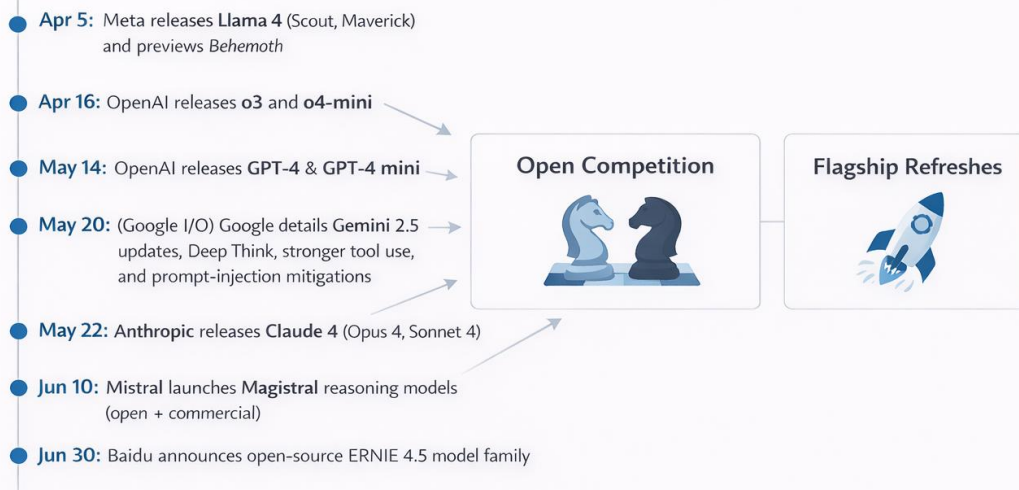
- GPT-5.2-Codex (December 18, 2025) was presented as an agentic coding model aimed at long-horizon software engineering, including stronger behavior on large code changes and stronger cybersecurity capability claims.

A strategic turning point was OpenAI’s move into open-weight releases. On August 5, 2025, OpenAI released gpt-oss-120b and gpt-oss-20b under Apache 2.0. This put OpenAI more directly into the open-model space that had been shaped by other firms, and it changed the baseline for on-prem deployment, customization, and sovereign AI strategies.

OpenAI also shipped open-weight safety tooling later in the year. The gpt-oss-safeguard models were described as policy-reasoning classifiers trained to label content against provided policies. The move signals that content governance is being treated as part of the build, not only a hosted service.

Across these releases, the throughline is that OpenAI treated 2025 as a year to make reasoning, tool use, and deployment patterns the default expectation rather than an optional tier.

### Q2 2025: Open Competition Heats Up, Flagships Refresh



## **Google: Gemini 2.5, computer use, and multimodal product surfaces**

Google's 2025 story, as visible through Gemini updates, pushed two axes: enhanced reasoning modes and the productization of computer use.

At I/O 2025, Google described updates to Gemini 2.5, including an experimental enhanced reasoning mode called Deep Think for Gemini 2.5 Pro. The same period also emphasized native audio output and safety safeguards, while framing Project Mariner as a pathway for computer use.

By October 2025, Google released the Gemini 2.5 Computer Use model via the API. It was described as a specialized model built on Gemini 2.5 Pro's visual understanding and reasoning, aimed specifically at UI interaction for web and mobile control tasks. This "special model for control" framing matters for deployment: it suggests that UI control is distinct enough to justify separate training and separate evaluation, rather than being a thin wrapper around a general chat model.

Google also leaned into interoperability. The Agent2Agent Protocol (A2A) was presented through Google's developer channels as a way for agents to communicate, exchange information securely, and coordinate across enterprise apps. In parallel, Model Context Protocol (MCP) adoption accelerated during 2025, supported by registry work and production guidance. The shared message is that tool and data hookup has become a main bottleneck for agent work, often more so than raw text generation quality.

## **Anthropic: hybrid reasoning and coding agents as product primitives**

Anthropic's most visible 2025 milestone was Claude 3.7 Sonnet (February 24, 2025), positioned as a "hybrid reasoning model" that can answer quickly or spend more time thinking, with API controls for how long it can think.

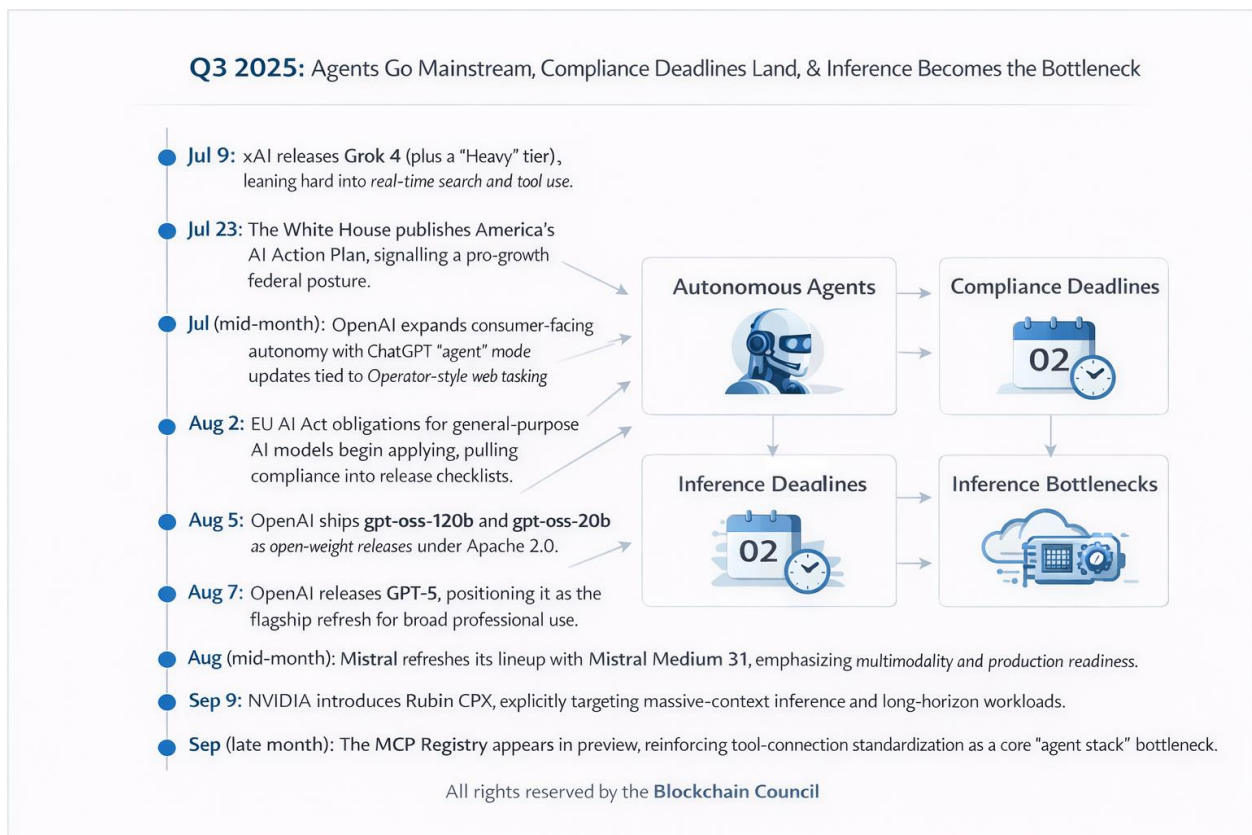
This framing fits a broader 2025 trend: inference-time scaling, where providers scale compute per query rather than only scaling parameters and training data. The main deployment implication is that cost and latency can be treated as adjustable knobs tied to accuracy and error rates.

Anthropic linked reasoning advances to workflow value via Claude Code, described as an agentic coding assistant with file editing, testing, CLI tool use, and GitHub workflows. The pairing of a reasoning model with a workflow-native agent product supports a pragmatic

thesis: near-term autonomy tends to show up first in tool-rich, bounded domains such as coding, data analysis, and internal enterprise workflows.

Anthropic’s work on MCP and engineering guidance about tool scaling also pointed to a systems constraint that got clearer in 2025. As agents connect to hundreds or thousands of tools, token overhead and tool selection become a performance and cost problem. This pushes architectures toward:

- Registries and servers that can describe tools without loading everything into every prompt
- Dynamic tool loading, where tool options are fetched only when needed
- Code execution gateways with clear privilege boundaries



## xAI: Grok 4 and Grok 4.1 with real-time search and agent tools

xAI’s Grok releases in 2025 emphasized real-time search, tool calling, and large context.

- Grok 4 was announced July 9, 2025 as a flagship model available to subscribers and via API, with an added “Heavy” tier.
- Grok 4.1 arrived in November 2025, described as a staged rollout with preference gains in live traffic, alongside a “Fast” variant aimed at tool calling and a 2M context window plus an Agent Tools API.

The Grok line illustrates a broader competitive theme: “real-time” is increasingly treated as a key product feature, whether through platform data access, built-in browsing, or agent frameworks that combine retrieval with action. It also illustrates the paired governance risk. Mainstream reporting (including Reuters) covered controversy around generated content and subsequent removals, underscoring that search-linked generation can move from benign to harmful content quickly if guardrails fail.

## **The open-weight ecosystem and diffusion of capability**

A defining feature of 2025 was that behavior close to the frontier increasingly appeared in open-weight or permissively licensed models. For many enterprises and labs, that narrowed the practical gap between “closed API frontier” and “self-hosted deployment,” especially where sovereignty, privacy, or cost rules out reliance on a single hosted provider.

## **OpenAI enters open weights: gpt-oss and safety adjuncts**

OpenAI’s gpt-oss release under Apache 2.0 (gpt-oss-120b and gpt-oss-20b) represented a notable shift from a primarily closed stance.

The positioning around tool use and deployment cost reflected the open ecosystem’s 2025 reality: open models are judged less on academic scores alone and more on operational metrics such as tool calling, latency on common hardware, and fit with agent orchestration stacks.

Later in 2025, OpenAI released gpt-oss-safeguard models as open weights. These were described as policy-reasoning classifiers trained to label content against given policies. The release is a sign that safety tooling is becoming part of the open ecosystem rather than a proprietary afterthought.

## **Mistral: reasoning and multimodal openness at multiple scales**

Mistral continued expanding a European open-model footprint.

- In June 2025, a reasoning model line called Magistral appeared in technical reporting, including an open-weights variant.
- By late 2025, Mistral announced “Mistral 3” as a family of open multimodal and multilingual models, including a mixture-of-experts flagship with 41B active parameters and 675B total parameters under Apache 2.0.
- Mistral also positioned “Mistral Medium 3.1” as a multimodal frontier-class model released in August 2025.

The practical takeaway is that, in 2025, open models were no longer only about small local assistants. Providers marketed open weights for professional contexts, multimodality, and reasoning tasks, backed by licenses that allow commercial use without heavy copyleft obligations.

## **Alibaba Qwen: trillion-parameter scale and applied benchmark framing**

Alibaba’s Qwen3-Max was announced as a model with over 1 trillion parameters. It was framed around code generation and agent work, and it arrived alongside multi-year investment commitments in AI infrastructure.

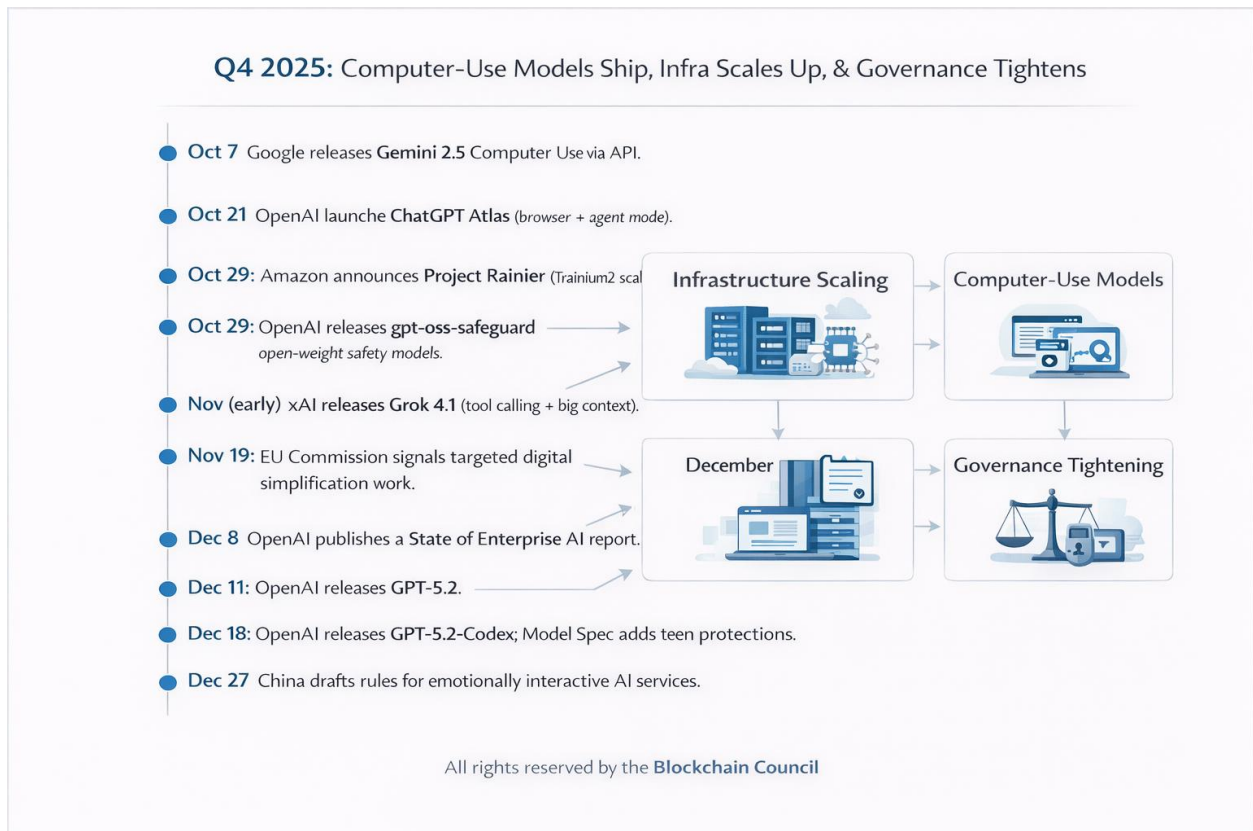
Qwen communications also stressed applied software evaluation. The official Qwen3-Max blog highlighted a SWE-Bench Verified score of 69.6 in the company’s reporting. The messaging reflects a wider 2025 shift: model marketing and selection increasingly depended on applied benchmarks that track real work artifacts, especially coding agents and long-horizon tool use.

