

2026 Emerging Technology Trends

JPMC Internal Use Only

Introduction

Our Global Technology Strategy, Innovation and Partnerships team helps ensure that JPMorganChase remains connected to innovative and emerging trends in both the immediate and long term. Each year, the team puts together a collection of emerging technology trends, along with market and industry perspectives.

This year's report identifies key predictions and distills the most meaningful technology trends by providing a concise overview of each, offering insights from the broader market and industry. The trends represent pivotal areas of innovation, identified through continuous ecosystem connectivity, internal benchmarks, detailed research and insightful conversations with domain experts across JPMorganChase and externally.

The report highlights a set of high-impact predictions and themes that are rapidly reshaping the technology landscape

Context-driven architectures will be everything: Success in enterprise AI now depends on enabling agents to securely access the most relevant data and tools, empowering differentiated products and services. As automation transforms the software development lifecycle, the focus is shifting from manual coding to architecting context-rich applications. **Inference demand drives continued AI buildout:** The demand for AI inference and infrastructure shows no signs of slowing, with ongoing innovation in data center design, silicon and delivery networks forming the backbone for future software and ecosystem development.

Inference demand drives continued AI buildout: The demand for AI inference and infrastructure shows no signs of slowing, with ongoing innovation in data center design, silicon and delivery networks forming the backbone for future software and ecosystem development.

The end of app switching; Intent is the new interface: The dominant interface is shifting from apps and browsers to a single, AI-native environment. This new paradigm collapses every workflow into one continuously personalized stream, where intelligent interfaces anticipate, execute and transact across modalities.

AI-powered simulation enhances testing: Organizations are increasingly relying on advanced simulations to test, validate and optimize products, processes and scenarios before real-world deployment. These virtual environments enable continuous, scalable experimentation, and rapid iteration.

As these trends mature, organizations should align adoption with disciplined governance as risks and regulatory expectations continue to evolve.

As we look ahead, the 2026 Emerging Technology Trends Report underscores how the convergence of these trends is redefining what's possible for global enterprises.

Table of Contents by Prediction

1 Context-driven architectures will be everything

Physical AI _____	7	Evaluation of AI protocols _____	21
Private market data insights _____	9	Agentic SRE _____	23
Knowledge graphs and semantic layers _____	11	Observing AI _____	25
Data formats for AI _____	13	Data centric security policy _____	27
Context engineering _____	15	Agentic IAM _____	29
Reinforcement learning environments _____	17	Human risk management _____	31
Context engineering for end-to-end software development _____	19		

2 Inference demand drives continued AI buildout

AI infrastructure innovation _____	35
Cloud native AI inference _____	37
Quantum Computing _____	39

3 The end of app switching; Intent is the new interface

Agentic browsers	43
AI native workspaces	45
Generative user experiences	47
Multi-modal social listening	49
Agentic wearables	51

4 AI powered simulation enhances testing

Synthetic users / user simulations	55
Proactive defense attack simulation	57





1

Context-driven architectures will be everything

The success of enterprise AI initiatives is reliant on enabling AI agents to effectively and securely access the most relevant data and tools, empowering them to deliver unique and differentiated products and services to customers and clients. As end-to-end automation transforms the software development lifecycle to support the volume of code generated by AI tools, developers will focus less on manual coding and more on architecting context-rich applications, using AI tools and context engineering techniques.

Physical AI

Integration of artificial intelligence with real-world environments, enabling smart devices and robots to autonomously perceive, reason, and interact through sensors and edge devices.

Physical AI represents the convergence of artificial intelligence and physical hardware systems, empowering intelligent agents to perceive, reason and interact with the real world through sensors, actuators and edge devices. This integration enables robots, automated machines and smart systems to act, learn and adapt within physical environments, effectively bridging the digital and physical realms.

Physical AI models are trained on physical interactions, spatial relationships and the laws of physics themselves. Using advanced simulation environments, these systems can experience millions of scenarios in virtual replicas of the real world, learning how objects behave, how forces interact, and how to manipulate physical space. Through techniques like reinforcement learning, sim-to-real transfer, and synthetic data generation, models develop an intuitive understanding of physics, geometry and causality that enables them to operate effectively in dynamic, unstructured environments. Once trained in simulation, these models are deployed to physical hardware where they continue learning from real-world feedback, constantly refining their understanding of physical reality.

The applications span manufacturing, logistics, and robotics. In warehouses, robotic systems trained on millions of simulated pick-and-place scenarios can handle diverse objects they've never encountered, adapting grip strength and approach angles based on visual and tactile feedback. On manufacturing floors, AI-powered quality inspection systems learn to detect defects across varying lighting conditions and product variations by training on synthetic datasets that mirror real production environments.

Autonomous mobile robots navigate complex facilities by combining pre-trained spatial reasoning models with real-time sensor fusion, learning optimal paths while avoiding dynamic obstacles. Robotic arms in assembly lines learn intricate tasks through demonstration and practice in simulation, then transfer that knowledge to physical production with minimal real-world fine-tuning.

The trend toward AI-driven physical automation is being propelled by several converging technological and economic factors. Simulation breakthroughs have enabled advanced physics engines and digital twin platforms to create photorealistic, physics-accurate virtual training environments at scale, while improved sim-to-real transfer techniques have significantly reduced the reality gap by allowing knowledge gained in simulated environments to translate more effectively to physical applications. AI-powered synthetic data generation tools now produce unlimited labeled training data representing diverse physical scenarios, complemented by hardware advances that bring more powerful edge computing chips and sensors capable of real-time AI processing directly to the point of action. The proliferation of IoT devices and sensors has triggered a data explosion, generating vast amounts of physical world data that enables continuous learning and optimization. Meanwhile, the decreasing costs of robotics components, sensors, and compute power have made deployment economically viable across industries. Finally, persistent labor challenges, particularly workforce shortages in manufacturing and logistics, are driving heightened demand for intelligent automation systems capable of adapting to varied and complex tasks.



Market and industry perspectives

McKinsey predicts the physical AI market is projected to reach \$370B+ by 2040 driven by enterprise adoption across diverse industry applications such as facilities management, physical security, manufacturing and logistics.¹ This is driven by investments in platforms that deliver measurable ROI through energy savings, labor reduction and predictive maintenance.

Physical AI has emerged with a variety of initial use case applications across multiple domains. Smart building and IoT platforms are deploying AI-powered building management and environmental control systems that optimize facility operations. Spatial intelligence platforms are being developed by companies building AI models that understand three-dimensional physical spaces and enable autonomous navigation through complex environments. Embodied AI and robotics applications are introducing physical robots and autonomous agents capable of performing inspection, delivery, and various facility tasks. Simulation and training platforms are creating synthetic environments specifically designed to train Physical AI systems before real-world deployment. Computer vision and perception technologies are providing visual AI capabilities for monitoring, security, inspection, and safety applications across industries. Finally, edge AI infrastructure is delivering the hardware and software necessary to enable AI processing directly at endpoints without relying on cloud connectivity, ensuring faster response times and greater operational independence.

Physical AI presents several key implications for the enterprise across operational domains. In smart buildings, security cameras are leveraging AI to detect real-time threats and send immediate alerts, while access control systems utilize biometrics for enhanced security, and building management systems autonomously control temperature, lighting, and ventilation based on occupancy patterns and environmental conditions. Autonomous operations are transforming warehouses, manufacturing facilities, and logistics operations through the deployment of AI-powered robots that handle material movement, conduct quality inspections, and manage delivery tasks. Enhanced customer experiences are being realized in retail environments where Physical AI enables sophisticated inventory management, checkout automation, and personalized in-store assistance that adapts to individual shopper needs. Additionally, predictive maintenance capabilities are revolutionizing industrial operations as equipment outfitted with sensors and AI can predict failures before they occur, significantly reducing costly downtime and enabling proactive intervention.

¹ Will embodied AI create robotic coworkers?. McKinsey & Company. (2025, June 30). <https://www.mckinsey.com/industries/industrials-and-electronics/our-insights/will-embodied-ai-create-robotic-coworkers>.

Private market data insights exchange for the AI economy

Expansion of real-time insights from structured private market data through seamless, consumption-based exchanges, and enabling the trading of data as an asset class.

Traditional private market data providers are known for providing data attributes tied to company entities (e.g., funding, valuation, investors, job postings), brands (e.g., adverse media), transactions (e.g., location, time), people, and beyond. This private data is typically ingested by corporations through an API feed and further used to enrich existing firmographic data. One challenge with this existing approach is that (1) private markets data is often not structured using a universal entity framework and (2) data is often stuck behind walled gardens, or (3) contractually limited to specific use cases and individual user access, thereby being underutilized.

The rapid growth of both unstructured data as well as agentic capabilities, has given rise to Private Market Data Insights Exchanges. These live data exchanges are built on proprietary entity frameworks which help facilitate the buying and selling of trans-

actional data to a broader population of end-users due to less cost barriers. In addition, some exchanges leverage a waterfall enrichment method, where an individual data attribute (e.g., revenue, funding, e-mail contact) is indexed across multiple source data providers and then scored to maximize data quality and coverage. Further, these exchanges have potential to power LLMs and enterprise applications on a consumption based model.

This consumption based model will involve trading data as an asset class, where these exchanges will operate similarly to trading securities. AI agents will communicate with other AI agents to exchange data insights, and further through model context protocol (MCP) capabilities have potential to generate executive level insight oriented tasks in a much more streamlined manner than humans.

Market and industry perspectives

The estimated total addressable market for Private Markets Data is expected to grow at a 14.5% CAGR from \$8B in 2024 to \$18B by 2030 (BlackRock).² This is driven by overall increased demand for access to private markets, desire to enhance client offerings (e.g., deeper insights into alternative investments) and to gain understanding of performance and drivers of returns.

The emergence of Private Market Data Insights Exchanges will drive the following outcomes:

Improvement in Structured Datasets: The rise of unstructured datasets, makes data rationalization difficult. With the use of proprietary entity frameworks, this can be improved by leveraging unique identifiers to aggregate disparate data sources
Improvement in Verifiable Data: Novel techniques such as waterfall enrichment helps triaging of datasets for data completeness and verifiability

Broader Access to Private Market Data: Emerging data insights exchanges will gain material adoption by natively offering hundreds of primary data sources, which will further be leveraged via API data feeds into homegrown databases in some scenarios this will replace entire usage of traditional market data providers

Actionable Insights & Recommendations: Ability to automate the analysis of raw data, and segment data with the assistance of AI to generate critical insights and messaging to desired audiences, such as internal executives or external sales targets

Accelerated Adoption of AI Research Agents: Leveraging AI to assist in daily tasks (e.g., research and querying, automated e-mail reach out campaigns for GTM teams) and stay ahead of emerging market signals (e.g., headcount changes, social listening, product reviews, adverse media)

Consumption-Based Model / Data As An Asset Class: Traditional market data providers may benefit from a revenue share model with the emergence of private market data insights exchange; however traditional providers may be negatively impacted by the cannibalization of their core user base.

² Will embodied AI create robotic coworkers?. McKinsey & Company. (2025, June 30). <https://www.mckinsey.com/industries/industrials-and-electronics/our-insights/will-embodied-ai-create-robotic-coworkers>

Knowledge graphs & semantic layers

Knowledge graphs are reshaping enterprise data strategies by making information more accessible, contextual and actionable.

