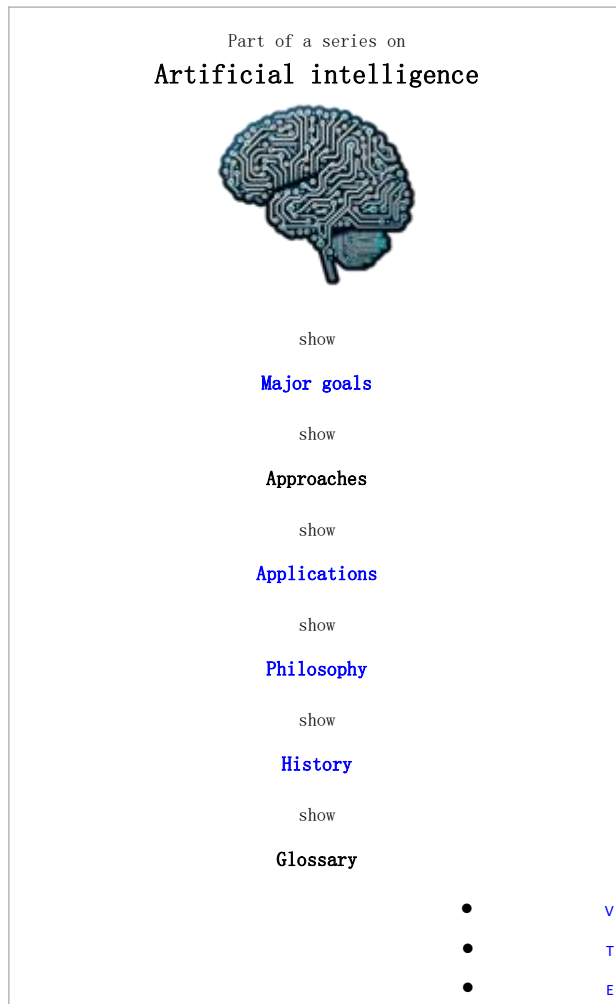


Artificial intelligence

151 languages

"AI" redirects here. For other uses, see [AI \(disambiguation\)](#), [Artificial intelligence \(disambiguation\)](#), and [Intelligent agent](#).



Artificial intelligence (AI), in its broadest sense, is *intelligence* exhibited by *machines*, particularly *computer systems*. It is a *field of research* in *computer science* that develops and studies methods and *software* that enable machines to *perceive their environment* and uses *learning* and intelligence to take actions that maximize their chances of achieving defined goals.^[1] Such machines may be called AIs.

AI technology is widely used throughout industry, government, and science. Some high-profile applications include advanced web search engines (e.g., Google Search); recommendation systems (used by YouTube, Amazon, and Netflix); interacting via human speech (e.g., Google Assistant, Siri, and Alexa); autonomous vehicles (e.g., Waymo); generative and creative tools (e.g., ChatGPT and AI art); and superhuman play and analysis in strategy games (e.g., chess and Go).^[2] 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."^{[3][4]}

Alan Turing was the first person to conduct substantial research in the field that he called machine intelligence.^[5] Artificial intelligence was founded as an academic discipline in 1956.^[6] The field went through multiple cycles of optimism,^{[7][8]} followed by periods of disappointment and loss of funding, known as AI winter.^{[9][10]} Funding and interest vastly increased after 2012 when deep learning surpassed all previous AI techniques,^[11] and after 2017 with the transformer architecture.^[12] This led to the AI boom of the early 2020s, with companies, universities, and laboratories overwhelmingly based in the United States pioneering significant advances in artificial intelligence.^[13]

The growing use of artificial intelligence in the 21st century is influencing a societal and economic shift towards increased automation, data-driven decision-making, and the integration of AI systems into various economic sectors and areas of life, impacting job markets, healthcare, government, industry, and education. This raises questions about the long-term effects, ethical implications, and risks of AI, prompting discussions

about [regulatory policies](#) to ensure the [safety and benefits of the technology](#).

The various subfields of AI research are centered around particular goals and the use of particular tools. The traditional goals of AI research include [reasoning](#), [knowledge representation](#), [planning](#), [learning](#), [natural language processing](#), [perception](#), and support for [robotics](#).^[a] [General intelligence](#)—the ability to complete any task performable by a human on an at least equal level—is among the field's long-term goals.^[14]

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](#).^[b] AI also draws upon [psychology](#), [linguistics](#), [philosophy](#), [neuroscience](#), and other fields.^[15]

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.^[a]

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](#).^[16] By the late 1980s and 1990s, methods were developed for dealing with [uncertain](#) or incomplete information, employing concepts from [probability](#) and [economics](#).^[17]

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.^[18] 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.^[19] Accurate and efficient reasoning is an unsolved problem.

