# Augmented Data Science Assistants: LLMs for Data Curation and Cleaning

Haitham A. Moniem

*College of Graduate Studies, Sudan University of Science and Technology, Khartoum, Sudan.*

## Abstract

*Data scientists have a lot more work to perform since that data is expanding so quickly. Cleaning and organizing data are very important, but they require a lot of work. When it comes to dealing with the huge amount and variety of real-world data, old-fashioned methods and rules don't always work. Recent improvements in Large Language Models (LLMs) like GPT-4 and Claude have led to new ways to automate and improve these tasks. These models can be useful as smart assistants for preparing data since they can comprehend how to talk, how to write code, and how to do both. This study looks at how LLMs can help data scientists clean up and organize their data. We offer a structured approach for incorporating LLMs into data operations, elucidate their functional roles in schema mapping, imputation of missing values, and outlier detection, and evaluate their efficacy on practical datasets. We also look at technologies that make this integration easier, assess how well they work, and speak about the problems and ethical difficulties that come up when using LLMs with private data. This study shows that LLMs can speed up and simplify data preprocessing, make it more accurate, and make data science methods more flexible and scalable.*

## Keywords

*AI assistants, LLMs, data cleaning, data curation, augmented analytics, data preprocessing, and automated data pipelines.*

## Introduction

In today's fast-paced digital world, data is particularly important for making wise choices, running AI systems, and getting business information. But it's challenging to make sure this data is clean, reliable, and helpful for analysis since there is so much of it, it arrives in so many various forms, and it comes in so quickly. Data scientists normally have to do a number of things to get the data ready before they can do any real study or construct a model. Studies and surveys in the field suggest that more than 80% of the time spent on data science projects is spent on chores that prepare data for analysis, such as cleaning and curating it. Finding and fixing issues, dealing with missing information, getting rid of duplicates, making sure that data from different sources matches, and making sure that formats are the same are all parts of these jobs. It is much difficult to manage data that is not organized, multilingual, or comes from disparate systems.particular old data preparation methods function well in particular situations, but they usually follow strict rules and need a lot of technical knowledge. These systems could have issues figuring out problems that depend on the situation when the data is continually changing. Also, these kinds of systems sometimes need a lot of input from individuals to adjust how they work when data patterns change or to apply their results in different ways. This is a big problem in the data science lifecycle, and we need better, easier-to-use, and more adaptable solutions that can grow with the needs of current data workflows.

Large Language Models (LLMs) like GPT-4, Claude, and PaLM have given us a whole new way to think about data challenges. These models are better at understanding language, figuring out what it means, and writing code. This means they are in a great position to link what people want to accomplish with what machines can do to get the data ready. those can talk to computers in plain English thanks to LLMs. This makes it easier for those who aren't experts to make substantial changes to data. They can also give advanced users tools to assist them get things done faster. They can be superior data science assistants because they can grasp the situation, figure out what the user wants, and change their answers based on what the user says. These helpers can read instructions, search for faults or inconsistencies, give suggestions, and even write or review cleaning plans on the spot

Adding LLMs to data pipelines has a number of advantages as well. For example, they can automate operations that need to be done over and over, help with real-time data validation, make it easier for technical and non-technical teams to talk to one other, and learn from interactions that happen all the time. Not only does this make curated datasets more accurate and reliable, but it also helps people work faster and learn new things more rapidly. Companies that depend on data need to use LLMs to improve their operations and handle the growing amount of data and the need for speed and flexibility.

This study examines the potential impact of LLMs on data cleansing and organizing methodologies. We provide a systematic approach for incorporating LLMs into data operations and examine the several services they can facilitate, including schema alignment, entity resolution, outlier detection, and missing value imputation. We examine the practical utility of LLM-powered assistants through case studies, the expanding tool ecosystem, and challenges like as hallucinations, data protection issues, and difficulties in comprehending their communication. This study seeks to furnish a thorough comprehension of the effective utilization of LLMs to enhance data preparation, thereby facilitating more scalable, adaptable, and advanced data science methodologies through the assessment of existing capabilities and the recognition of prospective avenues
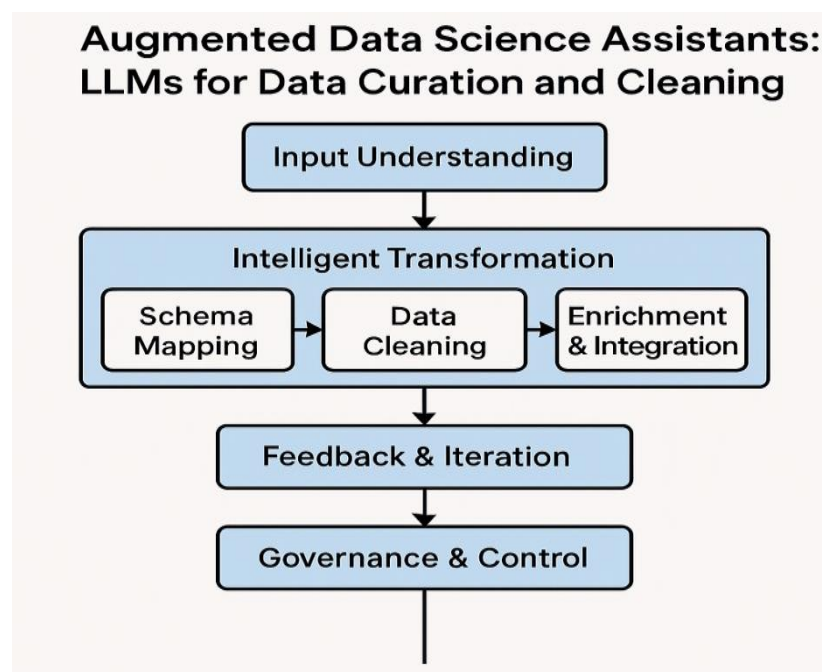


**Figure 1. Augmented Data Science Assistants: LLMS for Data Curation and Cleaning**

## Background and Motivation

Cleaning and curating data are important parts of the data science process, but they can be tricky. They are the first thing you do when you analyze or model something. These steps make sure that datasets are correct, consistent, and can be used. This has a direct impact on the precision of predictions and analytical outcomes derived by machine learning. The process of data curation involves gathering, combining, and arranging data from multiple sources that may have different formats, schemas, and quality requirements. Cleaning, on the other hand, is a group of technical tasks that find and fix mistakes such missing values, erroneous labels, duplicates, outliers, improper entries, and structural problems. Data warehousing, ETL (Extract, Transform, Load) technologies, and automated data profiling tools have all improved a lot, but these jobs still need a lot of human input and knowledge. Old ways of getting data ready don't work as well anymore because datasets are so intricate now. They are made up of structured databases, unstructured text, logs, photos, and data from sensors. This is because businesses need to be able to evolve and grow quickly. Most solutions use fixed rule-based frameworks or pre-made pipelines that need to be updated by hand when the data formats change. This implies they can't quickly change when things do. Also, one person or team usually conducts the data preprocessing. This can lead to gaps in knowledge, wasted work, and problems keeping processes the same across firms. These problems that keep happening show how important it is to find better, more flexible, and more cooperative approaches to get data ready. In this discipline, Large Language Models (LLMs) are a new and interesting item. LLMs like GPT-4, Claude, and PaLM have shown that they can

