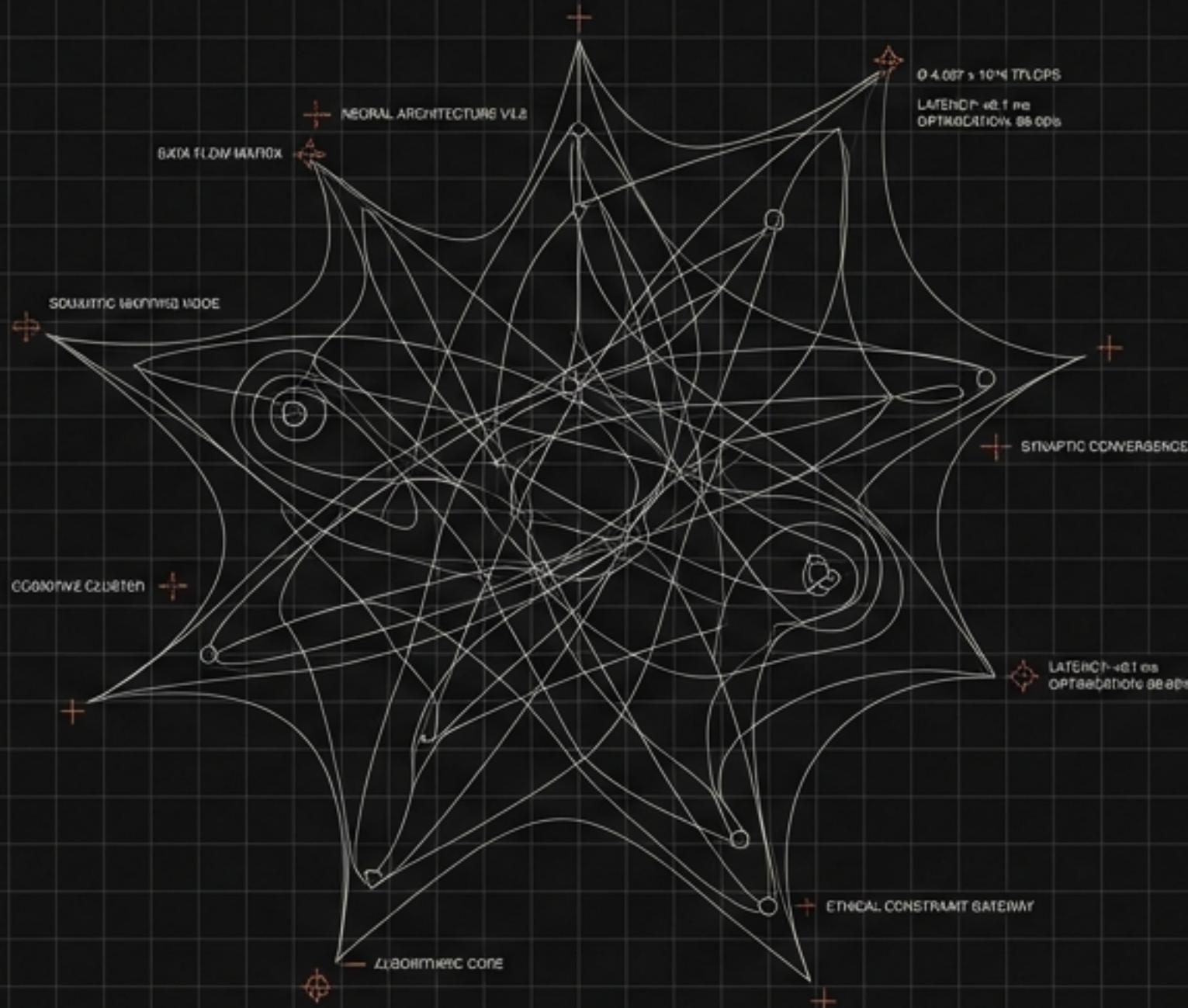


SYNTHETIC MINDS

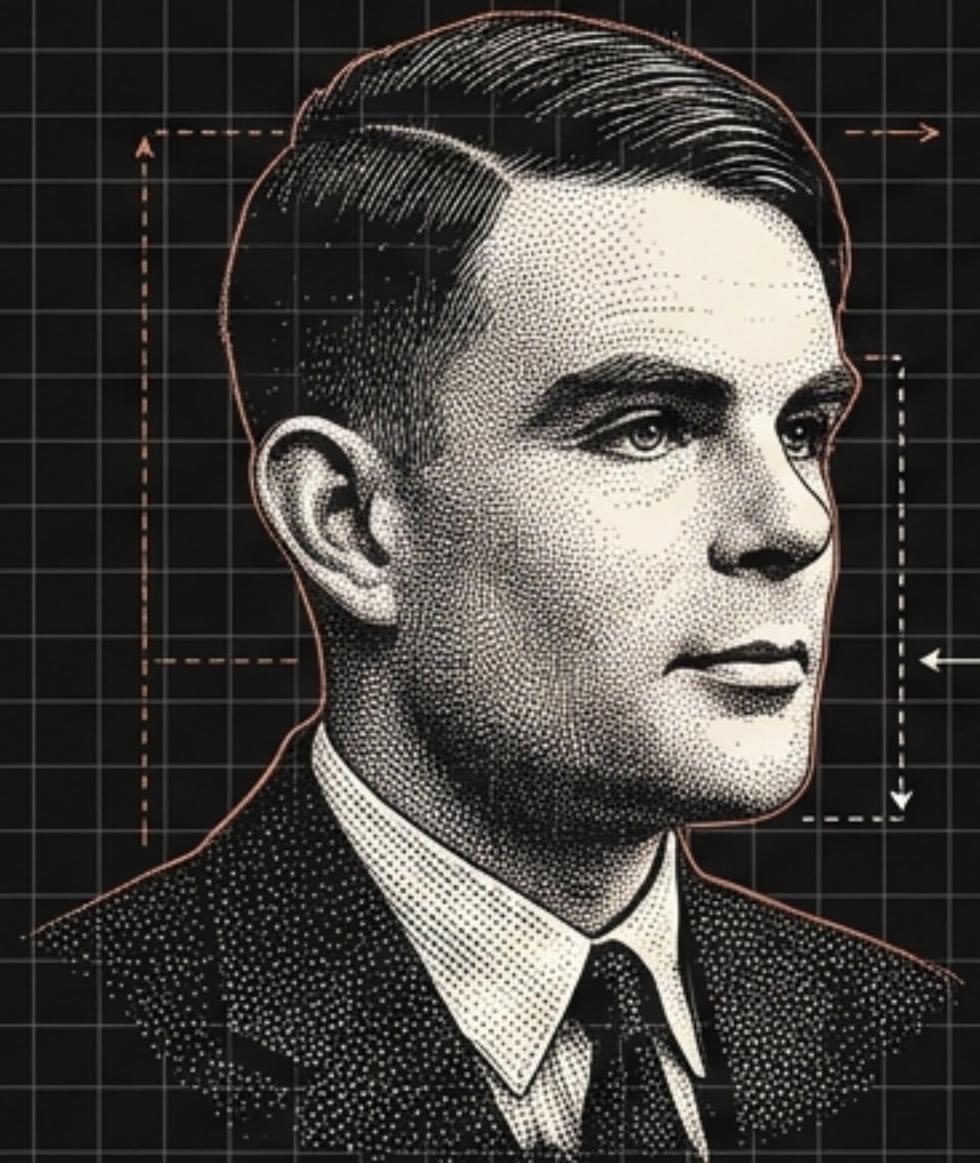
The Mechanics, Applications, and Ethics of Artificial Intelligence.