Knowledge representation



An ontology represents knowledge as a set of concepts within a domain and the relationships between those concepts.

Knowledge representation and knowledge engineering^[20] 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,^[21] scene interpretation,^[22] clinical decision support,^[23] knowledge discovery (mining "interesting" and actionable inferences from large databases),^[24] and other areas.^[25]

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.^[26] Knowledge bases need to represent things such as

objects, properties, categories, and relations between objects;^[27] situations, events, states, and time;^[28] causes and effects;^[29] knowledge about knowledge (what we know about what other people know);^[30] *default reasoning* (things that humans assume are true until they are told differently and will remain true even when other facts are changing);^[31] 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);^[32] 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).^[19] There is also the difficulty of *knowledge acquisition*, the problem of obtaining knowledge for AI applications.^[c]

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.^{[d][35]} In *automated planning*, the agent has a specific goal.^[36] 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.^[37]

In *classical planning*, the agent knows exactly what the effect of any action will be.^[38] 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.^[39]

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.^[40] [Information value theory](#) can be used to weigh the value of exploratory or experimental actions.^[41] 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.^[42]

[Game theory](#) describes the rational behavior of multiple interacting agents and is used in AI programs that make decisions that involve other agents.^[43]

Learning

[Machine learning](#) is the study of programs that can improve their performance on a given task automatically.^[44] It has been a part of AI from the beginning.^[6]

There are several kinds of machine learning. [Unsupervised learning](#) analyzes a stream of data and finds patterns and makes

predictions without any other guidance.^[47] *Supervised learning* requires a human to label the input data first, 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).^[48]

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".^[49] *Transfer learning* is when the knowledge gained from one problem is applied to a new problem.^[50] *Deep learning* is a type of machine learning that runs inputs through biologically inspired *artificial neural networks* for all of these types of learning.^[51]

Computational learning theory can assess learners by *computational complexity*, by *sample complexity* (how much data is required), or by other notions of *optimization*.^[52]

Natural language processing

Natural language processing (NLP)^[53] allows programs to read, write and communicate in human languages such as *English*. Specific problems include *speech recognition*, *speech synthesis*, *machine translation*, *information extraction*, *information retrieval* and *question answering*.^[54]

Early work, based on *Noam Chomsky's generative grammar* and *semantic networks*, had difficulty with *word-sense disambiguation*^[7] unless restricted to small domains called "*micro-worlds*" (due to the common sense knowledge problem^[32]). *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),^[55] *transformers* (a deep learning architecture using an *attention* mechanism),^[56] and others.^[57] In 2019, *generative pre-trained transformer* (or "GPT") language models began to generate coherent text,^{[58][59]} 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.^[60]

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.^[61]

The field includes *speech recognition*,^[62] *image classification*,^[63] *facial recognition*, *object recognition*,^[64] and *robotic perception*.^[65]

Social intelligence



Kismet, a robot head which was made in the 1990s; a machine that can recognize and simulate emotions.^[66]

Affective computing is an interdisciplinary umbrella that comprises systems that recognize, interpret, process, or simulate human *feeling, emotion, and mood*.^[67] 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.^[68] Moderate successes related to affective computing include textual *sentiment analysis* and, more recently, *multimodal sentiment analysis*, wherein AI classifies the affects displayed by a videotaped subject.^[69]

General intelligence

A machine with *artificial general intelligence* should be able to solve a wide variety of problems with breadth and versatility similar to human intelligence.^[14]

Techniques

AI research uses a wide variety of techniques to accomplish the goals above.^[6]

Search and optimization

AI can solve many problems by intelligently searching through many possible solutions.^[70] 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.^[71] 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*.^[72]

Simple exhaustive searches^[73] 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.^[18] "[Heuristics](#)" or "rules of thumb" can help prioritize choices that are more likely to reach a goal.^[74]

[Adversarial search](#) is used for [game-playing](#) programs, such as chess or Go. It searches through a [tree](#) of possible moves and counter-moves, looking for a winning position.^[75]

Local search

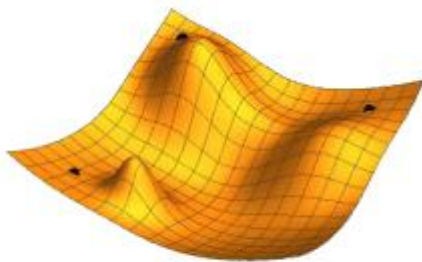


Illustration of [gradient descent](#) for 3 different starting points. Two parameters (represented by the plan coordinates) are adjusted in order to minimize the [loss function](#) (the height).

