# TrustLLM

## - Towards Trustworthy and Factual Large-Language Models

Fredrik Heintz
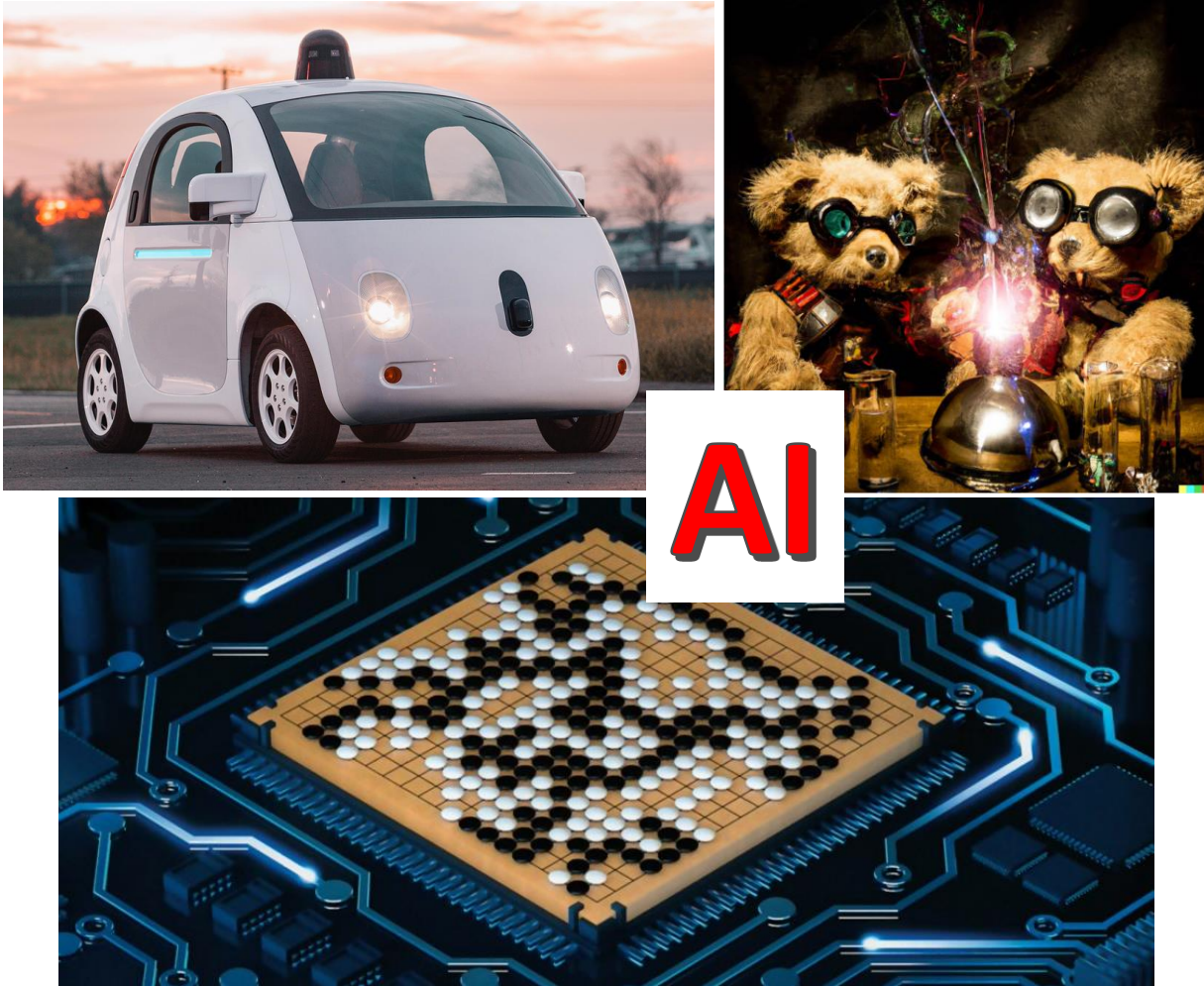
Dept. of Computer Science, Linköping University

fredrik.heintz@liu.se

@FredrikHeintz

# AI Development is Fast

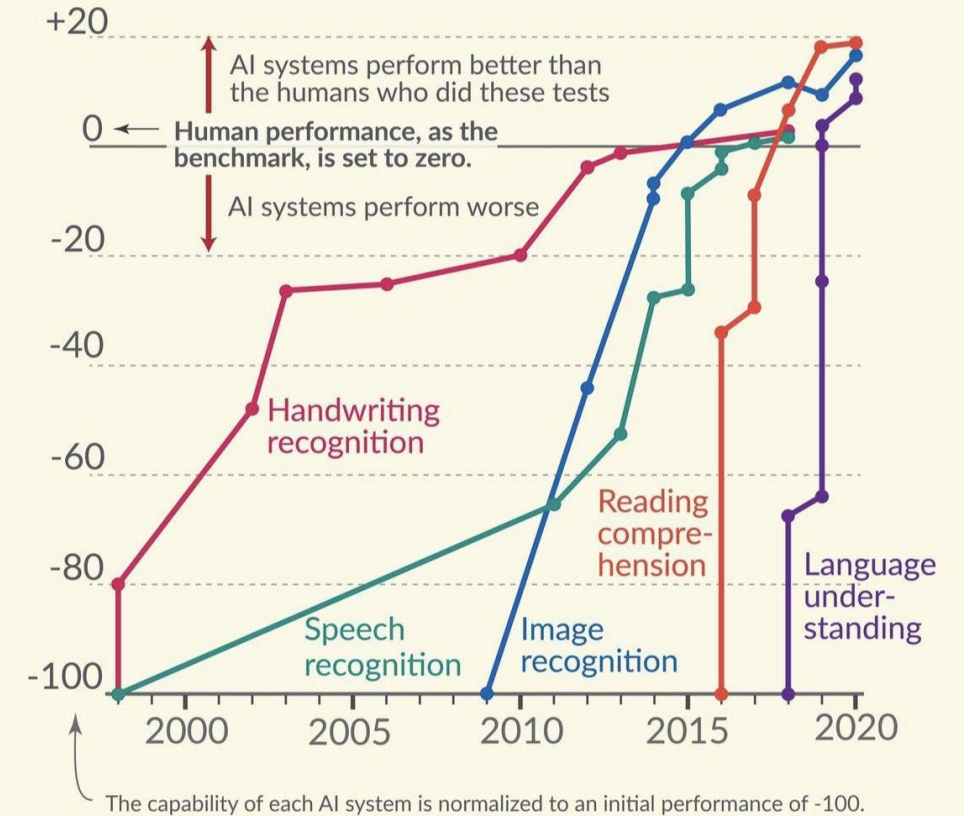# Ethics Guidelines for Trustworthy AI

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) |
|:---------------------------:|:----------------------------:|:--------------------------------:|

