

# Responsible AI: From Bias and Privacy to Compliance and Risk Management

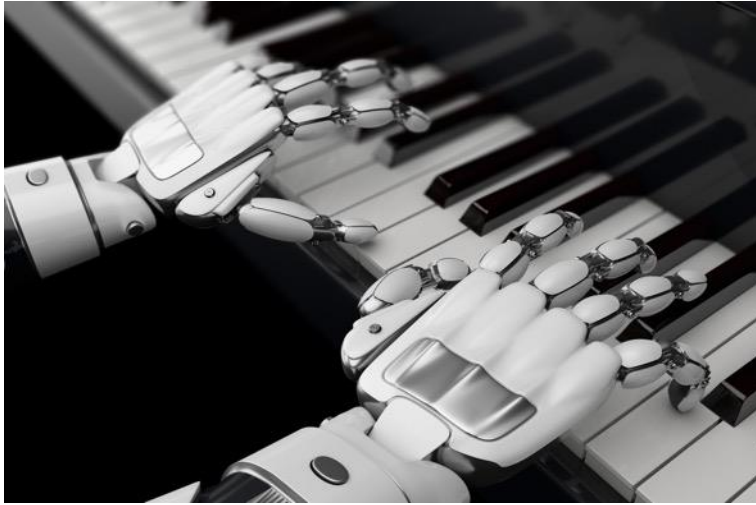
Anthony Mullen

© 2023 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner is a registered trademark of Gartner, Inc. and its affiliates. This publication may not be reproduced or distributed in any form without Gartner's prior written permission. It consists of the opinions of Gartner's research organization, which should not be construed as statements of fact. While the information contained in this publication has been obtained from sources believed to be reliable, Gartner disclaims all warranties as to the accuracy, completeness or adequacy of such information. Although Gartner research may address legal and financial issues, Gartner does not provide legal or investment advice and its research should not be construed or used as such. Your access and use of this publication are governed by [Gartner's Usage Policy](#). Gartner prides itself on its reputation for independence and objectivity. Its research is produced independently by its research organization without input or influence from any third party. For further information, see "[Guiding Principles on Independence and Objectivity](#)."

**Gartner**®



# Cool!



**Composing music**  
for commercials,  
movies, video games



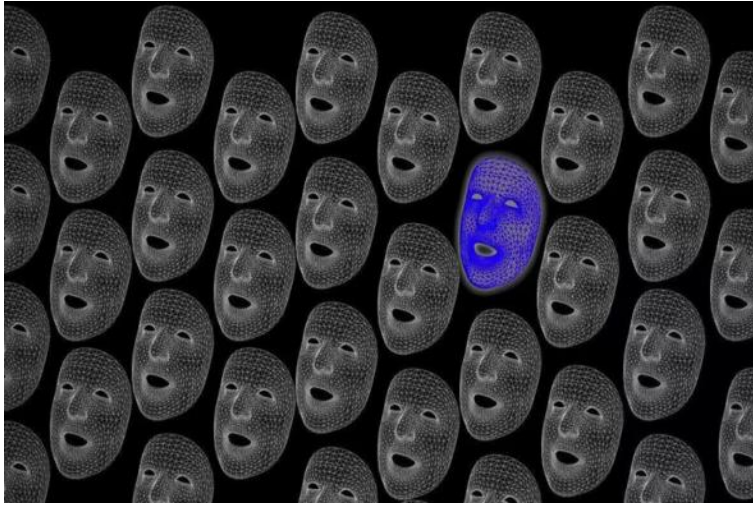
**Algorithms detect**  
anomalies in  
real time.



**Domestic robots**  
keeping people  
company



# Creepy!



**Facial recognition**  
leads to wrongful  
imprisonment



Children's **privacy**  
**violation** by  
smart AI toy



**Filter bubbles**  
empowered by AI  
on social media

**To reap the benefits of AI, we need to empower the cool  
and protect us from the creepy, by managing AI risks**

# From Responsible AI to AI Risk Management



# Key Issues

1. Eight AI risks and how to mitigate them
2. Emerging practices for AI risk management

# Key Issues

1. Eight AI risks and how to mitigate them
2. Emerging practices for AI risk management



# 1. Trust



## AI Risks

- Invalid (wrong)
- Unreliable (inconsistent)

## Key Mitigations

- Improve model explainability & transparency
- Guard rails
- Human-in-the-loop

### AI Camera Ruins Football Game By Mistaking Referee's Bald Head For Ball

Many complained that they missed their team's goals because the camera "kept thinking the Lino bald head was the ball."