[Local search](#) uses [mathematical optimization](#) to find a solution to a problem. It begins with some form of guess and refines it incrementally.^[76]

[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.^[77]

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.^[78]

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](#)).^[79]

Logic

Formal [logic](#) is used for [reasoning](#) and [knowledge representation](#).^[80] 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")^[81] 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 X s that are Y s").^[82]

[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](#)).^[83] 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.^[84] 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.^[85]

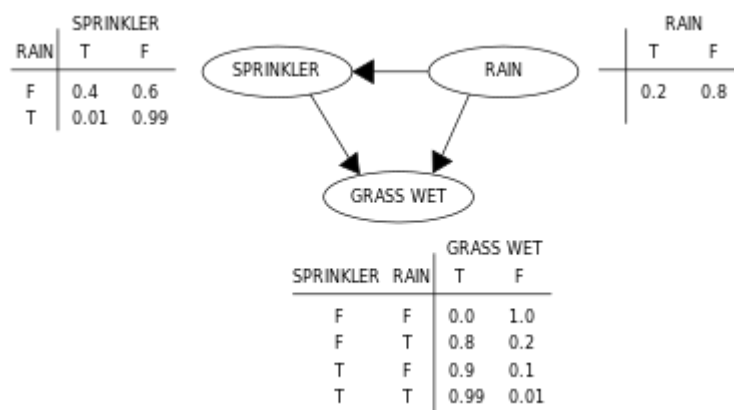
Inference in both Horn clause logic and first-order logic is [undecidable](#), and therefore [intractable](#). However, backward

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.^[86]

[Fuzzy logic](#) assigns a "degree of truth" between 0 and 1. It can therefore handle propositions that are vague and partially true.^[87]

[Non-monotonic logics](#), including logic programming with [negation as failure](#), are designed to handle [default reasoning](#).^[31] Other specialized versions of logic have been developed to describe many complex domains.

Probabilistic methods for uncertain reasoning



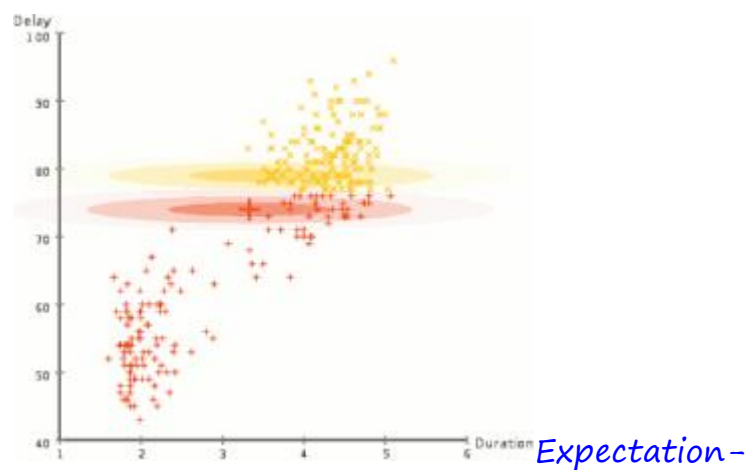
A simple [Bayesian](#)

[network](#), with the associated [conditional probability tables](#)

Many problems in AI (including in 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.^[88] Precise mathematical tools have been developed that analyze how an agent can make choices and plan, using [decision theory](#), [decision analysis](#),^[89] and [information value theory](#).^[90] These tools include models such as [Markov decision processes](#),^[91] dynamic [decision networks](#),^[92] [game theory](#) and [mechanism design](#).^[93]

Bayesian networks^[94] are a tool that can be used for reasoning (using the Bayesian inference algorithm),^{[9][96]} learning (using the expectation-maximization algorithm),^{[4][98]} planning (using decision networks)^[99] and perception (using dynamic Bayesian networks).^[92]

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).^[92]



Expectation-maximization clustering of Old Faithful eruption data starts from a random guess but then successfully converges on an accurate clustering of the two physically distinct modes of eruption.

