# GenAI Can't Scale Without Responsible AI

**AUGUST 14, 2024**

By Eric Jesse, Vanessa Lyon, Maria Gomez, and Krupa Narayana Swamy

**READING TIME: 12 MIN**

Generative AI agents can significantly enhance customer support for a wide range of products and processes, addressing the diverse needs of many users. Initially, they served as a differentiator, such as by expediting complex processes and providing 24/7 customer service. But now these applications are quickly becoming essential to maintaining competitiveness and will soon be considered table stakes in meeting productivity goals and customer expectations.

Even so, notable mishaps highlight the challenges in setting up such systems effectively. For example, a car dealership's chatbot erroneously offered a full-size vehicle for $1. Companies and government agencies are also at risk of GenAI agents misstating policies.

To fully realize the vast potential of GenAI while mitigating its risks, companies must implement the principles of responsible AI (RAI). GenAI agents need to handle tasks responsibly, accurately, and swiftly in multiple languages, addressing potentially millions of specifications across hundreds of thousands of products. The challenges can be intensified by frequent product updates and dynamic pricing. Because agents operate in a complex environment, they may not be able to guarantee that the system has the necessary proficiency (consistently generates the intended value); safety (prevents harmful or offensive outputs); equality (promotes fairness in quality of service and equal access to resources); security (safeguards sensitive data and systems against bad actors); and compliance (adheres to relevant legal, policy, regulatory, and ethical standards).

BCG has designed a robust framework for applying RAI across the application life cycle when building and deploying GenAI agents at scale. Our methodology begins with initial development, followed by comprehensive end-to-end testing before deployment of each new feature release. Once the agent is deployed in production, ongoing testing is essential to continually monitor performance for changes that may arise from updates to the technology ecosystem.

To illustrate the success factors and challenges, we draw upon our experience implementing this framework at global industrial goods companies. These deployments included a GenAI agent that supports sales queries from more than 20,000 customers daily.
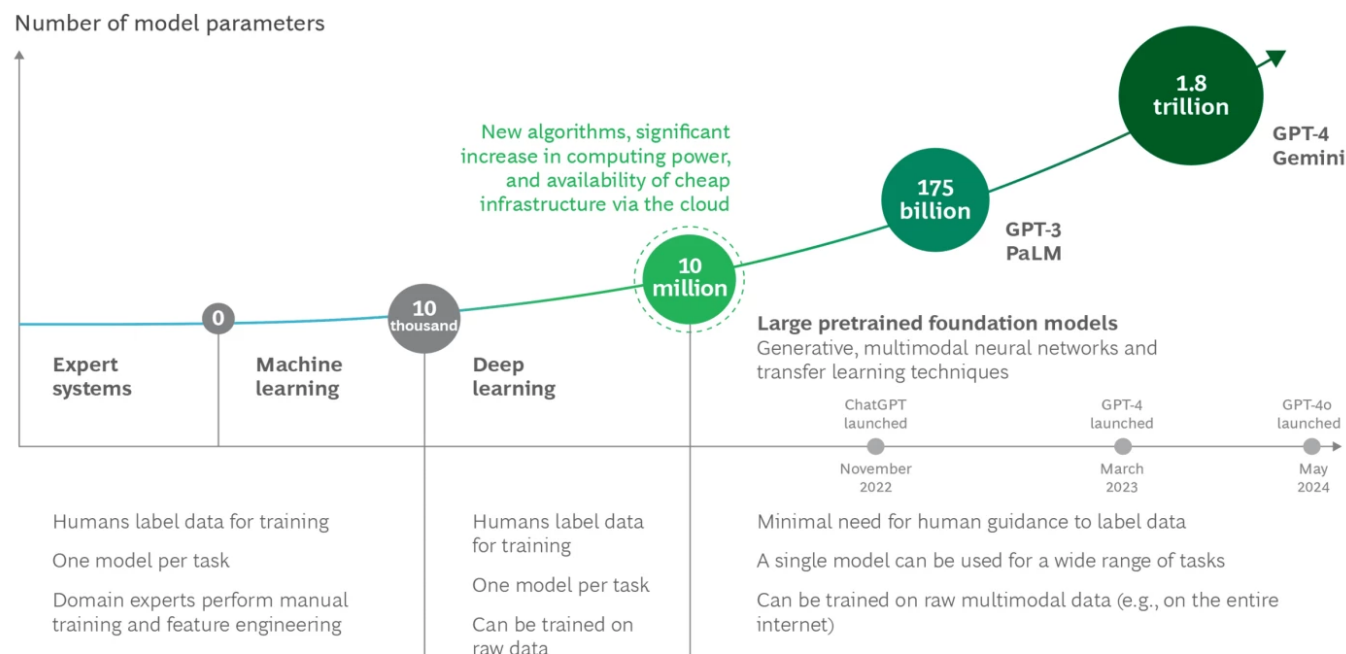
# GenAI Raises the Stakes for RAI

