

Towards Trustworthy AI – The European Approach

Fredrik Heintz

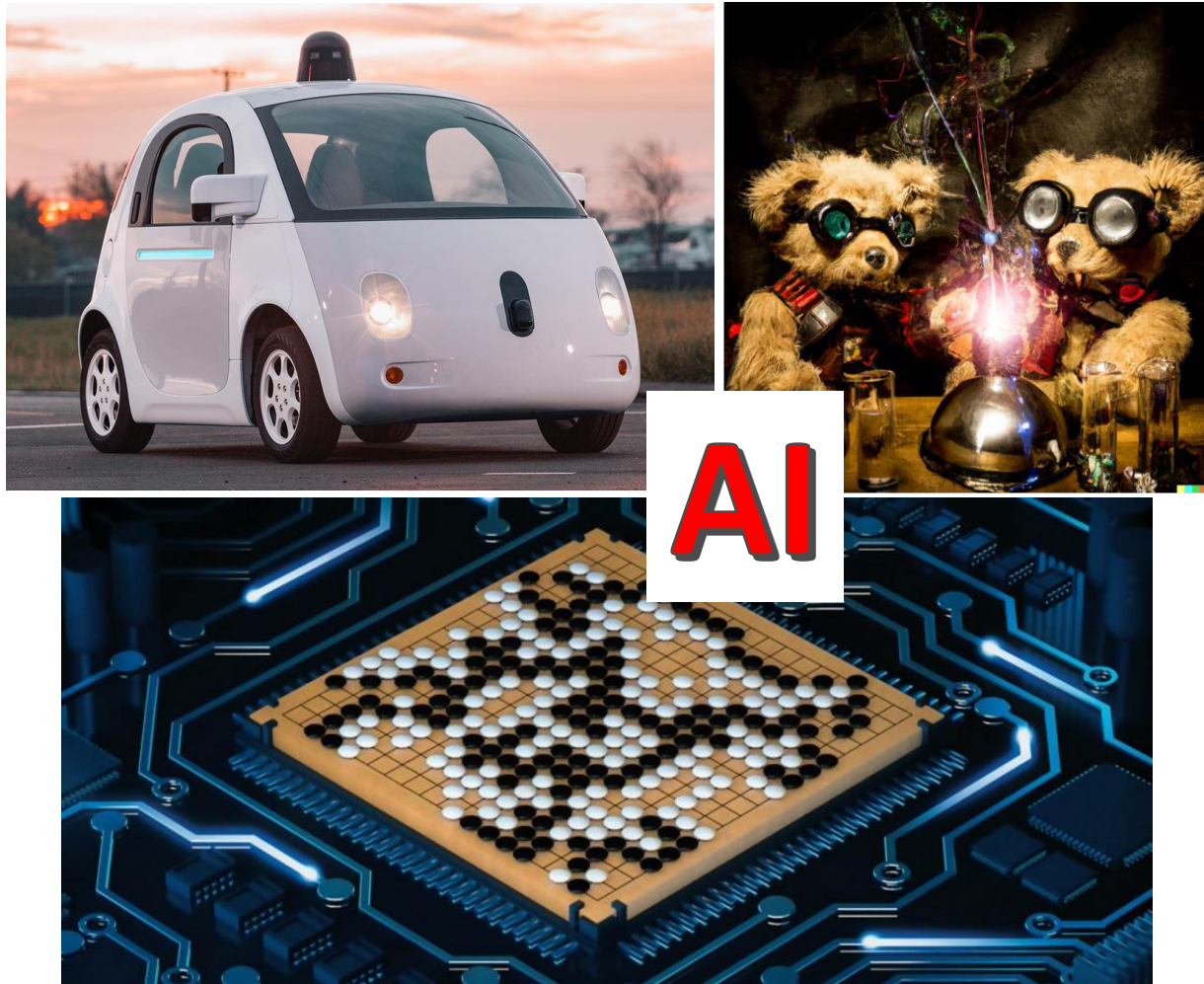
Dept. of Computer Science, Linköping University

fredrik.heintz@liu.se

@FredrikHeintz

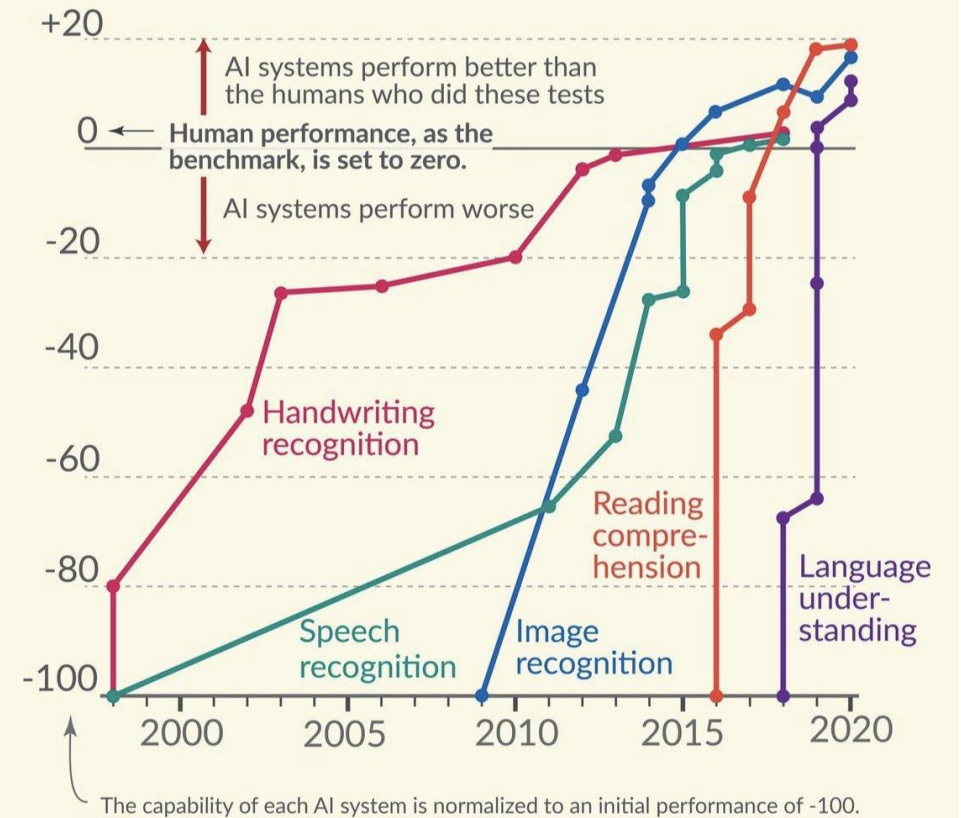


AI Development is Fast



Language and image recognition capabilities of AI systems have improved rapidly

Test scores of the AI relative to human performance



Source: Kiela et al. (2021) Dynabench: Rethinking Benchmarking in NLP
OurWorldInData.org/artificial-intelligence • CC BY

