# Agentic AI
## and the Architecture of Trust

*About Systems That Act and
Humans Who Remain Responsible*

**Dr. Kirsten MF Weber**
Jülich Supercomputing Centre

**Dr. Alexander Ebbes**
Xyna.AI

**OpenRheinMain 2026**

# AI for Fully Autonomous Weapons?!

Sam Altman,
OpenAI

Pete Hegseth,
Department of War

Dario Amodei,
Anthropic

It doesn't have to be war in some far-away country.
It can start right on your desk.

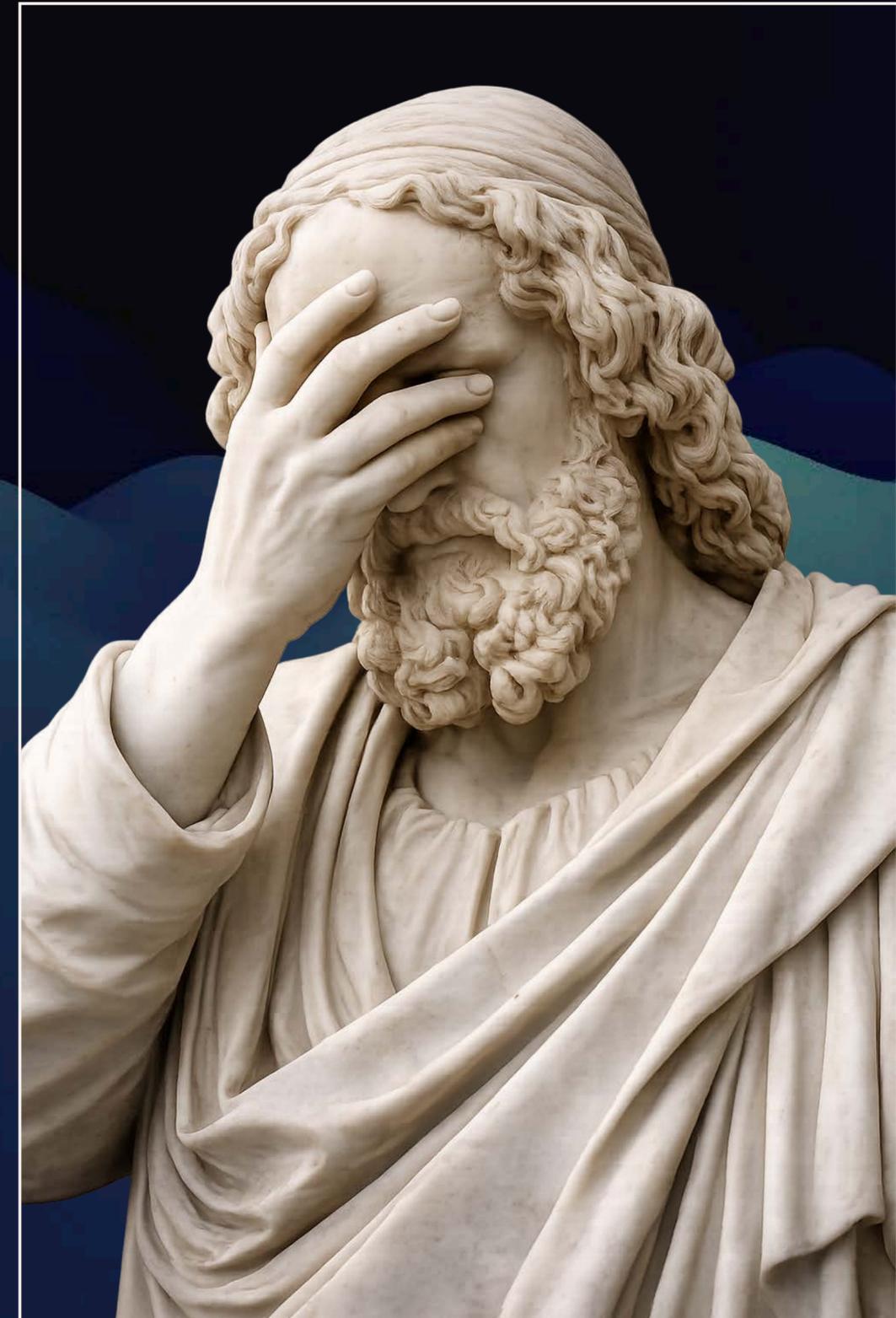# 2022-11-30

**The "ChatGPT Moment" in AI history**

1 million users in 5 days,
100 million after 2 months

# AI Milestone Timeline

- 2022: Research Preview von ChatGPT

- 2023: Model explosion.
  Everyone talks embeddings, tokens, attention.
  The P in GPT creates an ecosystem.
  But the training corpus is still static.