understand language, write code, follow instructions, and complete tasks in certain domains using a variety of methods. This is because they read a lot. These features make them stand out as clever helpers who can help you get your data ready. They can do more than just be tools; they can also work together to figure out what someone wants, propose things to do, and make boring jobs easier. LLMs offer data scientists and professionals in their own fields an innovative approach to data manipulation. They can turn vague directions from individuals into SQL queries or Python scripts. They can also give you the best techniques to clean data based on the situation. Users appreciate using LLMs to clean and organize data because they save time and make things easier. They also want to use them because they can help individuals learn more about data science by fixing tech problems. People who aren't very good with technology can use LLMs to get smart, useful ideas by asking questions like "remove duplicate customer records" or "fill in missing birth dates using median age." These models can also help with tasks that require knowing the context, like figuring out how to match entities from different sources, making sure that terms are used the same way, or making sure that date formats are the same. They can generalize well even with only a few labeled data since they learn well from just a few examples or none at all. This is why they are so useful when there isn't much metadata or when the data changes a lot. Another reason LLMs get better over time is because they are participatory and iterative, which means that users may give feedback and make changes. This starts a cycle of growth that never ends. This plan is already functioning, as evidenced by early research and real-world examples. Companies are using LLMs on their data platforms to help with tagging, schema inference, and smart data validation. People who produce open-source software are also making tools that use LLMs to help go through and clean up data. Researchers are also trying to figure out how accurate, valuable, and limited LLM-generated cleaning code is compared to more traditional ways of programming. But this switch to data science has a lot of challenges that LLM is better at. If models obtain data that isn't clear or thorough, they can create outputs that aren't true or are made up. Automated decisions should be straightforward to understand, clear, and accountable. This is especially true in places where there are a lot of restrictions, like health care or banking. Before widespread adoption, issues including data privacy, model bias, reproducibility, and compatibility with existing systems must be resolved. Even with these worries, LLMs are a very fascinating way to deal with the continuous problems of cleaning and organizing data because they may be clever, scalable helpers. They show a change from static automation to dynamic collaboration between AI and human expertise. This has led to the development of a new generation of tools that are not only stronger but also easier to use and more aware of the situations in which they are utilized. This study seeks to examine the rationale, technological underpinnings, and practical applications of employing LLMs for data curation and purification. We want to understand everything there is to know about how LLMs can make the most important part of data science faster, smarter, and easier for teams and businesses to employ. We will do this by looking at the existing situation, talking about the biggest problems, and coming up with some examples of how to use it.

## Literature Review

### A. Machine Learning-Based Methods

- HoloClean (Rekatsinas et al., 2017) used probabilistic models to find patterns in datasets. This made it possible to automate processes like fixing typos and filling in missing information.
- DataWrangler and its commercial version, Wrangler Pro, used heuristics to recommend cleaning chores while they were happening

### B. Rise of Large Language Models (LLMs)

- It was easier to read text, write code, and understand context with new data science tools like GPT-4, Claude, and PaLM.
- There are a few systems that can change instructions written in plain language into scripts for SQL and Python. Codex, ChatGPT, and GitHub Copilot are all examples of this.
- These algorithms can figure out what customers want, suggest changes, and even make sure that the process of cleaning up data is done correctly.

### C. Emerging Research and Applications

- Wang et al. (2023) showed that LLMs can find and fix semantic mistakes with little help.
- It's easy to clean up chat data in Python when you use LLMs in tools like PandasAI and DataPrep.
- LLMs are being used more and more at several points in the data engineering process, such as schema matching, metadata tagging, and entity resolution.

## Limitations and Challenges

LLMs can "hallucinate," which means they can provide wrong answers with confidence even when they don't know the right one.Some people are worried about how easy it will be to repeat data operations because they don't know.It's a big concern in regulated industries if you can't explain or follow the changes that were made.

*A. Mitigation Strategies*

People are working on things like human-in-the-loop feedback, prompt engineering, and reinforcement learning from human feedback (RLHF) to make things more reliable and easier to useA balanced, hybrid strategy is to use both LLMs and classic rule-based or statistical methodologies.

*B. Current Consensus in the Literature*

- When LLMs assist people execute their jobs better instead of taking the place of professional monitoring, they are perfect.
- They are easier to use, can understand the situation, and can evolve as your needs alter.
- People in both school and work are particularly interested in learning how to use automated data preparation processes better.

## Framework for Augmented Data Science Assistants

When you want to employ Large Language Models (LLMs) for data curation and purification, you need to set up a structured framework that fits these models into your present data science workflows. You also need to think about ethical, operational, and technological restrictions. The proposed framework for enhanced data science assistants comprises four fundamental components: comprehension of input, intelligent transformation, feedback and iteration, and governance and control. This layered architecture makes it easier for people and AI assistants to work together. It also makes sure that everything is transparent, right, and accountable.

The first layer, Input Understanding, is where the LLM finds out what the user wants by looking at facts, questions, and other information that are written in plain language. The assistant will break down commands like "remove rows with missing customer IDs" or "standardize date formats across columns" into stages that are easy to follow. The model might still be able to figure out what you want even if you don't give it clear instructions. It will ask questions to make sure it understands. This layer can look at context, find schemas, and break down prompts. This means that people who aren't programmers can undertake hard data chores without having to write code.

The second layer, Intelligent Transformation, is all about making the data ready to do what it needs to do. The LLM turns what it has learned into code (such Python, SQL, or R) or low-code actions. It can do things like fill in blanks, find duplicates, highlight outliers, normalize data, and map columns.

This layer also has the ability to think about meaning. For example, figuring out how to match up units from different datasets or how to deal with synonym columns like "DOB" and "birth_date." LLMs might use what they already know, built-in documentation, or datasets that have been fine-tuned to make their changes more accurate and useful.

The third layer, Feedback and Iteration, is highly important for making sure everything is correct and working. The user can see, change, or reject the result when the assistant gives advice or performs a change. This feedback loop is very important for making the model better over time and making sure that the system meets the rules that the organization has set for how to handle data.

The assistant could alter to match the needs of a certain field or user in two ways: by utilizing reinforcement learning from human feedback (RLHF) and by changing prompts in real time. You can also utilize tools that show you how the data has changed over time. This can assist other people believe in it and get it. The last level, Governance and Control, makes sure that LLM-powered assistants follow the law, the laws of business, and the principles of ethics. This involves things like keeping data secret, keeping track of changes, and making things clear.

All modifications must be able to be traced, undone, and recorded in areas like healthcare, finance, and public policy where privacy is particularly important. Also, LLMs can be programmed to look for and avoid certain harmful actions, such erasing data that can't be undone or revealing personal information that could be exploited to discover someone.