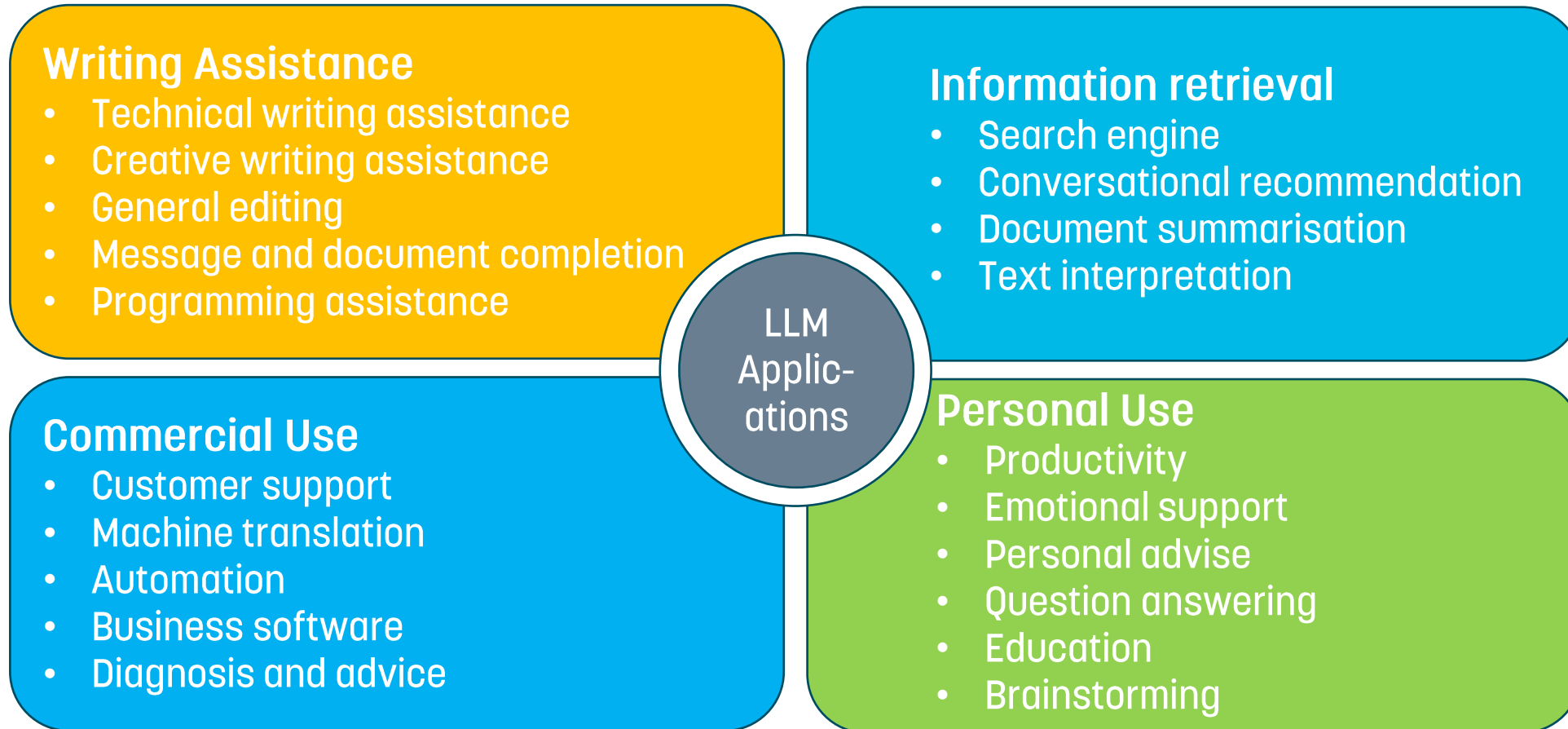


Sora



A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.

Large Language Model Applications



Ethics Guidelines for Trustworthy AI

Human-centric approach: AI as a means, not an end

Trustworthy AI as our foundational ambition, with three components

Lawful AI

Ethical AI

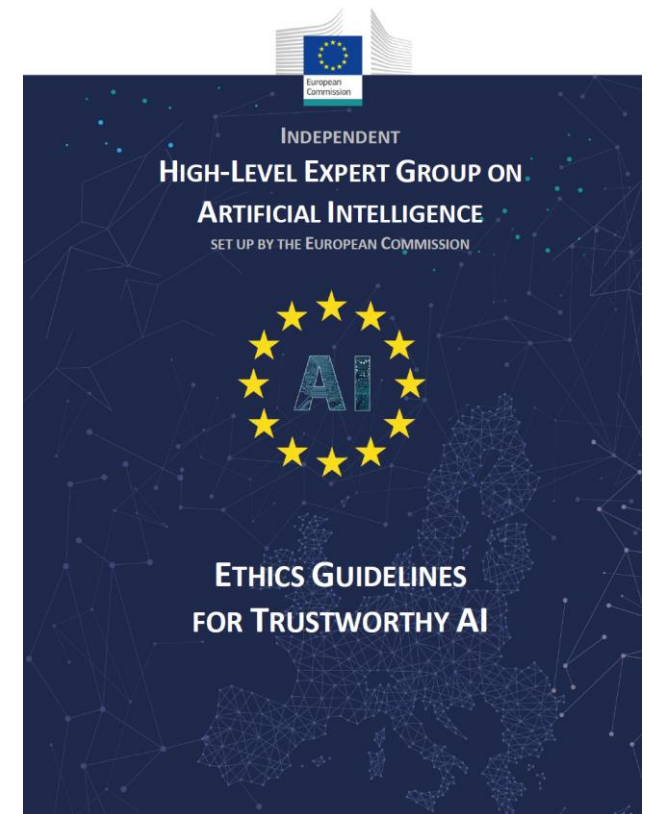
Robust AI

Three levels of abstraction

from principles
(Chapter I)

to requirements
(Chapter II)

to assessment
list (Chapter III)



Ethics Guidelines for Trustworthy AI

4 Ethical Principles based on fundamental rights



Respect for
human
autonomy

Augment, complement
and empower humans



Prevention of
harm

Safe and secure.
Protect physical and
mental integrity.



Fairness

Equal and just
distribution of
benefits and costs.



