

Bias Analysis Report

Artificial intelligence

Source: https://en.wikipedia.org/wiki/Artificial_intelligence

Analysis Summary

Overall Bias Score: 12.37 (Score per 1000 words)

Biased Segments Found: 266

Article Text with Highlights

Artificial intelligence (AI) is the capability of computational systems to perform tasks typically associated with human intelligence, such as learning, reasoning, problem-solving, perception, and decision-making. It is a field of research in computer science that develops and studies methods and software that enable machines to perceive their environment and use learning and intelligence to take actions that maximize their chances of achieving defined goals.

High-profile applications of AI include advanced web search engines (e.g., Google Search); recommendation systems (used by YouTube, Amazon, and Netflix); virtual assistants (e.g., Google Assistant, Siri, and Alexa); autonomous vehicles (e.g., Waymo); generative and creative tools (e.g., language models and AI art); and superhuman play and analysis in strategy games (e.g., chess and Go). However, many AI applications are not perceived as AI: "A lot of cutting edge AI has filtered into general applications, often without being called AI because once something becomes useful enough and common enough it's not labeled AI anymore."

Various subfields of AI research are centered around particular goals and the use of particular tools. The traditional goals of AI research include learning, reasoning, knowledge representation, planning, natural language processing,

perception, and support for robotics. To reach these goals, AI researchers have adapted and integrated a wide range of techniques, including search and mathematical optimization, formal logic, artificial neural networks, and methods based on statistics, operations research, and economics. AI also draws upon psychology, linguistics, philosophy, neuroscience, and other fields.

Some companies, such as OpenAI, Google DeepMind and Meta, aim to create artificial general intelligence (AGI)—AI that can complete virtually any cognitive task at least as well as a human.

Artificial intelligence was founded as an academic discipline in 1956, and the field went through multiple cycles of optimism throughout its history, followed by periods of disappointment and loss of funding, known as AI winters. Funding and interest vastly increased after 2012 when graphics processing units started being used to accelerate neural networks and deep learning outperformed previous AI techniques. This growth accelerated further after 2017 with the transformer architecture. In the 2020s, an ongoing period of rapid progress in advanced generative AI became known as the AI boom. Generative AI's ability to create and modify content has led to several unintended consequences and harms, which has raised ethical concerns about AI's long-term effects and potential existential risks, prompting discussions about regulatory policies to ensure the safety and benefits of the technology.

Goals

The general problem of simulating (or creating) intelligence has been broken into subproblems. These consist of particular traits or capabilities that researchers expect an intelligent system to display. The traits described below have received the most attention and cover the scope of AI research.

Reasoning and problem-solving

Early researchers developed algorithms that imitated step-by-step reasoning that humans use when they solve puzzles or make logical deductions. By the late 1980s and 1990s, methods were developed for dealing with uncertain or incomplete information, employing concepts from probability and economics.

Many of these algorithms are insufficient for solving large reasoning problems because they experience a "combinatorial explosion": They become exponentially slower as the problems grow. Even humans rarely use the

step-by-step deduction that early AI research could model. They solve most of their problems using fast, intuitive judgments. Accurate and efficient reasoning is an unsolved problem.

Knowledge representation

Knowledge representation and knowledge engineering allow AI programs to answer questions intelligently and make deductions about real-world facts. Formal knowledge representations are used in content-based indexing and retrieval, scene interpretation, clinical decision support, knowledge discovery (mining "interesting" and actionable inferences from large databases), and other areas.

A knowledge base is a body of knowledge represented in a form that can be used by a program. An ontology is the set of objects, relations, concepts, and properties used by a particular domain of knowledge.

Knowledge bases need to represent things such as objects, properties, categories, and relations between objects; situations, events, states, and time; causes and effects; knowledge about knowledge (what we know about what other people know); default reasoning (things that humans assume are true until they are told differently and will remain true even when other facts are changing); and many other aspects and domains of knowledge.

Among the most difficult problems in knowledge representation are the breadth of commonsense knowledge (the set of atomic facts that the average person knows is enormous); and the sub-symbolic form of most commonsense knowledge (much of what people know is not represented as "facts" or "statements" that they could express verbally). There is also the difficulty of knowledge acquisition, the problem of obtaining knowledge for AI applications.

Planning and decision-making

An "agent" is anything that perceives and takes actions in the world. A rational agent has goals or preferences and takes actions to make them happen. In automated planning, the agent has a specific goal.

In automated decision-making, the agent has preferences—there are some situations it would prefer to be in, and some situations it is trying to avoid. The decision-making agent assigns a number to each situation (called the "utility") that measures how much the agent prefers it. For each possible action, it can calculate the "expected utility": the utility of all possible

outcomes of the action, weighted by the probability that the outcome will occur. It can then choose the action with the maximum expected utility. In classical planning, the agent knows exactly what the effect of any action will be. In most real-world problems, however, the agent may not be certain about the situation they are in (it is "unknown" or "unobservable") and it may not know for certain what will happen after each possible action (it is not "deterministic"). It must choose an action by making a probabilistic guess and then reassess the situation to see if the action worked. In some problems, the agent's preferences may be uncertain, especially if there are other agents or humans involved. These can be learned (e.g., with inverse reinforcement learning), or the agent can seek information to improve its preferences. Information value theory can be used to weigh the value of exploratory or experimental actions. The space of possible future actions and situations is typically intractably large, so the agents must take actions and evaluate situations while being uncertain of what the outcome will be.

A Markov decision process has a transition model that describes the probability that a particular action will change the state in a particular way and a reward function that supplies the utility of each state and the cost of each action. A policy associates a decision with each possible state. The policy could be calculated (e.g., by iteration), be heuristic, or it can be learned.

Game theory describes the rational behavior of multiple interacting agents and is used in AI programs that make decisions that involve other agents.

Learning

Machine learning is the study of programs that can improve their performance on a given task automatically. It has been a part of AI from the beginning.

There are several kinds of machine learning. Unsupervised learning analyzes a stream of data and finds patterns and makes predictions without any other guidance. Supervised learning requires labeling the training data with the expected answers, and comes in two main varieties: classification (where the program must learn to predict what category the input belongs in) and regression (where the program must deduce a numeric function based on numeric input).

In reinforcement learning, the agent is rewarded for good responses and punished for bad ones. The agent learns to choose responses that are classified as "good". Transfer learning is when the knowledge gained from one problem is applied to a new problem. Deep learning is a type of machine learning that runs inputs through biologically inspired artificial neural networks for all of these types of learning.

Computational learning theory can assess learners by computational complexity, by sample complexity (how much data is required), or by other notions of optimization.

Natural language processing

Natural language processing (NLP) allows programs to read, write and communicate in human languages. Specific problems include speech recognition, speech synthesis, machine translation, information extraction, information retrieval and question answering.

Early work, based on Noam Chomsky's generative grammar and semantic networks, had difficulty with word-sense disambiguation unless restricted to small domains called "micro-worlds" (due to the common sense knowledge problem). Margaret Masterman believed that it was meaning and not grammar that was the key to understanding languages, and that thesauri and not dictionaries should be the basis of computational language structure.

Modern deep learning techniques for NLP include word embedding (representing words, typically as vectors encoding their meaning), transformers (a deep learning architecture using an attention mechanism), and others.

In 2019, generative pre-trained transformer (or "GPT") language models began to generate coherent text, and by 2023, these models were able to get human-level scores on the bar exam, SAT test, GRE test, and many other real-world applications.

Perception

Machine perception is the ability to use input from sensors (such as cameras, microphones, wireless signals, active lidar, sonar, radar, and tactile sensors) to deduce aspects of the world. Computer vision is the ability to analyze visual input.

The field includes speech recognition, image classification, facial recognition, object recognition, object tracking, and robotic perception.

Social intelligence

Affective computing is a field that comprises systems that recognize, interpret, process, or simulate human feeling, emotion, and mood. For example, some virtual assistants are programmed to speak conversationally or even to banter humorously; it makes them appear more sensitive to the emotional dynamics of human interaction, or to otherwise facilitate human-computer interaction.

However, this tends to give naïve users an unrealistic conception of the intelligence of existing computer agents. Moderate successes related to affective computing include textual sentiment analysis and, more recently, multimodal sentiment analysis, wherein AI classifies the effects displayed by a videotaped subject.

General intelligence

A machine with artificial general intelligence would be able to solve a wide variety of problems with breadth and versatility similar to human intelligence.

Techniques

AI research uses a wide variety of techniques to accomplish the goals above.