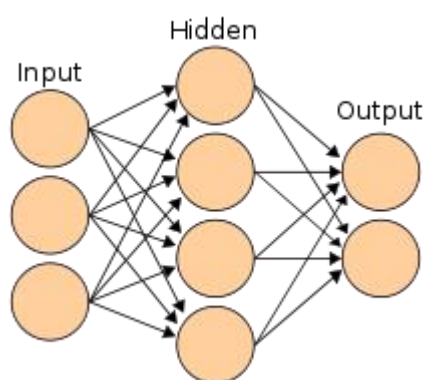
Classifiers and statistical learning methods

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^[100] 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.^[48]

There are many kinds of classifiers in use. The *decision tree* is the simplest and most widely used symbolic machine learning algorithm.^[101] *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.^[102] The *naive Bayes classifier* is reportedly the "most widely used learner"^[103] at Google, due in part to its scalability.^[104] *Neural networks* are also used as classifiers.^[105]

Artificial neural networks



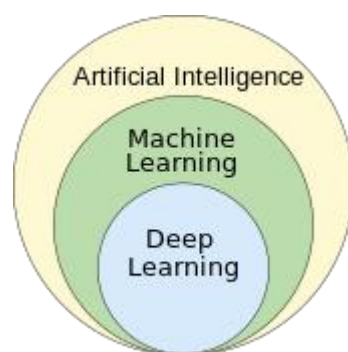
A neural network is an interconnected group of nodes, akin to the vast network of *neurons* in the *human brain*.

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.^[105]

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.^[106] 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.^[107]

In [feedforward neural networks](#) the signal passes in only one direction.^[108] [Recurrent neural networks](#) feed the output signal back into the input, which allows short-term memories of previous input events. [Long short term memory](#) is the most successful network architecture for recurrent networks.^[109] [Perceptrons](#)^[110] use only a single layer of neurons, deep learning^[111] uses multiple layers. [Convolutional neural networks](#) strengthen the connection between neurons that are "close" to each other—this is especially important in [image processing](#), where a local set of neurons must [identify an "edge"](#) before the network can identify an object.^[112]

Deep learning



Deep learning^[111] 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.^[113]

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](#),^[114] and others. The reason that deep learning performs so well in so many applications is not known as of 2023.^[115] 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)^[116] 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](#).^[117]

GPT

[Generative pre-trained transformers](#) (GPT) are [large language models](#) that are based on the semantic relationships between words in sentences ([natural language processing](#)). Text-based GPT models are pre-trained on a large corpus of text which can be from the internet. The pre-training consists in predicting the next [token](#) (a token being usually a word, subword, or punctuation). Throughout this pre-training, 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 still prone to generating falsehoods called "[hallucinations](#)", although this can be reduced with

RLHF and quality data. They are used in [chatbots](#), which allow you to ask a question or request a task in simple text.^{[124][125]}

Current models and services include: [Gemini](#) (formerly Bard), [ChatGPT](#), [Grok](#), [Claude](#), [Copilot](#) and [LLaMA](#).^[126] Multimodal GPT models can process different types of data ([modalities](#)) such as images, videos, sound, and text.^[127]

Specialized hardware and software

Main articles: [Programming languages for artificial intelligence](#) and [Hardware for artificial intelligence](#)

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.^[128] Historically, specialized languages, such as [Lisp](#), [Prolog](#), [Python](#) and others, had been used.

Applications

Main article: [Applications of artificial intelligence](#)

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 Face ID](#) or [Microsoft's DeepFace](#) and [Google's FaceNet](#)) and [image labeling](#) (used by [Facebook](#), [Apple's iPhoto](#) and [TikTok](#)).

Health and medicine

Main article: [Artificial intelligence in healthcare](#)

The application of AI in [medicine](#) and [medical research](#) has the potential to increase patient care and quality of life.^[129] 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.^[130] It has been suggested that AI can overcome discrepancies in funding allocated to different fields of research.^[130] 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](#).^[131] 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.^[132] 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.^{[133][134]}

Games

Main article: [Game artificial intelligence](#)

[Game playing](#) programs have been used since the 1950s to demonstrate and test AI's most advanced techniques.^[135] [Deep](#)

Blue became the first computer chess-playing system to beat a reigning world chess champion, Garry Kasparov, on 11 May 1997.^[136] 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.^[137] 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.^[138] Other programs handle imperfect-information games, such as the poker-playing program Pluribus.^[139] DeepMind developed increasingly generalistic reinforcement learning models, such as with MuZero, which could be trained to play chess, Go, or Atari games.^[140] 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.^[141] 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.^[142]

Military

Main article: Military artificial intelligence

Various countries are deploying AI military applications.^[143] The main applications enhance command and control, communications, sensors, integration and interoperability.^[144] Research is targeting intelligence collection and analysis, logistics, cyber operations, information operations, and semiautonomous and autonomous vehicles.^[143] 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 involving

manned and unmanned teams.^[144] AI was incorporated into military operations in Iraq and Syria.^[143]

In November 2023, US Vice President [Kamala Harris](#) disclosed a declaration signed by 31 nations to set guardrails for the military use of AI. The commitments include using legal reviews to ensure the compliance of military AI with international laws, and being cautious and transparent in the development of this technology.^[145]