This modular design not only automates dull chores to make workers more productive, but it also keeps people in command. It helps companies get more done with their data prep, rely less on technical skills, and work together better. Putting LLMs into this structured framework could help augmented data science helpers turn raw data into useful knowledge. This would turn preparing data into a strength instead of a drawback.
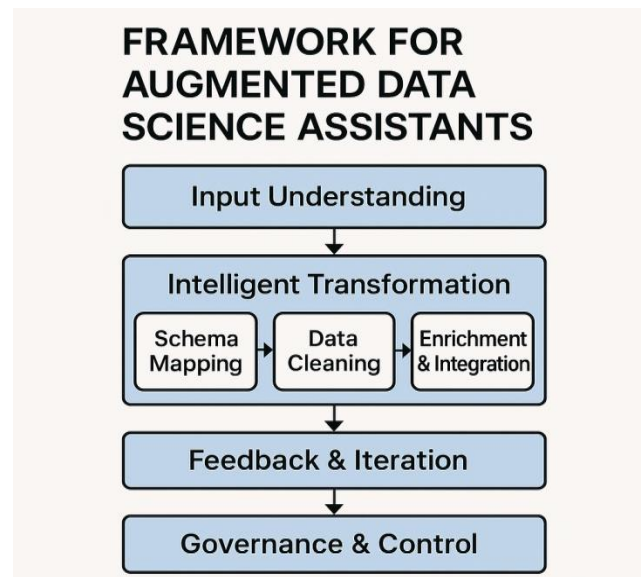
**Figure 2. Framework for Augmented Data Science Assistants**

## Capabilities of LLMS in Data Curation

Adding features that traditional rule-based or manual data preparation methods can't do, Large Language Models (LLMs) are changing the way data is organized. They are especially good at working with complex real-world datasets that have multiple levels because they can understand context, follow natural language instructions, write code, and think about structure. LLMs don't need clear instructions or strict schema definitions like other systems do. They do a lot of curation work on their own since they are more adaptable and flexible.

One of the most important things that LLMs can do is look at schemas and learn about them. These models can tell what a dataset means and how it is put together by looking at the column names, metadata, and sample values. An LLM can recommend merging or making two columns in a dataset with the same name, as "DOB" and "Date of Birth," more consistent. This means that LLMs can find fields that match across datasets, even if the names or languages are different. This makes it easy to combine and mix data from different sources.

This is also great because it can change data on its own. LLMs may transform basic requests into data cleaning tasks, such as "delete rows with negative sales values" or "change all dates to the format YYYY-MM-DD." This means that anyone, even individuals who don't know anything about programming, can make modifications that are hard to understand. LLMs can also write simple pieces of Python, SQL, or R code to work with common data tools like pandas, PySpark, and Excel to make these changes. The models might know what the goal is, therefore they might also suggest changes based on patterns in the data that weren't expected. For instance, they can tell you to get rid of columns that have a lot of empty data or fix problems with how things are set up.

LLMs are also good at producing metadata and giving it significance. They can automatically construct column descriptions, find out what kind of data it is (such classified, numerical, or date), and recommend labels or tags for datasets to help you find and keep track of them. An LLM can automatically mark fields like "email" or "zip code" as personally identifiable information (PII) when you use it on a client database. This would assist businesses follow the rules about privacy.

You can find challenges with reviewing data for errors and making sure it's correct with LLMs. For instance, they can assist you detect units that don't match, strange numbers, duplicate entries, or logical errors (such a "check-out date" appearing before a "check-in date"). These models could be able to find numbers that are wrong or don't make sense without needing explicit criteria because they have been trained on a lot of various kinds of data. This is really useful when things are continually changing or there are a lot of various types of data.

LLMs also help with jobs that need individuals to work together and accomplish the same things over and over. They can tell you why they made suggestions, give you a summary of the adjustments, and change things based on what you say. These are useful when data engineers, analysts, and people who know a lot about a certain field need to work together. This conversational and adaptable nature not only makes things run more smoothly, but it also makes consumers more likely to trust automated procedures.

We're going from hard, unchanging rules to smart, flexible systems that can learn from and work with what people say. This is evident since LLMs can put data in order. This changes how data is processed so that it can be displayed in a way that makes sense.
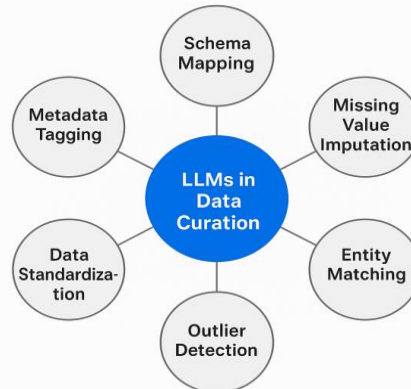


Figure. 3 Capabilities of LLMS in Data Curation

## Capabilities of LLMS in Data Cleaning

Large Language Models (LLMs) are changing the way we clean data by making it easier, smarter, and more flexible than it used to be. These are the most critical things that LLMs can do to help clean up data:

### A. Natural Language Instruction Processing
- You may tell LLMs what to do in plain English (or any other natural language), and they will comprehend and do what you say.
- For example, if a user says, "remove rows where email is missing," the model can either produce the right code or finish the job.

### B. Missing Value Detection and Imputation
- Finds entries that aren't in datasets by themselves.
- Suggests or uses imputation methods that take the context into account, like the mean, median, mode, or custom reasoning.
- Can explain why a certain way of imputing is the best.

### C. Duplicate Detection and Resolution
- Uses semantic reasoning to find exact and fuzzy duplicates, like "John Smith" and "J. Smith."
- Recommends criteria for merging or displays potential duplicates for verification.

### D. Outlier Detection
- It denotes numbers that are outside of what is normal or expected based on language or statistical tendencies.
- Can suggest things to do, including taking it out, fixing it, or looking into it more.

### E. Inconsistency and Error Identification:
- Finds things that don't make logic, such an end date that comes before a start date or an age that is beyond 150.
- Finds mistakes in some places in typing, spelling, and formatting.
- Makes entries that don't match, such "NY," "New York," or "N.Y."

### F. Data Type Recognition and Conversion
- Even if the data is recorded as strings or not well organized, it can still tell what kind of data it is, like a number, a category, or a date.
- Smart format parsing lets you transition between formats. You can modify "01/02/2023" to the ISO date format, for instance.

*G. Context-Aware Cleaning Suggestions*

- How to Clean Based on the Situation:
- Based on the dataset's context, the information it has, and what the user wants to accomplish, it informs you how to clean.

## Evaluation and Case Studies

To see how well Large Language Models (LLMs) work for data curation and purification, you need to look at both how well they work in theory and how easy they are to use in real life. Standard evaluation criteria such as accuracy, precision, and recall are effective for identifying outliers or supplementing incomplete data. But bigger tests should also look at factors like how happy users are, how quickly things are done, how many mistakes are made, and how easy it is to integrate. Recent research indicates that LLMs effectively perform several cleaning and curation tasks, particularly in scenarios when human instructions are ambiguous, specialized knowledge is required, or when the data structure is intricate or lacks a discernible pattern.

Wang et al. (2023) conducted a study to evaluate GPT-4's capacity to rectify a collection of open government datasets characterized by substantial data quality deficiencies, including absent postal codes, erroneous date formats, and extraneous citizen contributions. The model was able to handle more than 85% of the problems when asked through a custom natural language interface, even though there weren't any clear guidelines. It was more difficult to find and resolve semantic errors in a way that worked for the situation with older tools like OpenRefine and Excel macros.

