# ETHICS AND TRUST FROM A HUMAN-CENTRIC PERSPECTIVE

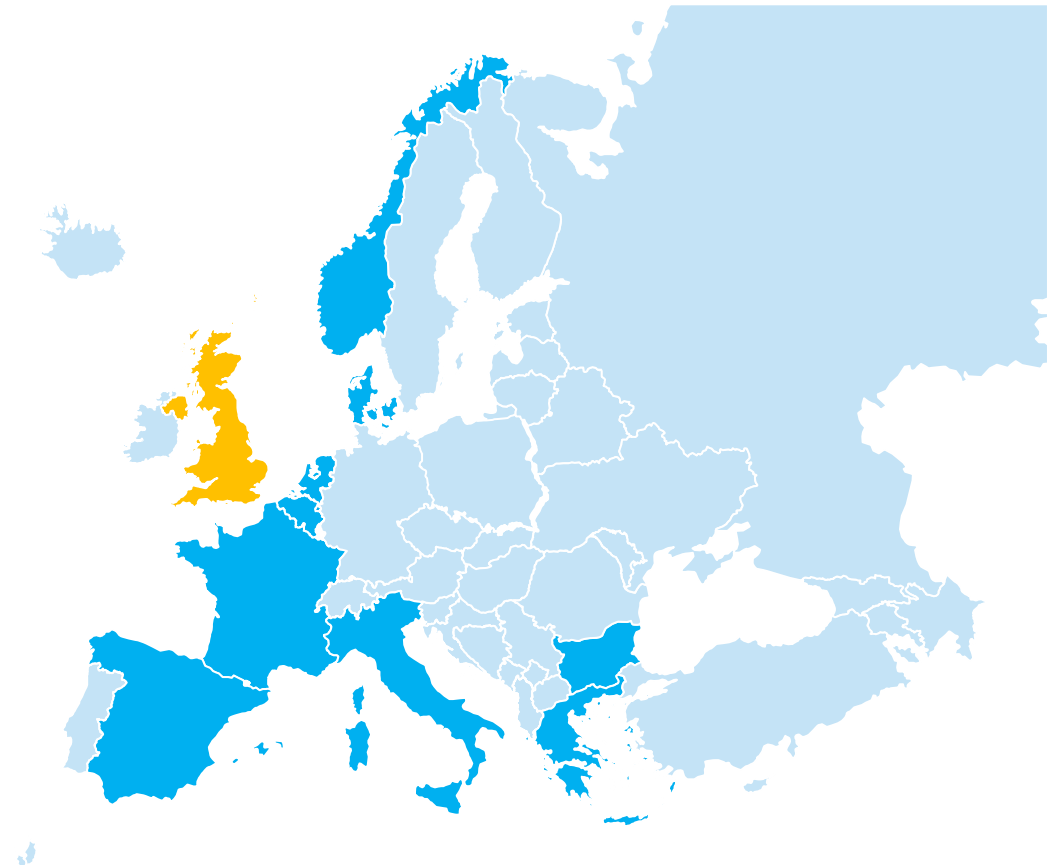ENFIELD WEBINAR, NOVEMBER 15, 2024

ASBJØRN FØLSTAD, SINTEF

THEMIS 5.0 (RIA): Human-centred trustworthiness optimization in hybrid decision making support

Call: HORIZON-CL4-2022-HUMAN-02 - 'a human-centred and ethical development of digital and industrial technologies'.