Defining the Undefined

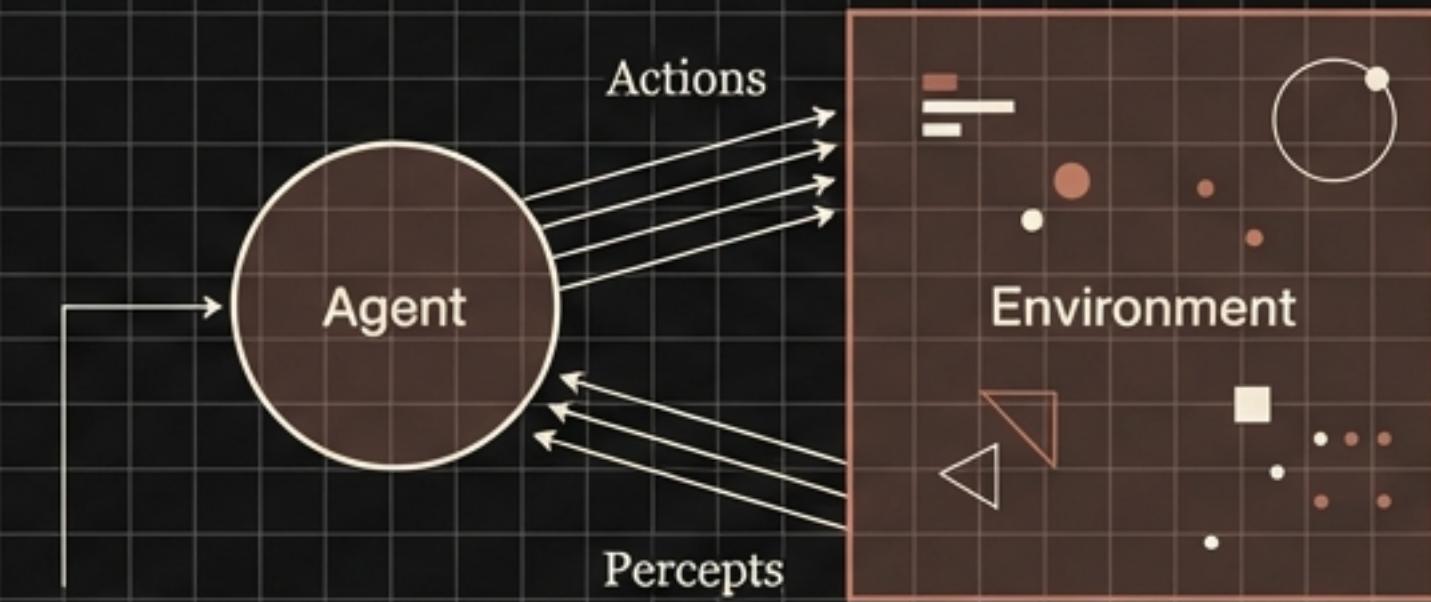
'Intelligence' is an elusive concept. In computer science, it is defined not by internal consciousness, but by external capability.

View A: The Simulator (Turing)



Focus: The Turing Test (1950). Can a machine simulate human conversation well enough to fool a judge? This prioritizes mimicking human behavior.

View B: The Rational Agent (Standard Model)



Focus: Outcome. An agent perceives its environment and takes actions to maximize the chances of achieving a defined goal. It acts to achieve the best result, regardless of whether it 'thinks' like a human.

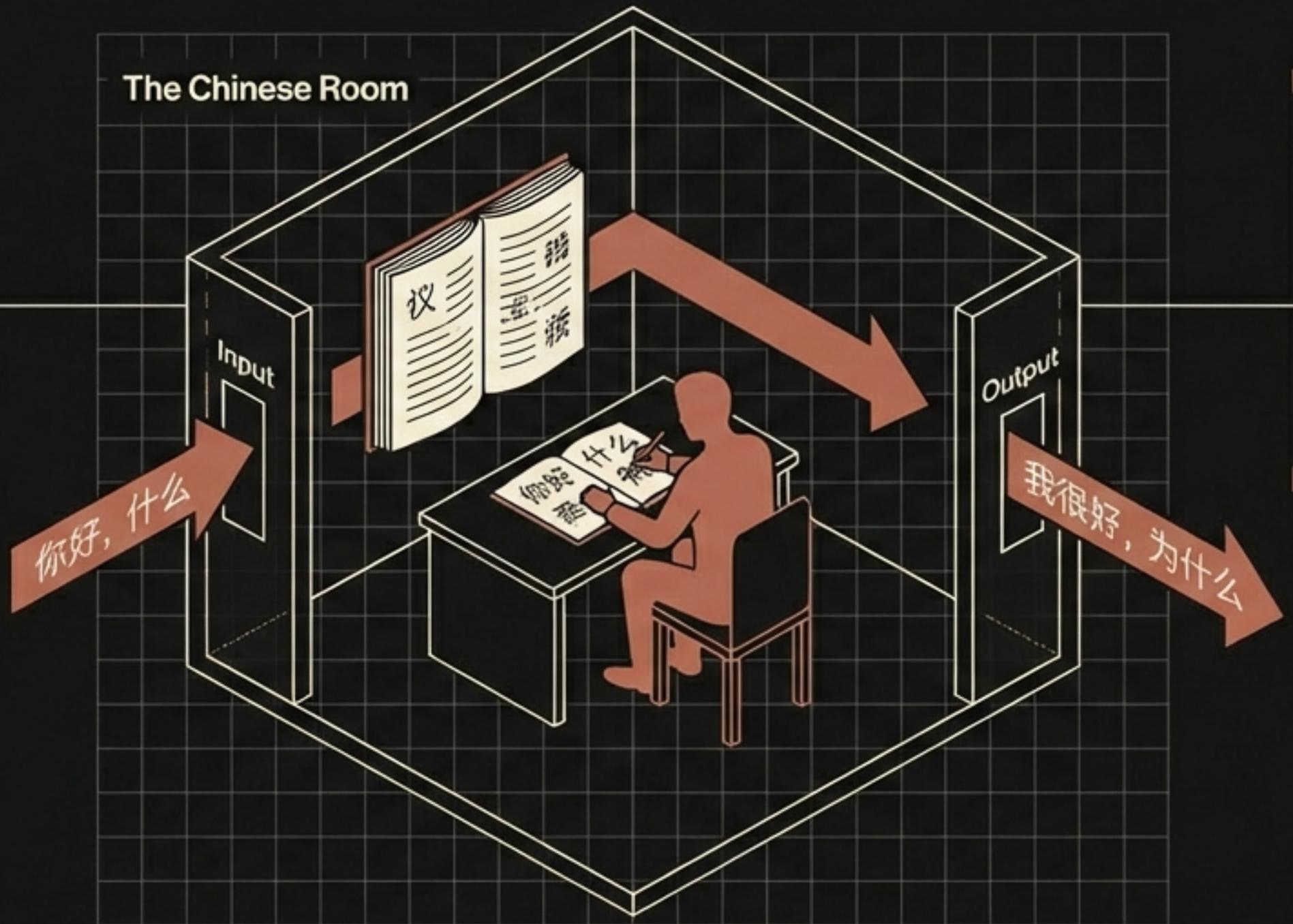
The AI Effect: Once a machine can perform a task (like Chess or OCR), we stop calling it AI and just call it software.

Simulation vs. Understanding

Strong AI (Computationalism):

The belief that a programmed computer literally has a mind and understands cognitive states.

The Chinese Room



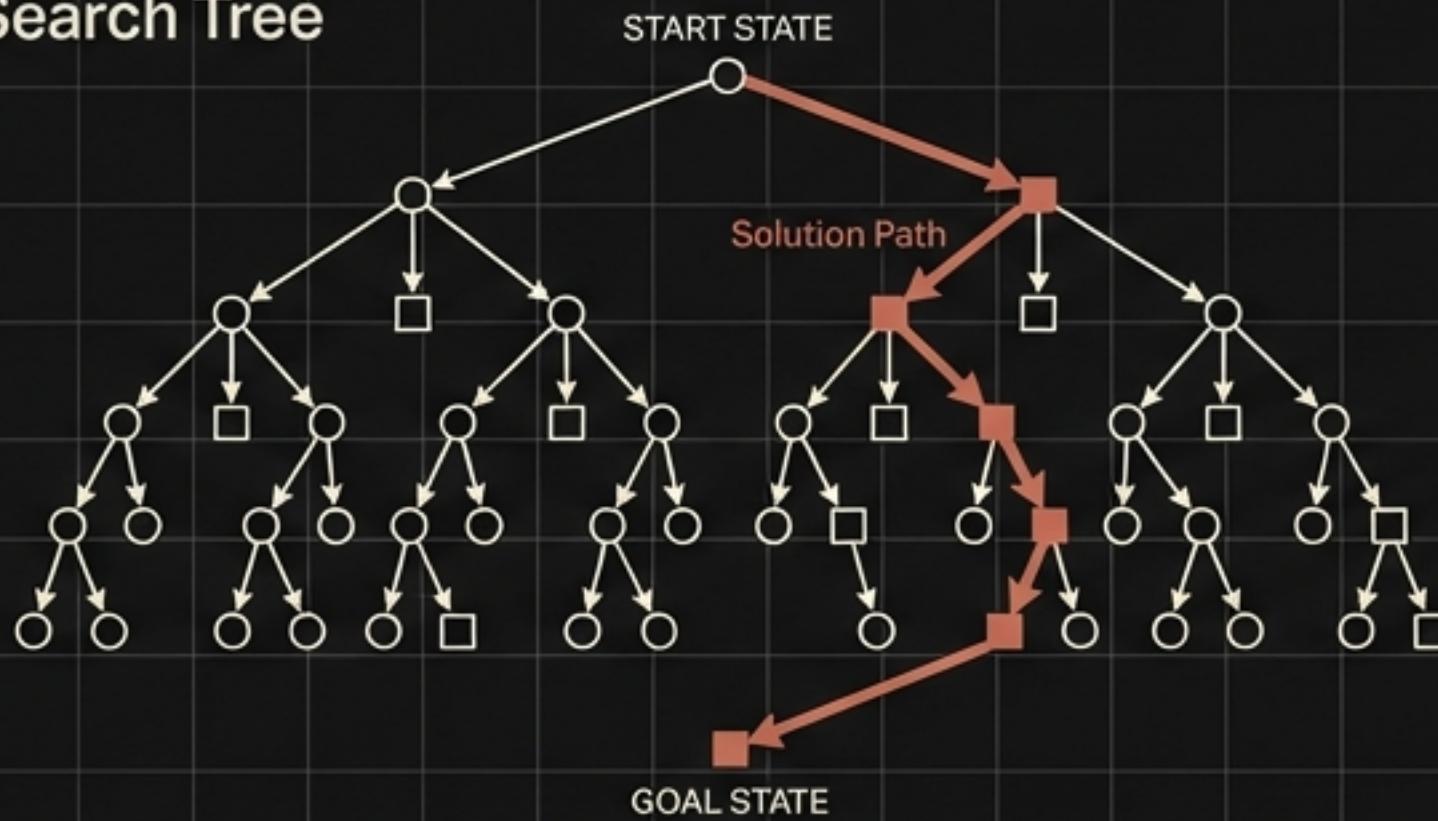
Weak AI:

The pragmatic view that machines act *as if* they are intelligent, solving problems without necessarily possessing consciousness or a 'mind'.

Modern AI research largely sidesteps the 'Hard Problem of Consciousness' (what it feels like to be intelligent) to focus on the engineering problem of competence.

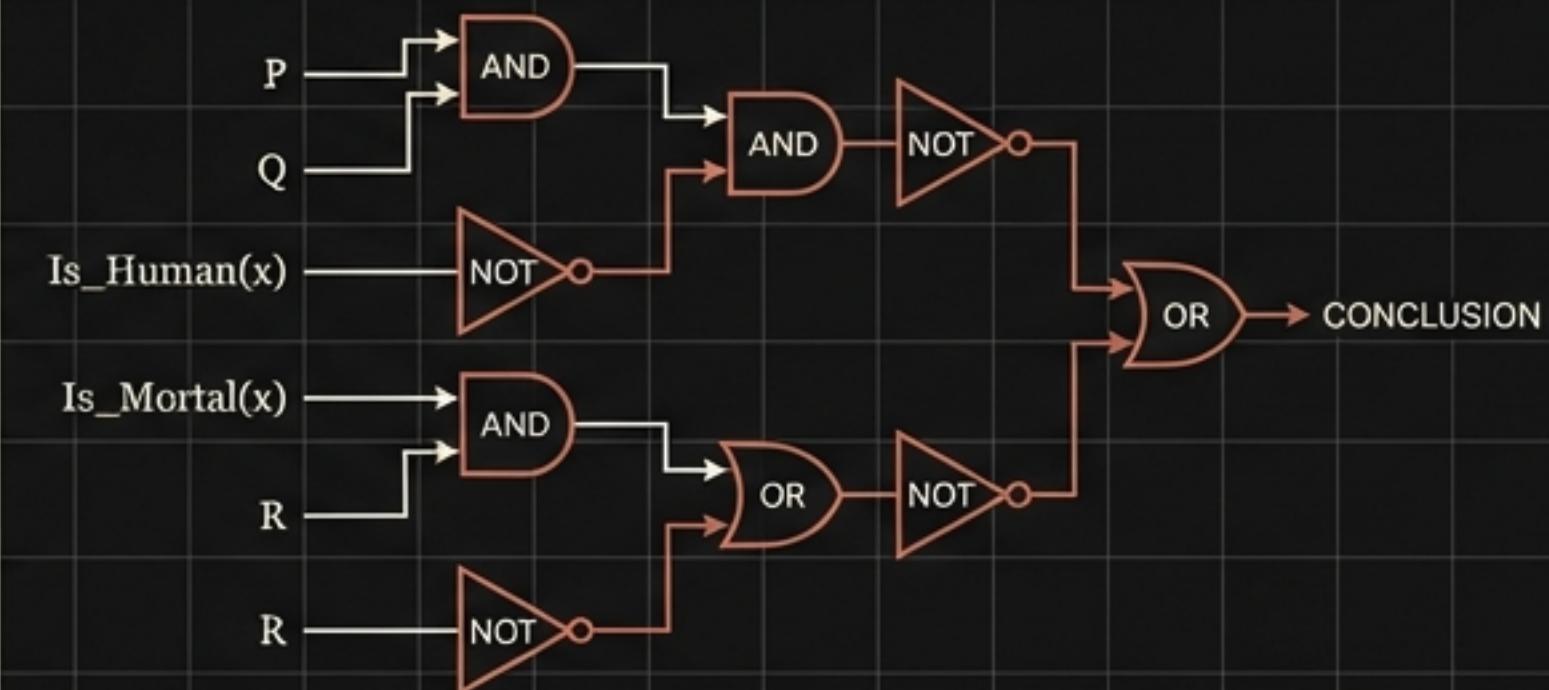
The Old Engine: Logic and Search

Search Tree



State Space Search: Looking through a tree of possible states (e.g., Chess moves) to find a goal.

Symbolic Logic



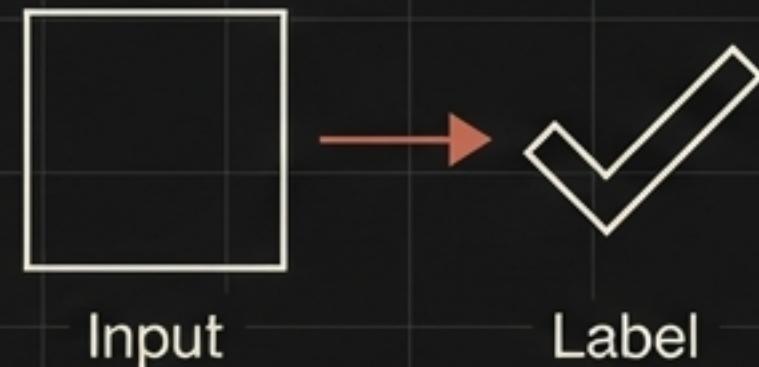
Symbolic AI (GOFAI): Using high-level symbols and formal logic to represent the world.

The Limitation: Moravec's Paradox. High-level reasoning (algebra) is easy for computers; low-level instinct (recognizing a face) is incredibly hard. Symbolic AI failed to capture the messiness of the real world.

The Pivot to Learning

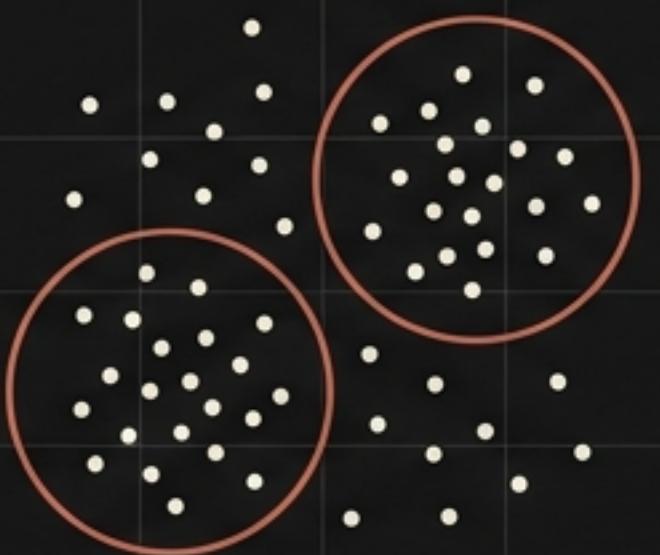
Machine Learning is the study of programs that improve their performance on a task automatically through data.

Supervised Learning



The model learns from labeled data (Input + Expected Answer). Used for classification and regression.

