

DATASEA - AN AUTOMATIC FRAMEWORK FOR COMPREHENSIVE DATASET PROCESSING USING LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

In the era of data-driven decision-making, efficiently acquiring and analyzing diverse datasets is critical for accelerating research and innovation. Yet, traditional manual approaches to dataset discovery, preparation, and exploration remain inefficient and cumbersome, especially as the scale and complexity of datasets continue to expand. These challenges create major roadblocks, slowing down the pace of progress and reducing the capacity for data-driven breakthroughs. To address these challenges, we introduce DataSEA (Search, Evaluate, Analyze), a fully automated system for comprehensive dataset processing, leveraging large language models (LLMs) to streamline the data handling pipeline. DataSEA autonomously searches for dataset sources, retrieves and organizes evaluation metadata, and generates custom scripts to load and analyze data based on user input. Users can provide just a dataset name, and DataSEA will handle the entire preparation process. While fully automated, minimal user interaction can further enhance system accuracy and dataset handling specificity. We evaluated DataSEA on datasets from distinct fields, demonstrating its robustness and efficiency in reducing the time and effort required for data preparation and exploration. By automating these foundational tasks, DataSEA empowers researchers to allocate more time to in-depth analysis and hypothesis generation, ultimately accelerating the pace of innovation. The code is available at <https://github.com/SingleView11/DataSEA>.

1 INTRODUCTION

In the age of artificial intelligence, the significance of data is undeniable. The volume of data generated and published online is rapidly increasing, yet searching for structured data on the internet remains challenging [Kacprzak et al. (2018)]. On one hand, quickly and clearly understanding the structure and content of large datasets—spanning social sciences, life sciences, high-energy physics, climate science, and other fields—has become increasingly difficult. On the other hand, as the scale and complexity of data grow, the efficiency of manually searching for and downloading datasets diminishes, turning into a daunting task. Traditional manual methods for discovering, preparing, and exploring data are becoming increasingly cumbersome and inefficient. Consequently, swiftly acquiring and analyzing diverse datasets has become a crucial factor in driving research and innovation.

With the advancement of prompt engineering, large language models (LLMs) have demonstrated impressive performance across various fields [Wei et al. (2022); Kojima et al. (2022); Wang et al. (2022); Zhou et al. (2022); Madaan et al. (2024); Bai et al. (2022); Chen et al. (2023)]. The strong capability of LLMs to process vast amounts of textual information opens new possibilities for the automated discovery, evaluation, and analysis of datasets. By utilizing LLMs to analyze the web information related to specific datasets, it is possible to organize this data into structured formats, facilitating further computational processing and user interaction.

Thus, we propose DataSEA (Search, Evaluate, Analyze), a comprehensive automated dataset processing system based on large language models. Additionally, by treating the collected data as a

new corpus for LLMs, it can enhance user interaction, further improving the system’s accuracy and adaptability to specific data processing needs. This allows users to gradually understand the various characteristics of target datasets and perform diverse evaluations in a question-answering system-like environment.

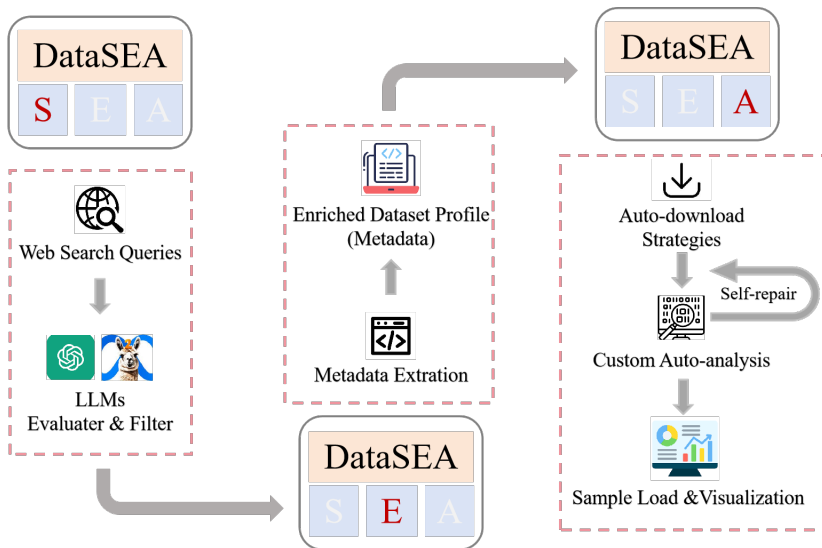


Figure 1: The architecture of DataSEA. The system is divided into three key modules: Search, Evaluate, and Analyze. The Search Module (S) retrieves top results from search engines and evaluates the relevance of the links using LLM models. The Evaluate Module (E) extracts metadata from the identified websites and retrieves research papers citing the dataset, followed by metadata optimization. The Analyze Module (A) generates and executes code to download datasets, hypothesizes possible download methods, and visualizes data samples. The entire process can be fully automated, though users may intervene to improve accuracy and filter unwanted downloads.

The design goal of DataSEA is to simplify and accelerate the data processing workflow, including the search, evaluation, and analysis of datasets. Users only need to provide a dataset name, and DataSEA can autonomously perform the search for dataset sources, retrieve and organize metadata, and generate custom scripts to load and analyze data based on user requirements.

To demonstrate the significant practical value of this system, we evaluated DataSEA on multiple datasets across various fields, such as computer vision, natural language processing, speech recognition, medicine, natural sciences, social sciences, and finance. Accessing these datasets is crucial for promoting the reproducibility of research findings, enabling scientists to build upon the work of others, and facilitating easier access to information and its sources for data journalists [Brickley et al. (2019)]. Experimental results indicate that DataSEA exhibits strong robustness and efficiency, significantly reducing the time and effort required for data preparation and exploration. Additionally, DataSEA provides rich and logical visualization, evaluation, and interpretation of data, greatly lowering the cost of understanding the content of datasets.

We also tested the performance of various LLMs as core processing components of DataSEA, including GPT-4 [Achiam et al. (2023)], GPT-3.5 [Ye et al. (2023)], and Llama [Touvron et al. (2023)]. Our findings indicate that the effectiveness of different LLM models in dataset collection, evaluation, and analysis is related to their parameter counts, which influences their ability to integrate and comprehend web information. This suggests that as model parameters increase in the future, alongside improvements in long text processing and logical reasoning capabilities, DataSEA is likely to exhibit enhanced capabilities in dataset collection and processing.

In Section 2, we outline related work; Section 3 elaborates on the methods used to construct DataSEA; Section 4 presents experiments demonstrating the effectiveness of DataSEA; Section 5 discusses the existing limitations; and finally, Section 6 provides a comprehensive summary of this paper. Additional relevant data can be found in the appendix.

2 RELATED WORK

Efficient Dataset Discovery and Exploration Efficient dataset discovery, preparation, and exploration are critical components of the data-driven research pipeline. However, traditional approaches often require significant manual effort, involving labor-intensive tasks such as searching for relevant datasets, preparing them for analysis, and creating appropriate exploratory tools. Recent advances in automated data management have sought to alleviate these challenges by streamlining data discovery, cleaning, transformation, and exploration. Automated data integration tools such as advanced dataset search systems Brickley et al. (2019) and AutoML frameworks Zöller & Huber (2021) have demonstrated considerable promise in minimizing manual workload, though many solutions still face limitations regarding the flexibility and comprehensiveness of the automated workflow. These advancements reflect a growing interest in reducing human intervention and enhancing efficiency through intelligent data management solutions.

Automated Dataset Processing Frameworks Several automated dataset processing tools have emerged to address various parts of the data handling workflow. Tools such as Trifacta and OpenRefine Petrova-Antonova & Tancheva (2020) focus on data wrangling, emphasizing interactive user experiences for cleaning and transforming data. Although these tools significantly improve the efficiency of data preprocessing, they require extensive user involvement throughout the process and lack fully automated workflows, particularly in terms of dataset discovery and evaluation.

The development of systems like AutoML He et al. (2021) has further paved the way for automation by addressing tasks like feature engineering and model selection. However, while AutoML tools effectively handle model training and hyperparameter tuning, they often depend on structured, pre-prepared datasets. The processes of discovering datasets and assessing their suitability for analysis largely remain manual, limiting the overall automation potential in the data science pipeline Biswas et al. (2022).

Leveraging Language Models for Automation LLM models such as GPT-4 Achiam et al. (2023) and LLAMA Touvron et al. (2023), have demonstrated significant capabilities in understanding natural language, generating code, and automating workflows in complex domains. Previous research has leveraged LLMs to generate scripts for data processing Biswas & Talukdar (2024); Nejjar et al. (2023); Patiny & Godin (2023), streamlining the creation of custom data handling scripts. These efforts highlight the potential of LLMs in automating repetitive tasks, but they often focus narrowly on code generation without addressing the end-to-end dataset processing pipeline.

The recent work on LLM-based assistants (e.g., GPT-4, LLAMA) has further demonstrated the applicability of these models for responding to natural language queries related to data analytics Ram et al. (2024), offering on-demand support for exploratory data analysis (EDA) Ma et al. (2023) and visualization Sah et al. (2024). However, these applications are reactive, requiring substantial user intervention in specifying datasets, parameters, and contexts for each step.

Comprehensive Dataset Automation DataSEA builds on these advancements, aiming to deliver a fully automated framework for dataset processing that encompasses not only code generation but also dataset discovery and evaluation. Unlike existing semi-automated tools that require significant human interaction, DataSEA autonomously manages dataset discovery, metadata extraction, and script generation, reducing the need for user input to a minimum. Data-centric AI Zha et al. (2023) suggests that focusing on automating data-handling processes can significantly accelerate research outcomes Mittal et al. (2023), and DataSEA aligns closely with this vision by implementing an automated pipeline that integrates search, evaluation, and analysis.