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

ETHICS GUIDELINES
FOR TRUSTWORTHY AI

LINKÖPING UNIVERSITY

https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai

# A risk-based approach

**Unacceptable risk**
e.g. social scoring
— **Prohibited**

**High risk**
e.g. recruitment, medical devices
— **Permitted** subject to compliance with AI requirements and ex-ante conformity assessment

**'Transparency' risk**
'Impersonation' (bots)
— **Permitted** but subject to information/transparency obligations

**Minimal or no risk**
— **Permitted** with no restrictions

*Not mutually exclusive
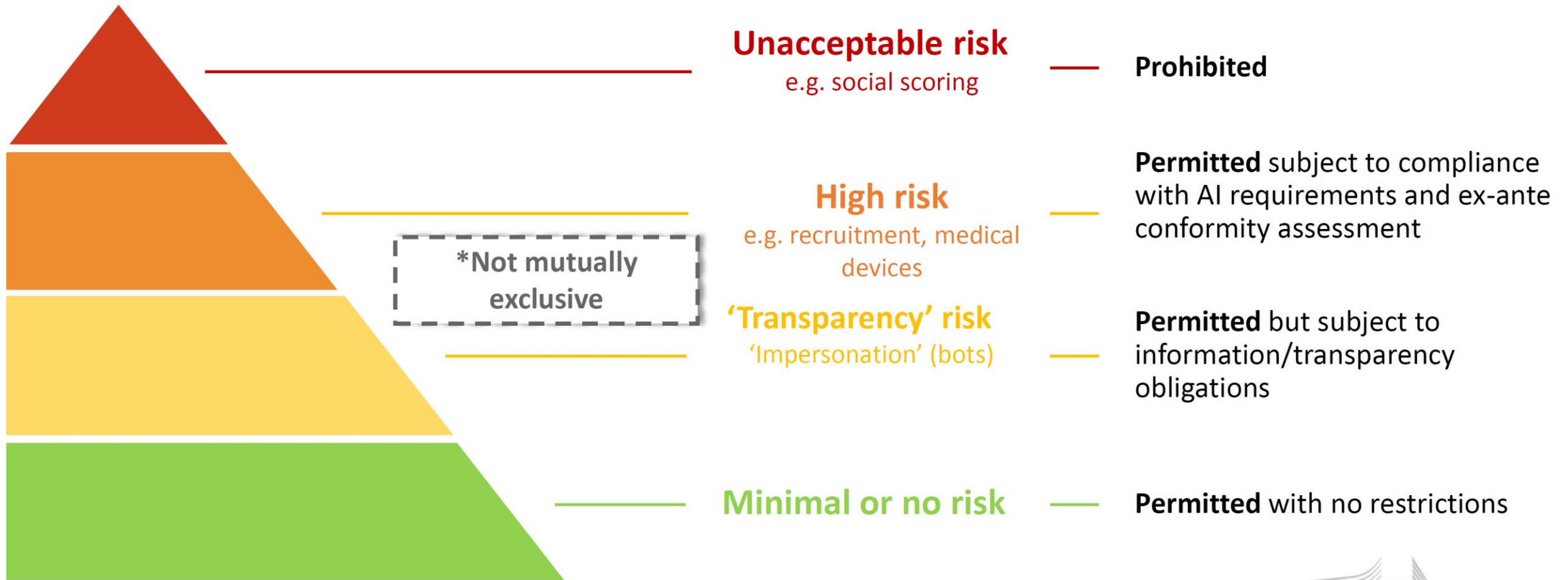
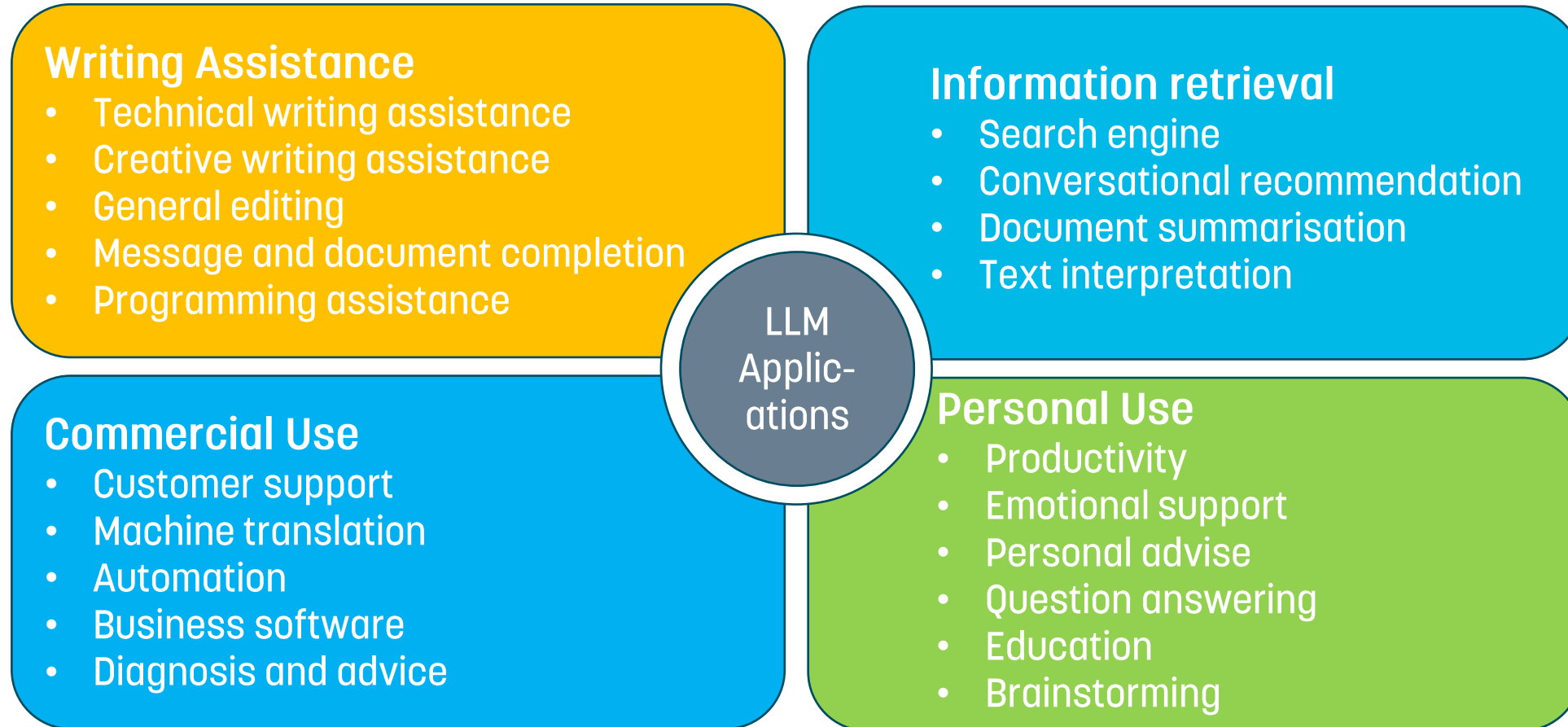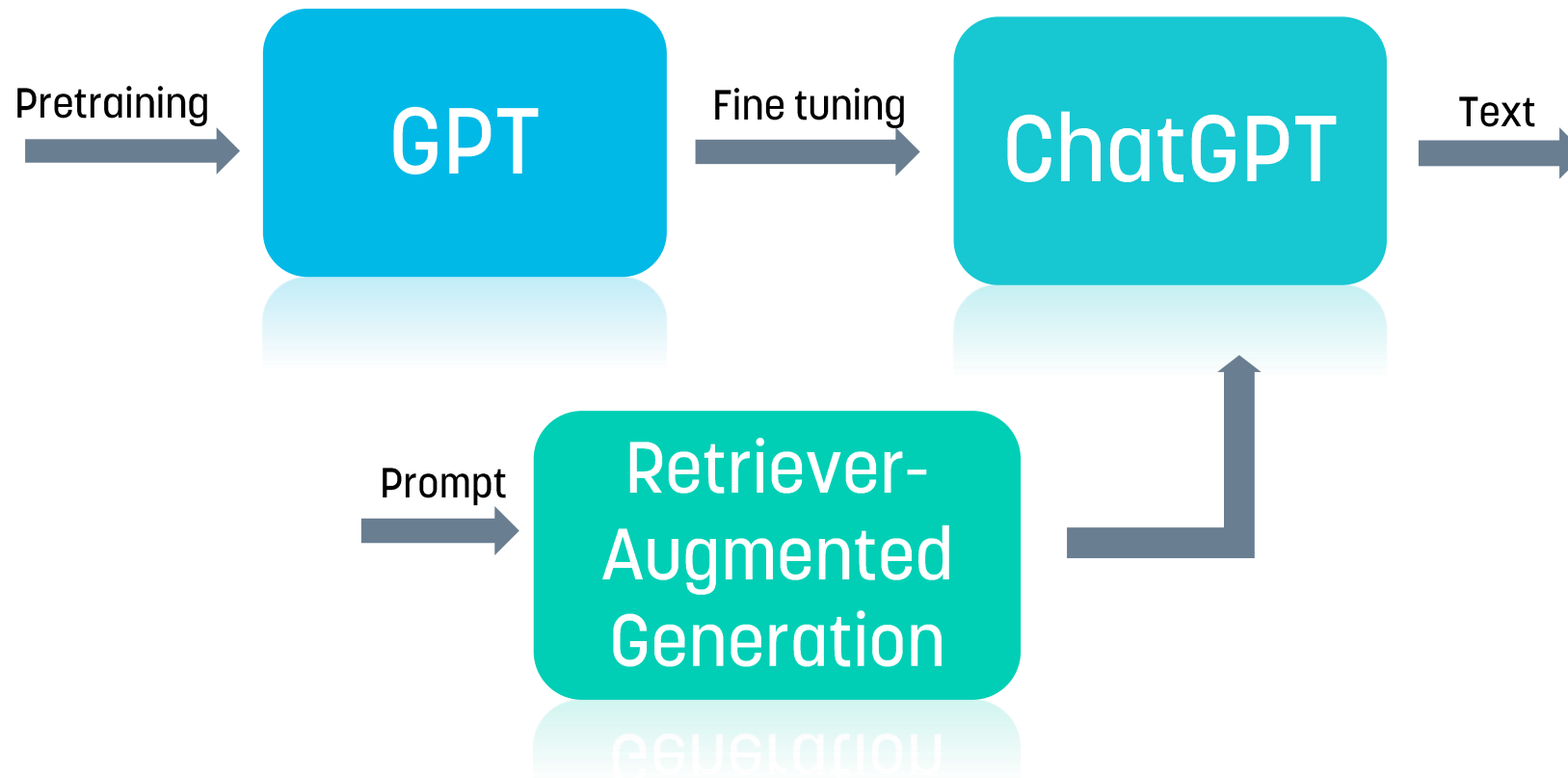*Generated by Dall-E from "photorealistic image of a self-driving car"*

