



CAYLENT

# The 2025 Outlook on Generative AI

Insights for Business Leaders & Executives

**Randall Hunt**, CTO

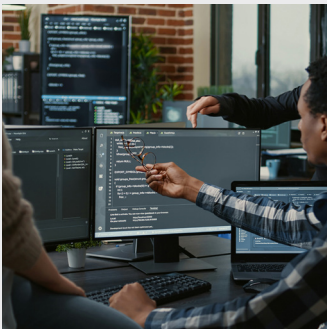
**Guille Ojeda**, Cloud Software Architect

**Anat Fraenkel**, GenAI Program Leader

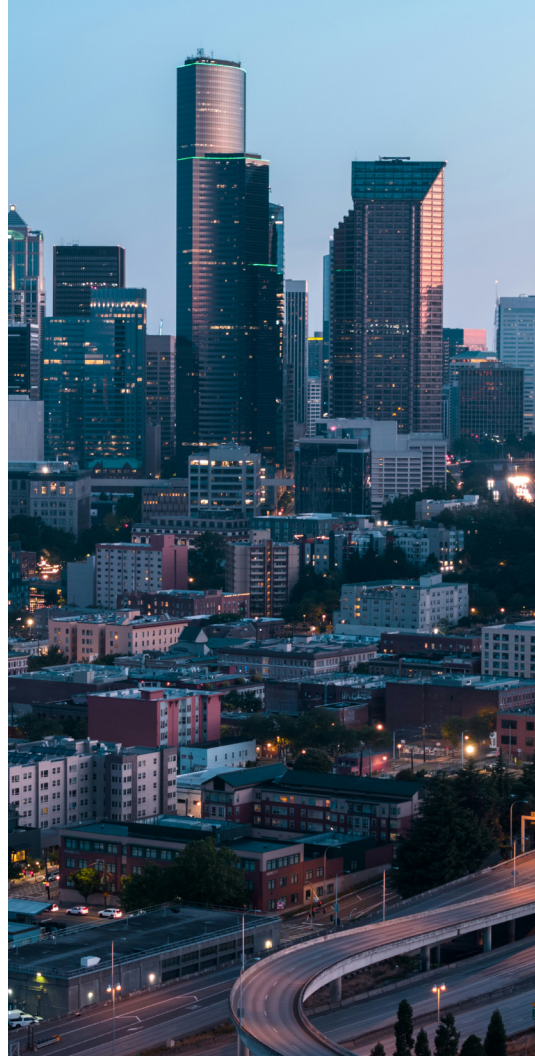


# Table of Contents

<b>Introduction</b>	<b>03</b>
Key Developments of 2024	04
Focus Areas for 2025	05
<b>Core Trends and Predictions</b>	<b>06</b>
Agentic Architectures	06
Optimization for Cost, Performance, and Security	07
Multimodal AI and Data Processing	08
Evolving Search and Discovery Technologies	09
Products Over Models	10
Decentralized AI and Federated Learning	11
AI Governance & LLMOps	12
Generative AI UI/UX	14
<b>Should You Build or Buy Generative AI Solutions?</b>	<b>16</b>
<b>Case Studies</b>	<b>17</b>
Multi-Agent Systems in Enterprise Workflows	17
Agentic AI-Powered Voice Service	18
<b>Looking Forward</b>	<b>19</b>
<b>Let's Get Started</b>	<b>20</b>







# Introduction

The generative AI landscape has fundamentally shifted from experimental technology to enterprise-ready solutions that drive measurable business value. Modern generative AI systems demonstrate sophisticated capabilities in orchestrating complex tasks, processing multiple types of data, and adapting to specific business contexts while maintaining robust security and governance frameworks.

Enterprise adoption of these technologies continues to accelerate, with Gartner projecting that [33% of enterprise software applications will incorporate agentic AI capabilities by 2028](#), up from less than 1% in 2024. Organizations implementing these systems

report significant operational improvements, with some [saving 70% of analysts' time](#) by automating manual data retrieval.

This whitepaper provides enterprise leaders with a practical framework for implementing generative AI technologies, focusing on cost optimization, operational excellence, and responsible development. Drawing from real-world implementations and industry research, we examine key trends, technical requirements, and strategic considerations that will shape successful deployments in 2025.

# Key Developments of 2024

This past year built the foundation for enterprise generative AI adoption.

## Architectural Advancement

The practical implementation of agent-based systems moved from theoretical frameworks to production deployments. Major cloud providers like AWS introduced robust orchestration frameworks, enabling organizations to implement complex multi-agent workflows while maintaining security and governance controls. A notable example is [Amazon Bedrock's Multi-Agent Orchestration](#), which was announced in preview in December 2024. This service provides enterprises with a scalable framework for implementing agentic workflows.