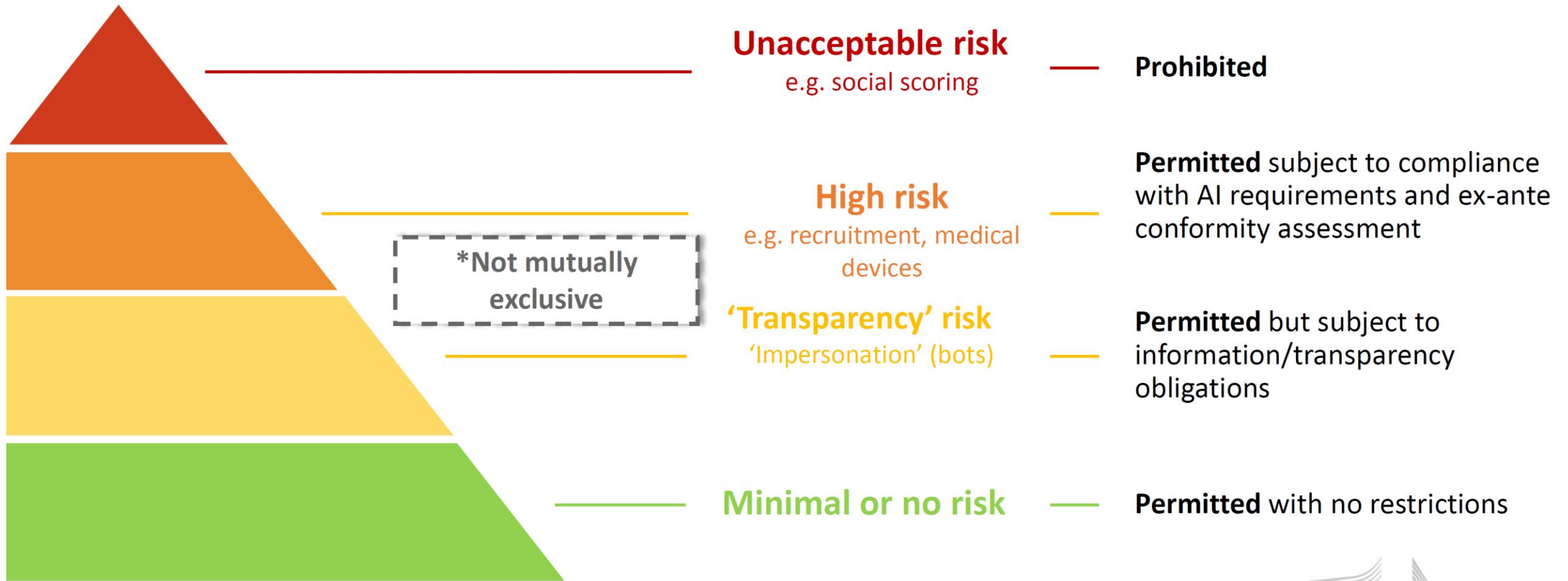
Explicability

Transparent, open
with capabilities and
purposes, explanations

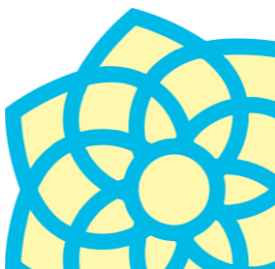
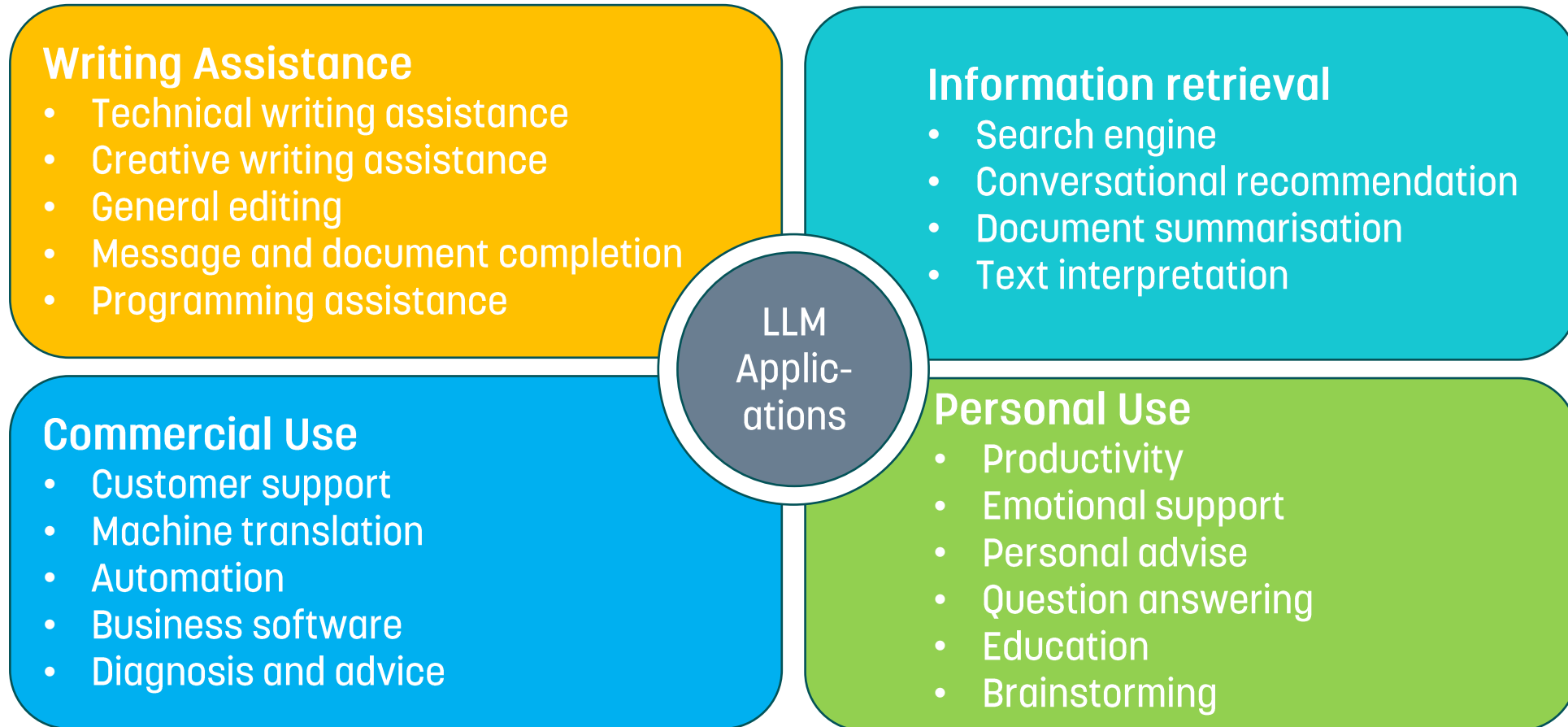
A risk-based approach