Generative AI

Main article: [Generative artificial intelligence](#)



Vincent van Gogh in watercolour
created by generative AI software

In the early 2020s, [generative AI](#) gained widespread prominence. In March 2023, 58% of U.S. adults had heard about [ChatGPT](#) and 14% had tried it.^[146] The increasing realism and ease-of-use of AI-based [text-to-image](#) generators such as [Midjourney](#), [DALL-E](#), and [Stable Diffusion](#) sparked a trend of [viral](#) AI-generated photos. Widespread attention was gained by a fake photo of [Pope Francis](#) wearing a white puffer coat, the fictional arrest of [Donald Trump](#), and a hoax of an attack on the [Pentagon](#), as well as the usage in professional creative arts.^{[147][148]}

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.^[149] A few examples are [energy storage](#), medical diagnosis, military logistics, applications that predict the result of judicial decisions, [foreign policy](#), or supply chain management.

In agriculture, AI has helped farmers identify areas that need irrigation, fertilization, pesticide treatments or increasing yield. 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 for discovering exoplanets, forecasting solar activity, and distinguishing between signals and instrumental effects in gravitational wave astronomy. It could also be used for activities in space such as space exploration, including analysis of data from space missions, real-time science decisions of spacecraft, space debris avoidance, and more autonomous operation.

Ethics

Main article: [Ethics of artificial intelligence](#)

AI has potential benefits and potential risks. AI may be able to advance science and find solutions for serious problems: [Demis Hassabis](#) of [Deep Mind](#) hopes to "solve intelligence, and then use

that to solve everything else".^[150] However, as the use of AI has become widespread, several unintended consequences and risks have been identified.^[151] 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.^[152]

Risks and harm

Privacy and copyright

Further information: [Information privacy](#) and [Artificial intelligence 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](#).

Technology companies collect a wide range of data from their users, including online activity, geolocation data, video and audio.^[153] 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.^[154] 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](#).^[155]

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](#).^[156] 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'".^[157]

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".^{[158][159]} Website owners who do not wish to have their content scraped can indicate it in a "[robots.txt](#)" file.^[160] In 2023, leading authors (including [John Grisham](#) and [Jonathan Franzen](#)) sued AI companies for using their work to train generative AI.^{[161][162]} 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.^[163]

Misinformation

See also: [YouTube & Moderation and offensive content](#)

[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.^[164] This convinced many users that the misinformation was true, and ultimately undermined trust in institutions, the media and the government.^[165] 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 steps to mitigate the problem.

In 2022, [generative AI](#) began to create images, audio, video and text that are indistinguishable from real photographs, recordings, films or human writing. It is possible for bad actors to use this technology to create massive amounts of misinformation or propaganda.^[166] AI pioneer [Geoffrey Hinton](#) expressed concern about AI enabling "authoritarian leaders to manipulate their electorates" on a large scale, among other risks.^[167]

Algorithmic bias and fairness

Main articles: [Algorithmic bias](#) and [Fairness \(machine learning\)](#)

Machine learning applications will be biased if they learn from biased data.^[168] The developers may not be aware that the bias exists.^[169] Bias can be introduced by the way [training data](#) is selected and by the way a model is deployed.^[170]^[168] 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](#).^[171] [Fairness](#) in machine learning is the study of how to prevent the harm caused by algorithmic bias. It has become serious area of academic study within AI. Researchers have discovered it is not always possible to define "fairness" in a way that satisfies all stakeholders.^[172]

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,^[173] a problem called "sample size disparity".^[174] 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.^[175]

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.^[176] In 2017, several researchers^[k] 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.^[178]

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".^[179] Moritz Hardt said "the most robust fact in this research area is that fairness through blindness doesn't work."^[180]

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.^[181] 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 necessarily descriptive and not proscriptive.^[l]

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.^[174]

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.^[183]

Lack of transparency

See also: [Explainable AI](#), [Algorithmic transparency](#), and [Right to explanation](#)



[Lidar](#) testing vehicle for autonomous driving

Many AI systems are so complex that their designers cannot explain how they reach their decisions.^[184] Particularly with [deep neural networks](#), in which there are a large amount of non-[linear](#) relationships between inputs and outputs. But some popular explainability techniques exist.^[185]

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.^[186] 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.^[187]

People who have been harmed by an algorithm's decision have a right to an explanation.^[188] 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.^[m] 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.^[189]

[DARPA](#) established the [XAI](#) ("Explainable Artificial Intelligence") program in 2014 to try and solve these problems.^[190]

