

How to prevent AI-solution abuse

Mitigate the risks of AI-solution abuse thorough human-centered strategies to ensure its ethical and responsible use

Mitigating the risk of AI-solution abuse:

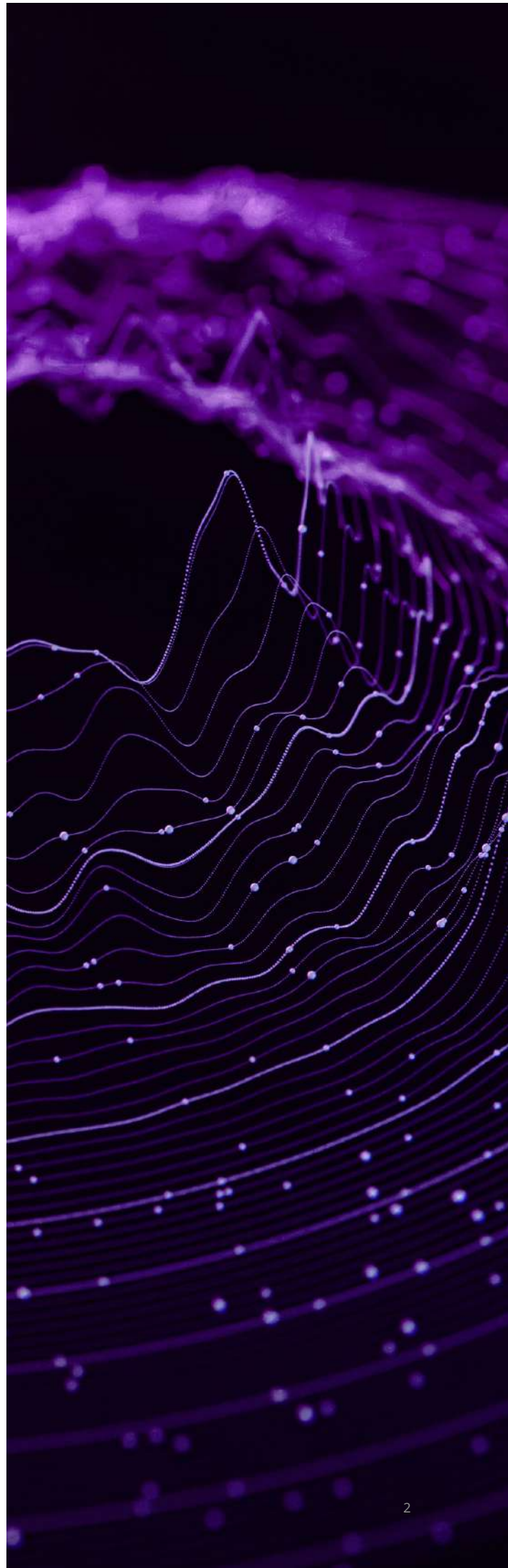
Human-centered strategies

In brief:

- In his latest blog, Jeroen Bet, Director of Strategy at Luxoft's Smashing Ideas, explores the growing risk of AI-solution abuse and highlights the importance of proactive measures in safeguarding ethical AI use
- The article presents four human-centered strategies: Creating bad actor personas, using the Six Thinking Hats technique, establishing guiding principles, and emphasizing continuous improvement and monitoring
- Jeroen underscores the need for organizations to anticipate potential abuse scenarios, foster a culture of accountability, and maintain an ongoing dialogue with users to ensure responsible AI adoption

Introduction

Artificial intelligence (AI) has revolutionized various industries, offering tremendous potential for innovation and transformation. However, as businesses adopt AI solutions, it is crucial to learn from past mistakes and take proactive measures to prevent abuse. This article offers four human-centered strategies that can effectively mitigate the risk of AI-solution abuse, ensuring ethical and responsible use of AI technologies.