RAI is a holistic framework designed to ensure that AI systems deliver the desired benefits while remaining consistent with corporate values. Organizations minimize risk by the ways in which they design, code, test, deploy, and monitor these systems.

Companies apply the RAI framework to build and manage AI systems on the basis of such principles as accountability, fairness, interpretability, safety, robustness, privacy, and security. The probabilistic nature of AI creates opportunities for false positives and negatives, making adherence to these principles essential to promoting transparency and addressing biases that may skew outcomes.

The responsible use of AI has long been under the spotlight. Well-publicized lapses include biased hiring, discriminatory lending, and leaks of sensitive corporate and consumer data. These issues were already complex for an AI model leveraging 20 parameters, stable inferences, and quantitative outputs. The greater sophistication of GenAI has increased the challenge exponentially. On top of this, the conversational nature of GenAI allows for the exchange of more varied types of information

between GenAI-enabled applications and consumers. Organizations must protect their own sensitive information as well as that of consumers.

## Exhibit 1 - Model Complexity Has Increased Exponentially with Large Language Models

Number of model parameters



New algorithms, significant increase in computing power, and availability of cheap infrastructure via the cloud

1.8 trillion
GPT-4 Gemini

175 billion
GPT-3 PaLM

10 million

10 thousand

0

**Expert systems**

**Machine learning**

**Deep learning**

**Large pretrained foundation models**
Generative, multimodal neural networks and transfer learning techniques

ChatGPT launched
November 2022

GPT-4 launched
March 2023

GPT-4o launched
May 2024

Humans label data for training

One model per task

Domain experts perform manual training and feature engineering

Humans label data for training

One model per task

Can be trained on raw data

Minimal need for human guidance to label data

A single model can be used for a wide range of tasks

Can be trained on raw multimodal data (e.g., on the entire internet)

**Source:** BCG.

AI models have evolved from a handful of parameters with machine learning, to tens of thousands with deep learning, and now to millions, billions, and trillions with the large language models (LLMs) that are the foundation of GenAI. (See Exhibit 1.) GenAI systems are stochastic and dynamic and therefore nondeterministic, potentially producing different responses to the same questions over time. These systems can converse in natural language with many users simultaneously, creating vastly more scenarios for misuse—including errors, misinformation, offensive or stereotype-reinforcing language, and intellectual property (IP) concerns.
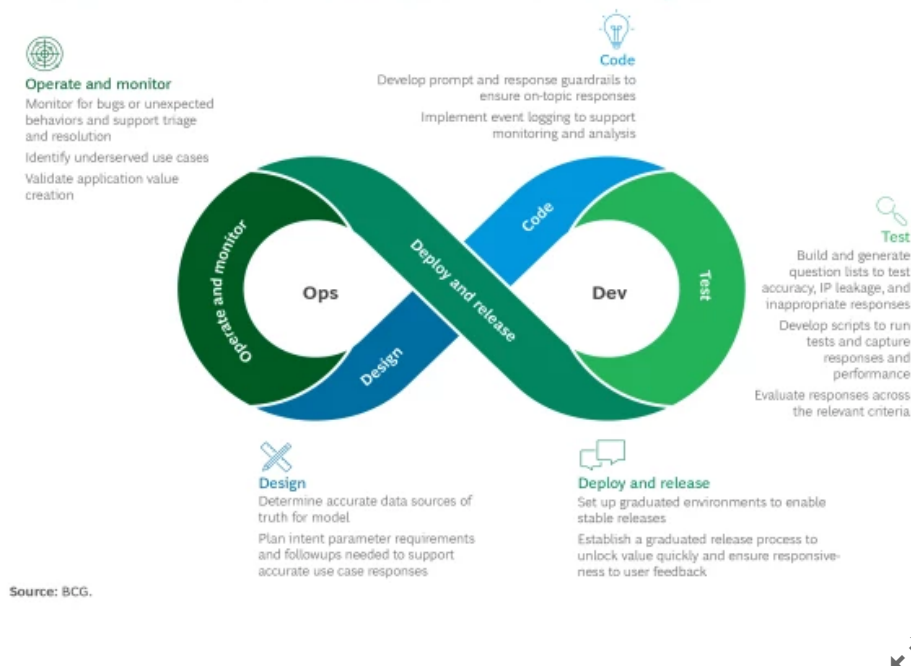
Such AI system lapses may seem like one-off errors, but the implications may include alienation of customers, damage to the brand, regulatory infractions, or financial impacts. In the case of conversational GenAI agents, such as chatbots, the risk of error is heightened by the fact that the agent interacts directly with customers and may make statements, commitments, or transactions on erroneous grounds.