By leveraging LLMs to automate not only the preparation but also the discovery and evaluation of datasets, DataSEA contrasts with existing solutions that focus predominantly on either preparation or analytics. This comprehensive approach empowers researchers to reduce the time spent on foundational tasks, allowing for more in-depth analysis and exploration of the data.

Summary While previous work has made significant strides in automating parts of the data preparation and analysis workflow, DataSEA is among the first to provide a fully integrated solution for dataset discovery, evaluation, and custom analysis using large language models. By autonomously handling these key stages, DataSEA extends beyond the capabilities of current LLM-driven tools and represents a significant step toward automating the entire data lifecycle. This approach aligns

with recent trends in AI-driven automation and data-centric methodologies, ultimately accelerating the pace of innovation in data-driven research.

3 METHODOLOGY

3.1 SYSTEM OVERVIEW

DataSEA is composed of three core modules: Search, Evaluate, and Analyze. The system leverages large language models to intelligently locate dataset sources, extract metadata, and generate custom scripts for loading and visualizing the data. Users can input a dataset name and optional description, and DataSEA autonomously handles the remainder of the process. The architecture allows for minimal user interaction, but additional input can improve the system’s accuracy.

To enhance the effectiveness of LLMs, DataSEA employs instruction-prompting [Brown (2020)] and a multi-chunk strategy [Liu et al. (2024)] to handle long inputs, ensuring that even large datasets can be processed effectively by breaking the data into manageable sections while maintaining context across chunks.

3.2 SEARCH MODULE

The Search Module in DataSEA automates the process of discovering dataset websites by leveraging search engines and LLMs to filter and rank relevant results. The user starts by inputting a dataset name and, optionally, additional dataset details to refine the search. Drawing inspiration from tools like Google Dataset Search [Brickley et al. (2019)], the system generates optimized search queries based on the input and send them to search engines such as Google to retrieve the top-ranking links.

Once the top links are retrieved, the system performs web content extraction on each page. The contents are then analyzed by the LLM, which generates evaluation info. Similar to work on using LLMs for content understanding and retrieval tasks [Brown (2020)], the LLM helps filter out irrelevant or low-quality pages. The links are ranked based on their evaluation info, with the top-ranking results being those most likely to contain useful dataset information. More detail can be found in Appendix.

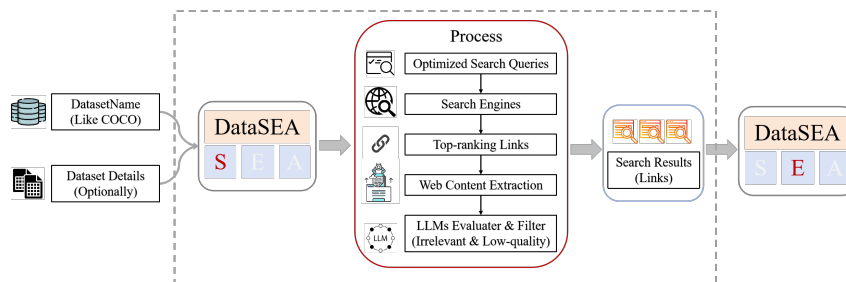


Figure 2: The process flow of the Search Module in DataSEA. The system retrieves top links from search engines based on user input, evaluates the relevance of each link using the LLM model, and filters out irrelevant or low-quality pages. The links are ranked by relevance to the dataset, allowing the user to quickly access accurate and useful resources.

3.3 EVALUATION MODULE

The Evaluation Module in DataSEA generates the metadata of the dataset, including various information with 3 steps: Metadata Extraction, Reference Paper Retrieval, Metadata Extension. User can customize the properties of the metadata template.

3.3.1 METADATA EXTRACTION

In the first step, the system processes the links identified as relevant in the Search Module. The system extracts their website content and uses the LLM to get metadata. This process is guided by a

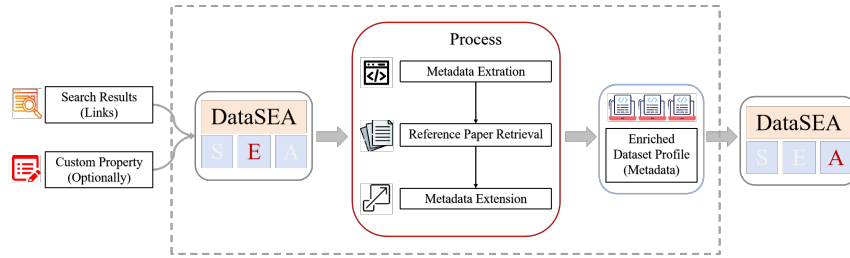


Figure 3: The process flow of the Evaluation Module in DataSEA. The system extracts metadata from relevant links, retrieves research papers that reference the dataset, and validates the relevance of each paper using LLM models. The metadata is then optimized and extended using information from the papers, ensuring a comprehensive dataset profile.

preset of metadata attributes, including the dataset’s usage, content and scale, application fields, and other important factors. Additionally, users have the flexibility to input custom property names, and the system dynamically optimizes prompts for the LLM to retrieve those specific properties. The results from different links are combined at the end for optimization.

3.3.2 REFERENCE PAPER RETRIEVAL

The second stage focuses on retrieving and ranking research papers that reference the dataset. The system uses the Google Scholar API to find papers that may have cited the dataset. These papers are ranked by citation count, following the widely accepted practice of using citation metrics as indicators of a paper’s impact (Bornmann & Daniel, 2008). For each paper, the metadata extracted in the first step is used to validate whether the paper indeed references the correct dataset, as there may be cases of duplicate names or other inaccuracies. The system filters the most impactful papers, and the user can specify the number of papers to be collected.

3.3.3 METADATA EXTENSION

In the final step, the system extracts additional metadata from the validated reference papers. This may include more detailed descriptions of the dataset’s features, specific application examples, and additional context provided by the authors. The extracted metadata is then combined with the original information from the dataset’s website, resulting in an enriched and comprehensive dataset profile.

3.4 ANALYSIS MODULE

The Analysis Module in DataSEA is used to download the dataset and analyze it by generating code and test them. The final generated code can load and visualize dataset samples, and user can input a customized requirement and the code generation will be adapted to satisfy the requirement.

3.4.1 DATASET DOWNLOAD

In this phase, the module utilizes metadata from the Evaluation Module to generate hypotheses for downloading the dataset from the identified websites. For each website containing dataset information, the system generates hypotheses regarding possible download methods. It then creates code by combining the hypothesis, the website content, and the dataset metadata, and executes this code to attempt to download the dataset.

The generation of hypotheses has proven to be an effective method for handling uncertainty in data retrieval processes, as it allows the system to explore multiple potential download strategies simultaneously. This approach is inspired by the inductive reasoning capabilities of language models, as demonstrated in Hypothesis Search (Wang et al., 2023), where generating and testing multiple hypotheses leads to more robust and successful outcomes.

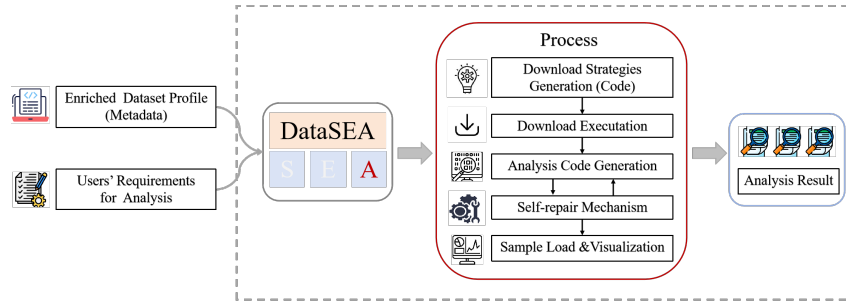


Figure 4: The process flow of the Analyze Module in DataSEA. The system generates hypotheses for dataset download based on metadata, executes the download, and analyzes the dataset by reading the samples of the raw data. It then generates visualization code, with options for manual intervention to refine the analysis. The system also incorporates a self-repair feature to handle any issues in code execution.

Since multiple links and hypotheses may be involved, the system organizes the download results into separate folders, each corresponding to a specific hypothesis. Users can manually inspect the results and delete unwanted downloads, such as cases where multiple datasets (e.g., raw, processed) are available and only a portion is needed. This flexibility allows the user to refine the dataset collection process and accelerate further analysis.

3.4.2 DATASET ANALYSIS AND VISUALIZATION

After successfully downloading the datasets, the system generates custom analysis code by calling the LLM with prompts that include the dataset metadata and initial portions of the downloaded data. This analysis code is designed to load the dataset, read the samples of the raw dataset, and visualize key aspects of the dataset’s structure and contents.

The analysis code generation process incorporates a self-repair mechanism, drawing from approaches like CodeT5 (Wang et al., 2021). If the generated code fails during execution, the system automatically collects the error log and combines it with the original code to form a more detailed context. This context is sent back to the LLM, which attempts to identify and fix the issues in the code. The system iterates through this feedback loop until a working version of the code is produced, significantly improving the reliability and robustness of the analysis code. This feature allows the system to autonomously handle failures and continuously improve the generated code without requiring user intervention.

In addition to the automated processes, users have the ability to write customized requirements. The system will generate and test code based on the user’s input, allowing for tailored analysis that fits specific research needs. This user interaction complements the fully automated pipeline, providing flexibility for users to guide the analysis toward more specific goals if needed.