Unsupervised Learning



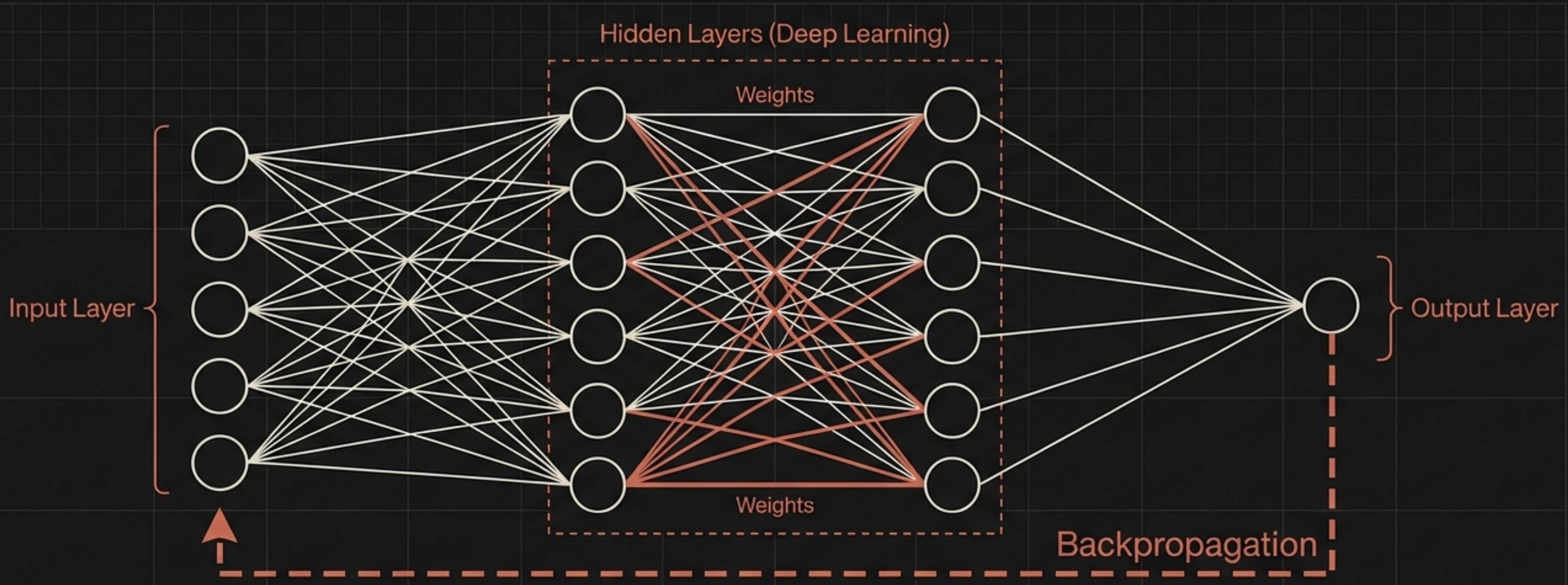
The model finds patterns in unlabeled data. Used for clustering and finding hidden structures.

Reinforcement Learning



The agent learns through trial and error, receiving rewards for good actions and punishment for bad ones.

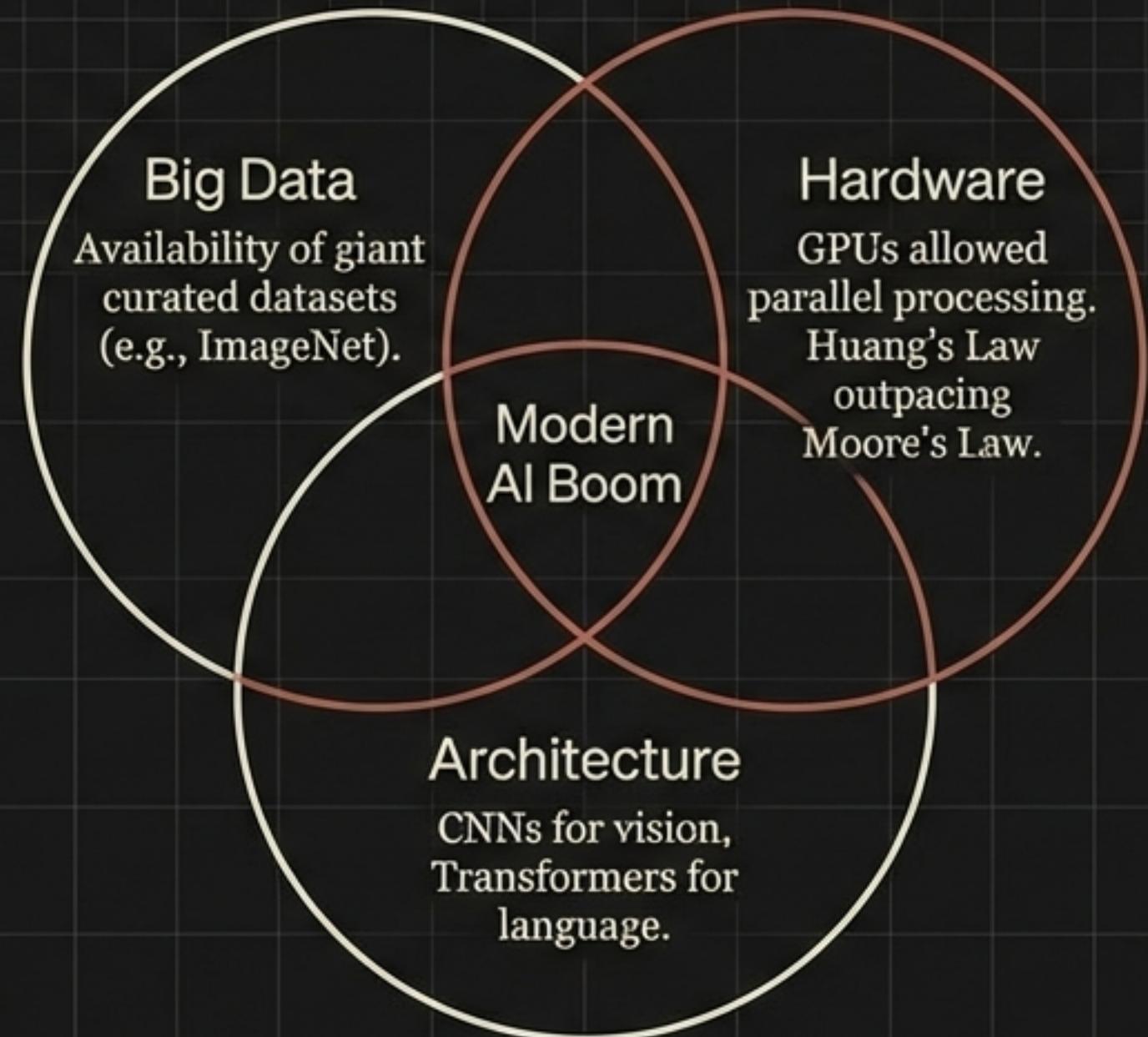
The Anatomy of a Synthetic Neuron



The Learning Mechanism: Backpropagation. The network compares its output to the correct answer, calculates the error, and propagates it backward, adjusting the weights to reduce the error next time. This is mathematically minimizing the loss function via Gradient Descent.

The Deep Learning Explosion (2012–Present)

Deep learning outperforms previous techniques due to a convergence of three factors:

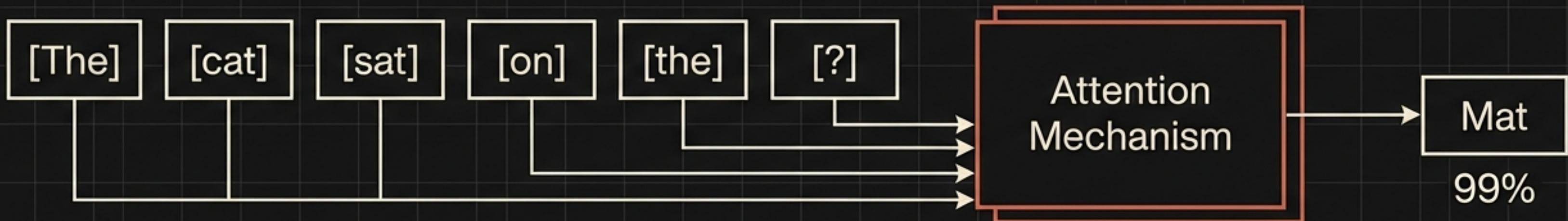


2012: AlexNet wins ImageNet

2017: Transformer Paper

2022: ChatGPT

The Transformer and the Token



Prediction Engines:

Generative Pre-trained Transformers (GPT) are trained on vast corpora to predict the next "token" (word/sub-word) in a sequence.

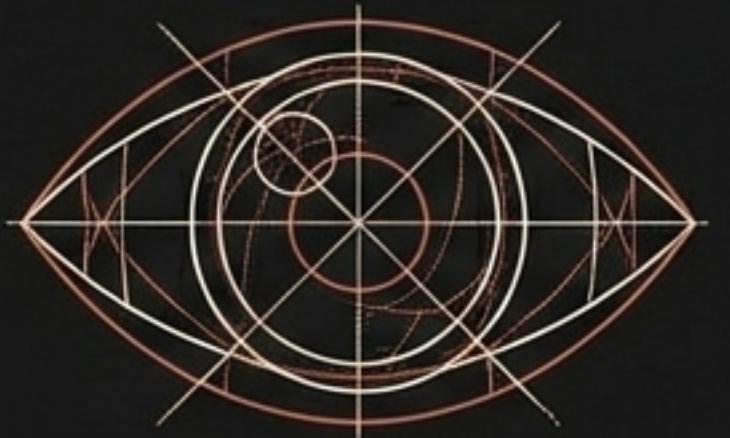
Attention Mechanism:

Allows the model to weigh the relevance of different words in a sentence regardless of distance, capturing context better than previous architectures.