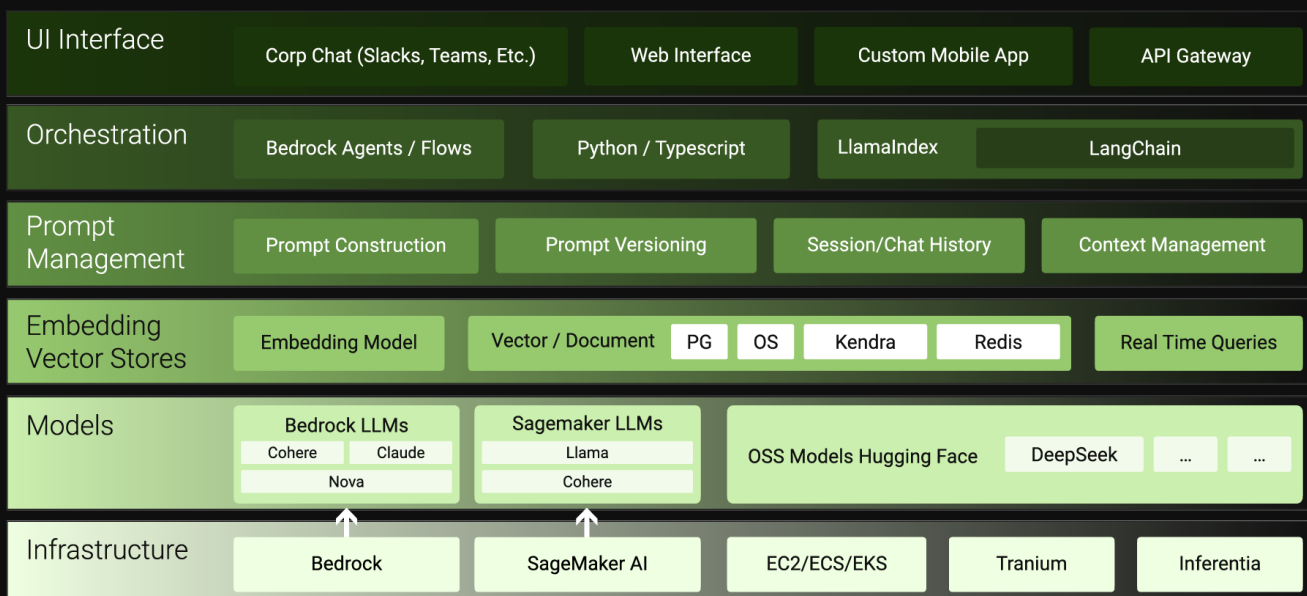
## Performance Optimization

Significant progress in model distillation and inference optimization addressed the computational and economic challenges of large-scale deployments. These advances make enterprise-scale deployment more practical, with [distilled models in Amazon Bedrock being up to 500% faster](#).

## Multimodal Integration

The integration of text, image, audio, and video processing capabilities has become more sophisticated and practical. This advancement enables more natural interactions between AI systems and users while expanding potential applications across industries. [Amazon Nova's multimodal models](#), which launched in December 2024, represent a significant step towards more cost-efficient multimodal processing.

## Generative AI Building Blocks





# Focus Areas for 2025

The transformation of generative AI from experimental technology to enterprise cornerstone brings both opportunities and challenges. Organizations that understand these dynamics and prepare accordingly will be best positioned to leverage AI's capabilities for a competitive advantage while maintaining operational excellence and ethical standards. These focus areas are:

## 01. Cost and Performance Optimization

Large models can be powerful but expensive - at times more powerful than is needed. Organizations should balance the capabilities of large models with the practicality of specialized, task-specific implementations that can operate within reasonable computational and economic constraints. This requires careful attention to model selection, model distillation, deployment architecture, and ongoing optimization strategies. Industry benchmarks are helpful when evaluating models, but organizations should evaluate models and optimization options based on their business needs.

## 02. Operational Excellence

Deploying generative AI systems requires robust operational frameworks that ensure reliability, security, and governance. These frameworks must implement:

- Comprehensive monitoring and observability systems that track both technical performance and business outcomes
- Clear accountability structures that define roles and responsibilities across technical and business teams
- Effective risk management protocols that ensure responsible deployment

## 03. Human-AI Collaboration

The evolving capabilities of AI systems are redefining how humans interact with technology, necessitating new frameworks for intuitive, transparent, and adaptable collaboration. Success in 2025 will require thoughtful interaction design that ensures AI enhances—rather than complicates—user experiences. This means establishing clear protocols for decision-making authority, user control, and explainability, ensuring AI-driven systems are both trustworthy and seamlessly integrated into workflows.