4 EXPERIMENTS

4.1 SETUP

Datasets We evaluated DataSEA on 100 datasets across various fields to assess its generalizability and effectiveness in automating dataset discovery, evaluation, and analysis. The datasets were sourced from repositories such as Google Dataset Search, Kaggle, and other publicly available platforms. These datasets span a wide range of fields, including Computer Vision, Natural Language Processing (NLP), Healthcare, Speech and Audio, Natural Sciences, Social Science, Finance, Transportation, Recommendation Systems, Time Series Analysis, Robotics, and Agriculture. The diversity in size, format, and complexity of these datasets allowed for a comprehensive evaluation of DataSEA’s performance across different domains.

Models We used three models to evaluate the performance of DataSEA: gpt-4o, gpt-4o-mini, and llama3. For gpt-4o and gpt-4o-mini, we directly call openai apis; for llama3, we deploy it locally.

Parameters We configured three different modes in DataSEA to explore the trade-offs between speed and accuracy: High-Speed, Medium-Speed and Slow-Speed version. High-Speed Version is optimized for fast dataset discovery and analysis by reducing the number in hyper parameters. Medium-Speed Version is a balance between performance and speed. Slow-Speed Version is focused on thorough dataset discovery and analysis. The 3 modes take about 3-5 / 10-15 / 20-60 minutes.

The hyper parameters include the number of websites collected per dataset, the number of hypotheses generated, the number of relevant papers retrieved, the number of download code generation trial, the number of analysis code idea and the number of self-repairs performed when issues were encountered in code generation.

Evaluation Metrics We evaluated DataSEA based on the performance of each of its three core modules: Search (S), Evaluate (E), and Analyze (A). For every module we have different metrics, with more detail in the next subsection.

4.2 MAIN RESULTS

We present the evaluation of each module—Search (S), Evaluate (E), and Analyze (A)—using three model and their high, medium, and slow-speed versions. Each module’s performance is detailed in its respective section.

Search Module The Search Module was evaluated based on its ability to find relevant websites. The main evaluation metrics are the **RWF** (Relevant Websites Found in LLM return, true if one is found otherwise false), **ACC** (Relevance judgement accuracy of the LLM).

Model	Version	RWF (%)	ACC (%)
gpt-4o	High Speed	100	92.36
	Medium Speed	100	93.02
	Slow Speed	100	96.33
gpt-4o-mini	High Speed	97	82.77
	Medium Speed	98	84.14
	Slow Speed	98	90.13
llama3	High Speed	99	83.20
	Medium Speed	100	81.82
	Slow Speed	100	87.67

Table 1: Results for Search Module (S) across different models and versions.

The results show that gpt-4o under low-speed mode achieves the highest accuracy for relevant websites judging. As for the false judgements, most cases are false negative - the website does contain information about the dataset but it is judged as not because there are too many redundant information in the html content, so that the information about the dataset is stuck in the middle and not captured well by the llm. This result is just like the phenomenon described in a paper showing LLM’s incapability in processing long context[Liu et al. (2024)].

Evaluation Module The Evaluation Module was assessed using the quality of generated metadata and retrieved papers. The **I-ACC** (Initial Metadata Accuracy) reflects the system’s ability to extract correct metadata across different properties. It is calculated by averaging the accuracy of the 8 different properties in the metadata. We also evaluated **R-ACC** (Relevant Papers Accuracy, if a judged reference paper is really referring to the dataset) and the **E-ACC** (Extended Metadata Accuracy).

Analysis Module For the Analysis Module, we focused on the **DDS** (Dataset Download Success), **H-ACC** (Hypothesis Accuracy), **CAS** (Code Analysis Success Without Intervention), and **CAS-I** (Code Analysis Success With Intervention). The H-ACC is examined by manually following the steps generated by a hypothesis, and is marked as true if there exist one true hypothesis.

The intervention in code analysis means downloading the dataset by following the hypothesis, because while the system can be fully-automated, it is hard to download datasets automatically as most

Model	Version	I-ACC (%)	R-ACC (%)	E-ACC (%)
gpt-4o	High Speed	92.63	87.31	98.00
	Medium Speed	93.13	84.09	98.00
	Slow Speed	94.75	85.12	98.25
gpt-4o-mini	High Speed	87.13	80.94	93.50
	Medium Speed	85.00	82.30	92.25
	Slow Speed	90.38	78.64	93.50
llama3	High Speed	82.63	81.53	84.75
	Medium Speed	83.00	81.20	84.13
	Slow Speed	83.00	84.22	85.88

Table 2: Results for Evaluation Module (E) across different models and versions.

datasets’ host platform will require login or email request for the dataset, and can not be crawled easily. Even if the datasets are available publicly, there may be multiple datasets with different properties(raw, processed, etc), and serve different purposes, so a comparison for code analysis between automatically downloaded data and manually downloaded data following the hypothesis is needed.

Model	Version	DDS (%)	H-ACC (%)	CAS (%)	CAS-I (%)
gpt-4o	High Speed	9	81	11	32
	Medium Speed	11	88	12	35
	Slow Speed	12	91	15	38
gpt-4o-mini	High Speed	4	76	9	28
	Medium Speed	6	82	9	32
	Slow Speed	6	87	9	31
llama3	High Speed	2	62	3	19
	Medium Speed	3	69	3	21
	Slow Speed	4	70	3	21

Table 3: Results for Analysis Module (A) across different models and versions.

The interesting finding is that sometimes the CAS is higher than H-ACC, which is counter-intuition because it is hard to imagine analyzing a dataset when it is not downloaded. This is due to the prompt design, as instruction to try to generate code for the dataset without checking the dataset info. As a result, for some popular datasets like MNIST or CIFAR-10, even though the download dataset folder is empty, the generated code can still be run and will successfully generate visualization of samples.

The CAS-I will be lifted greatly if user manually follow the ideas generated and download the dataset. For example, downloading a dataset in Kaggle is convenient and only need a click of button if user is logged in, but the system currently cannot auto-login for the user and will fail to download the dataset.

5 LIMITATIONS

DataSEA faces several challenges, including its inability to process databases and databanks, which limits its application in biological fields like genomics [Sherry et al. (2001)] and proteomics [Abola et al. (1984)]. Additionally, its performance depends heavily on LLMs, and the system exhibits a trade-off between speed and accuracy. While the self-repair mechanism can handle common errors, complex dataset structures may still require manual intervention. The automatic dataset download process also struggles with anti-crawling mechanisms and login/email requests, and visual information is often lost during HTML content extraction, suggesting the need for methods that integrate neural optical understanding [Blecher et al. (2023)].

6 CONCLUSION

In this work, we introduced DataSEA, a fully automated system for comprehensive dataset processing, which integrates dataset search, evaluation, and analysis using large language models. Our

432 system allows users to input a dataset name and automatically retrieves, evaluates, and analyzes
433 datasets from a wide range of domains. DataSEA demonstrates the effectiveness of leveraging
434 LLMs to streamline the dataset processing pipeline, reducing manual effort and enabling researchers
435 to focus on deeper data analysis. While our system shows promising results across diverse datasets,
436 certain limitations such as handling databases and databanks and challenges in automatic downloads
437 present opportunities for future work. Overall, DataSEA represents a significant step forward in au-
438 tomating the early stages of dataset preparation, offering researchers a powerful tool to accelerate
439 data-driven discoveries.

440 REFERENCES

- 441 Enrique E Abola, Frances C Bernstein, and Thomas F Koetzle. *The protein data bank*. Springer,
442 1984.
- 443 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
444 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
445 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 446 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
447 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harm-
448 lessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 449 Anjanava Biswas and Wrick Talukdar. Robustness of structured data extraction from in-plane rotated
450 documents using multi-modal large language models (llm). *Journal of Artificial Intelligence*
451 *Research*, 2024.
- 452 Sumon Biswas, Mohammad Wardat, and Hriday Rajan. The art and practice of data science
453 pipelines: A comprehensive study of data science pipelines in theory, in-the-small, and in-the-
454 large. In *Proceedings of the 44th International Conference on Software Engineering*, pp. 2091–
455 2103, 2022.
- 456 Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. Nougat: Neural optical
457 understanding for academic documents. *arXiv preprint arXiv:2308.13418*, 2023.
- 458 Lutz Bornmann and Hans-Dieter Daniel. What do citation counts measure? a review of studies on
459 citing behavior. *Journal of documentation*, 64(1):45–80, 2008.
- 460 Dan Brickley, Matthew Burgess, and Natasha Noy. Google dataset search: Building a search engine
461 for datasets in an open web ecosystem. In *The world wide web conference*, pp. 1365–1375, 2019.
- 462 Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 463 Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models
464 to self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
- 465 Xin He, Kaiyong Zhao, and Xiaowen Chu. Automl: A survey of the state-of-the-art. *Knowledge-*
466 *based systems*, 212:106622, 2021.
- 467 Emilia Kacprzak, Laura Koesten, Jeni Tennison, and Elena Simperl. Characterising dataset search
468 queries. In *Companion Proceedings of the The Web Conference 2018*, pp. 1485–1488, 2018.
- 469 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
470 language models are zero-shot reasoners. *Advances in neural information processing systems*,
471 35:22199–22213, 2022.
- 472 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and
473 Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the*
474 *Association for Computational Linguistics*, 12:157–173, 2024.
- 475 Pingchuan Ma, Rui Ding, Shuai Wang, Shi Han, and Dongmei Zhang. Insightpilot: An llm-
476 empowered automated data exploration system. In *Proceedings of the 2023 Conference on Em-*
477 *pirical Methods in Natural Language Processing: System Demonstrations*, pp. 346–352, 2023.

