

FOR EVERYONE

OpenAI offers generative AI services to individuals on a free and subscription basis. Free offerings include limited access to the flagship ChatGPT-5 and lesser models, while subscription offerings include extended access to more advanced GPT versions and features, increased message capacity, DALL-E image generation, and limited access to Sora video generation.



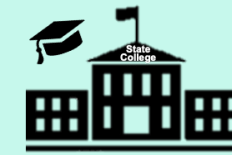
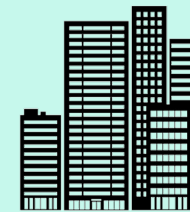
FOR BUSINESS

OpenAI offers generative AI services to organizations with small to moderate sized teams looking to supercharge their work through the integration of generative AI tools in a collaborative projects and workflows with administrative controls. These work-centered features assist teams with generating code, crafting emails, analyzing data, and brainstorming ideas, among other things.



FOR ENTERPRISES

OpenAI offers generative AI services to large organizations seeking to enable their entire workforce with generative AI tools like GPT-5, DALL-E, web browsing, and data analysis, but which require admin controls, domain verification and other security protocols, customization, scalability, analytics, and a dedicated account team for support.



FOR DEVELOPERS

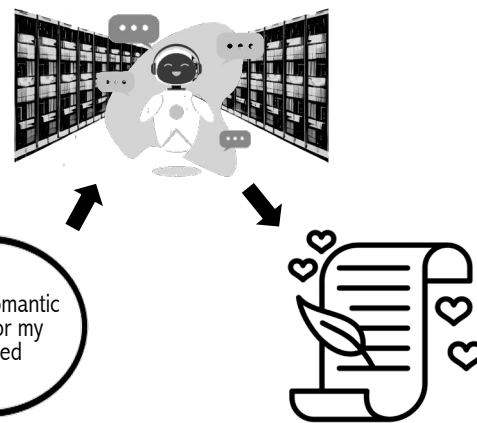
OpenAI's developer plan provides robust tools, APIs, and frameworks that enable developers to build, deploy, and scale AI applications. It includes access to flagship models (e.g., GPT-5) via its API, support for agentic and tool-integrated workflows, enterprise security features, and continuous updates to reduce costs and improve usability. The plan aims to make AI development more accessible, powerful, and commercially viable for individuals and organizations.



MODELS



Text to Text



TTS (text-to-speech)



Deep Research



Embeddings



Whisper



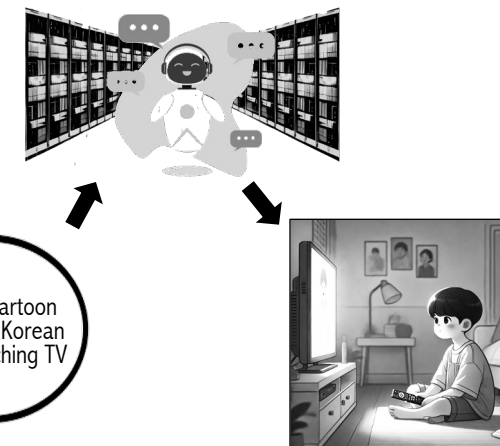
Moderation



Deprecated



Text to Image



Text to Video



AI Coding Agent for Software Development



GLOSSARY

AGI: AGI, or Artificial General Intelligence, refers to a theoretical form of AI that possesses the ability to understand, learn, and apply knowledge across a wide range of tasks at a level comparable to human intelligence. Unlike narrow AI, which is specialized for specific tasks, AGI would be capable of generalizing knowledge and reasoning, allowing it to perform any intellectual task that a human can do.

AI: Refers to "artificial intelligence," which is the simulation of the human brain by machines or computer systems to enable problem solving. Specific applications of AI include expert systems, natural language processing, speech recognition, and machine vision.

AI Chatbot: An AI chatbot refers to a type of AI-powered program capable of generating written content from a user's input prompt. AI chatbots can write anything from a rap song to an essay upon a user's request. The extent of what each chatbot can write about depends on its capabilities, including whether it is connected to a search engine. At the user's command, AI chatbots can write code, compose emails, draft a report, generate art, write Excel formulas, and much more.

AI Writer: The main difference between an AI chatbot and an AI writer is the type of output they generate and their primary function.

