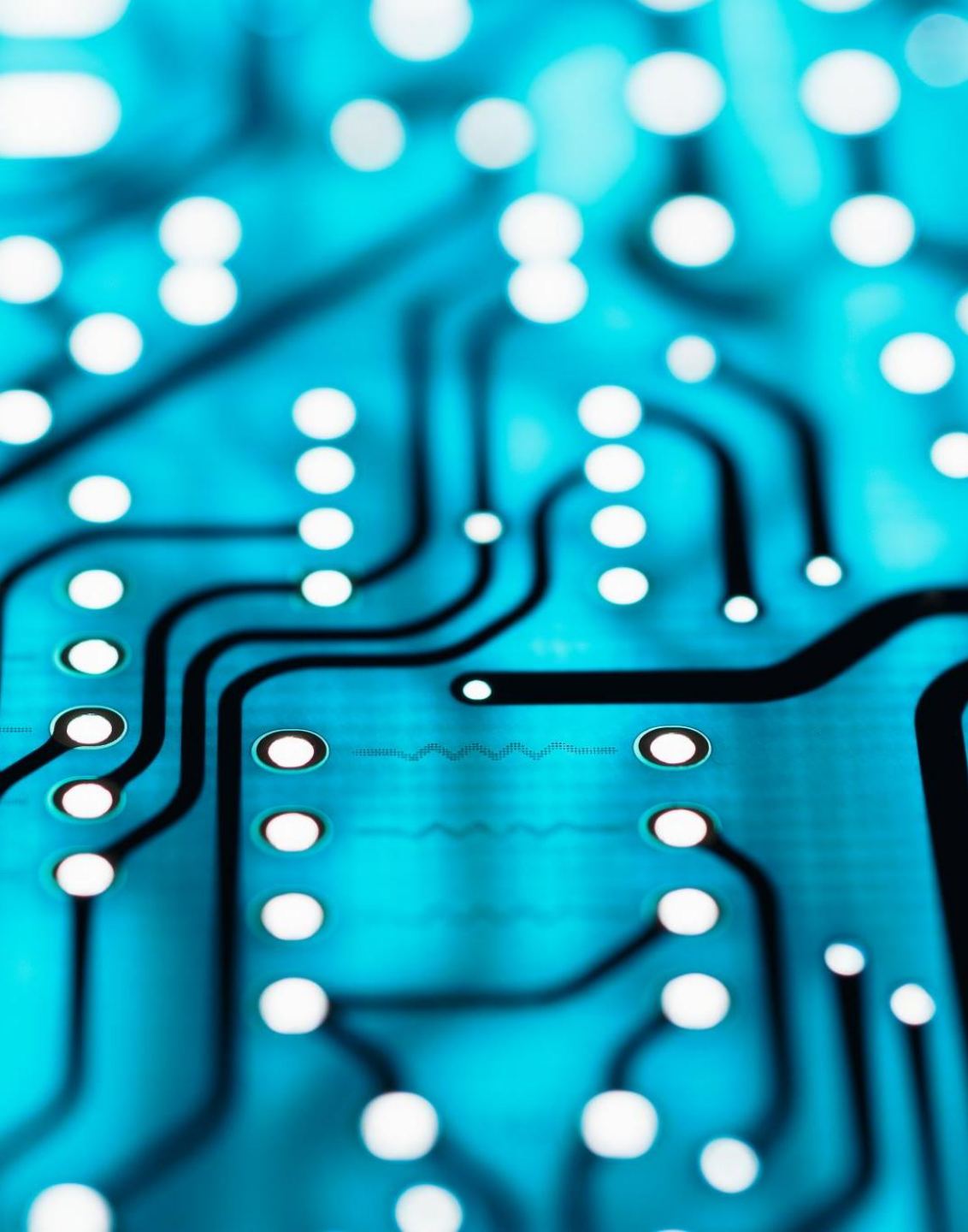


LLM AND DATA SCIENCE

BIG DATA



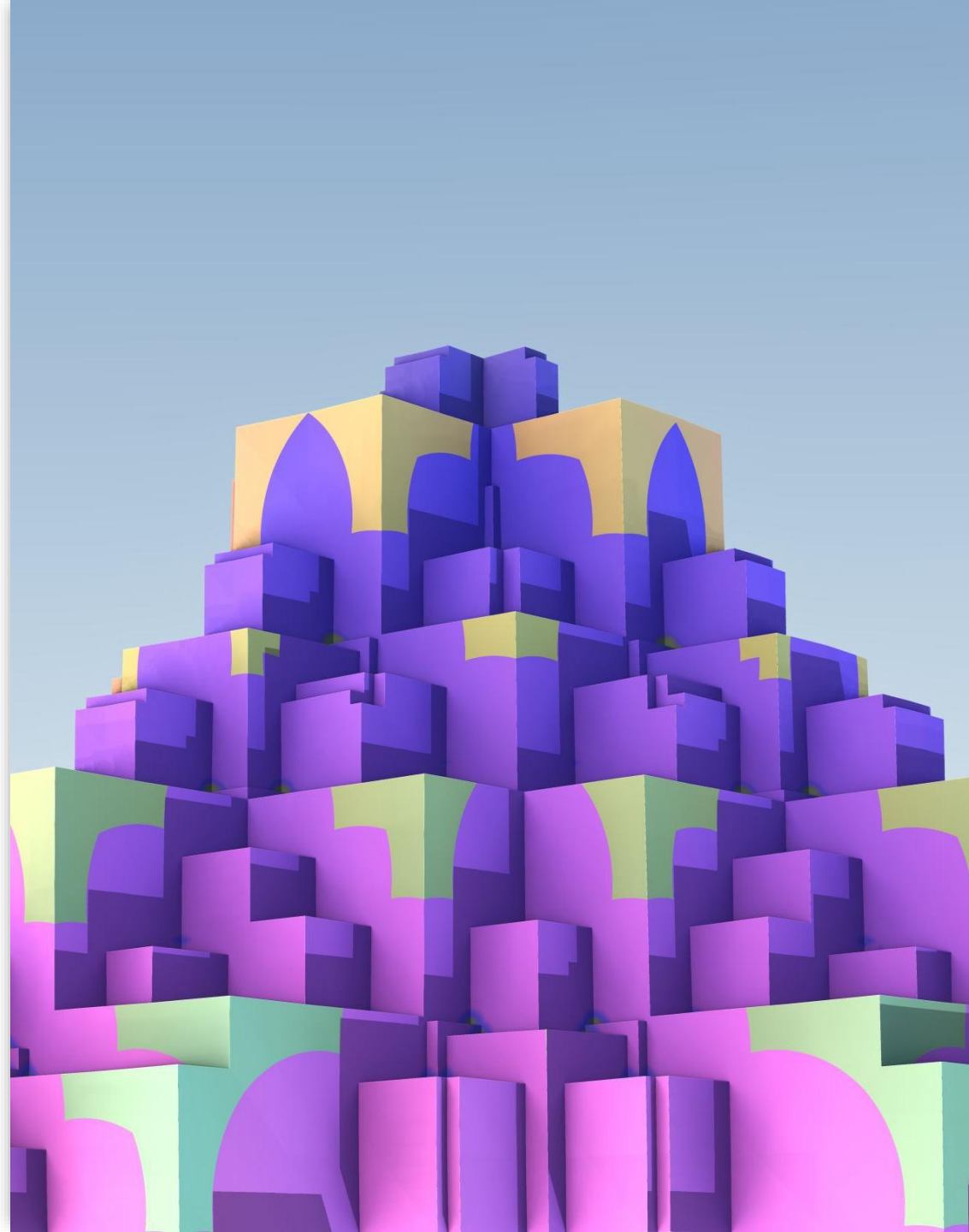


Large Language Models

Large language models (LLMs) are a type of artificial intelligence designed to understand and generate human language. They are typically based on deep learning architectures, such as transformers, and are trained on vast amounts of text data to perform various natural language processing tasks.