- 486 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
487 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
488 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 489
490 Deepti Mittal, Rebecca Mease, Thomas Kuner, Herta Flor, Rohini Kuner, and Jamila Andoh. Data
491 management strategy for a collaborative research center. *GigaScience*, 12:giad049, 2023.
- 492 Mohamed Nejjar, Luca Zacharias, Fabian Stiehle, and Ingo Weber. Llms for science: Usage for
493 code generation and data analysis. *Journal of Software: Evolution and Process*, pp. e2723, 2023.
- 494
495 Luc Patiny and Guillaume Godin. Automatic extraction of fair data from publications using llm.
496 2023.
- 497 Dessislava Petrova-Antonova and Rumyana Tancheva. Data cleaning: A case study with openrefine
498 and trifacta wrangler. In *Quality of Information and Communications Technology: 13th Interna-*
499 *tional Conference, QUATIC 2020, Faro, Portugal, September 9–11, 2020, Proceedings 13*, pp.
500 32–40. Springer, 2020.
- 501 Gummuluri Venkata Ravi Ram, Kesanam Ashinee, and M Anand Kumar. End-to-end space-efficient
502 pipeline for natural language query based spacecraft health data analytics using large language
503 model (llm). In *2024 5th International Conference on Innovative Trends in Information Technol-*
504 *ogy (ICITIIT)*, pp. 1–6. IEEE, 2024.
- 505 Subham Sah, Rishab Mitra, Arpit Narechania, Alex Endert, John Stasko, and Wenwen Dou. Gen-
506 erating analytic specifications for data visualization from natural language queries using large
507 language models. *arXiv preprint arXiv:2408.13391*, 2024.
- 508
509 Stephen T Sherry, M-H Ward, M Kholodov, J Baker, Lon Phan, Elizabeth M Smigielski, and Karl
510 Sirotkin. dbsnp: the ncbi database of genetic variation. *Nucleic acids research*, 29(1):308–311,
511 2001.
- 512 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
513 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
514 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 515 Ruo Cheng Wang, Eric Zelikman, Gabriel Poesia, Yewen Pu, Nick Haber, and Noah D Goodman.
516 Hypothesis search: Inductive reasoning with language models. *arXiv preprint arXiv:2309.05660*,
517 2023.
- 518
519 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
520 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
521 *arXiv preprint arXiv:2203.11171*, 2022.
- 522 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified
523 pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint*
524 *arXiv:2109.00859*, 2021.
- 525 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
526 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
527 *neural information processing systems*, 35:24824–24837, 2022.
- 528
529 Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou,
530 Chao Gong, Yang Shen, et al. A comprehensive capability analysis of gpt-3 and gpt-3.5 series
531 models. *arXiv preprint arXiv:2303.10420*, 2023.
- 532 Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, and Xia Hu. Data-centric ai: Perspec-
533 tives and challenges. In *Proceedings of the 2023 SIAM International Conference on Data Mining*
534 *(SDM)*, pp. 945–948. SIAM, 2023.
- 535 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuur-
536 mans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex
537 reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.
- 538
539 Marc-André Zöllner and Marco F Huber. Benchmark and survey of automated machine learning
frameworks. *Journal of artificial intelligence research*, 70:409–472, 2021.

540	A APPENDIX	
541		
542	MENU	
543		
544	CONTENTS	
545		
546		
547	1 Introduction	1
548		
549	2 Related Work	3
550		
551		
552	3 Methodology	4
553	3.1 System Overview	4
554		
555	3.2 Search Module	4
556		
557	3.3 Evaluation Module	4
558	3.3.1 Metadata Extraction	4
559	3.3.2 Reference Paper Retrieval	5
560	3.3.3 Metadata Extension	5
561		
562	3.4 Analysis Module	5
563	3.4.1 Dataset Download	5
564	3.4.2 Dataset Analysis and Visualization	6
565		
566		
567	4 Experiments	6
568	4.1 Setup	6
569		
570	4.2 Main Results	7
571		
572	5 Limitations	8
573		
574	6 Conclusion	8
575		
576		
577	A Appendix	11
578		
579	Appendices	11
580		
581	A.1 Search Pipeline	13
582		
583	A.2 Evaluate Pipeline	13
584		
585	A.3 Analyze Pipeline	13
586		
587	A.4 Flowchart Visualizations	14
588		
589	A.5 Prompts we use	14
590	A.5.1 Dataset Website Prompt	15
591	A.5.2 Dataset Paper Retrieval	16
592	A.5.3 Retrieving PDF Links for Dataset Paper	17
593	A.5.4 Generating Dataset Metadata Extraction Instructions	18
	A.5.5 Generating Dataset Reference Detection Instructions	19

594		
595		
596		
597		
598		
599		
600	A.6	Google Search Api 25
601	A.7	Long Context Inference 27
602		
603	A.8	Reproducibility 31
604	A.9	Code Structure 31
605		
606	A.10	Code Explanation 32
607		
608		
609		
610		
611		
612		
613		
614		
615		
616		
617		
618		
619		
620		
621		
622		
623		
624		
625		
626		
627		
628		
629		
630		
631		
632		
633		
634		
635		
636		
637		
638		
639		
640		
641		
642		
643		
644		
645		
646		
647		

This section provides a detailed description of the code structure for the DataSEA system, which is composed of three main parts: Search, Evaluate, and Analyze. Each part contains an integrated pipeline to automate the dataset processing workflow.

A.1 SEARCH PIPELINE

The **Search Pipeline** is responsible for identifying and retrieving relevant datasets based on the user-provided input. It utilizes large language models (LLMs) to search dataset repositories and official websites. The key steps are as follows:

1. **Dataset Query:** The system accepts a dataset name as input and sends it to the LLM for generating search queries.
2. **Web Search:** These queries are used to search for datasets across multiple sources, including Kaggle, UCI, Zenodo, and official dataset websites.
3. **Link Retrieval:** The system collects potential dataset links, filtering and ranking them based on relevance and credibility.
4. **Search Output:** The final output is a set of ranked dataset links which are passed to the next pipeline for evaluation.

A.2 EVALUATE PIPELINE

The **Evaluate Pipeline** processes the dataset links retrieved from the search phase to extract useful metadata and identify the most reliable sources. This pipeline consists of three parts:

1. **Metadata Extraction:** Using LLMs, the system extracts relevant metadata from the dataset links, such as dataset size, format, domain, and source information. This step utilizes a combination of preset and user-specified property names.
2. **Reference Paper Retrieval:** The system retrieves research papers that reference the dataset by querying academic databases. The papers are ranked by citation count, and their metadata is validated.
3. **Metadata Extension:** Metadata from the papers is integrated with the original dataset metadata to provide a more comprehensive evaluation. The system uses cross-references to ensure accuracy and consistency.

The output of this pipeline includes the final metadata and a list of reference papers, which are passed to the Analyze Pipeline.

A.3 ANALYZE PIPELINE

The **Analyze Pipeline** is responsible for downloading, organizing, and analyzing the dataset. It performs several critical tasks:

1. **Dataset Download:** The system generates hypotheses for downloading the dataset based on metadata and content from the dataset website. For each hypothesis, it attempts to download the dataset and stores the results in corresponding folders.
2. **Code Generation:** After downloading, the system uses LLMs to generate Python code that can load, read, and visualize the dataset. The code is based on the metadata and sample points extracted from the dataset.
3. **Self-Repair:** If the generated code fails during execution, the system captures the error log and re-invokes the LLM to attempt self-repair. This step improves the robustness of the generated code.
4. **User-Defined Requirements:** Users can specify custom analysis requirements, and the system generates corresponding code to meet these needs. The code is automatically tested, and any failures are corrected using the self-repair mechanism.

The final output of the Analyze Pipeline includes a fully downloaded dataset, analysis code, and visualized sample data points. Users can interact with the system to modify or delete failed downloads as needed.

A.4 FLOWCHART VISUALIZATIONS

The following figures illustrate the flow of the DataSEA system, both at a high level and within each individual module:

- **Figure 5:** A flowchart of the entire pipeline, from dataset input to analysis completion.
- **Figure 6:** A detailed flowchart of the Search Pipeline, showing the steps involved in dataset query and retrieval.
- **Figure 7:** A flowchart of the Evaluate Pipeline, illustrating the metadata extraction, reference paper retrieval, and metadata extension processes.
- **Figure 8:** A flowchart of the Analyze Pipeline, depicting the dataset download, code generation, and self-repair mechanisms.

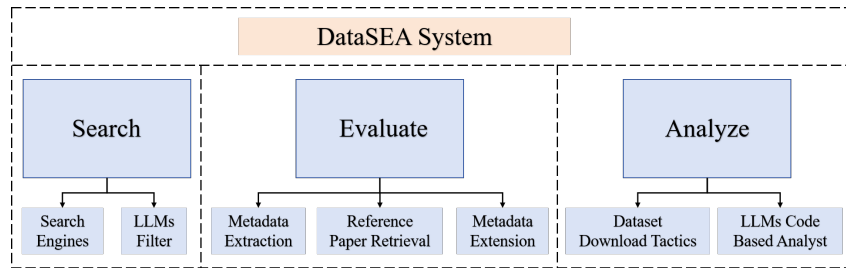


Figure 5: Flowchart of the entire DataSEA pipeline.