In the past, an AI writer was used specifically to generate written content, such as articles, stories, or poetry, based on a given prompt or input. An AI writer outputs text that mimics human-like language and structure. On the other hand, an AI chatbot is designed to conduct real-time conversations with users in text or voice-based interactions. The primary function of an AI chatbot is to answer questions, provide recommendations, or even perform simple tasks, and its output is in the form of text-based conversations.

While the terms AI chatbot and AI writer are now used interchangeably by some, the original distinction was that an AI chatbot was used for conversational purposes. However, with the introduction of more advanced technology such as ChatGPT, the line between the two has become increasingly blurred. Many AI chatbots are now capable of generating text-based responses that mimic human-like language and structure, similar to an AI writer.

Algorithm: An algorithm refers to a mathematical model or set of computational rules that enables machines to learn from data, make predictions, or perform specific tasks. These algorithms can range from simple decision trees to complex neural networks, and they are foundational in areas like machine learning, natural language processing, and computer vision.

API: An API, or Application Programming Interface, is a set of rules and protocols that allows different software applications to communicate and interact with each other. It defines the methods and data formats that programs use to request and exchange information. APIs enable developers to integrate third-party services and functionalities into their applications without needing to understand the underlying code.



Selected examples of how APIs are used to facilitate communication between different software applications are:

- **Data Retrieval:** APIs allow applications to fetch data from external sources, such as weather information, financial data, or social media feeds.
- **Integration of Services:** APIs enable different services to work together, such as integrating payment gateways (e.g., PayPal, Stripe) into e-commerce platforms.
- **Accessing Functionality:** Applications can use APIs to access specific functionalities of other software, like sending emails through an email service API (e.g., SendGrid).
- **Mobile and Web Applications:** APIs power mobile and web applications by allowing them to communicate with back-end servers for data processing and storage.
- **Automation:** APIs can automate workflows by connecting different applications. For example, an API can trigger an action in one application when a specific event occurs in another.
- **Third-party Services:** Many platforms offer APIs to allow developers to extend their services, such as integrating with customer relationship management (CRM) systems, analytics tools, or social media platforms.

Overall, APIs play a crucial role in enhancing interoperability, enabling developers to build more flexible, powerful, and interconnected applications.

Big Data: Big data refers to extremely large and complex datasets that cannot be easily managed, processed, or analyzed using traditional data processing tools. The processing of big data involves the use of advanced technologies and methodologies to capture, store, analyze, and derive insights from vast amounts of structured and unstructured data, often in real-time, to support decision-making and uncover trends.



CAPTCHA: CAPTCHA is the acronym for Completely Automated Public Turing test to tell Computers and Humans apart. They typically consist of a string of twisted letters or other visual symbols that humans can identify correctly but algorithms struggle with. CAPTCHA puzzles are used by websites to determine whether users are humans and to block bot attacks. In a test administered in 2023, OpenAI's GPT-4 model was able to overcome CAPTCHA visual puzzles.

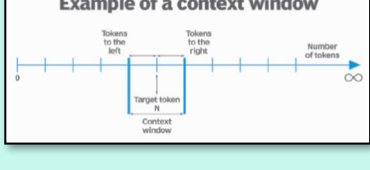


Chatbot: A chatbot is a software application designed to simulate conversation with human users, typically through text or voice interactions. Chatbots can be integrated into websites, messaging platforms, or mobile apps, and they utilize natural language processing (NLP) and artificial intelligence (AI) to understand and respond to user inquiries. They can provide information, assist with customer service, automate tasks, and enhance user engagement by offering instant responses and 24/7 availability. Chatbots vary in complexity, from simple rule-based systems to advanced AI-driven conversational agents, such as ChatGPT.

Cloud: Refers to servers that are accessed over the Internet, and the software and databases that run on those servers. Cloud servers are located in data centers all over the world. By using cloud computing, users and companies do not have to manage physical servers themselves or run software applications on their own machines.

Computer Vision: A field of artificial intelligence that enables computers to interpret and understand visual information from the world, such as images and videos. It involves developing algorithms and models that can analyze and make decisions based on visual input, facilitating tasks like object recognition, image classification, and scene understanding. By simulating human vision, computer vision applications are used in areas such as autonomous vehicles, facial recognition, and medical imaging.

Context Window: A context window refers to the amount of text or data that a language model, such as an LLM, can consider at once when generating or analyzing text. It defines the range of input that the model can use to understand the current situation, make predictions, or generate coherent responses. The size of the context window is typically measured in tokens, which can be words, parts of words, or symbols, and a larger context window allows the model to maintain a better understanding of long passages, complex instructions, or intricate conversations.