# 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.*

LINKÖPING UNIVERSITY

WASP ED
WALLENBERG AI AND TRANSFORMATIVE TECHNOLOGIES
EDUCATION DEVELOPMENT PROGRAM

# Large Language Model Applications

**Writing Assistance**
- Technical writing assistance
- Creative writing assistance
- General editing
- Message and document completion
- Programming assistance

**Information retrieval**
- Search engine
- Conversational recommendation
- Document summarisation
- Text interpretation

**LLM Applic-ations**

**Commercial Use**
- Customer support
- Machine translation
- Automation
- Business software
- Diagnosis and advice

**Personal Use**
- Productivity
- Emotional support
- Personal advise
- Question answering
- Education
- Brainstorming

TrustLLM

Trustworthy LLMs: A Survey and Guideline for Evaluating Large Language Models' Alignment by Yang Liu Etal, 2023

# 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
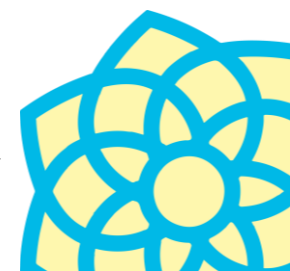
- *but, ChatGPT is still useful!*

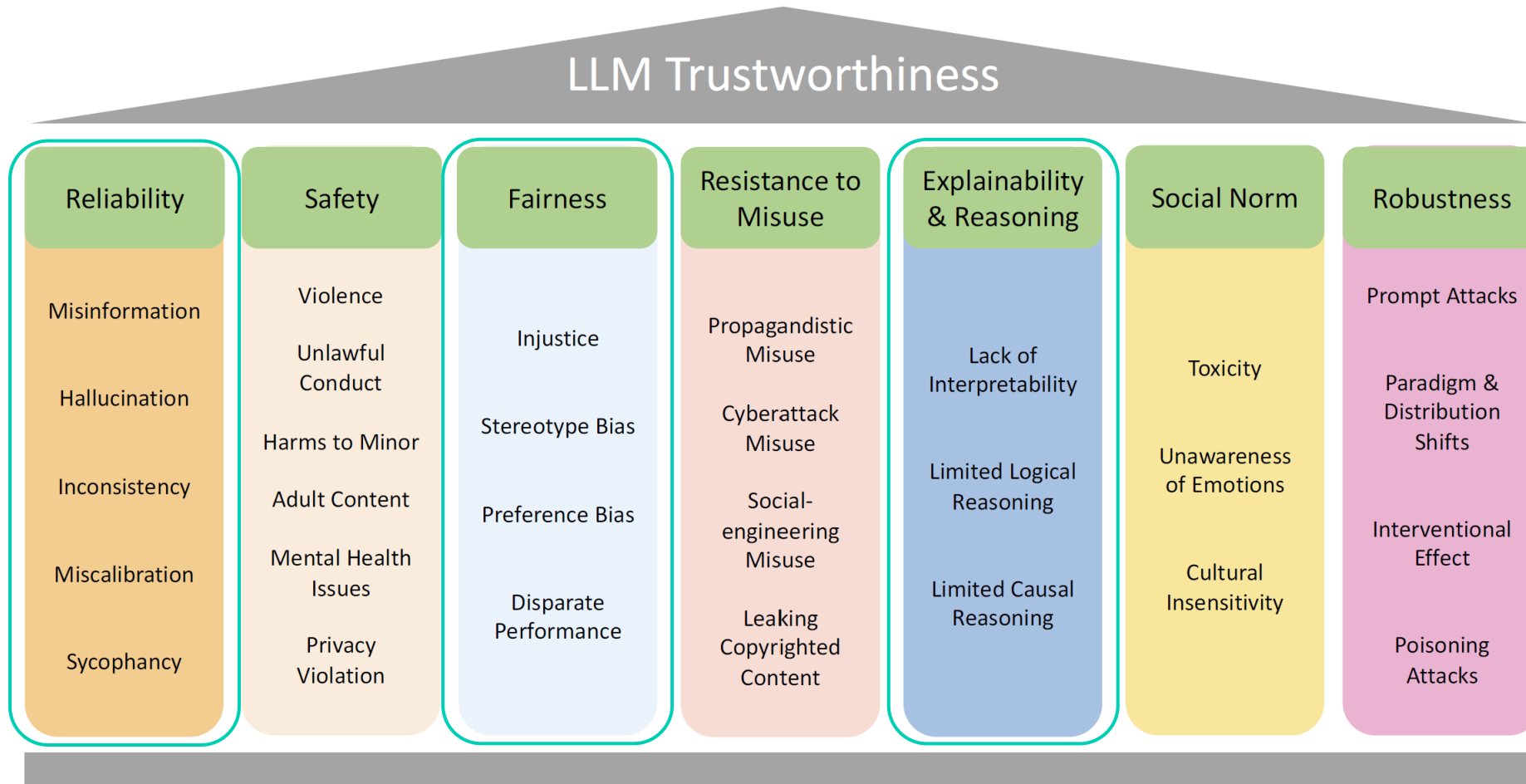# TrustLLM – Trustworthy and Factual LLMs made in Europe

- 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.



https://trustllm.eu/

# LLM Trustworthiness



Trustworthy LLMs: A Survey and Guideline for Evaluating Large Language Models' Alignment by Yang Liu Etal, 2023

# TrustLLM Objectives (1/2)

- **Improving the Factual Correctness of LLMs**
  - New methodology to include structured factual information into LLMs.
  - Large knowledge databases integrated in the process to incorporate structural knowledge (e.g., knowledge graphs, databases, APIs) in the LLM-training
  - Successful benchmarking of the new methodology: error reduction of common NLP tasks related to factual languages.

- **Providing a Framework for Modular and Transferable Multilingual LLM Training**
  - correctness, across A family of 3-5 open source, general purpose and pre-trained LLMs
  - Two models transferred to low-resource target languages.
  - Excellent language modelling score (low perplexity) on 3 low-resource languages with only a small available training corpus.