The challenge of AI-solution abuse

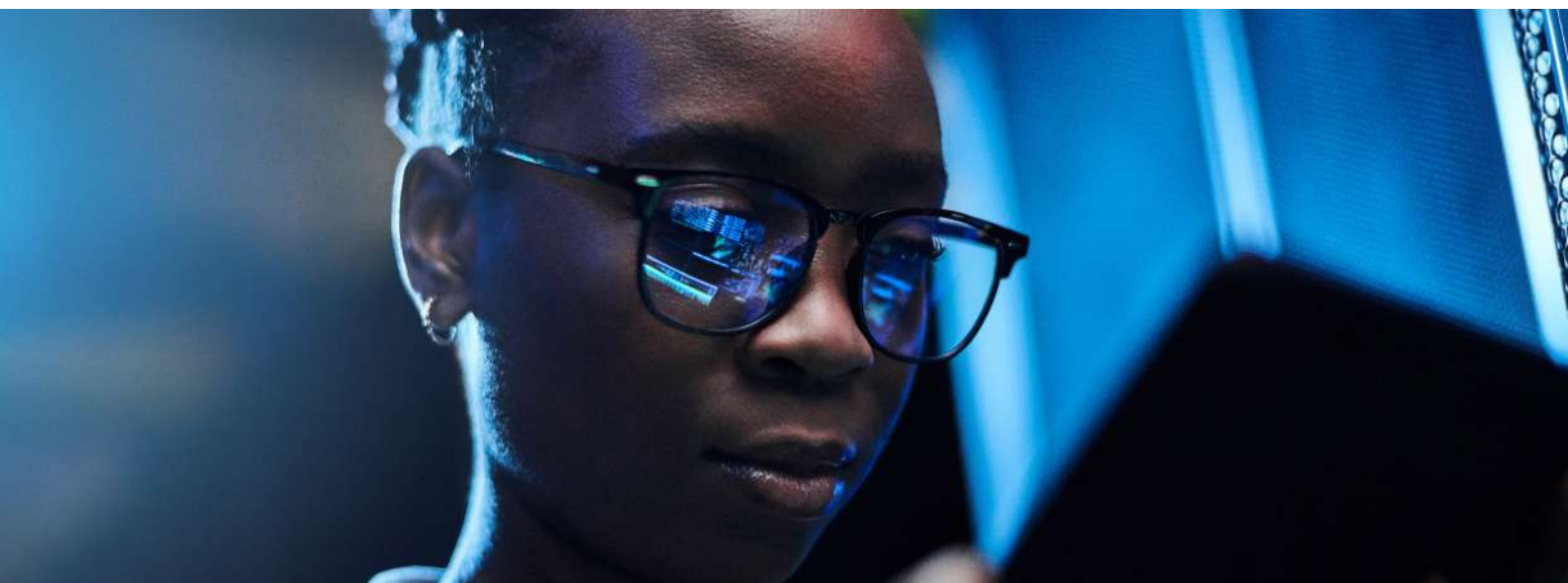
A couple of years ago, a client aimed to enhance diversity within their predominantly homogeneous workforce by implementing an AI engine. The concept was to utilize the AI engine for sorting and filtering resumes, generating merit-based hiring recommendations, and ultimately eliminating bias associated with traditional hiring methods. However, the AI model's training data set consisted solely of profiles from employees who had historically excelled within the company. Consequently, the training data set introduced inherent bias into the algorithm. As a direct consequence, the AI engine's recommendations unsurprisingly favored new hires who closely resembled the existing workforce, perpetuating the lack of diversity.

In the aforementioned scenario, mitigating unintentional bias in the system might have proven challenging without manipulation of the learning data set. From an AI perspective, the system itself did not commit any wrongdoing as it merely acted based on the defined success criteria of hiring candidates similar to the existing workforce. However, a comprehensive approach encompassing bias detection, auditing, and human evaluation of the data could have resulted in broader recommendations and helped steer the program away from biases. Although initial human intervention would be crucial, over time, the AI would progressively learn to act more responsibly and reduce bias tendencies.

AI-solution abuse can take various forms, ranging from unintentional biases in decision-making, as described above, to intentional manipulation for malicious purposes. These incidents highlight the need for strategies to address edge cases (extreme or unusual scenarios) and prevent abuse scenarios. Organizations must consider the potential risks associated with AI solutions and take steps to proactively mitigate them. The following are four human-centered strategies to get started:

1. Creating bad actor personas

To proactively identify potential abuse scenarios, organizations can create virtual bad actor personas. These personas represent users who intentionally misuse or exploit the AI solution. By understanding their motivations and behaviors, organizations can gain valuable insights that expose vulnerabilities in the system. This, in turn, enables them to develop specific edge cases that simulate scenarios and analyze how the personas interact with the system. Bad actor personas provide valuable insights into the deceptive practices, biases, security risks, privacy concerns, and system performance issues that could be exploited. This approach empowers teams to uncover potential areas of concern, anticipate abuse scenarios, and implement preventive measures.



About **the author**



Jeroen Bet

Director of Strategy,
Luxoft's Smashing Ideas

Jeroen is a strategy director with a solid background in customer experience. Over the course of his 25-year career, Jeroen has worked with companies such as Chempoint, Expedia, Costco, Amazon, and Microsoft. He has collaborated with a variety of stakeholders, including scientists, engineers, plant managers, marketers, conservationists, end-users, and business leaders. Jeroen has successfully led project teams in developing human experiences that incorporate AI and other cutting-edge technologies.

Ready to take a human-centered approach to your AI strategies?

At Luxoft, our strategy team can help you ask better questions to get to the right solutions.

Visit **luxoft.com** or **contact us** today to learn more about how we can help your organization innovate faster towards business objectives.

About Luxoft

Luxoft, a DXC Technology Company delivers digital advantage for software-defined organizations, leveraging domain knowledge and software engineering capabilities. We use our industry-specific expertise and extensive partnership network to engineer innovative products and services that generate value and shape the future of industries.

For more information, please visit **luxoft.com**