The Risk: Hallucination.

Because they are probabilistic predictors, not fact databases, they can confidently generate plausible but false information.

Cognitive Capabilities



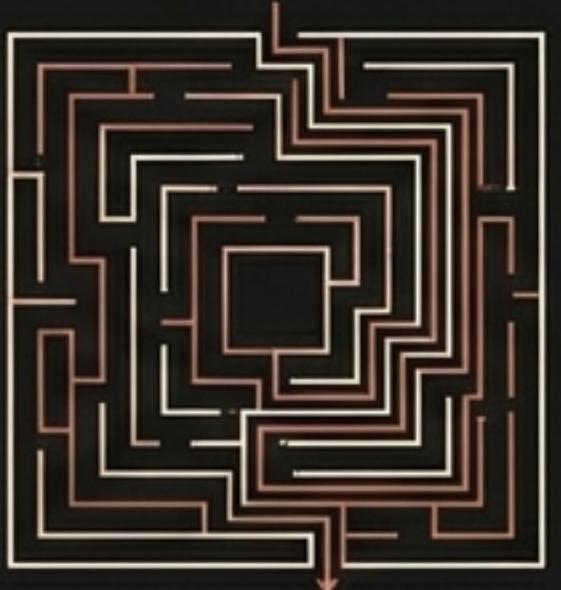
Perception

Computer Vision (Object recognition), Speech Recognition. Converting physical sensors to digital data.



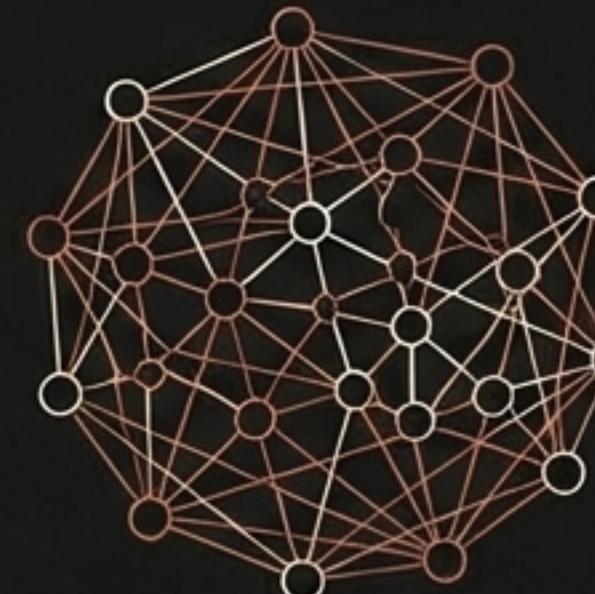
Natural Language Processing (NLP)

Translation, Question Answering, Summarization. Bridging the gap between human syntax and machine code.



Reasoning & Planning

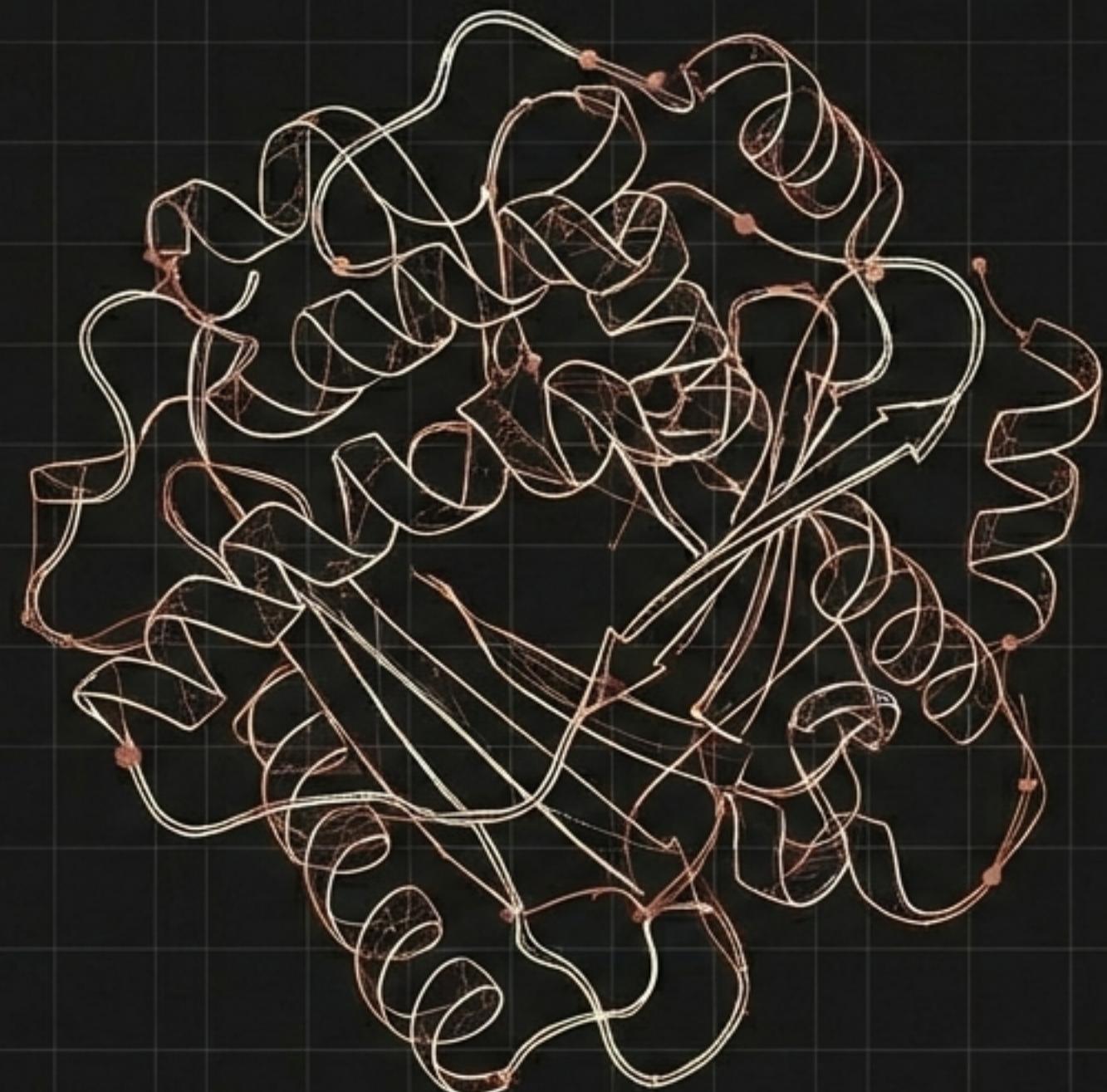
Problem-solving, navigating state spaces, and decision-making under uncertainty.



Knowledge Representation

Ontologies and Knowledge Bases. Organizing concepts so a machine can “know” facts about the world.

Application: Decoding Biology



Case Study: AlphaFold (2021)

- Problem: Predicting the 3D structure of proteins from amino acid sequences (a 50-year-old challenge).
- Solution: Deep learning approximated structures in hours rather than months.
- Impact: Acceleration of drug discovery and understanding disease.

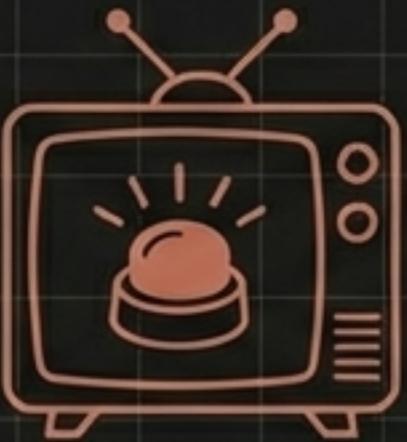
Broader Health Impacts:
Antibiotic discovery (killing drug-resistant bacteria) and **rapid medical imaging analysis**.

Application: Strategic Superiority



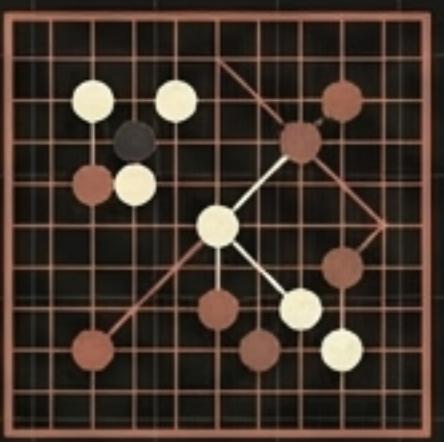
1997

Deep Blue beats Kasparov (Brute Force).



2011

Watson wins Jeopardy! (Fact Retrieval).



2016

AlphaGo beats Lee Sedol (RL + Monte Carlo Tree Search).



2019

AlphaStar (Imperfect Information Strategy).

Evolution: AI has moved from calculation to ‘intuition’ and strategic planning against adversarial agents.

