

## Highlights from the Responsible AI Leadership Consortium Event: AI-Driven Misinformation/Disinformation Seminar

**Date:** January 15, 2025

**Venue:** Permanent Mission of Liechtenstein to the United Nations, New York City

**Theme:** AI-Driven Misinformation/Disinformation and Possible Remedies

### Introduction

The Responsible AI Leadership Consortium (RAILC.org) convened a high-impact event that brought together global leaders, policymakers, and industry experts to address the accelerating challenges posed by AI-generated misinformation and disinformation. Hosted at the Permanent Mission of Liechtenstein to the United Nations, the discussions underscored the need for multi-stakeholder governance, ethical frameworks, and strategic foresight to navigate the evolving AI landscape.

### Keynote Highlights

The keynote provided an overview of the transformative journey of AI, from the advent of transformer models in 2017 to the exponential advancements in Large Language Models (LLMs). Emphasis was placed on the paradox of progress: while AI democratizes access to technology, it simultaneously exposes vulnerabilities, such as algorithmic "hallucinations" and synthetic data misuse.

### Critical Insights:

- **Misinformation Risks:** Deepfake technologies and open-source tools exacerbate political manipulation and fraud.
- **Data Challenges:** Despite strides in modeling, data management lags, necessitating improved governance frameworks.
- **Future Trends:** As AI models evolve to reason and self-learn, maintaining human oversight becomes pivotal.

### Panel Discussion 1: The Evolution of AI-Driven Misinformation & Emerging Threats

The panel discussion focused on how AI technologies, such as deepfakes, synthetic media, and large language models (LLMs), have significantly advanced disinformation strategies. Emerging threats, including real-time adversarial content creation, were highlighted. The conversation emphasized the necessity of proactive measures, such as regulatory frameworks, public education, and the development of robust detection tools.

The panel underscored the dual nature of AI - its potential to improve access to information while being misused by bad actors to propagate disinformation - and stressed the urgent need for safeguards to uphold truth in the digital age.

## **Panel Discussion 2: Policy, Regulation & Global Responses to AI-Generated Disinformation**

The panel addressed the human rights dimensions of AI governance, emphasizing the importance of partnerships to bridge the digital divide, capacity-building efforts and the establishment of multi-stakeholder decision-making bodies. The connection to the SDG 5 - Digital Literacy action items and the goals in the UN Global Digital Compact were also cited.

The discussion critiqued existing frameworks, highlighting their limitations in addressing environmental justice and corporate accountability. Preventive measures to safeguard democracy from the threats posed by unchecked AI were prioritized.

The panel highlighted the critical importance of cross-sector collaboration in enhancing data governance. Although AI models are advancing at an unprecedented pace, the supporting data infrastructure frequently lags and fails to meet best practices. To address this gap, the panel advocated for a multidisciplinary approach, emphasizing the need to prioritize data quality and uphold model credibility.

## **Key areas for implementation**

1. **Governance and Regulation:**
  - Establish independent agencies to oversee AI governance, focusing on data quality and not just technology.
  - Promote real-time fact-checking systems and transparent model training processes to combat misinformation.
2. **Ethical Data Management:**
  - Encourage tagging (labeling) sharing, and verifying data to build trustworthy and high quality AI systems.
  - Balance synthetic and organic data to ensure the reliability of foundational models.
3. **Sustainability and AI's Environmental Impact:**
  - Analyze and verify the facts surrounding the energy-intensive demands of AI infrastructure through investments in renewable technologies.
  - Foster global dialogue on sustainable AI practices, with the Global South as a key stakeholder.
4. **Human-Centric AI Development:**
  - Embed human oversight in AI systems while acknowledging the limits of scalability. What happens when the human expertise dwindles and AI systems have much more knowledge and expertise.
  - Leverage synthetic data with built-in governance to simulate edge-case scenarios.

## Future Direction

The event brought out the urgency of accelerating multi-stakeholder dialogues and collaborations. ***Some participants suggested the inevitability of an AI-related "event" triggering regulatory overcorrection, drawing parallels with the 2008 financial crisis.*** To avoid stifling innovation, frameworks must prioritize flexibility, inclusivity, and adaptability - Government regulators are you listening?

## Proposed Next Steps:

- Advocate for the development of a global compact on AI governance, incorporating ethical principles, environmental sustainability, and human rights protections to ensure responsible AI use worldwide.
- Launch regional and sector-specific pilot projects to test innovative AI governance models, providing insights to harmonize regulatory efforts across different jurisdictions.
- Establish collaborative platforms for real-time knowledge exchange among governments, private sectors, fostering transparency and shared learning in AI governance practices (UN University | World Bank | OECD).
- Promote public education initiatives to empower individuals with fact-checking skills, enabling them to verify important information independently and combat misinformation effectively.