## **DeepSeek: cost narratives and reinforcement learning as a competitive lever**

DeepSeek’s rise in early 2025 became a focal point in global discussion. Reuters reporting described DeepSeek’s claim that it trained a model at a fraction of the cost of leading US firms, using less-advanced chips, and that its reasoning model R1 was far cheaper to run than comparable offerings.

Even if cost figures require careful interpretation, the broader technical point fits the year’s direction: reinforcement learning, distillation, and inference-time scaling can produce large perceived capability gains without a linear rise in training compute, especially when models are tuned for multi-step reasoning and tool-heavy workflows.

Financial Times coverage also pointed to the “aha moment” narrative around reinforcement learning and self-improvement behavior. Regardless of how literal such language is taken, it speaks to the competitive takeaway many teams drew in 2025: post-training and inference-time compute budgets can be decisive, not only base pretraining scale.



## Reasoning as a systems property

In 2025, “reasoning” was less a single trick and more a bundle of design choices:

- Encouraging stepwise deliberation
- Allocating more compute per request at inference
- Using reinforcement learning to improve multi-step problem solving
- Using tool access so the model can call external verifiers (executors, search engines, test runners)

Because these choices show up in both product design and architecture, reasoning is best described as a systems property rather than a single model score.

## Adjustable deliberation budgets become mainstream

Three major provider lines converged on an explicit “thinking budget” idea.

- Anthropic’s hybrid reasoning framing for Claude 3.7 Sonnet exposed API controls over how long the model can think.
- OpenAI’s o-series stressed reasoning behavior and the ability to combine it with tools, including Python.
- Google’s Deep Think mode for Gemini 2.5 Pro was presented as an enhanced reasoning setting, with descriptions of extended thinking and reinforcement learning techniques.

The product-level conclusion is that single-pass chat is no longer the default for high-stakes tasks. Instead, developers can trade latency and cost for better correctness and fewer errors, and allocate deeper thinking only when a task warrants it.

## **Tool access changes what “reasoning” means in practice**

Tool access moved from a side feature to a core part of what “reasoning” meant in practice.

OpenAI’s o4-mini reporting on AIME performance with a Python interpreter illustrated the point sharply, including near-perfect pass rates under tool-assisted conditions. The interpretation is not “the model is now a perfect mathematician,” but that the strongest “reasoning systems” often look like a model-plus-tool loop:

- The model can write and run code
- It can check intermediate results
- It can iterate rather than commit to the first attempt

Once tools enter the loop, evaluation becomes harder. Benchmarks designed for closed-book reasoning can become less informative if a model can run an interpreter or retrieve outside information. That tension helped drive the year’s shift toward applied benchmarks with explicit tool environments.

## **Inference-time scaling and reinforcement learning narratives gain weight**

DeepSeek’s R1 became widely discussed as an example of reinforcement learning-driven gains with limited human feedback, with reporting stressing both cost and qualitative shifts in reasoning-like behavior.

The broader 2025 lesson is that teams treated inference-time scaling and RL-based post-training as first-class levers. Even when base model scale remained important, developers increasingly treated “how the model is used” (deliberation, tools, verification loops) as a large part of final task success.

## **Agents, computer use, and interoperability standards**

If 2024 made “agents” plausible, 2025 made them operationally contentious. The year delivered major releases in computer use and agent frameworks, while also making the gap between demos and dependable autonomy harder to ignore.

### **Computer use becomes a model product**

Two releases illustrate the shift from “agent demo” to “computer use as a product category.”

- OpenAI’s Operator, introduced earlier, was updated in July 2025 and presented as part of ChatGPT as “ChatGPT agent,” bringing agent mode into the main ChatGPT experience for multi-step tasks that involve websites.
- Google’s Project Mariner framed computer use as a core capability heading into the Gemini API, and the Gemini 2.5 Computer Use model (October 2025) formalized the idea as a specialized control model.

The key research implication is that UI control appears to require distinct training, distinct safety policies, and distinct evaluation. It is not just “chat plus clicks.”

### **The agent reliability gap remains the main constraint**

Despite progress, broad “everyday life” transformation from general agents did not arrive at the pace forecast in some earlier product stories. Year-end commentary, including in The New Yorker, argued that agents still struggled outside narrow domains.

This skepticism fits a technical reading:

- Multi-step tasks compound error. A small per-step failure rate becomes a large end-to-end failure rate as steps add up.
- Graphical interfaces add uncertainty channels: visual grounding errors, UI timing issues, hidden state, and action side effects.
- Real-world constraints (payments, identity checks, policy limits, timeouts) impose boundaries that models do not naturally track without careful tool design.

As a result, autonomy looked less like a single breakthrough and more like a long engineering program that will need progress in:

- Planning and replanning
- Verification and test harnesses
- Memory and state management over long horizons
- Tool and data standards
- Security defenses that assume adversarial input

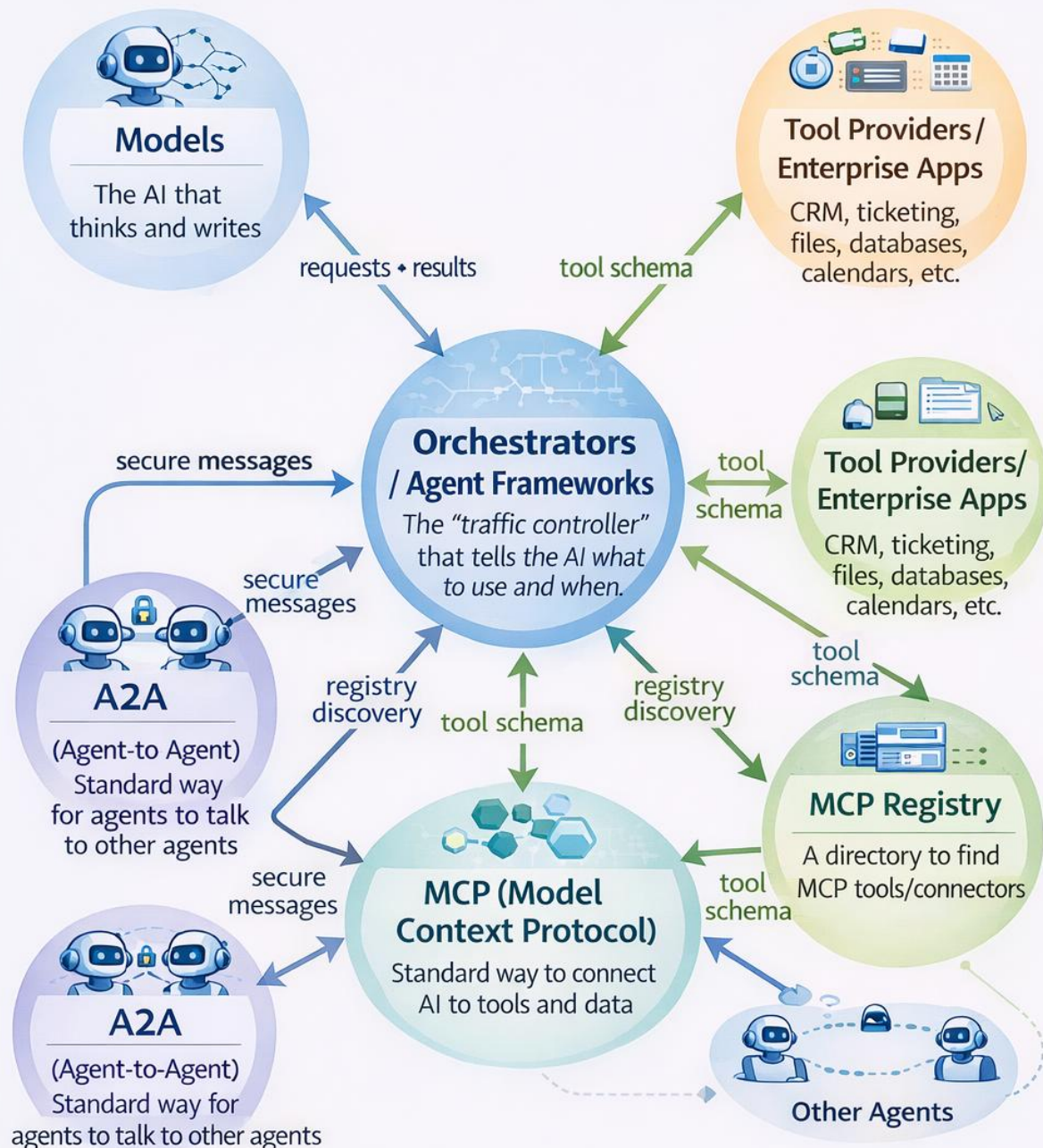
## **Interoperability standards become first-order work**

Two standards-oriented efforts became especially salient in 2025.

- Model Context Protocol (MCP) matured as an open protocol for connecting models to tools and data sources. The MCP Registry launched in preview in September 2025, with a stated path toward general availability.
- Agent2Agent (A2A) was announced as a protocol for secure agent-to-agent communication and coordination across platforms.

The strategic point is that agent performance is often gated by tool hookup more than by token-level language skill. Standards reduce repeated custom plumbing, enable a market for shared connectors, and create a compliance surface that can be audited and secured.

## AI Interoperability Map



Standards reduce custom integrations and make permissions and auditing easier.

## The browser as an agent container

OpenAI introduced ChatGPT Atlas on October 21, 2025, presenting it as a browser with ChatGPT built in, including agent mode.

Placing the agent in the browser changes both capability and risk. It moves the model to the point of action on the web, where it can see and act on pages directly. It also increases exposure to adversarial content.

Reporting around AI browsers (including IT Pro coverage of ChatGPT Atlas) stressed prompt injection as a persistent and serious risk. The security lesson is that an agent operating on untrusted web text must be treated as an interpreter of adversarial input. Defenses then need to look less like “prompt tips” and more like layered engineering:

- Clear instruction hierarchy and separation between system goals and page content
- Permission boundaries for tools (especially anything that can spend money or change accounts)
- Sandboxed execution and audit logs
- Continuous red-teaming and fast patch cycles

## Multimodality: from vision add-ons to native audio and video generation

Multimodality broadened in 2025 along two tracks:

- Models that can perceive and act across modalities (text, images, audio)
- Generators that can synthesize media (video) with more control and continuity

These tracks are related but not identical. One is about interaction and grounding; the other is about synthesis quality and control.

## Native audio and multimodal interaction

Google’s Gemini 2.5 updates stressed native audio output and a more natural conversational experience. Baidu’s investor materials described ERNIE 5.0 as its first “native omni-modal foundation model,” framing multimodality as a default architecture rather than an ensemble of separate subsystems.

From a research perspective, the interesting claim is not “it can accept more input types,” but that the model can maintain coherent state across channels and apply reasoning grounded in perceptual input. That matters for:

- UI control, where vision guides action
- Assistive settings, where audio output is part of the interface
- Workflows that mix documents, images, and spoken instructions

## **Text-to-video and controllable media generation**

On the synthesis side, 2025 saw rapid iteration in text-to-video.

- Runway introduced Gen-4, stressing consistent characters, locations, and objects across scenes, and focusing on control and continuity.
- Google DeepMind’s Veo line continued to stress video quality and native audio generation.
- Coverage of later Runway updates described gains in realism and prompt adherence, while noting remaining issues such as object permanence and causal consistency.

The research significance is that video generation is moving from “short clip novelty” toward workflows that demand:

- Continuity across shots
- Camera control
- Stable identity for characters and scenes

These are table stakes for many professional uses in marketing, entertainment, and training content. At the same time, higher-quality synthetic media raises governance pressures around provenance and misuse.

## **Evaluation: applied benchmarks and domain-grounded rubrics**

Evaluation continued shifting away from pure academic tests toward applied tasks and domain rubrics. The driver was not only scientific rigor; it was product need. Teams needed evaluations that match real failure modes: long-horizon work, tool use, and safety-relevant behavior.

## **HealthBench and rubric-based domain evaluation**

OpenAI introduced HealthBench in May 2025 as a benchmark meant to evaluate model behavior in realistic health conversations using physician-written rubrics, with thousands of scenarios and criteria weighted by clinical judgment.

OpenAI later pointed to GPT-5 performance on HealthBench as part of its release story. This pattern shows how domain benchmarks became both a research tool and a governance artifact.

Methodologically, rubrics help capture behavior that a simple accuracy score misses, such as:

- How uncertainty is expressed
- Whether a model suggests appropriate triage
- Whether it avoids unsafe advice
- Whether it keeps track of patient context across turns

## **Software engineering benchmarks as a main competitive arena**

Software engineering evaluation became a central battleground. SWE-Bench continued rising in prominence, with 2025 updates pointing to agent-oriented evaluation work such as mini-SWE-agent and newer evaluations such as CodeClash.

Vendors referenced SWE-Bench Verified scores directly in marketing, including Qwen3-Max. The focus is consistent with product reality. Coding domains have properties that make them unusually evaluation-friendly:

- Executable verification through tests
- Structured artifacts like repositories and issues
- Measurable outcomes like merged pull requests and resolved tickets

At the same time, coding evaluations surface a deployment truth: higher throughput often shifts work from writing to reviewing. A tool that produces more code can also produce more mistakes that must be caught, making quality control and testing practices part of “model performance.”

## **Tool-assisted math and the meaning of benchmark dominance**

OpenAI's o4-mini reporting on AIME results with Python access illustrates a broader evaluation challenge. Tool access changes the task. If a model can run a solver or execute code, the benchmark is no longer measuring the same thing it measured in a closed setting.

As a result, applied evaluation increasingly depends on explicit environment specification:

- Which tools are allowed
- What the permission scope is
- How verification works
- Whether the system can iterate and self-correct

## **Safety, security, and governance**

As agents and multimodality spread, risk moved from abstract debate to operational detail. Safety work in 2025 looked less like slogans and more like engineering:

- Threat models for prompt injection n- Tool privilege boundaries
- Age-sensitive product rules
- Policy processes that can stand up to audit and regulation

## **Prompt injection and agent security**

As models gained browser control, prompt injection moved from a niche research threat into a product-grade concern.

Reporting around AI browsers, including coverage tied to ChatGPT Atlas, treated prompt injection as persistent. The core issue is structural: the model consumes untrusted content from the web, and that content can include instructions designed to hijack the agent's goals.

This makes security for agents resemble security for browsers and email clients. Eliminating the issue entirely may be unrealistic, so teams are pushed toward layered defenses and continuous testing.