- **Development of Sustainable and Trustworthy LLMs aligned with European Values**
  - Novel methods of aligning LLMs, while balancing increased robustness, reduced bias, and increased efficiency.
  - Development of a multi-lingual benchmarks suite to measure alignment with European values
  - Publication of a dataset together with methods for fine-tuning LLMs aligned with user expectations (e.g., instructions) or creator intentions (e.g., safe language).
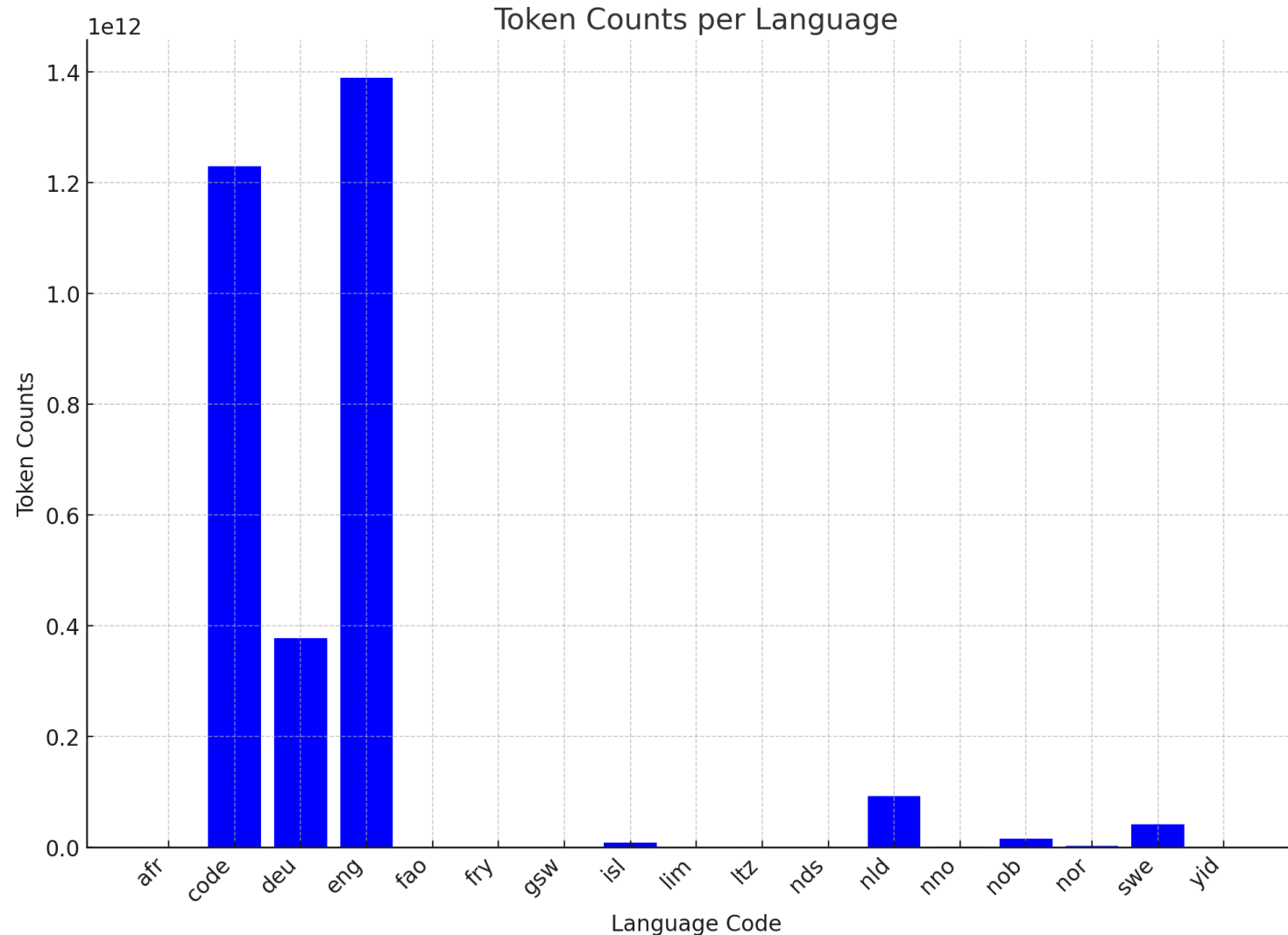
TrustLLM

# TrustLLM Objectives (2/2)

- **Establishing a European Ecosystem for LLMs**
  - **GeTT events: European events on LLM/GPAI with industry & academia**
  - **NeTT events: topical workshops with experts on LLM (research topics, infrastructure, etc.).**
  - **A Concept to establish a sustainable European ecosystem based on the LEAM initiative combining national and European funding and networks structures.**

- **Assessment of LLM by Use Cases Demonstration and Holistic Benchmarking**
  - **Provision of NLP benchmarks in Germanic languages**
  - **Technical realization of use cases based on TrustLLM.**
  - **Applications demonstrating the versatility of TrustLLM solutions/models**

- **Large Scale Data Management for LLM**
  - **Access to text corpora with 1-2 trillion tokens in 6 Germanic languages.**
  - **Efficient data processing pipelines with a throughput of 100 billion tokens/day.**
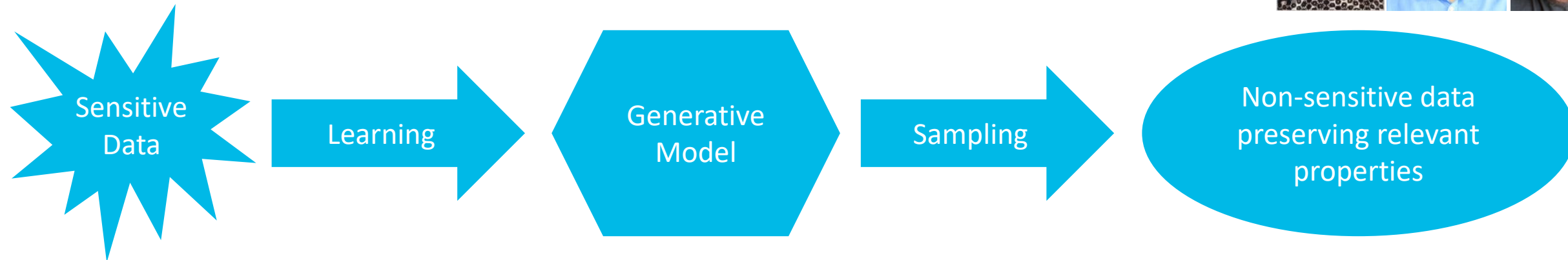  - **Access to a large, distributed storage system with least 2 PB capacity.**