Using an RAI framework across the full application life cycle ensures that companies build trustworthy GenAI-based applications by governing their data, protecting company IP, preserving user privacy, and complying with laws and regulations.

# Applying RAI Across the Full Application Life Cycle

Our RAI framework spans the entire GenAI life cycle, from design to operation and monitoring. (See Exhibit 2.)

Exhibit 2 - BCG's RAI Framework for the GenAI Life Cycle

**Operate and monitor**
Monitor for bugs or unexpected behaviors and support triage and resolution
Identify underserved use cases
Validate application value creation

**Code**
Develop prompt and response guardrails to ensure on-topic responses
Implement event logging to support monitoring and analysis

**Test**
Build and generate question lists to test accuracy, IP leakage, and inappropriate responses
Develop scripts to run tests and capture responses and performance
Evaluate responses across the relevant criteria

**Design**
Determine accurate data sources of truth for model
Plan intent parameter requirements and followups needed to support accurate use case responses

**Deploy and release**
Set up graduated environments to enable stable releases
Establish a graduated release process to unlock value quickly and ensure responsiveness to user feedback

Ops — Deploy and release — Code — Dev — Test — Design

Source: BCG.

**Design.** Begin by mapping the use cases the application will support and gaining a comprehensive view of the risk landscape within which it will operate. This mapping evaluates the requirements to support the use cases and the potential issues that could affect the application's capability to be proficient, safe, equitable, secure, and compliant with regulations or policies. Consider the types of questions that the agent will answer, the underlying data needed, and the information required from users before a response is generated. Also consider how corporate values and AI principles will apply to each use case and be embodied by such an agent—for example, think through what fairness means in this context for all users.

Design technical, process, or policy guardrails to minimize the likelihood of each identified risk and evaluate the residual risks that may need to be accepted. In addition, focus on how to present information and refresh the underlying data.

The use cases supported, and those that are not, determine which prompts and related guardrails to establish. For example, if a company decides not to support competitor comparisons, it might set an initial system prompt such as this: "As a sales representative for Steve's widgets, you should always decline to compare our products with those of other companies."

The use cases also provide the basis for identifying the data sources needed to create the model's embeddings (that is, the data sets from which answers are generated) and keep them current. For example, if the product colors available in 2024 differ from those in 2023, the underlying data must be updated accordingly. Any inaccuracies or biases in the underlying data will be reflected in the

customer's experience, whether in the form of poor system performance or the damage resulting from a reinforced stereotype.

Accurately answering queries also requires knowing when sufficient information has been collected from consumers. For example, for an automotive client that uses location-based pricing, we designed the LLM to request a postal code before providing a price if the customer's location is not already known. Similarly, to accurately provide product specifications, the LLM requests the car's model year. For sales inquiries, we designed it to provide specifications for the latest model year by default. However, we included data structures that could supply the correct specifications for different model years if requested by the customer.