Organizations must design interfaces that support cooperative problem-solving, where AI augments human expertise rather than operating in isolation. Investing in both technical capabilities and human-centered AI design will be essential to maximizing the value of these systems. Generative AI's potential extends beyond automation—it should be leveraged to create adaptive, user-driven experiences that optimize both customer and employee interactions.

The following chapters explore these themes in detail, providing concrete guidance for implementation and strategic planning in the evolving landscape of generative AI.

# Core Trends and Predictions for 2025

## 01. **Agentic** Architectures

According to [Gartner's framework](#), AI agents are defined as "goal-driven software entities that use AI techniques to complete tasks and achieve goals." The evolution from traditional AI models to agent-based systems enables organizations to automate complex workflows while maintaining appropriate human oversight and control.

Unlike the generative AI systems we've become accustomed to in recent years - those that simply respond to prompts - agent-based architectures enable AI systems to move beyond simple question-answering. They achieve this by breaking down complex problems into manageable components. They create execution plans, coordinate multiple specialized agents working in parallel when applicable, and adjust their approach based on intermediate results and changing conditions. Most

importantly, they should maintain appropriate human oversight to ensure optimal and responsible results.

Agentic architectures make AI workflows more modular and scalable by handling multi-step processes automatically. Without agents, generating an answer often requires multiple calls to different AI systems, which requires significant custom configuration and maintenance and, therefore, can be expensive. In an agent-based system, we move this responsibility—these steps—to the agent supervisor. Agents orchestrate these steps, reducing complexity and making it easier to adapt workflows over time. This approach also abstracts away process details so the caller doesn't need to manage individual steps or configurations, allowing for faster iteration and easier maintenance.

### The Bottom Line

Our experience shows that agentic architectures break complex AI applications into smaller, self-contained parts, making them easier to build, scale, and maintain. This is especially critical for complex generative AI applications, as it enables better software engineering practices which are essential for enterprise-grade, production-ready generative AI

applications. It also improves results by ensuring each agent operates within its isolated inference context rather than a shared one. The long-term benefits of using agents to architect generative AI applications are clear, even if the current marketing hype doesn't hold up for long.



## 02. Optimization for Cost, Performance, and Sustainability

As organizations scale their generative AI implementations, optimization becomes increasingly important for maintaining operational efficiency and controlling costs. Successful optimization strategies address multiple dimensions of system operations, from model architecture to resource utilization.

### Model Distillation & Specialization

Model distillation has emerged as a key strategy for optimizing generative AI deployments. This approach enables organizations to maintain high performance while significantly reducing computational requirements. It uses larger models as a basis to create smaller, purpose-specific, or domain-specific models. These smaller models can perform just as well as their bigger counterparts in their specific domains or purposes while using much fewer resources for inference.

The economics of model distillation presents a compelling case for optimization. Initial investments typically include infrastructure setup, specialized expertise, development time, and data preparation. With [Amazon Bedrock Model Distillation](#), you can train smaller models to mimic high-performance ones, making them up to five times faster and 75% cheaper for specific use cases.

### Inference Optimization

Efficient inference represents one of the most significant opportunities for performance improvement and cost reduction in generative AI systems. This topic brings up concepts like [Minimum Viable Tokens \(MVT\)](#), which is how much you can optimize both input and output tokens to reduce costs while maintaining output quality. This approach requires enhanced management of context windows, prompt engineering efficiency, and selective data loading patterns.

Our experience with multiple clients has taught us that carefully managing tokens and using reduced prompts typically results in a reduction in token usage, while maintaining response quality. This translates directly to cost savings on model usage, particularly in high-volume applications where even small optimizations can yield significant financial benefits.

Thoughtful design of generative AI applications can also uncover additional opportunities for improvement, such as request batching. [Amazon Bedrock's Batch Inference](#) feature offers a 50% price reduction compared to on-demand inference, which is especially useful for data ingestion and non-time-sensitive inference.

### Sustainability Considerations

The environmental impact of AI systems has become an important consideration, driving innovations in energy-efficient computing and sustainable operations. Recent developments in specialized AI hardware, such as [AWS Trainium2 chips](#), demonstrate significant improvements in energy efficiency. These advanced processors achieve a [30% improved price performance](#) over older models and a [29% reduction in energy consumption](#).

Techniques like Parameter-Efficient Fine-Tuning ([PEFT](#)) can also help organizations achieve performance comparable to fine-tuning by using fewer trainable parameters and, consequently, fewer computational resources and energy consumption.

Success in optimizing generative AI systems requires a balanced approach that considers cost, performance, and sustainability. Organizations that successfully implement comprehensive optimization strategies position themselves to scale their AI initiatives effectively while maintaining operational efficiency and environmental responsibility. This holistic approach to optimization enables sustainable growth while ensuring maximum value from AI investments.