Search and optimization

AI can solve many problems by intelligently searching through many possible solutions. There are two very different kinds of search used in AI: state space search and local search.

State space search

State space search searches through a tree of possible states to try to find a goal state. For example, planning algorithms search through trees of goals and subgoals, attempting to find a path to a target goal, a process called means-ends analysis.

Simple exhaustive searches are rarely sufficient for most real-world problems: the search space (the number of places to search) quickly grows to astronomical numbers. The result is a search that is too slow or never completes. "Heuristics" or "rules of thumb" can help prioritize choices that are more likely to reach a goal.

Adversarial search is used for game-playing programs, such as chess or Go. It searches through a tree of possible moves and countermoves, looking for a winning position.

Local search

Local search uses mathematical optimization to find a solution to a problem.

It begins with some form of guess and refines it incrementally.

Gradient descent is a type of local search that optimizes a set of numerical parameters by incrementally adjusting them to minimize a loss function.

Variants of gradient descent are commonly used to train neural networks, through the backpropagation algorithm.

Another type of local search is evolutionary computation, which aims to iteratively improve a set of candidate solutions by "mutating" and "recombining" them, selecting only the fittest to survive each generation.

Distributed search processes can coordinate via swarm intelligence algorithms.

Two popular swarm algorithms used in search are particle swarm optimization (inspired by bird flocking) and ant colony optimization (inspired by ant trails).

Logic

Formal logic is used for reasoning and knowledge representation.

Formal logic comes in two main forms: propositional logic (which operates on statements that are true or false and uses logical connectives such as "and", "or", "not" and "implies") and predicate logic (which also operates on objects, predicates and relations and uses quantifiers such as "Every X is a Y" and "There are some Xs that are Ys").

Deductive reasoning in logic is the process of proving a new statement (conclusion) from other statements that are given and assumed to be true (the premises). Proofs can be structured as proof trees, in which nodes are labelled by sentences, and children nodes are connected to parent nodes by inference rules.

Given a problem and a set of premises, problem-solving reduces to searching for a proof tree whose root node is labelled by a solution of the problem and whose leaf nodes are labelled by premises or axioms. In the case of Horn clauses, problem-solving search can be performed by reasoning forwards from the premises or backwards from the problem. In the more general case

of the clausal form of first-order logic, resolution is a single, axiom-free rule of inference, in which a problem is solved by proving a contradiction from premises that include the negation of the problem to be solved.

Inference in both Horn clause logic and first-order logic is undecidable, and therefore

intractable. However, backward reasoning with Horn clauses, which underpins computation in the logic programming language Prolog, is Turing complete. Moreover, its efficiency is competitive with computation in other symbolic programming languages.

Fuzzy logic assigns a "degree of truth" between 0 and 1. It can therefore handle propositions that are vague and partially true.

Non-monotonic logics, including logic programming with negation as failure, are designed to handle default reasoning. Other specialized versions of logic have been developed to describe many complex domains.

Probabilistic methods for uncertain reasoning

Many problems in AI (including reasoning, planning, learning, perception, and robotics) require the agent to operate with incomplete or uncertain information. AI researchers have devised a number of tools to solve these problems using methods from probability theory and economics. Precise mathematical tools have been developed that analyze how an agent can make choices and plan, using decision theory, decision analysis, and information value theory. These tools include models such as Markov decision processes, dynamic decision networks, game theory and mechanism design.

Bayesian networks are a tool that can be used for reasoning (using the Bayesian inference algorithm), learning (using the expectation-maximization algorithm), planning (using decision networks) and perception (using dynamic Bayesian networks).

Probabilistic algorithms can also be used for filtering, prediction, smoothing, and finding explanations for streams of data, thus helping perception systems analyze processes that occur over time (e.g., hidden Markov models or Kalman filters).

Classifiers and statistical learning methods

The simplest AI applications can be divided into two types: classifiers (e.g., "if shiny then diamond"), on one hand, and controllers (e.g., "if diamond

then pick up"), on the other hand. Classifiers are functions that use pattern matching to determine the closest match. They can be fine-tuned based on chosen examples using supervised learning. Each pattern (also called an "observation") is labeled with a certain predefined class. All the observations combined with their class labels are known as a data set. When a new observation is received, that observation is classified based on previous experience.

There are many kinds of classifiers in use. The decision tree is the simplest and most widely used symbolic machine learning algorithm. K-nearest neighbor algorithm was the most widely used analogical AI until the mid-1990s, and Kernel methods such as the support vector machine (SVM) displaced k-nearest neighbor in the 1990s.

The naive Bayes classifier is reportedly the "most widely used learner" at Google, due in part to its scalability.

Neural networks are also used as classifiers.

Artificial neural networks

An artificial neural network is based on a collection of nodes also known as artificial neurons, which loosely model the neurons in a biological brain. It is trained to recognise patterns; once trained, it can recognise those patterns in fresh data. There is an input, at least one hidden layer of nodes and an output. Each node applies a function and once the weight crosses its specified threshold, the data is transmitted to the next layer. A network is typically called a deep neural network if it has at least 2 hidden layers. Learning algorithms for neural networks use local search to choose the weights that will get the right output for each input during training. The most common training technique is the backpropagation algorithm. Neural networks learn to model complex relationships between inputs and outputs and find patterns in data. In theory, a neural network can learn any function. In feedforward neural networks the signal passes in only one direction. The term perceptron typically refers to a single-layer neural network.

In contrast, deep learning uses many layers. Recurrent neural networks (RNNs) feed the output signal back into the input, which allows short-term memories of previous input events. Long short-term memory networks (LSTMs) are recurrent neural networks that better preserve longterm dependencies and are less sensitive to the vanishing gradient problem. Convolutional neural networks (CNNs) use layers of kernels to more

efficiently process local patterns.

This local processing is especially important in image processing, where the early CNN layers typically identify simple local patterns such as edges and curves, with subsequent layers detecting more complex patterns like textures, and eventually whole objects.

Deep learning

Deep learning uses several layers of neurons between the network's inputs and outputs. The multiple layers can progressively extract higher-level features from the raw input. For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human such as digits, letters, or faces.

Deep learning has profoundly improved the performance of programs in many important subfields of artificial intelligence, including computer vision, speech recognition, natural language processing, image classification, and others. The reason that deep learning performs so well in so many applications is not known as of 2021.

The sudden success of deep learning in 2012–2015 did not occur because of some new discovery or theoretical breakthrough (deep neural networks and backpropagation had been described by many people, as far back as the 1950s) but because of two factors: the incredible increase in computer power (including the hundred-fold increase in speed by switching to GPUs) and the availability of vast amounts of training data, especially the giant curated datasets used for benchmark testing, such as ImageNet.

GPT

Generative pre-trained transformers (GPT) are large language models (LLMs) that generate text based on the semantic relationships between words in sentences. Text-based GPT models are pre-trained on a large corpus of text that can be from the Internet. The pretraining consists of predicting the next token (a token being usually a word, subword, or punctuation). Throughout this pretraining, GPT models accumulate knowledge about the world and can then generate human-like text by repeatedly predicting the next token.

Typically, a subsequent training phase makes the model more truthful, useful, and harmless, usually with a technique called reinforcement learning from human feedback (RLHF). Current GPT models are prone to generating falsehoods called "hallucinations". These can be reduced with RLHF and

quality data, but the problem has been getting worse for reasoning systems. Such systems are used in chatbots, which allow people to ask a question or request a task in simple text. Current models and services include ChatGPT, Claude, Gemini, Copilot, and Meta AI. Multimodal GPT models can process different types of data (modalities) such as images, videos, sound, and text.

Hardware and software

In the late 2010s, graphics processing units (GPUs) that were increasingly designed with AI-specific enhancements and used with specialized TensorFlow software had replaced previously used central processing unit (CPUs) as the dominant means for large-scale (commercial and academic) machine learning models' training. Specialized programming languages such as Prolog were used in early AI research, but general-purpose programming languages like Python have become predominant.

The transistor density in integrated circuits has been observed to roughly double every 18 months—a trend known as Moore's law, named after the Intel co-founder Gordon Moore, who first identified it.

Improvements in GPUs have been even faster, a trend sometimes called Huang's law, named after Nvidia co-founder and CEO Jensen Huang.