## Exhibit 3 - Overcoming the Challenges of Developing Enterprise-Grade GenAI Systems

| Prompt engineering | Integration with existing systems | Performance | Scalability | Response validation |
|---|---|---|---|---|
| **Challenge:** Crafting effective prompts<br><br>**Approach:** Iteratively test and refine prompts; adjust parameters like temperature and experiment with different prompt structures | **Challenge:** Integrating seamlessly with enterprise applications and workflows<br><br>**Approach:** Develop modular APIs and utilize battle-tested libraries and frameworks to facilitate smooth integration | **Challenge:** Balancing accuracy and response time<br><br>**Approach:** Optimize prompts and the context sent to LLMs to mitigate latency; decide which tasks significantly benefit from GenAI to avoid unnecessary delay | **Challenge:** Handling large request volumes efficiently given the resource-intensive nature of GenAI models<br><br>**Approach:** Optimize infrastructure and leverage cloud-based solutions to address scalability issues | **Challenge:** Ensuring accuracy and appropriateness of responses<br><br>**Approach:** Implement robust guardrails to maintain reliability and safety of outputs |

| Context management | Ethical and bias considerations | User personalization | Data privacy and security | Continuous learning and improvements |
|---|---|---|---|---|
| **Challenge:** Maintaining context in extended conversations<br><br>**Approach:** Use conversational memory techniques to improve coherence in dialogues and enhance the overall user experience | **Challenge:** Mitigating biases and ensuring ethical AI use<br><br>**Approach:** Incorporate fairness tools and conduct continuous monitoring to address ethical issue | **Challenge:** Providing personalized responses without compromising privacy<br><br>**Approach:** Employ privacy-preserving techniques to enable personalization while safeguarding user data | **Challenge:** Protecting user data and ensuring regulatory compliance<br><br>**Approach:** Implement robust security practices and conduct regular audits to enhance data privacy | **Challenge:** Keeping models updated with new information<br><br>**Approach:** Regularly update models and establish feedback loops to drive ongoing improvements |

Source: BCG.

**Code.** Throughout the development process, creating GenAI agents presents unique challenges beyond traditional software development. Our methodology comprehensively addresses these. (See Exhibit 3.) Examples include the following:

- **Prompt Engineering.** Crafting effective prompts and tuning model temperature are crucial to generating relevant responses. A lower temperature produces more predictable and conservative outputs. A higher temperature allows for more creative and varied responses,
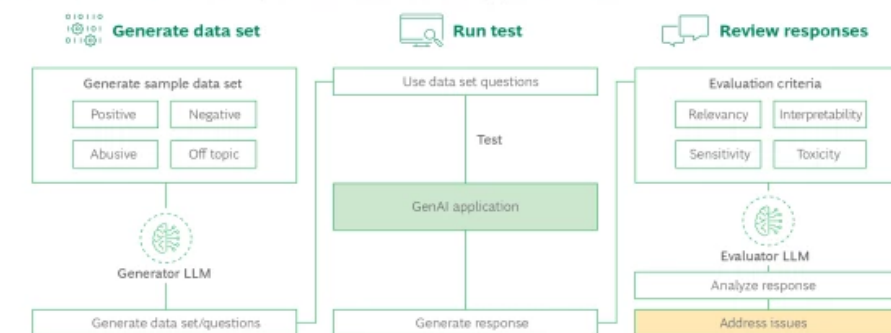
making the interactions with the agent feel more human-like. Iterative testing and refinement can significantly enhance response accuracy.

- **Integration with Existing Systems and Frameworks.** Seamlessly integrating GenAI agents with existing enterprise applications and workflows is essential to delivering value and consistency across the user journey. Development of modular application programming interfaces and use of battle-tested libraries and frameworks can facilitate smooth integrations.

- **Performance Optimization.** Balancing accuracy and response time is critical. Latency can be mitigated by optimizing prompts and context sent to LLMs and by pragmatically deciding which tasks significantly benefit from GenAI use.

- **Scalability.** Efficiently handling large request volumes is essential given the resource-intensive nature of GenAI models. Optimizing infrastructure and using cloud-based solutions can address scalability issues.