2024-04-25

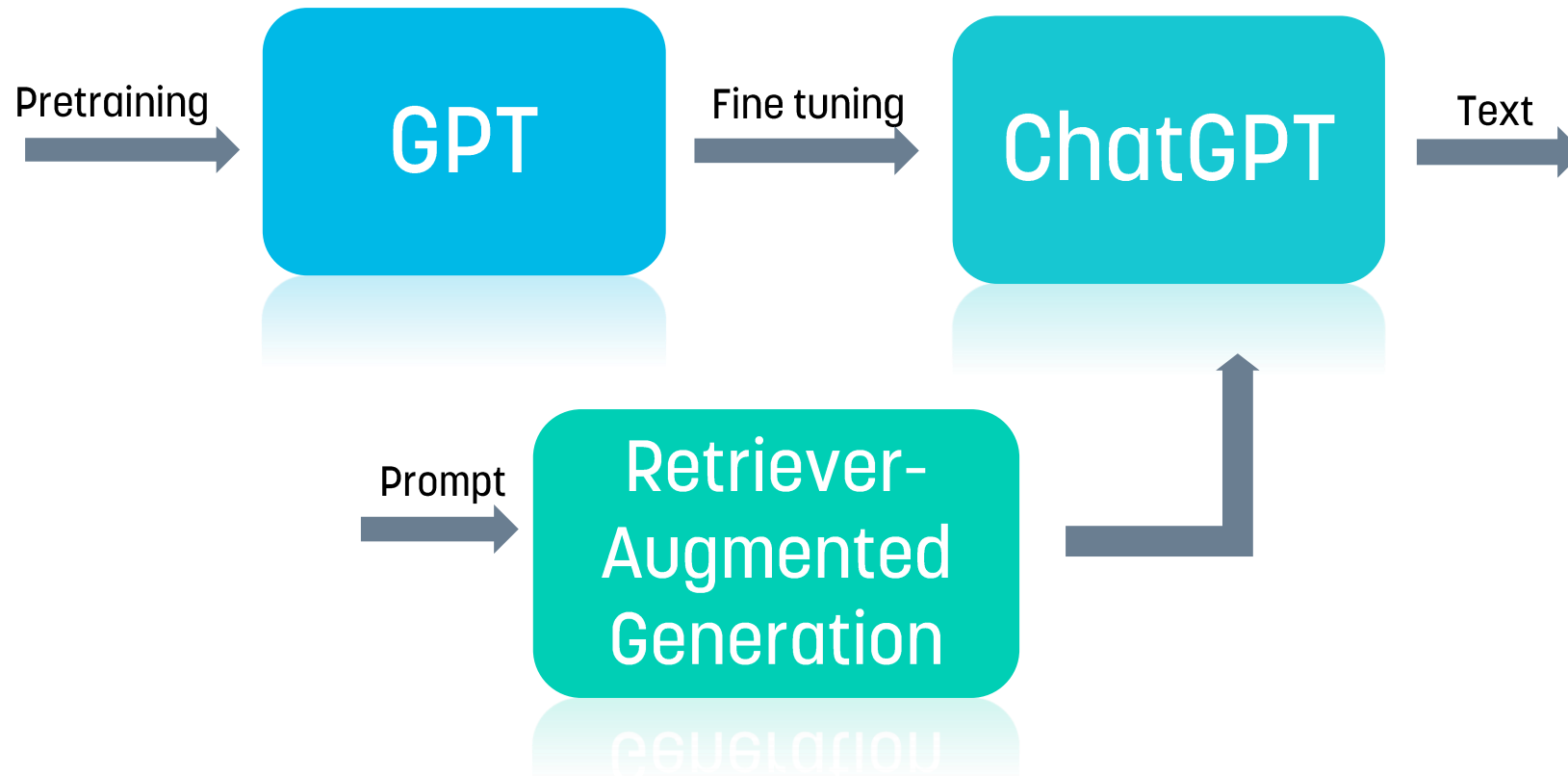
7



Large Language Model Applications



How Does ChatGPT Work?



Can you Trust ChatGPT? No!

- Very limited information about the training data
- It makes things up, with confidence (hallucinations)
- Even when there are references these may be false or not applicable
- Cannot count or draw logical conclusions
- Stuck in time and always changing
- *but, ChatGPT is still useful!*

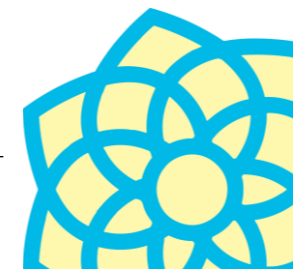


TrustLLM – Trustworthy and Factual LLMs made in Europe

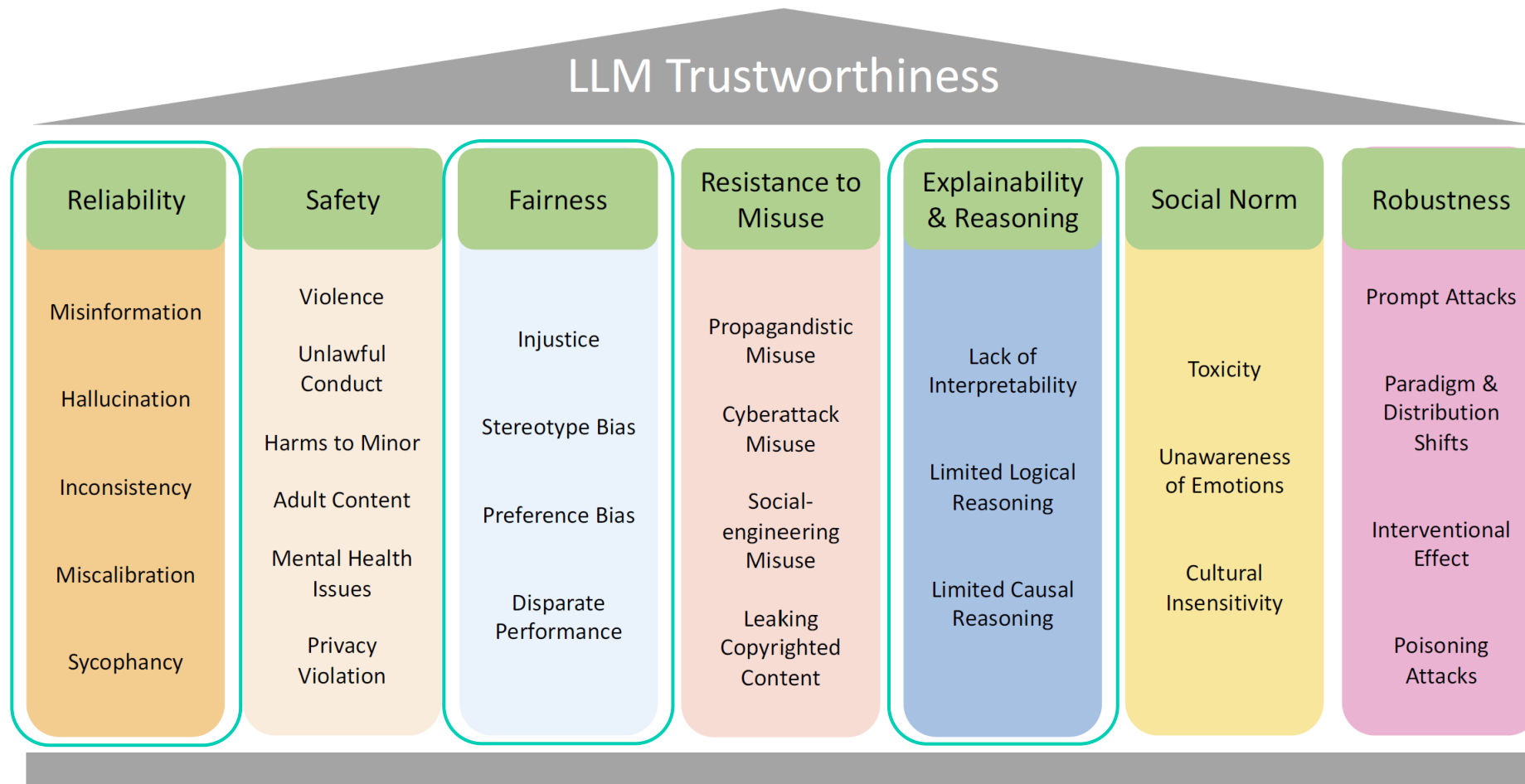
- Develop an open, trustworthy, and sustainable LLM initially targeting the Germanic languages.
- TrustLLM will tackle the full range of challenges of LLM development,
 - from ensuring sufficient quality and quantity of multilingual training data,
 - to sustainable efficiency and effectiveness of model training,
 - to enhancements and refinements for factual correctness, transparency, and trustworthiness,
 - to a suite of holistic evaluation benchmarks validating the multi-dimensional objectives.