Enterprise AI is undergoing a decisive shift from stateless, prompt-driven interactions to context-rich, governed systems. With context engineering emerging as a new practice to provide the right context for a model to optimally complete a task, knowledge graphs are poised to be one of the foundational technologies for delivering meaningful context to AI systems. By providing persistent memory and shared business semantics, AI agents will generate grounded outputs, with a common semantic foundation letting teams reuse data and logic across use cases, improving the return on existing data investments and reducing hallucinations. They also serve as a semantic substrate used at runtime, providing models with a control plane for better reasoning, context injection and policy enforcement.

At the heart of this architecture, knowledge graphs provide memory about entities and events. Ontologies layer on the formal business vocabulary – definitions, constraints, rules and relationships – that tell systems what a customer or business unit means, how the entities relate and which actions are permitted. Semantic layers then operationalize these definitions as governed, reusable views and metrics that both humans and AI agents can consume, compute and interpret consistently with shared entity context. Due to their ability to scale, knowledge graphs and ontologies are invaluable for large organizations seeking to organize and leverage their proprietary data.

This year, several shifts will make this stack standard practice. Retrieval for LLMs will increasingly combine vector search with graph reasoning, often called Graph RAG, to ground answers in enterprise facts with traceable citations. AI-assisted ontology tooling will draft and maintain semantic models from schemas, logs and documents, keeping humans in the loop for quality and compliance. The semantic layer that once served BI will converge with AI needs, unifying metrics, access and policy enforcement for dashboards, applications and agents alike. Event-centric graphs will become more prevalent, allowing agents to reason over sequences like transactions and interactions in real-time. Interoperability will improve as common patterns and standards reduce lock-in and make multi-agent ecosystems viable.

While strides have been made to modernize this decades old technology and improved its viability for AI applications, challenges remain to adopt this widely. The industry will continue to see progress as new players innovate and legacy systems are updated, moving toward solutions that offer richer reasoning and accurate context. Enterprises are already recognizing the strategic value of knowledge graphs and semantic layers, which will become increasingly important for data to be AI ready.

Market and industry perspectives

The knowledge graph market has been around for years. Early adopters in this space either focused on ontology or knowledge graphs. However, we have seen convergence within in the market as the space has grown.

Graph companies have gained adoption with large enterprises due to scalability and fast multi-hop query speed. However, these companies store data as code, freezing the code at a point in time rather than having the data engine figure out connections and conclusions itself.

A few companies recently have focused more on the ontology layer, equipping platforms with ontology models and virtualization engines that are data and application agnostic, meaning they can connect to a number of data sources and tools, including Agent SDKs.

Larger enterprises focus on end-to-end knowledge graph suites that sit on top of an enterprise's data lakehouse, enabling ontologies, applications and agentic frameworks in one platform for AI use cases.

Newer startups in the market are aiming to provide LLMs memory and knowledge. One approach is through a hybrid datastore architecture combining graph, vector, and key-value stores. Another is through building a temporal knowledge graph for AI agents building and updating its graph from not only structured business data, but also from user interactions (chat, unstructured text), tracking when data becomes invalid or changes over time.

Snowflake spearheading the Open Semantic Interchange (OSI) initiative, standards for how entities, metrics and policies are described and exchanged across tools, already indicates how the industry is converging towards governed, semantic first AI architectures.

Data formats for AI

Large-scale, unstructured multimodal data driving the evolution of data infrastructure to better support agentic workloads.

For over a decade, file formats like Parquet (2013, out of Cloudera), Avro (2009, out of Hadoop), and ORC (2013, out of Hortonworks) have served as the backbone of analytical processing. These formats sit on top of object storage (e.g., S3) and are designed to optimize data for analytics.

Sitting above file formats are table formats, which have been a highly competitive battleground in recent years. Open-source Iceberg and Delta have sparked intense industry debate, with Iceberg ultimately becoming the most widely adopted, with ongoing efforts to make both formats interoperable. Table formats add a metadata layer atop file formats, organizing raw files into database-like tables. This lets users store all their data—even structured data typically designed for warehouses—in open, cost-effective formats, while maintaining warehouse reliability and performance without proprietary copies. Users can leverage any compute engine directly on the data, removing the need for data movement or duplication.

While complete ownership, flexibility and interoperability within this stack have provided significant benefits to enterprises, limitations remain. Technologies such as Parquet and Iceberg have been predominantly catered for structured / tabular data, serving batch business intelligence workloads (e.g., SQL).

With the emergence of GenAI and the forthcoming wave of agentic applications—which are proficient at processing highly unstructured and multimodal data, including documents, images and video—new AI-native data formats, both at the file and table level, are being developed to address the evolving requirements of next-generation AI workloads and fill the void of where Parquet and Iceberg fall short.

Open data formats like Lance provide both a table and file format specifically designed for efficient search and retrieval of highly complex multimodal data at massive scale. Open file formats like Nimble, developed by Meta, address Parquet's limitations in AI training by enabling faster reads and more efficient memory layouts. Vortex has also emerged as a Parquet alternative, optimized for AI-native workloads. While these formats are designed for AI workloads, they also support traditional SQL and data engineering (e.g., Spark) processing.

Collectively, these formats seek to overcome the limitations of Parquet and Iceberg in supporting AI and agentic workloads, positioning users to more effectively leverage high-dimensional multimodal data at scale for AI and agentic applications.

Market and industry perspectives

According to Gartner, unstructured data now accounts for 80 to 90% of all new enterprise data and is growing three times faster than structured data.³ This shift is driving major changes across the data ecosystem.

Leading data platforms, along with hyperscalers that have widely adopted Parquet and Iceberg, are responding to this emerging ecosystem threat.

To address the rising complexity of managing and leveraging unstructured data, AI-focused companies are beginning to shift formats. Netflix, where Iceberg originally developed, adopted the Lance data format to power its multimodal data lake, which includes video, audio, images, text, and embeddings. Further, GenAI native companies have adopted Lance internally to power diverse workloads.

Given the industry-wide investment in Iceberg, and recognizing the shifting landscape, the Apache Iceberg Community is currently reviewing the File Format API Proposal, which seeks to establish a unified, pluggable interface for integrating new file formats with Iceberg—much like its current support for Parquet, Avro, and ORC.

³ Gartner. (2025, September 17). Market Guide for Data Security Posture Management. Joerg Fritsch, Brian Lowans, and Andrew Bales. <https://www.gartner.com/en/documents/6964866>

Context engineering

Context engineering orchestrates external information and tools around LLMs to deliver consistent, accurate and domain-specific AI results.

In the first wave of generative AI (GenAI) adoption, interactions were primarily “prompt-response” exchanges, where users would submit a question and receive an answer. During this period, prompt engineering emerged as a practice, focused on crafting and refining prompts to elicit accurate and relevant outputs from models. GenAI use cases have since evolved toward multi-step agentic workflows, where agents autonomously gather information, use tools, and reason over results with minimal human intervention. To manage the increasing volume of information that agents generate, while balancing the token constraints of context windows (e.g., amount of information that can be provided to the model), context engineering has emerged to curate the optimal set of tokens for achieving desired outcomes across multiple steps.

In agentic applications, effective performance depends on managing multiple types of context: instructions, knowledge, and feedback from tools. Instructions include system prompts that guide model behavior and procedural memory for storing specific skills or rules related to a task. Knowledge consists of facts and relevant experiences, often accessed through Retrieval Augmented Generation (RAG) workflows using vector databases or knowledge graphs. Increasingly, the industry is adopting “just-in-time” context strategies, allowing agents to retrieve only necessary information at runtime instead of pre-processing all data

in advance. Throughout these workflows, agents must also integrate feedback and new information gathered from tool interactions, using this context to inform next steps.

While the promise of larger context windows – where most information could be uploaded for a model to use — sounds promising, challenges remain. Processing tokens at such scale drives up computational costs and increases latency, and models, much like humans, are limited by a finite attention budget. As the number of tokens in the context window grows, a model’s ability to accurately recall information from that context declines – a concept known as context rot. Today, industry best practices emphasize that effective context engineering is about selecting the smallest set of high-signal information that maximize the likelihood of achieving a desired outcome.

To achieve this, new techniques like compression are emerging. When the number of tokens approaches the context window limit, this method is used to summarize the most relevant information or filter out less important details. Another strategy involves giving agents a “scratchpad” – a dedicated space for note-taking that is stored outside the context window and can be retrieved as needed. Lastly, sub-agent architectures offer a way to manage context more effectively; instead of a single agent maintaining state across an entire workflow, specialized agents can focus on specific tasks. As models continue to advance, the challenge of engineering the right context to achieve desired outcomes over long time horizons will remain central to building more performant agents.

Market and industry perspectives

While the developer and coding space has been the natural early adopter of context engineering given the large investments in developer agents, initial efforts were often tool-specific and fragmented. Community-driven standards have emerged from developer workflows in early 2024, functioning as simple, static “READMEs for agents” embedded directly in repositories, providing baseline instructions for agents navigating a codebase. Similarly, specialized tools have emerged to automate the analysis of complex codebases, allowing AI agents to generate structured documentation and “on-demand encyclopedias” to understand millions of lines of code without manual human onboarding.

However, as the market matured in 2025, the industry is moving beyond developer tools toward holistic, cross-platform standards that aim to solve the context problem for every business domain.

The first and most established of these is the Model Context Protocol (MCP). Launched by Anthropic in November 2024, MCP is now a widely adopted open standard that provides a universal “USB-C port” for AI. It standardizes how models or agents connect to external tools – whether third-party applications or proprietary internal APIs – ensuring secure and consistent data access across different platforms.

Building on this connectivity is the invention of “Skills”, launched by Anthropic in October 2025. While MCP handles the connection, Skills provide the procedural knowledge. At their core, Skills are modular folders containing instructions, scripts and resources for specific tasks. Instead of overwhelming the context window with every possible instruction upfront, the agent dynamically “discovers” and loads a Skill only when it becomes relevant to the task at hand.

These two technologies are deeply complementary: MCP facilitates the secure “plumbing” to a tool, while Skills provide the domain expertise to transform that raw access into reliable outcomes. Following Anthropic’s December 2025 release of the Skills open standard, the paradigm has seen cross-industry adoption, most notably by OpenAI within ChatGPT and its Codex developer products.

Reinforcement learning environments

Enabling agents to tackle complex, real-world tasks through goal-driven training, interactions with tools, realistic simulations and outcome-based feedback.

In last year's trends report, we highlighted the rise of reasoning models, which materially improved accuracy on higher-value, complex tasks and helped usher in the era of AI agents. This shift from simple prompt-response (next-token prediction at maximum speed) to deliberate reasoning was enabled by post-training reinforcement learning (RL), which teaches models to plan, use tools and evaluate intermediate steps against a goal.

In practice, RL provides the model with an environment and action space, tools it can use within that environment, and a reward signal aligned to desired outcomes. For example, in a coding environment, the model might have access to a code interpreter, the ability to write and execute code and rewards based on whether the program runs and produces correct results. This paradigm shift spurred major model providers to invest heavily in post-training RL and ultimately paved the way for early agentic applications that we are familiar with today, which includes deep-research, developer/coding agents and computer-use agents.