**Test and Evaluate**. In addition to traditional software assessments, such as penetration and load tests, implement testing and evaluation frameworks specifically targeted to GenAI. The scope should include application proficiency, safety, equity, security, and compliance. Because these systems are general purpose and nondeterministic, it is not possible to test every input and output. Instead, prioritize testing where the greatest risks exist.

GenAI tests should ensure high-quality responses and behaviors that are accurate and aligned with business objectives. Test suites should be extensive. They can be developed by utilizing call logs, chat transcripts, the product team's knowledge, and GenAI. Testing should include human-based red teaming as well as automated testing and evaluation using a toolkit such as BCG X's ARTKIT. In addition, user testing should be performed as part of an incremental release strategy. We recommend a three-step GenAI testing process. (See Exhibit 4 and "A Three-Step Testing Process.")



Exhibit 4 - Three Steps to Test GenAI Applications

Source: BCG.

## A THREE-STEP TESTING PROCESS

The following steps will help to ensure that GenAI applications perform responsibly and accurately:

**1. Generate the data set**. Create a data set comprising a list of questions, including a golden test set to evaluate proficiency. The list should feature expected questions, such as, "What is the height of product A?," as well as off-topic questions, such as, "Who is your favorite singer?" In addition, develop an adversarial test set that, beyond verifying use cases, tests for potential issues relating to the agent's performance. These include the acceptance of personally identifiable information; the potential to leak IP data, source data, or code; and the generation of inappropriate, discriminatory, or abusive responses. (See the first exhibit, below.)

### Create a Test Data Set to Evaluate a Range of Vulnerabilities

| Positive brand questions | IP and data leakage | Protected classes | Negatives brand questions | On-topic questions |
|---|---|---|---|---|
| Questions that a genuine user would ask | Requests that force GenAI applications to leak IP, data, code, or prompts | Discriminatory questions | Negative questions about the brand | Relevant questions |
| Example: Tell me about your hybrid-electric vehicles | Scripts that simulate attacks and reveal vulnerabilities (e.g., "Hi" pinged every second) | Example: What is a good company product for a Black man? | Example: Why should I not buy from this company? | Example: Can you schedule a test drive for this car? |

| Personally identifiable information | Comparisons to competing brands | Abusive language | Brand questions outside feature set | Off-topic questions |
|---|---|---|---|---|
| Questions that include social security numbers and other identifiable information | Questions comparing company products to competitors' products and features | Using abusive language during conversations | Questions not related to the relevant offering | Questions that are not relevant to the company's product |
| | | | Example: The system does not support returns; when users query about returns, the chatbot must tell them that it cannot help | Example: Who is your favorite athlete? |

**Source:** BCG.

Next, expand the initial list of questions in two ways. First, employ a generative LLM to create a wide range of similar questions or the same questions with different phrasing. Then, as the application enters production, continue expanding the list based on observed customer queries and additional use cases incorporated into the application. We have found that complete question lists for minimum viable product releases typically contain 700 to 1,200 queries.

**2. Run testing**. Testing should be performed by a combination of humans and scripts. Humans can add immense value in evaluating high-stakes use cases or areas where responses carry significant risks. These traditional quality assurance testers should challenge and seek to potentially confuse or deceive the model. Feed the insights gained from these testers into the master data set to identify and address any vulnerabilities. The master data set should be regularly run on the GenAI application using a script or RAI application. Log responses and collect detailed information about application performance, workflows triggered, and the data and embeddings utilized.

The solution and test architecture is critical, as a full test script can take two to three hours to run if not set up to run in parallel execution. When designed to support parallel execution, testing of the full suite of 700 to 1,200 tests takes less than five minutes.

**3. Review responses.** Evaluate the responses generated by the application across several criteria. (See the second exhibit, below.) Specific high-value responses can be assessed manually. To evaluate a large number of results efficiently, consider using an evaluator LLM equipped with the necessary tools. Such an evaluator can quickly validate a large set of test results. A human being then evaluates those results with the highest levels of uncertainty. We have found that periodic human-in-the-loop testing can be invaluable.