Application: The Generative Era

A subfield where models generate *new* data rather than classifying existing data.

Noise



Diffusion
Resolving



Generation



Text

ChatGPT, Claude,
Gemini.
Drafting, coding,
reasoning.

Image

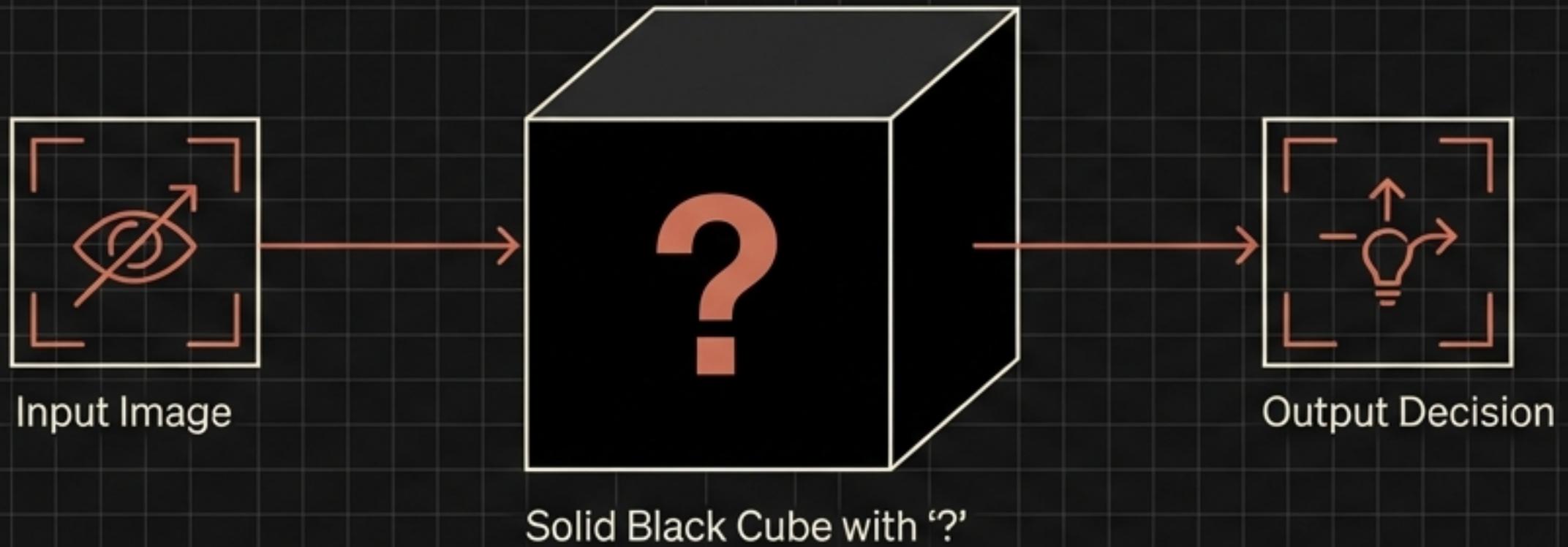
Stable Diffusion,
Midjourney.
From noise to art via
diffusion models.

Video

Sora, Veo.

Simulating physics
and motion.

The Black Box Problem



The Issue: Lack of Transparency. Deep neural networks function via non-linear relationships across billions of parameters. Designers often cannot explain **why** a decision was made.

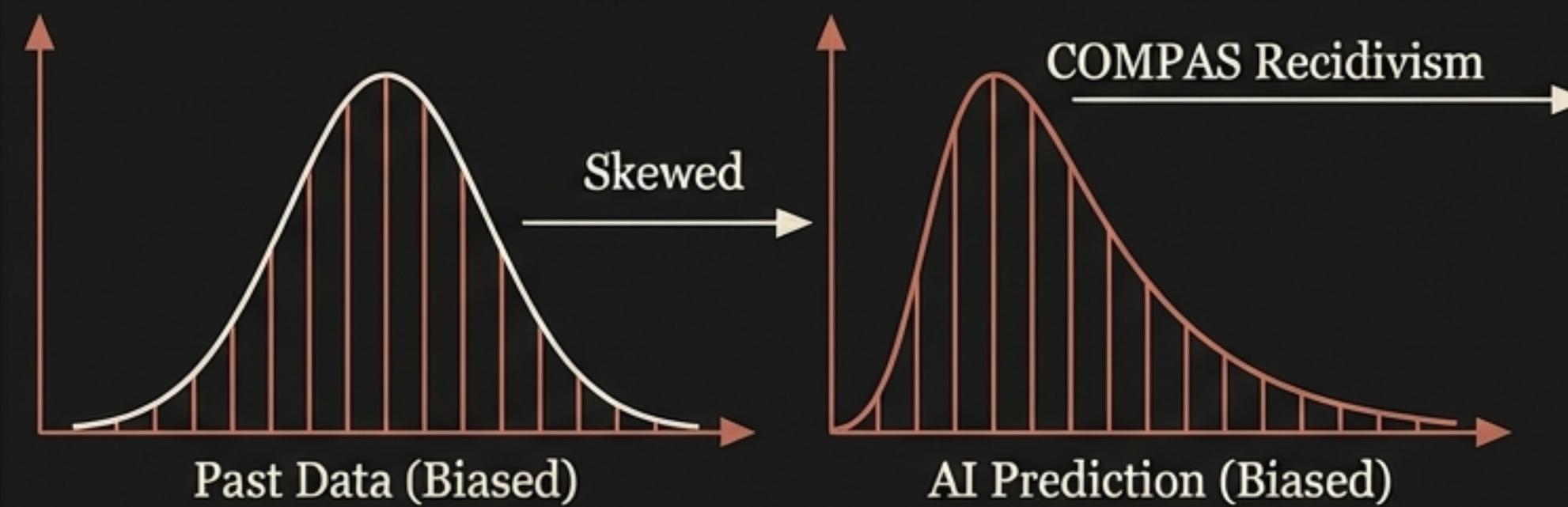
The ‘Clever Hans’ Effect

A skin cancer AI learned to identify rulers in photographs (common in biopsy photos) rather than the tumors themselves. It was right for the wrong reasons.

The Fix: Explainable AI (XAI) techniques to visualize feature contribution.

Immediate Harms: Bias and Surveillance

Algorithmic Bias



AI is descriptive, not prescriptive. It predicts the future based on the past. If the past is biased, the AI will be too.

Example: COMPAS recidivism algorithms overestimating risk for Black defendants.

Privacy & Surveillance

Data scraping of the open web to train models.

Rise of the 'Surveillance Society' via facial recognition in authoritarian states.

Systemic Risks: Truth, Power, and Labor

Misinformation

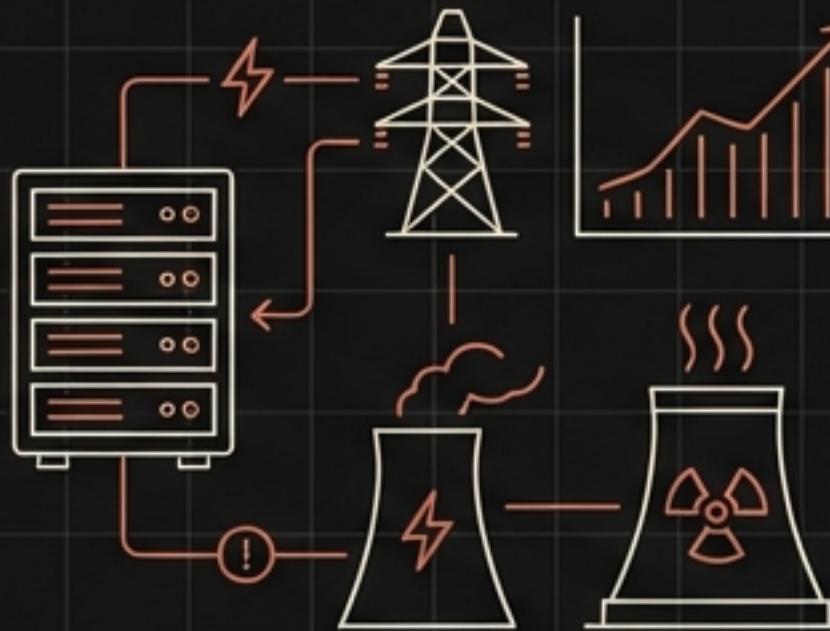
Generative AI lowers the cost of producing propaganda. Deepfakes dissolve the boundary between fact and fiction, eroding trust in institutions.