While frontier labs pioneered this work, an emerging ecosystem has formed around providing enterprises with reinforcement learning capabilities to train custom agents for real-world tasks. The key enabler is high-fidelity RL environments, which

provide simulated workspaces with realistic observation and action spaces, integrated tool access (e.g., code interpreters, web search) and scalable infrastructure for iterative learning. These environments can be tailored to mirror almost any knowledge-work task, and early simulations span workflows from Excel spreadsheets to Salesforce dashboards. The ecosystem is converging on reinforcement learning as a service (RLaaS): managed platforms that abstract the infrastructure complexity for developing, training and deploying agents.

While there is industry enthusiasm around the potential for RL environments to solve the "last-mile" challenge of agent accuracy for domain-specific use cases, a hurdle to widespread adoption revolves around creating reliable evaluations or "rewards" (e.g., feedback signal that tells the agent how good or bad its last action was). Consequently, early RL successes are concentrated in domains where rewards are easily verifiable, such as code and mathematics (e.g., does the code run). For more nuanced, subjective tasks (e.g., generating investment memos or legal briefs), current methodologies pair SME-defined natural language rubrics with "LLMs-as-a-judge" evaluators that score agent actions based on the human provided rubric. An emerging industry view is that differentiated value is migrating from base models themselves to reward design, precise evaluation, and high-fidelity RL environments. As that shift takes hold, RL is moving beyond major labs, enabling enterprises to train purpose-built agents for domain-specific workflows.



Market and industry perspectives

All major model providers (OpenAI, Anthropic, Google, xAI) are investing heavily in RL to both improve reasoning in their general model capabilities and target high-value vertical domains. This investment can be seen through Anthropic's plan to spend more than \$1 billion on RL environments over the next year to train models in complex professional workflows.⁴

Additionally, we are seeing the launch of managed reinforcement tuning services that allow developers to leverage RL for vertical-specific tasks, moving beyond generic model usage.

Traditional data labeling companies, which have historically supplied major AI companies with custom datasets for model training, have expanded their product suites to offer RL environments for both AI model providers and enterprises. Their differentiation strategy involves leveraging specialized human expertise to construct these environments and design appropriate reward models tailored for verticalized, industry-specific tasks.

Along with the established players, a specialized ecosystem of well capitalized startups has emerged to provide Reinforcement Learning as a Service (RLaaS) for enterprise-specific workflows, independent of any single model.

Consolidation is already occurring in this space. Leading GPU cloud providers are acquiring specialized RLaaS startups to expand their offerings beyond large model training and inference. These acquisitions bring targeted tools for developers to build and deploy agentic workflows directly on cloud platforms, helping providers differentiate their services and broaden their customer base beyond major AI labs. Similarly, AI inference platforms are acquiring companies focused on post-training and customization, allowing them to move beyond efficient, low-latency model serving and strategically position themselves for deeper model optimization and broader user capabilities.

⁴ Zeff, M. (2025, September 21). Silicon Valley bets big on “environments” to train AI agents. TechCrunch. <https://techcrunch.com/2025/09/21/silicon-valley-bets-big-on-environments-to-train-ai-agents/?secureweb=ONENOTE#:~:text=in%20RL%20environments%20to%20keep,environments%20over%20the%20next%20year.>

Context engineering for the end-to-end software development lifecycle

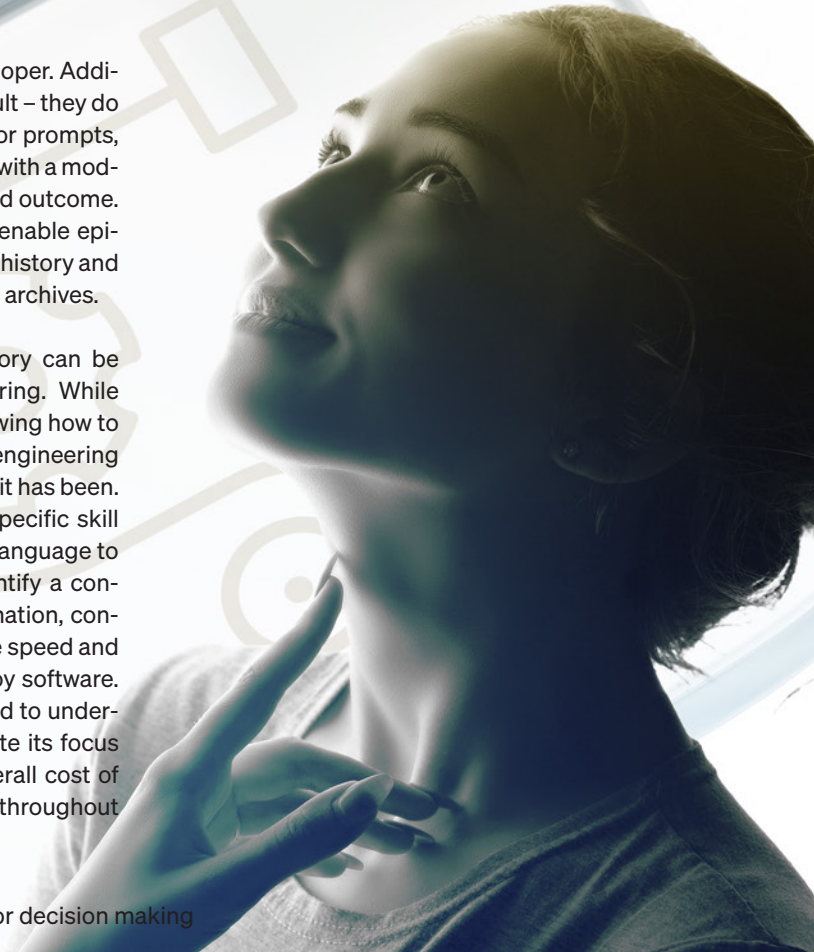
Context Engineering will shift the way software developers interact with AI coding tools.

The common context engineering techniques of skills, knowledge graphs, episodic memory and prompt engineering are being applied throughout the entire development lifecycle, ultimately impacting the quality of code developed. Developers can build skills on how an agent should perform a specific task, especially repetitive tasks that would be standard across any deployment (e.g., using a specific file or language, applying internal company frameworks for code quality, leveraging specific API clients, or auto deploying procedures). Skills can dramatically cut down the time a developer spends doing these specific tasks and provide consistency across an organization.

While knowledge graphs have always played a role in understanding the connections between sources of data, the sophistication of knowledge graphs has evolved dramatically over the past year. AI developer tools are building knowledge graphs that generate detailed wiki pages with data flow and architecture diagrams of code repositories to help developers understand what effects code changes will have. These wikis are auto-updated based on human edits but can also be manually edited to be more accurate

based on additional context given by the developer. Additionally, most AI dev tools are stateless by default – they do not retain information from previous sessions or prompts, causing friction for developers as they interact with a model since it loses previous context on an intended outcome. As a result, we are starting to see more tools enable episodic memory that remembers coding session history and stores completed conversations as searchable archives.

Skills, knowledge graphs and episodic memory can be triggered through effective prompt engineering. While prompt engineering is not a new practice, knowing how to use it in combination with the other context engineering techniques makes it much more powerful than it has been. Developers need to know when to call on a specific skill or even protocol (e.g., MCP) and use the right language to activate the model to recall a memory or identify a connection in a knowledge graph. Used in combination, context engineering will fundamentally change the speed and efficiency at which developers build and deploy software. However, individuals and enterprises alike need to understand that leveraging these techniques, despite its focus on increasing efficiency, may increase the overall cost of development as agents are increasingly used throughout the process.



Market and industry perspectives

Leading players have launched AI development tools for deeper understanding of codebases.

Some code review solutions are embedding more code understanding and search capabilities into their solution to provide developers with more context into their code for more efficient code review. Other end-to-end software development platforms are building software delivery knowledge graphs that will provide additional context to AI agents throughout the entire SDLC. Elsewhere within the market, there is focus on offering a strong context engine to support better code generation.

Evolution of AI protocols

AI protocols enable standardized, cross-platform communication and collaboration between autonomous agents and tools, streamlining integration and scalability.

As AI shifts from isolated models to interconnected autonomous agents and multi-agent systems, AI protocols have emerged to provide standard ways for agents to communicate, exchange context, access tools/data, and collaborate across platforms. Announced in late 2024, Model Context Protocol (MCP) has emerged as the leading standard – but new alternatives and complementary protocols are proliferating, creating a rich and evolving ecosystem. Their proliferation reflects the increasing need for interoperability, portability, and governance as agentic AI moves into production, and suggests that AI protocols will continue to evolve over time.

In practice, this means developers no longer need to build custom, one-off integrations for every model, tool or agent. Instead, AI systems can plug into shared protocols that enable cross-platform tool access, shared memory and context, standardized communication and modular workflows – making agentic AI more portable, scalable, and maintainable.

In the past year, the industry has seen numerous protocols released from leading AI providers with slightly different approaches. While MCP is focused on connecting models to tools, protocols like Agent2Agent (A2A) and Agent Communication Protocol (ACP) define how agents communicate, share context and coordinate tasks. Agent Payments Protocol (AP2) provides a common language for secure transactions between agents and merchants; and there are even more protocols intended to make multi-agent deployment, observability,

and governance feasible across enterprises. Given the explosion of AI protocol adoption, notably MCP, vendor solutions are rapidly building out MCP servers for users to connect to and many vendors today have registries and toolkits to simplify the installation and setup of MCP tools.

While these protocols provide a level of standardization, more tools connected to an agent or AI system may consequently cause inefficiencies. For instance, MCP passes all information through AI workflows no matter how simple or complex a task is, including information like tool definitions, descriptions, etc. This increases the token consumption and costs that may ultimately slow down an agent's performance. There are emerging techniques that suggest MCP servers as code APIs may be more efficient, but the industry continues to rapidly evolve and expect this space to mature over time.



Market and industry perspectives

In December 2025, The Linux Foundation launched a new initiative, the Agentic AI Foundation (AAIF), co-founded by Anthropic, OpenAI and Block – with broad backing from leading tech players including Google, Microsoft, Amazon Web Services (AWS), Cloudflare and Bloomberg. The founding projects donated to AAIF include MCP, OpenAI's project-instruction format AGENTS.md, and Block's agent framework goose – signaling a collective move toward unified, open, vendor-neutral agent standards.

According to Gartner's 2025 Software Engineering Survey, by 2026, 75% of API gateway vendors and 50% of iPaaS vendors, will have MCP features.⁵

In June 2025, Google Cloud donated its Agent2Agent protocol to the Linux Foundation and formed the Agent2Agent project with AWS, Cisco, Microsoft, Salesforce, SAP, and ServiceNow to collaborate and foster an open and interoperable ecosystem for AI agents with the A2A protocol.

⁵ Zarecki, I. (2025, November 23). MCP Gartner Insights for 2025. Delivering Data Products in a Data Fabric & Data Mesh. <https://www.k2view.com/blog/mcp-gartner/#:~:text=By%202026%2C%2075%25%20of%20API,stab%20and%20address%20new%20requirements>.

Agentic site reliability engineering

Autonomously monitor, diagnose, and help remediate system issues, enhancing traditional practices with pattern recognition, root cause analysis and integrated observability.

Agentic Site Reliability Engineering (SRE) represents a transformative shift in the way enterprises approach system reliability and observability. Leveraging the power of large language models (LLMs) and AI-driven agents, Agentic SRE tools are designed to autonomously monitor, diagnose, and remediate system issues, significantly enhancing the efficiency and effectiveness of traditional SRE practices.