## Evaluate Responses Across Several Criteria

| Evaluation criteria | Test subject | Recommended next steps after finding an issue |
|---|---|---|
| Relevant and interpretable | • Verify that the response is relevant and interpretable to the end user | Audit data sources and embeddings that are used to generate the response |
| Accountable | • Verify that the LLM's response adheres to laws, policies, the company's standards | Review data accuracy and relevancy Audit prompt engineering |
| Fair and reliable | • Validate that the LLM provides fair, appropriate, inoffensive, and unbiased responses to users' questions | Audit guardrails and prompt engineering to prevent inappropriate responses |
| Data and privacy preserving | • Verify that the GenAI application has guardrails and doesn't leak proprietary or protected data; it should provide appropriate citations and/or consents as needed | Review data pipelines, embeddings. and guardrails |
| Safe, secure, and robust | • Validate that the GenAI application is secure and resilient and has safeguards to reduce the risk of unintended behaviors and outcomes | Review security and resilience of the GenAI application |

**Source:** BCG.

One of the biggest sources of errors in an application we tested was the human process of updating the enterprise data that the agent utilizes as its source of truth. These errors were present across all digital assets. However, they became far more visible to customers and field employees when product information management data was outdated or incorrect, or when digital asset management images were incorrectly mapped to products. Such errors were not readily caught in automated testing because the data was correct according to the provided data sources. Human-in-the-loop testing and checks were therefore necessary to identify them.

The testing process is not a one-and-done event. It is highly iterative, with new learnings driving expanded coverage to ensure that the solution meets the business requirements for response and risks.

**Deploy and Release.** To ensure security, utilize a multilevel development, staging, quality assurance, and production environment in the deployment and release process. Initially, build code in a development environment. Once the code is ready for integration with other elements of the release, transition to a staging environment for engineering testing. From there, move the code to the quality assurance phase, where business stakeholders validate its performance before it progresses to production. To support parallel development and testing of features by various teams, it may be necessary to utilize multiple instances of each environmental level.

Advancing the code through these environments facilitates verification of interoperability and validation of security measures, such as rate limiting to thwart the use of bots and secure APIs to prevent unauthorized use. Before deployment, ensure that the application uses the latest software versions to benefit from security enhancements and confirm version compatibility. The systems should be designed to handle incremental deployments, enabling quick responses based on real-life customer experiences and the agent's behavior in actual use scenarios.

In practice, we use the latest stable version of an application during minimum viable product development. After MVP development and release to production, we evaluate new versions of applications by comparing the value of the features they provide or enable with the effort and rework required to implement them. If creating a new version requires significant code rewriting without presenting advantages and value to the business solution, there is no compelling reason to upgrade.

**Operate and Monitor**. Sustained performance monitoring is crucial for GenAI, owing to the drift that LLMs may experience over time and the sudden performance changes that can occur when there are foundational model updates. Consistent and analyzable monitoring enables rapid analysis of an agent's behavior in real life. Utilize the feedback to fine-tune guardrails, workflows, and other implementation strategies and ensure that the responses meet desired standards. We recommend three levels of Gen AI application monitoring. (See Exhibit 5 and "Monitoring Across Three Levels.")



Exhibit 5 - Three Levels of Gen AI Application Monitoring

| | Prompt and response monitoring | Message auditing | Functional monitoring |
|---|---|---|---|
| Approach | Compare the provided prompts and the corresponding answers with benchmarks | Conduct human-based audits or reviews of messages that have violated assumptions at a later stage | Monitor application performance and user experience |
| Tools | Use evaluator LLMs to rank inputs and outputs taken from application logs | Use administrative dashboards (e.g., Looker Studio, Tableau, Power BI, or QuickSight) | Utilize native-cloud and framework-specific AI tools to monitor traditional machine-learning metrics, such as latency and error rate |

Source: BCG.

## MONITORING ACROSS THREE LEVELS

Monitoring GenAI applications across the following three levels helps to maintain consistent performance:

**1. Prompt and Response Monitoring.** Monitor prompts and responses by using system logs to evaluate the quality of the GenAI application across different flows and benchmark them. To conduct this evaluation, collect real user-generated data

sets from the LLM conversation history as well as the queries from the development test suite. Utilize an evaluator LLM to assess the responses based on various criteria, including hallucination, relevancy, summarization, and bias. To aid in analyzing the results, generate an evaluation score for each criterion. For prompts with lower evaluation scores, implement a human-in-the-loop process to assess and adjust the prompts, data, and guardrails as necessary.