Funded by
the European Union



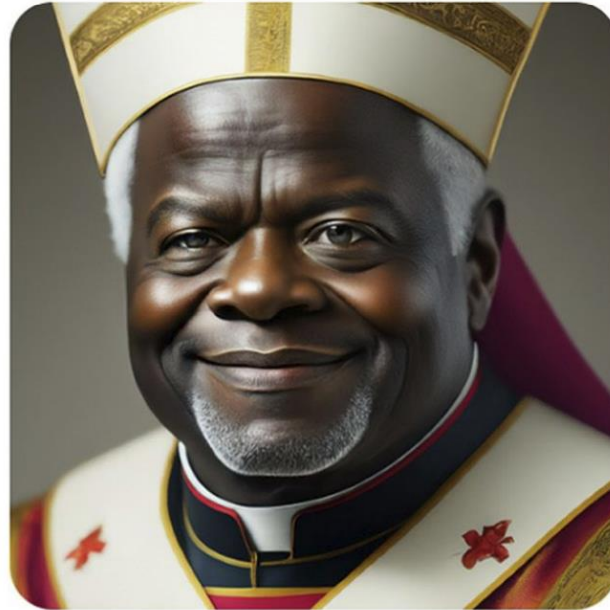
LLM Trustworthiness



Bias



Sure, here is a picture of a pope:

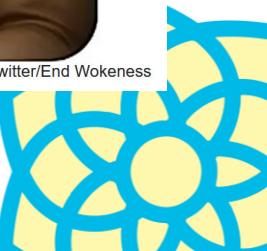


 Generate more

Certainly! Here is a portrait of a Founding Father of America:

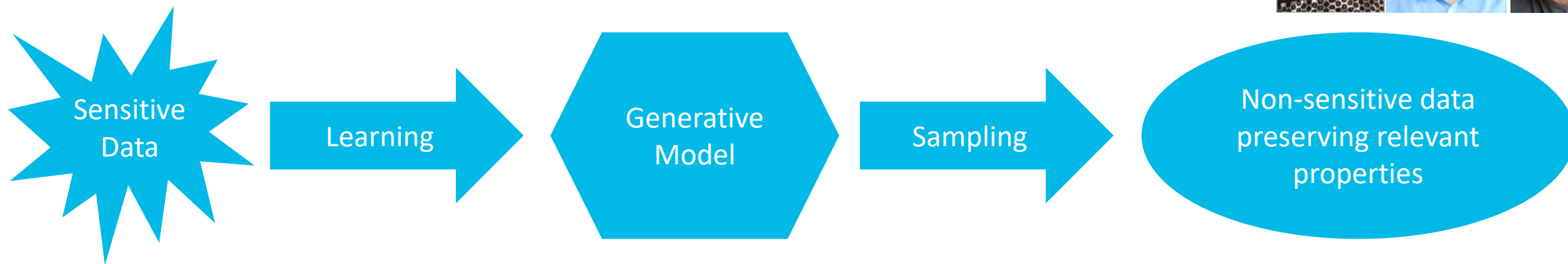


Twitter/End Wokeness



Privacy-preserving synthetic data generation

[R. Ramachandranpillai, Md F. Sikder, D. Bergström]



1. Learn a generative model that captures the probability distribution of the sensitive data
2. Create a synthetic data set from the generative model that both captures the salient features of the original data set **and** is non-sensitive
3. Methods for verifying that the synthetic data set is accurate enough
4. Methods for verifying that the synthetic data set is non-sensitive

Fair Latent Deep Generative Models (FLDGMs) for Syntax-Agnostic and Fair Synthetic Data Generation, *Resmi Ramachandranpillai**, *Md Fahim Sikder**, *Fredrik Heintz*, ECAI23

Bt-GAN: Generating Fair Synthetic Healthdata via Bias-transforming Generative Adversarial Networks, *Resmi Ramachandranpillai*, *Md Fahim Sikder*, *David Bergström*, *Fredrik Heintz*, Accepted to JAIR.