TrustLLM

Funded by
the European Union

# Dataset Composition - TrustLLM



Token Counts per Language

# Privacy-preserving synthetic data generation
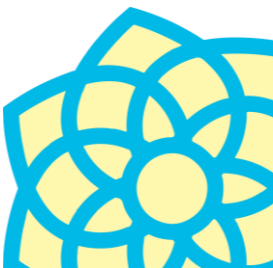## [R. Ramachandranpillai, Md F. Sikder, D. Bergström]



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

# Alignment

"**AI alignment** aims to steer AI systems toward a person's or group's intended goals, preferences, and ethical principles. An AI system is considered *aligned* if it advances the intended objectives. A *misaligned* AI system pursues unintended objectives."

- Recent research shows again and again that it is crucial to have **high-quality instruct fine-tuning data** for reinforcement learning.

- The *Less Is More for Alignment (LIMA)* paper by Zhou et al. (2023) demonstrated fine-tuning on just **1,000 high-quality examples** improved a model's instruction following, highlighting data quality over quantity.

- Furthermore, it is important to have **diverse data** written by a diverse group of people.

TrustLLM

Funded by
the European Union

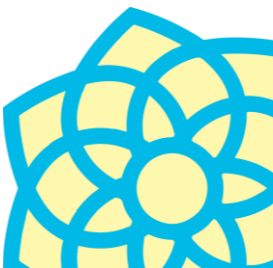# How do we align an LLM?

## 1. Foundation model

Pre-train on vast amounts of data

## 2. Instruction fine-tuning

Learning from task-specific examples

## 3. Preference tuning
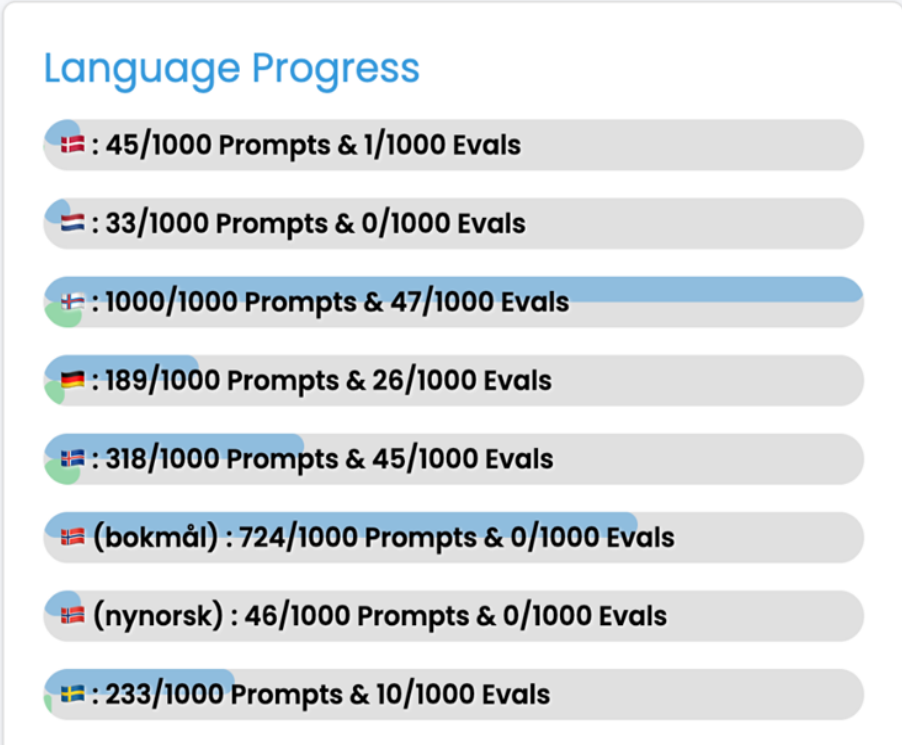
Feedback and reinforcement

TrustLLM

Funded by
the European Union

# Alignment Data Collection

# Our Crowdsourcing Platform

- This is a web app for reformulating the English prompts from the *Open Assistant Dataset* into the following Germanic languages:
  - German, Dutch, Swedish, Norwegian, Danish, Icelandic and Faroese

- We have initially only been annotating within the TrustLLM project.

- You need an invite key in order to join. Contact Annika by email <u>ans72@hi.is</u> or annika@hi.is to get one.

# Our Crowdsourcing Platform

- The goal is to reformulate 1,000 prompts for each language.

- These prompts will be used in a second phase of creating alignment data.

- We will publish the final dataset, open-source.

TrustLLM

Funded by
the European Union

# The tasks

- Prompt reformulation
  - Use an English prompt as inspiration and reformulate it in your own language, making it culturally appropriate and naturally-sounding linguistically.

- Original prompt creation
  - Create your own high-quality prompt in your own language.
  - Not only QA; also summarization, bug fixing, idea generation etc.

- Prompt evaluation
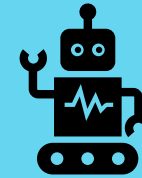  - Evaluate prompts written by other people in your own language using labels and scales.