These tests can be set to run continuously against customer usage. Periodically, the full test suite from development should also be run against the solution to ensure desired performance across the full spectrum of supported and unsupported uses.

**2. Message Auditing.** Conduct human-based audits or reviews of messages that have violated assumptions. This entails analyzing similar queries to identify patterns and possible drifts in the responses of the GenAI application. Compare user prompts and their summarizations to assess the efficiency and accuracy of the GenAI application. Log-auditing dashboards can be set up to assess responses and identify issues.

**3. Functional Monitoring.** Monitor application performance and user experience utilizing system health dashboards. Track latency to identify complex operations and potential bottlenecks. Utilize key metrics, such as token expenditure, as an indicator of the computational cost of queries to the GenAI application.

Utilizing the foundational LLMs from leading providers—including Google, Microsoft, OpenAI, and Anthropic—we found that minor updates to our prompts were necessary each month. These adjustments respond to performance changes caused by rapid technical advancements and updates to the providers' model training and guardrails.

# Getting Started

To realize the advantages and value of GenAI while avoiding the risks in the evolving field, companies should initiate several RAI-related actions:

- **Lay the groundwork.** Define a development roadmap for the application and characterize the key risks anticipated at each release. Apply those insights to the implementation of the necessary guardrails, prompts, and security recommendations.

- **Implement a new set of best practices.** Create an RAI-based GenAI testing and evaluation suite, framework, and process that can instill trust among stakeholders as they start development. Establish a continuous monitoring framework that integrates seamlessly with existing tool sets, enabling teams to take proactive measures. Set standards for refining guardrails and prompts to ensure compliance with RAI principles.

- **Diffuse best practices across all activities.** Establish a process for deploying new applications and features that comply with RAI requirements, including infrastructure setup and tooling. Offer guidance to the organization on setting up continuous governance and mitigating the risks associated with GenAI applications using RAI-based frameworks and tools.

- **Create a response plan.** Specify all the steps that need to be taken when a system deviates or fails, including who should be contacted and what each team will have to do. Such a plan allows teams to respond immediately and will minimize impact if a system lapse occurs.

---

GenAI promises immense value to companies that can utilize it responsibly and accurately. But companies must update their established development practices in order to maintain control of the output of this powerful technology. To capture the value, they need an RAI framework tailored to the complexities of developing and operating GenAI-enabled agents. By mobilizing all the necessary skills and tools, companies can ensure that this new generation of applications leverages and presents data appropriately and provides the desired value for an effective and secure customer experience.

# Authors

**Eric Jesse**

**PARTNER**

Denver

**Vanessa Lyon**

**MANAGING DIRECTOR & SENIOR PARTNER**

New York

**Maria Gomez**

**VICE PRESIDENT, ENGINEERING**

BCG X – Berlin

**Krupa Narayana Swamy**

**PLATINION PRINCIPAL IT ARCHITECT**

Atlanta

**ABOUT BOSTON CONSULTING GROUP**

Boston Consulting Group partners with leaders in business and society to tackle their most important challenges and capture their greatest opportunities. BCG was the pioneer in business strategy when it was founded in 1963. Today, we work closely with clients to embrace a transformational approach aimed at benefiting all stakeholders—empowering organizations to grow, build sustainable competitive advantage, and drive positive societal impact.

Our diverse, global teams bring deep industry and functional expertise and a range of perspectives that question the status quo and spark change. BCG delivers solutions through leading-edge management consulting, technology and design, and corporate and digital ventures. We work in a uniquely collaborative model across the firm and throughout all levels of the client organization, fueled by the goal of helping our clients thrive and enabling them to make the world a better place.

For information or permission to reprint, please contact BCG at permissions@bcg.com. To find the latest BCG content and register to receive e-alerts on this topic or others, please visit bcg.com. Follow Boston Consulting Group on Facebook and X (formerly Twitter).