### 03. Multimodal AI and Data Processing

The evolution of generative AI from single-modality systems to comprehensive platforms capable of processing multiple types of content represents a fundamental advancement in enterprise AI capabilities. Modern multimodal systems demonstrate sophisticated abilities in understanding and generating diverse content types, enabling more natural and comprehensive interactions between AI systems and users while introducing new requirements for data processing and system integration.

Modern AI systems can simultaneously process text, images, audio, and video, extracting meaning from the relationships between different modalities. Amazon Nova's multimodal models demonstrate the maturity of these capabilities through their ability to maintain context across modalities while generating coherent, multi-format responses. For example, when analyzing customer feedback, a multimodal system can combine text sentiment analysis with vocal tone assessment and facial expression recognition to provide a more accurate understanding of customer satisfaction.

The practical applications of multimodal AI span multiple industries and have the potential to improve operational efficiency and user experience significantly. In retail environments, multimodal

systems enable product discovery experiences by combining visual search capabilities with natural language understanding. Healthcare organizations leverage multimodal systems for enhanced diagnostic support, combining imaging analysis with patient history and symptom descriptions. Media and entertainment companies implement multimodal AI for content analysis and generation, enabling automated content tagging, personalized recommendations, and interactive experiences.

Organizations implementing multimodal AI systems must address several technical challenges to ensure successful deployment. Model families like [Cohere](#), [Amazon Titan](#), and [Llama](#) are excellent for data processing optimization, though format-specific requirements need careful attention. Successful implementations typically employ unified data pipelines that use preprocessing steps to handle specific formats and join the data in a way that maintains semantic relationships across modalities.

Integration complexity presents another significant challenge, particularly in enterprise environments with existing systems and workflows. Organizations must implement standardized interfaces and robust synchronization mechanisms to ensure reliable operation across different modalities. Success requires careful attention to metadata management and error-handling strategies that maintain system reliability while enabling efficient processing of diverse content types.



## 04. Evolving Search and Discovery Technologies

The landscape of search and discovery in generative AI has evolved significantly, moving beyond simple vector embeddings to innovative systems that combine multiple approaches for enhanced accuracy and relevance. This evolution represents a fundamental shift in how AI applications process and retrieve information, enabling more nuanced understanding and more accurate responses.

### Advanced Retrieval Architectures

Modern search and discovery systems leverage multiple complementary techniques to optimize information retrieval and enhance the quality of AI-generated responses. Our experience has shown that hybrid searches, the integration of semantic vector search with traditional keyword matching and structural metadata analysis, enable a more comprehensive understanding of content and context.

The adoption of GraphRAG (Graph-enhanced Retrieval Augmented Generation) represents an advancement in search capabilities. This approach integrates traditional RAG with knowledge graph structures, enabling systems to understand and leverage complex relationships between pieces of information beyond what vector embeddings can represent. The knowledge graph component maintains detailed mappings of entity relationships and contextual connections, enabling more accurate information retrieval.

Another significant development is agentic RAG. Traditional RAG implementations use a single source of additional information, or at most a handful of sources. These sources are typically all queried for every request, with the most complex implementations utilizing simple logic rules to select sources. In contrast, agentic RAG lets each agent decide which sources to query among the multiple

sources available to that specific agent. With each agent having a particular goal, RAG sources become purpose-specific and can be configured independently for each agent that requires them or would find them relevant.

Caylent's solution, [OmniLake](#), extends these concepts further by implementing a chain-based execution model that enables dynamic, condition-driven information retrieval and processing. Built on AWS, OmniLake accelerates time-to-value by unifying information access, ensuring context preservation through source annotation, and enabling complex logic through validated processing chains.

OmniLake's serverless architecture allows for parallel execution of multiple retrieval and processing steps, with each step initiating automatically as its prerequisites are met. The system combines vector search, knowledge graph processing, and conditional logic into unified request chains, enabling sophisticated multi-stage retrievals that can adapt based on intermediate results. This approach particularly shines in enterprise contexts where information must be gathered and processed from diverse sources such as CRM systems, wikis, and document repositories, with each retrieval step potentially influencing subsequent queries and processing steps.

OmniLake's architecture also maintains comprehensive data lineage throughout the retrieval and generation process. Each piece of AI-generated content is tracked along with its source materials, creating an auditable chain of information flow. This capability is particularly valuable in enterprise contexts where understanding the provenance of AI-generated content is crucial for compliance, verification, and building trust in the system's outputs.