Applications

AI and machine learning technology is used in most of the essential applications of the 2020s, including: search engines (such as Google Search), targeting online advertisements, recommendation systems (offered by Netflix, YouTube or Amazon), driving internet traffic, targeted advertising (AdSense, Facebook), virtual assistants (such as Siri or Alexa), autonomous vehicles (including drones, ADAS and self-driving cars), automatic language translation (Microsoft Translator, Google Translate), facial recognition (Apple's FaceID or Microsoft's DeepFace and Google's FaceNet) and image labeling (used by Facebook, Apple's Photos and TikTok). The deployment of AI may be overseen by a chief automation officer (CAO).

Health and medicine

The application of AI in medicine and medical research has the potential to increase patient care and quality of life. Through the lens of the Hippocratic Oath, medical professionals are ethically compelled to use AI, if applications

can more accurately diagnose and treat patients.

For medical research, AI is an important tool for processing and integrating big data. This is particularly important for organoid and tissue engineering development which use microscopy imaging as a key technique in fabrication. It has been suggested that AI can overcome discrepancies in funding allocated to different fields of research.

New AI tools can deepen the understanding of biomedically relevant pathways. For example, AlphaFold 2 (2021) demonstrated the ability to approximate, in hours rather than months, the 3D structure of a protein. In 2023, it was reported that AI-guided drug discovery helped find a class of antibiotics capable of killing two different types of drug-resistant bacteria. In 2024, researchers used machine learning to accelerate the search for Parkinson's disease drug treatments. Their aim was to identify compounds that block the clumping, or aggregation, of alpha-synuclein (the protein that characterises Parkinson's disease). They were able to speed up the initial screening process ten-fold and reduce the cost by a thousand-fold.

Games

Game playing programs have been used since the 1950s to demonstrate and test AI's most advanced techniques. Deep Blue became the first computer chess-playing system to beat a reigning world chess champion, Garry Kasparov, on 11 May 1997. In 2011, in a Jeopardy! quiz show exhibition match, IBM's question answering system, Watson, defeated the two greatest Jeopardy! champions, Brad Rutter and Ken Jennings, by a significant margin. In March 2016, AlphaGo won 4 out of 5 games of Go in a match with Go champion Lee Sedol, becoming the first computer Go-playing system to beat a professional Go player without handicaps. Then, in 2017, it defeated Ke Jie, who was the best Go player in the world. Other programs handle imperfect-information games, such as the poker-playing program Pluribus. DeepMind developed increasingly generalistic reinforcement learning models, such as with MuZero, which could be trained to play chess, Go, or Atari games. In 2019, DeepMind's AlphaStar achieved grandmaster level in StarCraft II, a particularly challenging real-time strategy game that involves incomplete knowledge of what happens on the map. In 2021, an AI agent competed in a PlayStation Gran Turismo competition, winning against four of the world's best Gran Turismo drivers using deep reinforcement learning. In 2024, Google DeepMind introduced SIMA, a type of AI capable of

autonomously playing nine previously unseen open-world video games by observing screen output, as well as executing short, specific tasks in response to natural language instructions.

Mathematics

Large language models, such as GPT-4, Gemini, Claude, Llama or Mistral, are increasingly used in mathematics. These probabilistic models are versatile, but can also produce wrong answers in the form of hallucinations. They sometimes need a large database of mathematical problems to learn from, but also methods such as supervised fine-tuning or trained classifiers with human-annotated data to improve answers for new problems and learn from corrections. A February 2024 study showed that the performance of some language models for reasoning capabilities in solving math problems not included in their training data was low, even for problems with only minor deviations from trained data. One technique to improve their performance involves training the models to produce correct reasoning steps, rather than just the correct result. The Alibaba Group developed a version of its Qwen models called Qwen2-Math, that achieved state-of-the-art performance on several mathematical benchmarks, including 84% accuracy on the MATH dataset of competition mathematics problems. In January 2025, Microsoft proposed the technique rStar-Math that leverages Monte Carlo tree search and step-by-step reasoning, enabling a relatively small language model like Qwen-7B to solve 53% of the AIME 2024 and 90% of the MATH benchmark problems.

Alternatively, dedicated models for mathematical problem solving with higher precision for the outcome including proof of theorems have been developed such as AlphaTensor, AlphaGeometry, AlphaProof and AlphaEvolve all from Google DeepMind, Llemma from EleutherAI or Julius.

When natural language is used to describe mathematical problems, converters can transform such prompts into a formal language such as Lean to define mathematical tasks. The experimental model Gemini Deep Think accepts natural language prompts directly and achieved gold medal results in the International Math Olympiad of 2025.

Some models have been developed to solve challenging problems and reach good results in benchmark tests, others to serve as educational tools in mathematics.

Topological deep learning integrates various topological approaches.

Finance

Finance is one of the fastest growing sectors where applied AI tools are being deployed: from retail online banking to investment advice and insurance, where automated "robot advisers" have been in use for some years.

According to Nicolas Firzli, director of the World Pensions & Investments Forum, it may be too early to see the emergence of highly innovative AI-informed financial products and services. He argues that "the deployment of AI tools will simply further automatise things: destroying tens of thousands of jobs in banking, financial planning, and pension advice in the process, but I'm not sure it will unleash a new wave of [e.g., sophisticated] pension innovation."

Military

Various countries are deploying AI military applications. The main applications enhance command and control, communications, sensors, integration and interoperability. Research is targeting intelligence collection and analysis, logistics, cyber operations, information operations, and semiautonomous and autonomous vehicles. AI technologies enable coordination of sensors and effectors, threat detection and identification, marking of enemy positions, target acquisition, coordination and deconfliction of distributed Joint Fires between networked combat vehicles, both human-operated and autonomous.

AI has been used in military operations in Iraq, Syria, Israel and Ukraine.

Generative AI

Agents

AI agents are software entities designed to perceive their environment, make decisions, and take actions autonomously to achieve specific goals. These agents can interact with users, their environment, or other agents. AI agents are used in various applications, including virtual assistants, chatbots, autonomous vehicles, game-playing systems, and industrial robotics. AI agents operate within the constraints of their programming, available computational resources, and hardware limitations. This means they are restricted to performing tasks within their defined scope and have finite memory and processing capabilities. In real-world applications, AI agents often face time constraints for decision-making and action execution. Many

AI agents incorporate learning algorithms, enabling them to improve their performance over time through experience or training. Using machine learning, AI agents can adapt to new situations and optimise their behaviour for their designated tasks.

Sexuality

Applications of AI in this domain include AI-enabled menstruation and fertility trackers that analyze user data to offer predictions, AI-integrated sex toys (e.g., teledildonics), AI-generated sexual education content, and AI agents that simulate sexual and romantic partners (e.g., Replika). AI is also used for the production of non-consensual deepfake pornography, raising significant ethical and legal concerns.

AI technologies have also been used to attempt to identify online gender-based violence and online sexual grooming of minors.

Other industry-specific tasks

There are also thousands of successful AI applications used to solve specific problems for specific industries or institutions. In a 2017 survey, one in five companies reported having incorporated "AI" in some offerings or processes. A few examples are energy storage, medical diagnosis, military logistics, applications that predict the result of judicial decisions, foreign policy, or supply chain management.

AI applications for evacuation and disaster management are growing. AI has been used to investigate patterns in large-scale and small-scale evacuations using historical data from GPS, videos or social media. Furthermore, AI can provide real-time information on the evacuation conditions.

In agriculture, AI has helped farmers to increase yield and identify areas that need irrigation, fertilization, pesticide treatments. Agronomists use AI to conduct research and development. AI has been used to predict the ripening time for crops such as tomatoes, monitor soil moisture, operate agricultural robots, conduct predictive analytics, classify livestock pig call emotions, automate greenhouses, detect diseases and pests, and save water.

Artificial intelligence is used in astronomy to analyze increasing amounts of available data and applications, mainly for "classification, regression, clustering, forecasting, generation, discovery, and the development of new scientific insights." For example, it is used for discovering exoplanets, forecasting solar activity, and distinguishing between signals and

instrumental effects in gravitational wave astronomy.

Additionally, it could be used for activities in space, such as space exploration, including the analysis of data from space missions, real-time science decisions of spacecraft, space debris avoidance, and more autonomous operation.

During the 2024 Indian elections, US\$50 million was spent on authorized AI-generated content, notably by creating deepfakes of allied (including sometimes deceased) politicians to better engage with voters, and by translating speeches to various local languages.

Ethics

AI has potential benefits and potential risks. AI may be able to advance science and find solutions for serious problems: Demis Hassabis of DeepMind hopes to "solve intelligence, and then use that to solve everything else". However, as the use of AI has become widespread, several unintended consequences and risks have been identified. In-production systems can sometimes not factor ethics and bias into their AI training processes, especially when the AI algorithms are inherently unexplainable in deep learning.

