# Deloitte.



# Risk and Trust in the Age of Agentic Al



Trust has been a concern ever since artificial intelligence (AI) first operated. Deloitte established the Trustworthy AI Framework to identify and address the many aspects of that challenge, and the expanded capabilities of Generative AI placed a new emphasis on it.



Now, another evolution is prompting a fresh look at trust: Al agents are reasoning engines that can plan workflows, connect to external tools and data, and execute actions to achieve a defined goal. Instead of merely interacting with the user, they are designed to reason and act on behalf of the user.

What does this new capability, known as agentic AI, mean for trust? Our framework still holds: This new architecture involves many of the same issues, such as reliability, transparency, fairness, and others. But here, the dimensions of trust are more complex. Humans remain in the loop but manage it in different ways. With many operations difficult to track inside an AI agent's "black box," organizations don't only need to examine their AI plans. They need to reexamine the processes those agents may soon control.

The structure and function of Al agents has been examined in print before. In some ways, it represents a new arrangement of existing capabilities. This is more a change in trust than a change in technology--but that change in trust might be profound. An organization that implements Al agents should take careful measures to maintain confidence not only in the system's outputs, but in its methods.

# Familiar AI risks, only larger and more numerous

Instead of "only" generating answers based on human-fed prompts, agentic Al makes decisions and directs other processes autonomously. That puts humans in a different relationship with the step-by-step function of the system. Trust issues an organization may already have considered in the use of machine learning or Generative Al may reappear, in a familiar form but at a larger scale.

Al agents may never be fully autonomous, but the degree of autonomy they do represent means organizations likely need a fresh take on the human-machine relationship.

#### Points to consider:

If the output of one internal process becomes the input for another, it can be difficult to identify where errors or hallucinations originated. An error early in a sequence may become magnified by the end.

Governance doesn't only control processes; it is a process. Al can make it more effective. But an insufficient approach to trust and control has the potential to amplify error—or conceal it.

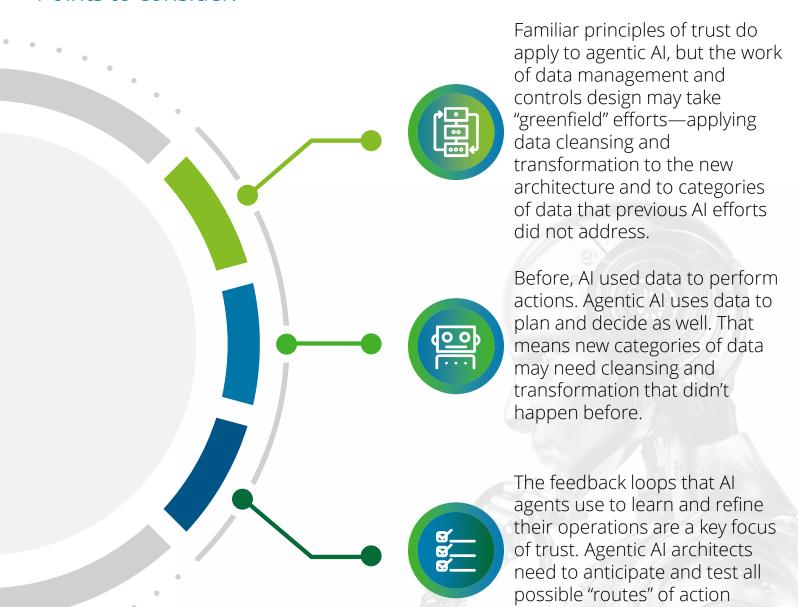
The learning and continuous adjustment inherent to Al makes the system non-deterministic. Each operation and its output is different from the last. That's one of Al's many valuable attributes. But risks multiply when it becomes harder to answer the question: Different why?

When humans have a smaller role in the loop, they have less need for some of the AI skills they have honed in recent years. But they will need new ones—not only to manage agentic AI, but also to take on the higher-level strategic jobs the new architecture frees them to address.

## The basis for trust is there

Humans built the processes you use today, and established the protocols that make them reliable. That means humans know where the critical points and guardrails are—and when it's time to implement agentic AI for those processes, it will be humans who design those systems and know where to focus the required controls.

### Points to consider:



within a system.

# Right now

The risks that may imperil trust in a future agentic AI system are likely to concentrate in the critical points of the processes as they exist today. For any operation your organization may assign to AI agents, catalog the hinges and decision points where control and monitoring will make the greater difference.

#### Points to consider:

Set rules for human oversight that correspond to the Al agent's new role. For example, if an Al agent is automating many of the tasks a junior financial analyst might perform, can the junior analyst's job evolve to include monitoring and analyzing the agent's outputs?

Be realistic and transparent in what you expect from agentic Al and the ways you use it. The fastest way to erode trust is to set an expectation and fail to meet it.

Consider an end-to-end evaluation of your enterprise Al policy.

To begin the work of establishing governance and controls for later agentic AI design and deployment, consider three broad steps:

Assessment: Inventory applicable use cases and determine the risks of each one.

Mitigation: Apply risk logic and devise controls and guardrails for each use case, including the use of simulations.

Monitoring: Continually track the model's performance and retrain it as needed. Don't only flag anomalies; optimize processes wherever possible.

# Standing on the shoulders of past understanding

Because of the way it works, agentic AI might not capture the public imagination the way Generative AI has done. But inside your organization, you and your stakeholders will see and feel the change. This is ultimately a realignment of existing AI tools into new process structures.

That means it will require a renewed application of trust awareness and safeguards that you already understand. A careful mapping of the systems agentic AI will run can help inform a re-examination of the critical points where monitoring and control are most important. From that starting point, organizations can look with confidence toward a new era of advanced automation.

Look out for our upcoming white paper on agentic Al's trust and risk implications, with practical examples and use cases for organizations to consider.



#### Contacts



Beena Ammanath

Managing Director

Global Deloitte Al Institute Leader

Deloitte Consulting LLP

bammanath@deloitte.com



# Acknowledgments

Amandeep Singh, Ashley Reichheld, Derek Snaidauf, Jim Rowan, Mike Crowthers, Prakul Sharma, Scott Holcomb



This article contains general information only and Deloitte is not, by means of this article, rendering accounting, business, financial, investment, legal, tax, or other professional advice or services. This article is not a substitute for such professional advice or services, nor should it be used as a basis for any decision or action that may affect your business. Before making any decision or taking any action that may affect your business, you should consult a qualified professional adviser. Deloitte shall not be responsible for any loss sustained by any person who relies on this article.

#### **About Deloitte**

Deloitte refers to one or more of Deloitte Touche Tohmatsu Limited, a UK private company limited by guarantee ("DTTL"), its network of member firms, and their related entities. DTTL and each of its member firms are legally separate and independent entities. DTTL (also referred to as "Deloitte Global") does not provide services to clients. In the United States, Deloitte refers to one or more of the US member firms of DTTL, their related entities that operate using the "Deloitte" name in the United States and their respective affiliates. Certain services may not be available to attest clients under the rules and regulations of public accounting. Please see www.deloitte.com/about to learn more about our global network of member firms.

Copyright © 2025 Deloitte Development LLC. All rights reserved.