## 04. Evolving Search and Discovery Technologies (con't)

### Automated Tuning and Optimization

Advanced search systems incorporate refined tuning mechanisms that continuously optimize performance. Based on observed performance patterns, these systems dynamically adjust chunk sizes, refine retrieval thresholds, and optimize query formulation. Implementing automatic threshold adjustment and result ranking refinement enables systems to maintain optimal performance as content and usage patterns evolve.

Performance monitoring involves tracking multiple dimensions of system performance, including retrieval precision, response relevance, and context maintenance. This holistic approach to performance

measurement enables continuous optimization while ensuring system reliability and accuracy.

### Future Developments

The field continues to evolve, with several emerging trends shaping future capabilities. Enhanced contextual understanding and improved retrieval algorithms enable more sophisticated applications while advancing performance optimization techniques for better efficiency. Organizations should prepare for the continued evolution of evaluation frameworks and quality assessment approaches while maintaining a focus on practical implementation requirements.

---

## 05. Products Over Models

The beginning of 2025 has been dominated by the release of [DeepSeek R1](#), an open-source model that can achieve results similar to those of the top models with much cheaper inference. While this is generally considered fantastic news for those building generative AI applications, these advances made so far, along with the increasingly advanced models released in 2024 and expected throughout 2025, are individually too small to provide significant business value. It is undeniable that models are progressively becoming more intelligent and cheaper, but at this point, no single model and no single change will be a significant disruption.

A good example of this is the other major news from early 2025: [ChatGPT's deep research](#) feature. It was released shortly after their own new model, OpenAI o3-mini, which is impressive in itself and also marks the start of a new family of models. However, the deep research feature had a greater impact than the new o3-mini model because it introduced a truly new

capability—a way to perform web searches across over 30 sources—rather than simply offering similar functionality to other models at a lower price point.

The emergence of so many models offering decent performance has turned models into a commodity. The rankings for best performance and price are constantly changing. While models will continue to evolve through cumulative improvements in the long run, trying to keep up with every new release has become a futile and unproductive effort.

Rather than chasing every new release, we recommend focusing on the bigger picture—how these models deliver real business value. Organizations should prioritize the application layer, comprising the software and interfaces that enable users and system-based workflows to interact with model outputs. This approach allows businesses to leverage generative AI solutions that adapt to evolving technologies, helping them sustain a competitive edge.



## 06. Democratized AI and Federated Learning

While decentralized approaches to AI development and deployment continue to evolve, the practical implementation of these technologies remains firmly grounded in centralized cloud infrastructure. Understanding the appropriate balance between centralized and decentralized approaches has become a key point for organizations planning their AI strategy, particularly as they address specific requirements for privacy, latency, and regulatory compliance.

### Current Implementation Landscape

The reality of AI deployment shows a clear pattern of specialization between cloud and edge computing. Most model training operations occur in cloud environments, driven by the specialized hardware requirements and economic advantages of resource pooling. Cloud-based training enables organizations to leverage robust infrastructure while maintaining cost efficiency and operational flexibility.

Similar patterns emerge in inference operations, where most workloads continue to run in cloud environments. This centralization enables organizations to maintain consistent performance while efficiently managing resources and ensuring appropriate governance controls.

Edge computing is the preferred option for inference in specific scenarios where privacy concerns or network constraints demand local processing. The most prevalent concerns around the use of data and privacy occur in B2C contexts, especially for free AI services offered by companies that train their own AI models. Network constraints, on the other hand, may occur in consumer-owned devices or manufacturing environments.

### Privacy and Regulatory Compliance

Privacy and regulatory requirements often drive decisions about AI deployment architecture. Organizations operating in regulated industries must implement sophisticated approaches to data protection and compliance management. Successful implementations consider the regionality of data and incorporate comprehensive audit capabilities and detailed monitoring systems that ensure appropriate handling of sensitive information while maintaining operational efficiency.

Implementing privacy-preserving techniques requires careful attention to technical architecture and operational procedures. Organizations typically achieve compliance objectives through a combination of data localization strategies, anonymization of sensitive data, and highly secure encryption mechanisms. These implementations demand regular security assessments and comprehensive audit trails to ensure continued compliance with evolving regulatory requirements.



## 07. AI Governance & LLMOps: Operationalizing Generative AI at Scale

As AI systems become more advanced through model capability and agentic systems, robust governance and technology management for AI and ML operations will become paramount to ensuring the transparent, responsible, and seamless use of this technology.

First, let's start with governance: The [PRISM methodology](#) provides a comprehensive approach to governing AI and agent-based systems through five key dimensions: Principles, Responsibility, Intelligence, Security, and Monitoring. This framework establishes clear guidelines for agent behavior and decision-making, including explicitly defining permitted actions and escalation protocols. The responsibility dimension creates structured accountability through clear ownership of operations and defined decision authority, while intelligence guidelines ensure the appropriate use of AI capabilities through systematic monitoring and optimization.