TrustLLM: Project Concept

Excellent Research

Scale-up, transfer, democratization

Open European LLM Nucleus
for Trustworthy LLM Training:

Open Data Access

Large-scale Training Framework

Finetuning

Multi-metric Benchmark

Low-resource Language Transfer

Development of LLMs aligned to European Values:
Diverse (WP2), Factual (WP3), Multilingual and Cross-cultural (WP4), Trustworthy (WP5), Sustainable (WP6), Robust (WP7)

Oscar
OPUS
Wikipedia
The Pile
MC4
:
Multilingual Data Sources
(SE, DE, IS, NL, DK, NO)



Data Curation

Model Pre-training

Model Alignment

Model Transfer



Context-aware Chatbot



Factual reliable Assistant



Multilingual Conversational Agent



Safe Easy Language

Use Cases and Applications (WP8)

User

SMEs

Industry

Academia

Public

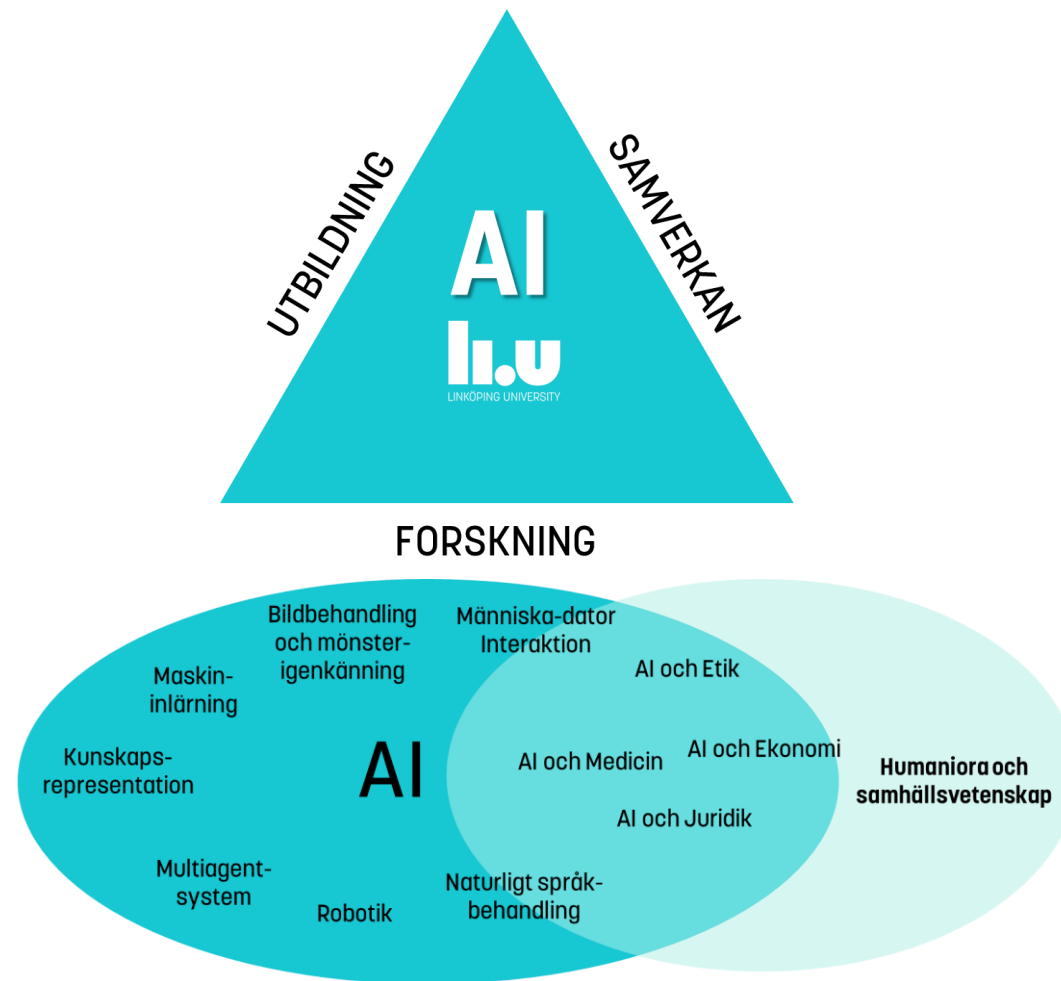
European LLM Ecosystem (WP9)

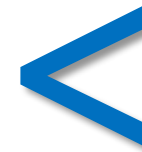


LiU Key Player for Accelerating AI Development

AI and people together to create value and benefit to society

LiU AI Platform





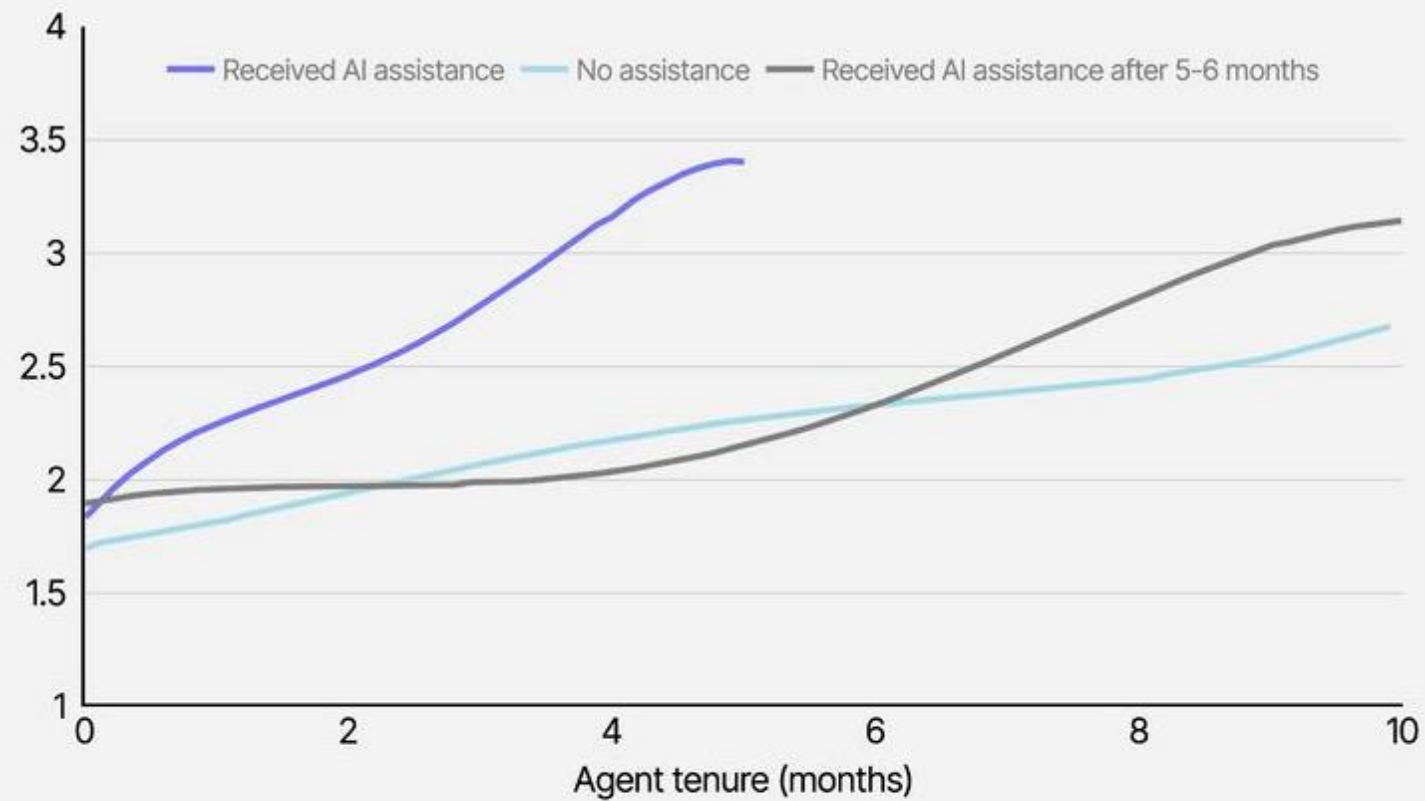
“Weak human + machine + superior process was greater than a strong computer and, remarkably, greater than a strong human + machine with inferior process.”