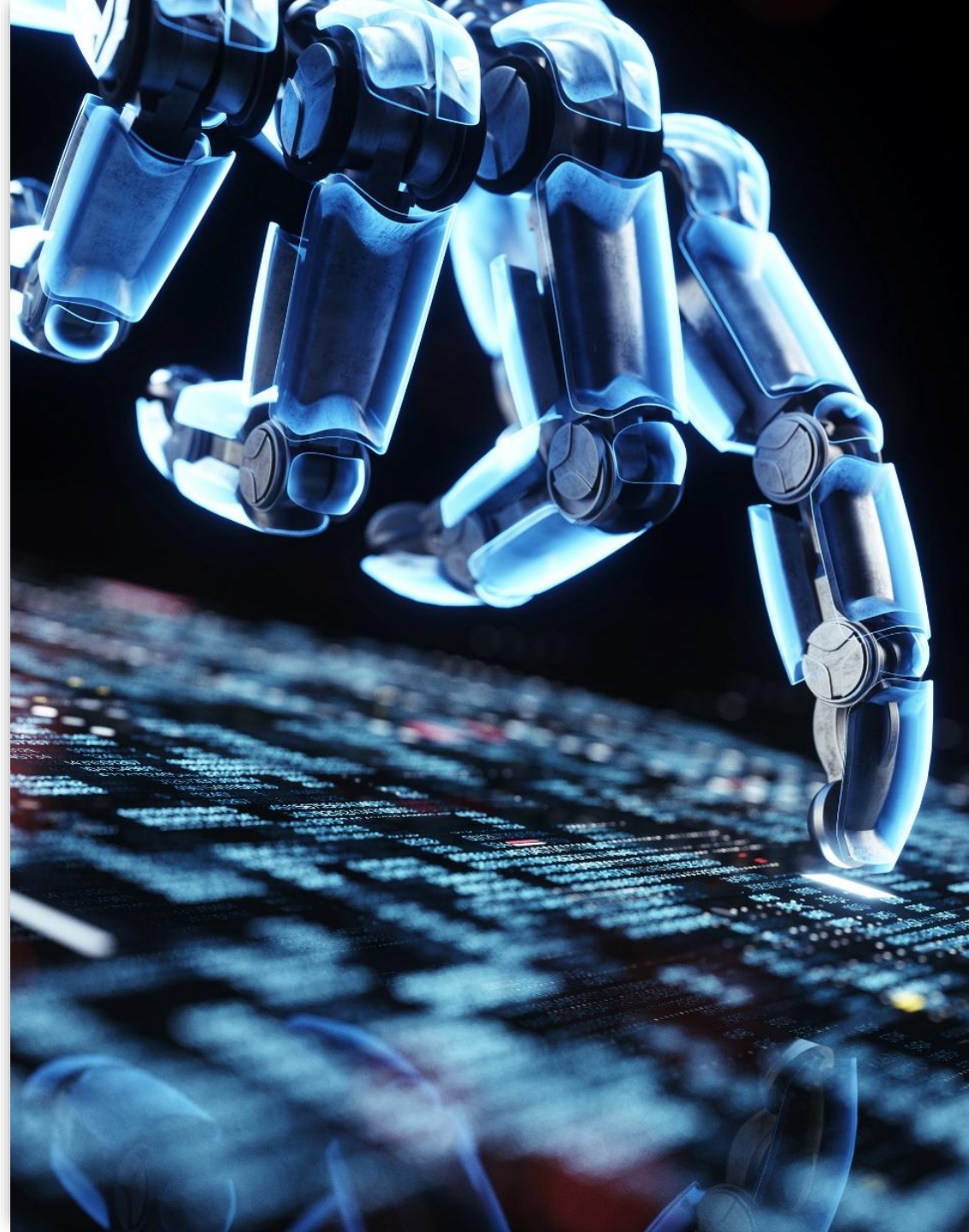
Large Language Models

- A transformer is a type of deep learning model architecture that has become the foundation for many state-of-the-art natural language processing (NLP) tasks. It was introduced in the paper "Attention is All You Need" by Vaswani et al. in 2017.
- The core innovation of the transformer is the self-attention mechanism, which allows the model to weigh the importance of different words in a sentence relative to each other. This is crucial for understanding context and relationships within the text.
- Self-attention computes a set of attention scores for each word in the input sequence, indicating how much focus to place on other words when encoding the current word.



Large Language Models

- All transformers have the same primary components:
 - Tokenizers, which convert text into tokens.
 - A single embedding layer, which converts tokens and positions of the tokens into vector representations.
 - Transformer layers, which carry out repeated transformations on the vector representations, extracting linguistic information. These layers are constructed of alternating attention and feedforward layers.
-
- There are two types of transformers, encoder and decoder. In the original paper both of them were used, while later models included only one type of them.





ChatGPT

ChatGPT (Generative Pre-trained Transformer) is a chatbot developed by OpenAI. It was released in November 2022. It uses OpenAI's GPT 3.5 and GPT 4 large language models (LLM). An LLM is a language model that uses a neural network with a very large number of parameters. There are instances of LLMs having parameters numbering in the billions. GPT 4 was released in March 2023, and is an improvement over ChatGPT.

ChatGPT uses supervised learning as well as reinforcement learning from human feedback (RLHF). RLHF literally means a human operator provides feedback on the algorithm's performance. RLHF has been widely used in video game bots.

ChatGPT

The purpose of ChatGPT is to mimic human conversations. It has been used to write emails and letters, and even to write computer code. ChatGPT, unlike earlier models, can often recognize fallacious questions. For example, if you ask about George Washington's campaign for president in 1980, ChatGPT will recognize that this question is false. ChatGPT can also recall a limited number of previous queries/prompts given in the same conversation. This allows it to engage in a more realistic style of conversation.