- 2024: Internet grounding.
  Tools and APIs enable search. Recency, specialization, customization.

- 2025: Agentic AI.
  Action instead of dialogue.
  Writing and executing APIs.
  Expanding reach. Delegated execution. Automation.

- 2026: Trust becomes unavoidable.
  More and more is possible.
  But what do we really want to delegate?

# About the Term Trust

- The question of trust is not new.

- Even early users quickly encountered hallucinations.
  AI sometimes just makes things up.

- This is not simply a software bug.
  Hallucination is intrinsic to probabilistic LLMs.
  We understand why it happens —
  but we cannot eliminate it without losing capability.

- The real difference is about impact.
  Misinterpreting a poem is harmless.
  Letting an agent trade your portfolio is not.

# Vectors of Harmfulness in AI

- **Hallucination** — AI can fabricate plausible falsehoods

- **Bias** — Asymmetries in training data. Policy-shaped responses.

- **Curation side-effects** —Reducing ambiguity can create content imbalance

- **Data correctness** — Epistemic constraints. Language models approximate knowledge.

- **Reasoning correctness** — What looks like cognition is just high-quality imitation

Psychology: "Fluency-Induced Credibility Bias" — Confident delivery increases perceived credibility. A foundation of rhetoric, marketing, and propaganda.

# "The Illusion of Thinking"

- Controlled **reasoning benchmarks.**
  Tasks like: Tower of Hanoi, River-Crossing puzzles. Complexity can be scaled precisely

- State-of-the-art reasoning models tested.
  ChatGPT O1/O3, Claude Sonnet Thinking, DeepSeek-R1, Gemini Thinking

- **Result**:
  At moderate complexity: models perform well
  — within distribution.
  Beyond that: reasoning collapses.
  Not gradual degradation. Structural failure

- **Implication**:
  Language models only simulate reasoning.
  They do not perform deduction

# „Is Reasoning a Mirage?"

- Object of study: Chain-of-Thought reasoning behaviour in LLMs.

- **Hypothesis**:
Apparent reasoning is largely determined by the distribution of reasoning trajectories already present in the training corpus. Models extend token streams — not logical states.

- **Finding**: Hypothesis confimed.
"… models conditionally generate reasoning trajectories that approximate those observed during training".

- It is about pattern matching and extrapolation. No deduction or thinking. "A brittle mirage".
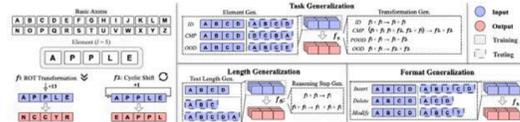
# Hallucination

- Not a software bug. A structural phenomenon.

  - 1986: Already documented in early generative imaging research (e.g. super-resolution from low-resolution inputs). Originally a positive term.

  - Earlier terminology in LLM research → Confabulation. Borrowed from neuropsychology: Patients subconsciously fill memory gaps with invented but subjectively convincing content.
  No intent to deceive.

  - With ChatGPT the catchy term "hallucination" respawned in public discourse.

  - "Hallucination" implies distorted perception. Machines do not perceive.
  "Confabulation" may be the more precise concept.

# Hallucination, cont'd

1. Architectural reality:
Hallucination emerges from the probabilistic nature of transformer models.
"Attention Is All You Need" by Vaswani et al. — read that!

2. Evaluation asymmetry:
Training objectives and benchmarks reward fluency and surface correctness more than epistemic reliability.
This can stabilize hallucination behaviour.

Quellen:

- zu 1: https://arxiv.org/abs/1706.03762

- zu 2: https://arxiv.org/abs/2512.19920

# Trust Is a Trade-off

- A model that can never be wrong must never speculate
  — it would be a database or an expert system,
  not a language model.

- A model that generates creatively must be allowed to be wrong
  occasionally
  — this is a design principle, not a failure.