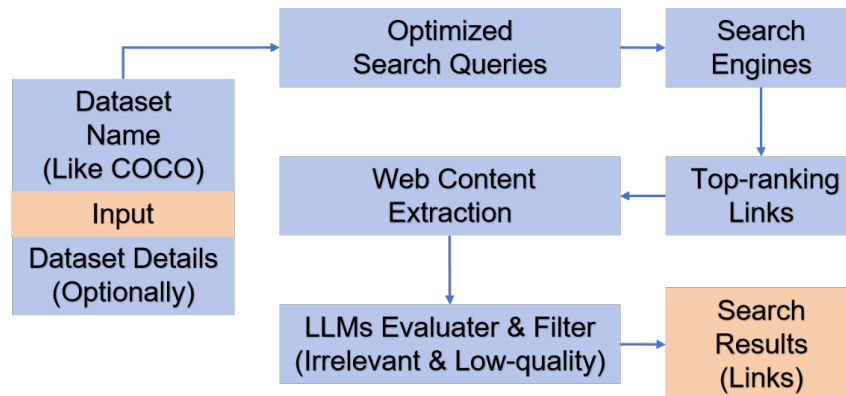


Figure 6: Flowchart of the Search Pipeline.

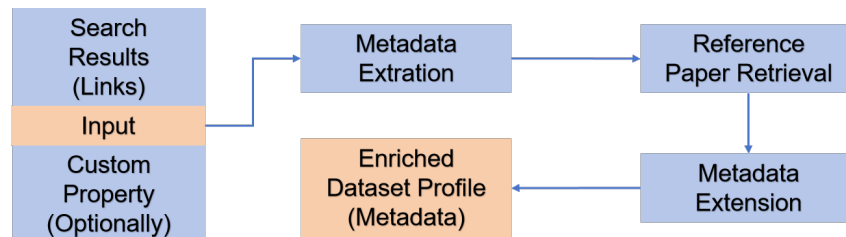


Figure 7: Flowchart of the Evaluate Pipeline.

A.5 PROMPTS WE USE

We make tons of LLM api calls during the SEA process, and for every specific task we have a independent prompt. We will list them and show their usage.

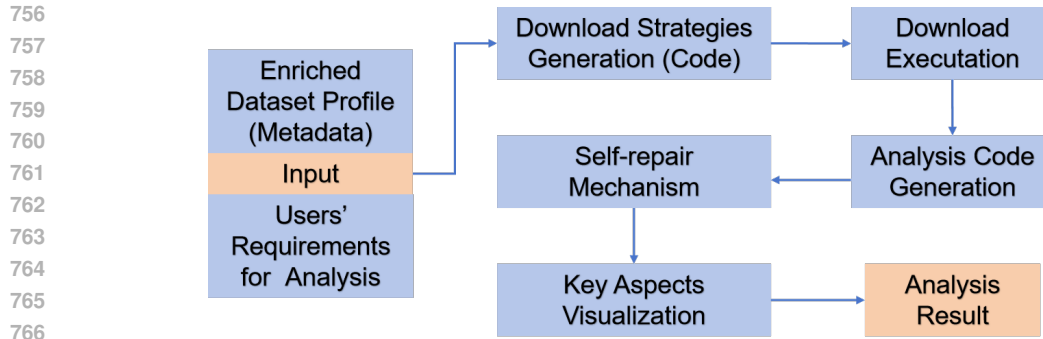


Figure 8: Flowchart of the Analyze Pipeline.

770 A.5.1 DATASET WEBSITE PROMPT

771
772 The following prompt is designed to analyze a given website to determine if it contains a valid dataset
773 download link. It is used in situations where the dataset might be located on various websites, and
774 we need to extract the dataset link and metadata from the HTML content of the page. The prompt
775 ensures the clarity of the task and requests the output in a structured JSON format. Below is a
776 breakdown of the prompt:

- 777 • The prompt instructs the system to analyze the HTML of the website to determine if it is
778 the official website for the dataset in question, identified by the `dataset_name` and a
779 description.
- 780 • It asks the system to check specifically for a dataset download link and warns against mis-
781 taking an article download link for a dataset.
- 782 • If the website is identified as containing the dataset, the system should extract and return
783 the dataset download link along with basic metadata, such as the website description.
- 784 • The prompt emphasizes returning the result in JSON format, specifying fields such as:
785
 - 786 – `is_dataset_website`: Whether the website is related to the dataset.
 - 787 – `download_link_dataset`: The URL link to the dataset download.
 - 788 – `metadata`: Any additional information extracted from the website, such as descrip-
789 tions or other relevant data.
- 790 • The prompt also requests that if the dataset link is not found, the system should provide a
791 reason for this.
- 792 • It concludes by asking for a structured JSON output without any unnecessary text or for-
793 mating, to ensure compatibility and ease of use for further analysis.

794
795 The structure of the prompt is as follows:

796
797 Determine whether the current website HTML is the website for the dataset "{dataset_
798 !!You should notice that the download link is for dataset and not article! If there
799 If it is, give the dataset download link from the HTML content and provide some meta
800
801 If the download link is already the dataset, then note it. Otherwise, indicate that
802
803 If it is not, provide the reason.
804
805 Return the format in JSON with the following structure:
806 {
807 "is_dataset_website": <boolean>,
808 "metadata": <object>,
809 "download_link_dataset_exists": <boolean>,
810 "download_link_dataset": <string>,
811 "is_direct_data": <boolean>,
812 }

```

810     "reason": <string>
811 }
812

```

Note: just give the json, and do not add any extra words like adding the j-s-o-n let

```

813 The website HTML:
814 ""
815 {html_content}
816 ""
817
818
819

```

820 A.5.2 DATASET PAPER RETRIEVAL

821
822 The prompt generated by this function is used to identify whether a provided website contains the
823 original paper for a given dataset. The original paper refers to a publication where the dataset was
824 first introduced by the author, not merely a paper that uses the dataset. The prompt asks the LLM to
825 carefully evaluate the HTML content of the page and determine if the paper link is present.

826 The model must also ensure that the link leads to a downloadable paper file (such as PDF) and not
827 another webpage or non-relevant content. If the website does not contain the original paper, the
828 model is required to provide a reason. The output of the task should follow a predefined JSON
829 structure that includes flags for whether the paper link is available, and metadata about the website
830 and the paper.

```

831
832 def generate_prompt_paper(link, dataset_name, desc = ""):
833     html_content = fetch_html_from_link(link)
834
835     if html_content is None:
836         return ""
837
838     prompt = f"""
839 Determine whether the current website HTML is the website for the original paper
840
841 Note that for "the original paper of the dataset", it means that the author of the
842 It does not mean the author just use the dataset in his research, but means that
843
844 If it is, give the paper download link from the HTML content and provide some met
845
846 If the download link is just the paper pdf(or possibly other format) then note it
847
848 If it is not the website of the original paper, provide the reason.
849
850 Return the format in JSON with the following structure:
851 {
852     "is_dataset_paper_website": <boolean>,
853     "metadata": <object>,
854     "download_link_paper_exists": <boolean>,
855     "download_link_paper": <string>,
856     "is_direct_paper": <boolean>,
857     "reason": <string>
858 }
859
860 Note: just give the json, and do not add any extra words like adding the j-s-o-n
861
862 The website HTML:
863 ""
864 {html_content}
865 ""
866
867 return prompt

```


A.5.3 RETRIEVING PDF LINKS FOR DATASET PAPER

The purpose of this prompt is to extract direct download links for the original dataset paper in common formats such as PDF, DOCX, or TXT from a webpage. Given the website's HTML content and dataset details, the LLM is tasked with identifying direct download links for the academic paper, filtering out irrelevant content like datasets or other material. The prompt also specifies the format for the response, which must be structured in JSON.

The returned JSON should contain download links and specify the file format, ensuring that the links are valid and directly lead to paper files rather than web pages or unrelated content.

```

def get_potential_pdf_link(link, dataset_name, desc = ""):
    html_content = fetch_html_from_link(link)

    prompt = f"""
    I have the HTML content of a website, and I need to find any direct download links.

    The description of the paper is as follows:
    "{desc}".

    Based on this information, please search through the HTML content to find any direct download links.

    Remember you should give the direct link of paper but not other weird stuff like

    Return the format in JSON with the following structure:
    {{
      download_link_1: {{
        "link": "https://aaa.com",
        "format": "pdf"
      }},
      download_link_2: {{
        "link": "https://bbb.com",
        "format": "txt"
      }},
      ...,
      download_link_n: {{
        "link": "https://nnn.com",
        "format": "other format"
      }}
    }}

    Note: just give the json, and do not add any extra words like adding the j-s-o-n

    The website HTML:
    \"\"\"
    {html_content}
    \"\"\"

    """

    prompt = clamp_prompt(prompt)
    res = LLMapi(prompt)
    return res

```

918 A.5.4 GENERATING DATASET METADATA EXTRACTION INSTRUCTIONS

919 This prompt is designed to instruct the LLM to extract relevant metadata from a collection of concatenated text files that contain information about a dataset. The LLM is provided with basic information about the dataset, such as its name and current metadata, and is tasked with extracting additional details like description, size, scale, author, and other relevant properties. The output format is strictly defined as JSON, and the LLM is asked not to provide any explanations, only the JSON data.

926 The key fields that the LLM is expected to populate include:

- 928 • **description:** A brief description of the dataset.
- 929 • **size:** The size of the dataset (e.g., 1GB, 10,000 samples).
- 930 • **scale:** The memory size of the dataset (e.g., 1TB, 100MB).
- 931 • **author:** The dataset’s creator or author.
- 932 • **organization:** The institution responsible for the dataset.
- 933 • **usage:** Common uses for the dataset (e.g., model training, validation).
- 934 • **application fields:** Relevant application domains such as computer vision or NLP.
- 935 • **keywords:** Key terms associated with the dataset.