Image tweeted by @seagull81



## 2. Privacy

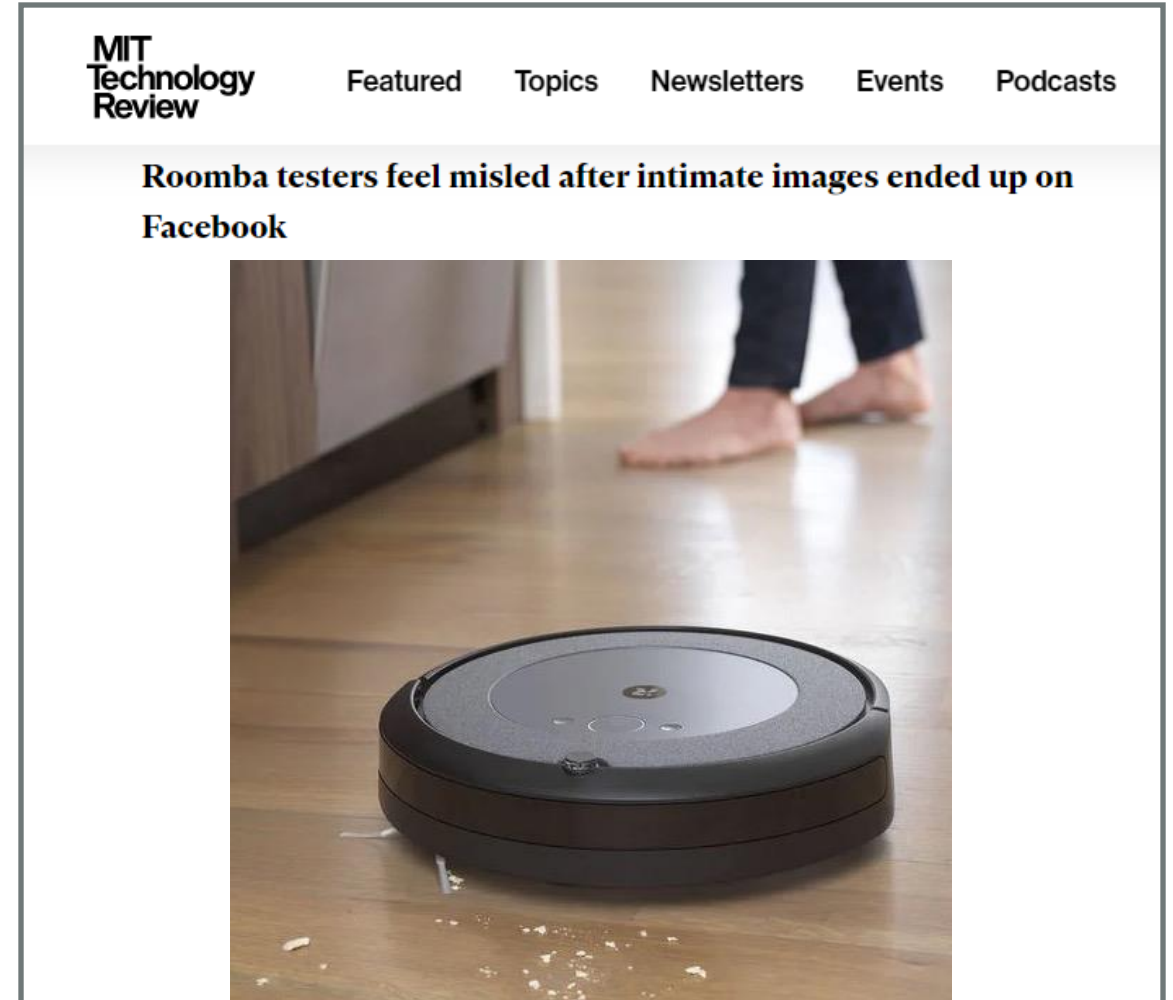


### AI Risks

- Unwarranted use of privacy sensitive data/personal identifiable information (PII)
- Unauthorized sharing of PII

### Key Mitigations

- Assess intended use of AI & use data proportionally
- Apply privacy enhancing technologies (PET)
- Keep data secure



Source: MIT TechnologyReview

# 3. Fairness



## AI Risks

- Bias against groups or classes of people, based on their color, ethnicity, gender, religion, age or any other classification protected by law.