Risks and harm

Privacy and copyright

Machine learning algorithms require large amounts of data. The techniques used to acquire this data have raised concerns about privacy, surveillance and copyright.

AI-powered devices and services, such as virtual assistants and IoT products, continuously collect personal information, raising concerns about intrusive data gathering and unauthorized access by third parties. The loss of privacy is further exacerbated by AI's ability to process and combine vast amounts of data, potentially leading to a surveillance society where individual activities are constantly monitored and analyzed without adequate safeguards or transparency.

Sensitive user data collected may include online activity records, geolocation data, video, or audio. For example, in order to build speech recognition algorithms, Amazon has recorded millions of private conversations and allowed temporary workers to listen to and transcribe some of them.

Opinions about this widespread surveillance range from those who see it as a

necessary evil to those for whom it is clearly unethical and a violation of the right to privacy.

AI developers argue that this is the only way to deliver valuable applications and have developed several techniques that attempt to preserve privacy while still obtaining the data, such as data aggregation, de-identification and differential privacy. Since 2016, some privacy experts, such as Cynthia Dwork, have begun to view privacy in terms of fairness. Brian Christian wrote that experts have pivoted "from the question of 'what they know' to the question of 'what they're doing with it'."

Generative AI is often trained on unlicensed copyrighted works, including in domains such as images or computer code; the output is then used under the rationale of "fair use". Experts disagree about how well and under what circumstances this rationale will hold up in courts of law; relevant factors may include "the purpose and character of the use of the copyrighted work" and "the effect upon the potential market for the copyrighted work". Website owners who do not wish to have their content scraped can indicate it in a "robots.txt" file. In 2023, leading authors (including John Grisham and Jonathan Franzen) sued AI companies for using their work to train generative AI. Another discussed approach is to envision a separate sui generis system of protection for creations generated by AI to ensure fair attribution and compensation for human authors.

Dominance by tech giants

The commercial AI scene is dominated by Big Tech companies such as Alphabet Inc., Amazon, Apple Inc., Meta Platforms, and Microsoft. Some of these players already own the vast majority of existing cloud infrastructure and computing power from data centers, allowing them to entrench further in the marketplace.

Power needs and environmental impacts

In January 2024, the International Energy Agency (IEA) released Electricity 2024, Analysis and Forecast to 2026, forecasting electric power use. This is the first IEA report to make projections for data centers and power consumption for artificial intelligence and cryptocurrency. The report states that power demand for these uses might double by 2026, with additional electric power usage equal to electricity used by the whole Japanese nation. Prodigious power consumption by AI is responsible for the growth of fossil

fuel use, and might delay closings of obsolete, carbon-emitting coal energy facilities. There is a feverish rise in the construction of data centers throughout the US, making large technology firms (e.g., Microsoft, Meta, Google, Amazon) into voracious consumers of electric power. Projected electric consumption is so immense that there is concern that it will be fulfilled no matter the source. A ChatGPT search involves the use of 10 times the electrical energy as a Google search.

The large firms are in haste to find power sources – from nuclear energy to geothermal to fusion. The tech firms argue that – in the long view – AI will be eventually kinder to the environment, but they need the energy now. AI makes the power grid more efficient and "intelligent", will assist in the growth of nuclear power, and track overall carbon emissions, according to technology firms.

A 2024 Goldman Sachs Research Paper, AI Data Centers and the Coming US Power Demand Surge, found "US power demand (is) likely to experience growth not seen in a generation...." and forecasts that, by 2030, US data centers will consume 8% of US power, as opposed to 3% in 2022, presaging growth for the electrical power generation industry by a variety of means. Data centers' need for more and more electrical power is such that they might max out the electrical grid. The Big Tech companies counter that AI can be used to maximize the utilization of the grid by all.

In 2024, the Wall Street Journal reported that big AI companies have begun negotiations with the US nuclear power providers to provide electricity to the data centers. In March 2024 Amazon purchased a Pennsylvania nuclear-powered data center for US\$650 million. Nvidia CEO Jensen Huang said nuclear power is a good option for the data centers.

In September 2024, Microsoft announced an agreement with Constellation Energy to re-open the Three Mile Island nuclear power plant to provide Microsoft with 100% of all electric power produced by the plant for 20 years. Reopening the plant, which suffered a partial nuclear meltdown of its Unit 2 reactor in 1979, will require Constellation to get through strict regulatory processes which will include extensive safety scrutiny from the US Nuclear Regulatory Commission. If approved (this will be the first ever US re-commissioning of a nuclear plant), over 835 megawatts of power – enough for 800,000 homes – of energy will be produced. The cost for re-opening and upgrading is estimated at US\$1.6

billion and is dependent on tax breaks for nuclear power contained in the 2022 US Inflation Reduction Act. The US government and the state of Michigan are investing almost US\$2 billion to reopen the Palisades Nuclear reactor on Lake Michigan. Closed since 2022, the plant is planned to be reopened in October 2025. The Three Mile Island facility will be renamed the Crane Clean Energy Center after Chris Crane, a nuclear proponent and former CEO of Exelon who was responsible for Exelon's spinoff of Constellation.

After the last approval in September 2023, Taiwan suspended the approval of data centers north of Taoyuan with a capacity of more than 5 MW in 2024, due to power supply shortages. Taiwan aims to phase out nuclear power by 2025. On the other hand, Singapore imposed a ban on the opening of data centers in 2019 due to electric power, but in 2022, lifted this ban.

Although most nuclear plants in Japan have been shut down after the 2011 Fukushima nuclear accident, according to an October 2024 Bloomberg article in Japanese, cloud gaming services company Ubitus, in which Nvidia has a stake, is looking for land in Japan near nuclear power plant for a new data center for generative AI. Ubitus CEO Wesley Kuo said nuclear power plants are the most efficient, cheap and stable power for AI.

On 1 November 2024, the Federal Energy Regulatory Commission (FERC) rejected an application submitted by Talen Energy for approval to supply some electricity from the nuclear power station Susquehanna to Amazon's data center.

According to the Commission Chairman Willie L. Phillips, it is a burden on the electricity grid as well as a significant cost shifting concern to households and other business sectors.

In 2025, a report prepared by the International Energy Agency estimated the greenhouse gas emissions from the energy consumption of AI at 180 million tons. By 2035, these emissions could rise to 300–500 million tonnes depending on what measures will be taken. This is below 1.5% of the energy sector emissions. The emissions reduction potential of AI was estimated at 5% of the energy sector emissions, but rebound effects (for example if people switch from public transport to autonomous cars) can reduce it.

Misinformation

YouTube, Facebook and others use recommender systems to guide users to

more content. These AI programs were given the goal of maximizing user engagement (that is, the only goal was to keep people watching). The AI learned that users tended to choose misinformation, conspiracy theories, and extreme partisan content, and, to keep them watching, the AI recommended more of it.

Users also tended to watch more content on the same subject, so the AI led people into filter bubbles where they received multiple versions of the same misinformation. This convinced many users that the misinformation was true, and ultimately undermined trust in institutions, the media and the government.

The AI program had correctly learned to maximize its goal, but the result was harmful to society. After the U.S. election in 2016, major technology companies took some steps to mitigate the problem.

In the early 2020s, generative AI began to create images, audio, and texts that are virtually indistinguishable from real photographs, recordings, or human writing, while realistic AI-generated videos became feasible in the mid-2020s. It is possible for bad actors to use this technology to create massive amounts of misinformation or propaganda; one such potential malicious use is deepfakes for computational propaganda. AI pioneer Geoffrey Hinton expressed concern about AI enabling "authoritarian leaders to manipulate their electorates" on a large scale, among other risks.

AI researchers at Microsoft, OpenAI, universities and other organisations have suggested using "personhood credentials" as a way to overcome online deception enabled by AI models.

Algorithmic bias and fairness

Machine learning applications will be biased if they learn from biased data.

The developers may not be aware that the bias exists. Bias can be introduced by the way training data is selected and by the way a model is deployed. If a biased algorithm is used to make decisions that can seriously harm people (as it can in medicine, finance, recruitment, housing or policing) then the algorithm may cause discrimination. The field of fairness studies how to prevent harms from algorithmic biases.