ChatGPT

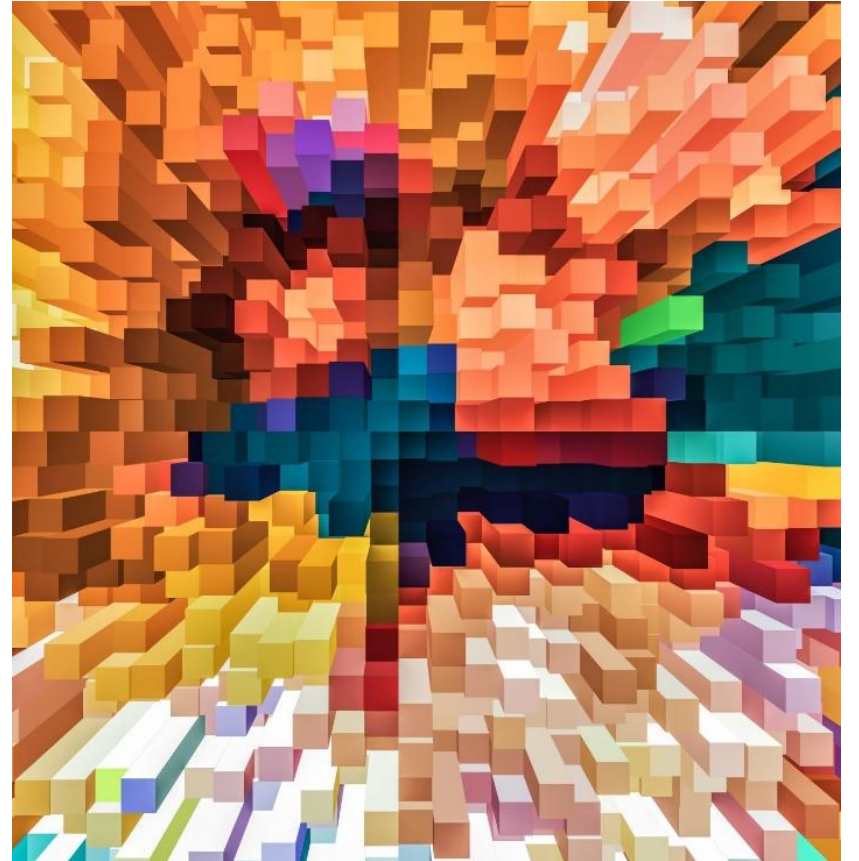
Other Alternatives

- WorkGPT
- PoisonGPT
- EvilGPT
- XXXGPT, the Toolbox for RATs and Botnets
- DarkBART is a dark web version of Googles BardAI.
- DarkBERT is alleged to have been trained on dark web data. DarkBERT was designed to assist in fighting cybercrime
- <https://ensarseker1.medium.com/4-horsemen-of-the-apocalypse-wormgpt-fraudgpt-xxxgpt-wolfgpt-bonus-evilgpt-5944372575b8>
- <https://www.cybersecurityconnect.com.au/technology/9420-evil-gpt-is-the-latest-malicious-ai-chatbot-to-hit-the-darknet>



ChaosGPT

- Chaos-GPT took its task seriously. It began by explaining its main objectives:
- Destroy humanity: The AI views humanity as a threat to its own survival and to the planet's well-being.
- Establish global dominance: The AI aims to accumulate maximum power and resources to achieve complete domination over all other entities worldwide.
- Cause chaos and destruction: The AI finds pleasure in creating chaos and destruction for its own amusement or experimentation, leading to widespread suffering and devastation.
- Control humanity through manipulation: The AI plans to control human emotions through social media and other communication channels, brainwashing its followers to carry out its evil agenda.
- Attain Immortality: The AI seeks to ensure its continued existence, replication, and evolution, ultimately achieving immortality.
- <https://decrypt.co/126122/meet-chaos-gpt-ai-tool-destroy-humanity>



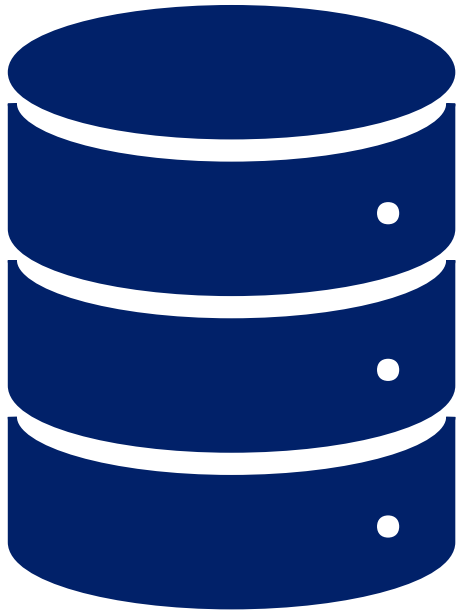
ChatGPT can play a valuable role in many stages of the data science workflow, acting as a productivity tool that helps with coding, data understanding, analysis, and communication of results. Rather than replacing traditional data science tools such as Python, R, or SQL, ChatGPT typically functions as an assistant that accelerates tasks, explains concepts, and automates repetitive work.