Copyright: Copyright is a legal protection granted to the creators of original works, including literary, musical, and artistic creations, giving them exclusive rights to use, distribute, and profit from their work. This protection typically lasts for a set period, after which the work enters the public domain, allowing others to use it freely. Copyright laws aim to balance the creator's rights with the public's interest in accessing and building upon creative works.

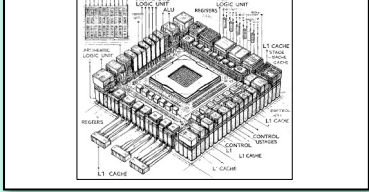
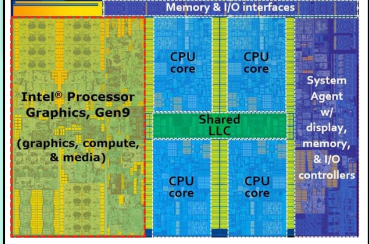
CPU: Refers to "central processing unit" which is a chip that functions as the brains of computing devices. Also called the central processor, main processor or processor, the CPU is made up of a set of electronic circuits that execute instructions comprising a computer program to perform basic arithmetic, logic, controlling, and I/O (input/output) operations specified in a computer program.

Creative AI: Creative AI refers to artificial intelligence systems designed to produce original and innovative content, such as art, music, writing, or design, by mimicking human creativity. These systems use machine learning algorithms, often trained on vast datasets, to generate new ideas, solve complex problems, and assist in creative processes. Creative AI has the potential to augment human creativity, offering new tools and perspectives in various artistic and creative fields.

Cybersecurity: Cybersecurity is the practice of protecting computers, networks, data, and systems from digital attacks, unauthorized access, damage, or theft. It involves implementing measures like firewalls, encryption, and intrusion detection systems, as well as educating users on safe practices, to safeguard sensitive information and ensure the integrity, confidentiality, and availability of digital assets.



Cores: Cores are individual processing units within a CPU or GPU that execute instructions. Each core can independently perform tasks, allowing for parallel processing and improved performance. More cores enable a chip to handle multiple tasks simultaneously, enhancing multitasking capabilities and overall efficiency. Cores are essential for running applications more quickly and efficiently, particularly in modern computing environments that require high processing power for tasks like gaming, data analysis, and machine learning (ML).

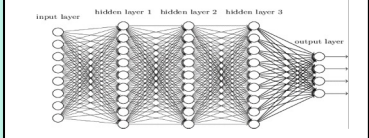


Data Engineering: Data engineering is the practice of designing, building, and maintaining the infrastructure and systems that enable the collection, storage, and processing of large volumes of data. It involves creating data pipelines, managing databases, and ensuring data is clean, accessible, and ready for analysis, supporting data scientists and analysts in their work to extract insights and make data-driven decisions.

Data Science: Data science is an interdisciplinary field that combines statistical analysis, machine learning, and domain expertise to extract insights and knowledge from structured and unstructured data. It involves collecting, processing, analyzing, and interpreting large datasets to inform decision-making and solve complex problems across various industries.

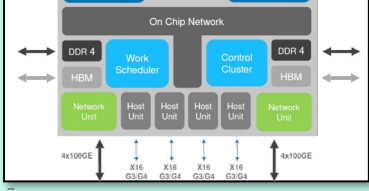
Deepfake: Deepfake technology uses artificial intelligence (AI) to create or manipulate audiovisual content, often replacing faces or voices in videos to make them appear authentic but are actually synthetic.

Deep Learning: A subset of machine learning, deep learning is a computer science approach where neural networks (1) are trained to recognize patterns from massive amounts of data (in the form of images, sounds and text), often better than humans, and in turn (2) provide predictions in a production process. Deep learning is often used to make non-linear, complex correlations, and needs to run on a specialized graphics processing unit (GPU).



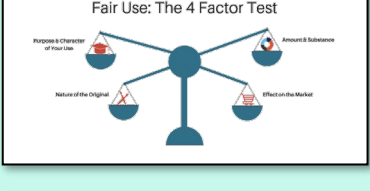
DevOps: DevOps is a set of practices that combines software development (Dev) and IT operations (Ops) to shorten the software development lifecycle and deliver high-quality software more reliably. It emphasizes collaboration, automation, continuous integration, and continuous delivery (CI/CD), enabling the public's interest in accessing and building upon creative works.