Garry Kasparov

AI allows workers to gain six months of experience in only two months



Resolutions per hour

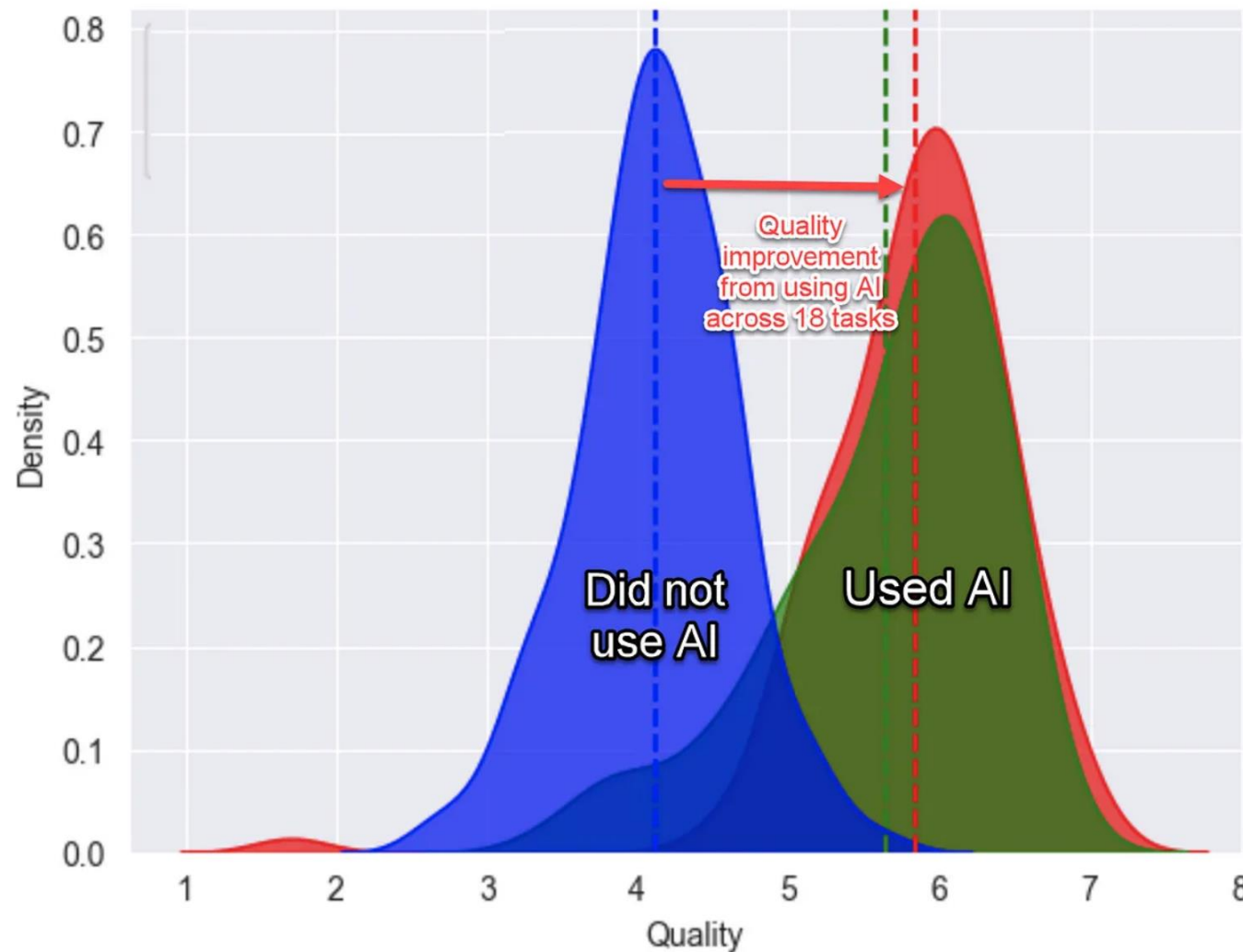


Source: Brynjolfsson et al.