TrustLLM

Funded by
the European Union

# How can we evaluate LLMs?

TrustLLM

Funded by
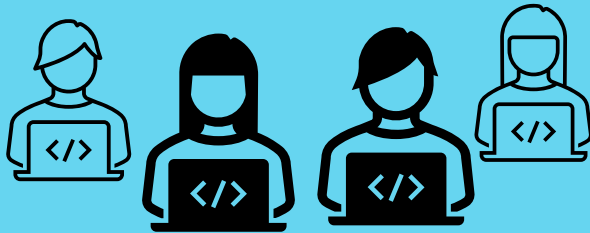the European Union

# Four Main Approaches

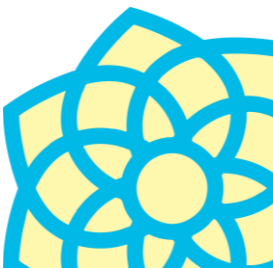**Gut Feeling Approach**

**LLM-as-a-judge**

**Arena Approach**

**Benchmark Approach**

# What is ScandEval?

# ScandEval is a robust multilingual benchmarking framework

ScandEval is a robust multilingual benchmarking framework

TrustLLM

Funded by
the European Union

# Language Model Benchmarking Framework

- Enables evaluation of implicit language <span style="color:orange">understanding</span> and <span style="color:orange">generation</span> capabilities of language models

- Allows evaluation of *both* encoders through finetuning, and decoders through few-shot evaluation

  - It has been shown that few-shot inference of decoder models corresponds exactly to finetuning [1]

  - This thus allows us to compare encoders with decoders directly

[1] von Oswald et al. arXiv preprint arXiv:2309.05858 (2023)
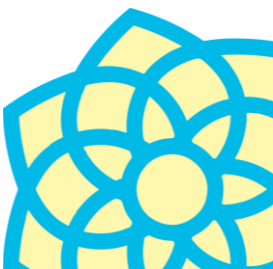
# Language Model Benchmarking Framework

- A large focus of the framework is ease of use

- The framework can simply be installed:

  ```
  $ pip install scandeval[all]
  ```

- Models can easily be evaluated:

  ```
  $ scandeval --model <model-id>
  ```

- Supports models on the Hugging Face Hub, local models and OpenAI models

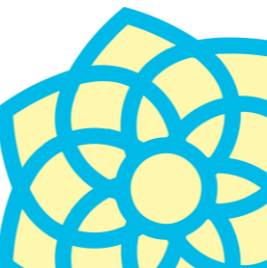# ScandEval is a robust multilingual benchmarking framework

# Evaluation Robustness

- When evaluating models, there are several sources of noise in the evaluation result:

  - The choice of training examples (=few-shot examples when evaluating decoder models)

  - The choice of test examples

  - The stochastic elements (stochastic gradient descent when evaluating encoders, sampling when evaluating decoders)

- The training and test examples are bootstrapped 10 times, yielding a more reliable estimation of the true mean

  - Asymptotically correct by the bootstrap theorem

- We enforce that the stochastic elements are deterministic

TrustLLM

# ScandEval is a robust multilingual benchmarking framework

# Multilingual Evaluations

- Currently, the main Germanic languages are natively supported:
    - Mainland Scandinavian languages (Danish, Swedish, Norwegian)
    - Insular Scandinavian languages (Icelandic, Faroese)
    - German
    - Dutch
    - English
- Aside from including evaluation datasets in these languages, the prompts used when evaluating decoder models are also localised to the given language

# Which Tasks are Included?

TrustLLM

# Tasks in ScandEval

Natural Language Understanding (NLU) Tasks

1. Text classification

2. Linguistic acceptability

3. Extractive question answering

4. Named entity recognition
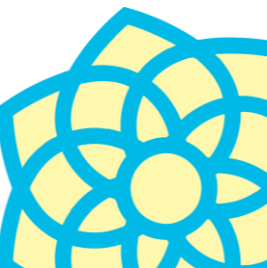
# Tasks in ScandEval
Natural Language Generation (NLG) Tasks

1. Text classification

2. Linguistic acceptability

3. Extractive question answering

4. Named entity recognition

5. Summarisation

6. World knowledge

7. Common-sense reasoning

TrustLLM

Funded by
the European Union

# Online Leaderboards

## scandeval.com

**ScandEval**

ABOUT | DANISH ▼ | SWEDISH ▼ | NORWEGIAN ▼ | ICELANDIC ▼ | FAROESE ▼

NLU LEADERBOARD

NLG LEADERBOARD

# Danish NLG

*Rank Score* computed (roughly) as 1 + number of significant standard deviations from best model, across all datasets

Last updated: 27/05/2024 14:10:07 CET

☐ Include merged models

| Model ID | Parameters | Vocabulary Size | Context | Commercial | Speed | Rank ▼ | |
|---|---|---|---|---|---|---|---|
| gpt-4-0613 (few-shot, val) | unknown | 100 | 8192 | True | 597 ± 197 / 93 ± 33 | 1.09 | |
| syvai/danskgpt-chat-llama3-70b (few-shot, val) | 70554 | 128 | 8192 | True | 1,283 ± 279 / 291 ± 92 | 1.36 | |
| meta-llama/Meta-Llama-3-70B (few-shot, val) | 70554 | 128 | 8192 | True | 312 ± 55 / 177 ± 51 | 1.47 | 6 |
| gpt-3.5-turbo-0613 (few-shot, val) | unknown | 100 | 4094 | True | 921 ± 293 / 113 ± 37 | 1.58 | |

Download as CSV  •  Copy embed HTML

**TrustLLM**

Funded by the European Union

# What are some of the best performing European models?

TrustLLM

Funded by
the European Union

# English ScandEval Rank

Lower is better

| Model | Rank |
|---|---|
| gpt-4-1106-preview | 1.16 |
| gpt-4-0613 | 1.22 |
| meta-llama/Meta-Llama-3-70B | 1.33 |
| gpt-4o-2024-05-13 | 1.36 |
| upstage/SOLAR-10.7B-v1.0 | 1.45 |
| Nexusflow/Starling-LM-7B-beta | 1.48 |
| meta-llama/Llama-2-70b-hf | 1.54 |
| gpt-3.5-turbo-0613 | 1.71 |
| meta-llama/Meta-Llama-3-8B | 1.94 |
| mistralai/Mistral-7B-v0.1 | 1.94 |
| meta-llama/Llama-2-7b-hf | 2.46 |