There are several possible solutions to the transparency problem. SHAP tried to solve the transparency problems by visualising the contribution of each feature to the output.^[191] LIME can locally approximate a model with a simpler, interpretable model.^[192] [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.^[193] [Deconvolution](#), [DeepDream](#) and other [generative](#) methods can allow developers to see what different layers of a deep network have learned and produce output that can suggest what the network is learning.^[194]

Bad actors and weaponized AI

Main articles: [Lethal autonomous weapon](#), [Artificial intelligence arms race](#), and [AI safety](#)

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.^[n] 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](#).^[196] Even when used in conventional warfare, it is unlikely that they will be unable to reliably choose targets and could potentially [kill an innocent person](#).^[196] 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.^[197] By 2015, over fifty countries were reported to be researching battlefield robots.^[198]

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*.^[199] All these technologies have been available since 2020 or earlier—AI *facial recognition systems* are already being used for *mass surveillance* in China.^{[200][201]}

There many other ways that 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.^[202]

Reliance on industry giants

Training AI systems requires an enormous amount of computing power. Usually only *Big Tech* companies have the financial resources to make such investments. Smaller startups such as *Cohere* and *OpenAI* end up buying access to *data centers* from *Google* and *Microsoft* respectively.^[203]

Technological unemployment

Main articles: Workplace impact of artificial intelligence and 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.^[204]

In the past, technology has tended to increase rather than reduce total employment, but economists acknowledge that "we're in uncharted territory" with AI.^[205] 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*.^[206] 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".^{[10][208]} 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.^[204] In April 2023, it was reported that 70% of the jobs for Chinese video game illustrators had been eliminated by generative artificial intelligence.^{[209][210]}

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".^[211] 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.^[212]

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.^[213]

Existential risk

Main article: [Existential risk from artificial general intelligence](#)

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](#)".^[214] 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.^[p] 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 factory manager](#)).^[216] [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."^[217] 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".^[218]

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 made of [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.^[219]

The opinions amongst experts and industry insiders are mixed, with sizable fractions both concerned and unconcerned by risk from eventual superintelligent AI.^[220] Personalities such as [Stephen Hawking](#), [Bill Gates](#), and [Elon Musk](#) have expressed concern about existential risk from AI.^[221] AI pioneers including [Fei-Fei Li](#), [Geoffrey Hinton](#), [Yoshua Bengio](#), [Cynthia Breazeal](#), [Rana el Kaliouby](#), [Demis Hassabis](#), [Joy Buolamwini](#), and [Sam Altman](#) have expressed concerns about the risks of AI. In 2023, many leading AI experts

issued 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".^[222]

Other researchers, however, spoke in favor of a less dystopian view. AI pioneer [Juergen 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."^[223] 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."^{[224][225]} [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."^[226] [Yann LeCun](#) "scoffs at his peers' dystopian scenarios of supercharged misinformation and even, eventually, human extinction."^[227] 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.^[228] However, after 2016, the study of current and future risks and possible solutions became a serious area of research.^[229]

Ethical machines and alignment

Main articles: [Machine ethics](#), [AI safety](#), [Friendly artificial intelligence](#), [Artificial moral agents](#), and [Human Compatible](#)

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.^[230]

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.^[231] The field of machine ethics is also called computational morality,^[231] and was founded at an [AAAI](#) symposium in 2005.^[232]

Other approaches include [Wendell Wallach's](#) "artificial moral agents"^[233] and [Stuart J. Russell's three principles](#) for developing provably beneficial machines.^[234]

Open source

Active organizations in the AI open-source community include [Hugging Face](#),^[235] [Google](#),^[236] [EleutherAI](#) and [Meta](#).^[237] Various AI models, such as [Llama 2](#), [Mistral](#) or [Stable Diffusion](#), have been made open-weight,^{[238][239]} 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.^[240] 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 can't be deleted everywhere if needed. They recommend pre-release audits and cost-benefit analyses.^[241]

Frameworks

Artificial Intelligence projects can have their ethical permissibility tested while designing, developing, and implementing an AI system. An AI framework such as the [Care and Act Framework](#) containing

the SUM values—developed by the [Alan Turing Institute](#) tests projects in four main areas:[242][243]

- 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;[244] however, these principles do not go without their criticisms, especially regards to the people chosen contributes to these frameworks.[245]

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.[246]

The AI Safety Institute in the UK has released a testing toolset called 'Inspect' for AI safety evaluations available under a 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.[247]

Regulation

Main articles: [Regulation of artificial intelligence](#), [Regulation of algorithms](#), and [AI safety](#)



The first global [AI Safety Summit](#) was held in 2023 with a declaration calling for international co-operation.