Topic: HORIZON-CL4-2022-HUMAN-02-01   AI for human empowerment (AI, Data and Robotics Partnership)
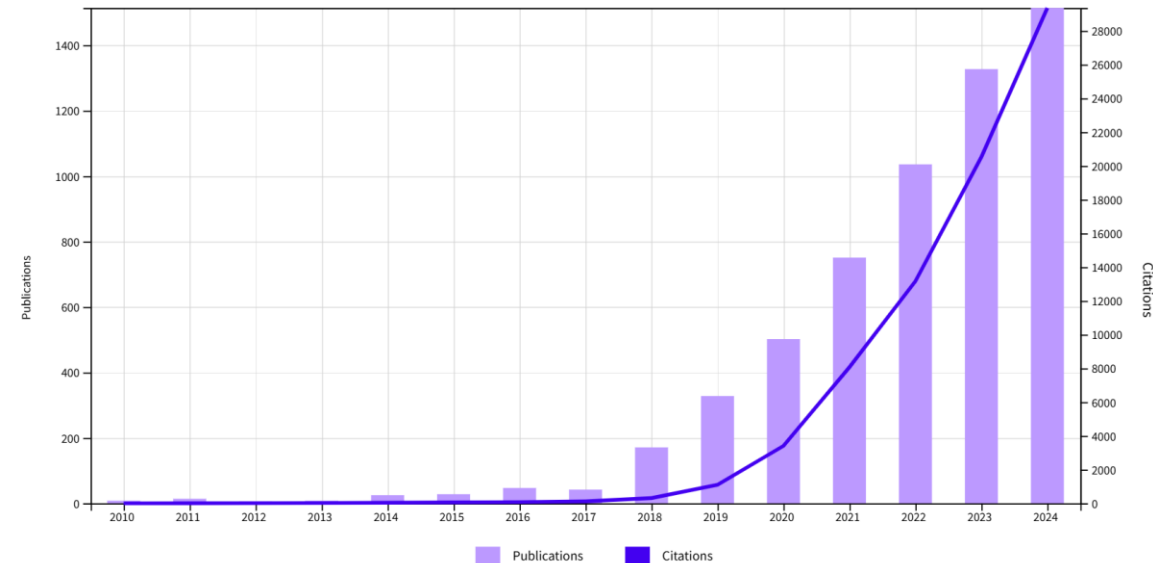
# ETHICS, TRUST, AND TRUSTWORTHY AI

- **AI ethics and trust** key issues as AI becomes increasingly advanced and taken up in society

- **Trustworthy AI** key to meeting ethical requirements and establish a sound basis for trust in AI

Citation report: ethics AND «artificial intelligence»
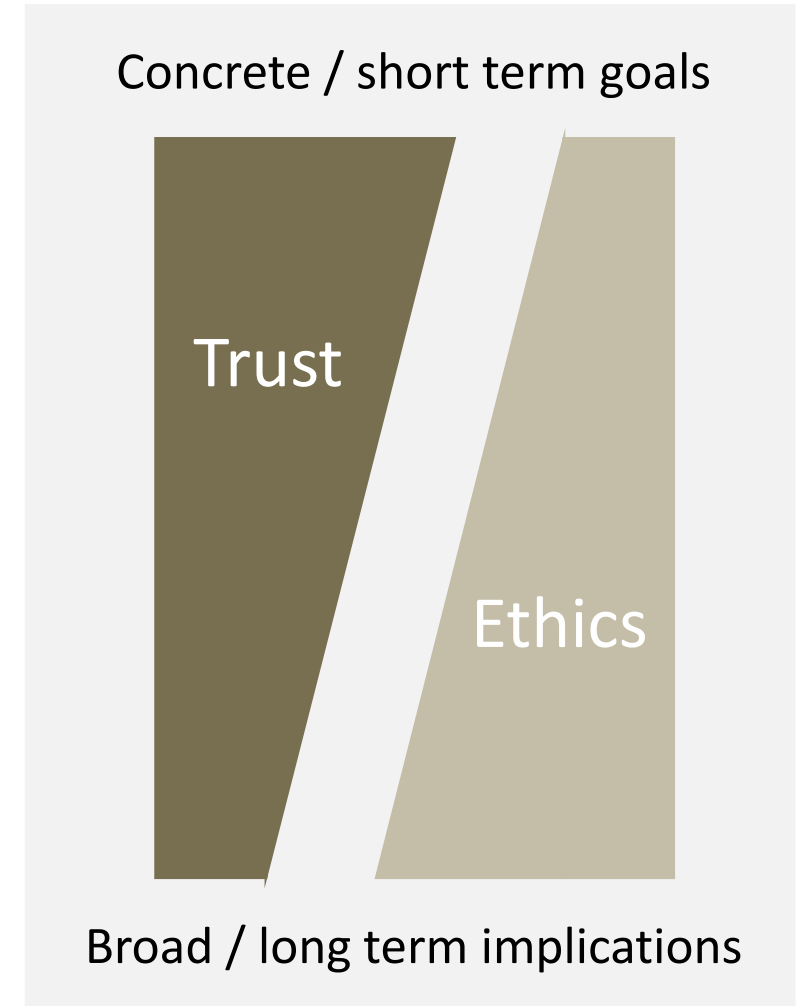


https://www.webofscience.com/
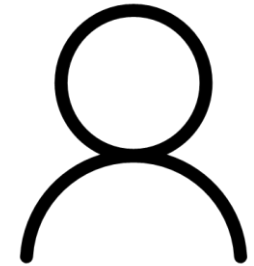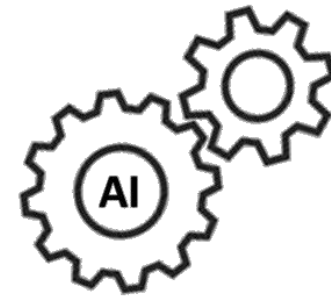
# ETHICS, TRUST, AND TRUSTWORTHY AI