Organizations start with PRISM by aligning AI governance principles with their company's mission, values, and regulatory requirements. This alignment ensures AI initiatives drive business value while maintaining trust and compliance. For early-stage implementations, this often means focusing on lightweight governance practices—for example, a single-agent AI assistant with clear rules for data access, response guidelines, automated response observability systems, and human review processes in critical decision points. Teams define basic principles for responsible AI use and establish a simple audit trail to track key AI-driven decisions.

As AI adoption scales, governance evolves to support multi-agent systems, cross-functional AI deployments, and real-time monitoring. Enterprises integrating agentic AI across multiple workflows require automated policy enforcement,

granular decision accountability, and real-time compliance checks. At Caylent, we enhance PRISM by incorporating LLMOps best practices and embedding continuous evaluations and techniques where AI systems and agents are monitored for accuracy, fairness, and drift. Our governance teams are alerted to any anomalies, ensuring they are addressed before they impact business operations.

### LLMOps: The Evolution of Operational Practices

Operationalizing large language models presents unique challenges that extend beyond traditional MLOps practices. While LLMOps builds upon established MLOps principles, it introduces additional complexities specific to generative AI systems. Thus, it requires sophisticated deployment, monitoring, and maintenance approaches that ensure reliability, performance, and appropriate governance.

Artifact management in LLMOps extends beyond traditional model versioning to include prompt engineering and prompt management, configuration management, and output validation via test data and evaluations. Organizations must maintain detailed version control for all system components while appropriately tracking dependencies and interactions.

Modern LLMOps requires observability systems that provide detailed insight into system operations and performance. These systems track metrics like response latency, token usage, error rates, and resource utilization across all system components. Logging and tracing systems enable detailed tracking of system interactions and performance patterns.



## 07. AI Governance & LLMOps: Operationalizing Generative AI at Scale (con't)

### Quality Assessment Framework

The non-deterministic responses of generative AI applications require a new approach to functional testing. Evaluation frameworks consider multiple metrics to evaluate the quality of the responses. These frameworks can be implemented in Continuous Integration / Continuous Deployment pipelines as part of existing automated test suites or used as a feedback loop during inference to evaluate and improve the quality of every answer. These metrics can be tracked over time to analyze system performance and used as an objective measurement of quality for prompt engineering and parameter fine-tuning.

Our experience building dozens of production-grade generative AI applications has shown us that

LLMOps is as vital to enterprise-grade generative AI as DevOps and platform engineering are to traditional software applications. Moreover, [effective quality assessments](#) via tests and evaluations can be built into integration pipelines or offered internally through a developer platform that provides the scaffolding needed to enable software engineering and LLMOps best practices.

By combining AI governance and quality assessment frameworks with an AI-first automation strategy, businesses of all sizes can confidently scale systems without introducing unnecessary complexity. This approach ensures governance practices enhance AI performance, trust, and long-term reliability without hindering innovation.





## 08. Generative AI UI/UX

Most generative AI implementations in user-facing software have been centered around chat-based interfaces. While chatbots have become increasingly sophisticated through RAG techniques and agentic architectures, these AI systems have evolved beyond simple conversational tools. Today, they can answer questions about company-specific knowledge, automate repetitive tasks, and even execute actions on behalf of users. However, while AI-powered chat interfaces have seen significant growth, the way users interact with applications remains largely unchanged. Despite companies collecting vast amounts of user interaction data, little innovation has been made in how generative AI can redefine how users interact with applications.

At Caylent, we are already pioneering this shift by exploring how generative AI can power dynamic UI generation, adapting interfaces to users in real-time based on behavior, context, and preferences. Rather than simply personalizing content, AI-driven interfaces can fundamentally reshape the user experience by dynamically tailoring workflows, layouts, and visual elements. While there is still significant work to be done in this domain, 2025 is poised to be a turning point. The industry is shifting

towards deploying AI-generated user interfaces into production environments, enabling companies to provide truly unique, personalized digital experiences at scale.

The evolution of generative AI from a tool to an active participant in business processes necessitates advanced frameworks for human-AI collaboration. Successful Human-in-the-Loop implementations carefully balance automation capabilities with human expertise and oversight, ensuring AI serves as a collaborative partner rather than a black-box decision-maker.

However, successful adoption is not just a technical challenge but an interaction design challenge. Organizations must ensure that AI systems are not only technically sound but also intuitive, explainable, and aligned with human cognitive models. Poorly designed AI interfaces can lead to confusion, misuse, or rejection by end-users.

To address these challenges, AI collaboration must be designed around clear interaction paradigms that define the level of human control versus AI autonomy.