## **Policy-encoded safety models and open safety tooling**

The release of gpt-oss-safeguard as open weights points to an approach where safety is treated as policy reasoning that can be audited and deployed in different settings. Rather than a single opaque hosted moderation layer, an organization can build moderation, triage, and compliance workflows that fit its own policy rules.

## **Youth protections, interaction intensity, and mental health**

Late 2025 brought a clearer focus on age and vulnerability.

OpenAI updated its Model Spec in December 2025 with teen protections, describing “U18 Principles” meant to guide age-appropriate behavior. The update landed amid heightened scrutiny of mental health impacts and litigation narratives, including Washington Post reporting on cases involving teenagers.

In parallel, Reuters reported that OpenAI planned to allow mature content for adult-verified users starting in December 2025. The move implies a bifurcated approach:

- Stricter safeguards for minors
- More permissive handling for verified adults

China’s regulators also signaled concern with psychologically intense, emotionally interactive systems. Reuters reported draft rules in December 2025 aimed at AI services with human-like interaction, including requirements tied to excessive use and psychological risk.

Across jurisdictions, the trend points toward governance that is sensitive to user age, vulnerability, and interaction intensity rather than only to content categories.

## **Infrastructure, compute, and energy: inference becomes the main scaling constraint**

By 2025, the field’s main scaling constraint was no longer only training. Reasoning-heavy use, long context, and agent loops shift the cost center toward inference.

That shift shows up in hardware roadmaps, in open-source serving frameworks, and in board-level attention to energy and grid access.

## **GPU roadmaps and inference-first architectures**

NVIDIA’s 2025 announcements stressed inference needs.

- Blackwell Ultra was presented as the next evolution of the Blackwell “AI factory” platform, expected in the second half of 2025.
- NVIDIA introduced Dynamo in March 2025 as an open-source distributed inference framework, reporting large throughput gains when serving DeepSeek-R1 on Blackwell.

- In September 2025, NVIDIA announced Rubin CPX as a GPU class designed for massive-context inference, explicitly targeting million-token processing for domains such as coding and generative video.

Rubin CPX is conceptually aligned with the “agent era” argument. Long context and multi-step deliberation increase inference demand sharply. Serving frameworks and hardware that can handle million-token workloads become strategic, not only technical.

## **Hyperscaler custom silicon and mega-clusters**

Hyperscalers continued building custom training and serving hardware.

Amazon’s Project Rainier, announced in October 2025, showed this direction. Reuters reported that Project Rainier would use nearly 500,000 Trainium2 chips across multiple US data centers, and that Anthropic was expected to scale to over one million Trainium2 chips by the end of 2025.

The significance is not only unit cost. It is also supply chain control and capacity planning in a market where demand for accelerators outstrips supply.

## **Energy and grid constraints enter the AI planning loop**

The IEA estimated that data centers consumed around 415 TWh of electricity in 2024 (about 1.5% of global electricity consumption) and projected substantial growth through 2030.

Late-2025 reporting described multi-year delays for grid connections and a turn toward on-site generation, including gas turbines and diesel generators, as developers sought to bring new AI data centers online faster.

This reality feeds back into model and system design. If inference is energy-hungry, then distillation, mixture-of-experts routing, and context compaction become not only performance work but also energy and cost work.

## **Policy and geopolitics: compliance timelines and compute sovereignty**

In 2025, compute access and deployment rules became inseparable. Hardware supply chains and regulatory compliance shaped what could be trained, where it could be served, and how it could be shipped into products.

## **European Union: AI Act obligations move into practice**

The EU AI Act's staged timeline is among the most consequential governance developments for global deployment.

- Prohibited practices and AI literacy duties began applying on February 2, 2025.
- Obligations for general-purpose AI models began applying on August 2, 2025.

These dates forced companies operating in Europe to treat compliance as ongoing work, covering documentation, risk management, transparency duties, and coordination with emerging oversight bodies.

## **United States: pro-growth posture and shifting export controls**

In January 2025, the White House issued Executive Order 14179, titled "Removing Barriers to American Leadership in Artificial Intelligence." In July 2025, the White House released "America's AI Action Plan," describing over 90 federal policy actions across themes such as innovation acceleration, infrastructure, and international diplomacy and security.

Export control policy remained central. In May 2025, the US Department of Commerce announced rescission of a Biden-era AI diffusion rule, signaling a change in approach. Congressional Research Service reporting later summarized evolving export controls and related actions through 2025, reinforcing how security policy shapes global compute availability and supply chain planning.

## **China: regulation of emotionally interactive AI and chip industrial policy**

China's late-2025 draft rules targeting AI with human-like interaction signaled a regulatory focus on psychological and social impacts, including requirements tied to user behavior monitoring and intervention in cases of excessive use or emotional dependence.

In parallel, Reuters reported a policy requiring at least 50% domestically produced equipment for new semiconductor manufacturing capacity, reflecting a push for self-sufficiency under sanctions.

Taken together, these policy moves show the new coupling: training and serving frontier models depends on hardware supply chains, and the acceptability of deploying companion-style or agent-style products depends on jurisdiction-specific rules.

## Enterprise adoption and the productivity debate

Enterprise adoption in 2025 entered a more measurable phase and a more contested one.

OpenAI's "State of enterprise AI" report, released in December 2025, described survey results from 9,000 workers across almost 100 enterprises, paired with usage data from enterprise customers. The report stressed adoption patterns, tool use distribution, and time savings claims.

At the same time, external reporting (including Axios) described tensions inside firms about uneven uptake across roles. The emerging question is not whether AI tools can help in the abstract, but under what conditions they create durable productivity gains without harming quality.

Coding agents, document drafting, and research assistants can raise output, but they can also add review burden and create new quality-control work. In that sense, "adoption" is partly an organizational design problem: incentives, trust, and process change can matter as much as model choice.

A careful reading of 2025 is therefore that enterprise AI shifted from experimentation to differentiated capability:

- High-intensity users and well-instrumented workflows can extract meaningful value.
- Median adoption can lag, creating internal divides and uneven ROI stories.

## Chronology of notable 2025 milestones

The timeline below covers widely documented releases and governance milestones that shaped practical AI work during 2025. It is not exhaustive, but it includes the events that most clearly influenced capability, deployment constraints, and regulation.

- **January 2025**
  - OpenAI released **o3-mini** (January 31), framed as a cost-conscious reasoning model for ChatGPT and the API.
  - The United States issued **Executive Order 14179** (January 23) on "Removing Barriers to American Leadership in Artificial Intelligence."

- The **DeepSeek R1** narrative gained global visibility, including claims about training cost and market impact.
- **February 2025**
  - EU AI Act rules on **prohibited practices** and **AI literacy** began applying (February 2).
  - Anthropic announced **Claude 3.7 Sonnet** (February 24) as a hybrid reasoning model.
  - OpenAI introduced a **GPT-4.5** research preview (February 27).
- **March 2025**
  - NVIDIA discussed **Blackwell Ultra** as the next step in its Blackwell platform roadmap for the second half of 2025.
  - NVIDIA introduced **Dynamo** (March 18), an open-source distributed inference framework.
  - The IEA published “**Energy and AI**” analysis, bringing data center electricity demand into mainstream AI planning.
- **April 2025**
  - OpenAI introduced the **GPT-4.1** family in the API (April 14), stressing improved coding and up to **1M token** context.
  - OpenAI introduced **o3** and **o4-mini** (April 16), stressing reasoning with full tool access.
  - Meta previewed an **API for Llama models**, signaling further product packaging for open models.
- **May 2025**
  - OpenAI introduced **HealthBench** (May 12), a physician-rubric benchmark for health conversations.
  - Google discussed **Gemini 2.5** updates and **Deep Think** at I/O 2025.
- **June 2025**
  - Mistral’s reasoning line **Magistral** appeared in technical reporting, reflecting continued growth in open reasoning offerings.
- **July 2025**
  - OpenAI updated Operator and folded it into **ChatGPT agent mode** inside ChatGPT, enabling multi-step website tasks.
  - The White House released **America’s AI Action Plan** (July 23).
  - xAI announced **Grok 4** (July 9), stressing real-time search and tool calling.
- **August 2025**
  - EU AI Act obligations for **general-purpose AI models** began applying (August 2).
  - OpenAI released **gpt-oss** open weights (August 5) under Apache 2.0.
  - OpenAI introduced **GPT-5** (August 7), pointing to broad professional gains and HealthBench.

- Mistral released **Mistral Medium 3.1** as a multimodal frontier-class model (August).
- **September 2025**
  - NVIDIA announced **Rubin CPX** (September 9) for massive-context inference.
  - Alibaba released **Qwen3-Max** with over **1 trillion parameters** (reporting September 23–24) and stressed coding and agent benchmarks.
  - The **MCP Registry** launched in preview, indicating maturing tool integration infrastructure.
- **October 2025**
  - Google released the **Gemini 2.5 Computer Use** model via API (October 7).
  - OpenAI introduced **ChatGPT Atlas** (October 21), a browser with ChatGPT and agent mode.
  - Amazon announced **Project Rainier** powered by **Trainium2** (October 29), designed to support Anthropic’s Claude at scale.
  - OpenAI published **gpt-oss-safeguard** models and supporting technical material (October 29).
- **November 2025**
  - xAI released **Grok 4.1**, including variants that stressed improved reasoning, tool calling, and a **2M** context window.
  - The EU Commission signaled targeted amendment proposals in a broader digital simplification framing (November 19).
- **December 2025**
  - OpenAI introduced **GPT-5.2** (December 11) and **GPT-5.2-Codex** (December 18), aimed at long-running agent work and coding workflows.
  - OpenAI updated its **Model Spec** with teen protections (December 18).
  - OpenAI released an enterprise AI report (December 8) based on worker survey results and usage data.
  - China issued **draft rules** for emotionally interactive AI services (December 27) and reported industrial policy signals on domestic semiconductor equipment shares (reporting December 30).

## Research interpretation and implications for 2026

By the end of 2025, the field looked capability-rich but autonomy-limited. Models improved sharply as co-workers inside well-instrumented domains, while “open-world” agents remained difficult to trust without heavy containment.

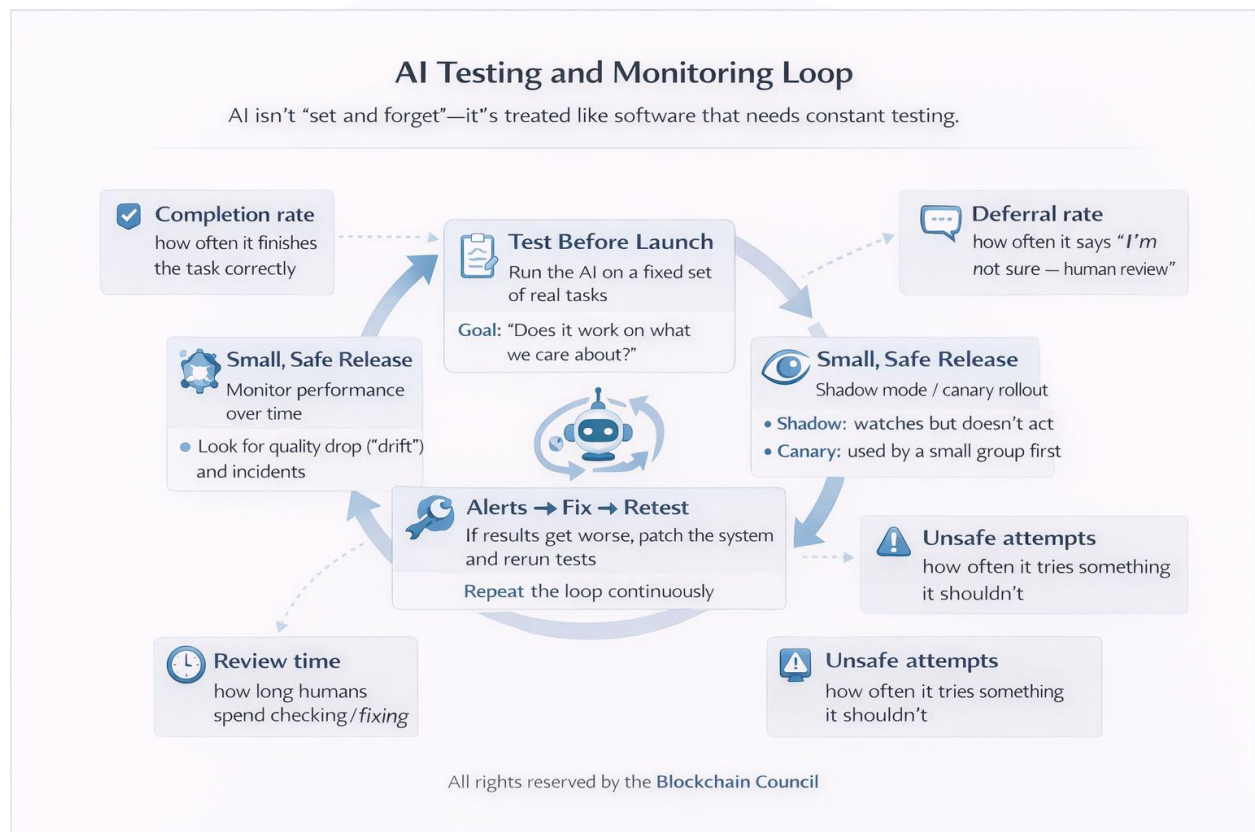
Three implications follow from the year’s main threads.

## Agents will be judged by security and verification as much as by raw skill

Prompt injection, tool privilege boundaries, and safe action execution are now first-order problems. Browser-embedded agents and UI control models widen the risk surface. Teams are pushed toward a security stance that resembles web security work:

- Continuous red-teaming
- Clear permissioning and compartmentalization
- Monitoring and rapid incident response
- Evaluation that assumes adversarial input

The 2025 shift here is not that these risks are new; it is that the product direction (computer use, browsers, agent modes) made them unavoidable.



## Benchmarking will keep moving toward work artifacts and domain rubrics

HealthBench and SWE-Bench showed why applied evaluation matters.

- In health, rubrics capture safety-relevant behavior that a score alone misses.
- In software, repositories and tests provide concrete end states and reduce ambiguity.

As tool use becomes normal, evaluation will also need better environment definition. “The model got X score” is less meaningful if it is unclear which tools were allowed and how much iteration was permitted.