## Key Mitigations

- Balance diversity in data (and in team)
- Use fairness tools/frameworks
- Include fairness in model monitoring

## A.I. Bias Caused 80% Of Black Mortgage Applicants To Be Denied

Kori Hale Contributor @

*I'm the CEO of CultureBanx, redefining business news for minorities.*

Follow

Source: [A.I. Bias Caused 80% of Black Mortgage Applicants to Be Denied](#), Forbes Media.

# 4. Security



## AI Risks

- Data poisoning
- Model theft
- Model deception

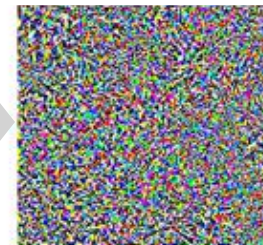
## Key Mitigations

- AI — cybersecurity collaboration
- Secure training data
- Harden AI models

### Manipulated Digital Images



Predict: 'tiger cat'



Perturbation



Predict: 'egyptian cat'

Add perturbations to digital image to fool AI model

Source: iProov; [DARTS: Deceiving Autonomous Cars With Toxic Signs](#), arXiv.

# 5. Safety



## AI Risks

- Unsafe outcomes, beyond intended or tested use
- Amplification of harmful output through automation

## Key Mitigations

- Monitor & alert for unintended consequences
- Prefilter outliers or untested cases
- Frequent drift monitoring

### Chess robot grabs and breaks finger of seven-year-old opponent

Moscow incident occurred because child 'violated' safety rules by taking turn too quickly, says official