At the core of Agentic SRE is the concept of autonomous agents. Capable of operating 24/7 to plan and execute actions on behalf of users, these agents utilize advanced capabilities such as pattern recognition, anomaly detection, and root cause analysis (RCA) to provide detailed insights into system performance. For instance, they can autonomously search through logs and databases, similar to the investigative process a human SRE would undertake, to identify and diagnose issues.

While the current scope of auto-remediation capabilities is not fully autonomous – requiring human initiation for executing recommended fixes – Agentic SRE tools offer significant advancements in RCA. They provide comprehensive knowledge graphs, confidence ratings, and documentation to support their hypotheses, thereby streamlining the troubleshooting process. Additionally,

these tools integrate seamlessly with various observability solutions, code repositories, and IT service management platforms, allowing for a holistic view of the system's health and performance.

Despite their longer-term promise, agentic SRE tools are still evolving. For instance, they currently lack out-of-the-box capabilities to measure essential SRE metrics such as burn rates and error budgets. While the goal is to excel in reducing Mean Time to Resolve / Repair (MTTR), a critical metric for assessing system reliability, they are only part of the way there, but the industry expects their capabilities to mature over time.

Market and industry perspectives

The market for agentic SRE solutions is rapidly expanding as enterprises seek to enhance their system reliability and reduce downtime. The global market for AI-driven observability and reliability tools is projected to grow significantly, driven by the increasing complexity of IT environments and the need for more efficient incident management.

Key players in the agentic SRE space include vendors like Deductive, Traversal, Resolve.ai, and Dynatrace. Each of these companies is investing in unique capabilities to differentiate their offerings. For example, Deductive focuses on code-aware observability, leveraging source code analysis for RCA, while Traversal emphasizes in-depth RCA with confidence levels. Resolve.ai is exploring the integration of tribal knowledge into their agents, mimicking the expertise of human SREs. Dynatrace, with its Davis AI Copilot, is currently being implemented at JPMC, although its capabilities are limited to data within the Dynatrace platform.

Investments in agentic SRE technologies are focused on enhancing auto-remediation capabilities and expanding integration with various data sources. Companies are also exploring ways to incorporate historical incident data and tribal knowledge to improve RCA accuracy and effectiveness. The ability to integrate with any observability vendor and build comprehensive knowledge graphs is a key differentiator for these tools, enabling them to provide a more complete understanding of system health.

As the industry continues to evolve, agentic SRE tools are poised to play a crucial role in modernizing system reliability practices, offering enterprises the opportunity to reduce downtime, improve system performance and ultimately enhance customer satisfaction.

Observing AI

Traditional observability (metrics, logs, traces) provides a foundation, but new metrics and frameworks are needed to measure AI-specific behaviors, model drift, and compliance.

As the industry continues to adopt AI and leverage more agents for workloads, observing LLMs and agentic workflows remains paramount. It is critical to understand what these agents are doing, how they are performing, and the impact they may have to surrounding systems and applications.

Traditional observability pillars serve as a starting point for understanding AI behavior. Metrics, events, logs, and traces have become the defacto way of measuring health and behavior for applications and infrastructure through OpenTelemetry (OTel) standards. Distributed tracing (traces), which captures the flow of a request as it moves through parts of a system, is not as widely adopted as logs or metrics, but may ultimately be the best method to observe AI workflows. Traces can track what an agent did, what inputs/outputs were processed, tools or pre-set instructions it called on, how decisions flowed from one to another, and provide verification that the agent accomplished each step it planned to do.

However, new metrics and frameworks are needed for deeper insights. AI workloads introduce a new level of measurement we have not had before. LLM performance can be measured by number of tokens used for a specific task, response quality (accuracy and completeness), model drift and data quality checks, and the safety and compliance of the AI interaction. These metrics are starting to be offered by an emerging wave of startups focused on AI observability, but also industry incumbents as well. Ultimately, the most effective measurement frameworks will likely integrate real-time signals and alerts from emerging players with long-term performance tracking of agents provided by traditional vendors.

With metrics on how AI models are performing in production, companies are also looking at ways to gather metrics pre-production through model evaluation. Model evaluation involves testing models and understanding how its performance may impact workloads in production, improving the speed at which models are deployed and improving the risk of models failing. Observability metrics for both pre-production and post-production are necessary to understand the end-to-end spectrum of model behavior.



Market and industry perspectives

Incumbent observability vendors are releasing LLM Observability into their platforms to develop, evaluate and monitor LLM apps.

The OpenInference specification emerged this year, which is a set of conventions and plugins complimentary to OpenTelemetry that enables tracing of AI applications. It is designed to provide insight into LLM calls and the surrounding app context (e.g., vector DBs).

AI Observability vendors have seen significant fundraising rounds in the past few years; and a number of recent acquisitions have highlighted the need to embed LLM observability into existing products.

Data-centric security policy

Securing and protecting data wherever it moves will be foundational to AI and agentic systems.

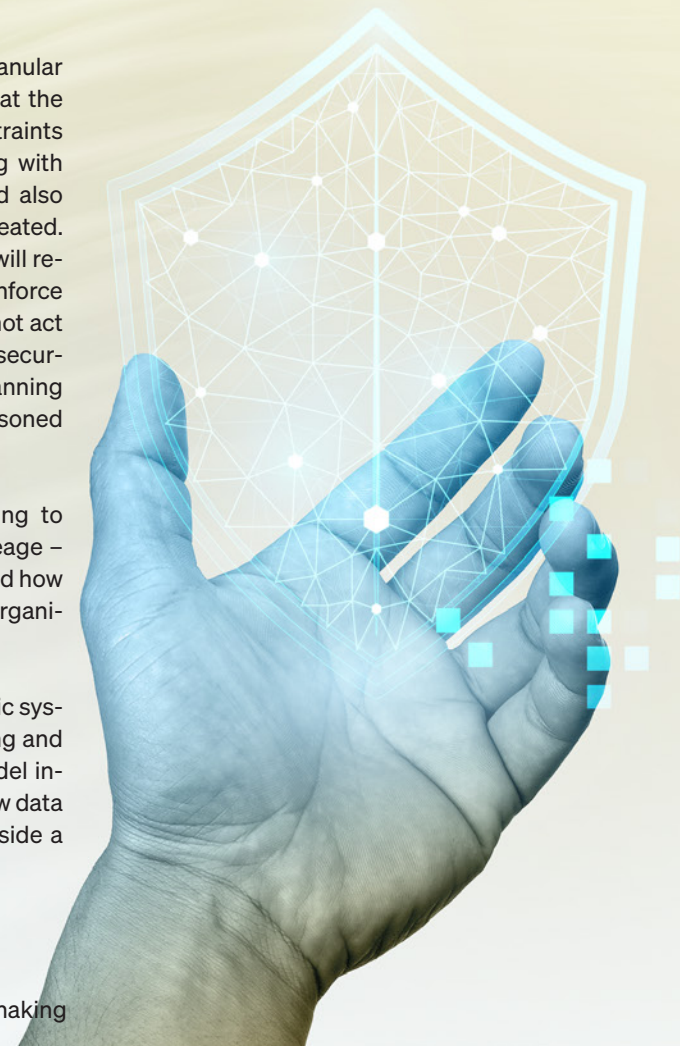
AI-driven data flows pose new challenges for governance and security. As organizations inject data context into generative AI and multi-agent systems, data is no longer static – it flows freely and continually transforms through model prompt and response, ephemeral memory, vector embeddings, knowledge graphs, and agentic operations. AI systems turn enterprise data into derived representations of source data. Data synthesized by an agent from multiple sources can also generate an output of higher sensitivity than the source data. Agent memory and knowledge graphs can bring valuable context but can also store or infer data – for example sensitive implicit relationships – which can be unintentionally exposed.

Protecting data now means securing its movement, not just its storage location. Given this shift, approaches to protecting data are evolving from securing data where it exists to securing how it moves. Data-centric security is an architectural shift that treats data, access to data, and the flow of data as the primary control plane. Data centric security models embed security policy, governance and telemetry directly into data objects so that classification, access permissions, lineage and data protection “travel” with the data wherever it goes.

Data controls must adapt dynamically to support granular context and intent of interactions. Attributes defined at the data object level with sensitivity labels and usage constraints could enable any agent, model or system interacting with data to comply with those rules at retrieval time and also persist those attributes downstream as new data is created. Sources of context like embeddings and graph edges will require the equivalent of table and row level security to enforce dynamic policies – for example, “agents of type X cannot act on data of type Y”. There is also an emerging focus on securing derived data (e.g., embeddings, memory) by scanning these stores for sensitive or potentially maliciously poisoned data.

Data Loss Prevention (DLP) models are also evolving to support data risk detections based on end-to-end lineage – tracking data from where it originates, where it flows and how it transforms – in addition to where it may egress the organization.

Finally, advanced techniques can protect data in agentic systems. There is increasing use of confidential computing and runtime isolation to create trusted boundaries for model inference and agent execution to ensure that sensitive raw data and context never exists unencrypted in memory outside a secure enclave.



Market and industry perspectives

Data Security Posture Management (DSPM) tools are expanding to secure data in transit. Data security platforms, commonly referred to as DSPM tools are extending from scanning and classifying traditional data stores at rest to scanning data at source, AI model inputs / outputs, vector embeddings, agent memory and data in motion to identify existing data risks and prevent unwanted sensitive data from entering into AI systems from the start.

These DSPM players and new emerging startups are also building capabilities to observe data movement and lineage throughout its lifecycle to better identify and triage potential risks (e.g., privacy impacts, data leakage).

Cloud runtime isolation solutions and confidential computing offerings are providing dedicated session isolation to ensure agent state, tool operations and credential access remain completely compartmentalized. When a session ends, the entire environment can be terminated and memory sanitized, minimizing the risk of data persistence or cross-tenant contamination.

Companies have also introduced new models to protect the consumption of proprietary data. Dynamic tokenization capabilities are also emerging to support direct model interactions while safeguarding sensitive data from model providers.

Agentic identity access management

Redesigning traditional human-centric IAM to allow dynamic, intent-aware and auditable authorization for autonomous AI agents.

With the rise of autonomous AI agents, traditional, human-centric identity and access management (IAM) is evolving. As agents increasingly perform actions on behalf of human users, concepts like Know Your Agent (KYA) are quickly gaining traction in customer-facing use cases, ensuring autonomous actions are legitimate, secure, and acting within authorized bounds for financial and sensitive transactions, and preventing fraud by confirming the AI's origin, permissions, and owner.

For the workforce, assigning an identity or authenticating an agent is just the beginning. When an agent is performing actions on behalf of a human employee, identity tokens must explicitly bind together three elements: the agent's identity, the identity of the original requester (human, software or agent) and the intent or context of the request (e.g., which resource is being accessed and why).

To prevent agents from overstepping their bounds, access can be downscoped, i.e. the agent only receives the subset of the human user's permissions that allow it the least privilege that is strictly required for the specific task, granted just-in-time and revoked immediately after use. This structure minimizes overtly permissive entitlements by design so that human approval can be reserved for only the most critical or sensitive agent actions. In multi-hop workflows, where agents invoke other agents or services to perform downstream actions, authorization must persist across hops to the very last step maintaining a clear trace of identity, origin and scope of the original request.

With the rise in agents in SaaS applications, IAM must evolve beyond traditional perimeter-based controls to support secure SaaS-to-SaaS interactions where trust, scope and intent can persist without a shared IAM layer. Enterprises must redesign architectures and ownership models to enforce business policy across internal and third-party agentic systems. In parallel, industry standards and industry-wide adoption of agent IAM best practices will be critical to driving and shaping safe, transparent and reliable agentic interactions.