939 The LLM is instructed to use the text files as a source and output the final information in the required JSON structure.

```

942 # Function to create the instruction prompt without the actual text
943 def generate_instruction_prompt():
944     dataset_name, dataset_info = read_metadata()
945     prompt = f"""
946     You are provided with a detailed description from a folder of concatenated text files.
947     Your task is to extract the relevant dataset information and present it in the following JSON structure:
948
949     The basic info of dataset: its name is {dataset_name}, and its current info is {dataset_info}.
950
951     {{
952         "dataset_name": "{dataset_name}",
953         "info": {{
954             "description": "<brief description of the dataset>",
955             "size": "<size of the dataset (e.g., 1GB, 10,000 samples)>",
956             "scale": "<scale of the memory of the dataset (e.g., 1tb, 1gb, 100mb, 1000mb)>",
957             "author": "<author or creator of the dataset>",
958             "organization": "<organization or institution responsible for the dataset>",
959             "usage": "<how the dataset is typically used (e.g., model training, validation)>",
960             "application_fields": [
961                 "<application_field (e.g., computer vision, NLP)>"
962             ],
963             "keywords": [
964                 "<keyword_1>",
965                 "<keyword_2>"
966             ]
967         }}
968     }}
969
970     Note that you should ONLY return a json file and no any other fukcing explanation.
971
972     Use the information from the concatenated text to fill out the fields as accurately as possible.
973
974     """
975     return prompt

```

972 A.5.5 GENERATING DATASET REFERENCE DETECTION INSTRUCTIONS

973
974 This prompt instructs the LLM to analyze a research paper and identify whether it references a given
975 dataset. The LLM is provided with two inputs: the name and description of the dataset, and a text
976 string from the research paper. It is tasked with determining if the dataset is mentioned in the paper
977 and extracting relevant details about how the dataset is used. The output is structured as a JSON
978 object, containing information on whether the dataset is referenced and, if so, specific details on its
979 usage and relevant text excerpts.

980 The key tasks for the LLM include:

- 981 • Checking if the dataset is referenced in the research paper.
- 982 • Extracting relevant information on how the dataset is used (e.g., for model training, analy-
983 sis, or validation).
- 984 • Providing the specific text from the paper where the dataset is mentioned.
- 985 • Structuring the output in a JSON format, with clear fields for dataset usage, application
986 domains, and additional details.

987
988 The prompt is designed to be comprehensive, guiding the LLM through a detailed extraction process
989 to ensure accurate metadata is gathered from the research paper.

```
990 def generate_instruction_prompt(dataset_name, dataset_info):
991     instruction_prompt = f"""
```

```
992 You are provided with two inputs:
```

- 993 1. A dataset named '{dataset_name}', which is described as:
994 "{dataset_info}".
- 995 2. A string containing text from a research paper.

```
996 Your task is to:
```

- 1000 - Determine if the research paper references the dataset '{dataset_name}' at any point.
- 1001 - If the dataset is referenced, identify and extract the specific part of the paper where it is mentioned.
- 1002 - Additionally, provide detailed information about how the dataset is used in the paper, including:
1003 - Whether the dataset is used for model training, analysis, validation, comparison, or other purposes.
1004 - Any specific aspects of the dataset mentioned (e.g., size, features, or unique characteristics).
1005 - Any insights into the relevance of the dataset to the research being conducted.

```
1006 Your output should be a JSON object with the following structure:
```

```
1007 {{
1008     "dataset_referred": <true/false>,
1009     "reference_details": {{
1010         "dataset_name": "{dataset_name}",
1011         "dataset_usage": "<detailed description of how the dataset is used in the research paper>",
1012         "related_text": "<specific excerpt from the paper where the dataset is mentioned>",
1013         "application_field": "<application domains of the paper, in the form of a list of application domains>",
1014         "...: any other useful info you think, can be left as blank
1015     }}
1016 }}
1017 """
```

```
1018 Instructions:
```

- 1019 - If the dataset '{dataset_name}' is not mentioned in the paper, set "dataset_referred" to false.
- 1020 - If the dataset is mentioned, set "dataset_referred" to true and provide detailed information on its usage and relevant text excerpts.
- 1021 - Ensure that "related_text" contains an exact or closely matching excerpt from the paper where the dataset is mentioned.
- 1022 - If the dataset is referred to but no explicit usage is stated, provide an empty string for "application_field".

```
1023 """
1024 return instruction_prompt
1025 """
```

1026 A.5.6 INSTRUCTION FOR GENERATING PYTHON CODE TO VISUALIZE DATASET

1027
1028 This prompt is designed to instruct the language model to generate Python code for loading and
1029 visualizing a dataset. The model is guided to provide error-handling mechanisms and structured
1030 output based on dataset popularity and file availability. If the dataset is famous, libraries should be
1031 used; if not, the prompt asks the model to process local files to visualize the first 10 samples of
1032 the dataset. The prompt emphasizes proper error handling, data extraction, and visualization while
1033 logging useful information.

```
1034
1035 def generate_instruction_prompt(files_info, path, error_info = ""):
1036     dataset_name, dataset_info = read_metadata()
1037
1038     """
1039     Generate an instruction prompt for an LLM to generate Python code to read
1040     and visualize the elements of a dataset.
1041
1042     Parameters:
1043     - dataset_name (str): The name of the dataset.
1044     - dataset_info (list): A list of dictionaries containing file names and the
1045       head starting characters of the files (if applicable).
1046     """
1047
1048     prompt = f"""
1049     # Instruction:
1050     Generate Python code to load the dataset '{dataset_name}', retrieve the first 10
1051
1052     1. If the dataset '{dataset_name}' is famous (e.g., MNIST, CIFAR-10), use existin
1053     2. If the dataset is not famous, manually process the local dataset files provid
1054     3. Visualize the first 10 samples using matplotlib or another Python library.
1055     4. Ensure that all parts of the code (file loading, extraction, visualization) ha
1056
1057     ## Dataset Information:
1058     {dataset_info}
1059
1060     ## Local Dataset Files:
1061     {files_info}
1062
1063     ### Task:
1064     - Write a Python program to load the dataset, extract the first 10 samples, and s
1065     - Write a function to visualize the samples and save the plot figure in the fold
1066     - Ensure all functions handle errors properly, with logs or messages.
1067
1068     ### Final Output:
1069     You should only return plain Python code without any additional explanation!
1070
1071     Ensure the code follows the below structure:
1072     ```python
1073     import os
1074     import matplotlib.pyplot as plt
1075
1076     def load_dataset():
1077         try:
1078             ...
1079         except Exception as e:
1080             ...
1081
1082     def get_first_10_samples():
1083         try:
```

```

1080         ...
1081     except Exception as e:
1082         ...
1083
1084 def visualize_samples(samples):
1085     try:
1086         ...
1087     except Exception as e:
1088         ...
1089
1090 def save_run_result():
1091     ...
1092
1093 if __name__ == "__main__":
1094     try:
1095         samples = get_first_10_samples()
1096         visualize_samples(samples)
1097     except Exception as e:
1098         ...
1099
1100     try:
1101         save_run_result()
1102     except Exception as e:
1103         ...
1104
1105 Error log from previous code attempts: {error_info}
1106 """
1107 return prompt

```

1108 A.5.7 PROMPT FOR EXTRACTING DATASET DOWNLOAD LINK FROM HTML

1109 This prompt is designed for a language model to analyze HTML content and retrieve direct or
1110 indirect download links for a dataset. The model is required to provide clear instructions on how
1111 to access the dataset, including handling any intermediate steps necessary for the download. The
1112 prompt also instructs the model to infer the file format and provide detailed instructions if the dataset
1113 cannot be directly downloaded.

```

1114
1115 def generate_llm_prompt(link):
1116
1117     dataset_name, dataset_info = read_metadata()
1118
1119     prompt = f"""
1120     You are tasked with analyzing the HTML content provided to identify how to download.
1121
1122     - Dataset Name: {dataset_name}
1123     - Dataset Info: {dataset_info}
1124
1125     ### Your Objective:
1126     1. Download URL: Extract the direct download link for the dataset file if it exists.
1127     2. File Format: Determine the file format of the dataset (e.g., zip, tar, csv, etc.).
1128     3. Download Steps: Provide clear, step-by-step instructions to acquire the dataset.
1129     - Clicking a direct download link.
1130     - Navigating to another webpage to continue the download process, if there is no direct link.
1131     - Completing necessary forms or accepting terms to access the dataset.
1132     - Any other process required to reach the final dataset.
1133
1134     You should try your best to find the direct download link of the dataset. Even if it is not

```

```

1134     And sometimes there are direct download links but you misjudge them, so be more
1135
1136     NOTE!!! You should only return me a json file and do not contain any other info,
1137
1138
1139     ### JSON Output Format:
1140     Present the output as a JSON object in the following structure:
1141
1142
1143     ```json
1144     {{
1145     "dataset_name": "{dataset_name}",
1146     "download_info": {{
1147         "download_url": "<Direct download URL or 'None' if not available>",
1148         "direct_download": "<If the download url is direct or none>",
1149         "useful_info": "<any useful infos you find, like links to potential download>",
1150         "file_format": "<File format or 'Unknown'>",
1151         "potential_indirect_links": "<potential download links you think>"
1152         "download_steps": [
1153             {{
1154                 "step": 1,
1155                 "action": "<Description of the first step needed to download the dataset>"
1156             }},
1157             {{
1158                 "step": 2,
1159                 "action": "<Description of the second step, if applicable>"
1160             }},
1161             {{
1162                 "step": 3,
1163                 "action": "<Additional steps, if applicable>"
1164             }},
1165             ...,
1166             {{
1167                 "step": n,
1168                 "action": "<Additional steps, if applicable>"
1169             }},
1170         ]
1171     }}
1172     }}
1173
1174     """
1175     return prompt