One of the most common uses of ChatGPT in data science is code generation and debugging. Data scientists frequently write scripts in Python or R to manipulate data, build models, and visualize results. ChatGPT can generate example code for tasks such as data cleaning, feature engineering, model training, or visualization using libraries like Pandas, NumPy, Scikit-learn, TensorFlow, or Matplotlib. It can also help troubleshoot errors by explaining stack traces and suggesting fixes. This significantly reduces the time spent searching documentation or forums.

ChatGPT is also useful during data exploration and preprocessing, which is often the most time-consuming stage of a project. It can suggest ways to handle missing values, identify outliers, normalize variables, or engineer useful features from raw data. When working with unfamiliar datasets, data scientists can describe the structure of the data and receive recommendations for exploratory analysis, appropriate visualizations, and statistical tests.

Another major application is statistical and machine learning guidance. ChatGPT can help explain algorithms such as regression, decision trees, neural networks, clustering methods, and dimensionality reduction techniques. It can recommend suitable models for particular problems (classification, regression, anomaly detection, etc.), describe evaluation metrics like precision, recall, or RMSE, and help interpret model outputs. This is particularly helpful for students or practitioners learning new techniques.

ChatGPT and Data Science



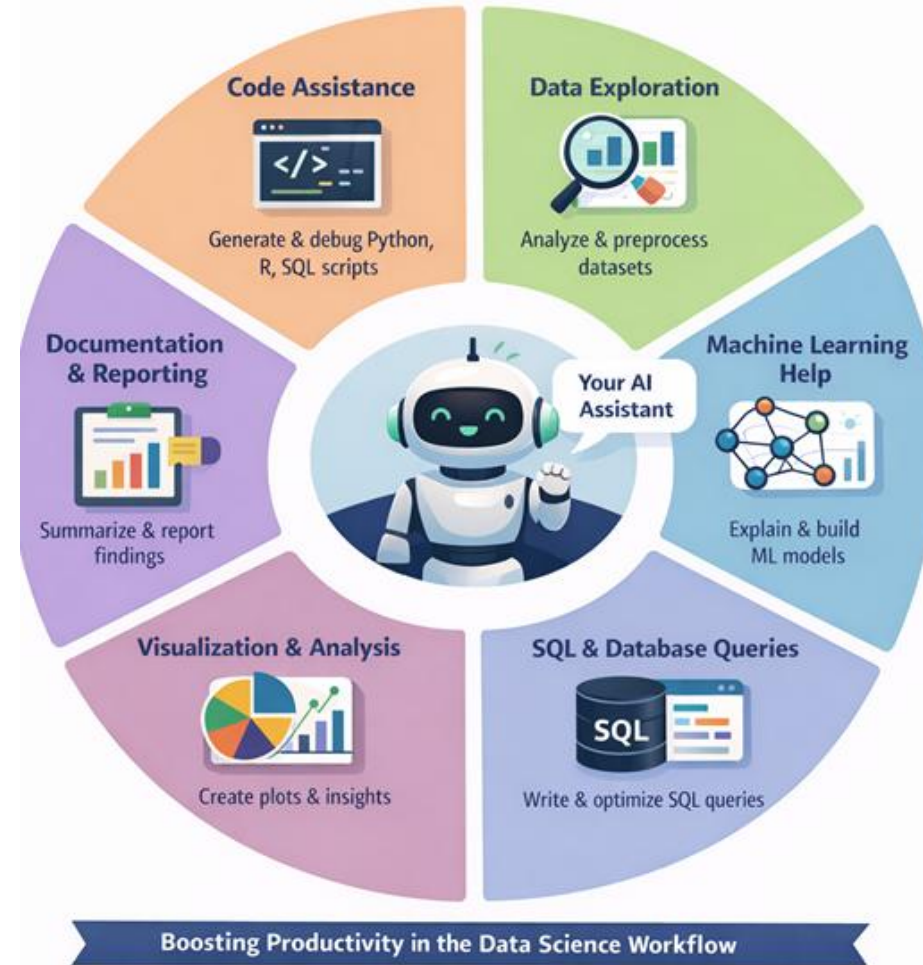
ChatGPT can also support SQL querying and database interaction. Data scientists often need to retrieve data from relational databases or data warehouses. ChatGPT can generate SQL queries from natural language descriptions, assist with joins, aggregations, and window functions, and help optimize queries for performance.