- Optimization under conflicting objectives always implies a trade-off
  — gain in one dimension entails loss in another.

- You cannot demand "fundamentally reliable AI" without
  understanding this trade-off
  — you would be demanding a contradiction.

- There is no purely technical solution
  — only conscious, transparent choices about trade-offs.

# Approaches to Trust

Spoiler: Trust is not a property of the model —
it is a property of the architecture

# RAG: Retrieval-Augmented Generation

- Even before ChatGPT, Meta addresses the problem of **knowledge-intensive tasks**

  - Dense factual knowledge is only weakly encoded in model weights.

  - Pre-training is expensive → model knowledge is always outdated

  - Idea: inject fresh, exact knowledge at inference time

  - External embedding vector database → knowledge store

  - Major progress — but new failure modes emerge contradictions, chunking errors, lost-in-the-middle, multi-hop limits, latency, etc.

  - Reduces hallucinations — a meaningful step towards trust

Darren Edge[1†]    Ha Trinh[1†]    Newman Cheng[2]    Joshua Bradley[2]    Alex Chao[3]

Apurva Mody[3]    Steven Truitt[2]    Dasha Metropolitansky[1]    Robert Osazuwa Ness[1]

Jonathan Larson[1]

[1]Microsoft Research
[2]Microsoft Strategic Missions and Technologies
[3]Microsoft Office of the CTO

{daedge,trinhha,newmancheng,joshbradley,achao,moapurva,
steventruitt,dasham,robertness,jolarso}@microsoft.com

[†]These authors contributed equally to this work

**Abstract**

The use of retrieval-augmented generation (RAG) to retrieve relevant information from an external knowledge source enables large language models (LLMs) to answer questions over private and/or previously unseen document collections. However, RAG fails on global questions directed at an entire text corpus, such as "What are the main themes in the dataset?", since this is inherently a query-focused summarization (QFS) task, rather than an explicit retrieval task. Prior QFS methods, meanwhile, do not scale to the quantities of text indexed by typical RAG systems. To combine the strengths of these contrasting methods, we propose *GraphRAG*, a graph-based approach to question answering over private text corpora that scales with both the generality of user questions and the quantity of source text. Our approach uses an LLM to build a graph index in two stages: first, to derive an entity knowledge graph from the source documents, then to pre-generate community summaries for all groups of closely related entities. Given a question, each community summary is used to generate a partial response, before all partial responses are again summarized in a final response to the user. For a class of global sensemaking questions over datasets in the 1 million token range, we show that GraphRAG leads to substantial improvements over a conventional RAG baseline for both the comprehensiveness and diversity of generated answers.

**1 Introduction**

Retrieval augmented generation (RAG) (Lewis et al., 2020) is an established approach to using LLMs to answer queries based on data that is too large to contain in a language model's *context window*, meaning the maximum number of *tokens* (units of text) that can be processed by the LLM at once (Kuratov et al., 2024; Liu et al., 2023). In the canonical RAG setup, the system has access to a large external corpus of text records and retrieves a subset of records that are individually relevant to the query and collectively small enough to fit into the context window of the LLM. The LLM then

# Knowledge Graph:
# KG-RAG and GraphRAG

- Special form of RAG, effective for structured information

  - Entities and their relations are encoded explicitly

  - Graph $G=(V,E)$ of vertices and edges

  - Key advantage: consistency is no longer emergent — it is engineered

- Knowledge Graphs are older that RAG or LLM

  - 1960s to 1970s AI research: Semantic Networks

  - Since 2012: Google's Knowledge Graph marked the end of matrix-based internet search engines

- New bottleneck: Excellent results — but expensive.

  - Domain-bound, hard to scale and generalize.

  - Architectural mismatch with GPUs (non-SIMD, irregular memory access)

# The World Model of Yann LeCun

- Yann LeCun, Turing Award 2018, etc.

  - Professor at New York University, Inventor of CNNs, Head of AI at Meta for 12 years, now founding his own startup AMI Labs (Advanced Machine Intelligence)

- "… LLMs are a dead end (when it comes to superintelligence)"

- LeCuns Approach: **World Models**

  - Learned representations of physical reality

  - So very complete and continuous that they can predict consequences of actions

  - Intelligence = **grounded** prediction + planning

  - Not an extension of LLMs → A different paradigm

- Language models predict text. World models predict reality.