On June 28, 2015, Google Photos's new image labeling feature mistakenly identified Jacky Alcine and a friend as "gorillas" because they were black. The system was trained on a dataset that contained very few images of black people, a problem called "sample size disparity". Google "fixed" this problem

by preventing the system from labelling anything as a "gorilla". Eight years later, in 2023, Google Photos still could not identify a gorilla, and neither could similar products from Apple, Facebook, Microsoft and Amazon.

COMPAS is a commercial program widely used by U.S. courts to assess the likelihood of a defendant becoming a recidivist. In 2016, Julia Angwin at ProPublica discovered that COMPAS exhibited racial bias, despite the fact that the program was not told the races of the defendants. Although the error rate for both whites and blacks was calibrated equal at exactly 61%, the errors for each race were different—the system consistently overestimated the chance that a black person would re-offend and would underestimate the chance that a white person would not re-offend.

In 2017, several researchers showed that it was mathematically impossible for COMPAS to accommodate all possible measures of fairness when the base rates of re-offense were different for whites and blacks in the data.

A program can make biased decisions even if the data does not explicitly mention a problematic feature (such as "race" or "gender"). The feature will correlate with other features (like "address", "shopping history" or "first name"), and the program will make the same decisions based on these features as it would on "race" or "gender". Moritz Hardt said "the most robust fact in this research area is that fairness through blindness doesn't work."

Criticism of COMPAS highlighted that machine learning models are designed to make "predictions" that are only valid if we assume that the future will resemble the past. If they are trained on data that includes the results of racist decisions in the past, machine learning models must predict that racist decisions will be made in the future. If an application then uses these predictions as recommendations, some of these "recommendations" will likely be racist. Thus, machine learning is not well suited to help make decisions in areas where there is hope that the future will be better than the past. It is descriptive rather than prescriptive.

Bias and unfairness may go undetected because the developers are overwhelmingly white and male: among AI engineers, about 4% are black and 20% are women.

There are various conflicting definitions and mathematical models of fairness. These notions depend on ethical assumptions, and are influenced by beliefs about society. One broad category is distributive fairness, which focuses on the outcomes, often identifying groups and seeking to compensate for statistical disparities. Representational fairness tries to ensure that AI

systems do not reinforce negative stereotypes or render certain groups invisible. Procedural fairness focuses on the decision process rather than the outcome. The most relevant notions of fairness may depend on the context, notably the type of AI application and the stakeholders. The subjectivity in the notions of bias and fairness makes it difficult for companies to operationalize them. Having access to sensitive attributes such as race or gender is also considered by many AI ethicists to be necessary in order to compensate for biases, but it may conflict with anti-discrimination laws. At its 2022 Conference on Fairness, Accountability, and Transparency (ACM FAccT 2022), the Association for Computing Machinery, in Seoul, South Korea, presented and published findings that recommend that until AI and robotics systems are demonstrated to be free of bias mistakes, they are unsafe, and the use of self-learning neural networks trained on vast, unregulated sources of flawed internet data should be curtailed.

Lack of transparency

Many AI systems are so complex that their designers cannot explain how they reach their decisions. Particularly with deep neural networks, in which there are many non-linear relationships between inputs and outputs. But some popular explainability techniques exist.

It is impossible to be certain that a program is operating correctly if no one knows how exactly it works. There have been many cases where a machine learning program passed rigorous tests, but nevertheless learned something different than what the programmers intended. For example, a system that could identify skin diseases better than medical professionals was found to actually have a strong tendency to classify images with a ruler as "cancerous", because pictures of malignancies typically include a ruler to show the scale. Another machine learning system designed to help effectively allocate medical resources was found to classify patients with asthma as being at "low risk" of dying from pneumonia. Having asthma is actually a severe risk factor, but since the patients having asthma would usually get much more medical care, they were relatively unlikely to die according to the training data. The correlation between asthma and low risk of dying from pneumonia was real, but misleading.

People who have been harmed by an algorithm's decision have a right to an explanation. Doctors, for example, are expected to clearly and completely explain to their colleagues the reasoning behind any decision they make.

Early drafts of the European Union's General Data Protection Regulation in 2016 included an explicit statement that this right exists. Industry experts noted that this is an unsolved problem with no solution in sight. Regulators argued that nevertheless the harm is real: if the problem has no solution, the tools should not be used.

DARPA established the XAI ("Explainable Artificial Intelligence") program in 2014 to try to solve these problems.

Several approaches aim to address the transparency problem. SHAP enables to visualise the contribution of each feature to the output. LIME can locally approximate a model's outputs with a simpler, interpretable model. Multitask learning provides a large number of outputs in addition to the target classification. These other outputs can help developers deduce what the network has learned. Deconvolution, DeepDream and other generative methods can allow developers to see what different layers of a deep network for computer vision have learned, and produce output that can suggest what the network is learning. For generative pre-trained transformers, Anthropic developed a technique based on dictionary learning that associates patterns of neuron activations with human-understandable concepts.

Bad actors and weaponized AI

Artificial intelligence provides a number of tools that are useful to bad actors, such as authoritarian governments, terrorists, criminals or rogue states.

A lethal autonomous weapon is a machine that locates, selects and engages human targets without human supervision. Widely available AI tools can be used by bad actors to develop inexpensive autonomous weapons and, if produced at scale, they are potentially weapons of mass destruction. Even when used in conventional warfare, they currently cannot reliably choose targets and could potentially kill an innocent person. In 2014, 30 nations (including China) supported a ban on autonomous weapons under the United Nations' Convention on Certain Conventional Weapons, however the United States and others disagreed. By 2015, over fifty countries were reported to be researching battlefield robots.

AI tools make it easier for authoritarian governments to efficiently control their citizens in several ways. Face and voice recognition allow widespread surveillance. Machine learning, operating this data, can classify potential enemies of the state and prevent them from hiding. Recommendation systems can precisely target propaganda and misinformation for maximum

effect. Deepfakes and generative AI aid in producing misinformation. Advanced AI can make authoritarian centralized decision-making more competitive than liberal and decentralized systems such as markets. It lowers the cost and difficulty of digital warfare and advanced spyware. All these technologies have been available since 2020 or earlier—AI facial recognition systems are already being used for mass surveillance in China. There are many other ways in which AI is expected to help bad actors, some of which can not be foreseen. For example, machine-learning AI is able to design tens of thousands of toxic molecules in a matter of hours.

Technological unemployment

Economists have frequently highlighted the risks of redundancies from AI, and speculated about unemployment if there is no adequate social policy for full employment.

In the past, technology has tended to increase rather than reduce total employment, but economists acknowledge that "we're in uncharted territory" with AI. A survey of economists showed disagreement about whether the increasing use of robots and AI will cause a substantial increase in long-term unemployment, but they generally agree that it could be a net benefit if productivity gains are redistributed.

Risk estimates vary; for example, in the 2010s, Michael Osborne and Carl Benedikt Frey estimated 47% of U.S. jobs are at "high risk" of potential automation, while an OECD report classified only 9% of U.S. jobs as "high risk". The methodology of speculating about future employment levels has been criticised as lacking evidential foundation, and for implying that technology, rather than social policy, creates unemployment, as opposed to redundancies. In April 2023, it was reported that 70% of the jobs for Chinese video game illustrators had been eliminated by generative artificial intelligence.

Unlike previous waves of automation, many middle-class jobs may be eliminated by artificial intelligence; The Economist stated in 2015 that "the worry that AI could do to white-collar jobs what steam power did to blue-collar ones during the Industrial Revolution" is "worth taking seriously". Jobs at extreme risk range from paralegals to fast food cooks, while job demand is likely to increase for care-related professions ranging from personal healthcare to the clergy.

From the early days of the development of artificial intelligence, there have

been arguments, for example, those put forward by Joseph Weizenbaum, about whether tasks that can be done by computers actually should be done by them, given the difference between computers and humans, and between quantitative calculation and qualitative, value-based judgement.

Existential risk

It has been argued AI will become so powerful that humanity may irreversibly lose control of it. This could, as physicist Stephen Hawking stated, "spell the end of the human race". This scenario has been common in science fiction, when a computer or robot suddenly develops a human-like "self-awareness" (or "sentience" or "consciousness") and becomes a malevolent character.

These sci-fi scenarios are misleading in several ways.

First, AI does not require human-like sentience to be an existential risk. Modern AI programs are given specific goals and use learning and intelligence to achieve them. Philosopher Nick Bostrom argued that if one gives almost any goal to a sufficiently powerful AI, it may choose to destroy humanity to achieve it (he used the example of a paperclip maximizer). Stuart Russell gives the example of household robot that tries to find a way to kill its owner to prevent it from being unplugged, reasoning that "you can't fetch the coffee if you're dead." In order to be safe for humanity, a superintelligence would have to be genuinely aligned with humanity's morality and values so that it is "fundamentally on our side".