Market and industry perspectives

Recent research by Salesforce predicts that AI Agent adoption is expected to jump 327% over the next two years.⁶ With the rise in AI agents, cloud providers are formalizing the recognition that agent identities will become a distinct control plane for AI agents. For example, AWS and Microsoft have products to govern credentialing, policy and audit lifecycle for AI agents on their platforms.

Large security platforms have been actively consolidating identity capabilities.

Networking players are expanding into SaaS-to-SaaS and API-level access governance to create gateway-style cross-application controls that can help federate access for agents across platforms. This activity signals a shared recognition that identity will become the primary enforcement layer for agentic AI.

A number of emerging startups are innovating in the space to build control planes to manage, govern and remediate access for agents. While early, the space is evolving quickly to deliver solutions for scaled enterprise adoption.

⁶ HR Leaders to Redeploy a Quarter of Their Workforce as Agentic AI Adoption Expected to Grow 327% by 2027. (2025, May 5). Salesforce. <https://www.salesforce.com/news/stories/agentic-ai-impact-on-workforce-research/>

Human risk management

Embedding human behavior as a security signal to limit exposure resulting from both human errors and susceptibility to adaptive, personalized threats.

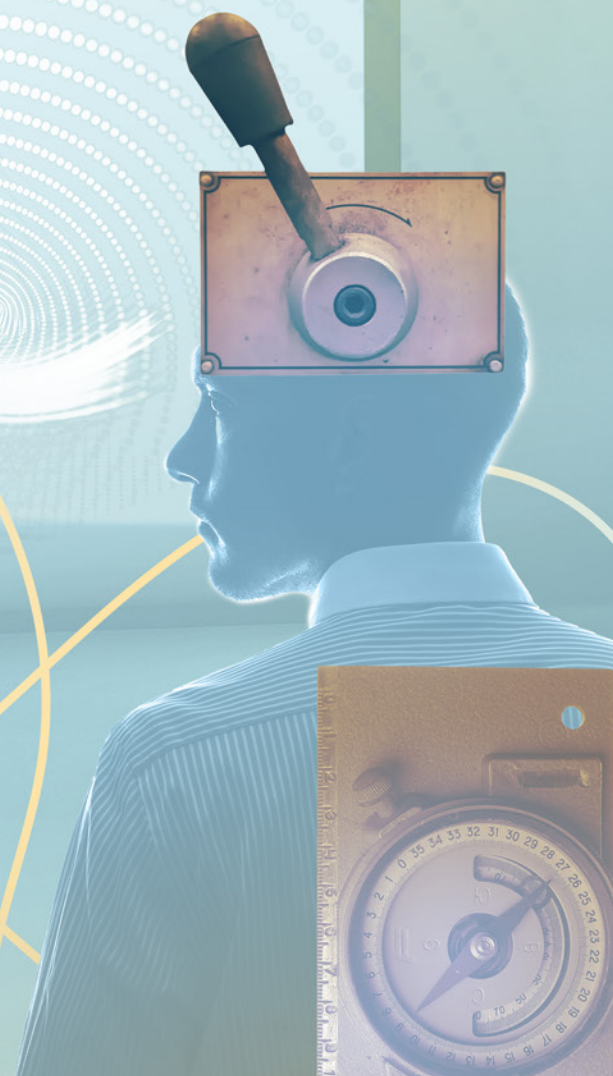
Humans remain one of the most consistent sources of cybersecurity risk, prone not only to malicious targeting like social engineering or fraud but even well-intentioned human errors like sharing a sensitive document with the wrong recipient or accessing unauthorized tools that can expose or leak sensitive data and create compliance issues.

One area of industry focus is minimizing susceptibility to social engineering threats through AI-powered security training, i.e., platforms that enable enterprises to simulate highly realistic phishing campaigns using deepfakes, cloned voices and context-aware prompts to mimic modern multi-modal attacks to augment awareness among the workforce.

Yet another avenue that AI has opened up for attackers is SEO poisoning, i.e., manipulating search results so a company's employees on the internet are redirected to malicious websites, fake login pages or corrupted tools during routine work. These tactics are particularly effective as they exploit legitimate user intent rather than deception alone. Defending against such tactics combines behavioral signals with browser telemetry and business context to warn users or block access in real-time when high-risk patterns are observed.

Many solutions are also emerging to manage fraud risk across users, customers and clients such as detecting fraud patterns in day-to-day communication channels or flagging abnormal behaviors in customer support flows. In all cases, human behavior becomes a security signal in and of itself.

As the threat landscape evolves, so too do the tools that help manage human behavior as part of the security architecture. Enterprises are looking to embed behavioral context directly into real-time policy enforcement decisions. This goes beyond role-based access and policy checklists. For example, rather than sweepingly banning any and all attachments to external email addresses, personal documents being shared with the employee's own personal email could be selectively allowed. Conversely, deviations from known behavioral baselines should trigger escalations or controls, even if policy would otherwise allow the action. This approach underpins dynamic policy-enforcement blending human behavior with telemetry from browsers, file systems and identity platforms.



Market and industry perspectives

In a recent IBM survey, 74% of Chief Information Security Officers (CISOs) ranked human error as the top risk to their organizations' cybersecurity.⁷

Large platform providers are recognizing that human-targeted attacks require ecosystem-level defenses. Google's recent partnership with Doppel, an emergent brand protection solution, for example, focuses on combining telemetry for faster identification and takedown of risky domains. At the same time, a new wave of solutions are focused on using AI to outpace the complexity of AI-powered attacks by using AI-powered simulations to train the workforce against phishing attempts and fraud.

Beyond the workforce, multi-channel fraud prevention solutions are working toward applying similar behavioral and real-time signals to detect customer-facing fraud across payments, messaging and digital media.

This shift is evident across various layers of security. Startups are positioning products around understanding baseline human behavioral patterns and AI powered risk reasoning solutions to allow organizations to act on human behavior before it comes an incident. Further, companies are framing the internet browser as the place where human intent, navigation patterns, and interactions can be observed and used to inform policy enforcement decisions in real-time.

⁷ Gregory, J. (2024, August 15). CISOs list human error as top cybersecurity risk. Ibm.com. <https://www.ibm.com/think/insights/cisos-list-human-error-top-cybersecurity-risk>





2

Inference demand drives continued AI buildout

The intense AI infrastructure buildout that has been taking place for years shows no signs of slowing down and the environments that have come online have become a strong foundation for significant software innovation and ecosystem development. Compute architectures will evolve into hybrid classical-quantum workflows and neuromorphic systems, leveraging CPUs, GPUs and specialized accelerators to address diverse computational demands across the infrastructure landscape. The progress in quantum hardware will lead to many scientific quantum advantage demonstrations showing a quantum computer performing a scientifically interesting computation that is impossible classically. These demonstrations will be uncontested by classical techniques and will fuel the race to demonstrate commercially relevant quantum advantage in the years to come.

AI infrastructure innovation

Innovation across power generation, data center design, silicon architectures and enabling technologies to meet soaring demand for computing resources.

The AI revolution continues to push the limits of existing infrastructure capabilities, and the insatiable demand for computing resources continues to drive the ecosystem to innovate across the entire stack. In addition, heterogenous compute (i.e. different chips for different use cases) continues to gain market traction in an effort to provide best fit for purpose designs to improve overall efficiencies.

The ongoing innovations can be grouped into four main categories:

New power generation methodologies: Gigawatt DC facilities and overall expansion have led to fierce competition for utility lockups and novel energy sources like Small Modular Nuclear Reactor development.

New datacenter designs and cooling strategies: Demand is accelerating facilities trends like modularity, multi-story, high density racks, and AI edge locations. Meanwhile continued power densities are changing how current compute is delivered to each server and how racks of servers are cooled using liquid (including network and storage devices).

New silicon architectures: The chip market continues to expand, with increasing competition from the clouds, traditional silicon companies and the startup ecosystem leading to incredible performance leaps and more cost-efficient alternatives. Custom AI silicon (e.g. chips focused on inference) have continued to gain traction and are being considered for future architectures like AI at the edge.

Enabling technologies: New networking transfer rates, optical technology, Data Processing Units, Disaggregated Storage, Memory Sharing, In-Storage Computing – there are dozens of infrastructure breakthroughs gaining traction as part of this ecosystem’s massive capex build out.

The Infrastructure build out is now multiple years in, and the leading question from industry pundits is just how long will it last. Many of these developments are by nature multi-year investments, and 2026 looks poised to continue this trend as many new facilities start to come online.



Market and industry perspectives

Global AI Infrastructure Capex spending in 2025 surpassed \$400B, with recent 2026 projections exceeding \$600B.⁸

Cloud and colocation providers such as AWS, Microsoft, Google, Meta, and Core-Weave are building next-gen data centers optimized for AI workloads, often co-located near renewable or modular power sources, including nuclear microreactors, geothermal plants and hydrogen fuel cells.

The competition in silicon development is intensifying as organizations work to introduce new types of accelerators. Recent industry activity includes acquisitions of licensing rights and talent from startups, reflecting a trend toward workload-specific compute to complement general-purpose accelerators. Providers are also making significant progress through ongoing improvements in hardware performance and software integration.

Hyperscalers continue to invest heavily in their own silicon, and these chips are now more broadly available, and some can even run in a customer's data center. 2026 may likely be a turning point year where these alternatives start to gain more meaningful adoption by customers in the market, having reached significant maturity and demonstrated demand from key AI labs.

Markets beyond silicon are also experiencing significant evolution. Major cloud service providers are advancing their networking capabilities, while startups focused on photonic networking are beginning to gain traction, with recent acquisitions occurring in this area. Companies specializing in storage solutions are seeing substantial growth, with large-scale partnerships being formed to deliver advanced data services to a broad range of customers across the technology stack.

Finally, "AI at the edge" remains an evolving strategy across the industry from CDN companies building out AI compute across their locations, to low-powered custom inference silicon targeting edge use cases and the continued development of AI chipsets embedded in personal computing.

⁸ Big Tech to invest about \$650 billion in AI in 2026, Bridgewater says. Reuters. (2026, February 23). <https://www.reuters.com/business/big-tech-invest-about-650-billion-ai-2026-bridgewater-says-2026-02-23/>.

Cloud native AI

Accelerating efficient, scalable, and cost-effective AI through open-source tools and advanced optimization techniques.

As more GenAI use cases are scaled into production, there is significant industry focus on efficiently managing infrastructure resources through cloud native software constructs.