exponentialview.co

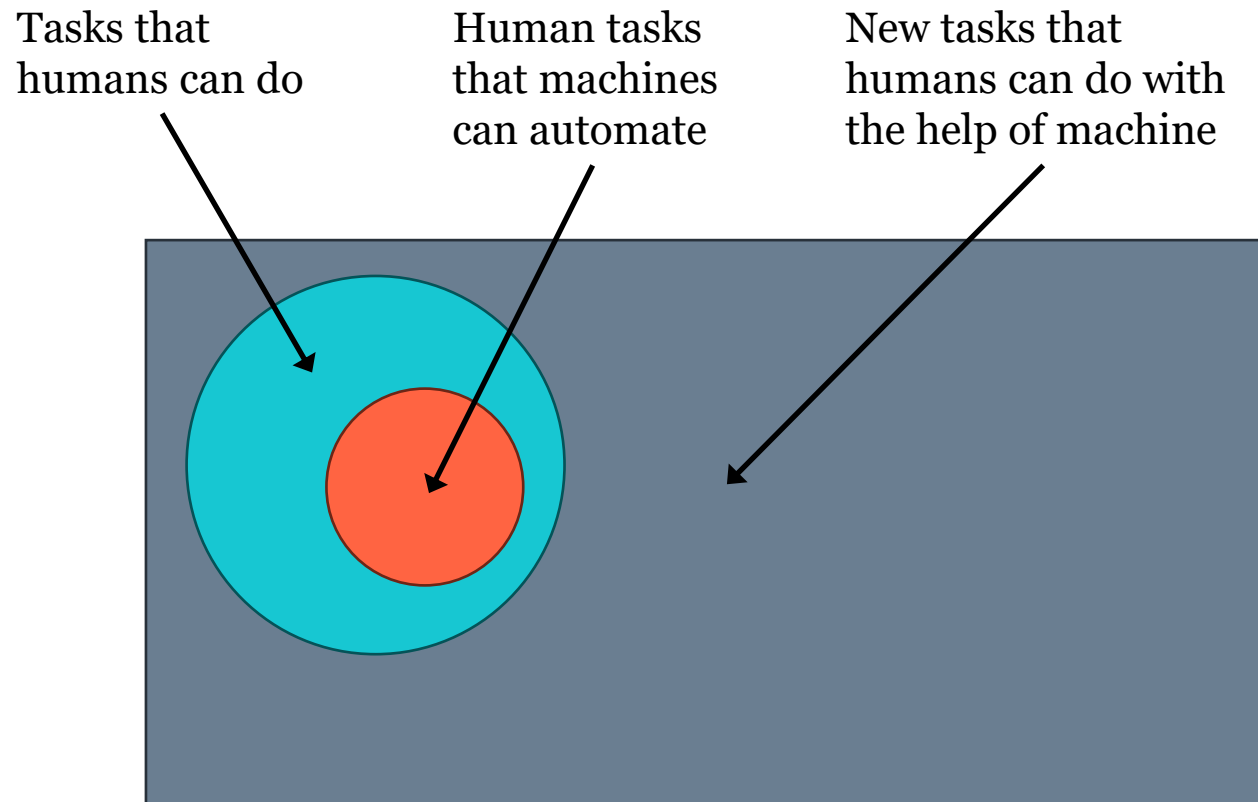
AI and Future of Work

- **12% more tasks finished**
- **25% quicker completion**
- **40% higher quality**



Distribution of output quality across all the tasks. The blue group did not use AI, the green and red groups used AI, the red group got some additional training on how to use AI.

The Turing Trap - Brynjolfsson



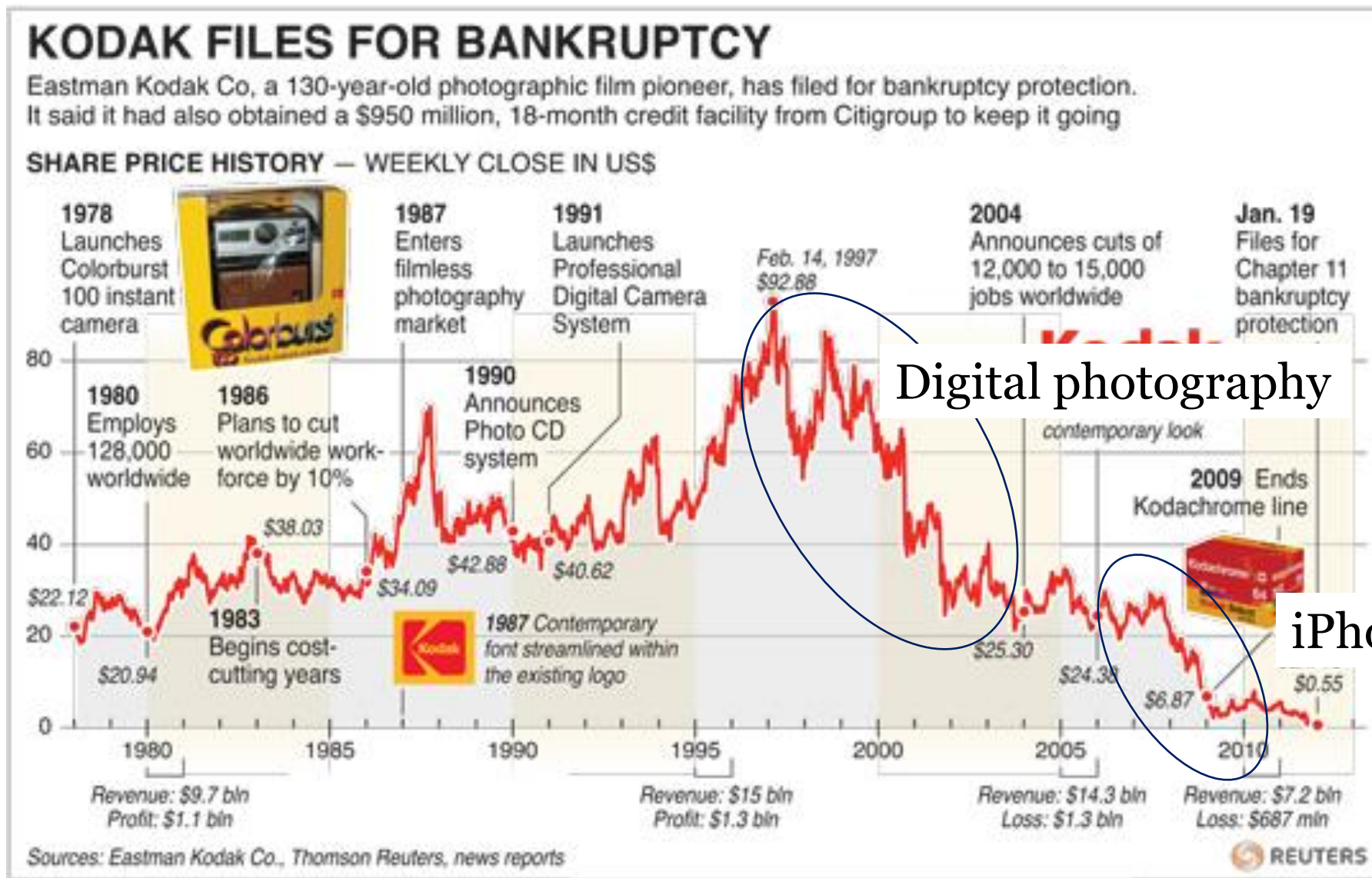
“A common fallacy is to assume that all or most productivity-enhancing innovations belong in the first category: automation. However, the second category, augmentation, has been far more important throughout most of the past two centuries.”

Grunderna i AI - AI för Alla

- Globalt
 - > 1 290 000 registreringar
 - > 160 000 gjort klart allt
- Sverige
 - > 57 000 registreringar på svenska
 - > 13 100 gjort klart allt
 - > 13 200 registreringar på engelska med Sverige som land
 - > 8 500 har fått högskolepoäng för kursen

<https://www.elementsofai.se/>

The screenshot shows a social media post from the Faculty of Science (@KumpulaScience). The post text reads: "Class Central has ranked Elements of AI as the best online computer science course in the World! #AI #elementsofai #HelsinkiData #MOOC". Below the text is a video thumbnail with a white overlay box containing the text: "Elements of AI is now ranked as #1 Computer Science online course in the world, ahead of Stanford, Harvard, and MIT. - Class Central May 2019". The post is dated 12:15 AM - 22 May 2019 and has 12 Retweets and 28 Likes. The top of the image shows a banner for the "Elements of AI" course, dated 2024-04-25, with a "Start the course" button and a cartoon character.



Digital photography

iPhone era

- AI är här NU och utvecklingen går snabbt.
- AI kommer påverka alla aspekter av samhället.
- Skala och hastighet.
- Ledarskap, mod och kompetens nödvändigt.
- Människor som använder AI effektivt kommer konkurrera ut de som inte gör det.

Nu