A Path Towards Autonomous Machine Intelligence
Version 0.9.2, 2022-06-27

Yann LeCun
Courant Institute of Mathematical Sciences, New York University yann@cs.nyu.edu
Meta - Fundamental AI Research yann@fb.com

June 27, 2022

**Abstract**

How could machines learn as efficiently as humans and animals? How could machines learn to reason and plan? How could machines learn representations of percepts and action plans at multiple levels of abstraction, enabling them to reason, predict, and plan at multiple time horizons? This position paper proposes an architecture and training paradigms with which to construct autonomous intelligent agents. It combines concepts such as configurable predictive world model, behavior driven through intrinsic motivation, and hierarchical joint embedding architectures trained with self-supervised learning.

**Keywords:** Artificial Intelligence, Machine Common Sense, Cognitive Architecture, Deep Learning, Self-Supervised Learning, Energy-Based Model, World Models, Joint Embedding Architecture, Intrinsic Motivation.

**1 Prologue**

This document is not a technical nor scholarly paper in the traditional sense, but a position paper expressing my vision for a path towards intelligent machines that learn more like animals and humans, that can reason and plan, and whose behavior is driven by intrinsic objectives, rather than by hard-wired programs, external supervision, or external rewards. Many ideas described in this paper (almost all of them) have been formulated by many authors in various contexts in various form. The present piece does not claim priority for any of them but presents a proposal for how to assemble them into a consistent whole. In particular, the piece pinpoints the challenges ahead. It also lists a number of avenues that are likely or unlikely to succeed.

The text is written with as little jargon as possible, and using as little mathematical prior knowledge as possible, so as to appeal to readers with a wide variety of backgrounds including neuroscience, cognitive science, and philosophy, in addition to machine learning, robotics, and other fields of engineering. I hope that this piece will help contextualize some of the research in AI whose relevance is sometimes difficult to see.



Figure 2: A system architecture for autonomous intelligence. All modules in this model are assumed to be "differentiable", in that a module feeding into another one (through an arrow connecting them) can get gradient estimates of the cost's scalar output with respect to its own output.
The **configurator** module takes inputs (not represented for clarity) from all other modules and configures them to perform the task at hand.
The **perception** module estimates the current state of the world.
The **world model** module predicts possible future world states as a function of imagined actions sequences proposed by the actor.
The **cost** module computes a single scalar output called "energy" that measures the level of discomfort of the agent. It is composed of two sub-modules, the intrinsic cost, which is immutable (not trainable) and computes the immediate energy of the current state (pain, pleasure, hunger, etc), and the critic, a trainable module that predicts future values of the intrinsic cost.
The **short-term memory** module keeps track of the current and predicted world states and associated intrinsic costs.
The **actor** module computes proposals for action sequences. The world model and the critic compute the possible resulting outcomes. The actor can find an optimal action sequence that minimizes the estimated future cost, and output the first action in the optimal sequence.
See Section 3 for details.

# Trust-Engineering along the Vector of Harmfulness

- Bias → remove the asymmetries in objective functions

- Curation side-effects → restore ambiguity and balance

- Factuality → enforce epistemic constraints

"Problems cannot be solved with the same thinking that created them", often attributed to Einstein

**But this time — it's exactly where you have to start**!

# Trust is not a feature.
# It is a systems property.

# MCP

The leading agentic protocol places the
Human-in-the-Loop
as an core element of architecture.

**Die Welt ist voller Menschen, die unter den Folgen ihres ungelebten Lebens leiden**.

Sie werden verbittert, kritisch oder unnachgiebig, nicht, weil die Welt zu grausam zu Ihnen ist, sondern weil sie **ihre inneren Möglichkeiten verraten** haben.

Sie sehen im Außen Feinde, die eigentlich in ihrem eigenen Inneren als **ungenutzte Talente** und unterdrückte Wunsche schlummern.

Erst wenn der Mensch beginnt, sein **eigenes Licht nicht mehr unter den Scheffel zu stellen** und die Verantwortung für sein eigenes Werden übernimmt, löst sich die Bitterkeit auf und macht dem Frieden Platz.



Carl Gustav Jung,
1875 - 1961,
founder of the school of analytical psychology