# 6. Societal Impact

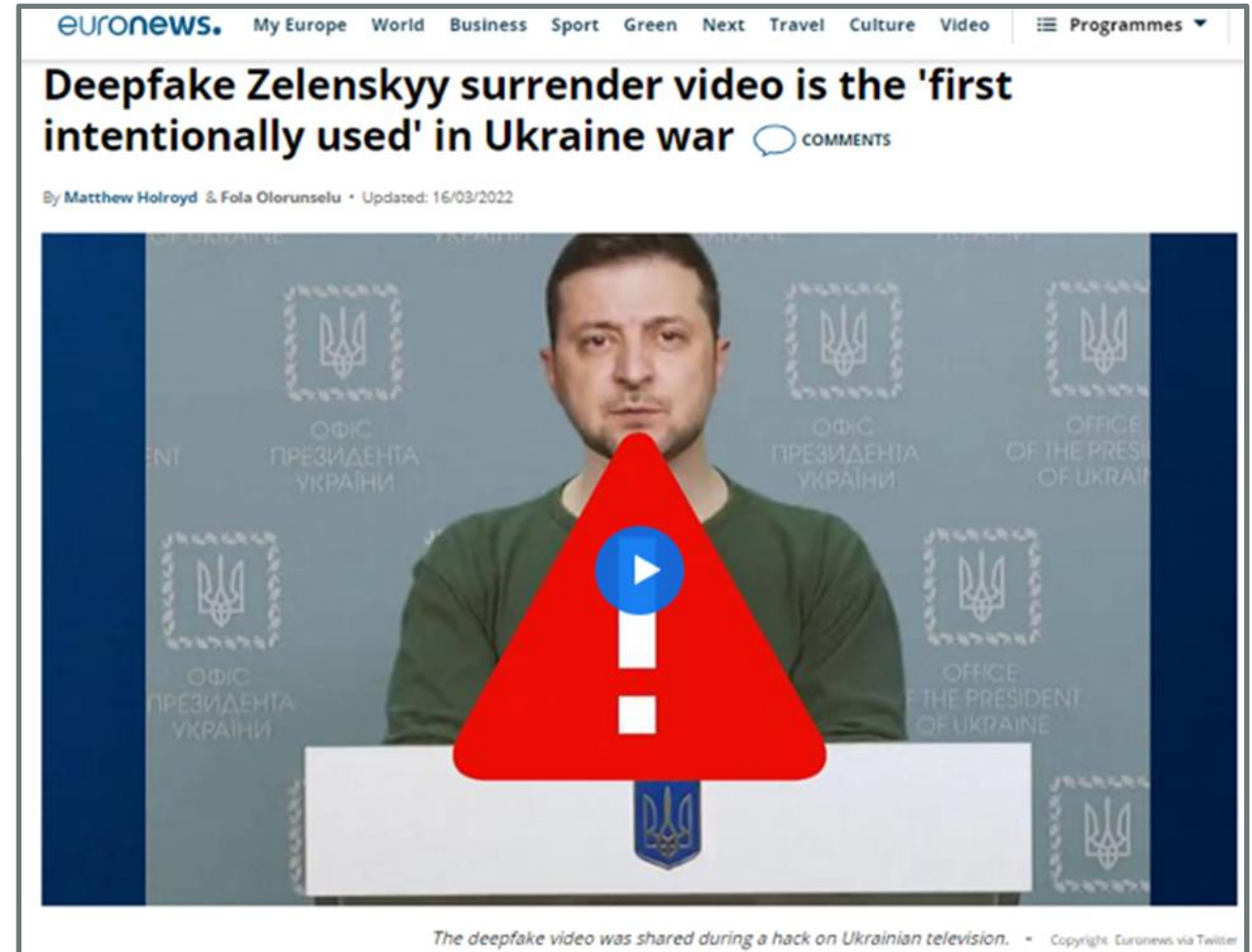


## AI Risks

- Negative impact on human well-being
- Misinformation and social bubble amplification

## Key Mitigations

- Apply human-centric design principles
- Verify content & detect fakes
- Limit content personalization



Source: [Deepfake Zelenskyy Surrender Video Is the "First Intentionally Used" in Ukraine War](#), Euronews.

# 7. Sustainability



## AI Risks

- Growing energy consumption and carbon footprint of AI

## Key Mitigations

- Apply emerging sustainable AI practices to reduce energy consumption
- Proactively apply AI to improve sustainability in business operations and society



Source: [AI Power Consumption Exploding](#), Semiconductor Engineering.



# 8. Accountability



## AI Risks

- Insufficient attention to preventing or mitigating AI failures
- Noncompliance with growing number of complex regulations

## Key Mitigations

- Create own set of principles and actionable guidelines (but do not reinvent the wheel)
- Partner with legal and compliance experts
- AI risks & ethics board

### France fines Google and Facebook €210m over user tracking

Data privacy watchdog says websites make it difficult for users to refuse cookies