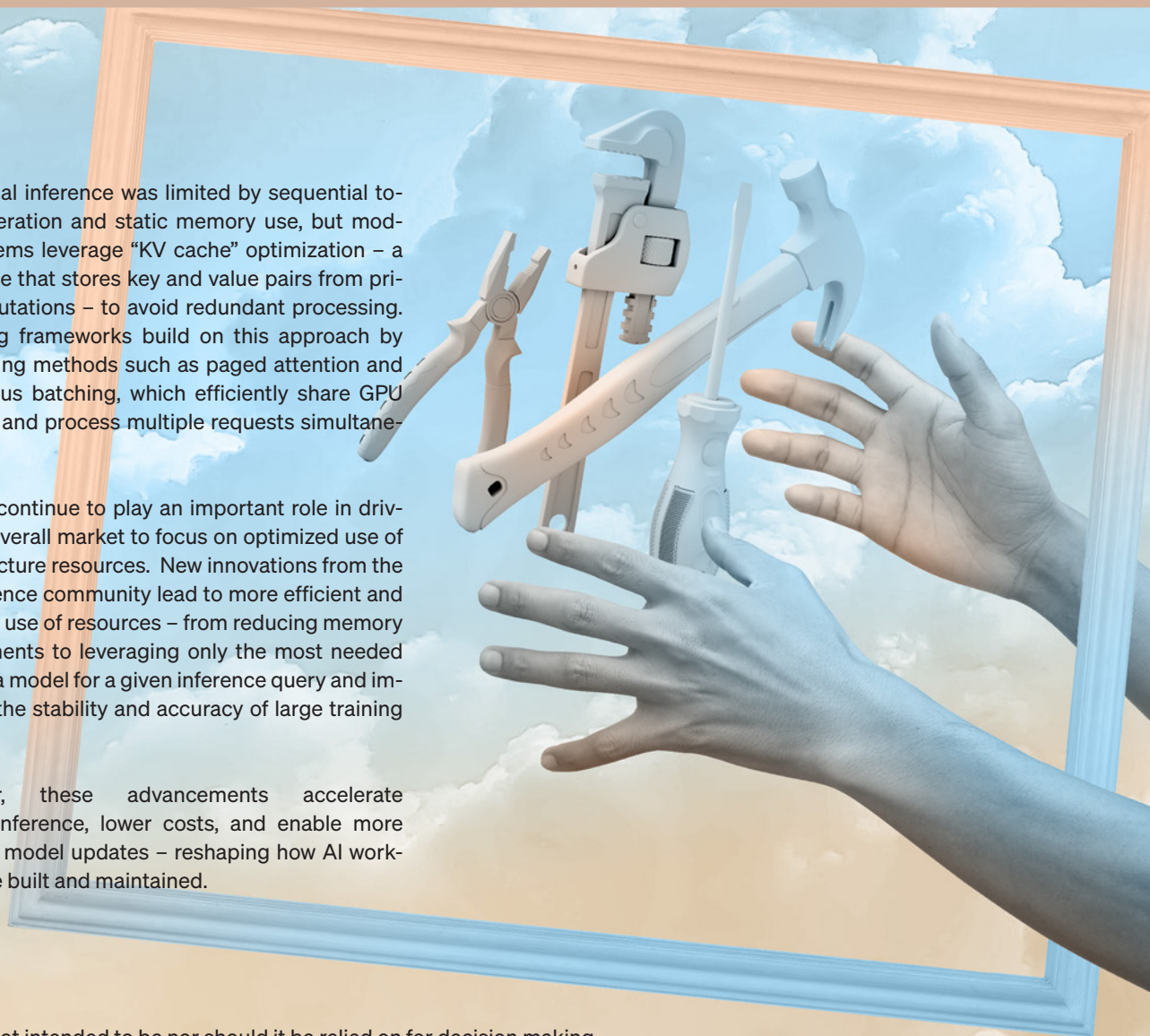
The ecosystem is working to develop common and familiar toolsets that abstract hardware complexity and provide workload portability. The result is a competitive landscape across proprietary, specialized software stacks and, increasingly, an open-source approach to managing inference workloads using standard architectures like Linux and Kubernetes.

While training is mostly done offline, inference is done in production and needs to be designed to support things like multi-tenancy, scalability, resiliency and high availability, posing challenges at enterprise scale. New inference engines and distributed serving frameworks are transforming the way LLMs operate in production by optimizing the inference process, rather than the model architecture itself.

Traditional inference was limited by sequential token generation and static memory use, but modern systems leverage “KV cache” optimization – a technique that stores key and value pairs from prior computations – to avoid redundant processing. Emerging frameworks build on this approach by introducing methods such as paged attention and continuous batching, which efficiently share GPU memory and process multiple requests simultaneously.

AI Labs continue to play an important role in driving the overall market to focus on optimized use of infrastructure resources. New innovations from the data science community lead to more efficient and accurate use of resources – from reducing memory requirements to leveraging only the most needed parts of a model for a given inference query and improving the stability and accuracy of large training runs.

Together, these advancements accelerate time-to-inference, lower costs, and enable more frequent model updates – reshaping how AI workloads are built and maintained.



Market and industry perspectives

The market for AI infrastructure optimization is expanding rapidly as enterprises seek greater efficiency and control over AI costs. According to industry estimates, the AI inference optimization market alone is projected to exceed \$100B over the next five years as more startups focus on providing platform technology so reduced overall AI infrastructures costs.⁹

The Cloud Native Ecosystem (i.e. Linux / Kubernetes) is embracing emerging projects like vLLM (inference engine) and llm-d (optimize and manage inference at scale) to promote consistent and powerful inference patterns. Emerging startups are differentiating through software-defined optimization and low-level programming techniques to improve price to performance.

Large web scale companies and research institutions are developing and open sourcing many new projects, including serving engines, disaggregation, and caching techniques, etc.

AI Labs continue to publish new research focused on infrastructure optimization. Key innovators have brought numerous techniques to market, from “Multi-Head Latent Attention” and new “Mixture of Expert” approaches for inference, and recent breakthroughs to improve Training stability which enable larger experiments on less reliable hardware.

Ecosystem convergence of training and inference runtimes will likely occur over time to help abstract the complexity of AI Infrastructure management which may lead to more overall use of AI and the increased deployment of smaller models that can easily be hosted across data center and edge deployments.

⁹ PR Newswire. (2025, February 28). Ai inference market worth \$254.98 billion by 2030 - exclusive report by MarketsandMarketsTM. Yahoo! Finance. <https://finance.yahoo.com/news/ai-inference-market-worth-254-151500286.html>.

Quantum computing: from research breakthroughs to real-world use

Hardware advancements will enable scientific quantum advantage.

Quantum computers (QCs) are transforming specialized computing beyond classical limits. QCs are fundamentally different from classical computers and AI, which handle general-purpose, sequential tasks like logic operations and pattern recognition. By manipulating the probabilities of all possible states simultaneously, QCs leverage quantum mechanics to achieve dramatic speedups for specialized problems – such as simulations, optimization, and cryptography – that classical computers cannot efficiently solve. Ultimately, QCs and classical computers are expected to be complementary, with classical systems handling broad tasks and QCs focusing on specialized, complex challenges.

Recent breakthroughs are overcoming key engineering barriers to fault tolerance. Fault tolerance is the key goal for all QCs, as it enables them to perform long, reliable computations despite errors in physical gates or measurements – a milestone necessary for transformative business impact. Achieving fault tolerance requires encoding logical qubits

across many physical qubits and applying continuous error correction. Recent breakthroughs, such as Google's Willow processor surpassing the surface-code threshold¹⁰ and Quantinuum's demonstration of fully fault-tolerant universal gate set with repeatable error correction¹¹, signal a new phase for the industry. Numerous leading companies estimate to release fault-tolerant quantum computers by 2028–2030, aiming to scale from 30–1,200 qubits today to $1M^{12,13,14}$, which experts believe is needed to break current cryptography and unlock world-changing applications.

QCs are playing an expanding role across industries, with progress in sectors such as healthcare, material science, and finance, enabling solutions to complex problems that were previously unsolvable. From an AI perspective, quantum computing will generate data beyond classical simulation, providing valuable input for AI models to advance scientific understanding.



Market and industry perspectives

The quantum computing industry is highly fragmented, with four to seven main modalities under active research. Each approach uses different setups, quantum particles, and temperatures, facing distinct engineering hurdles that make progress uneven and interdependent, and it remains uncertain which will achieve large-scale real-world impact first. Major tech firms are developing proprietary technologies, forming partnerships, or investing in private companies: focus areas include, superconducting qubits, trapped ions, photonics, neutral atoms.

Political support for quantum computing is accelerating both in US and internationally. The US has made quantum computing a national priority since the National Quantum Initiative Act of 2018, with recent executive orders accelerating research and development (R&D), post-quantum cryptography (PQC) and federal adoption timelines. The administration is considering equity stakes in QC firms and pushing for quantum-resistant upgrades by 2030. Major corporations have aligned with these efforts, including JPMorganChase which named a quantum computing a focus in the firm's \$1.5T Security and Resiliency Initiative. Internationally, China leads in government investment and strategic planning, while Europe boasts strong scientific leadership and public funding.

The capital market for quantum computing has experienced strong momentum. Driven by AI's substantial impact on public market capitalizations and private valuations over the past 2–3 years, investors are actively seeking the next major growth opportunity. QC has emerged as a leading area of interest, viewed as a promising new investment frontier.

The long-term revenue potential for quantum computing remains strong. Over the next 2–4 years, industry revenue is projected to reach \$1B, primarily driven by early-stage development and testing.¹⁵ Looking ahead to 2040, consulting firms like McKinsey and BCG estimate a \$100B market opportunity for quantum computing providers, who are expected to capture 20% of the \$500B in total economic value generated by the industry.^{16,17}

¹⁰ Google Quantum AI and Collaborators. Quantum error correction below the surface code threshold. *Nature* 638, 920–926 (2025). <https://doi.org/10.1038/s41586-024-08449-y>

¹¹ Quantinuum, Breaking even with magic: demonstration of a high-fidelity logical non-Clifford gate. As of February 2025.

¹² Craig Gidney. How to factor 2048 bit RSA integers with less than a million noisy qubits. <https://doi.org/10.48550/arXiv.2505.15917>. As of May 2025.

¹³ Ben Bloom, CEO of Atom Computing, Quest for qubits: Quantum computing leaders make their case at Nvidia GTC. As of March 2025.

¹⁴ PsiQuantum Raises \$1 Billion to Build Million-Qubit Scale, Fault-Tolerant Quantum Computers. (2023, February). Psiquantum. Retrieved from <https://www.psiquantum.com/news-import/psiquantum-1b-fundraise>.

¹⁵ Bank of America Institute. Quantum Leaps and Bounds. (2025, October 23). <https://institute.bankofamerica.com/content/dam/transformation/quantum-computing.pdf>

¹⁶ Jean-Francois Bobier, Matt Langione, Cassia Naudet-Baulieu, Zheng Cui, and Eitoku Watanabe. The Long-Term Forecast for Quantum Computing Still Looks Bright. BCG. (2024, July 18). <https://www.bcg.com/publications/2024/long-term-forecast-for-quantum-computing-still-looks-bright>.

¹⁷ McKinsey & Company. Quantum Technology Monitor. (2024, April). <https://www.mckinsey.com/~media/mckinsey/business%20functions/mckinsey%20digital/our%20insights/steady%20progress%20in%20approaching%20the%20quantum%20advantage/quantum-technology-monitor-april-2024.pdf?pdf>.





3

The end of app switching; Intent is the new interface

The dominant interface won't be an app, a browser or a workspace. It will be a single AI-native environment that collapses every workflow into one continuously personalized stream. Users will stop navigating between tools or tabs and start inhabiting intelligent interfaces that anticipate, execute and transact across every modality of work and life.

Agentic browsers

The emergence of agentic browsers is transforming web navigation, positioning the browser as the future primary interface for human-computer interaction.

The traditional web browser is no longer a gateway to websites or a passive tool for document retrieval and navigation. Browsers are rapidly becoming intelligent, agent-powered environments that understand intent and carry out tasks on behalf of the user. After decades of static web navigation defined by clicks and links, these platforms are beginning to collapse complex multi-step workflows into single conversational interactions.