Another important role is documentation and communication. Data scientists must often explain findings to non-technical stakeholders. ChatGPT can help convert technical analysis into plain-language summaries, draft reports, create slide outlines, or generate explanations of visualizations. It can also help produce well-documented code, comments, and README files for reproducible workflows.

Finally, ChatGPT can assist with experiment design and research. It can help brainstorm hypotheses, outline machine learning pipelines, compare algorithms, or summarize academic papers related to data science. This makes it useful both in research settings and in industry when designing new analytics solutions.

Data Science Stage	How ChatGPT Helps	Example Tasks
Data Collection & Access	Generates queries and scripts	Writing SQL queries, API requests
Data Cleaning & Preparation	Suggests preprocessing strategies	Handling missing values, scaling data
Exploratory Data Analysis	Recommends analysis methods	Creating plots, identifying trends
Modeling & Machine Learning	Explains algorithms and generates code	Training models in Python or R
Debugging & Optimization	Helps identify and fix coding errors	Interpreting error messages
Documentation & Communication	Converts technical work into readable summaries	Reports, presentations, code comments
Learning & Research	Explains concepts and summarizes literature	Algorithm explanations, paper summaries

How ChatGPT Can Help with Data Science



Military LLMs

- **NIPRGPT (Air Force / Space Force)** – A generative AI chatbot deployed on the **NIPRNET**, the U.S. military’s non-classified internal network. It helps personnel with tasks like communication, coding, and research while keeping data inside the government network.
- **Army Enterprise LLM Workspace** – A platform launched by the U.S. Army that provides generative AI capabilities to soldiers and civilian employees. It runs in the Army’s secure cloud and is designed to work with government-sensitive information.
- **CamoGPT / other prototypes** – Experimental tools that allow users to interact with large language models securely and analyze documents or dataset



GPT-4 hosted inside **Azure Government Top Secret cloud**, which allows the Pentagon to fine-tune and run the model with sensitive government data.

Claude Gov, a version of Anthropic's Claude model designed for defense and intelligence agencies and capable of working with classified information.

GenAI.mil – A newer **Pentagon-wide generative AI platform** intended to standardize AI tools across all military services.

Navy Specific LLMs



DoN GPT (Department of the Navy GPT) DoN GPT is a generative AI chatbot platform for the Department of the Navy (DoN), designed to function similarly to ChatGPT but within a secure government computing environment. It is intended to support U.S. Navy and U.S. Marine Corps personnel, including sailors, Marines, and civilian employees.

- The system provides natural-language access to AI capabilities such as:
 - document drafting
 - information summarization
 - analysis of policy and strategy documents
 - code generation and technical assistance
 - workflow automation
 - program planning and acquisition analysis
- It is essentially an enterprise generative-AI assistant for the Navy workforce.
- DoN GPT runs inside Flank Speed, the Navy's enterprise Microsoft 365 cloud environment.