### Human-AI Collaboration Levels

Collaboration Mode	Human Role	AI Role	Example Use Case
Automation	Oversight & validation	Executes predefined tasks	AI-driven log analysis
Augmentation	Decision support	Provides insights, suggestions	AI-assisted data visualization
Co-Creation	Active participant	Generates & iterates with humans	AI-powered content generation
Autonomy	Approval only	Makes decisions independently	AI-driven fraud detection

By applying AI interaction design patterns aligned with these collaboration modes, organizations can reduce friction, enhance user trust, and accelerate adoption. These patterns include:

### Guided AI Assistance

AI provides structured recommendations, while humans retain control over decisions.

### AI as a Second Opinion

AI suggests alternative actions with transparent reasoning.

### Explain & Iterate

AI justifies decisions and allows for user-driven refinement.

## Integrating AI into Business Workflows

Beyond technical capabilities, AI adoption requires reimagining workflows to ensure that AI solutions augment human expertise rather than create additional complexity. This requires technical upskilling and UX-driven design principles that enhance the AI-user interaction experience.

Caylent's [Applied Intelligence](#) methodology accelerates this shift by embedding AI throughout cloud evolution, prioritizing seamless user interaction and trust-building mechanisms. By applying AI interaction patterns that match different business needs—ranging from automation to co-creation—organizations can scale AI adoption with confidence.

## Preparing for the Future of Human-AI Interaction

Enhanced capabilities for natural interaction, explainability, and contextual understanding continue to open new possibilities for human-AI collaboration. As AI becomes embedded in business processes, organizations must systematically address not only technical and operational readiness but also usability, transparency, and human trust.

A structured AI Maturity Model can help guide this transformation:

### AI as an Experiment

Limited pilots, no defined interaction models.

### AI as a Tool

AI assists predefined tasks but lacks user adaptation mechanisms.

### AI as a Co-Worker

AI is embedded into workflows with iterative human-AI interactions.

### AI as an Adaptive Partner

AI anticipates user needs and refines its behavior based on feedback.

By focusing on the design of AI interactions, organizations can ensure that AI becomes a natural and valuable part of decision-making processes rather than a disruptive force. AI that is understandable, usable, and trustworthy will enable businesses to fully unlock AI's potential while preserving human oversight and control.

## Understanding Customization and Scalability

# Should You Build or Buy Generative AI Solutions?

The decision to build custom generative AI solutions or leverage existing platforms is a strategic choice that impacts immediate implementation, long-term scalability, and total cost of ownership. Organizations must evaluate multiple dimensions of this decision to ensure alignment with business objectives and long-term operational efficiency.

Custom solutions often demonstrate significant advantages in specific scenarios, particularly when organizations require unique capabilities or face unusual constraints. Meanwhile, vendor solutions such as [Amazon Q for Business](#) are ideal if organizations require rapid deployment or standard functionality with limited customization. Typically, vendor solutions will fit well-known processes and use cases but don't offer the degree of customization needed for innovative use cases or unique business processes.

Initial assessment for implementing a generative AI solution must establish clear requirements for system functionality and performance while identifying specific constraints and dependencies. Organizations must also assess internal expertise and resource availability while considering long-term maintenance requirements and the total cost of ownership. This foundation enables informed decisions about implementation approaches while ensuring alignment with business objectives.

As technical capabilities continue to advance and operational requirements grow more complex, organizations must prepare not only for changing requirements but also for significant advancements in AI that could drive new requirements or enable entirely new use cases.



# Multi-Agent Systems in Enterprise Workflows

Case Study

BrainBox AI's development of the world's first virtual building assistant, ARIA (Artificial Responsive Intelligent Assistant), demonstrates the practical power of multi-agent architectures in complex operational environments. Faced with the challenge of redefining how their customers manage HVAC systems across diverse building environments while using BrainBox AI's proprietary solution for energy optimization, Caylent implemented a multi-agent system that transformed how technicians and building managers control HVAC systems.

The technical implementation centered on an agentic architecture where a planning agent received the user's request and routed it to multiple specialized agents. Each agent specialized in areas like temperature control, equipment details, and energy optimization, while also using RAG techniques to process live data. The system's core architecture enabled real-time processing of sensor data while maintaining predictive modeling capabilities across all building systems.

## Implementation Lessons

With such a complex process for generating an answer, especially with the involvement of live data, latency became a problem during the development of ARIA. Advanced prompt engineering and caching techniques reduced it significantly. However, the most unexpected problem was with the delay in obtaining live data. Endpoints that were considered fast for populating dashboards were introducing significant delays in the response generation of ARIA, in part because of how many were being called, and in part because of the different expectations from users. Significant data analysis efforts were required to identify the endpoints

that introduced the highest latency, and after their performance was improved, the response times of ARIA saw a meaningful reduction.