GPU: Data Processing Unit, which is a programmable computer processor designed to efficiently handle data-centric workloads, such as data transfer, reduction, security, compression, analytics and encryption, at scale in data centers. A GPU tightly integrates a general-purpose CPU with network interface hardware.



Edge Server: Edge servers are servers that run the processing at an edge location. In a centralized network, client devices are connected to one server whose job is to process the information requested by the users and hand it back to them. While a centralized network is good enough to interact with simpler websites, companies have found that complex projects involving a great bank of customers work better with an edge server. Unlike a centralized network, because an edge server sits at the edge of a network, it improves latency, reduces loading times and removes the load from the origin server. Instead of sending unprocessed data to the data center, edge servers process it themselves and send it back to the client machines.

Fair Use Doctrine: The Fair Use Doctrine is a legal principle that allows limited use of copyrighted material without the owner's permission under certain circumstances. It typically applies to uses for purposes such as criticism, commentary, news reporting, teaching, scholarship, or research. Fair use is determined by considering factors like the purpose and character of the use, the nature of the copyrighted work, the amount used, and the effect on the market value of the original work. This doctrine seeks to balance the rights of creators with the need for freedom of expression and access to information.



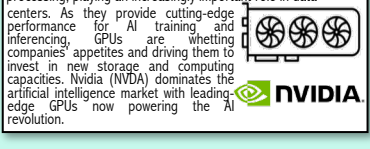
FPGA: Refers to Field-Programmable Gate Array, an integrated circuit designed to be configured by a customer or a designer after manufacturing—hence "field-programmable." FPGAs contain an array of programmable logic blocks and a hierarchy of reconfigurable interconnects, allowing the blocks to be wired together to perform complex combinational and sequential logic functions. This flexibility makes FPGAs highly valuable in applications that require parallel processing, high-speed data handling, and real-time processing, such as digital signal processing, aerospace, defense systems, and increasingly in data centers and artificial intelligence (AI) applications. Unlike fixed-function application-specific integrated circuits (ASICs), FPGAs can be reprogrammed to adapt to changing requirements, offering a cost-effective and versatile solution for many industries.



GANs: Generative Adversarial Networks (GANs) are a type of artificial intelligence framework consisting of two neural networks—the generator and the discriminator—competing against each other to create realistic data. The generator creates fake data, such as images or text, while the discriminator evaluates their authenticity, distinguishing between real and generated data. Over time, this adversarial process improves the generator's ability to produce increasingly convincing outputs, making GANs powerful tools for generating high-quality synthetic content in fields like art, design, and deepfake technology.

GPT: GPT stands for "Generative Pre-trained Transformer," a type of language model developed by OpenAI that uses deep learning techniques to generate human-like text based on the input it receives. The "Generative" aspect indicates its ability to produce text, "Pre-trained" means it is initially trained on a large corpus of text data before being fine-tuned for specific tasks, and "Transformer" refers to the underlying architecture that enables it to process and understand language effectively.

GPU: Refers to "graphics processing unit," a type of microprocessor capable of rendering graphics display on an electronic device, such as a gaming console or vehicle infotainment system. GPUs also have application for video editing and content creation. As a result, GPUs are also known as "video cards" or "graphics cards." However, GPUs are the backbone of machine and deep learning and artificial intelligence (AI) centers. As they provide cutting-edge performance for AI training and inferencing, GPUs are whetting companies' appetites and driving them to invest in new storage and computing capacities. Nvidia (NVDA) dominates the artificial intelligence market with leading-edge GPUs now powering the AI revolution.



Hallucination: An AI hallucination is when an AI model produces inaccurate, misleading, or biased information. This can happen when AI models generate content that is not based on real-world data, but rather on the model's own imagination.

Human-Computer Interaction (HCI): Human-Computer Interaction (HCI) is the study and design of how people interact with computers and technology, focusing on creating user-friendly and efficient interfaces. It involves understanding user behavior, designing intuitive interfaces, and evaluating usability to enhance the overall user experience with digital products and systems.

Inferencing: AI inferencing refers to the process of applying a trained machine learning model to new, unseen data to make predictions or decisions. It involves using the learned patterns and relationships encoded in the model to process input data and produce an output. In practical terms, inferencing is where an AI system takes input, applies what it has learned during training, and generates a response or output, such as classifying an image, translating text, or recommending a product.