The regulation of artificial intelligence is the development of public sector policies and laws for promoting and regulating artificial intelligence (AI); it is therefore related to the broader regulation of algorithms.^[248] The regulatory and policy landscape for AI is an emerging issue in jurisdictions globally.^[249] 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.^{[250][251]} Between 2016 and 2020, more than 30 countries adopted dedicated strategies for AI.^[252] 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.^[252] 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.^[252] [Henry Kissinger](#), [Eric Schmidt](#), and [Daniel Huttenlocher](#) published a joint statement in November 2021 calling for a government commission to regulate AI.^[253] In 2023, OpenAI leaders published recommendations for the governance of superintelligence, which they believe may happen in less than 10 years.^[254] In 2023, the United Nations also launched an advisory

body to provide recommendations on AI governance; the body comprises technology company executives, governments officials and academics.^[255]

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".^[250] A 2023 [Reuters/Ipsos](#) poll found that 61% of Americans agree, and 22% disagree, that AI poses risks to humanity.^[256] 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".^{[257][258]}

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.^[259] 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.^{[260][261]}

History

Main article: [History of artificial intelligence](#)

For a chronological guide, see [Timeline of artificial intelligence](#).

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.^{[262][5]} This, along with concurrent discoveries

in [cybernetics](#), [information theory](#) and [neurobiology](#), led researchers to consider the possibility of building an "electronic brain".^[9] They developed several areas of research that would become part of AI,^[264] such as [McCullouch](#) and [Pitts](#) design for "artificial neurons" in 1943,^[265] and Turing's influential 1950 paper '[Computing Machinery and Intelligence](#)', which introduced the [Turing test](#) and showed that "machine intelligence" was plausible.^{[266][5]}

The field of AI research was founded at [a workshop](#) at [Dartmouth College](#) in 1956.^{[7][6]} The attendees became the leaders of AI research in the 1960s.^[5] They and their students produced programs that the press described as "astonishing":^[8] computers were learning [checkers](#) strategies, solving word problems in algebra, proving [logical theorems](#) and speaking English.^{[u][7]} Artificial intelligence laboratories were set up at a number of British and U.S. Universities in the latter 1950s and early 1960s.^[5]

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.^[270] [Herbert Simon](#) predicted, "machines will be capable, within twenty years, of doing any work a man can do".^[271] [Marvin Minsky](#) agreed, writing, "within a generation ... the problem of creating 'artificial intelligence' will substantially be solved".^[272] They had, however, underestimated the difficulty of the problem.^[v] In 1974, both the U.S. and British governments cut off exploratory research in response to the [criticism](#) of [Sir James Lighthill](#)^[274] and ongoing pressure from the U.S. Congress to [fund more productive projects](#).^[275] [Minsky's](#) 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.^[276] The "AI winter", a period when obtaining funding for AI projects was difficult, followed.^[9]

In the early 1980s, AI research was revived by the commercial success of *expert systems*,^[277] 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*.^[8] 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.^[10]

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*,^[278] and began to look into "sub-symbolic" approaches.^[279] *Rodney Brooks* rejected "representation" in general and focussed directly on engineering machines that move and survive.^[w] *Judea Pearl*, *Lofti Zadeh* and others developed methods that handled incomplete and uncertain information by making reasonable guesses rather than precise logic.^{[88][284]} But the most important development was the revival of "*connectionism*", including neural network research, by *Geoffrey Hinton* and others.^[285] In 1990, *Yann LeCun* successfully showed that *convolutional neural networks* can recognize handwritten digits, the first of many successful applications of neural networks.^[286]

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](#)).^[287] By 2000, solutions developed by AI researchers were being widely used, although in the 1990s they were rarely described as "artificial intelligence".^[288] 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.^[14]

[Deep learning](#) began to dominate industry benchmarks in 2012 and was adopted throughout the field.^[11] For many specific tasks, other methods were abandoned.^[x] Deep learning's success was based on both hardware improvements ([faster computers](#),^[290] [graphics processing units](#), [cloud computing](#)^[291]) and access to [large amounts of data](#)^[292] (including curated datasets,^[291] such as [ImageNet](#)). Deep learning's success led to an enormous increase in interest and funding in AI.^[9] The amount of machine learning research (measured by total publications) increased by 50% in the years 2015–2019.^[252]

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.^[229]

In the late teens 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 was taught only the rules of the game and developed 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.^[293] These programs, and others, inspired an aggressive [AI boom](#), where large companies began investing billions in AI research. According to AI Impacts, about \$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".^[294] About 800,000 "AI"-related U.S. job openings existed in 2022.^[295]

Philosophy

Main article: [Philosophy of artificial intelligence](#)