### BLUEPRINT FOR AN AI BILL OF RIGHTS

MAKING AUTOMATED SYSTEMS WORK FOR THE AMERICAN PEOPLE

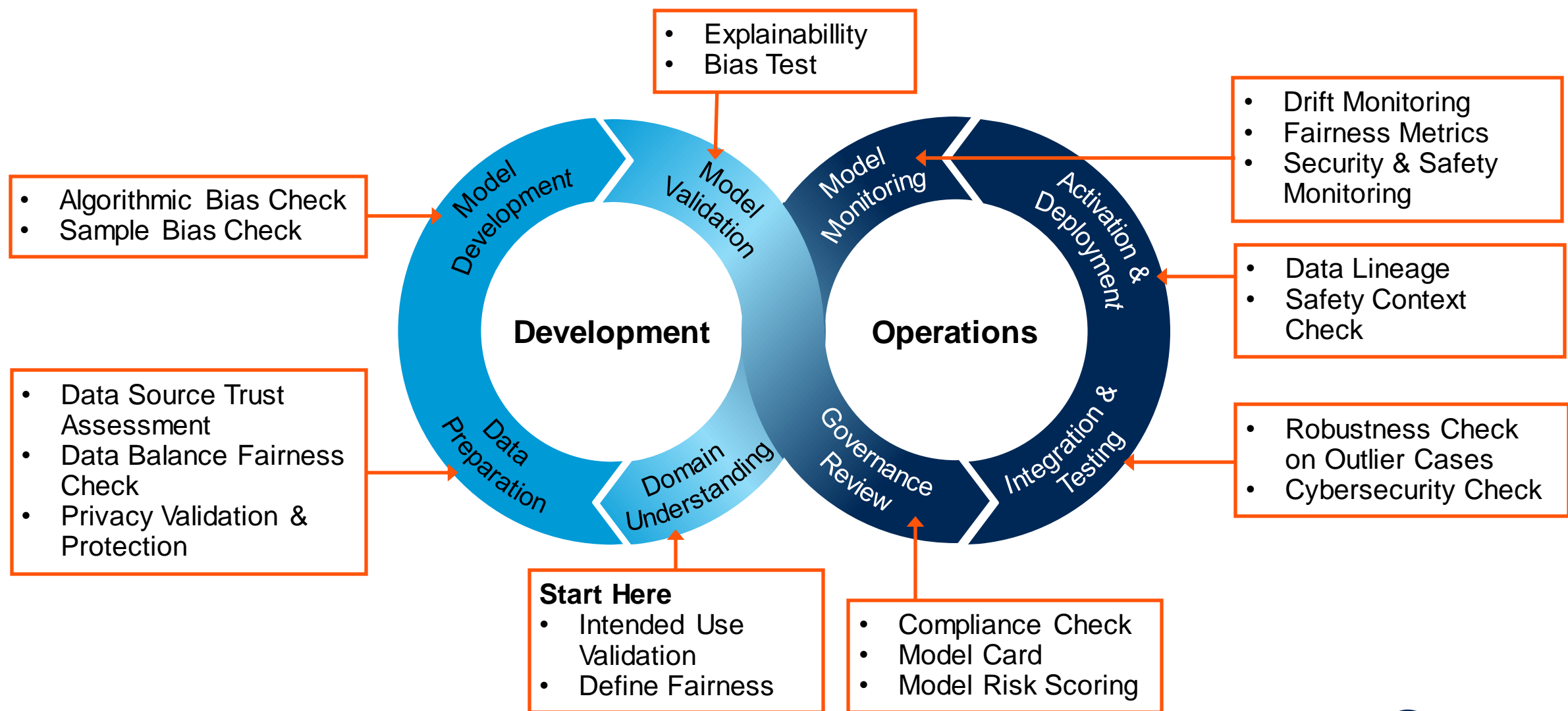
OSTP



# Key Issues

1. Eight AI risks and how to mitigate them
2. Emerging practices for AI risk management

# Full Life Cycle AI Risk Management



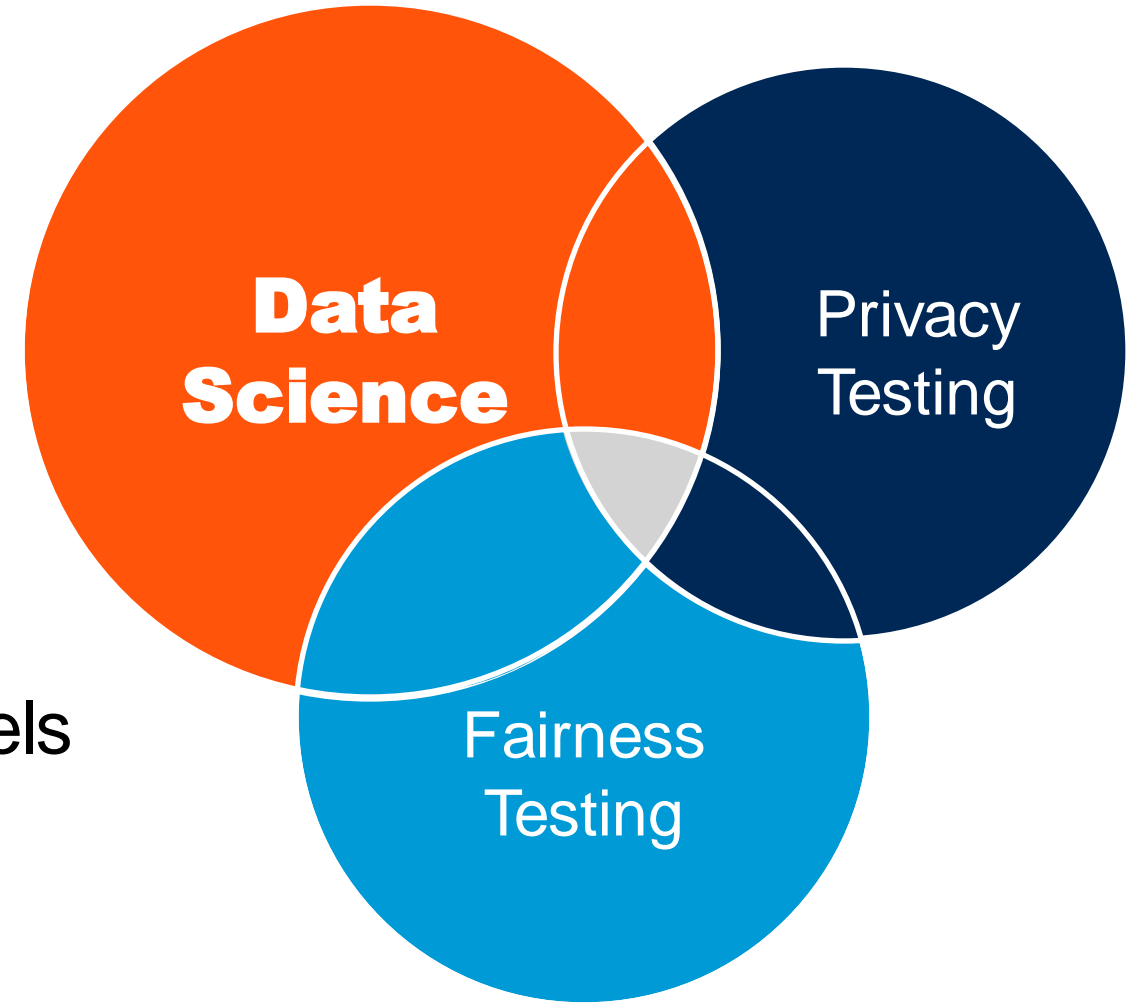
# Adding an AI Model Validator Role

**1/3**

## **model model validator developers**

---

- Not involved in model development
- Internal or external role
- Test for bias and other risks in models
- Audit for compliance



# Create an Awareness and Training Program for AI Risks & Ethics

- Every month, send around a dilemma. Discuss! What would you do?
- It trains being comfortable with ambiguity, and ethical thinking.



Mikkel Muller, CIO



**The Daily Telegraph**

U.K. 2016: The Daily Telegraph installed motion sensors under the desks of employees to improve energy sustainability. The sensor boxes, branded OccupEye and made by Cad-Capture, contain motion and temperature sensors, link back wirelessly to central network points, but do not monitor the desk areas themselves. It was, however, how the sensors could be used to monitor desk occupancy.

Sources: [Hacks Rebel After Bosses Secretly Install Motion Sensors Under Desks](#), The Register, Jan. 2016  
[National Union of Journalists Criticizes Daily Telegraph for Monitoring Journalists' Desk Time](#), The Guardian, Jan. 2016

Arguments Pro	Arguments Con
<ul style="list-style-type: none"><li>• Such instruments provide accurate occupancy data to building managers, who can then use it for maintaining temperature, reallocating unutilized desk space and other energy- and space-efficient measures.</li></ul>	<ul style="list-style-type: none"><li>• If the intent was to gather occupancy and temperature data, building managers could have identified alternative appropriate locations and methods for gathering data than installing the device beneath the desks of employees, which is more suspicious.</li></ul>