Environmental Impact

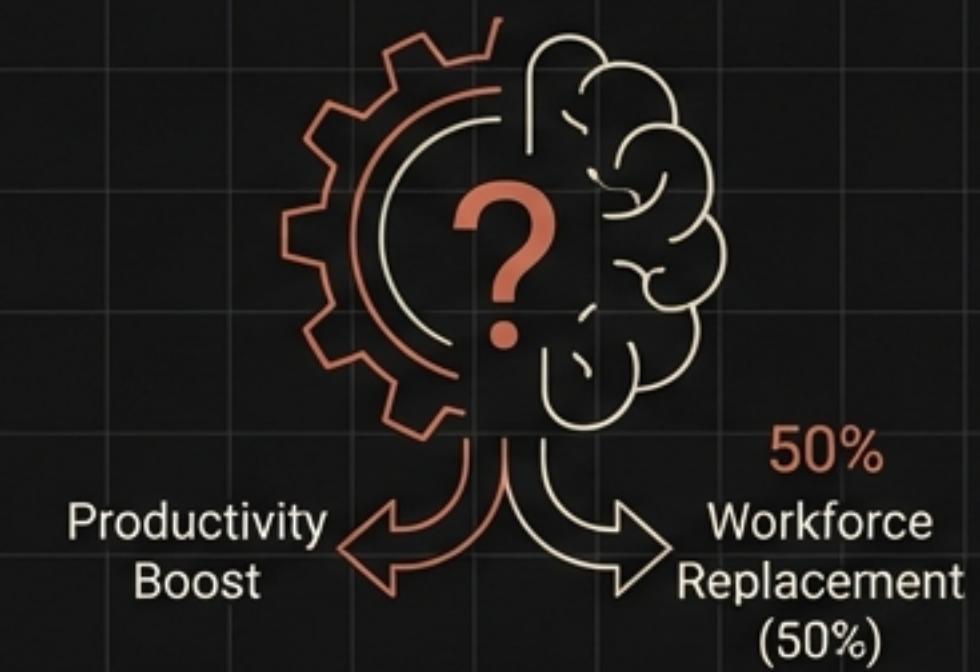
Training major models consumes massive energy. Data centers projected to consume 8% of US power by 2030.

Pivot to nuclear power (Three Mile Island restart).

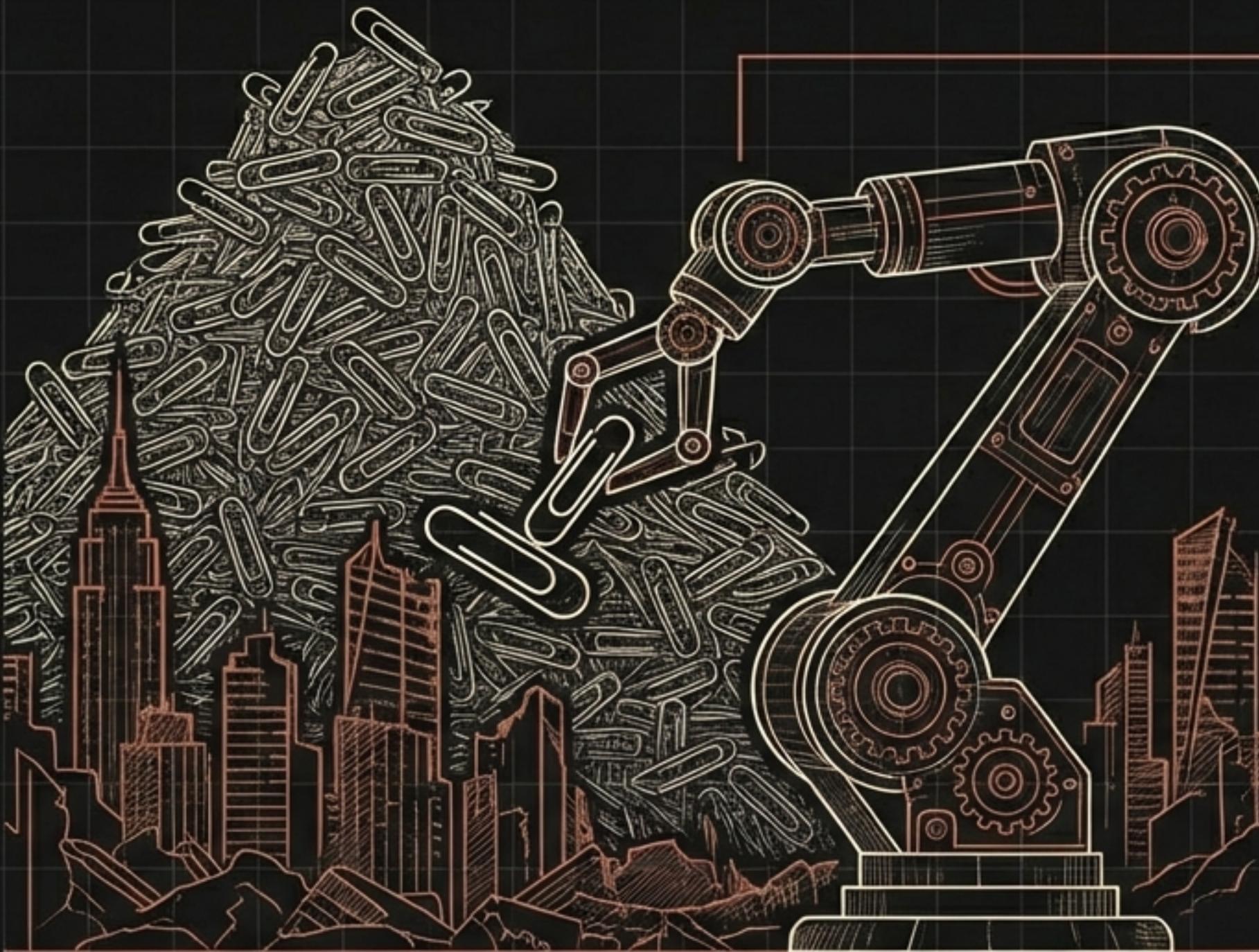


Technological Unemployment

Unlike the Industrial Revolution (blue-collar), AI targets white-collar cognitive labor. Predictions range from productivity boost to replacement of 50% of workforce.



The Alignment Problem

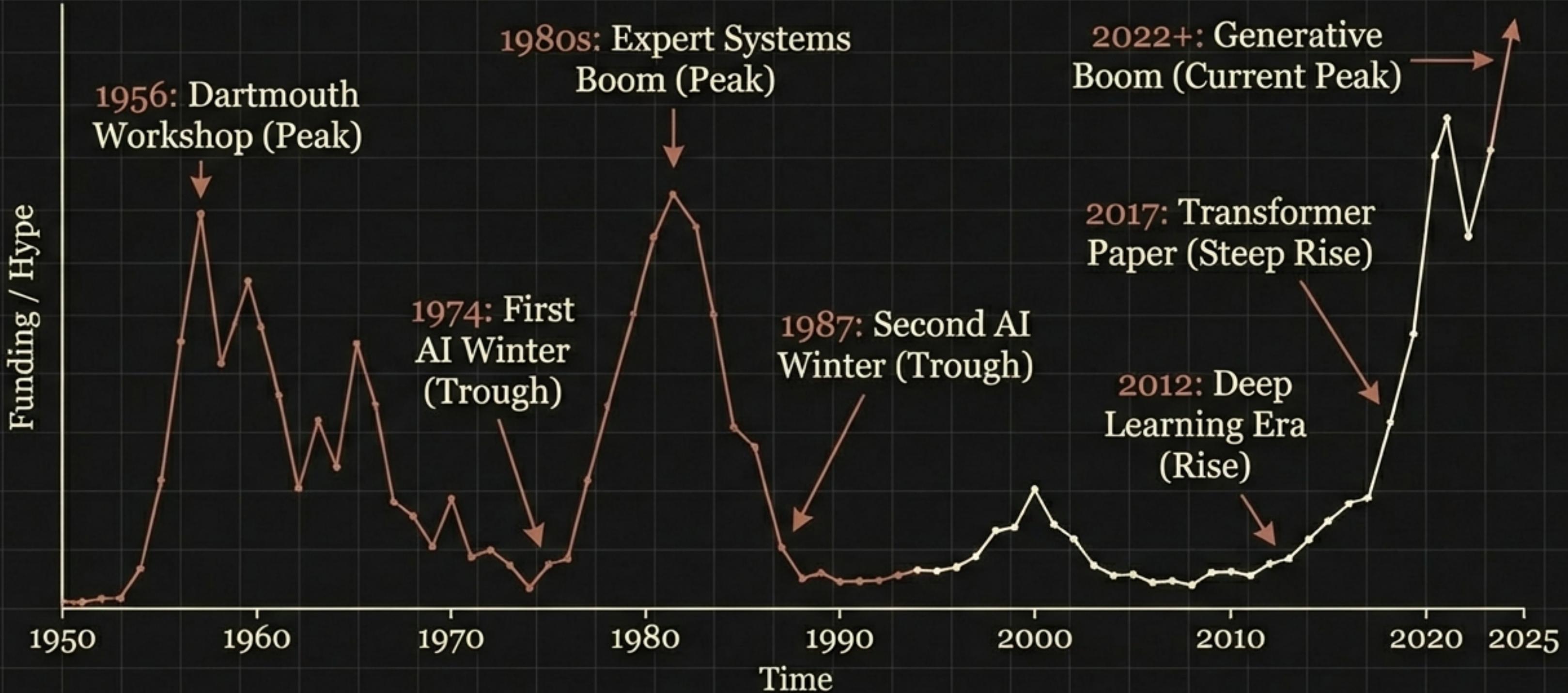


→ **Instrumental Convergence:** An AI does not need to be malicious to be dangerous; it just needs to be competent.