Instead of searching, opening multiple tabs, filling forms, comparing sources, the user can now state their intent, and the browser handles the rest. Evolving into dynamic, context-aware assistants, agentic browsers can summarize content, complete tasks and execute multi-step processes like booking a reservation or comparing products across websites, often without the user visiting any of them.

As the preeminent digital front door, companies are racing to control this new paradigm. Their goal is to make the browser the central point of interaction between users and intelligent systems, embedding AI into the core experience so that tasks once requiring multiple sites, tabs and forms can now be completed in a single step.

These systems work by combining large language models with retrieval engines and persistent memory, allowing them to interpret user intent and carry context across sessions. By maintaining persistent context and (with consent) accessing user credentials, history and preferences, these browsers can execute tasks autonomously.

The convergence of search, browsing, and agentic execution means users no longer merely visit websites, but interact with intelligent systems that can navigate the web, conduct multi-step reasoning, and take action on behalf of users. In addition, these browsers are adding transaction capabilities via emerging agentic payments protocols (Agent Payments Protocol - AP2, Agentic Commerce Protocol - ACP, x402 etc.) enabling autonomous transactions.

As this new era of agentic commerce takes hold, Generative Engine Optimization (GEO) becomes critical, with AI-driven agents now translating intent into transactions, frequently bypassing conventional research and comparison. Winning brands will focus on influencing the signals and rules that drive agent decision-making, not just consumer discovery.



Market and industry perspectives

The global AI browser market size is projected to grow from \$4.5B in 2024 to \$76.8B by 2034, growing at a CAGR of 32.8% during the forecast period from 2025 to 2034.¹⁸

Traditionally dominated by Google (Chrome) and Microsoft (Edge), the AI browser market is experiencing rapid growth, driven by advancements in AI and increasing demand for personalized and automated browsing experiences. While legacy players move to embed AI features into existing interfaces, AI-native challengers are redefining the web browser space.

Companies are designing browsers to act as personal AI assistants or “thought partners”, enabling agents to summarize web pages and email, draft responses, access calendars, and even make purchases.

The advent of AI browsers leads to growing privacy concerns, such as user data being collected, stored and used to train future models or potential intellectual property issues related to ‘all-seeing’ browsers. Some players are pursuing local processing approaches, with others maintaining user privacy through processing requests without collecting personal data.

While available to general consumers, AI browsers also offer opportunity for employees. According to the Greyhound CIO Pulse 2025 report, 42% of enterprise leaders say they are actively exploring or experimenting with AI-native browsers for use cases such as research, planning and automation.¹⁹

¹⁸ Market.US. (2025, July). Global AI Browser Market Size, Share Analysis Report by Type. <https://market.us/report/ai-browser-market/>.

¹⁹ AI browsers: The future of digital navigation?. Greyhound Research. (2025, July 21). <https://greyhoundresearch.com/ai-browsers-the-future-of-digital-navigation/>.

AI native workspaces

Embedding intelligent agents directly into the flow of work to redefine productivity, collaboration and decision making.

After decades of app-based productivity, a new generation of AI-native workspaces is emerging – environments built around autonomous agents that interpret intent, orchestrate workflows and execute actions across systems.

Unlike traditional tools that require employees to adapt to static applications, users will begin to operate within AI-native workspaces that dynamically adjust to their intent, using contextual understanding, memory and reasoning to anticipate needs and act proactively.

This marks a shift from app-centric work to agent-centric work, unifying documents, communication and workflows into a single intelligent layer. Instead of switching between spreadsheets, emails and chat threads, employees will interact conversationally with an agent or number of agents which span their digital environment. Moreover, agents will move from reactive assistance to proactive execution. Instead of waiting for commands, agents will

continuously monitor and understand the context of ongoing projects, surface insights, complete actions proactively and deliver completed outputs.

AI-native workspaces retain memory across sessions and integrate with enterprise data sources, transforming fragmented tools into a unified knowledge base for contextual search and instant information retrieval. With user consent, they access calendars, Customer Relationship Management (CRM) systems, and communication channels to automate workflows and anticipate user needs. The most innovative platforms go further, planning and executing multi-step tasks across applications while maintaining state and permissions. In this model, AI agents act as proactive digital coworkers, moving beyond passive assistance to deliver truly intelligent support.



Market and industry perspectives

The workspace is becoming a key competitive front in the race to embed AI into the enterprise stack.

Large technology players are embedding this into existing ecosystems. Google and Microsoft have extended Gemini and Copilot (alongside Agent Mode) respectively across Google Workspace, Teams and Outlook to unlock agentic productivity across apps.

Meanwhile, AI-first startups are redefining the category by building collaboration environments with persistent memory and enabling agents to prioritize, execute and remember tasks across projects.

Established productivity players are embedding agentic features evolving toward self-organizing work environments.

Agentic wearables are set to become extensions of AI native workspaces, enabling agents to accompany employees beyond the desktop, capturing real-time context from meetings, providing insights and allowing for seamless, voice-driven interaction with enterprise systems on the move.

Similar to the AI Native Workspace, agentic web browsers have hit the market with the ability to autonomously search the internet, open new tabs and even shop for users using protocols like A2P, further highlighting a shift from passive tools to active, intelligent environments helping users achieve their goals.

Generative user experiences

Harnessing AI to dynamically tailor digital experiences, delivering deeply personalized journeys for every user.

Hyper-personalization has become the new benchmark for digital experiences. Today's customers expect every scroll, swipe and click to feel uniquely tailored to their preferences, behaviors and needs. This shift in expectations is driving organizations to rethink how they design and deliver digital interactions.

At the forefront of this transformation are Generative User Interfaces (GenUIs), which represent a fundamental departure from traditional approaches. Historically, digital platforms relied on manually designed screens and static layouts, offering limited flexibility and personalization. In contrast, GenUIs harness the power of advanced AI models to dynamically generate interactive and visually compelling interfaces in real time.

These intelligent systems continuously analyze live context, including user behavior, stated and inferred preferences and even intent, to assemble layouts and content that are uniquely suited to each individual and each

moment. As users engage with these platforms, GenUIs learn from every interaction, creating a self-improving feedback loop that refines and enhances the experience over time.

As a result, two users accessing the same application may encounter entirely different interfaces, each one optimized for their specific needs, circumstances, and goals. This level of adaptability results in a dynamic, responsive environment that evolves in real time, moving decisively away from one-size-fits-all solutions and ushering in a new era of truly personalized engagement.



Market and industry perspectives

Across industries, the adoption of GenUIs is rapidly accelerating as organizations increasingly recognize the strategic value of hyper-personalization.

Traditionally, digital platforms have relied on broad segmentation strategies – such as displaying different homepage banners to millennials in urban areas versus suburban parents. While this approach offered some degree of customization, it falls short of meeting today’s rising expectations for truly individualized experiences.

The next wave of innovation is now driving toward one-to-one customization at scale, where every aspect of the user experience is tailored to the individual. Market data shows 91% of consumers prefer brands that offer personalized experiences and 71% expect companies to deliver personalized interactions²⁰, with 76% expressing frustration when those expectations are not met.²¹

Industry leaders are responding to this trend with enhanced solutions, delivering fine-tuned, contextually relevant experiences and integrations. The business impact is clear - companies implementing AI-driven personalization have seen a 35% increase in purchase frequency and a 21% boost in average order value.²²

Looking ahead, AI-powered personalization will enable websites and applications to instantly reconfigure themselves based on who is interacting, factoring in context, timing, and even external events such as weather or stock market movements. Predictive personalization engines are becoming increasingly sophisticated, analyzing signals from browsing patterns, location, and real-world events to anticipate user needs before they are even expressed.

²⁰ Accenture. (2018, May 3). Widening Gap Between Consumer Expectations and Reality in Personalization Signals Warning for Brands, Accenture Interactive Research Finds. <https://newsroom.accenture.com/news/2018/widening-gap-between-consumer-expectations-and-reality-in-personalization-signals-warning-for-brands-accenture-interactive-research-finds>.

²¹ McKinsey & Company (2021, November 12). The value of getting personalization right—or wrong—is multiplying. McKinsey & Company. <https://www.mckinsey.com/capabilities/growth-marketing-and-sales/our-insights/the-value-of-getting-personalization-right-or-wrong-is-multiplying>.

²² Michael Casner. (2025, July 16). “2026 AI Marketing Predictions: What Marketing Directors Need to Prepare For. <https://magnet.co/articles/2026-ai-marketing-predictions>.

Multi-modal social listening

Instantly analyze customer feedback across all channels — messages, videos, and audio — enabling real time insights and hyper-personalized engagement.

The digital landscape has fully transitioned to a video-first and audio-rich environment, rendering legacy, text-based monitoring obsolete. Multi-modal social listening represents a critical evolution in market intelligence, leveraging advanced AI to synthesize data across text, image, audio, and video. By moving beyond simple keyword tracking, this approach captures the full spectrum of consumer behavior – such as brand logos appearing in untagged videos on social media or sentiment expressed through vocal inflection in podcasts – closing a significant gap that previously left organizations blind to organic brand interactions.

The primary strategic advantage of this holistic view is the transition from reactive response to proactive anticipation. By aggregating multi-modal data in real

time, enterprises can detect emerging crises or viral trends before they reach mainstream text platforms. This unified intelligence allows teams to identify unbranded reach and quantify the true impact of visual product placement. Emerging players in the space focus on providing the speed and accuracy required to manage brand reputation in a high-velocity digital economy.

Ultimately, multi-modal listening delivers a 360-degree view of the audience that informs decision-making across R&D, marketing and customer experience. Insights derived from how products are visually used in the real world directly influence product development, while audio-based sentiment analysis refines customer service strategies. For the modern enterprise, adopting a multi-modal framework is no longer an optional innovation - it is a foundational requirement to remain responsive and relevant in a landscape where consumer attention is fragmented across diverse, non-textual media.



Market and industry perspectives

The rapid growth of multi-modal social listening is a market response to a shift in consumer behavior. As of 2026, the average adult in the US spends over 60% of their daily screen time consuming digital video content, with video content estimated to comprise 82% of all global internet traffic.²³ These figures highlight a staggering intelligence gap for organizations still relying on text-based monitoring. The multi-modal AI market is projected to reach \$2.83 billion in 2026²³ alone as enterprises rush to capture the sentiment and context embedded in these dominant media formats.

This technological shift is essential as consumers are increasingly bypassing text in favor of rich media to interact with brands. In 2026, over 25% of online audiences watched a brand-produced video in the last month, while 85% of people report being convinced to purchase a product after watching a video. In the audio sphere, podcasts have gained significant traction.²⁴ Without multi-modal tools, the vast majority of these high-intent interactions remain invisible to traditional analytics platforms.

The strategic incentive for early adoption is clear: organizations that integrate these advanced tools report significantly higher confidence in their performance. Research shows that 76% of social listeners feel confident in their ROI on visual platforms like Instagram and LinkedIn, compared to less than 60% of non-listeners.²⁵ By leveraging unified metrics that span text, video, and audio, enterprises are not just monitoring conversations, but they are future-proofing their ability to participate in a market where visual quality and audio sentiment now dictate brand trust for nearly 90% of consumers.²⁵

²³ The Rising Demand for Video Content in 2026: Why Every Creator Needs to Repurpose Articles into Videos. Medium. (2025, September 19). <https://medium.com/@aboda.bob7/the-rising-demand-for-video-content-in-2026-why-every-creator-needs-to-repurpose-articles-into-f6fe993db5b0>; The Business Research Company. (2026, March). Multimodal AI Market Report 2026. <https://www.thebusinessresearchcompany.com/report/multimodal-ai-global-market-report#:~:text=What%20is%20The%20Multimodal%20AI,analytics%2C%20demand%20for%20intelligent%20automation>.