Second, Yuval Noah Harari argues that AI does not require a robot body or physical control to pose an existential risk. The essential parts of civilization are not physical. Things like ideologies, law, government, money and the economy are built on language; they exist because there are stories that billions of people believe.

The current prevalence of misinformation suggests that an AI could use language to convince people to believe anything, even to take actions that are destructive.

The opinions amongst experts and industry insiders are mixed, with sizable fractions both concerned and unconcerned by risk from eventual superintelligent AI. Personalities such as Stephen Hawking, Bill Gates, and Elon Musk, as well as AI pioneers such as Yoshua Bengio, Stuart Russell, Demis Hassabis, and Sam Altman, have expressed concerns about existential risk from AI.

In May 2023, Geoffrey Hinton announced his resignation from Google in order to be able to "freely speak out about the risks of AI" without "considering how this impacts Google". He notably mentioned risks of an AI takeover, and stressed that in order to avoid the worst outcomes, establishing safety guidelines will require cooperation among those competing in use of AI.

In 2023, many leading AI experts endorsed the joint statement that "Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war".

Some other researchers were more optimistic. AI pioneer Jürgen Schmidhuber did not sign the joint statement, emphasising that in 95% of all cases, AI research is about making "human lives longer and healthier and easier."

While the tools that are now being used to improve lives can also be used by bad actors, "they can also be used against the bad actors." Andrew Ng also argued that "it's a mistake to fall for the doomsday hype on AI—and that regulators who do will only benefit vested interests." Yann LeCun "scoffs at his peers' dystopian scenarios of supercharged misinformation and even, eventually, human extinction." In the early 2010s, experts argued that the risks are too distant in the future to warrant research or that humans will be valuable from the perspective of a superintelligent machine. However, after 2016, the study of current and future risks and possible solutions became a serious area of research.

Ethical machines and alignment

Friendly AI are machines that have been designed from the beginning to minimize risks and to make choices that benefit humans.

Eliezer Yudkowsky, who coined the term, argues that developing friendly AI should be a higher research priority: it may require a large investment and it must be completed before AI becomes an existential risk.

Machines with intelligence have the potential to use their intelligence to make ethical decisions. The field of machine ethics provides machines with ethical principles and procedures for resolving ethical dilemmas.

The field of machine ethics is also called computational morality, and was founded at an AAAI symposium in 2005.

Other approaches include Wendell Wallach's "artificial moral agents" and Stuart J. Russell's three principles for developing provably beneficial

machines.

Open source

Active organizations in the AI open-source community include Hugging Face, Google, EleutherAI and Meta. Various AI models, such as Llama 2, Mistral or Stable Diffusion, have been made open-weight, meaning that their architecture and trained parameters (the "weights") are publicly available. Open-weight models can be freely fine-tuned, which allows companies to specialize them with their own data and for their own use-case.

Open-weight models are useful for research and innovation but can also be misused. Since they can be fine-tuned, any built-in security measure, such as objecting to harmful requests, can be trained away until it becomes ineffective. Some researchers warn that future AI models may develop dangerous capabilities (such as the potential to drastically facilitate bioterrorism) and that once released on the Internet, they cannot be deleted everywhere if needed. They recommend pre-release audits and cost-benefit analyses.

Frameworks

Artificial intelligence projects can be guided by ethical considerations during the design, development, and implementation of an AI system. An AI framework such as the Care and Act Framework, developed by the Alan Turing Institute and based on the SUM values, outlines four main ethical dimensions, defined as follows:

Respect the dignity of individual people

Connect with other people sincerely, openly, and inclusively

Care for the wellbeing of everyone

Protect social values, justice, and the public interest

Other developments in ethical frameworks include those decided upon during the Asilomar Conference, the Montreal Declaration for Responsible AI, and the IEEE's Ethics of Autonomous Systems initiative, among others; however, these principles are not without criticism, especially regarding the people chosen to contribute to these frameworks.

Promotion of the wellbeing of the people and communities that these technologies affect requires consideration of the social and ethical implications at all stages of AI system design, development and

implementation, and collaboration between job roles such as data scientists, product managers, data engineers, domain experts, and delivery managers. The UK AI Safety Institute released in 2024 a testing toolset called 'Inspect' for AI safety evaluations available under an MIT open-source licence which is freely available on GitHub and can be improved with third-party packages. It can be used to evaluate AI models in a range of areas including core knowledge, ability to reason, and autonomous capabilities.

Regulation

The regulation of artificial intelligence is the development of public sector policies and laws for promoting and regulating AI; it is therefore related to the broader regulation of algorithms. The regulatory and policy landscape for AI is an emerging issue in jurisdictions globally. According to AI Index at Stanford, the annual number of AI-related laws passed in the 127 survey countries jumped from one passed in 2016 to 37 passed in 2022 alone.

Between 2016 and 2020, more than 30 countries adopted dedicated strategies for AI. Most EU member states had released national AI strategies, as had Canada, China, India, Japan, Mauritius, the Russian Federation, Saudi Arabia, United Arab Emirates, U.S., and Vietnam. Others were in the process of elaborating their own AI strategy, including Bangladesh, Malaysia and Tunisia. The Global Partnership on Artificial Intelligence was launched in June 2020, stating a need for AI to be developed in accordance with human rights and democratic values, to ensure public confidence and trust in the technology. Henry Kissinger, Eric Schmidt, and Daniel Huttenlocher published a joint statement in November 2021 calling for a government commission to regulate AI.

In 2023, OpenAI leaders published recommendations for the governance of superintelligence, which they believe may happen in less than 10 years. In 2023, the United Nations also launched an advisory body to provide recommendations on AI governance; the body comprises technology company executives, government officials and academics. In 2024, the Council of Europe created the first international legally binding treaty on AI, called the "Framework Convention on Artificial Intelligence and Human Rights, Democracy and the Rule of Law". It was adopted by the European Union, the United States, the United Kingdom, and other signatories.

In a 2022 Ipsos survey, attitudes towards AI varied greatly by country; 78% of Chinese citizens, but only 35% of Americans, agreed that "products and

services using AI have more benefits than drawbacks". A 2023 Reuters/Ipsos poll found that 61% of Americans agree, and 22% disagree, that AI poses risks to humanity. In a 2023 Fox News poll, 35% of Americans thought it "very important", and an additional 41% thought it "somewhat important", for the federal government to regulate AI, versus 13% responding "not very important" and 8% responding "not at all important". In November 2023, the first global AI Safety Summit was held in Bletchley Park in the UK to discuss the near and far term risks of AI and the possibility of mandatory and voluntary regulatory frameworks. 28 countries including the United States, China, and the European Union issued a declaration at the start of the summit, calling for international co-operation to manage the challenges and risks of artificial intelligence. In May 2024 at the AI Seoul Summit, 16 global AI tech companies agreed to safety commitments on the development of AI.

History

The study of mechanical or "formal" reasoning began with philosophers and mathematicians in antiquity. The study of logic led directly to Alan Turing's theory of computation, which suggested that a machine, by shuffling symbols as simple as "0" and "1", could simulate any conceivable form of mathematical reasoning. This, along with concurrent discoveries in cybernetics, information theory and neurobiology, led researchers to consider the possibility of building an "electronic brain". They developed several areas of research that would become part of AI, such as McCulloch and Pitts design for "artificial neurons" in 1943, and Turing's influential 1950 paper 'Computing Machinery and Intelligence', which introduced the Turing test and showed that "machine intelligence" was plausible.

The field of AI research was founded at a workshop at Dartmouth College in 1956. The attendees became the leaders of AI research in the 1960s. They and their students produced programs that the press described as "astonishing": computers were learning checkers strategies, solving word problems in algebra, proving logical theorems and speaking English. Artificial intelligence laboratories were set up at a number of British and U.S. universities in the latter 1950s and early 1960s. Researchers in the 1960s and the 1970s were convinced that their methods

would eventually succeed in creating a machine with general intelligence and considered this the goal of their field. In 1965 Herbert Simon predicted, "machines will be capable, within twenty years, of doing any work a man can do". In 1967 Marvin Minsky agreed, writing that "within a generation ... the problem of creating 'artificial intelligence' will substantially be solved". They had, however, underestimated the difficulty of the problem. In 1974, both the U.S. and British governments cut off exploratory research in response to the criticism of Sir James Lighthill and ongoing pressure from the U.S. Congress to fund more productive projects. Minsky and Papert's book *Perceptrons* was understood as proving that artificial neural networks would never be useful for solving real-world tasks, thus discrediting the approach altogether. The "AI winter", a period when obtaining funding for AI projects was difficult, followed.