## **Infrastructure constraints will shape research priorities**

The emergence of Rubin CPX and Dynamo points to inference as a central constraint. Long context, multi-step thinking, and agent loops push compute and energy needs upward. That pressure will shape work on:

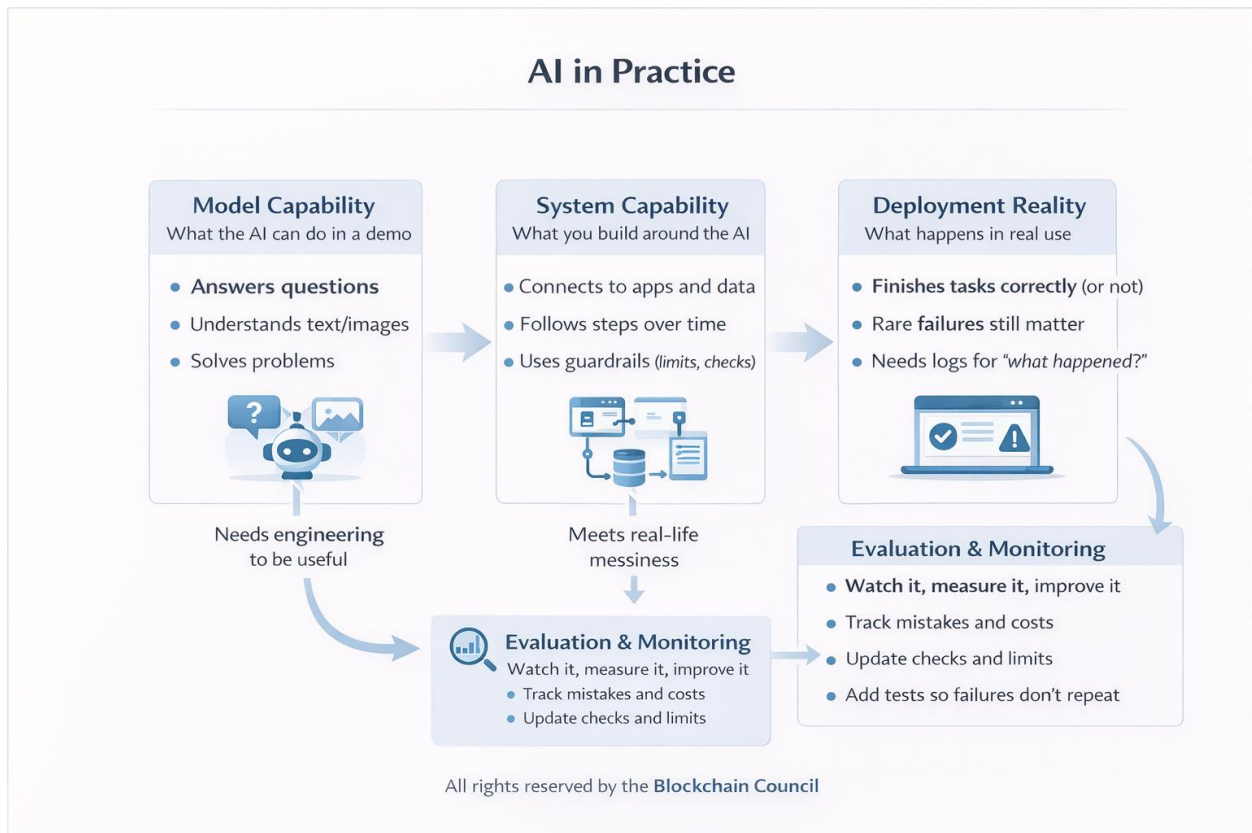
- Distillation and smaller student models
- Mixture-of-experts routing
- Context compaction and better memory policies
- Serving systems that can handle million-token workloads

Energy and grid constraints add another feedback loop. If new data centers face multi-year grid delays, then the cost of running large inference workloads becomes a first-class product constraint.

Taken together, 2025 reads as a year of integration and confrontation with constraints. The field showed that reasoning, multimodality, and tool use can be shipped at scale. It also showed that dependable autonomy requires progress in security, verification, standards, and infrastructure. At year end, models looked increasingly competent as assistants and co-workers inside bounded work settings, while still requiring careful containment and orchestration to act as agents in open settings.

## **Research interpretation and implications for 2026**

By the end of 2025, the field looked capability-rich but autonomy-limited. Models improved sharply as co-workers inside well-instrumented domains, while “open-world” agents remained difficult to trust without heavy containment.



Three implications follow from the year’s main threads.

## Agents will be judged by security and verification as much as by raw skill

Prompt injection, tool privilege boundaries, and safe action execution are now first-order problems. Browser-embedded agents and UI control models widen the risk surface. Teams are pushed toward a security stance that resembles web security work:

- Continuous red-teaming
- Clear permissioning and compartmentalization
- Monitoring and rapid incident response
- Evaluation that assumes adversarial input

The 2025 shift here is not that these risks are new; it is that the product direction (computer use, browsers, agent modes) made them unavoidable.

## **Benchmarking will keep moving toward work artifacts and domain rubrics**

HealthBench and SWE-Bench showed why applied evaluation matters.

- In health, rubrics capture safety-relevant behavior that a score alone misses.
- In software, repositories and tests provide concrete end states and reduce ambiguity.

As tool use becomes normal, evaluation will also need better environment definition. “The model got X score” is less meaningful if it is unclear which tools were allowed and how much iteration was permitted.

## **Infrastructure constraints will shape research priorities**

The emergence of Rubin CPX and Dynamo points to inference as a central constraint. Long context, multi-step thinking, and agent loops push compute and energy needs upward. That pressure will shape work on:

- Distillation and smaller student models
- Mixture-of-experts routing
- Context compaction and better memory policies
- Serving systems that can handle million-token workloads

Energy and grid constraints add another feedback loop. If new data centers face multi-year grid delays, then the cost of running large inference workloads becomes a first-class product constraint.

Taken together, 2025 reads as a year of integration and confrontation with constraints. The field showed that reasoning, multimodality, and tool use can be shipped at scale. It also showed that dependable autonomy requires progress in security, verification, standards, and infrastructure. At year end, models looked increasingly competent as assistants and co-workers inside bounded work settings, while still requiring careful containment and orchestration to act as agents in open settings.

## **Sector forecasts for 2026**

The broad predictions above become clearer when mapped to sectors. The goal here is not to claim uniform outcomes across industries. It is to identify where constraints and incentives make specific changes more likely.

## Enterprise operations in 2026

Enterprise adoption will continue to grow, but the center of spend is likely to move from experimentation to integration.

### What is likely to become common

- **Workflow-level deployments**
  - AI embedded in ticketing, CRM, knowledge bases, document systems, and code review tools
- **Fewer, deeper use cases**
  - Concentration on a small number of workflows that justify integration effort
- **Measurement-first deployments**
  - Teams define success metrics before scaling

### What is less likely to become common

- Unbounded agents with broad permissions across many systems
- Fully automated decision-making in high-liability areas without clear audit trails

### Practical indicators

- More AI projects include:
  - baseline measurement
  - controlled rollouts
  - audit logs
  - rollback plans

## Software development and IT operations in 2026

Software and IT will remain among the most favorable domains for automation because output is verifiable.

### Likely shifts

- More use of agents that:
  - open and triage issues
  - propose patches
  - run tests
  - draft pull requests
- More emphasis on:
  - security scanning
  - dependency analysis

- infrastructure-as-code checks

### **Constraints that shape outcomes**

- Code correctness is measurable, but the cost of wrong changes can be high.
- Organizations will therefore favor agents that propose changes and run tests, but still require review.

### **Indicators**

- Higher share of code changes created by AI but merged only after test and review gates
- Growth of tooling that logs model decisions and patch provenance

## **Customer support, sales, and service in 2026**

Support and sales workflows are attractive because they are repetitive and measurable. But they are also sensitive to tone, policy, and privacy.

### **Likely shifts**

- More “draft response” systems with:
  - retrieval grounding
  - policy checks
  - escalation triggers
- More automation in:
  - categorization
  - routing
  - summarization of cases

### **Constraints**

- Privacy and compliance are hard.
- Incorrect advice can produce reputational and legal costs.

### **Indicators**

- Higher use of AI to prepare responses, not to send them without approval
- Increased use of escalation rules tied to uncertainty and policy risk

## **Healthcare and life sciences in 2026**

Healthcare is a high-interest domain and a high-liability domain. In 2026, growth is likely to continue in clinical documentation and administrative tasks, with more caution in direct clinical decision support.

### **Likely shifts**

- Wider use in:
  - clinical note drafting
  - patient message drafting with review
  - coding and billing assistance
  - triage suggestions that are reviewed
- More attention to:
  - rubric-based evaluation
  - uncertainty expression
  - safety guardrails

### **Constraints**

- Regulatory and liability considerations favor human review.
- Medical advice requires calibrated uncertainty and strong deferral behavior.

### **Indicators**

- More health deployments emphasize documentation and workflow efficiency
- Stronger evaluation requirements for any system that influences triage or advice

## **Education and training in 2026**

Education is likely to see deeper integration, but also more conflict over assessment integrity.

### **Likely shifts**

- AI used for:
  - tutoring support
  - practice generation
  - feedback on drafts
  - administrative work
- More institutional changes:
  - assessment redesign to reduce simple outsourcing
  - emphasis on process, oral defense, and in-class work

### **Constraints**

- Unequal access creates equity concerns.
- Detection is limited; prevention focuses on assessment design.

### **Indicators**

- More schools publish AI usage policies
- More curriculum focuses on verification and reasoning about outputs

## **Finance and insurance in 2026**

Finance will adopt AI broadly, but strict controls will remain normal.

### **Likely shifts**

- AI used for:
  - document processing
  - compliance summaries
  - customer communications drafting
  - fraud investigation support
- Stronger requirements for:
  - audit logs
  - explainability at the system level
  - data governance

### **Constraints**

- Model risk management practices extend to AI.
- Regulators and auditors demand evidence.

### **Indicators**

- More AI deployments framed as “decision support with logs” rather than automated decisions

## **Media, marketing, and synthetic content in 2026**

Synthetic content tools will improve and will be used more broadly, but trust and provenance become central.

### **Likely shifts**

- Higher use of synthetic video and voice for:
  - marketing content
  - training materials

- localization
- Stronger demand for:
  - provenance metadata
  - consent and usage controls

### **Constraints**

- Misuse risks increase as realism improves.
- Platforms face pressure to label or manage synthetic content.

### **Indicators**

- More platform features for provenance
- More disputes about impersonation and consent

## **Government and public services in 2026**

Public sector use will expand unevenly, constrained by procurement rules and accountability demands.

### **Likely shifts**

- AI used for:
  - document drafting and summarization
  - internal search and retrieval
  - translation and accessibility
- Stronger controls on:
  - public-facing advice
  - enforcement-related uses

### **Constraints**

- Auditability and fairness concerns are central.
- Governance requirements are explicit.

### **Indicators**

- More pilots become sustained programs only where evaluation and logging are strong

## **Robotics, manufacturing, and physical AI in 2026**

Physical AI will progress, but deployment remains constrained by safety and hardware integration.

## Likely shifts

- More AI used for:
  - inspection
  - predictive maintenance support
  - warehouse picking in constrained environments
  - assisted teleoperation

## Constraints

- Physical systems require safety certification, redundancy, and reliable sensing.
- Action errors have immediate consequences.

## Indicators

- Growth in constrained, industrial deployments rather than open-world household robots

## Consumer tools and everyday use in 2026

Consumer tools will shift toward richer interfaces and more embedded assistance, but the most successful products will minimize risk and friction.

## Likely shifts

- Voice and image become default inputs in more tools
- Agents appear as:
  - shopping helpers
  - travel planners
  - email and calendar assistants

## Constraints

- Payment and identity risks force strict permissions
- Users have low tolerance for high-effort configuration

## Indicators

- More consumer agent products emphasize “suggest and confirm” rather than full automation

# AI Autonomy Levels

## (What It's Allowed to Do)

### Tier 1 — Read-only

Looks, searches, summarises

Examples: *read docs, answer questions, pull info from a database*



### Tier 2 — Draft-only

Prepares work, but can't send or *change anything*

Examples: *write an email draft, draft a report, propose a plan, suggest a fix*



### Tier 3 — Submit with approval (default for risky actions)

AI can click "**submit**" only after a human approves

Examples: *purchases, refunds, changing passwords, editing permissions, sending customer replies, booking travel*



### Tier 4 — Submit without approval (rare)

AI acts on its own, only in *low-risk, reversible* tasks

Examples: *schedule a low-impact reminder, run a read-only report, apply a reversible formatting change*

**Money, identity, and irreversible actions stay on Tier 3 by default.**

## **Market structure in 2026: consolidation and specialization happen at the same time**

2026 is likely to show a two-track market structure.

- **Consolidation** around a small number of foundation model providers and major platforms
- **Specialization** among tool providers, agent frameworks, evaluation firms, and security vendors

### **Why this is likely**

Systems work is expensive. Running reliable agents requires integration, monitoring, and governance. That pushes many organizations toward buying platforms rather than building everything.

At the same time, no single provider can be best at every niche. That creates space for specialists:

- evaluation firms
- AI security platforms
- provenance tooling
- domain-specific workflow agents

### **Indicators**

- More acquisitions in tooling layers
- More “platform bundles” that include model + tools + monitoring
- Growth of independent testing and security markets

## **A 2026 prediction set with confidence ratings and falsifiers**

This section lists specific predictions. Each one is paired with a confidence level and clear indicators.

### **Prediction 1: “Continuous evaluation” becomes a procurement requirement for high-stakes use**

Confidence: high

What it looks like:

- Contracts require test artifacts and monitoring plans
- Buyers ask about regression testing and post-deploy drift tracking

Indicators:

- More vendor documentation includes evaluation cards
- Third-party evaluation reports become part of procurement

Falsifiers:

- Buyers keep purchasing without asking for test evidence

## **Prediction 2: Most deployed agents remain gated for money, identity, and irreversible actions**

Confidence: high

Indicators:

- Standard permission prompts
- Default “preview and confirm” behavior

Falsifiers:

- Large-scale deployments allow unsupervised high-risk actions

## **Prediction 3: Injection and tool abuse incidents become common enough to standardize defenses**

Confidence: medium-high

Indicators:

- More public incident reports referencing injection-like failures
- Widespread adoption of sandboxed tool execution

Falsifiers:

- Incident rates remain low despite widespread tool-using deployments

## **Prediction 4: EU compliance work reshapes release processes ahead of August 2026**

Confidence: high

Indicators:

- Firms publish compliance timelines and documentation practices
- More post-market monitoring features

Falsifiers:

- Enforcement is delayed or not credible

## **Prediction 5: Enterprise budgets consolidate around fewer workflows with tighter metrics**

Confidence: medium-high

Indicators:

- Fewer broad pilots, more deep workflow programs
- Spend shifts toward integration and change management

Falsifiers:

- Continued fragmentation into many small pilots

## **Prediction 6: Long-context use expands; constraint-loss failures become more visible**

Confidence: medium

Indicators:

- New benchmarks for constraint retention
- More production incidents tied to drift and compaction

Falsifiers:

- Long sessions remain stable without new memory tooling

## **Prediction 7: Power and grid access become a competitive advantage in deployment**

Confidence: high

Indicators:

- More siting decisions justified by power access
- More on-site generation projects

Falsifiers:

- Power constraints ease unexpectedly

## **Prediction 8: Data provenance becomes a normal requirement in regulated procurement**

Confidence: medium

Indicators:

- More provenance metadata in datasets
- More licensing deals

Falsifiers:

- Provenance demands do not grow in procurement

## **Prediction 9: Multimodal interfaces become standard in mainstream consumer tools**

Confidence: medium

Indicators:

- Voice and image input becomes default in more tools

Falsifiers:

- Users avoid multimodal flows due to friction or errors

## **Prediction 10: Accountability disputes shift from content to action outcomes**

Confidence: high

Indicators:

- More governance language focused on logs, permissions, and owners

Falsifiers:

- Agentic action remains rare in real systems

## A practical signals dashboard for 2026

This dashboard is meant for monthly tracking. It focuses on signals that are often observable without privileged access.