Input Processing: Input processing refers to the method by which raw data or user inputs are received, interpreted, and prepared for further analysis or action by a system. This involves various steps, such as data cleaning, normalization, feature extraction, and transformation, to convert the inputs into a format that the system can effectively utilize. Effective input processing is crucial for ensuring accurate, efficient, and meaningful outputs in tasks like machine learning, natural language processing, and other automated systems.

Intellectual Property (IP): A category of property that includes intangible creations of the human intellect. There are many types of intellectual property, and some countries recognize more than others. The best-known types are patents, copyrights, trademarks, and trade secrets.

Licensing: Licensing of intellectual property involves a legal agreement between two parties - the IP rights owner (licensor) and another party (licensee) who is authorized to use those rights for contractually agreed purposes.

LLM: A Large Language Model (LLM) is an advanced type of artificial intelligence model designed to understand and generate human-like text based on vast amounts of training data. These models use deep learning techniques, particularly transformers, to analyze and predict text, making them capable of performing a wide range of language-related tasks such as translation, summarization, and conversation. This is in contrast to models designed to generate speech, sounds, or images.

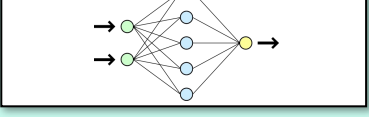
Machine Learning: Machine learning, or ML, is a branch of artificial intelligence that enables systems to learn from data, identify patterns, and make decisions with minimal human intervention. It involves training algorithms on datasets to improve their performance over time in tasks such as prediction, classification, and optimization. In other words, ML is an application of AI where machines/computer programs are given access to data and then use algorithms to find patterns in data (i.e., learn from the data). ML applications are often used to make simple, linear correlations, and can run on a CPU instead of a GPU that excels at running many smaller tasks at once.

With ML, a machine can automatically learn new information and improve itself from experiences without having to be programmed in a certain way since the machine will be able to teach itself new information. However, the machine does not sense things (yet) and must be (1) told when they get something wrong or (2) given new data to figure out what the correct answer is. If new data or feedback is not given to the machine, it will not organically "learn" from its mistakes and adjust accordingly. As a result, ML does not provide "feedback"; it needs to be told explicitly how to fix a problem. AI, on the other hand, operates off a feedback loop. By contrast to ML where new data must be input to make corrections, AI is able to seek out new sources of data on its own and rebuild itself based on those feedbacks.

Model Application: A model application refers to the practical use of a trained machine learning or AI model to solve real-world problems or perform specific tasks. This involves integrating the model into software systems or platforms where it can analyze data, make predictions, automate decisions, or generate content based on its training. Model applications span a wide range of industries, including healthcare, finance, marketing, and technology, where they are used for tasks like diagnosing diseases, predicting market trends, personalizing customer experiences, and optimizing business operations.

Natural Language Processing (NLP): A field of artificial intelligence focused on the interaction between computers and human language. It involves developing algorithms and models that enable machines to understand, interpret, and generate human language in a way that is both meaningful and useful. NLP applications include tasks such as language translation, sentiment analysis, and text summarization.

Neural Network: The Neural Network conducts correlation analysis on all tokens, characterized by the number of nodes and layers. Nodes in one layer connect to nodes in another layer, with connections represented by lines. Weights and bias terms are assigned to these connections. The bias term is calculated as the product of nodes in the layer multiplied by the number of connections, plus the number of nodes in the next layer. The total bias terms across layers define the parameters.



Nodes: Nodes are individual devices or points that participate in a network or system. In data structures, nodes represent fundamental elements that hold data and may link to other nodes, forming structures like trees or graphs. Nodes play a critical role in organizing and managing information, facilitating communication, and executing processes within systems.

NPU: Refers to "neural processing unit," also known as a neural processor, tensor processing unit, or AI accelerator. The NPU is a specialized processor that implements all necessary control and arithmetic logic necessary to execute machine learning (ML) algorithms. Typically by operating on predictive models such as artificial neural networks or random forests. In short, NPUs, such as the M4 chip used in latest versions of Apple's laptops and tablets, specialize in the acceleration of ML algorithms (10,000 times less time than a GPU).

Output Generation: Output generation refers to the process by which a system produces a response or result based on the processed inputs and underlying algorithms. In AI and computing, this could involve generating text, images, predictions, or actions that align with the intended goal of the application. Output generation is the final step in a data processing pipeline, where the system's internal decisions or computations are translated into a form that is understandable or usable by humans or other systems.