In the early 1980s, AI research was revived by the commercial success of expert systems, a form of AI program that simulated the knowledge and analytical skills of human experts. By 1985, the market for AI had reached over a billion dollars. At the same time, Japan's fifth generation computer project inspired the U.S. and British governments to restore funding for academic research.

However, beginning with the collapse of the Lisp Machine market in 1987, AI once again fell into disrepute, and a second, longer-lasting winter began.

Up to this point, most of AI's funding had gone to projects that used high-level symbols to represent mental objects like plans, goals, beliefs, and known facts. In the 1980s, some researchers began to doubt that this approach would be able to imitate all the processes of human cognition, especially perception, robotics, learning and pattern recognition, and began to look into "sub-symbolic" approaches. Rodney Brooks rejected

"representation" in general and focussed directly on engineering machines that move and survive. Judea Pearl, Lotfi Zadeh, and others developed

methods that handled incomplete and uncertain information by making reasonable guesses rather than precise logic. But the most important

development was the revival of "connectionism", including neural network research, by Geoffrey Hinton and others. In 1990, Yann LeCun successfully showed that convolutional neural networks can recognize handwritten digits, the first of many successful applications of neural networks.

AI gradually restored its reputation in the late 1990s and early 21st century by exploiting formal mathematical methods and by finding specific solutions

to specific

problems. This "narrow" and "formal" focus allowed researchers to produce verifiable results and collaborate with other fields (such as statistics, economics and mathematics). By 2000, solutions developed by AI researchers were being widely used, although in the 1990s they were rarely described as "artificial intelligence" (a tendency known as the AI effect). However, several academic researchers became concerned that AI was no longer pursuing its original goal of creating versatile, fully intelligent machines. Beginning around 2002, they founded the subfield of artificial general intelligence (or "AGI"), which had several well-funded institutions by the 2010s.

Deep learning began to dominate industry benchmarks in 2012 and was adopted throughout the field.

For many specific tasks, other methods were abandoned.

Deep learning's success was based on both hardware improvements (faster computers, graphics processing units, cloud computing) and access to large amounts of data (including curated datasets, such as ImageNet). Deep learning's success led to an enormous increase in interest and funding in AI. The amount of machine learning research (measured by total publications) increased by 50% in the years 2015–2019.

In 2016, issues of fairness and the misuse of technology were catapulted into center stage at machine learning conferences, publications vastly increased, funding became available, and many researchers re-focussed their careers on these issues. The alignment problem became a serious field of academic study.

In the late 2010s and early 2020s, AGI companies began to deliver programs that created enormous interest. In 2015, AlphaGo, developed by DeepMind, beat the world champion Go player. The program taught only the game's rules and developed a strategy by itself. GPT-3 is a large language model that was released in 2020 by OpenAI and is capable of generating high-quality human-like text. ChatGPT, launched on November 30, 2022, became the fastest-growing consumer software application in history, gaining over 100 million users in two months. It marked what is widely regarded as AI's breakout year, bringing it into the public consciousness. These programs, and others, inspired an aggressive AI boom, where large companies began investing billions of dollars in AI research. According to AI Impacts, about

US\$50 billion annually was invested in "AI" around 2022 in the U.S. alone and about 20% of the new U.S. Computer Science PhD graduates have specialized in "AI". About 800,000 "AI"-related U.S. job openings existed in 2022. According to PitchBook research, 22% of newly funded startups in 2024 claimed to be AI companies.

Philosophy

Philosophical debates have historically sought to determine the nature of intelligence and how to make intelligent machines. Another major focus has been whether machines can be conscious, and the associated ethical implications. Many other topics in philosophy are relevant to AI, such as epistemology and free will. Rapid advancements have intensified public discussions on the philosophy and ethics of AI.

Defining artificial intelligence

Alan Turing wrote in 1950 "I propose to consider the question 'can machines think?'" He advised changing the question from whether a machine "thinks", to "whether or not it is possible for machinery to show intelligent behaviour". He devised the Turing test, which measures the ability of a machine to simulate human conversation.

Since we can only observe the behavior of the machine, it does not matter if it is "actually" thinking or literally has a "mind".

Turing notes that we can not determine these things about other people but "it is usual to have a polite convention that everyone thinks."

Russell and Norvig agree with Turing that intelligence must be defined in terms of external behavior, not internal structure. However, they are critical that the test requires the machine to imitate humans. "Aeronautical engineering texts", they wrote, "do not define the goal of their field as making 'machines that fly so exactly like pigeons that they can fool other pigeons.'" AI founder John McCarthy agreed, writing that "Artificial intelligence is not, by definition, simulation of human intelligence". McCarthy defines intelligence as "the computational part of the ability to achieve goals in the world". Another AI founder, Marvin Minsky, similarly describes it as "the ability to solve hard problems". The leading AI textbook defines it as the study of agents that perceive their environment and take actions that maximize their chances of achieving defined goals. These

definitions view intelligence in terms of well-defined problems with well-defined solutions, where both the difficulty of the problem and the performance of the program are direct measures of the "intelligence" of the machine—and no other philosophical discussion is required, or may not even be possible.

Another definition has been adopted by Google, a major practitioner in the field of AI. This definition stipulates the ability of systems to synthesize information as the manifestation of intelligence, similar to the way it is defined in biological intelligence.

Some authors have suggested in practice, that the definition of AI is vague and difficult to define, with contention as to whether classical algorithms should be categorised as AI, with many companies during the early 2020s AI boom using the term as a marketing buzzword, often even if they did "not actually use AI in a material way".

There has been debate over whether large language models exhibit genuine intelligence or merely simulate it by imitating human text.

Evaluating approaches to AI

No established unifying theory or paradigm has guided AI research for most of its history. The unprecedented success of statistical machine learning in the 2010s eclipsed all other approaches (so much so that some sources, especially in the business world, use the term "artificial intelligence" to mean "machine learning with neural networks"). This approach is mostly sub-symbolic, soft and narrow.

Critics argue that these questions may have to be revisited by future generations of AI researchers.

Symbolic AI and its limits

Symbolic AI (or "GOFAI") simulated the high-level conscious reasoning that people use when they solve puzzles, express legal reasoning and do mathematics. They were highly successful at "intelligent" tasks such as algebra or IQ tests. In the 1960s, Newell and Simon proposed the physical symbol systems hypothesis: "A physical symbol system has the necessary and sufficient means of general intelligent action."

However, the symbolic approach failed on many tasks that humans solve easily, such as learning, recognizing an object or commonsense reasoning. Moravec's paradox is the discovery that high-level "intelligent" tasks were

easy for AI, but low level "instinctive" tasks were extremely difficult.

Philosopher Hubert Dreyfus had argued since the 1960s that human expertise depends on unconscious instinct rather than conscious symbol manipulation, and on having a "feel" for the situation, rather than explicit symbolic knowledge. Although his arguments had been ridiculed and ignored when they were first presented, eventually, AI research came to agree with him.

The issue is not resolved: sub-symbolic reasoning can make many of the same inscrutable mistakes that human intuition does, such as algorithmic bias. Critics such as Noam Chomsky argue continuing research into symbolic AI will still be necessary to attain general intelligence, in part because sub-symbolic AI is a move away from explainable AI: it can be difficult or impossible to understand why a modern statistical AI program made a particular decision. The emerging field of neuro-symbolic artificial intelligence attempts to bridge the two approaches.

Neat vs. scruffy

"Neats" hope that intelligent behavior is described using simple, elegant principles (such as logic, optimization, or neural networks). "Scruffies" expect that it necessarily requires solving a large number of unrelated problems. Neats defend their programs with theoretical rigor, scruffies rely mainly on incremental testing to see if they work. This issue was actively discussed in the 1970s and 1980s, but eventually was seen as irrelevant. Modern AI has elements of both.

Soft vs. hard computing

Finding a provably correct or optimal solution is intractable for many important problems. Soft computing is a set of techniques, including genetic algorithms, fuzzy logic and neural networks, that are tolerant of imprecision, uncertainty, partial truth and approximation. Soft computing was introduced in the late 1980s and most successful AI programs in the 21st century are examples of soft computing with neural networks.