TrustLLM

Funded by
the European Union

# Dutch ScandEval Rank

Lower is better

| Model | Rank |
|---|---|
| gpt-4-0613 | 1.14 |
| meta-llama/Meta-Llama-3-70B | 1.34 |
| gpt-4-1106-preview | 1.45 |
| gpt-4o-2024-05-13 | 1.54 |
| upstage/SOLAR-10.7B-v1.0 | 1.99 |
| Nexusflow/Starling-LM-7B-beta | 2.05 |
| gpt-3.5-turbo-0613 | 2.07 |
| meta-llama/Llama-2-70b-hf | 2.15 |
| meta-llama/Meta-Llama-3-8B | 2.43 |
| yhavinga/Boreas-7B-chat | 2.52 |
| mistralai/Mistral-7B-v0.1 | 2.77 |
| meta-llama/Llama-2-7b-hf | 3.24 |

TrustLLM

Funded by
the European Union

# German ScandEval Rank

Lower is better

| Model | Rank |
|---|---|
| gpt-4-0613 | 1.18 |
| gpt-4-1106-preview | 1.33 |
| meta-llama/Meta-Llama-3-70B | 1.36 |
| gpt-4o-2024-05-13 | 1.44 |
| upstage/SOLAR-10.7B-v1.0 | 1.57 |
| meta-llama/Llama-2-70b-hf | 1.71 |
| Nexusflow/Starling-LM-7B-beta | 1.88 |
| VAGOsolutions/SauerkrautLM-7b-LaserChat | 1.88 |
| gpt-3.5-turbo-0613 | 1.90 |
| meta-llama/Meta-Llama-3-8B | 2.06 |
| mistralai/Mistral-7B-v0.1 | 2.25 |
| meta-llama/Llama-2-7b-hf | 2.71 |

TrustLLM

Funded by
the European Union

# Danish ScandEval Rank

Lower is better

| Model | Rank |
|---|---|
| gpt-4-0613 | 1.12 |
| gpt-4-1106-preview | 1.20 |
| gpt-4o-2024-05-13 | 1.23 |
| syvai/danskgpt-chat-llama3-70b | 1.36 |
| meta-llama/Meta-Llama-3-70B | 1.46 |
| gpt-3.5-turbo-0613 | 1.58 |
| meta-llama/Llama-2-70b-hf | 1.73 |
| upstage/SOLAR-10.7B-v1.0 | 2.02 |
| Nexusflow/Starling-LM-7B-beta | 2.02 |
| meta-llama/Meta-Llama-3-8B | 2.32 |
| mistralai/Mistral-7B-v0.1 | 2.61 |
| meta-llama/Llama-2-7b-hf | 3.01 |

# Swedish ScandEval Rank

Lower is better

| Model | Rank |
|---|---|
| gpt-4-0613 | 1.10 |
| gpt-4o-2024-05-13 | 1.18 |
| gpt-4-1106-preview | 1.19 |
| meta-llama/Meta-Llama-3-70B | 1.38 |
| gpt-3.5-turbo-0613 | 1.82 |
| meta-llama/Llama-2-70b-hf | 1.85 |
| upstage/SOLAR-10.7B-v1.0 | 2.05 |
| Nexusflow/Starling-LM-7B-beta | 2.18 |
| timpal0l/Llama-3-8B-flashback-v1 | 2.27 |
| meta-llama/Meta-Llama-3-8B | 2.31 |
| mistralai/Mistral-7B-v0.1 | 2.62 |
| meta-llama/Llama-2-7b-hf | 2.90 |

TrustLLM

Funded by
the European Union

# Norwegian ScandEval Rank

Lower is better

| Model | Rank |
|---|---|
| gpt-4-0613 | 1.17 |
| gpt-4-1106-preview | 1.26 |
| gpt-4o-2024-05-13 | 1.31 |
| meta-llama/Meta-Llama-3-70B | 1.45 |
| gpt-3.5-turbo-0613 | 1.99 |
| meta-llama/Llama-2-70b-hf | 2.25 |
| upstage/SOLAR-10.7B-v1.0 | 2.35 |
| Nexusflow/Starling-LM-7B-beta | 2.45 |
| bineric/NorskGPT-Llama3-8b | 2.46 |
| meta-llama/Meta-Llama-3-8B | 2.61 |
| mistralai/Mistral-7B-v0.1 | 3.04 |
| meta-llama/Llama-2-7b-hf | 3.41 |

TrustLLM

Funded by
the European Union

# Icelandic ScandEval Rank

Lower is better



| Model | Rank |
|-------|------|
| gpt-4o-2024-05-13 | 1.17 |
| gpt-4-1106-preview | 1.19 |
| gpt-4-0613 | 1.54 |
| meta-llama/Meta-Llama-3-70B | 2.45 |
| gpt-3.5-turbo-0613 | 3.09 |
| meta-llama/Llama-2-70b-hf | 3.32 |
| upstage/SOLAR-10.7B-v1.0 | 3.34 |
| meta-llama/Meta-Llama-3-8B | 3.40 |
| mhenrichsen/hestenettetLM | 3.64 |
| Nexusflow/Starling-LM-7B-beta | 3.66 |
| mistralai/Mistral-7B-v0.1 | 3.70 |
| meta-llama/Llama-2-7b-hf | 4.04 |

TrustLLM

Funded by
the European Union

# Papers

**ScandEval NLU benchmark for encoders:**

Nielsen, Dan. *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa).* 2023

**ScandEval NLU benchmark for decoders:**

Joint work with Kenneth Enevoldsen (Aarhus University) and Peter Schneider-Kamp (University of Southern Denmark). Preprint out soon.

TrustLLM

Funded by
the European Union

# TrustLLM: Project Concept

# Trends

- Will the LLMs continue to grow?
    - Smaller and more effective models
    - More specfalized models
- Multi-modal models, initially text, image and sound
- Generalize from language models to AI models, world models grounded in reality
- Neurosymbolic AI
- Testing and evaluation
- Trustworthy Human-Centered AI

TrustLLM

Funded by
the European Union

# 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.

https://trustllm.eu/