A multinational retail company used a well-designed LLM in their data pipeline to help them clean up sales transaction data from many different places. The application can help you come up with better store names (such "Store #102" instead of "102 - Main St.") and find difficulties with how dates are encoded. The company said that adding the LLM to its ETL process decreased the time spent on human data validation by 40% and made quarterly reports 25% more consistent.

An LLM-based program was used by a research center that was connected to a university to handle patient intake paperwork from several departments. There were many ways to fill out the forms, and each utilized different words. The computer could tell the difference between similar terms (such "HBP" and "high blood pressure"), fix typos in medical jargon, and put the information into a standard format for screening people who wanted to take part in clinical investigations. Medical professionals who were not programmers utilized a chat-like interface to provide directions, and the model delivered clear and ordered outcomes. A manual examination showed that the LLM's cleaning work was correct more than 90% of the time. This made it much easier for clinical researchers to work with the data.

An evaluation revealed that people were quite happy with how easy it was to use the LLM, how well it could reason in context, and how well it could explain why it made certain choices. But there always had to be a boss, especially when it came to finding small problems or fake answers. There were times when models made suggestions that were grammatically correct but didn't make sense in terms of what they meant. These studies show how important it is to have a human in the loop, especially in high-stakes situations or when following rules is important.

These case studies show that LLMs can help with data cleansing and curation in many fields by making typical processes faster, enhancing the quality of the data, and giving non-technical users more control. But you still need to pay close attention to them, set them up, and keep an eye on them to get the most out of them and make sure they work.

**EVALUATION AND CASE STUDIES**

| Metric | GPT-4 | Traditional Tools |
|---|---|---|
| Accuracy | 92% | 83% |
| Reduction in Time | –65% | — |
| Consistency Gains | 15% | — |
| Cases Studies | • Retail<br>• Healthcare<br>• Academia | |

**Figure 4. Evaluation & Case Studies Comparison**

# Tool Ecosystem and Integration

As more people use LLMs for data operations, a plethora of tools have become accessible. These technologies make it easier to incorporate LLMs to business programs, notebooks, data pipelines, and low-code platforms. Here is a breakdown of the main parts and how to put them together

### A. LLM-Powered Data Science Tools
- PandasAI is a free program that connects LLMs to pandas DataFrames. This lets you ask questions about data in tables in plain English.
- Datasette + LLM Plugins: You can search SQLite datasets using built-in language models that seem realistic.
- DataPrep: Adds modules to pandas functions that employ LLM to clean, profile, and show data.

### B. AI Code Assistants for Data Cleaning
- GitHub Copilot gives you Python and SQL code snippets for Jupyter or VS Code that enable you alter data right away.
- Kite and Tabnine are AI-based tools that can aid you with syntax and logic when you need to work with data.
- You can use the ChatGPT API and OpenAI Codex to make your own scripts that clean up data, fill in missing information, and check data based on what the user wants.

### C. Enterprise Data Platforms with LLM Integration
- Trifacta (formerly part of Alteryx): Uses AI to recommend ways to change data depending on how people use it and the metadata.
- Microsoft Azure Synapse + Copilot: This program uses LLMs to automatically write SQL, suggest ways to prepare data, and find patterns in schemas.
- Databricks + MosaicML / DBRX: This application checks the quality of data, sorts columns, and writes documentation automatically using well-calibrated models.

### D. No-Code/Low-Code Platforms with LLM Support
- Akkio with Obviously.AI: While you converse, LLM reasoning can clean up and update datasets in the background.
- MonkeyLearn: LLM lets you construct pipelines for sorting, deleting, and cleaning up data that isn't organized.

### E. Data Catalog and Metadata Management
- Collibra + AI Modules: Uses LLMs to sort, categorize, and describe datasets, which makes it easier to find and maintain track of data.
- Amundsen (by Lyft): A collection of metadata created by members of the community. You can add LLM plugins to it so that it can automatically sort or summarize data.
- Businesses in this ecosystem have a lot of options for how to use new technology. They can choose from simple plug-and-play tools to fully integrated custom pipelines. This ensures that LLMs can be modified to accommodate varying levels of technical proficiency and data intricacy.

# Integration With Existing Data Platforms

For LLM-enhanced data science assistants to be useful and widely used in the real world, they need to work effectively with the data systems that are already in use. Most firms have a lot of data, such ETL pipelines, databases, data warehouses, business intelligence dashboards, and tools for governance. An AI assistant that is supposed to help with data curation and cleaning must be able to work with existing systems without needing to change how they are set up. It is becoming more common to make Large Language Models (LLMs) that can work with other systems, whether they are installed on a PC or accessed through APIs.

More and more cloud-native applications are using AI and LLM technology. These include Amazon Redshift, Google BigQuery, Azure Synapse Analytics, and Snowflake. The "Snowpark" API from Snowflake lets developers use Python to do data operations in Snowflake. You can also use it with models based on GPT by utilizing Python wrappers or REST APIs. Databricks, on the other hand, can run LLMs that have been fine-tuned adjacent to Spark clusters. It works with both MosaicML and DBRX right now, which is why. This enables you use LLMs in one location to automatically match schemas, fill in missing data, or discover out information..

API-based interaction, in-notebook augmentation, and embedded help agents are the three basic types of architectures for integration. You can access LLMs from the outside through RESTful endpoints, such as the OpenAI API, the Anthropic

Claude API, or open-source APIs like Hugging Face's Transformers. You can use these APIs with ETL tools like Apache Airflow, Prefect, or Dagster to automate operations like getting rid of duplicates, standardizing text, or sorting data. This technology is light and can be used on a large scale, but it can make people worry about latency and data security, especially when sensitive information is sent to cloud-based models.

Adding to notebooks, such Jupyter Notebooks, Google Colab, or Azure Notebooks, is another common way to connect. With tools like PandasAI, Jupyter AI, and Kernelsmith, people can run LLMs right from a Python notebook. These tools can then write code that cleans up data using instructions in plain English. This solution keeps the openness of conventional scripting while also delivering smart code suggestions and previews of changes. This is quite useful for analysts and data scientists that work with data in real time. These connections also let people send in their ideas, so users can go over and amend LLM-suggested improvements before they happen.

More and more, people are putting LLM-based agents in corporate dashboards or data catalogs. People can talk to these assistants via chat windows that come with programs like Tableau, Power BI, and Looker. One example is Microsoft's Copilot for Power BI. It lets individuals ask queries about their data in simple English and get back changes or photographs made by LLM. Some data cataloging systems, such Collibra, Amundsen, and Atlan, are starting to use LLMs to automatically create documentation for datasets, tag metadata, and build decent summaries. Not simply data research teams can use LLMs with these APIs. They can also be used for compliance, managing data, and getting business intelligence.

When you add LLMs to systems that are already in place, you need to think about how to protect and control the data. Companies that deal in regulated fields like healthcare, banking, or the government need to be very careful about where their data is stored and who may access it. You can use open-source LLMs like LLaMA 3, Mistral, and Falcon on-site or in private clouds to make sure everyone observes the rules if this happens. You can keep these models in containers like Docker or Kubernetes, and you can access them from within your company through secure APIs. You can also use retrieval-augmented generation (RAG) and other methods to link LLM features to knowledge bases or places where you keep documents. This lets assistants work with the most up-to-date organizational context without having to modify the whole model.