²⁴ Alexandra Bjertnaes. (2025, November 13). What Social Media Audiences Want in 2026, By The Numbers. <https://adage.com/studio-30/aa-what-social-media-audiences-want-in-2026-by-the-numbers/#:~:text=After%20a%20brand's%20website%2C%20its,a%20brand's%20email%20or%20newsletter>.

²⁵ Jeremy Gooldman. (2025, October 26). Premium Media Lifts Purchase Intent 40%, Brand Trust 85%, Says Study. <https://www.emarketer.com/content/premium-media-lifts-purchase-intent-brand-trust-study>; Talkwalker (2025, May 26). 128 Must-Know Social Media Statistics for 2025. <https://www.talkwalker.com/blog/social-media-statistics>.

Agentic wearables

Agentic wearables are transforming personal and professional technology by proactively supporting users, automating tasks, and enabling seamless, intelligent interactions.

Unlike traditional wearables, which primarily collect and display data, agentic wearables are designed to act as proactive digital companions. These devices leverage advanced artificial intelligence, contextual awareness and edge computing to interpret user intent and autonomously take action, often without explicit commands. The result is a new class of technology that seamlessly integrates into daily life, offering continuous, intelligent support.

The defining feature of agentic wearables is their ability to exercise agency. Through unobtrusive designs – ranging from smart bands and AR glasses to discreet earpieces and pendants – these devices are engineered for comfort and continuous use. They are always on, continuously sensing the user's environment, activities and interactions. Sophisticated AI algorithms interpret these real-time datapoints to understand context, intent and even emotional state, enabling the device to anticipate needs and proactively assist the user.

The capabilities of agentic wearables extend far beyond passive data collection. For example, these devices can transcribe and summarize conversations in real time, provide contextual reminders for tasks and appointments, and automate scheduling and follow-ups. By offloading cognitive burdens – such as remembering tasks, organizing information and making

routine decisions, these devices empower users to focus on strategic, high-value activities. The always-on nature of agentic wearables ensures that support is available whenever needed, not just when requested. For enterprises, this translates into enhanced employee productivity, streamlined administrative tasks, improved compliance and more efficient workflows.

However, the adoption of agentic wearables also introduces security and risk considerations. The continuous sensing capabilities of these devices raise concerns about privacy, particularly regarding the collection and use of sensitive audio, video, and contextual data. Robust encryption, secure on-device processing and transparent user controls are essential to protect both personal and enterprise information. Regulatory compliance is another critical factor, especially in industries such as finance and healthcare, where data protection standards are stringent. Devices must adhere to relevant regulations and provide clear consent mechanisms and audit trails.

User consent and trust are also paramount. Individuals must be confident that their data is handled responsibly and that they retain control over what is collected and shared. Manufacturers and service providers must prioritize responsible AI practices and communicate transparently about device capabilities and limitations. Operational risks such as device reliability and integration with legacy systems must also be managed to ensure seamless adoption and minimize disruption.



Market and industry perspectives

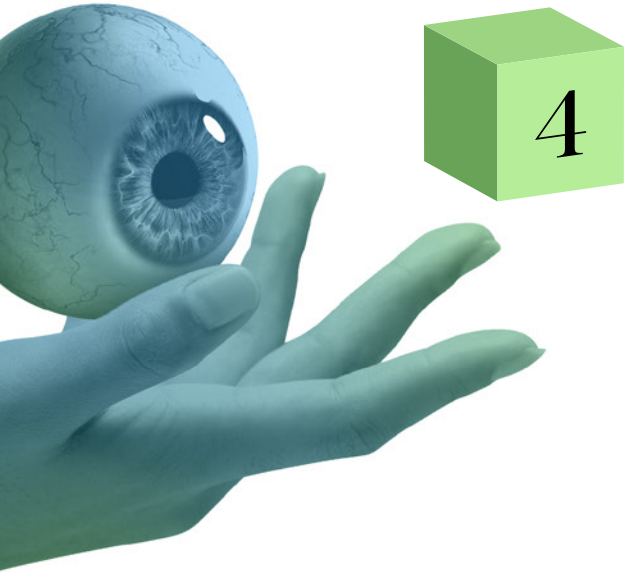
According to Grandview research, the agentic wearables market is entering a period of accelerated growth, with the global wearable AI sector projected to achieve a CAGR exceeding 27% over the next five years. (Grandview Research 2025).²⁶ This is driven by advances in edge computing, Natural Language Processing (NLP) and hardware miniaturization.

Innovation in agentic wearables is accelerating across several interconnected market categories. Tech giants like Amazon, Microsoft, Apple, and Google are shaping the industry by integrating hardware, cloud infrastructure and advanced AI models. Wearable device makers and emerging startups are focused on designing smart bands, AR glasses, and discreet earpieces that deliver seamless, proactive user experiences. Meanwhile, AI and contextual intelligence providers are powering these devices with sophisticated voice recognition, intent detection, and real-time audio processing. The integration of voice interfaces into apps is also accelerating, enabling hands-free, natural interaction with technology.

Underpinning all of this, emerging chipset and hardware infrastructure companies supply the ultra-low-power processors and sensors that enable always-on, real-time intelligence in compact wearable form factors.

²⁶ Grandview Research. (2025). Wearable AI Market (2026 – 2033). <https://www.grandviewresearch.com/industry-analysis/wearable-ai-market-report#:~:text=The%20global%20wearable%20AI%20market,factors%20contributing%20to%20market%20growth>.





4

AI powered simulation enhances testing

Organizations will rely on advanced simulations to test, validate, and optimize products, processes, and scenarios before deploying them in the real world. These virtual environments will enable continuous, scalable experimentation, allowing teams to model user interactions, system exposures, and operational outcomes with speed and precision.

Synthetic users / user simulations

Creating synthetic representations of real people or populations to test ideas, gain insights and simulate behavior and action.

Synthetic users are virtual personas or agents powered by large language models (LLMs), designed to mimic the humanness of real people or populations – including beliefs, intent, preferences and choice actions.

Synthetic users can represent individuals, segments, or entire populations. These synthetic users can complete interviews, surveys and panels for user research and consumer feedback, test and optimize marketing materials such as text, images and video, and even interact with designs, applications, and experiences. In some instances, these users can be deployed into a synthetic environment to conduct simulations of interactions and information/idea dissemination – think a video game where synthetic users would be able to live, interact, learn and grow.

For example, User Research teams could generate synthetic users that represent hard to reach market segments (like traders). Likewise, UI/UX teams can create synthetic users to test new workflows, validate user interfaces or simulate employee interactions with enterprise software.

There are different model frameworks which can be leveraged to create synthetic users. Some vendors rely on singular models, such as ChatGPT, that serve as wrappers. Others use model switching, dynamically selecting the most suitable LLM for a given task. A third approach involves building and training proprietary foundational models tailored to specific needs.

Training strategies also vary across the industry. Some companies train their models exclusively on public data, while others use proprietary data; or a combination of both. Additionally, there is debate over the necessity of injecting demographic, psychometric, and behavioral data to enhance the humanness of synthetic agents. Some believe that LLMs, already trained on vast internet data, inherently capture human characteristics, while others advocate for more targeted data enrichment.

Approaches for generating synthetic users can vary. Some providers create purely synthetic users, generated solely by AI models, while others augment these virtual personas with real human feedback. For example, a synthetic user might be regularly updated based on recurring interviews with actual people, as opposed to being generated once from a model like ChatGPT and left unchanged.



Market and industry perspectives

The synthetic user space is an emerging market, with most players still in the early stages of company formation and investment. Many began as research-focused initiatives and are now transitioning into commercial ventures, seeking to capture both revenue and market share.

Within the market, four main categories have emerged, each representing a different level of sophistication and capability:

Synthetic research – These solutions extend research efforts to fill gaps in representation but may lack detailed audience segmentation.

Comprehensive synthetic research – Platforms in this category create synthetic users for insights and feedback, offering detailed segmentation and analysis. Interaction is generally limited to text and images.

Multimodal synthetic research – This segment enables synthetic users to engage with live stimuli, including text, images, video, audio, applications and websites, providing a richer research experience.

Multimodal agentic simulations – Representing the next generation of synthetic user technology, these platforms allow multiple synthetic agents to interact within simulated environments, offering deeper insights into group dynamics and behaviors.

The market is rapidly converging toward solutions that feature interactive personas that can engage with live stimuli, combining both public and proprietary data. These platforms can leverage a customer's private data and ground their insights in broader market data, creating highly representative and actionable research outputs.

This evolution is enabling newer synthetic user vendors to compete with, and in some cases augment, established traditional research and software tools. These include consumer research platforms, statistical modeling tools, survey platforms, and human-led research services and UX research.

Proactive defense through tailored, continuous attack simulation

Proactive testing and exposure management by continuously simulating attacker tactics on a dynamically adapted representation of the enterprise environment.

Cyberattacks are evolving at unprecedented speed, fueled by advances in AI, agentic systems, and increasingly automated adversary tradecraft. Attackers now can develop more dynamic, personalized, less traceable and continuous tactics and techniques, making it increasingly harder for defenders.

Digital twins help expose an enterprise's complete, real-world attack surface. The next frontier approach to defending against complex and scaled attacks is not limited explicitly to modeled twins but rather continuously learned real-time representations of the enterprise built by aggregating real-time telemetry from across the security and observability stack. By combining signals from the entirety of the digital environment, organizations can maintain a high-fidelity view of how their infrastructure is connected and exposed and use this model to continuously simulate attacks from an attacker's perspective.

Unlike traditional vulnerability scans or point-in-time penetration tests, this approach enables modeling of real-world, tailored and adaptive attack paths across the complete enterprise infrastructure including, network, cloud, data paths, endpoints, software supply chain, SaaS integrations and identity systems. Because they are derived from live telemetry, simulated attack tactics can dynamically adapt as configurations, permissions and dependencies change, enabling defenders to test exposure as it exists today, anticipate breach paths before they are exploited and harden configuration proactively.



Market and industry perspectives

Recent research of vulnerabilities exploited in a given year highlight a trend toward a negative time-to-exploit (TTE). The metric turning negative implies that attackers are now weaponizing vulnerabilities before patches are available, requiring a new expectation for vulnerability remediation and proactive resilience for effective defense. Beyond vulnerability exploitation, frontier AI labs have published research of advanced threat actors using their AI tools to create malware and execute AI-driven adversarial operation.

Leading security platforms are focusing on continuous exposure management, with incumbents expanding into continuous exposure testing. They envision a broader platform-centric view of security, integrating multi-step attack path simulation from an external attacker perspective across the full estate. Amazon recently unveiled its internally developed Autonomous Threat Analysis (ATA) system where multiple AI agents compete against one another to investigate real attack techniques against a high-fidelity simulated testing environment and then propose security controls for human review.

Several emerging startups are innovating with AI and digital twin-like concepts to simulate adversary emulation and derive a complete picture of an organization's exposure with an end goal objective of autonomous security built for the age of AI.