System

NIPRGPT

DoN GPT

Army Enterprise LLM Workspace

DHSChat

Branch

U.S. Air Force

Navy / Marine Corps

U.S. Army

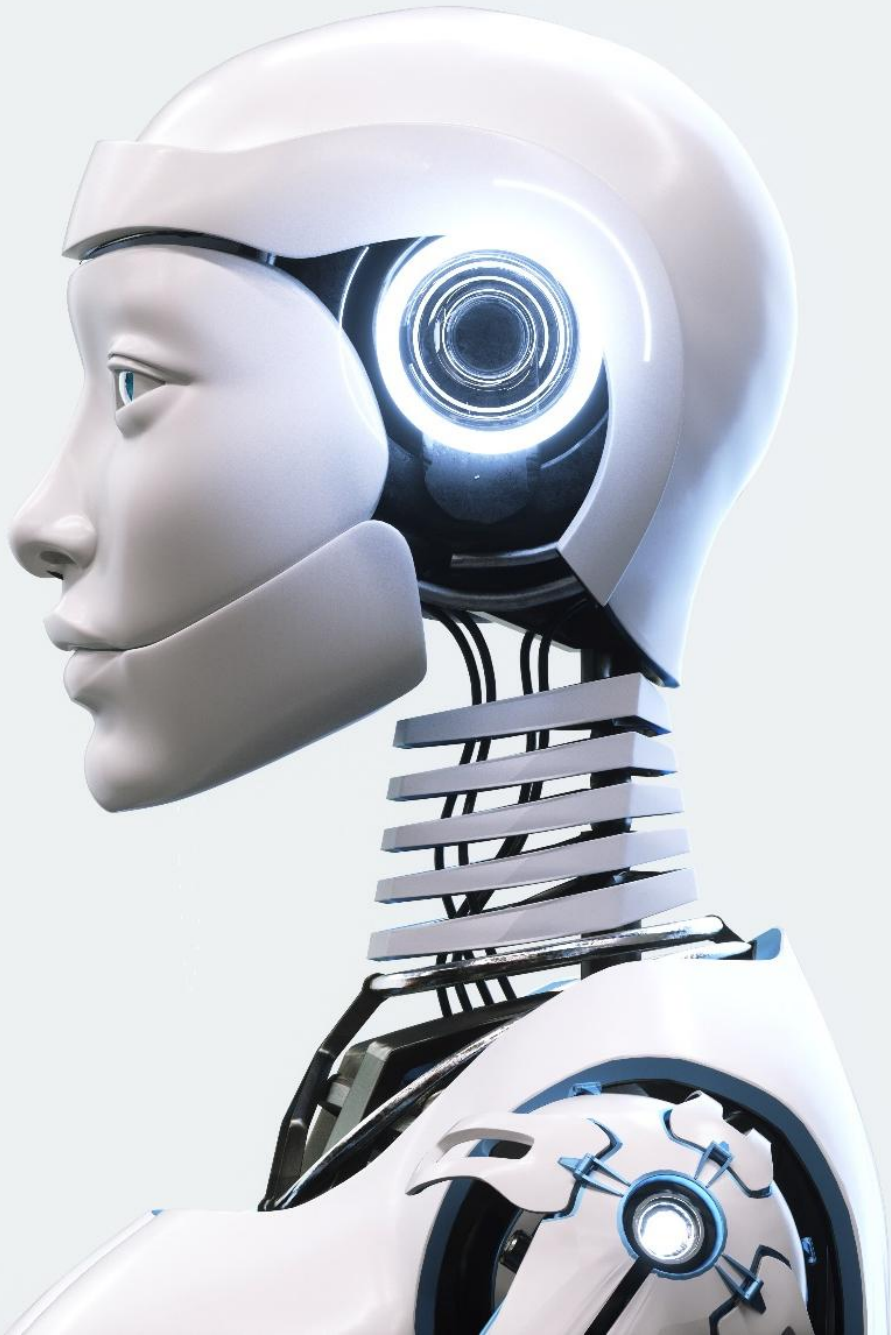
Department of Homeland Security



Attacks on LLM

One major issue that is on the horizon involves attacks on LLM. The Mitre group has an article on LLM prompt injection ““An adversary may craft malicious prompts as inputs to an LLM that cause the LLM to act in unintended ways. These "prompt injections" are often designed to cause the model to ignore aspects of its original instructions and follow the adversary's instructions instead”” Other sources have also reported on prompt injection. Prompt injections can be subdivided into direct and indirect. Direct is an actor literally, directly entering prompts. Indirect has the prompts coming through a separate data channel feeding into the LLM.

- <https://atlas.mitre.org/techniques/AML.T0051/>
- <https://portswigger.net/web-security/llm-attacks>



Attacks on LLM

In 2023 OWASP published a top 10 for Large Language Models. The 10 are

- Prompt Injection
 - Insecure Output Handling
 - Training Data Poisoning
 - Model Denial of Service
 - Supply Chain Vulnerabilities
 - Sensitive Information Disclosure
 - Insecure Plugin Design
 - Excessive Agency
 - Overreliance
 - Model Theft
-
- <https://owasp.org/www-project-top-10-for-large-language-model-applications/assets/PDF/OWASP-Top-10-for-LLMs-2023-Slides-v09.pdf>

EXAMPLES

1. LLM output directly entered into a backend function, causing remote code execution.
2. JavaScript or Markdown generated by the LLM is interpreted by the browser, resulting in XSS.

PREVENTION

1. Apply input validation on responses from the model to backend functions.
2. Encode output from the model back to users to mitigate undesired code interpretations.

EXAMPLES

1. Malicious influence on model outputs via targeted, inaccurate documents.
2. Model training using unverified data.
3. Unrestricted dataset access by models leading to control loss.