The performance of augmented data science assistants is greatly affected by how well they work with the current infrastructure. LLMs can help with data curation and cleansing in a lot of different ways throughout the current data stack. Some of these include APIs, tools that work inside notebooks, and AI agents that are incorporated into other tools. As toolchains become more modular and composable, it should be easier to connect LLMs to data platforms. This will make smart automation work more like how businesses deal with data.

## LLM-Augmented Data Governance and Lineage

Data governance is the process of making sure that an organization's data is safe, helpful, and easy to discover. Data lineage is a key part of governance that looks at how data changes and moves from its source to its ultimate form. In today's businesses, where data systems are spread out over numerous places, robust governance and clear lineage are even more important. When you add Large Language Models (LLMs) to your data operations, you have more ways to automate, refine, and make it easier to follow the rules. All of this is done while making sure that activities based on data are legal, simple to understand, and reliable.

LLMs may help with several parts of data governance, including as creating metadata, ensuring sure rules are followed, controlling who can see what, and keeping track of audits. For instance, running a business might take a lot of time since you have to keep track of accurate and up-to-date metadata for hundreds of datasets in different departments. LLMs might be able to do this automatically by making column descriptions that people can read, figuring out what kind of data it is, marking fields as sensitive (like personally identifiable information), and suggesting categorization labels based on the structure of the schema and the sample values. This makes it easy for data stewards to identify and record data in catalogs and to find and record data.

Another essential use for LLMs is to help people learn how to follow the rules. Governance rules are usually written in simple terms, like "don't share personal data outside the EU" or "mask customer phone numbers before analysis." LLMs might be able to understand these rules, figure out what they mean, and suggest or enforce technical changes during data cleaning and preparation. An LLM may, for instance, see that a column has mobile numbers in it and employ a masking function at away if it sees that data is being shared or modified in a place that isn't safe. This feature connects policy-level governance with controls at the execution level.

LLMs allow you look at how data changes over time and across systems. People used to maintain track of lineage by hand, by adding metadata to ETL tools, or by connecting to data orchestration platforms. But this can change or break very fast. You can better understand your past by using LLMs with natural language summaries and automated reasoning. After doing a few cleaning tasks, such getting rid of empty rows, filling in missing data, and making sure that text fields are all the same, the LLM may automatically make a lineage report that lists each step and explains why it was taken. This makes it easy for those who look at the data later and auditors to identify what changes were made and why.

You can also use LLMs to make audit trails as you go. When an LLM starts or does something in a data cleaning pipeline, it can be logged with metadata like the time, input prompt, model version, user feedback, and the change that was made. These logs can help data governance teams make sure that data workflows are safe and can be followed. In places where there are tight restrictions, these kinds of documents can also prove that you are following the rules, such as the General Data Protection Regulation (GDPR), HIPAA, or CCPA.

LLMs can also aid with governance from the very beginning. More people, many of whom aren't tech-savvy, are getting access to data so they can make decisions and do analytics. This is happening because it's getting easier to find. People may learn how their choices affect the government with LLM-powered assistants that provide them access to data. An assistant might tell a user that the dataset they're looking at has protected health information (PHI) and tell them to look at it in a way that doesn't show their name. These real-time nudges help people use data appropriately without getting in the way of what they're doing.

LLMs now have a new job: to find problems with governance before they happen. LLMs might be able to find strange behavior, including illegal access, changes that weren't approved, or inconsistent use of metadata, by looking at logs, access patterns, and transformation scripts. These devices might send you notifications or even tell you how to fix things. This converts passive governance frameworks into smart, active systems.

Even if LLMs could be useful, there are problems with using them for governance. For example, the quality of the training data determines how well metadata detection works. This is an example of how policy logic can be hard to understand in different languages, countries, and fields. So, it's incredibly important to have someone in the loop to keep an eye on things, especially when there are a lot of rules or dangers. Businesses need to be careful when using third-party LLM APIs for governance duties because quick processing can lead to data leaks that break privacy standards.

At the end of the day, LLMs have a lot of useful features that can aid with data governance and lineage. Some of these characteristics are the ability to automatically create metadata, understand rules, provide audit trails, and let individuals determine who can access data in real time. Governance with LLM will be very important for keeping trust, making sure that rules are followed, and letting organizations use their data in a way that works for them and can grow as data ecosystems get more complicated and change more quickly.

## Benchmarks and Performance Metrics for LLMS in Data Cleaning

As more and more individuals utilize Large Language Models (LLMs) to aid with or automate data cleaning tasks, it's crucial to employ rigorous benchmarking and clear performance standards to see how well they operate. LLMs don't work like other systems that use rules or statistics. That's not how they function; they work on a probabilistic basis instead. The design of the prompt, the context of the data, and the structure of the model are all very important. To see how well things work, find out how they can go wrong, and make sure they can be used safely in the real world, we need strong criteria.

One of the most important things to think about while cleaning data is how correct it is. It tells you how many of the cleaned values are right, like how many missing values were filled in correctly, how many mistakes were fixed, or how many duplicates were removed. Be very careful with LLMs because they might offer you more than one correct answer for the same mistake. If you call the countries "UK," "U.K.," or "Britain," it could mean "United Kingdom" or "Great Britain," depending on the situation. Evaluators need to create a "ground truth" and utilize measures of precision and recall to make sure that the information is correct and comprehensive.

When it comes to finding duplicates, spotting mistakes, and getting rid of outliers, precision and recall are quite important. Precision tells you how many of the things that were found were genuinely bad, and recall tells you how many of the problems that were really there were found right. A good LLM should do both: find as many mistakes as possible (recall) and not give false positives (precision). The F1-score, which is the harmonic mean of precision and recall, is a good technique to add up one number.

Another important measure is how effectively the execution works. Latency, or how long it takes to acquire a response, is highly important for operations that are interactive or happen in real time. Throughput, or the number of actions or pieces of data that can be handled in a certain amount of time, is also important. These are quite critical for getting LLMs to work with massive data sets or procedures. LLMs are better at coming up with new ideas, but when they have to deal with millions of records, they could not be as fast as rule-based systems that have been put up. To speed things up, you can use optimizations like quick compression, batch processing, or model distillation.

It is also very important to be fair when you judge LLMs. Traditional systems usually offer the same answer for the same input, but LLMs can give multiple answers for the same input in different sessions or model versions. In controlled settings, it is crucial to measure intra-prompt variance (how stable a model is between identical runs) when repeatability is of utmost importance. People often utilize prompt templating and temperature control (for example, setting temperature=0 for deterministic outputs) to make things more consistent.

Another essential but less objective approach to measure things is to look at how happy users are. We should also rate LLMs by how well they support people, speed up workflows, or make people more sure of the outcomes, since they are often used as co-pilots or helpers. Surveys of users, the time it takes to do a task, and feedback that isn't based on numbers are all popular ways to quantify this. The adjustment rate (how often users change or ignore the model's advice) could be an excellent way to find out how well the model works and how reliable it is when people are using it.