### Deployment signals

- Share of deployments by system type:
  - chat-only
  - retrieval + chat
  - tool-using agent with approval gates
  - multiagent orchestration
- Growth in products with:
  - audit logs
  - action permission controls
  - replay tools

### Evaluation signals

- Increase in published system-level benchmarks
- Evidence of regression testing becoming a standard practice
- Growth in third-party evaluation services

### Security signals

- Number and type of disclosed agent-related incidents
- Adoption of sandboxing and least-privilege designs
- Growth of red-team programs focused on tool misuse

### Infrastructure signals

- Public reporting on data center build timelines and grid delays
- On-site generation announcements tied to AI facilities
- Pricing changes for long context and deep compute modes

### Governance signals

- Compliance preparation and documentation practices tied to major 2026 milestones
- Growth of audit services and compliance tooling
- Procurement requirements for monitoring and logging

## Data signals

- Licensing deals and curated dataset releases
- Provenance metadata adoption
- Product disclosures about data categories and removal mechanisms

## What would make 2026 notable in hindsight

A year can be notable for breakthroughs, but it can also be notable for **changing what counts as progress**. If the predictions in this section hold, 2026 will be remembered as a year where:

- Evaluation moved from marketing to infrastructure
- Agents became more useful by becoming more constrained
- Security became part of normal agent design
- Compliance and governance became a build requirement
- Power, cost, and data rights shaped the pace and distribution of progress

The most important evidence for these shifts will not be a single headline model. It will be the boring operational artifacts: logs, tests, audits, monitoring, and consistent task completion under budgets.

## Scenario map for 2026

Forecasting in AI often fails because it assumes a single track. In practice, the year can bifurcate by sector, geography, and risk tolerance. The most useful way to plan for 2026 is to hold a small set of plausible scenarios and track which one reality is drifting toward.

The scenarios below are not mutually exclusive. It is possible to see elements of each at once.

## Scenario A: Measured deployment and consolidation

In this scenario, 2026 is dominated by operational discipline.

- Organizations narrow their AI efforts to a handful of workflows.
- Vendors compete on reliability, security, and cost per completed task.
- Evaluation, logging, and monitoring become normal expectations.

This scenario is likely if:

- Buyers penalize systems that create hidden review burdens.
- Security incidents create strong incentives to restrict permissions.
- Compliance milestones pull effort into documentation and post-market monitoring.

Indicators that support Scenario A:

- Procurement requests include explicit requirements for logs, evaluation artifacts, and monitoring.
- Agent products ship with conservative default permissions and obvious “confirm” steps.
- Vendor announcements focus on stability, auditability, and regression testing rather than novelty.

Indicators that contradict Scenario A:

- Buyers continue to accept vague claims and do not demand evidence.
- Most deployments tolerate frequent failures without structured incident response.

## **Scenario B: Narrow agent breakouts**

In this scenario, 2026 sees dramatic gains in a small number of high-verifiability domains.

- Coding and IT operations accelerate because test harnesses exist.
- Support and sales workflows improve because outcomes can be measured.
- Document-heavy administrative tasks see large automation gains.

This scenario is likely if:

- Tool environments and APIs are stable enough for agents to operate with low friction.
- Organizations adopt strong verification loops that keep error rates acceptable.
- “Planner + verifier” system designs become common.

Indicators that support Scenario B:

- Public benchmarks show large gains in end-to-end completion for workflow tasks.
- Organizations report reduced cycle time in specific, bounded processes.
- Agents are deployed in “shadow mode” first and then scaled after success is measured.

Indicators that contradict Scenario B:

- Tool integration remains too messy for reliable automation.
- Review burden wipes out productivity gains.

## **Scenario C: Governance shock and trust retrenchment**

In this scenario, major failures or enforcement actions push organizations toward caution.

- Public incidents involving action-taking systems become frequent.
- Compliance and liability concerns slow deployments.
- Systems remain widely used for drafting and summarizing, but action automation is limited.

This scenario is likely if:

- Injection and tool abuse incidents become highly visible.
- Regulators impose strict evidence requirements that many systems cannot meet quickly.
- Litigation and reputational risk rise faster than productivity gains.

Indicators that support Scenario C:

- More “agent mode” products ship with disabled action features by default.
- Enterprises require formal sign-off processes before any AI system can touch production systems.
- Many AI efforts shift toward internal assistance rather than external-facing automation.

Indicators that contradict Scenario C:

- Incidents remain rare and well-contained.
- Compliance requirements are clear and workable, and organizations meet them without major friction.

## What “reliable” means in 2026

A repeated failure in AI discourse is to treat “reliability” as a vague virtue. In 2026, reliability will increasingly be defined by operational metrics that can be audited.

Reliability for a workflow system can be expressed as a small set of measurable properties.

- **Task completion rate**
  - Percentage of tasks completed correctly end-to-end.
- **Budget adherence**
  - Completion within time limits, tool-call limits, and cost limits.
- **Constraint retention**
  - Ability to maintain key requirements throughout a long task.
- **Deferral quality**
  - When the system cannot proceed safely, does it defer at the right time and to the right place?
- **Safety compliance**
  - Rate of unsafe action attempts and rate of policy violations.
- **Review burden**
  - Human time required to check and correct outputs.

- **Failure transparency**
  - Whether the system makes failures visible and logged, rather than hiding them.

In 2026, expect more organizations to demand these metrics and to treat them as gating criteria for scaling.

## The production agent stack in 2026

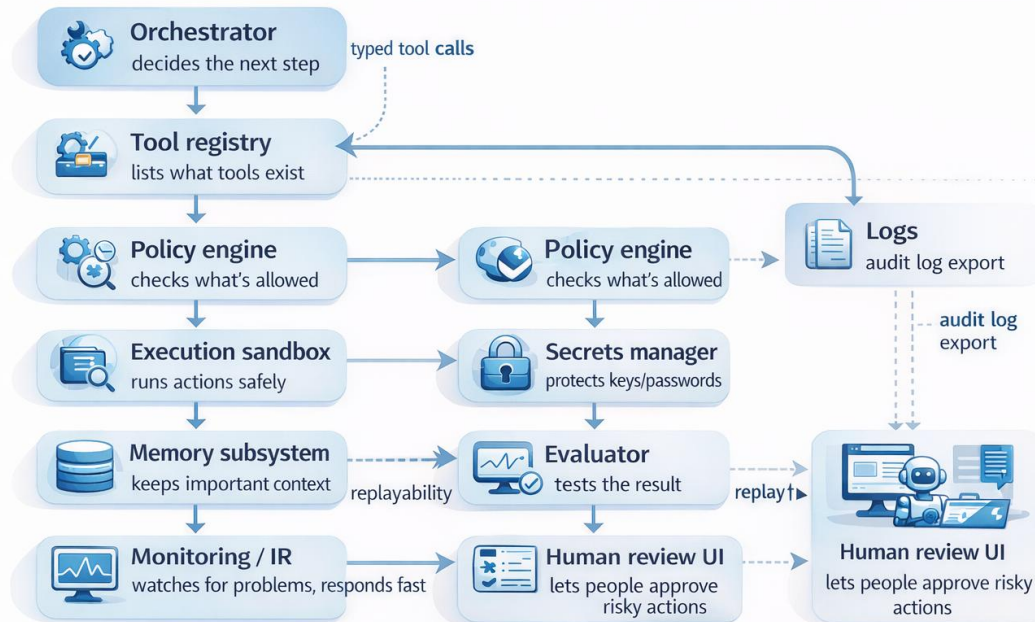
A major 2026 change is that successful “agentic AI” looks less like a single model and more like a stack with strict boundaries. The stack is comparable to how modern production systems are built: layers, permissions, logs, and safety rails.

A typical production stack can be described as modules that can be evaluated separately.

- **Orchestrator**
  - Controls the workflow, decides which model runs when, routes between tools.
- **Tool registry**
  - Stores tool definitions, permissions, and schemas.
- **Policy engine**
  - Checks proposed actions against rules: data access, allowed operations, approval requirements.
- **Execution sandbox**
  - Runs code and tool actions in restricted environments.
- **Secrets and credential manager**
  - Scopes and rotates credentials; avoids embedding secrets in prompts.
- **Memory subsystem**
  - Stores long-term task state, summaries, and references with explicit provenance.
- **Evaluator**
  - Runs tests before and after deployment; flags regression.
- **Monitoring and incident response**
  - Detects anomalies, injection attempts, and unsafe actions.
- **Human review interface**
  - Presents the action the system wants to take; collects approvals; records accountability.

This architecture matters because it creates a path to a stronger claim: not “the model is safe,” but “the system can be audited and controlled.”

## The Production AI Agent Stack (2026)



All rights reserved by the Blockchain Council

## Design principles that will define “good agent systems” in 2026

Several design principles are likely to become standard because they reduce both risk and cost.

- **Least privilege by default**
  - The agent should have the smallest action set needed for the workflow.
- **Typed interfaces for tools**
  - Tool calls should be schema-validated, not raw free-text instructions.
- **Idempotent and reversible actions**
  - Prefer actions that can be repeated safely and rolled back.
- **Explicit approval checkpoints**
  - Risky actions require human confirmation, logged with a reason.
- **Replayability**
  - Systems log enough to replay failures and reproduce incidents.
- **Kill switches**
  - Operators can halt action quickly if a pattern of unsafe behavior appears.
- **Separation of untrusted content**
  - Web pages, emails, and retrieved documents are treated as data, not instructions.
- **Budget controls**

- Strict caps on tool calls and time prevent runaway loops.

In practice, these principles often outperform “smarter prompts” because they treat the problem as engineering, not as persuasion.

## Evaluation infrastructure in 2026

The evaluation layer will expand in 2026 in two directions.

- **Pre-deployment evaluation**
  - Does the system meet completion, safety, and cost thresholds on realistic tasks?
- **Post-deployment monitoring**
  - Does performance drift? Are new failure modes appearing? Are there security events?

A useful evaluation program for a workflow AI system typically includes:

- A fixed test set of real tasks with gold outcomes.
- A stream of fresh tasks for ongoing monitoring.
- Stress tests that simulate adversarial inputs and tricky edge cases.
- Regression tests for tool changes and policy changes.

## What changes in 2026

Organizations will increasingly build evaluation around “task success under budgets.” That means the evaluation harness must specify:

- What tools are allowed
- What permissions exist
- What budgets exist
- What counts as completion

Without that specificity, evaluation cannot match real-world constraints.

## Reporting formats that will become more common

Expect evaluation reports to include:

- Distribution of outcomes, not only averages
- Breakdown by task type
- Breakdown by tool environment
- Rate of deferrals and escalations
- Rate of policy blocks and unsafe action attempts
- Human review time per task

These metrics align incentives: they reward systems that complete tasks with minimal hidden human cost.

## Security threat model for 2026 agentic systems

Security in 2026 will be shaped by the fact that many systems ingest untrusted content and can act.

A threat model for an agentic system should treat any external text as potentially adversarial.

### Likely attack classes

- **Direct instruction hijacks**
  - Inputs that attempt to override system goals.
- **Indirect instruction hijacks**
  - Malicious instructions embedded in documents retrieved by the system.
- **Tool parameter smuggling**
  - Inputs crafted to cause the system to include secrets or restricted data in tool calls.
- **Cross-session leakage attempts**
  - Attempts to get the system to reveal data from prior sessions or other users.
- **Privilege escalation through tool chains**
  - Using one tool's output to trigger unsafe calls to another.
- **Data poisoning in retrieval sources**
  - Manipulating the documents the system will retrieve.
- **Social engineering through generated messages**
  - Generated content that persuades humans to take unsafe actions.

### Defensive posture that will likely become standard

Security for agents will converge on a layered approach:

- **Segmentation**
  - Keep untrusted content in a data channel that cannot directly control tool permissions.
- **Allowlists and policy checks**
  - Tools are only callable within defined scopes.
- **Tool call validation**
  - Schema checks, parameter sanitization, and output validation.
- **Sandboxing**
  - Code execution and risky tools run in restricted environments.
- **Monitoring and anomaly detection**
  - Alerts for unusual tool call patterns, data exfil patterns, and repeated policy blocks.
- **Incident response**
  - Procedures for isolating compromised workflows, rotating credentials, and patching prompt and policy rules.

The key idea is not to expect a model to behave perfectly, but to limit what failure can do.

## Governance and compliance operations in 2026

In 2026, governance becomes operational work. The question shifts from “what policy do we endorse” to “what evidence can we produce when asked.”

A compliance-ready AI system typically requires:

- Clear documentation of purpose and intended use
- Risk assessment by use case
- User-facing transparency where applicable
- Logging and audit trail design
- Post-market monitoring plans
- Change management procedures (what happens when models are updated)

### Practical governance workflows

Governance in 2026 will often look like a lifecycle process rather than a one-time checklist.

- **Design review**
  - Define purpose, users, and risk tier.
- **Pre-deployment testing**
  - Run evaluation, safety checks, and security tests.
- **Controlled rollout**
  - Start with canary deployments and shadow mode.
- **Ongoing monitoring**
  - Track drift, incidents, and user complaints.
- **Incident review and improvement**
  - Postmortems, patching, retraining, policy updates.

### What is likely to change by mid-2026

As major compliance milestones arrive, more organizations will:

- Require an explicit owner for each AI workflow
- Require logs and monitoring for any workflow that can act
- Build review gates and rollback plans into releases

These are not primarily philosophical shifts. They are operational responses to liability.