Patent: A type of intellectual property that gives its owner the legal right to exclude others from making, using, or selling an invention for a limited period of time in exchange for publishing an enabling disclosure of the invention. The patent period in the U.S. is 20 years from the date of filing in the case of utility patents, and 15 years from the date of grant in the case of design patents.

Parameters: Parameters are variables or values passed into functions, methods, or procedures to influence their behavior or output. They define the specific details or data that a function or algorithm needs to operate correctly. Parameters allow for customization and flexibility, enabling the same function to be used in different contexts with varying inputs.

Real-time and Batch Processing: Real-time and batch processing are two approaches to handling and processing data within computing systems. Real-time processing involves the immediate processing of data as it is received, enabling systems to provide instant responses or updates. This is essential in applications like online transactions, live monitoring systems, and interactive applications where time-sensitive decisions or actions are required. Batch processing refers to the processing of large volumes of data in groups or batches at scheduled intervals. It is typically used for tasks that do not require immediate results, such as payroll processing, end-of-day reports, or data backups, allowing for efficient handling of large datasets over time.

Robotics and Autonomous Systems: Robotics and autonomous systems involve the design and development of machines that can perform tasks independently, with minimal or no human intervention. Robotics focuses on creating physical robots capable of interacting with the environment, often equipped with sensors, actuators, and AI to carry out complex tasks. Autonomous systems extend this concept by enabling these robots or software agents to make decisions, adapt to new situations, and operate in dynamic environments, from manufacturing and logistics to exploration and healthcare.



Semiconductors: Also known as chips or processors, semiconductors are made by imprinting a network of electronic circuits/components onto a wafer that partially conducts electricity, which can then perform various functions—e.g., (a) processing, amplifying and selectively filtering electronic signals (b) controlling, (c) system functions, and (c) storing and transmitting data.

Simulation and Gaming: A field of artificial intelligence that utilizes generative AI for creating virtual environments, game characters, procedural content generation, and interactive simulations.

Software Development: Software development is the process of designing, coding, testing, and maintaining applications, frameworks, or other software components. It involves various stages, including requirement analysis, system design, programming, testing, and deployment, to create functional and reliable software that meets user needs and business goals.

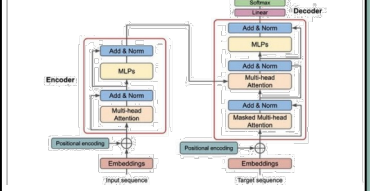


Speech Synthesis and Recognition: Speech Synthesis is the technology that converts text into spoken language, enabling machines to generate human-like voice outputs. Speech Recognition, on the other hand, involves processing and interpreting spoken language to convert it into text or actionable commands. These technologies enable applications such as virtual assistants, voice-controlled devices, and automated transcription services.



Tokens: Tokens refer to individual units of text, such as words or characters, that a model processes during training or inference. These tokens serve as the basic building blocks for understanding and generating human language.

Transformer Architecture: The LLM undergoes pre-training with a substantial amount of information or documents within a context window. The Transformer Architecture serves as the foundational structure for language processing, encompassing both decoding and encoding components. The decoder interprets inputs, while the encoder generates outputs.



LANGUAGE MODEL TRAINING PROCESS

Step 1: Establish a Goal

The first step in training an AI system is to establish an "objective function." Most LLMs will have the same objective function or goal: given a sequence of text, guess which words come next.

Step 2: Collect/Tokenize Data

To train an LLM effectively, a massive amount of data will be required. Once data are collected, they will be broken down or subdivided into units called tokens. These subdivisions will help the model analyze it more easily.

Step 3: Build the Neural Network

After the data have been collected and tokenized, the next step is to create the model's neural network. The neural network (consisting of layer upon layer of logic and decision trees) serves as the foundational structure for language processing.

Step 4: Train the Neural Network

Next, the tokenized data are fed into the neural network in a step called "training the model." The model trains by identifying patterns and relationships in the data.

Step 5: Fine-Tune the Model

Once the model is trained, it will be fine-tuned with human feedback (or reinforcement learning). For example, humans could rate a model's initial responses and feed the ratings back into the model until it improves to an acceptable level.

Step 6: Launch and Monitor

After the model has been fine-tuned, it is ready to be launched. However, because AI systems can hallucinate or be erratic and unpredictable, the developer will monitor the system and follow up with additional rounds of changes as needed before the model is fully functional.