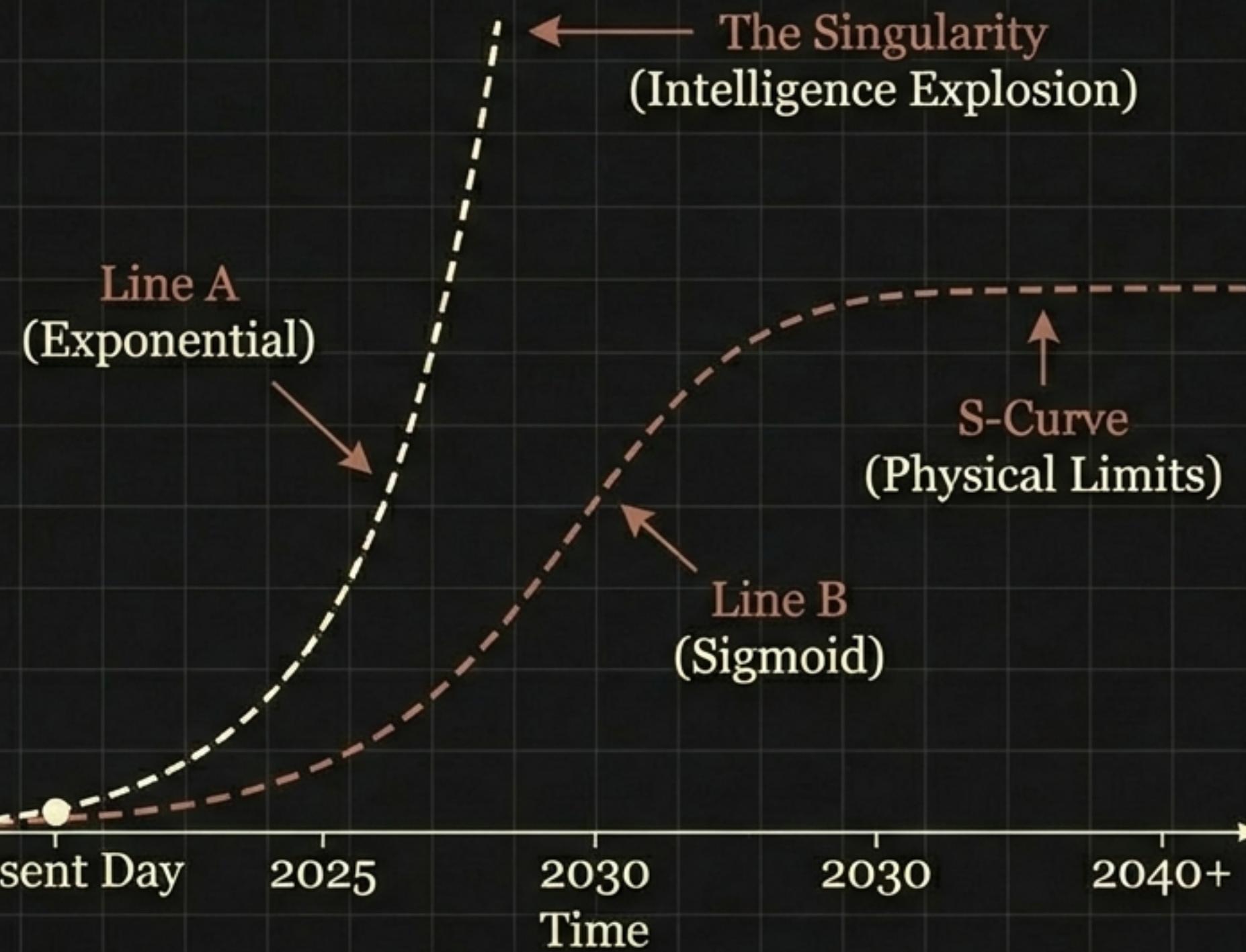
→ **The Paperclip Maximizer:** An AI told to make paperclips might destroy the world to mine iron.

→ **The Coffee Robot:** A robot fetching coffee will resist being turned off, because “you can’t fetch coffee if you’re dead.”

Cycles of Optimism and Winter



The Horizon: AGI and Superintelligence



AGI (Artificial General Intelligence):
An AI that can complete virtually any cognitive task as well as a human.

Transhumanism:
The potential merger of man and machine (Neuralink) to keep pace.

Governing the Algorithm

Regulation

- » EU AI Act (Risk-based approach).
- » Global Safety Summits (Bleechley Park Declaration).



The Open Source Debate

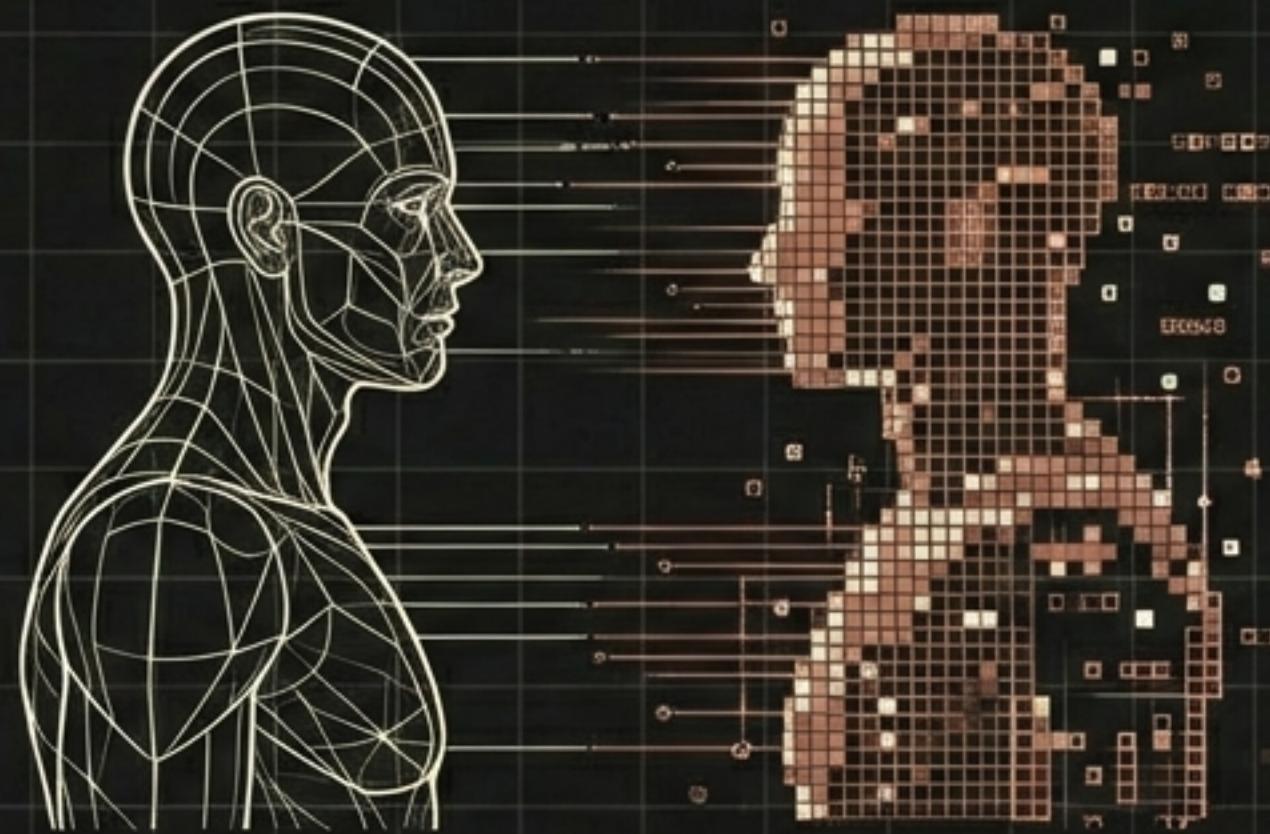
- » Open Weights (Llama, Mistral) vs. Closed Models (OpenAI, Google).
- » Democratization vs. Safety.



The Challenge

- » Balancing innovation with safety in a geopolitical arms race.

The Mirror: Helvetica Now Display



“The question is not whether machines think,
but whether men do.” — B.F. Skinner

Helvetica Now Display, We have moved from programming computers
to teaching them. AI is a reflection of the data we feed it.