- **AI ethics and trust** key issues as AI becomes increasingly advanced and taken up in society

- **Trustworthy AI** key to meeting ethical requirements and establish a sound basis for trust in AI

?

Concrete / short term goals

Trust

Ethics

Broad / long term implications

# TRUST AND TRUSTWORTHINESS

- Trust – in the eye of the beholder

- Trustworthiness – characteristic of the AI system

Trusting?

Trustworthy?

## TRUST AND TRUSTWORTHINESS

- Trust – in the eye of the beholder

- Trustworthiness – characteristic of the AI system

- Beneficial use of AI depends on **adequate configuration of trust** to the actual trustworthiness of the AI system



The New York Times

A.I. and Chatbots ›    Testing a Tutorbot    Chatbot Prompts to Try    A.I.'s Literary Skills    Spot the A.I. I

### The ChatGPT Lawyer Explains Himself

In a cringe-inducing court hearing, a lawyer who relied on A.I. to craft a motion full of made-up case law said he "did not comprehend" that the chat bot could lead him astray.
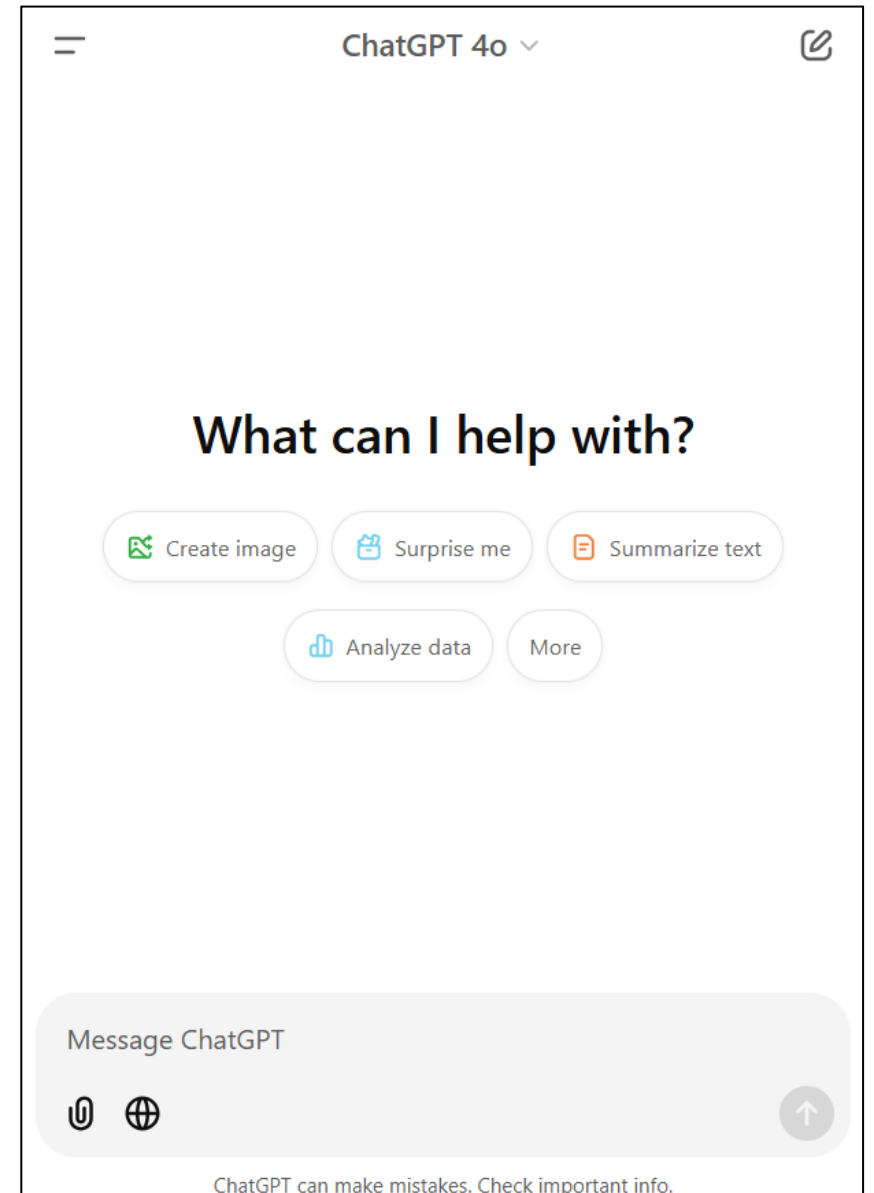
Give this article    267