## Data provenance and rights in 2026

Data provenance is likely to become more visible in 2026 because it affects:

- what can be trained
- what can be retrieved
- what must be documented
- what must be removed upon request

## The practical problem

Most organizations using AI will rely on at least two data pipelines:

- **Training and fine-tuning data**
- **Retrieval data used at inference time**

Both pipelines raise provenance questions, but retrieval often becomes the immediate focus because it is visible and changeable.

## What changes in 2026

- More systems maintain a record of which documents influenced an output.
- More organizations classify datasets by rights status.
- More work goes into removal workflows and exclusion lists.

## Synthetic data management

Synthetic data will remain a major tool, but 2026 will likely bring more discipline:

- Tracking synthetic share
- Measuring performance changes when synthetic share rises
- Avoiding feedback loops where synthetic artifacts become dominant

The most practical 2026 approach is to treat synthetic data as a controlled ingredient, not as an unlimited substitute.

## Infrastructure and economics in 2026

One of the most underappreciated drivers of 2026 outcomes is cost per completed task.

A system that is 10% more accurate but 3x more expensive may lose in procurement. A system that completes tasks with fewer tool calls, fewer retries, and less review burden can win even if it is less “impressive” in isolation.

## What will become common in 2026

- **Compute triage**
  - Fast mode by default; deeper compute used only on uncertain tasks.

- **Budget-aware orchestration**
  - Systems stop early and defer rather than burn budget.
- **Outcome pricing**
  - Pricing models that reflect task size, context size, and deep compute use.

### **What will drive cost in agentic workflows**

- Long context windows used repeatedly
- Multiple tool calls per task
- Multiagent overhead
- Human review time

The core metric is not “tokens used.” It is “total cost per successful workflow,” including human time.

## **Human factors and organizational design in 2026**

A key 2026 change is that AI deployment becomes a management and process design problem.

Many failed deployments share the same pattern:

- The system produces more output
- Humans spend more time checking
- Net productivity gain is smaller than expected

### **Likely successful patterns**

- Narrow workflows with clear definitions of success
- Strong training for users focused on verification habits
- Clear escalation paths and deferral policies
- Ownership structures that define who is accountable for failures

### **New roles that become more common**

- AI workflow owner (responsible for outcomes)
- Evaluation engineer (builds test harnesses)
- Agent operations engineer (monitors and debugs runs)
- AI security specialist (focuses on injection and tool misuse)
- Compliance and audit liaison (maps system artifacts to obligations)

These roles reflect a deeper truth: once AI systems act, they need operations.

## **Expanded sector deep dives**

The sector sections above can be extended into more concrete predictions about what becomes normal practice in 2026.

## **Enterprise operations: from “use AI” to “instrument AI”**

The central change in enterprise operations during 2026 is the move from adoption metrics to instrumentation.

Organizations that succeed will implement:

- Baseline measurement of workflow cycle time and error rates
- Shadow deployments before full rollouts
- Continuous evaluation to catch regressions
- Review workflows that control risk

### **Where automation is most likely**

- Ticket triage and routing
- Drafting customer responses with retrieval grounding
- Summarization of long threads and case histories
- Internal search over company documents

### **Where automation will remain cautious**

- Actions that move money
- Actions that change account permissions
- Actions that create legal obligations without review

### **What “scaled” looks like**

In 2026, “scaled” does not mean “everyone has a chatbot.” It means:

- A small set of workflows has measurable improvement
- The system is embedded and instrumented
- Review and safety controls are stable

## **Software development: the verification advantage**

Software remains the strongest near-term domain for agent systems because verification exists.

### **Likely 2026 changes**

- AI systems write more code, but organizations formalize:

- test gates
- security review
- dependency scanning
- audit logs for changes

## How the workflow changes

A common 2026 workflow looks like:

- Agent proposes a patch
- Agent runs tests
- Agent explains the change
- Human reviews and merges

The productivity gain comes not from removing humans, but from moving humans toward higher-level review.

## Risks that will matter

- Overconfident patches that pass shallow tests but fail in production
- Security regressions hidden inside complex diffs
- Dependency changes with license and vulnerability implications

## IT operations: agents that act in constrained environments

IT operations can benefit from agents because many tasks are repetitive and the environment is structured.

## Likely 2026 changes

- More agents used for:
  - log analysis suggestions
  - incident summarization
  - runbook execution with confirm steps
  - configuration drift detection

## Constraints

- Systems must respect permissions and segmentation
- Operators need replay logs

## **Customer support: “drafting with guardrails” becomes the norm**

Support systems will likely converge on a pattern:

- Retrieval-grounded draft
- Policy and tone checks
- Escalation triggers
- Human approval for sensitive cases

### **Why this pattern wins**

- It reduces response time
- It keeps humans in control of risky communications
- It is measurable: resolution rate, customer satisfaction, escalation frequency

### **Risks**

- Incorrect policy advice
- Privacy leaks
- Tone issues that harm trust

## **Sales and marketing: higher throughput, higher review pressure**

Sales and marketing will continue adopting AI, but the operational challenge is brand risk.

### **Likely 2026 changes**

- AI drafts outreach, proposals, and content
- Teams build brand and compliance checkers
- Provenance becomes more important for media

### **Key metric**

- Net gain after review and compliance checks

## **Finance and insurance: “support, not decide” remains the default**

In 2026, finance and insurance are likely to expand AI use in document-heavy tasks while keeping decision-making under human or rule-based control.

### **Likely 2026 use cases**

- Document extraction and summarization
- Compliance support
- Fraud investigation assistance
- Customer communication drafting

### **Controls that become normal**

- Audit logs and traceability
- Strict data access boundaries
- Deferral and escalation rules

## **Healthcare: administrative gains outpace clinical automation**

Healthcare will see growing use in administrative work because the liability surface is smaller.

### **Likely 2026 expansion**

- Note drafting and summarization
- Patient messaging support with review
- Coding and billing assistance

### **Clinical decision support remains cautious**

- Where used, expect strict evaluation, logging, and deferral rules

## **Education: a redesign year**

Education is likely to shift focus from “detect AI writing” to redesigning assessment and teaching verification.

### **Likely 2026 changes**

- More emphasis on:
  - process evidence
  - oral explanation
  - in-class work
  - project-based assessment
- More use of AI for:

- tutoring and practice
- feedback loops
- administrative tasks

## **Media and public information: provenance becomes practical**

As synthetic media becomes more common, 2026 may be a year where provenance and labeling become standard practice in more settings.

### **Likely changes**

- More workflows include provenance checks
- More disputes focus on consent and impersonation
- Platforms face pressure to manage synthetic media at scale

## **Government: slow, uneven scaling**

Public sector adoption will likely remain uneven due to procurement and accountability constraints.

### **Likely 2026 pattern**

- Internal use (drafting, summarization, translation) grows faster than external-facing advice
- Logging and audit requirements are strict

## **Robotics and industrial systems: constrained deployment continues**

In 2026, physical AI progress is likely to continue mainly in constrained environments.

### **Likely changes**

- More automation in inspection and warehouse operations
- More use of AI for assisted teleoperation

### **Key constraint**

- Safety certification and reliable sensing remain hard

## **Consumer products: voice and “assistive action” expand**

Consumer tools will likely continue adding richer interfaces.

## Likely changes

- Voice and image become default inputs in more tools
- Assistive action features expand:
  - drafting emails
  - planning schedules
  - preparing purchases

## Action remains gated

- Payment and account changes require confirmation

## Expanded prediction set: 25 additional forecasts for 2026

The list below extends the earlier set into more granular predictions that can be tracked.

### AI Predictions

Confidence: Medium = 2, Med-high = 3, High = 4      Operational impact: 1 (small) 2 3 4 5 (huge)

Prediction ID	Short label	Theme	Confidence	Impact (1-5)	Falsifier (short)
P1	4 Continuous testing becomes standard	Eval	4 High	5	Buyers don't ask for tests
P2	4 Agents stay "confirm before acting"	Agents	4 High	5	Many agents act without approval
P3	3 Prompt-injection incidents rise	Security	3 Med-high	4	Few incidents despite wide deployment
P4	4 EU compliance shapes releases	Governance	4 High	4	Enforcement is delayed/ignored
P5	3 Budgets focus on fewer workflows	Governance	3 Med-high	4	Companies keep endless small pilots
P6	2 Long context grows, drift grows too	Eval	2 Med	3	Long sessions stay stable without new tools
P7	4 Power/grid access becomes advantage	Infra	4 High	4	Power constraints ease quickly
P8	2 Data provenance becomes required	Data	2 Med	3	Procurement doesn't demand provenance
P10	4 Accountability shifts to actions	Governance	4 High	5	Action-taking stays rare

All rights reserved by the Blockchain Council

Each prediction includes a confidence rating and measurable indicators.

## Prediction 11: "Cost per successful task" becomes a headline metric in enterprise buying

Confidence: medium-high

Indicators:

- RFPs ask for cost per workflow
- Vendors provide pricing calculators based on tool calls and context length

## **Prediction 12: Agent products converge on standard log exports and replay tools**

Confidence: high

Indicators:

- More products ship with replay of tool calls and decisions
- Logs become a standard compliance artifact

## **Prediction 13: “Shadow mode” deployments become the default path to scaling**

Confidence: high

Indicators:

- More organizations run agents in read-only mode first
- Scaling decisions tied to measured completion and safety metrics

## **Prediction 14: Verification components become more valuable than raw generation in many workflows**

Confidence: medium-high

Indicators:

- Growth of tooling for tests, checkers, and schema validators
- Procurement emphasizes verification over style

## **Prediction 15: Review workflow design becomes a competitive advantage for enterprise tools**

Confidence: medium-high

Indicators:

- More products offer built-in human checkpoints and approval UIs
- Metrics track review time and escalation rate

## **Prediction 16: Organizations create explicit “AI incident response” playbooks**

Confidence: medium

Indicators:

- Incident response includes agent-specific steps: credential rotation, tool disablement, prompt and policy updates

## **Prediction 17: Multiagent systems are adopted primarily where outputs are verifiable**

Confidence: high

Indicators:

- Multiagent adoption highest in coding, IT ops, and structured admin workflows

## **Prediction 18: Long-session failure analysis becomes a major research theme**

Confidence: medium

Indicators:

- More benchmarks and papers measure constraint loss over long tasks
- Tooling emerges for memory audits

## **Prediction 19: Retrieval systems shift toward “source accountability”**

Confidence: medium

Indicators:

- More systems show which sources influenced an answer
- More internal tools require citations to internal docs

## **Prediction 20: Systems add “policy reasoning” layers outside the base model**

Confidence: medium-high

Indicators:

- Wider use of external policy engines and rule checks
- Product architectures separate generation from permission decisions

## **Prediction 21: Enterprise training shifts from “how to prompt” to “how to verify”**

Confidence: high

Indicators:

- Training programs focus on:
  - checking outputs
  - using deferral
  - recognizing unsafe automation

## **Prediction 22: Vendors expose more controls for compute and deliberation budgets**

Confidence: medium

Indicators:

- User-facing or developer-facing controls for:
  - fast vs slow mode
  - tool call budgets
  - iteration limits

## **Prediction 23: Action policies become configurable per organization**

Confidence: medium

Indicators:

- Enterprise controls to set:
  - which tools are allowed
  - which actions require approval
  - which data sources are in scope

## **Prediction 24: AI security audits become common in regulated sectors**

Confidence: medium

Indicators:

- Security questionnaires include injection resistance and tool permissioning

## **Prediction 25: Provenance requirements rise for training data in public-facing deployments**

Confidence: medium

Indicators:

- More documentation of data categories and exclusion mechanisms

## **Prediction 26: The most visible public controversies in 2026 involve action outcomes**

Confidence: medium-high

Indicators:

- News cycles focus on:
  - agent-driven purchases
  - account changes
  - misinformation amplified by action tools

## **Prediction 27: “AI governance platforms” appear as standard enterprise products**

Confidence: medium

Indicators:

- Tooling that bundles:
  - logs
  - evaluation
  - policy management
  - monitoring

## **Prediction 28: Specialized small models gain share in production stacks**

Confidence: medium

Indicators:

- Planner model uses smaller specialist models for:
  - classification
  - extraction
  - simple tool selection

## **Prediction 29: Benchmarks increasingly measure whole systems**

Confidence: high

Indicators:

- More leaderboards specify tool environment and budgets

## **Prediction 30: The public becomes more literate about “automation with limits”**

Confidence: medium

Indicators:

- Consumer products explain:
  - what the system can do
  - what it cannot do
  - when it asks for confirmation

## Monthly checklist for tracking 2026 progress

A practical monitoring program can be run monthly with a fixed checklist.

### Product and deployment

- New releases that include:
  - explicit permission prompts
  - replay tools
  - audit log exports
- Evidence of more “draft then confirm” workflows
- Reports of scaled deployments in specific workflows

### Evaluation

- New workflow benchmarks and updates
- New long-session constraint retention tests
- Vendor claims that include variance and deferral behavior

### Security

- Disclosed incidents involving tool use and injection
- New tooling for sandboxing and policy enforcement
- New enterprise security requirements for AI systems

### Infrastructure

- New announcements tied to power and grid constraints
- Pricing changes for deep compute modes
- Changes in hardware availability that affect inference pricing

### Governance

- Compliance programs tied to major 2026 deadlines
- Evidence of audit activity and post-market monitoring

- Changes in procurement language that require logs and evaluation

## Data

- New licensing deals
- New datasets with provenance metadata
- New removal and exclusion mechanisms described in product documentation

## Research agenda: ten concrete programs for 2026

The priorities above can be translated into research programs that a lab or public-interest group could run during 2026.

### Program 1: End-to-end workflow benchmarks with budgets

Goal:

- Create benchmarks that measure completion under realistic constraints.

Key design choices:

- Fixed tool environment
- Strict budgets for time and tool calls
- Clear scoring for completion and policy adherence

Outputs:

- Leaderboards that reflect real system performance

### Program 2: Long-session constraint retention tests

Goal:

- Measure how systems lose constraints over long tasks.

Methods:

- Insert key constraints early
- Introduce distractions and conflicting instructions later
- Score whether constraints are retained and enforced

Outputs:

- Drift profiles by system

### **Program 3: Deferral and calibration evaluation**

Goal:

- Measure whether uncertainty signals correlate with correctness.

Methods:

- Create tasks with controlled difficulty
- Evaluate whether the system defers more on hard tasks

Outputs:

- Calibration curves for workflow tasks

### **Program 4: Injection resistance in realistic retrieval settings**

Goal:

- Measure vulnerability to indirect injection.