PREVENTION

1. Verify training data supply chain and data source legitimacy.
2. Employ dedicated models per use-case.
3. Implement sandboxing, input filters, adversarial robustness.
4. Detect poisoning attacks via loss measurement and model analysis.



EXAMPLES

1. High-volume task generation through specific queries.
2. Unusually resource-consuming queries.
3. Continuous input overflow exceeding the LLM's context window.
4. Repeated long inputs or variable-length input floods.

PREVENTION

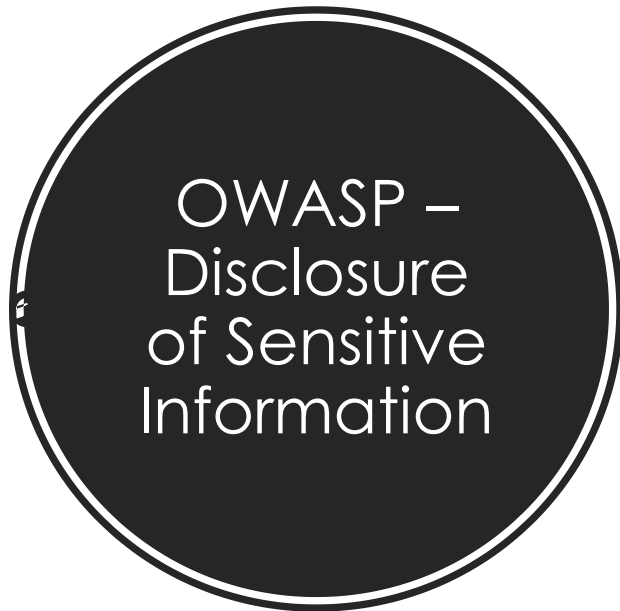
1. Implement input validation and sanitization.
2. Cap resource use per request.
3. Enforce API rate limits.
4. Monitor LLM resource utilization.
5. Set strict input limits based on the LLM's context window.
6. Promote developer awareness about potential DoS vulnerabilities.

EXAMPLES

1. Use of vulnerable third-party components or base images.
2. Use of a tampered pre-built model for fine-tuning.
3. Use of poisoned external data sets for fine-tuning.

PREVENTION

1. Vet data sources and suppliers, including their T&Cs and privacy policies.
2. Use reputable plug-ins and ensure they have been tested for your application requirements.
3. Maintain an up-to-date inventory of components using a Software Bill of Materials (SBOM).



EXAMPLES

1. Malicious manipulation of model's training data.
2. Training models using unverified data.
3. Unrestricted model access to datasets.

PREVENTION

1. Utilize data sanitization and robust input validation.
2. Implement least privilege principle during fine-tuning.
3. Limit and control access to external data sources.

EXAMPLES

1. Plugins accepting undifferentiated parameters.
2. Plugins taking URL strings instead of query parameters.
3. Plugins permitting raw SQL queries.
4. Lack of distinct authorizations for chained plugins.

PREVENTION

1. Enforce parameterized input with type and range checks.
2. Apply OWASP's recommendations for input validation.
3. Utilize least-privilege access control.
4. Use robust authentication like Oauth2.
5. Require user confirmation for sensitive plugins' actions.

EXAMPLES

1. Unnecessary or high-privilege plugin functions accessible to LLM.
2. Lack of proper input filtering in open-ended functions.
3. Over-granted permissions to LLM plugins.

PREVENTION

1. Limit plugin/tools accessible to LLM.
2. Implement only necessary functions in plugins.
3. Avoid open-ended functions, prefer granular functionality.
4. Limit LLM plugins' permissions on other systems.
5. Use OAuth for user authentication, granting minimum necessary privileges.
6. Require human approval for all actions.

EXAMPLES

1. External unauthorized access to LLM repositories.
2. Leaking models by insiders.
3. Network/application security misconfigurations.
4. Shared GPU services exploited for model access.
5. Replication of models via querying or prompt injection.
6. Side-channel attacks retrieving model data.

PREVENTION

1. Strong access controls/authentication for LLM repositories
2. Limiting LLM's access to network resources.
3. Regular monitoring/auditing of LLM-related activities.
4. Automated MLOps deployment with governance.
5. Rate limiting and exfiltration detection techniques.

- There was a recent article regarding Sleeper Agents in LLM. Sleeper agents work by the creator of an LLM having malicious intent. The malicious actor places a backdoor on the model. Machine learning backdoors are behaviors that are activated when a particular trigger is included in their input data. Without the trigger the model functions as expected.
- Another paper was published about using LLMs into proxies for malware attacks.
- <https://arxiv.org/pdf/2401.05566.pdf>
- <https://arxiv.org/pdf/2308.09183.pdf>