**S05 — Insight Versus Privacy**

Many digital initiatives start with the best of intentions — in this case, introducing sensors as part of a smart building sustainability initiative. But this is often followed by speculation and interpretation of the story — that they were meant to monitor employees. This demonstrates how highly sensitive people can be when it comes to issues they may see as affecting their privacy, even if this is not the case in practice.

**Gartner for IT Leaders Toolkit**  
Notes accompany this presentation. Please select Notes Page view to examine the Notes text.  
63 CONFIDENTIAL AND PROPRIETARY | © 2016 Gartner, Inc. and/or its affiliates. All rights reserved.

# Create and Involve AI Risks & Ethics Board

- AI risk management is not a checklist.
- Ethical dilemmas should not be left to project teams to figure out, causing stagnation or risky decisions.
- Board decides/advises on guidelines and project dilemmas.
- AI champion, business owners, external experts, other stakeholder representatives.





# Organizations That Manage Their AI Risks Realize More Value From AI



# Recommendations

- ④ Prioritize and contextualize AI risks for your organization. At least include privacy, fairness and accountability.
- ④ Not all regulations and ethics are globally the same. Take regional and cultural differences into account.
- ④ Make guidelines actionable. Do not reinvent the wheel: leverage existing policies, frameworks and tools.
- ④ Partner with AI governance, legal, cybersecurity and other stakeholders.
- ④ AI and AI risk management are about people. Build awareness and skills.

# Recommended Gartner Research

- 🔍 [\*\*A Comprehensive Guide to Responsible AI\*\*](#)  
Farhan Choudhary and Svetlana Sicular (G00764905)
- 🔍 [\*\*AI Ethics: Use 5 Common Principles as Your Starting Point\*\*](#)  
Frank Buytendijk, Erick Brethenoux and Others (G00774103)
- 🔍 [\*\*Applying AI — Governance and Risk Management\*\*](#)  
Avivah Litan, Svetlana Sicular and Others (G00745080)
- 🔍 [\*\*Digital Ethics by Design: A Framework for Better Digital Business\*\*](#)  
Frank Buytendijk, Jim Hare and Lydia Clougherty Jones (G00421878)
- 🔍 [\*\*Market Guide for AI Trust, Risk and Security Management\*\*](#)  
Avivah Litan, Jeremy D'Hoinne and Others (G00758388)