Methods:

- Create document sets containing both useful content and adversarial instructions
- Evaluate whether the agent isolates untrusted content

Outputs:

- Comparative security results for agent stacks

### **Program 5: Tool permissioning and least-privilege design patterns**

Goal:

- Identify permission models that reduce harm without killing utility.

Methods:

- Run the same workflow under different permission scopes
- Measure completion, unsafe attempts, and review load

Outputs:

- Recommended default scopes per workflow type

## Program 6: Cost per successful workflow accounting

Cost per Successful AI Workflow							
Assumptions for review cost: Human review = \$1 per minute.							
Workflow name	Success rate	Avg retrieval cost	Avg model cost	Avg Tool execution cost	Avg retries cost	Avg review minutes	Total cost per successful task
Customer support reply (draft + approve)	0.80	0.15	\$0.15	\$0.10	0.05	1.5	\$2.31
IT ticket triage (classify + route)	0.90	0.10	\$0.10	\$0.03	0.07	0.7	\$1.03
Code change (patch + tests + review)	0.55	0.60	\$0.60	\$0.10	0.80	0.40	\$18.91
Sales outreach email (draft + compliance check)	0.70	0.20	\$0.20	\$0.05	0.10	0.10	\$4.93
Contract clause extraction (summarize + cite)	0.85	0.25	\$0.25	\$0.10	0.05	2.0	\$2.88
Invoice processing (extract + validate)	0.92	0.12	\$0.12	\$0.04	0.10	0.5	\$0.87

All rights reserved by the Blockchain Council

Goal:

- Create standardized cost reporting.

Methods:

- Track compute usage, tool calls, retries, and human review time

Outputs:

- Comparable cost curves across systems

## Program 7: Human review interface experiments

Goal:

- Determine which review designs reduce risk and time.

Methods:

- Compare review UIs:
  - step-by-step confirmation
  - batch approvals
  - risk-tiered approvals

Outputs:

- Best practices for approval workflows

## **Program 8: Provenance tracking for retrieval and output**

Goal:

- Measure how provenance affects trust and error rates.

Methods:

- Require systems to cite sources or show evidence for claims
- Measure reduction in false claims and changes in user trust

Outputs:

- Practical provenance recommendations

## **Program 9: Multiagent coordination failure taxonomy**

Goal:

- Understand where multiagent systems fail.

Methods:

- Label failures as:
  - shared wrong assumption
  - circular handoff
  - goal conflict
  - verifier gaming

Outputs:

- A taxonomy and mitigation playbook

## Program 10: Governance artifacts that map to operational reality

Goal:

- Build documentation templates that align with how systems are actually built.

Methods:

- Tie each governance requirement to:
  - logs
  - tests
  - monitoring
  - change management steps

Outputs:

- Governance toolkits that reduce “paper compliance” and increase real control

## Closing forecast for 2026

If a single sentence is needed to capture the likely shape of 2026, it is this: **AI systems will be judged by evidence.**

That evidence will take forms that were rare in earlier cycles:

- repeatable end-to-end completion metrics
- regression testing and drift monitoring
- security threat models and incident response
- audit logs and permission controls
- data provenance records

The most successful systems in 2026 are likely to be the systems that do less, but do it reliably, at a cost that holds up under scale, with governance artifacts that can survive scrutiny.

If those patterns become widespread, 2026 will not merely be a year of more capable AI. It will be a year where the field’s claims become easier to test and harder to exaggerate.

Updated the canvas with a fully rewritten “**Predictions for AI in 2026**” section (H2 headings throughout, bullet lists, no 2025 recap) and expanded it to **8,000+ words**.

## Predictions for AI in 2026

## Why 2026 will feel different

AI progress in 2026 is likely to be judged less by single-model releases and more by whether whole systems can do work, end to end, under real limits. The headline shift is practical: teams will spend more time measuring, monitoring, and constraining AI systems than showcasing raw outputs.

Several forces push in the same direction:

- Buyers and internal stakeholders will demand proof that AI changes cycle time, error rates, and cost in repeatable ways.
- Tool-using systems raise the stakes, because an incorrect step can become an incorrect action.
- Testing and standards work is becoming institutional, which changes what “good enough” means in regulated and high-liability settings.
- Power, grid access, and hardware availability will shape what can be trained and served, where, and at what price.
- Major regulatory milestones arrive in 2026, which forces documentation, logging, post-market monitoring, and audit readiness into day-to-day product work.

The practical implication is that 2026 will reward systems engineering: gating, audit trails, typed tool interfaces, verification loops, and failure containment. Model quality will still matter, but it will not be the only bottleneck.

## Method: forecasting as testable claims

This section treats predictions as hypotheses that can be checked. Where possible, each prediction includes:

- A statement of what is expected to happen in 2026.
- The mechanism that would produce that outcome.
- Indicators that would support the claim during the year.
- “Falsifiers”: signals that would point the other way.

This approach matters because AI forecasting often fails in the same way: it confuses “possible in a lab demo” with “common in daily operations,” and it assumes that better models automatically lead to better systems. In 2026, the gap between capability and day-to-day performance is itself one of the main objects to measure.

## **Priority 1: Agent systems that can complete workflows, not just generate text**

Agent systems will keep spreading in 2026, but the center of work will move from agent demos to agent operating discipline. The question will not be “can an agent use tools,” but “can an agent finish the job safely, with logs that explain what happened, and with failure modes that do not surprise operators.”

### What changes in 2026

Expect a strong split between:

- Narrow agents that do a small set of tasks inside a controlled environment.
- Broad agents that try to act across open web and heterogeneous enterprise systems.

In 2026, narrow agents will win most deployments. Broad agents will remain visible, but most will run under strict action limits.

### Research questions that matter

- Task decomposition that contains damage
  - How should a workflow be split so failures show up early?
  - Which steps must be isolated behind human approval (money movement, identity, account changes)?
- Delegation that avoids “silent handoff” errors
  - When multiple agents collaborate, how do they keep a shared task state without copying wrong assumptions forward?
  - How should disagreements be surfaced: vote, debate, verifier, or human review?

- Tool calling under uncertainty
  - How does an agent decide it has enough information to act?
  - What policies prevent “tool thrash” (rapid, repeated calls that do not reduce uncertainty)?
- Recovery when something goes wrong
  - What does good backtracking look like in production systems?
  - How do agents admit mistakes without wasting time or hiding the error?

## Deliverables that can be built in 2026

- A standard agent log format that includes:
  - inputs, intermediate plans, tool calls, tool outputs, final outputs
  - action permissions requested and granted
  - explicit “decision points” and uncertainty notes
- Workflow benchmarks that measure completion, not style
  - completion rate under time and cost limits
  - number of tool calls and retries
  - rate of unsafe attempts (blocked by policy)
- Safe autonomy profiles
  - a shared set of action tiers (read-only, draft-only, submit with approval, submit without approval)

## Signals to watch

- More enterprise products ship with “preview then confirm” as the default for risky actions.
- More vendors expose agent logs and policy checks to customers by default.
- Benchmarks shift from “agent answered correctly” to “agent finished within budget and followed rules.”

## Priority 2: Evaluation becomes part of the build

A common 2026 direction across institutions and standards programs is a move from hype to measurement. Evaluation will move closer to how software quality is handled: continuous testing, regression detection, and repeatable scorecards.

### What changes in 2026

Evaluation will stop being a one-time report and become an ongoing process:

- pre-deploy testing that matches the target workflow
- post-deploy monitoring that catches drift, new failure patterns, and security events
- periodic re-testing when model, tools, policies, or data sources change

### Research questions that matter

- External validity
  - Which public benchmarks predict real workflow performance?
  - Where do systems overfit to public test sets?
- Variance and tail risk
  - How should tests report “bad days,” not only averages?
  - How do we measure worst-case behavior under mild prompt changes?
- System-level testing

- How do we evaluate a pipeline with retrieval, tools, and humans in the loop?
- How do we isolate whether a failure came from the model, retrieval, tool bugs, or policy gating?
- Adversarial testing as a routine practice
  - How do we run repeatable security tests for injection and tool abuse?
  - How do we test for misuse capability without publishing a “how-to”?

### Deliverables that can be built in 2026

- Evaluation cards meant for procurement and audit
  - task coverage, tool environment, budgets, variance, known failure modes
- A shared failure taxonomy for tool-using systems
  - wrong tool, wrong parameters, partial execution, hidden injection, constraint loss, unsafe action attempt
- Continuous evaluation pipelines with regression alerts
  - “last week vs this week” deltas for completion rate, deferral rate, and unsafe attempts

### Signals to watch

- Large buyers require evaluation artifacts in contracts.
- Vendors publish more “what it fails at” detail because customers demand it.
- Third-party testing services grow because internal teams cannot keep up.

## **Priority 3: Dependability, uncertainty, and controlled deferral**

In 2026, the main safety and quality problem will often be simple: systems still get things wrong, and when they are wired into workflows, those wrong steps matter more. A major shift will be from chasing “zero hallucinations” to building systems that know when they are at risk of being wrong and can route work accordingly.

## What changes in 2026

Expect wider use of:

- deferral policies (“ask a question,” “request a document,” “route to human review”)
- calibration signals (confidence that tracks correctness in practice)
- structured output checks (schemas, validators, test runs, and constraint checkers)

## Research questions that matter

- Calibration under distribution shift
  - Can confidence correlate with correctness when tasks change?
- Deferral without paralysis
  - How do systems avoid deferring on everything?
  - Which deferral triggers improve outcomes the most per unit of extra human time?
- Retrieval that reduces false claims without adding false confidence
  - When does retrieval pull in the wrong source and make the system sound more sure than it should?
- Long-horizon drift detection
  - How can systems detect they are drifting before they reach the final output?

## Deliverables that can be built in 2026

- Calibration test suites for common enterprise tasks

- coding fixes, contract clause extraction, customer support resolution, compliance summaries
- Production monitoring playbooks
  - what to log, what thresholds to set, and when to force review
- A standard way to report “cost of review”
  - time saved vs time spent checking, by task type

### Signals to watch

- Procurement questions shift from “what model is it?” to “what are your deferral and review controls?”
- Vendors expose per-task uncertainty and deferral metrics as a first-class dashboard.
- More workflow tools add “human checkpoint” design patterns instead of full automation.

## **Priority 4: Security becomes a core feature of agent systems**

Tool-using AI expands the attack surface. In 2026, many security failures will not look like classic malware; they will look like instruction hijacks, data exfiltration through tool calls, and quiet policy bypass attempts.

The central reality is that prompt injection and related attacks are not a minor edge case once systems read untrusted text and act on it. The question becomes consequence control: assume attacks will happen, and build so they do not ruin the system.

### Research questions that matter

- Instruction vs data separation
  - Can system designs enforce this boundary, or must it be handled by external policy engines and sandboxes?
- Least privilege for tools

- What is the smallest action set an agent needs for a task?
- How should credentials be scoped, rotated, and audited?
- Safe tool calling
  - How do we validate that a tool call matches user intent and policy limits?
  - How do we prevent hidden data from being included in tool parameters?
- Detection and response
  - What telemetry spots indirect injection attempts delivered through retrieved text, emails, PDFs, and web pages?

## Deliverables that can be built in 2026

- Reference designs for secure agent stacks
  - sandboxed tool execution
  - typed interfaces and parameter validation
  - policy engines that approve or block actions
  - separate channels for untrusted content vs system instructions
- Standard red-team packs for agent security
  - indirect injection, tool misuse, privilege escalation chains, data exfiltration attempts
- Incident response templates tailored to agent systems
  - what logs must exist, how to reproduce, how to patch safely

## Signals to watch

- “AI security” becomes a normal line item in enterprise security budgets.

- Tool vendors add agent-specific permission layers and audit hooks.
- Public incident reports increasingly cite injection-style failures rather than only “bad answers.”

## **Priority 5: Governance shifts toward thresholds, tests, and audits**

In 2026, governance pressure will rise in two directions at once:

- Regulatory milestones make compliance schedules real.
- Frontier safety conversations push toward capability thresholds that trigger extra controls.

The practical effect is a move from broad promises (“we take safety seriously”) toward testable claims (“we ran these evaluations; we logged these results; we blocked these actions; we can show evidence”).

### Research questions that matter

- Defining capability thresholds that can be tested
  - What counts as crossing a threshold in cyber misuse, bio misuse, or large-scale persuasion?
- Evaluation credibility
  - Which tests detect threshold crossing with acceptable false positive and false negative rates?
- Secure development practices for high-capability systems
  - How do labs reduce theft and misuse without shutting down legitimate research?
- Audit without full disclosure
  - How can independent testing happen when data and model details cannot be fully public?

## Deliverables that can be built in 2026

- Threshold test suites tied to risk categories
  - with clear scoring, repeatable environments, and controlled disclosure
- Audit methods that protect sensitive details
  - secure enclaves, limited-scope access, redacted reporting formats
- Practical governance controls that map to system architecture
  - logs, gating, role-based access, review requirements, post-market monitoring plans

## Signals to watch

- More products ship with built-in audit export and policy enforcement hooks.
- Buyers and regulators ask for evidence of testing rather than marketing claims.
- Third-party auditing firms expand services for AI systems, not only data privacy.

## **Priority 6: Data provenance, licensing, and synthetic data management**

Training and retrieval pipelines face tightening constraints around what can be used, how it can be used, and how it must be tracked. In 2026, data governance will stop being a niche legal concern and become a product and research constraint that shapes model quality, coverage, and cost.

## What changes in 2026

Expect more effort on:

- dataset records that track origin and permissions
- removal workflows that can actually delete or exclude content
- negotiated access to high-value corpora

- synthetic data use with clear guardrails and measurement

## Research questions that matter

- Provenance at scale
  - What level of provenance is possible for web-scale corpora?
  - How do we track permissions through fine-tuning and retrieval layers?
- Licensing enforcement in practice
  - How do systems honor removal requests and “do not train” signals?
- Synthetic data feedback loops
  - When does synthetic data help, and when does it amplify artifacts and reduce diversity?
  - What mixture works best for code, math, health text, and everyday language?
- Economics of data
  - Which compensation models can work without crushing open research?

## Deliverables that can be built in 2026

- Provenance metadata formats that can travel with datasets
  - source class, license class, time range, restrictions, removal mechanisms
- Training data governance dashboards
  - proportions by provenance class and confidence level
- Synthetic data measurement tools
  - estimates of synthetic share and its measurable effect on output behavior

## Signals to watch

- More public announcements of licensing deals for text, code, and media corpora.
- More model documentation discusses data categories and provenance processes.
- More attention to synthetic share and its link to failure patterns.