Steven A. Schwartz told a judge considering sanctions that the episode had been "deeply embarrassing."  Jefferson Siegel for The New York Times

By Benjamin Weiser and Nate Schweber
June 8, 2023

hhttps://www.nytimes.com/2023/06/08/nyregio n/lawyer-chatgpt-sanctions.html

# TRUSTWORTHY AI

- Trustworthy AI addressed through **detectable trustworthiness characteristics**,

- … that is, characteristics on which adherence to ethical requirements can be assessed

# TRUSTWORTHY AI

- Trustworthy AI addressed through **detectable trustworthiness characteristics**,

- … that is, characteristics on which adherence to ethical requirements can be assessed

Accuracy!

## TRUSTWORTHY AI

- Trustworthy AI addressed through **detectable trustworthiness characteristics**,

- ... that is, characteristics on which adherence to ethical requirements can be assessed

- **Balancing different approaches and frameworks**

- Ethics: Broad range of ethical guidelines and requirements

- Trustworthy AI: Broad range of trustworthiness characteristics. Partially overlapping frameworks

- Balancing different approaches and frameworks

- **Ethics: Broad range of ethical guidelines and requirements**

- Trustworthy AI: Broad range of trustworthiness characteristics. Partially overlapping frameworks



Table 1 Overview of AI ethics

Hagendorff, T. (2020). The ethics of AI ethics: An evaluation of guidelines. Minds and machines, 30(1), 99-120.

## ETHICS AND TRUSTWORTHINESS – CONCEPTUAL BALANCING

- Balancing different approaches and frameworks

- Ethics: Broad range of ethical guidelines and requirements

- **Trustworthy AI: Broad range of trustworthiness characteristics. Partially overlapping frameworks**



Human agency
Robustness and safety
Privacy
Transparency
Fairness
Well-being
Accountability

INDEPENDENT HIGH-LEVEL EXPERT GROUP ON ARTIFICIAL INTELLIGENCE
SET UP BY THE EUROPEAN COMMISSION

ETHICS GUIDELINES FOR TRUSTWORTHY AI

Validity/reliability
Safety
Security/resilience
Explainability
Privacy
Fairness
Accountability

NIST AI 100-1

Artificial Intelligence Risk Management Framework (AI RMF 1.0)

NIST NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY U.S. DEPARTMENT OF COMMERCE

- **Balancing for specific contexts and situations**

- Use case specific ethical requirements, with implications for AI trustworthiness assessment and optimization

# ETHICS AND TRUSTWORTHINESS – CONTEXTUAL BALANCING

THEMIS 5.0
USE CASES:



**News media:** Mitigate disinformation

- Balancing for specific contexts and situations



**Maritime:** Optimise port operations

- **Use case specific ethical requirements, with implications for AI trustworthiness assessment and optimization**



**Healthcare:** Personalised risk prediction

# TRUSTWORTHINESS ASSESSMENT AND OPTIMIZATION – THE THEMIS 5.0 APPROACH

- **Assessment of AI system in socio-technical environment**
  - Use-cases: healthcare, news media, port management