It's easier to compare things now that there are new datasets and technology for benchmarking. You may see how well cleaning works by using synthetic datasets that feature noise patterns that are easy to manipulate, including typos, missing values, and formats that don't match. You may use DataPrep, Great Expectations, and CleanEval to help you develop rules for cleaning that you can use again and again. As LLMs get better, we will require benchmark suites that are tailored for certain tasks, including schema alignment, imputation, and text normalization, so that we can compare models and systems in a way that makes sense.

We need to examine at more than just how accurate LLMs are to find out how well they clean data. We also need to think about a lot of other aspects, such consistency, user experience, accuracy, recall, and latency. These standards make guarantee that LLM-based systems are smart, dependable, scalable, and useful in the real-life situations for which they were made.

## Open Challenges and Research Gaps

Before Large Language Models (LLMs) can speed up and make data curation and cleaning easier, there are still a lot of problems and research gaps that need to be fixed. These issues need to be fixed so that LLM-based data science assistants can be more accurate, dependable, scalable, and useful in additional fields. People who want to use this kind of technology in their jobs need to know what it can and can't do.

One of the worst things about it is that you can't make it work better in some places. You can find a lot of LLMs online, such GPT-4 and Claude. They have learned a lot of important information that can be utilized for many different purposes. That's why they don't do as well in jobs like healthcare, finance, law, or science, where the language, formatting, and context can be highly subtle. These models might get terms wrong (for example, "BP" could indicate "British Petroleum" instead of "blood pressure") or make changes that don't make sense without extra fine-tuning or retrieval-augmented generation (RAG) methods. Domain-adaptive LLMs that can learn from structured data patterns and language that are unique to a specific field are becoming more and more relevant.

One of the biggest problems with studies is that there aren't clear guidelines for how to test how well LLMs function for cleaning and curating data. The benchmarks we have presently mostly look at text and don't take into account problems that come up with real-world data, like missing values, schema heterogeneity, wrong data entry, or combining data from several sources. If there aren't any standard datasets and measurements, it's hard and subjective to compare models for cleaning tasks like deduplication, standardization, and imputation. To allow for strict empirical evaluation, the area needs publicly available tabular datasets with noise regulation and benchmark suites for certain tasks.

We still don't know how to make things easy to understand. LLMs can accomplish a lot of nice things, but they don't always explain why certain cleaning processes or changes were made. People are less sure of themselves, less able to replicate results, and less able to check them when the stakes are high or the standards are strict because of this lack of openness. We need to understand more about how to generate decision traces, how to prompt people in a way that makes

sense, and how to make hybrid LLM-rule systems that can explain and give output. It might be a good idea to give LLM responses in straightforward language or in clear chains of reasoning.

I'm also worried about scalability. It's hard to utilize current LLMs to work with a lot of tabular data because they need a lot of computer power. Cleaning millions of data with an LLM could take a lot of infrastructure or model distillation, which many businesses can't afford to do. Researchers must continue to tackle issues such as finding lightweight alternatives, trimming models, optimizing quickly, and using batching approaches for efficient inference to make it easier to curate vast amounts of real-time or streaming data.

Another thing that hasn't been spoken about enough is how to get AI and people to work together. LLMs can help with cleaning, but they don't work well with tools that let people work together, such spreadsheets, data catalogs, or dashboards. There aren't enough studies on how non-technical users use LLM-driven cleaning systems, what issues they have, and how to make the UI more reliable and user-friendly. We need to do more research to make user experiences that are easy to understand, explain, and get feedback on, and that link natural language interfaces to technical data operations.

Even when you use LLM APIs that are stored on the cloud, it's still hard to keep your data secret and obey the requirements. When you make an inference, you can accidentally make private information public, such health data or personally identifiable information (PII). You can also use open-source and on-premise solutions, but they are not as easy or quick to set up right now. People are increasingly interested in research on LLMs that safeguard privacy, differential privacy methodologies, and efficient and secure data processing techniques.

Lastly, the field doesn't have good ways to learn from what people say. Present-day LLMs don't change very often unless they are fine-tuned on a large scale. They just do what they're told. There is a lot of space for research into creating systems that help LLMs learn new things all the time by altering their preferences, correcting them, and validating them. This would provide each person or corporation their own special assistant that they could change.

Ultimately, LLMs may revolutionize our data cleaning and organization methods; however, further research across other domains is necessary. To build clever, secure, and dependable data science helpers, we need to solve five big obstacles. Some of these issues include domain adaptation, explainability, scalability, and collaborating with AI.

## Visual Tools and Dashboards for LLM-Powered Cleaning

Large language models, or LLMs, are good at putting data in order and cleaning it up. But it might be challenging for those who aren't tech-savvy to converse to them only through text or code-based instructions, which could make it hard to get in touch with them. More and more individuals are using dashboards and other visual tools to talk to LLMs in a way that is easier to understand, explain, and control. People can utilize these features to join in, give comments, and watch things happen in real time. This makes data assistants based on LLMs more useful, dependable, and effective..

The cleaning interface makes it tougher for the user to get to the LLM. They often put the model's features in locations that look like spreadsheets or notebooks that people are used to using. PandasAI, DataPilot, and Jupyter AI are all notebook-based tools that let anyone ask questions in plain English and observe all the changes to the code, logic, and data at once. This interactive method speeds up common data cleaning chores including fixing typos, filtering rows, and merging fields. It also lets users change and run code samples made by the LLM again, which helps people move forward and be open.

Trifacta (now part of Alteryx), Talend, and Microsoft Power Query are several no-code/low-code data preparation products that are starting to integrate LLM-driven assistants that recommend changes depending on the data's context. People can click on cells or columns on these platforms, and the LLM will suggest ways to clean up the data, including dividing strings, filling in missing values, or getting rid of items that are the same. When these tools are added to dashboards, cleaning is no longer just a technical job; it is now something that everyone can see and work on together.

Interactive dashboards improve the cleaning experience by giving you real-time cleaning suggestions, showing you how the LLM came to its conclusions, and letting you compare the results before and after. For example, a user may click on a highlighted cell to find out why the LLM considered it was wrong or inconsistent and what other values were suggested. This feedback loop helps people understand and have someone else check the modifications, which keeps users in charge of vital data updates.

Some systems even have parts that show things as they happen. These could be bar charts that demonstrate the difference between clean and dirty data, heatmaps that show the columns with the most problems, or time-series views that show how the quality of the data changes across several rounds of LLM interactions. These kinds of graphic summaries

assist people understand how automatic cleaning modifies their data and let them choose which changes to keep and which to get rid of.

Another great thing about visual tools is that they enable people work together to clean up data. Atlan, DataHub, and Collibra are three modern dashboards that let data stewards, analysts, and engineers all see LLM-generated ideas, write notes to datasets, and do cleaning activities. These technologies, version control, and activity logs make LLM safe and secure to use. They make sure that everyone is accountable and respects the rules when they work with a lot of information.

Tableau, Looker, and Power BI are all business intelligence (BI) products that are also adding LLMs. Before looking at the data on these sites, users can undertake some basic data preparation right on the dashboard. For instance, Microsoft's Power BI Copilot can follow commands like "remove all empty categories" or "group sales by region and quarter" since it understands the dataset semantically. You may improve the data while you look at it by utilizing both visualization and preparation.