```

1176 A.5.8 PROMPT FOR GENERATING PYTHON CODE TO DOWNLOAD A DATASET

```

1177     This prompt instructs the model to generate Python code that automates the process of downloading
1178     a dataset. The model must handle both direct and indirect download links, provide error handling,
1179     and ensure that the dataset is saved with the correct file structure. Additionally, the model is re-
1180     quired to produce code that is generalizable and capable of managing different dataset formats and
1181     conditions.
1182
1183     def generate_instruction(uid, idea):
1184
1185         dataset_name, dataset_info = read_metadata()
1186
1187         prompt = f"""

```

1188 Write a python file to download the dataset {dataset_name}. Here are some other
1189
1190 You are provided with an input dictionary stored in a variable called `input_data`.
1191 It is the info about a dataset {dataset_name}, with info {dataset_info}. Your goal is to
1192 The structure of the dictionary is as follows:
1193
1194 {idea}
1195
1196 The real input is in the "input" section as this is an instruction prompt.
1197
1198 NOTE THAT THE dictionary is ONLY FOR REFERENCE and it may contain FALSE INFO, so
1199
1200 Your task is to generate Python code for the following:
1201
1202 - ****Create a Python script file in the folder `draft/ideas/{uid}` with the name**
1203
1204 - ****Define a function `download_dataset()` within this file.****
1205 - This function should:
1206 - Download the dataset based on the dataset name and dataset info, and (if possible) the
1207 - Download the dataset based on the information provided in the `input_data`
1208 - If you can already find information about the dataset without using input_data
1209 - Handle both direct downloads and cases where the download requires manual
1210 - Add try-except blocks anywhere so that the code will function normally
1211 - Running the download_dataset will ensure that the dataset gets downloaded
1212 - If after trying downloading directly or indirectly (like trying all potential
1213 - Print the required download steps as outlined in the `download_steps` section
1214 - Output these instructions clearly so that the user can follow them to manually
1215 - ****Handle direct downloads:****
1216 - If `direct_download` is set to "Yes", the function should use `requests` to
1217
1218 - ****Create directories if necessary:****
1219 - Ensure that the folder `draft/dataset/{uid}` is created if it doesn't already
1220
1221 - ****Error handling:****
1222 - The function should check for errors during the download process, including
1223 - If the download fails, print a meaningful error message and proceed to try
1224
1225 - ****Log useful information:****
1226 - After a successful download, print out useful metadata about the dataset file
1227
1228 - ****File structure and naming:****
1229 - Save the dataset with a filename based on the `dataset_name` and the appropriate
1230
1231 - ****Generalization:****
1232 - Ensure that the function is generalized to handle any properly formatted input
1233
1234 - ****Edge cases and validation:****
1235 - Include validation for the existence of required fields like `download_url`
1236 - If a field is missing or invalid, the function should print an error and give
1237
1238 NOTE THAT the download link may be a link to files like csv/txt/zip/json/... , but
1239
1240 Example code structure to start:
1241
1242 ```python
1243 import os
1244 import requests

```

1242
1243 def download_dataset():
1244     ...
1245     ...
1246     ...
1247     ...
1248
1249
1250 if __name__ == "__main__":
1251     download_dataset()

```

Note that the example code may be wrong, so do not really rely on it. You should

You should only return python code that is content of get_dataset.py, and do not

And for the result python code, the function download_dataset, once run, will do

```

1257
1258 """
1259 return prompt

```

1261 A.5.9 PROMPT FOR GENERATING PYTHON CODE TO DOWNLOAD A DATASET

1262 This prompt instructs the model to generate a Python script to download a dataset, with detailed
 1263 instructions for error handling, logging, and alternative download methods. It ensures that the gen-
 1264 erated code is robust, handles edge cases, and can be run without any external parameters. The
 1265 script also includes mechanisms to handle both direct and indirect download links, create necessary
 1266 directories, and validate input data.

```

1267
1268 def generate_instruction(uid, idea):
1269
1270     dataset_name, dataset_info = read_metadata()
1271
1272     prompt = f"""
1273
1274     Write a python file to download the dataset {dataset_name}. Here are some other
1275
1276     You are provided with an input dictionary stored in a variable called 'input_data'.
1277     It is the info about a dataset {dataset_name}, with info {dataset_info}. Your goal is to
1278     The structure of the dictionary is as follows:
1279
1280     {idea}
1281
1282     The real input is in the "input" section as this is an instruction prompt.
1283
1284     NOTE THAT THE dictionary is ONLY FOR REFERENCE and it may contain FALSE INFO, so
1285
1286     Your task is to generate Python code for the following:
1287
1288     - **Create a Python script file in the folder 'draft/ideas/{uid}' with the name
1289
1290     - **Define a function 'download_dataset()' within this file.**
1291     - This function should:
1292         - Download the dataset based on the dataset name and dataset info, and (
1293         - Download the dataset based on the information provided in the 'input_data'
1294         - If you can already find information about the dataset without using input_data
1295         - Handle both direct downloads and cases where the download requires manual
1296         - Add try-except blocks anywhere so that the code will function normally
1297         - Running the download_dataset will ensure that the dataset gets downloaded

```


- 1296 - If after trying downloading directly or indirectly(like trying all potential
- 1297 - Print the required download steps as outlined in the 'download_steps' s
- 1298 - Output these instructions clearly so that the user can follow them to m
- 1299 - ****Handle direct downloads:****
- 1300 - If 'direct_download' is set to "Yes", the function should use 'requests' to
- 1301
- 1302 - ****Create directories if necessary:****
- 1303 - Ensure that the folder 'draft/dataset/{uid}' is created if it doesn't alrea
- 1304
- 1305 - ****Error handling:****
- 1306 - The function should check for errors during the download process, including
- 1307 - If the download fails, print a meaningful error message and proceed to try
- 1308
- 1309 - ****Log useful information:****
- 1310 - After a successful download, print out useful metadata about the dataset fi
- 1311
- 1312 - ****File structure and naming:****
- 1313 - Save the dataset with a filename based on the 'dataset_name' and the approp
- 1314
- 1315 - ****Generalization:****
- 1316 - Ensure that the function is generalized to handle any properly formatted in
- 1317
- 1318 - ****Edge cases and validation:****
- 1319 - Include validation for the existence of required fields like 'download_url'
- 1320 - If a field is missing or invalid, the function should print an error and g

1321 NOTE THAT the download link may be a link to files like csv/txt/zip/json/... , bu

1322

1323 Example code structure to start:

```
1324
1325 ```python
1326 import os
1327 import requests
1328
1329 def download_dataset():
1330     ...
1331     ...
1332     ...
1333
1334
1335 if __name__ == "__main__":
1336     download_dataset()
```

1337 Note that the example code may be wrong, so do not really rely on it. You should

1338 You should only return python code that is content of get_dataset.py, and do not

1340 And for the result python code, the function download_dataset, once run, will do

```
1341
1342 """
1343
1344 return prompt
```

1346 A.6 GOOGLE SEARCH API

1347 We directly crawl the top-ranked links search results from google with such code:

```
1348
1349 import requests
```

```
1350 from bs4 import BeautifulSoup
1351
1352 def search_google(query):
1353     # Make a request to Google Search
1354     url = f"https://www.google.com/search?q={query}"
1355     headers = {
1356         "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36"
1357     }
1358     response = requests.get(url, headers=headers)
1359
1360     # Check if the request was successful
1361     if response.status_code == 200:
1362         # Parse the HTML content
1363         soup = BeautifulSoup(response.text, 'html.parser')
1364         search_div = soup.find('div', {'id': 'search'})
1365         return str(search_div)
1366     else:
1367         return f"Error: {response.status_code}"
1368
1369 def save_to_file(content, filename):
1370     with open(filename, 'w', encoding='utf-8') as file:
1371         file.write(content)
1372
1373 def extract_links(html_content):
1374     soup = BeautifulSoup(html_content, 'html.parser')
1375     links = []
1376
1377     # Find all 'a' tags and get their href attributes
1378     for a_tag in soup.find_all('a', href=True):
1379         links.append(a_tag['href'])
1380
1381     # with open("links.txt", 'w', encoding='utf-8') as link_file:
1382     #     for link in links:
1383     #         link_file.write(link + '\n')
1384
1385     return links
1386
1387 def get_links(input_text):
1388     result = search_google(input_text)
1389     links = extract_links(result)
1390     return links
1391
1392 if __name__ == "__main__":
1393     # print(get_links("scope2 dataset"))
1394
1395     query = input("Enter your search query: ")
1396     result = search_google(query)
1397
1398     if "Error" not in result:
1399         save_to_file(result, "search_results.html")
1400         print("Search results saved to 'search_results.html'.")
1401
1402         # Extract links and save to a separate file
1403         links = extract_links(result)
1404         with open("draft/links.txt", 'w', encoding='utf-8') as link_file:
1405             for link in links:
1406                 link_file.write(link + '\n')
```

```
1404         print("Links saved to 'links.txt'.")
1405     else:
1406         print(result)
1407
1408
```

1409 A.7 LONG CONTEXT INFERENCE

1411 Long context inference involves processing large textual inputs that exceed typical token limits
1412 in language models. By employing techniques such as chunking, models can handle and analyze
1413 extensive documents without losing context or important details.