- **Detectable trustworthiness characteristics for assessment**
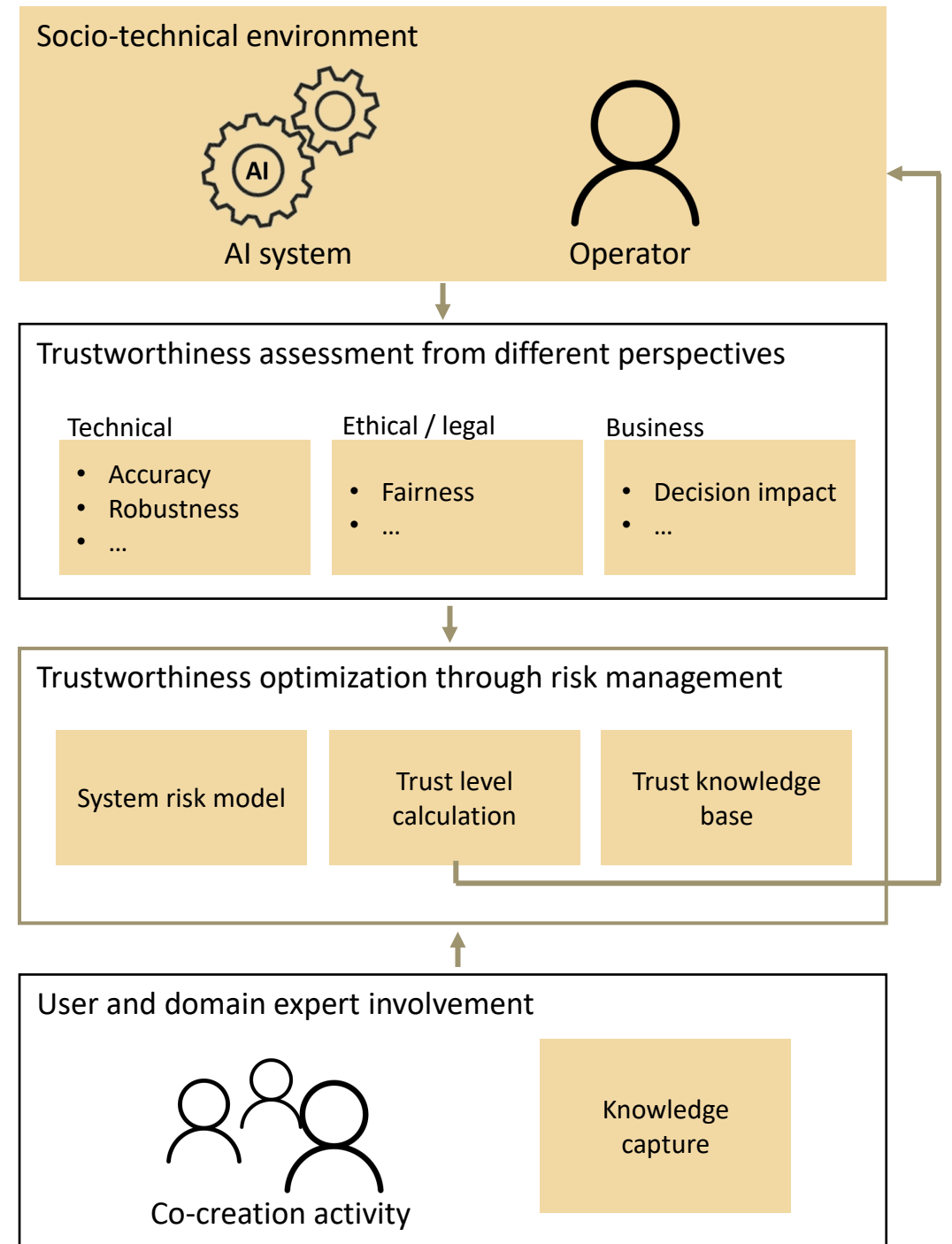
- **Trustworthiness optimization through risk management**
  - Optimized with regards to user group preferences

- **Human-centred assessment**
  - Identification of trustworthiness characteristics with users and domain experts
  - Explanations for assessment by users and domain experts

**Socio-technical environment**

AI system          Operator

**Trustworthiness assessment from different perspectives**

| Technical | Ethical / legal | Business |
| --- | --- | --- |
| • Accuracy<br>• Robustness<br>• … | • Fairness<br>• … | • Decision impact<br>• … |

**Trustworthiness optimization through risk management**

| System risk model | Trust level calculation | Trust knowledge base |
| --- | --- | --- |

**User and domain expert involvement**

Co-creation activity          Knowledge capture

- AI ethics and trust with high attention, in society and academia alike

- Trustworthy AI key to meeting ethical requirements and establish a sound basis for trust in AI

- Ethical requirements and corresponding trustworthiness characteristics dependent on user case

- Human-centre approach needed to assess and optimize AI trustworthiness

THEMIS 5.0