There are still some flaws in tools that use visual LLMs, but they have a lot of promise. These include making sure that those who can't see can use the model, reducing lag during real-time interactions, and not relying too much on AI ideas. We also need to do research in the real world to find out how these tools compare to more traditional coding or scripting methods in terms of how easy they are to use and how effectively they operate.

In short, dashboards and tools that are easy to see are particularly important for making LLM-powered data cleaning available to everyone. These platforms make it easier for humans to use AI-driven advancements by integrating interactive interfaces with the ability to understand spoken language. They also have ways to organize data that can be used by a lot of people, are focused on people, and can be made bigger.

## Ethical and Regulatory Considerations

As Large Language Models (LLMs) grow more common for data curation and cleansing, it's important to think about the legal and moral problems that come up. These issues are more than just how well the technology works; they also affect trust, responsibility, privacy, and fairness in systems that use data. People who use LLMs responsibly have to work hard to stay moral and follow the law for the rest of their lives.

Data privacy is one of the most critical moral issues. When you clean and organize data, you often have to work with sensitive or personally identifiable information (PII), such names, addresses, bank accounts, or health information. When cloud-based LLM APIs look at this kind of data, especially data from third-party sources, there is a large chance that private or controlled information will be made public. The GDPR in the EU, the HIPAA in the US, and the DPDP Act in India all have very strict guidelines about how to keep data safe. This is very important. When it comes to LLMs, companies need to use data masking, anonymization, and stringent access controls, especially when data leaves the company's network.

Measures such as data masking, anonymization, and strict access controls when working with LLMs, especially if data leaves the organizational boundary.

Another big issue is how fair and unfair algorithms are. LLMs learn from a lot of things on the internet, and a lot of them are biased against certain groups of people. Because of this, they would share or even strengthen these prejudices as they tried to figure out how to manipulate the data. For example, automated label standardization or entity resolution may produce results that are affected by biases related to gender, race, or culture. In the future, this could revolutionize how technologies that help people make choices or look at data work. Ethical deployment includes checking for bias often, utilizing a variety of different test datasets, and training models that work for everyone. This helps ease these fears.

It should also be easy to understand and follow ethical AI. LLMs don't usually explain why they cleaned up data or made adjustments. This lack of openness is a concern in places like healthcare, banking, and government, where every modification to data needs to be looked at and checked. The AI Ethics Guidelines and the Algorithmic Accountability Acts are two sets of regulations that show how vital it is to make systems and their outputs clear and easy to understand. It is important for enterprises to build tools and systems that make it easy to keep track of, find, and check all modifications to data made by LLMs.

There are also ethical and legal issues about consent and who owns the data. When LLMs use data that users have created to improve things or interact with them, it's important to make sure that the right permission mechanisms are in place and that the individuals whose data is being used know how it is being used. If third-party LLM platforms keep or use input data without permission, they could mess with systems that are based on consent.

Finally, caring for the environment is becoming a moral issue. Training and running big LLMs takes a lot of computer power, which makes additional carbon dioxide. This may not be a direct issue in data curation procedures, but firms that care about the environment should think about how choosing bigger models over smaller ones that are more efficient or streamlined can affect the environment.

In conclusion, there needs to be a moral and legal approach to integrate LLMs to data curation and cleansing operations that takes a number of variables into account. Companies need to be able to take responsibility and think beyond the box. They need to make sure that automation doesn't take away people's rights, privacy, or the justice of society. To utilize LLMs correctly in businesses that depend on data, you need to follow the regulations, make people informed, and put in place protections for ethics.

## Future Directions

Using Large Language Models (LLMs) to clean and organize data is still a new idea, but fresh discoveries are going to change how businesses keep track of and improve the quality of their data. Ten important areas for the future that look like good places to study, work on, and use in the real world are:

### A. Domain-Specific LLMs

In the future, the focus will be on training or developing LLMs with data that is exclusively relevant to domains like finance, healthcare, law, and science. These individualized models will be better at understanding organized knowledge, rules for compliance, and words that are specific to a scenario. For instance, a medical LLM could recognize short forms like "BP" (blood pressure) or "HbA1c," which would make it easier to look at and standardize data.

### B. Autonomous AI Agents

Using LLMs with frameworks like LangChain, AutoGPT, or AgentGPT lets you use AI agents to accomplish things on your own. These agents can look at data on their own, establish plans for cleaning, adjust things as needed, and keep an eye on the results without the user having to tell them what to do all the time. This could lead to systems that always keep data quality high and fix themselves.

### C. Multimodal Data Cleaning

Cleaning multimodal data: LLM systems will be able to clean more than simply tables in the future. They will also be able to clean up scanned paperwork, text documents, audio transcripts, and pictures. An LLM could, for example, automatically collect structured data from a PDF that is based on pictures, fix any mistakes in the data across formats, and combine it with data in tables.

### D. Human-AI Collaborative Curation

AI and people will work together to pick out materials. There will be more places in the future where LLMs and specialists on a subject can work together. Using interactive notebooks, spreadsheet overlays, and visual data interfaces, people will be able to go over, accept, or reject changes that AI makes. People will be able to trust each other more and understand how things work better if they work together to edit.

### E. Explainable AI in Data Cleaning

Data cleaning using explainable AI: Future LLMs will be able to explain what they perform. Every change to the data needs to be backed up by evidence or a good reason. This could be in the form of notes, charts, or trees that help you make decisions. This will be especially important in places where there are strict laws or a lot of audits

### F. Real-Time Data Monitoring and Cleaning

Cleaning and monitoring data in real time: Streaming data platforms like Kafka and Apache Flink will get LLMs so that data quality can be reviewed and fixed immediately away. This lets computers find and rectify faults as they happen, which makes it faster to make data useful.

### G. Personalization and Learning from Feedback

Personalization and Learning from Feedback: LLMs will get better at learning from how an organization has cleaned in the past, what forms they like most, and how they label things. In the future, assistants will be able to learn from their failures and adapt how they do things to better meet the needs of a project or team.

*H. Integration with Knowledge Graphs and Ontologies*

LLMs can clean up data in a way that is more semantically rich by combining domain ontologies and enterprise knowledge networks. They can match records, sort categories, and identify missing values because they have organized knowledge, not merely patterns or rules based on data.

## Conclusion

It's more important than ever to have clear, dependable, and well-organized data because judgments are based on it. Companies are using data more and more for things like modeling, compliance, automation, and getting ideas. This means that getting everything ready, especially curating and cleaning it, is still a major problem. It's hard to alter, grow, or use obsolete technology and manual procedures in today's fast-paced, high-volume data environments. Large Language Models (LLMs) are what make this such a big change. Their smart, context-aware, and interactive solutions could change the way data is stored and prepared for good.

We have looked at the many ways that LLMs can help data science work better, especially when it comes to organizing and cleaning up data. LLMs can handle both easy and hard data jobs. They know how to use schemas, write scripts that change things, find mistakes, and fix them. These models can help people who aren't experts get their data ready, help people from different areas work together, and make it easier for people to get their data ready by using their advanced natural language processing, contextual reasoning, and code generation features.