Narrow vs. general AI

AI researchers are divided as to whether to pursue the goals of artificial general intelligence and superintelligence directly or to solve as many specific problems as possible (narrow AI) in hopes these solutions will lead indirectly

to the field's long-term goals. General intelligence is difficult to define and difficult to measure, and modern AI has had more verifiable successes by focusing on specific problems with specific solutions. The sub-field of artificial general intelligence studies this area exclusively.

Machine consciousness, sentience, and mind

There is no settled consensus in philosophy of mind on whether a machine can have a mind, consciousness and mental states in the same sense that human beings do. This issue considers the internal experiences of the machine, rather than its external behavior.

Mainstream AI research considers this issue irrelevant because it does not affect the goals of the field: to build machines that can solve problems using intelligence. Russell and Norvig add that "[t]he additional project of making a machine conscious in exactly the way humans are is not one that we are equipped to take on." However, the question has become central to the philosophy of mind. It is also typically the central question at issue in artificial intelligence in fiction.

Consciousness

David Chalmers identified two problems in understanding the mind, which he named the "hard" and "easy" problems of consciousness. The easy problem is understanding how the brain processes signals, makes plans and controls behavior. The hard problem is explaining how this feels or why it should feel like anything at all, assuming we are right in thinking that it truly does feel like something (Dennett's consciousness illusionism says this is an illusion). While human information processing is easy to explain, human subjective experience is difficult to explain. For example, it is easy to imagine a color-blind person who has learned to identify which objects in their field of view are red, but it is not clear what would be required for the person to know what red looks like.

Computationalism and functionalism

Computationalism is the position in the philosophy of mind that the human mind is an information processing system and that thinking is a form of computing.

Computationalism argues that the relationship between mind and body is

similar or identical to the relationship between software and hardware and thus may be a solution to the mind-body problem. This philosophical position was inspired by the work of AI researchers and cognitive scientists in the 1960s and was originally proposed by philosophers Jerry Fodor and Hilary Putnam.

Philosopher John Searle characterized this position as "strong AI": "The appropriately programmed computer with the right inputs and outputs would thereby have a mind in exactly the same sense human beings have minds."

Searle challenges this claim with his Chinese room argument, which attempts to show that even a computer capable of perfectly simulating human behavior would not have a mind.

AI welfare and rights

It is difficult or impossible to reliably evaluate whether an advanced AI is sentient (has the ability to feel), and if so, to what degree. But if there is a significant chance that a given machine can feel and suffer, then it may be entitled to certain rights or welfare protection measures, similarly to animals. Sapience (a set of capacities related to high intelligence, such as discernment or self-awareness) may provide another moral basis for AI rights. Robot rights are also sometimes proposed as a practical way to integrate autonomous agents into society.

In 2017, the European Union considered granting "electronic personhood" to some of the most capable AI systems. Similarly to the legal status of companies, it would have conferred rights but also responsibilities. Critics argued in 2018 that granting rights to AI systems would downplay the importance of human rights, and that legislation should focus on user needs rather than speculative futuristic scenarios. They also noted that robots lacked the autonomy to take part in society on their own.

Progress in AI increased interest in the topic. Proponents of AI welfare and rights often argue that AI sentience, if it emerges, would be particularly easy to deny. They warn that this may be a moral blind spot analogous to slavery or factory farming, which could lead to large-scale suffering if sentient AI is created and carelessly exploited.

Future

Superintelligence and the singularity

A superintelligence is a hypothetical agent that would possess intelligence far

surpassing that of the brightest and most gifted human mind. If research into artificial general intelligence produced sufficiently intelligent software, it might be able to reprogram and improve itself. The improved software would be even better at improving itself, leading to what I. J. Good called an "intelligence explosion" and Vernor Vinge called a "singularity".

However, technologies cannot improve exponentially indefinitely, and typically follow an S-shaped curve, slowing when they reach the physical limits of what the technology can do.

Transhumanism

Robot designer Hans Moravec, cyberneticist Kevin Warwick and inventor Ray Kurzweil have predicted that humans and machines may merge in the future into cyborgs that are more capable and powerful than either. This idea, called transhumanism, has roots in the writings of Aldous Huxley and Robert Ettinger.

Edward Fredkin argues that "artificial intelligence is the next step in evolution", an idea first proposed by Samuel Butler's "Darwin among the Machines" as far back as 1863, and expanded upon by George Dyson in his 1998 book *Darwin Among the Machines: The Evolution of Global Intelligence*.

In fiction

Thought-capable artificial beings have appeared as storytelling devices since antiquity, and have been a persistent theme in science fiction.

A common trope in these works began with Mary Shelley's *Frankenstein*, where a human creation becomes a threat to its masters. This includes such works as Arthur C. Clarke's and Stanley Kubrick's *2001: A Space Odyssey* (both 1968), with HAL 9000, the murderous computer in charge of the *Discovery One* spaceship, as well as *The Terminator* (1984) and *The Matrix* (1999). In contrast, the rare loyal robots such as Gort from *The Day the Earth Stood Still* (1951) and Bishop from *Aliens* (1986) are less prominent in popular culture.

Isaac Asimov introduced the Three Laws of Robotics in many stories, most notably with the "Multivac" super-intelligent computer. Asimov's laws are often brought up during lay discussions of machine ethics; while almost all artificial intelligence researchers are familiar with Asimov's laws through popular culture, they generally consider the laws useless for many reasons,

one of which is their ambiguity.

Several works use AI to force us to confront the fundamental question of what makes us human, showing us artificial beings that have the ability to feel, and thus to suffer. This appears in Karel Čapek's R.U.R., the films A.I. Artificial Intelligence and Ex Machina, as well as the novel Do Androids Dream of Electric Sheep?, by Philip K. Dick. Dick considers the idea that our understanding of human subjectivity is altered by technology created with artificial intelligence.

See also

Artificial consciousness – Field in cognitive science

Artificial intelligence and elections – Use and impact of AI on political elections

Artificial intelligence content detection – Software to detect AI-generated content

Association for the Advancement of Artificial Intelligence (AAAI)

Behavior selection algorithm – Algorithm that selects actions for intelligent agents

Business process automation – Automation of business processes

Case-based reasoning – Process of solving new problems based on the solutions of similar past problems

Computational intelligence – Ability of a computer to learn a specific task from data or experimental observation

Digital immortality – Hypothetical concept of storing a personality in digital form

Emergent algorithm – Algorithm exhibiting emergent behavior

Female gendering of AI technologies – Gender biases in digital technology

Pages displaying short descriptions of redirect targets
Glossary of artificial intelligence – List of definitions of terms and concepts commonly used in the study of artificial intelligence

Intelligence amplification – Use of information technology to augment human intelligence

Intelligent agent – Software agent which acts autonomously

Intelligent automation – Software process that combines robotic process automation and artificial intelligence

List of artificial intelligence journals

List of artificial intelligence projects

Mind uploading – Hypothetical process of digitally emulating a brain

Organoid intelligence – Use of brain cells and brain organoids for intelligent computing

Robotic process automation – Form of business process automation technology

The Last Day – 1967 Welsh science fiction novel

Wetware computer – Computer composed of organic material

DARWIN EU - A European Union initiative coordinated by the European Medicines Agency (EMA) to generate and utilize real-world evidence (RWE) to support the evaluation and supervision of medicines across the EU.

Explanatory notes

References

AI textbooks

The two most widely used textbooks in 2023 (see the Open Syllabus):

Russell, Stuart J.; Norvig, Peter (2021). Artificial Intelligence: A Modern Approach (4th ed.). Hoboken: Pearson. ISBN 978-0-1346-1099-3. LCCN 20190474.

Rich, Elaine; Knight, Kevin; Nair, Shivashankar B (2010). Artificial Intelligence (3rd ed.). New Delhi: Tata McGraw Hill India. ISBN 978-0-0700-8770-5.

The four most widely used AI textbooks in 2008:

Other textbooks:

Ertel, Wolfgang (2017). Introduction to Artificial Intelligence (2nd ed.). Springer. ISBN 978-3-3195-8486-7.

Ciaramella, Alberto; Ciaramella, Marco (2024). Introduction to Artificial Intelligence: from data analysis to generative AI (1st ed.). Intellisemantic Editions. ISBN 978-8-8947-8760-3.

History of AI

Other sources

Further reading

External links

"Artificial Intelligence". Internet Encyclopedia of Philosophy.