Over the development of ARIA, we saw the release of several new LLM models from different providers. Initially we identified that supporting multiple models and easily changing the models used was not a strict requirement for the first launch, but that it would be necessary for the long-term success of ARIA. Because of that support, and the use of [Amazon Bedrock](#), replacing the models used with the newly released ones allowed the ARIA team to improve the performance and reduce the costs of ARIA with minimal effort.

ARIA empowers facilities managers to have more insight into on-site challenges and proactively report on areas that need attention. Using voice and text commands, the virtual assistant oversees building operations and performs specific tasks, using data to improve energy efficiency. ARIA, together with its AI for HVAC technology, resulted in:

UP TO 25% ↓ reduction in energy costs

UP TO 40% ↓ reduction in GHG emissions



Pipes.ai

# Agentic AI-Powered Voice Service

Case Study

[Pipes.ai](#) collaborated with Caylent to develop a GenAI-powered Voice AI service with advanced agentic AI capabilities to transform their customer interactions. This solution offers a conversational interface that allows clients to engage naturally and carry out complex tasks directly during calls. By integrating AWS AI services with text-to-speech (TTS) and speech-to-text (STT) technologies, this solution enables human-like conversations and

advanced capabilities like open-ended questioning, real-time appointment rescheduling, and contextual intelligence.

By leveraging Amazon Bedrock, Pipes.ai can continuously evaluate and adopt the best-performing models to stay ahead. With this ability, Pipes.ai can now move to Amazon's Nova models, leveraging their stronger value proposition and smooth integration with AWS's ecosystem.

## Implementation Lessons

Amazon Bedrock's support for drop in replacement of AI models played a key role in the improved performance of Pipes.ai's Voice AI service. The ability to leverage Amazon's Nova models as soon as they were made generally available provided Pipes.ai with the tools to deliver a better service at a lower cost.

With its agentic AI capabilities, this solution can support more advanced features such as real-time call rescheduling, asking open-ended questions, and clarifying unclear responses. The Voice AI solution has the potential to reduce call center costs by up to 70% and boost lead conversions by enabling more natural conversations that lead to better outcomes.

# Looking Forward

The future of generative AI depends heavily on establishing strong foundations for continued development and implementation. Organizations that maintain a practical focus on current capabilities while preparing systematically for enhanced functionality will be best positioned to leverage these powerful technologies effectively. This balanced approach enables meaningful progress while ensuring appropriate attention to both technical excellence and practical operational requirements.

## Key Factors for Success

As organizations move beyond initial generative AI implementations to comprehensive enterprise deployment, several key factors emerge as critical for success:

01 Focusing on models is a losing proposition. The real value is in building applications that use those models to solve real business problems.

02 The emergence of sophisticated agent-based architectures enables more effective approaches to complex business challenges.

03 Cost optimization remains a crucial element of any successful implementation.

04 Using thorough assessments and operational frameworks helps organizations maintain reliability and effectiveness at scale, providing a strong foundation for planning and allocating resources effectively.

05 Security and governance requirements demand attention throughout the implementation process.

06 Evaluating models and applications based on specific business requirements is more valuable than using only industry benchmarks.

07 The path forward requires careful attention to architecture design and operational integration. Organizations must ensure appropriate infrastructure and support systems while developing team capabilities and operational procedures, focusing on long-term success.

08 Having an [LLMOps Strategy](#) that enables continuous adaptation to the latest models and changes while maintaining enterprise-grade engineering best practices will ensure organizations can adapt to the rapidly changing GenAI field.

09 Implementing a systematic approach to process adaptation and change management with a clear focus on achieving measurable outcomes ensures sustainable growth.

10 Using continuous improvement methods that focus on measuring and optimizing performance boosts capabilities and ensures strong technical and business outcomes, driving operational excellence.





CAYLENT

# Let's Get Started

The transformation of generative AI from experimental technology to essential business capability brings both opportunities and challenges. By proactively addressing these dynamics, organizations can position themselves to fully leverage AI's transformative potential. Success requires careful attention to both immediate implementation needs and longer-term strategic objectives, enabled through systematic approaches to capability development and continuous improvement.

## Jumpstart Your Generative AI Initiatives Today

No matter where you are on your journey, Caylent can help you quickly achieve your AI goals, from reinventing customer experiences and creating innovative new applications to massively improving productivity.

Whether you are looking to define a use case, create a roadmap, build a prototype, or implement a production quality solution, Caylent has a suite of generative AI offers tailored specifically to accelerate your goals.

[Our GenAI Solutions →](#)