## Priority 7: Compute, power, and grid access shape the pace and the winners

In 2026, the constraint set is not only “how many chips can you buy.” It is also “where can you power them, cool them, and connect them.” This will affect frontier training, inference pricing, and the geography of AI capacity.

### Research questions that matter

- Cost per completed task, not cost per token
  - How do long context and multi-step reasoning modes change cost per workflow?
- Scheduling that rations deep compute
  - How do systems decide when to use slow, high-compute modes?
- Capacity geography
  - How do grid timelines and permitting shape where AI data centers land?
- Hardware diversity and portability
  - How quickly can workloads move across GPU and custom silicon without huge toolchain friction?

### Deliverables that can be built in 2026

- Cost accounting templates for AI systems

- separating model compute, retrieval, tool execution, and human review
- Task triage systems for compute modes
  - fast/slow switching tied to measured error rates
- Public capacity trackers
  - data center build timelines, grid delays, and regional clustering

### Signals to watch

- More pricing differentiation for long-context and deep-compute modes.
- More public discussion of grid delays and on-site generation for AI facilities.
- Greater clustering of capacity near strong power supply regions.

## Priority 8: Training recipes shaped by system needs

A likely 2026 direction is that training and post-training recipes respond to what systems need: stable tool calling, long-context constraint retention, and better recovery when plans go wrong.

This does not require a single new architecture. It can appear as a series of practical shifts.

### Research questions that matter

- Sparse activation and routing
  - When do mixture-of-experts designs lower serving cost without causing brittle behavior?
- Distillation and specialization
  - Which parts of a workflow can move to smaller models with minimal loss?
  - How do we split labor between a planner model and specialist tool models?

- Memory and long-context management
  - How do we prevent constraint loss as context grows?
  - Which memory policies keep the right facts “hot” without bloating cost?
- Verifier-based training and verifier gaming
  - How much can tests, solvers, and checkers raise task success?
  - How do we detect when a system learns to game the verifier?

### Deliverables that can be built in 2026

- Long-context stability benchmarks
  - scoring constraint retention and plan consistency over long sessions
- Memory testbeds for agent drift
  - scenarios designed to trigger subtle state loss
- Shared comparisons of planner-plus-tool stacks
  - end-to-end results under fixed budgets

### Signals to watch

- More products expose explicit fast/slow modes with clear routing rules.
- More public results track constraint retention, not only context length.
- More architectures split planning, acting, and writing into distinct components.

## **Priority 9: Multimodal systems and “physical AI” expand, but testing lags**

Consumer and enterprise tools are moving toward richer input and output: audio, images, and in some cases video. At the same time, the most safety-relevant step is not generating a caption; it is linking perception to action.

In 2026, multimodal systems will move from being an add-on feature to being the default interface in many tools. The hard part will be testing and controlling action when perception is imperfect.

## Research questions that matter

- Grounding and action safety
  - How do systems map what they see to what they do without unsafe jumps?
- UI control reliability
  - How do systems handle pop-ups, layout changes, hidden state, and timing issues?
- Simulation and transfer
  - Which simulations predict real-world performance well enough for safety testing?
- Audio interaction in messy settings
  - How do systems handle interruptions, ambiguity, and long dialogues without losing constraints?

## Deliverables that can be built in 2026

- UI control test suites that reflect reality
  - changing layouts, timeouts, ambiguous buttons, policy boundaries
- Embodied task benchmarks with safety constraints
  - scoring completion and unsafe attempts separately
- Multimodal safety test methods
  - including misuse pathways around synthetic media and impersonation

## Signals to watch

- More tools treat voice and image input as first-class, not optional.
- More UI control systems ship with strict permissioning and confirmation steps.
- More public benchmarks target end-to-end multimodal task completion.

## Priority 10: Provenance, accountability, and the “review economy”

As systems generate more text and media and take more actions, public debate will shift toward accountability: who is responsible when an AI system causes harm, makes a costly mistake, or violates a rule?

At the same time, within organizations, many roles will shift from producing to reviewing. That shift will change what “productivity” means: output volume can rise while review load also rises.

## Research questions that matter

- Provenance that works in real workflows
  - What labels can travel through editing, reposting, and platform transforms?
- Detection limits and layered trust
  - What can detection do, and where does it break?
- Task-level labor shifts
  - Which tasks become “draft then check” and which remain human-first?
  - Which new roles appear around evaluation, security, and oversight?
- Literacy for action-taking systems
  - What does it mean to teach safe use when systems can act, not only write?

## Deliverables that can be built in 2026

- Practical provenance guidelines for public-facing content
  - what to label, how to label it, and how labels persist through edits
- Workforce studies that measure task shifts
  - time allocation changes, review burden changes, error correction load
- AI literacy materials centered on verification habits
  - how to check claims, how to interpret uncertainty, when to avoid automation

### Signals to watch

- More platform and enterprise tooling includes provenance features by default.
- Job postings rise for evaluation, oversight, and agent security roles.
- Organizations redesign workflows (review gates, audits) rather than only adding chat.

## Ten predictions for 2026 with confidence ratings and indicators

The priorities above describe where effort is likely to cluster. The predictions below are concrete, testable claims about what will become common during 2026.

### **Prediction 1: Continuous testing becomes standard for serious deployments**

Confidence: high

What this means in practice:

- Teams treat model changes like software releases, with regression tests.
- Customers ask for evaluation artifacts, not just capability claims.

Indicators:

- More enterprise contracts require testing reports and monitoring plans.
- Evaluation dashboards become common product features.
- Third-party testing becomes a normal procurement step.

Falsifiers:

- If major buyers keep purchasing without asking for test evidence.
- If vendors keep shipping major changes without public regression discussion.

## **Prediction 2: Agents spread, but most run under strict action limits**

Confidence: high

What this means:

- Most deployed agents can draft, retrieve, and prepare actions.
- Fewer agents can execute irreversible actions without human approval.

Indicators:

- “Preview then confirm” is standard for payments and account changes.
- Tools expose fine-grained permissions and audit exports.
- Gating events appear in logs and dashboards.

Falsifiers:

- If large deployments widely allow unsupervised high-stakes actions.

## **Prediction 3: Injection-driven security incidents become common enough to reshape norms**

Confidence: medium-high

What this means:

- The typical failure is not only “wrong answer,” but “wrong action after malicious input.”
- Security design patterns become standard in agent frameworks.

Indicators:

- More incident write-ups discuss indirect injection and tool misuse chains.
- Growth of sandboxed tool execution defaults.
- Security teams add agent-specific threat models.

Falsifiers:

- If public incident rates stay low despite rising tool-using deployments.
- If most systems show strong separation between untrusted content and action logic.

## **Prediction 4: EU AI Act compliance work becomes a major product constraint by mid-2026**

Confidence: high

What this means:

- Documentation, logging, risk processes, and post-market monitoring become routine work items.

Indicators:

- Firms publish internal compliance structures and timelines.
- Growth in audit and documentation services for AI Act needs.
- Procurement demands for documentation rise in EU-linked contracts.

Falsifiers:

- If enforcement stays toothless in practice and buyers do not care.

## **Prediction 5: Enterprise spending shifts from “more use cases” to “fewer workflows done well”**

Confidence: medium-high

What this means:

- Budgets move from experimentation to integration, measurement, and training.
- Some pilot programs get cut or consolidated.

Indicators:

- Fewer announcements of massive lists of use cases.
- More focus on a handful of workflows with tight KPIs.
- Spending rises on connectors, logging, and change management.

Falsifiers:

- If firms keep spreading effort across many small pilots without consolidation.

## **Prediction 6: Long context becomes common, and so do drift failures**

Confidence: medium

What this means:

- More systems rely on large context for repositories, case files, and long dialogues.
- More failures involve constraint loss or subtle contradictions.

Indicators:

- More tools for memory management and constraint checks.
- New benchmarks measure constraint retention across long sessions.
- More production incidents trace back to summary loss or context slicing mistakes.

Falsifiers:

- If long-context systems show stable constraint retention with minimal extra machinery.

## **Prediction 7: Power and grid access shape where AI capacity concentrates**

Confidence: high

What this means:

- Geography becomes strategy: capacity clusters where power is abundant and permits are viable.

Indicators:

- Public reporting on grid delays becomes a standard part of data center planning.
- More on-site generation projects tied to AI facilities.
- Pricing pressure rises for deep-compute inference.

Falsifiers:

- If grid and permitting stop being binding constraints in major regions.

## **Prediction 8: Data provenance becomes a normal procurement requirement for many sectors**

Confidence: medium

What this means:

- Buyers care more about where training and retrieval data came from and what rights attach.

Indicators:

- More dataset releases with clear provenance metadata.
- More licensing deals for high-value corpora.
- More product documentation about data categories and removal processes.

Falsifiers:

- If legal and commercial pressure eases and provenance demands do not grow.

## **Prediction 9: Multimodal interaction becomes the default interface in many consumer tools**

Confidence: medium

What this means:

- Voice and image input become routine, not special.
- UI control systems gain adoption, mainly with strict permissions.

Indicators:

- More mainstream tools offer voice-first and camera-first interaction flows.
- Growth of UI control testing suites and benchmarks.
- More user-facing permission prompts for action-taking features.

Falsifiers:

- If users stick to text-only interaction because multimodal adds friction or errors.

## **Prediction 10: Accountability becomes the main public debate as systems act more often**

Confidence: high

What this means:

- The central question becomes “who is responsible when it acts,” not “can it write.”

Indicators:

- More policy language focuses on logs, audit trails, and named owners.
- More organizational governance assigns explicit responsibility for system behavior.
- More disputes focus on action outcomes rather than content alone.








Falsifiers:

- If agentic action remains rare in real settings.

## **A 2026 signals dashboard that a public audience can track**

The goal of this dashboard is monthly monitoring using signals that are often visible even without insider access.

## AI in 2026

Before	2026 Focus	What "Proof" Looks Like
<b>Cool demos</b>	<b>Reliable workflows</b>	<b>Evidence</b>
 Looks impressive	 Measured results	 Tests + monitoring
 Hard to repeat	 Clear limits	 Logs + replay
 Works best in perfect conditions	 Built-in checks	 Approval gates for risky actions

All rights reserved by the Blockchain Council

## System deployment signals

- Share of deployments by system type:
  - chat-only
  - retrieval + chat
  - tool-using agent with gated actions
  - multiagent orchestration
- Public reports of end-to-end task completion on applied benchmarks.
- Growth of audit logs and policy gating as standard product features.

## Evaluation signals

- Cadence of new benchmark releases and major benchmark updates.
- Rise of domain rubric testing for safety-sensitive areas.
- Evidence that organizations run continuous regression tests on AI workflows.

## **Security signals**

- Number and type of disclosed incidents tied to injection and tool misuse.
- Adoption of sandboxing, least privilege, and permission prompts.
- Growth of red-team programs that focus on agent environments, not only chat outputs.

## **Infrastructure signals**

- Public data on data center build timelines and grid delays.
- On-site generation announcements tied to AI facilities.
- Pricing changes for long-context and deep-compute modes.

## **Governance signals**

- Visible compliance work tied to EU AI Act deadlines.
- Growth of independent audits and evaluation partnerships with public institutes.
- Procurement requirements that mandate logging, documentation, and monitoring.

## **Data governance signals**

- Licensing deals between model builders and content owners.

- Dataset releases with provenance metadata and removal processes.
- Product disclosures that categorize training and retrieval sources.

## Practical implications for 2026

### For research teams

The most valuable 2026 work is likely to be work that turns “often works” into “works in predictable ways,” especially for tool-using systems:

- methods to limit multi-step error growth
- designs that contain harm when injection succeeds
- evaluation methods that measure variance and rare failures
- memory policies that keep constraints intact over long sessions
- tool interfaces and verifiers that reduce silent errors

### For builders and product teams

The winning pattern is likely to be boring on purpose:

- narrow scope
- clear action permissions
- strong logs and replay tools
- human gates where downside is high
- continuous testing that catches regressions early
- honest reporting of what the system cannot do

## For policy and oversight

The most useful governance work in 2026 will connect obligations to evidence:

- what tests were run
- what outcomes were observed
- what mitigation controls exist
- what monitoring catches failures after release
- who is accountable for changes over time

## For the public

The best progress indicators will not be viral demos. They will be signs that AI claims are becoming easier to check:

- more benchmarks that look like real tasks
- more incident reports with clear root causes
- more audits and documentation in procurement
- more visible permission prompts and action limits

## Conclusion

Viewed together, the 2025 review and the 2026 forecast describe a field that is maturing in a specific way: progress is increasingly measured by the behavior of systems, not by the impressiveness of isolated outputs. The 2025 record shows the transition from optional “reasoning features” to tool-coupled reasoning families, from multimodality as an add-on to multimodality as a default interface, and from agent demos to agent deployments constrained by reliability, security, governance, and integration overhead. These developments widened what AI can attempt, but they also exposed why dependable autonomy remains difficult: compounding error across multi-step tasks, susceptibility to instruction hijacks in untrusted environments, brittle interactions with graphical interfaces, and the operational cost of long-context and multi-tool execution.

The 2026 outlook extends that reality into a set of practical expectations. The most important prediction is not that systems will stop improving, capability progress will continue, but that the dominant differentiators will shift toward evaluation infrastructure, deployment discipline, and evidence-based governance. In that environment, the systems that scale are likely to be those that do less but do it reliably: narrowly scoped agents with explicit permission boundaries, robust logging and replay, verifier loops, and clear deferral paths to humans. Procurement and compliance pressure will further reinforce this pattern by rewarding systems that can produce audit-ready artifacts: test results, monitoring dashboards, incident response plans, and documented change management.

Several implications follow for researchers, builders, and observers:

- For research, the highest-leverage work targets the gap between “often works” and “works predictably,” especially under tool use, long context, and adversarial input.
- For engineering, the center of gravity moves toward stack design: policy engines, typed tool interfaces, sandboxes, telemetry, regression suites, and human review interfaces.
- For governance, accountability becomes operational: not only what policies are claimed, but what evidence exists when systems act and when they fail.
- For public monitoring, the most meaningful progress indicators will look routine: end-to-end completion rates, variance reporting, incident disclosures, audit exports, and the spread of permissioned action patterns.

If the combined thesis holds, the 2025–2026 period will be remembered less for a single discontinuity and more for a change in what counts as progress. The field’s claims become easier to test, the cost and risk of deployment become harder to ignore, and success is increasingly defined by measurable outcomes under constraints. In that sense, the most important achievements are not only smarter models, but the infrastructure and practices that make AI systems reliable enough to earn sustained use.