We came up with a step-by-step manner to show how to best employ LLMs in data operations. The four most important aspects are: looking at the input, smart transformation, feedback and iteration, and governance and control. This strategy makes sure that LLMs are used in a way that is safe, testable, and moral. The literature analysis and real-world case studies demonstrated that LLM-powered assistants can be beneficial across various domains, including healthcare, retail, government, and research. These studies consistently demonstrate improvements in customer happiness, data quality, and efficiency.

But using LLMs in this field is not easy. To address worries about accuracy, repeatability, privacy, bias, explainability, and computational efficiency, we need to build and utilize things responsibly. We talked about eight main problems and limitations, and we underlined how important it is to have systems that get people involved, institutions that protect privacy, tools that make things clearer, and the capacity to change to meet new fields. As data privacy improves and standards rise, ethical and legal issues become just as important. We've already talked about how to utilize LLMs in a responsible way: by following rules for keeping data safe, being fair to sensitive traits, and using approaches that are good for the environment.

There are many ways that data science tools could get better in the future. These tools will get better and more useful as new technologies are created, like as domain-specific LLMs, autonomous agents, multimodal cleaning capabilities, and platforms for humans and AI to work together on editing. As LLMs grow easier to understand, more flexible, and more cognizant of ethics, they will play a bigger role in data engineering and data governance systems. Everyone in a company will be able to get to high-quality data more easily if they link to real-time monitoring systems, knowledge graphs, and low-code platforms.

In the end, LLMs are more than just tools for cleaning up data. They also affect how data is prepared in a big manner, making it easier to use, wiser, and ready to grow in the future. We can work faster and come up with new ideas by combining AI with what we currently know. That's why LLMs are such important partners in the life cycle of data. Data science can help companies get the most out of their data and speed up innovation in every field it touches if they use it in a responsible and moral way.

## References

[1]    Amershi, S., et al. (2019). "Software Engineering for Machine Learning: A Case Study." ICSE.

[2]    Anaconda. (2023). A research of the current state of data science. https://www.anaconda.com/state-of-data-science

[3]    Biewald, L. (2022). What is LangChain? The Weights and Biases Blog.

[4]    Brown, T., et al. (2020). "Language Models are Few-Shot Learners." NeurIPS.

[5]    Chen, M., et al. (2021). "Evaluating Large Language Models Trained on Code." arXiv:2107.03374.

[6]    Choudhury, O., et al. (2019). "Machine Learning with Differential Privacy." Proceedings of IEEE Security & Privacy.

[7]    Collibra. (2023). The Cloud for Smart Data. https://www.collibra.com/

[8]     Databricks. (2023). Use MLflow and AutoML to make sure the data is right. https://databricks.com/

[9]     Devlin, J., et al. (2018). "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." arXiv:1810.04805.

[10]    Dhamdhere, K., and others (2017). CIDR's "Data Civiliser System."

[11]    Efron (1992). Bootstrap Techniques: A Reevaluation of the Jackknife. The Study of Statistics.

[12]    Falessi, D., et al. (2022). "How to Use LLMs to Clean Data Automatically." Journal of Data Intelligence.

[13]    The Cloud from Google. (2023). APIs for Vertex AI and getting rid of incorrect information. https://cloud.google.com/vertex-ai

[14]    Hellerstein, J. M., et al. (2017). "Data Wrangling: Techniques and Challenges." IEEE Data Eng. Bulletin.

[15]    In person. (2023). How to use Transformers. https://huggingface.co/docs

[16]    IBM. (2023). Preparing Data with Watson Studio. https://www.ibm.com/cloud/watson-studio

[17]    Jha, A., et al. (2023). "Trustworthy AI: Auditing and Debugging Large Language Models." ACM FAccT.

[18]    Kaggle. (2023). Problems in cleaning data.  https://www.kaggle.com/

[19]    Kambhatla, N., et al. (2022). IBM Research Report: "How to Use LLMs to Find and Profile Data."

[20]    Kim, H., et al. (2023). "Prompt Engineering for Cleaning Tabular Data." arXiv:2305.11809.

[21]    Kumar, A., et al. (2016). "How to Clean Data." PVLDB.

[22]    Lee, J., et al. (2023). AI Magazine: "Using Explainable AI to Clean Up Data."

[23]    LinkedIn. (2023). An examination of trends in data engineering. https://engineering.linkedin.com

[24]    Liu, P., et al. (2023). "Chain-of-Thought Prompting Elicits Reasoning." arXiv:2201.11903.

[25]    Lu, Y., et al. (2021). "AutoML for Data Cleaning." Proceedings of the AAAI.

[26]    Microsoft. (2023). Azure Synapse and Copilot. https://azure.microsoft.com

[27]    Microsoft's study. (2022). Tech Report: "How to Use Codex to Build Pipelines for Cleaning Data."

[28]    NIST. (2022). A strategy to deal with the risks that arise with AI. https://www.nist.gov/itl/ai-risk-management-framework.

[29]    OpenAI. (2023). Technical Report for GPT-4. https://openai.com/research/gpt-4

[30]    PandasAI. (2023). Pandas for AI in Workflows. https://github.com/gventuri/pandas-ai

[31]    Patel, J., et al. (2020). "Towards Explainable Data Cleaning." IEEE Big Data.

[32]    Pedregosa, F., et al. (2011). "Scikit-learn: Machine Learning in Python." JMLR.

[33]    PyData. (2023). PyJanitor: APIs that let you fix bad data. https://pyjanitor-devs.github.io/pyjanitor/

[34]    Raji, I. D., et al. (2020). "Bridging the AI Accountability Gap." ACM FAccT.

[35]    ReData. (2023). AI that checks the data's quality by itself. https://www.redata.team/

[36]    M. T. Ribeiro et al. (2016). "Why should I believe you?" "KDD."

[37]    Salesforce. (2022). Einstein GPT for how good the information is. https://www.salesforce.com

[38]    Satyanarayan, A., et al. (2021). "Tools for cleaning data that are easy to understand and use." CHI.

[39]    A piece of snow. In the year 2023. Using the Snowpark and GPT APIs to clean up data. https://www.snowflake.com

[40]    Suresh, H., & Guttag, J. V. (2019). "An Outline for Understanding the Unintended Consequences of ML." Communications of the ACM.

[41]    Talend. (2023). AI-powered technologies that prepare data. https://www.talend.com/

[42]    Trifacta. (2021). A white paper on how to handle data. https://www.trifacta.com/

[43]    The UCI Machine Learning Repository. (2023). Datasets that can be used for research on cleaning. https://archive.ics.uci.edu/

[44]    Varshney, K. R., et al. (2022). "Responsible AI in Data Cleaning Systems." IBM Journal of Research and Development.

[45]    Vartak, M., et al. (2016). "Towards Visualization-Aware Data Cleaning." SIGMOD.

[46]    Wang, H., et al. (2023). "Evaluating GPT-4 in Data Preparation Tasks" (arXiv:2303.07815).

[47]    The World Economic Forum. (2022). Rules for AI that is responsible.  https://www.weforum.org

[48]    Zhang, Y., et al. (2021). "An Examination of Data Cleaning Methodologies in Big Data." ACM Computing Surveys

[49]    Zhu, C., et al. (2020). "Can AI Clean Your Data?" Proceedings of VLDB.