Defining artificial intelligence

Main articles: [Turing test](#), [Intelligent agent](#), [Dartmouth workshop](#), and [Synthetic intelligence](#)

[Alan Turing](#) wrote in 1950 "I propose to consider the question 'can machines think?'"^[296] He advised changing the question from whether a machine "thinks", to "whether or not it is possible for machinery to show intelligent behaviour".^[296] He devised the Turing test, which measures the ability of a machine to simulate human conversation.^[266] 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"^[297]

[Russell](#) and [Norvig](#) agree with Turing that intelligence must be defined in terms of external behavior, not internal structure.^[1] 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.'" ^[298] AI founder John McCarthy agreed, writing that "Artificial intelligence is not, by definition, simulation of human intelligence". ^[299]

McCarthy defines intelligence as "the computational part of the ability to achieve goals in the world". ^[300] Another AI founder, Marvin Minsky similarly describes it as "the ability to solve hard problems". ^[301] 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. ^[1] 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, ^[302] 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.

Evaluating approaches to AI

No established unifying theory or paradigm has guided AI research for most of its history. ^[2] 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 (or "*GOFAI*")^[304] 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."^[305]

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.^[306] 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.^[307] Although his arguments had been ridiculed and ignored when they were first presented, eventually, AI research came to agree with him.^{[aa][19]}

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,^{[309][310]} 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

Main article: [Neats and scruffies](#)

"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,^[311] but eventually was seen as irrelevant. Modern AI has elements of both.

Soft vs. hard computing

Main article: [Soft computing](#)

Finding a provably correct or optimal solution is [intractable](#) for many important problems.^[18] 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

Main articles: [Weak artificial intelligence](#) and [Artificial general intelligence](#)

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.^{[312][313]} 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 experimental sub-field of artificial general intelligence studies this area exclusively.

Machine consciousness, sentience, and mind

Main articles: [Philosophy of artificial intelligence](#) and [Artificial consciousness](#)

The [philosophy of mind](#) does not know 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."^[314] 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

Main articles: [Hard problem of consciousness](#) and [Theory of mind](#)

[David Chalmers](#) identified two problems in understanding the mind, which he named the "hard" and "easy" problems of consciousness.^[315] 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*.^[316]

Computationalism and functionalism

Main articles: [Computational theory of mind](#), [Functionalism \(philosophy of mind\)](#), and [Chinese room](#)

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](#).^[317]

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."^[ab] Searle counters this assertion with his Chinese room argument, which attempts to show that, even if a machine perfectly simulates human behavior, there is still no reason to suppose it also has a mind.^[321]

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.^[322] 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.^{[323][324]} [Sapience](#) (a set of capacities related to high intelligence, such as discernment or [self-awareness](#)) may provide another moral basis for AI rights.^[323] [Robot rights](#) are also sometimes proposed as a practical way to integrate autonomous agents into society.^[325]

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.^[326] 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 to society on their own.^{[327][328]}

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.^{[324][323]}

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.^[313]

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](#)".^[329]

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.^[330]

Transhumanism

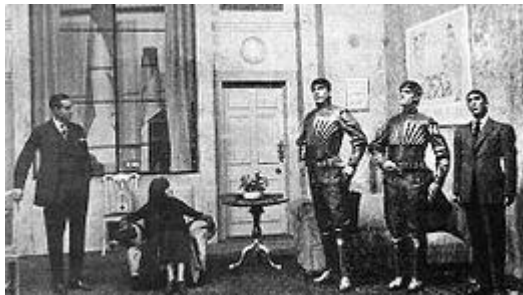
Robot designer [Hans Moravec](#), cyberneticist [Kevin Warwick](#), and inventor [Ray Kurzweil](#) have predicted that humans and machines will merge in the future into [cyborgs](#) that are more capable and

powerful than either. This idea, called transhumanism, has roots in [Aldous Huxley](#) and [Robert Ettinger](#).^[331]

[Edward Fredkin](#) argues that "artificial intelligence is the next stage 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 book of the same name in 1998.^[332]

In fiction

Main article: [Artificial intelligence in fiction](#)



The word "robot" itself was coined by [Karel Čapek](#) in his 1921 play [R.U.R.](#), the title standing for "Rossum's Universal Robots".

Thought-capable artificial beings have appeared as storytelling devices since antiquity,^[333] and have been a persistent theme in [science fiction](#).^[334]

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.^[335]

Isaac Asimov introduced the *Three Laws of Robotics* in many books and stories, most notably the "Multivac" series about a super-intelligent computer of the same name. Asimov's laws are often brought up during lay discussions of machine ethics;^[336] 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.^[337]

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.^[338]