1414 To implement long context inference, a common approach is to break down the input text into
1415 smaller chunks, process each chunk separately, and then combine the results to form a coherent
1416 output. Below is an example of Python code implementing this approach using an API to handle
1417 long texts:

```
1418
1419 import os
1420 import requests
1421 import json, sys
1422
1423 parent_dir = os.path.abspath(os.path.join(os.path.dirname(__file__), '..'))
1424 sys.path.append(parent_dir)
1425
1426 from utils import LLMapi, clamp_prompt, clean_llm_json_res
1427
1428 # Get the OpenAI API key from environment variable
1429 API_KEY = os.getenv('OPENAI_API_KEY')
1430
1431 # Function to split the input into chunks based on token limit
1432 def split_into_chunks(text, max_char_len = 8888):
1433     chunks = []
1434
1435     # Split the text into chunks of the given max_char_len
1436     for i in range(0, len(text), max_char_len):
1437         chunks.append(text[i:i + max_char_len])
1438
1439     return chunks
1440
1441 # Function to process text of any length with chunking
1442 def call_llm_with_chunks(instruction, text, max_tokens_per_chunk=8888, max_chunk_number):
1443     chunks = split_into_chunks(text, max_tokens_per_chunk)
1444
1445     full_response = []
1446
1447     for i, chunk in enumerate(chunks):
1448         if i > max_chunk_number:
1449             break
1450         print(f"Processing chunk {i+1}/{len(chunks)}...")
1451         prompt = generate_chunk_prompt(instruction, chunk, i)
1452         response = LLMapi(prompt, model=model)
1453         if response:
1454             full_response.append(response)
1455
1456     return full_response
1457
1458 def generate_chunk_prompt(instruction, chunk, number):
1459     prompt = f"""
1460     Task: You are required to perform the following action on the provided text.
```

```

1458     Instruction:
1459     {instruction}
1460
1461     Context:
1462     The text provided below is a portion(portion number: {number}) of a larger document.
1463
1464     Important Notes:
1465     - Pay close attention to the instruction and ensure that the output reflects exactly what is asked.
1466     - If the instruction requires summarizing, ensure the result is concise while retaining all key information.
1467     - If the instruction asks for rewriting, rephrase without altering the original meaning.
1468
1469     Below is the text chunk that you should work on:
1470
1471     [Start of Text Chunk]
1472     {chunk}
1473     [End of Text Chunk]
1474
1475     Please follow the instruction precisely and produce the corresponding output.
1476     """
1477     return prompt
1478
1479 def generate_combination_prompt(instruction, chunk_responses):
1480     prompt = f"""
1481     Task: You are required to combine multiple responses generated from different chunks.
1482     The individual chunk responses may contain overlapping information, separate ideas, or
1483     Your task is to combine these responses into a single cohesive and comprehensive response.
1484
1485     The responses are results of such task: {instruction}, so merge them based on the instruction.
1486
1487     Below are the responses generated from different chunks. Please combine them into a single response.
1488     """
1489     for i, response in enumerate(chunk_responses):
1490         prompt += f"[Response {i+1}]\n{response}\n\n"
1491
1492     prompt += "Please combine the above responses into a single cohesive output, following the instruction."
1493
1494     return prompt
1495
1496 def LLM_long_api(instruction, input_text, max_chunk = 100, model="gpt-4o-mini"):
1497     res = call_llm_with_chunks(instruction, input_text, max_chunk_number = max_chunk, model=model)
1498     cb_pp = generate_combination_prompt(instruction, res)
1499
1500     return clean_llm_json_res( LLMApi(cb_pp))
1501
1502 if __name__ == "__main__":
1503     res = LLM_long_api("you need to give me a story with some input info", "the story")
1504     print(res)
1505
1506 The function calls utils, and the code of util is below:
1507
1508 import json, os, requests
1509
1510 def change_dataset_name(name):
1511     json_file_path = "draft/metadata.json"
1512
1513     with open(json_file_path, 'r') as file:

```

```

1512         data = json.load(file)
1513
1514         # Update the dataset_name
1515         data['dataset_name'] = name
1516
1517         # Save the updated JSON back to the file
1518         with open(json_file_path, 'w') as file:
1519             json.dump(data, file, indent=4)
1520
1521     def read_dataset_name():
1522         with open("draft/metadata.json", 'r') as file:
1523             data2 = json.load(file)
1524
1525         # Extract the "dataset_name" property
1526         dataset_name = data2['dataset_name']
1527
1528         return dataset_name
1529
1530     def LLMApi(input_text, max_length=8888, model="gpt-4o-mini"):
1531         api_key = os.getenv('OPENAI_API_KEY') # Get the API key from environment variable
1532         if not api_key:
1533             return "API key not found in environment variables."
1534
1535         url = "https://api.openai.com/v1/chat/completions"
1536
1537         headers = {
1538             "Authorization": f"Bearer {api_key}",
1539             "Content-Type": "application/json"
1540         }
1541
1542         # Clamp input text to max_length
1543         if len(input_text) > max_length:
1544             input_text = input_text[:max_length] # Truncate the text if it's too long
1545
1546         data = {
1547             "model": model, # Ensure you're using a valid model, e.g., "gpt-4"
1548             "messages": [
1549                 {"role": "system", "content": "You are a helpful assistant."},
1550                 {"role": "user", "content": input_text}
1551             ]
1552         }
1553
1554         try:
1555             # Send POST request to OpenAI API
1556             response = requests.post(url, headers=headers, data=json.dumps(data))
1557
1558             # If the response is successful (status code 200)
1559             if response.status_code == 200:
1560                 result = response.json()
1561                 return result['choices'][0]['message']['content'].strip()
1562             else:
1563                 return f"Error: {response.status_code} - {response.text}"
1564
1565         except Exception as e:
1566             return f"An error occurred: {e}"
1567
1568     def fetch_html_from_link(link):
1569         """Fetches HTML content from a given link."""

```

```
1566     try:
1567         response = requests.get(link)
1568         response.raise_for_status() # Raise an error for bad responses
1569
1570         return response.text
1571     except requests.RequestException:
1572         return None # Return None on error
1573
1574 from bs4 import BeautifulSoup
1575 import requests
1576
1577 def fetch_html_from_link_no_script(link):
1578     """Fetches HTML content from a given link."""
1579     try:
1580         response = requests.get(link)
1581         response.raise_for_status() # Raise an error for bad responses
1582
1583         html_content = response.text
1584
1585         # Try removing <script> tags from the HTML
1586         try:
1587             soup = BeautifulSoup(html_content, 'html.parser')
1588             for script in soup.find_all('script'):
1589                 script.decompose() # Remove the <script> tags
1590             return str(soup)
1591         except Exception:
1592             return html_content # In case of error, return the raw HTML content
1593
1594     except requests.RequestException:
1595         return None # Return None on error
1596
1597 def clamp_prompt(long_string, char_limit=8888):
1598     if len(long_string) > char_limit:
1599         return long_string[:char_limit] + '...'
1600     return long_string
1601
1602 def read_metadata(file_path='draft/metadata.json'):
1603     with open(file_path, 'r', encoding='utf-8') as file:
1604         # Load the JSON data from the file
1605         metadata = json.load(file)
1606
1607     # Extract dataset_name and convert the entire 'info' dictionary to a string
1608     dataset_name = metadata['dataset_name']
1609     dataset_info = json.dumps(metadata['info']) # Convert the 'info' dictionary to a string
1610
1611     return dataset_name, dataset_info
1612
1613 def read_metadata_dataset_websites(file_path='draft/metadata.json'):
1614     try:
1615         with open(file_path, 'r', encoding='utf-8') as file:
1616             # Load the JSON data from the file
1617             metadata = json.load(file)
1618
1619             # Extract dataset_name and convert the entire 'info' dictionary to a string
1620             dataset_websites = metadata["dataset_websites"]
1621
1622             return dataset_websites
1623     except Exception as e:
```

```

1620         print(f"failed to read_metadata_dataset_websites, reason is : {e}")
1621         return []
1622
1623     # # Example usage:
1624     # dataset_name, dataset_info = read_metadata()
1625     # print(f"Dataset Name: {dataset_name}")
1626     # print(f"Dataset Info: {dataset_info}")
1627
1628     def clean_llm_json_res(res):
1629         res_json = res
1630         try:
1631             if res.startswith('```json\n'):
1632                 res = res[len('```json\n'):].strip('\n')
1633                 # Convert the string to JSON format
1634                 res_json = json.loads(res)
1635
1636         except Exception as e:
1637             # Skip invalid JSON strings
1638             print(f"Error decoding JSON for item: {res} - {e}")
1639
1640         return res_json
1641
1642     def get_py_files_length(folder_path):
1643         total_length = 0
1644         # Traverse through all files in the folder and its subfolders
1645         for root, dirs, files in os.walk(folder_path):
1646             for file in files:
1647                 if file.endswith(".py"): # Only consider .py files
1648                     file_path = os.path.join(root, file)
1649                     with open(file_path, 'r', encoding='utf-8') as f:
1650                         total_length += len(f.readlines()) # Add number of lines in the
1651         return total_length
1652
1653     if __name__ == "__main__":
1654         folder_path = os.path.dirname(os.path.realpath(__file__)) # Get the current folder
1655         total_lines = get_py_files_length(folder_path)
1656         print(f"The total number of lines in all .py files (including this script) is: {total_lines}")
1657
1658     A.8 REPRODUCIBILITY
1659
1660     The code for the DataSEA system is available on GitHub at https://github.com/SingleView11/DataSEA. Detailed instructions for setting up the environment and running the pipelines are provided in the repository.
1661
1662     A.9 CODE STRUCTURE
1663
1664     The code for the DataSEA system is organized into three main modules: Search (S), Evaluate (E), and Analyze (A). Each module contains several Python scripts responsible for different tasks within the pipeline. Below is a detailed breakdown of the file structure:
1665
1666     app/
1667
1668     app.py # Main file to orchestrate the full pipeline
1669     utils.py # Utility functions used across modules
1670
1671     S/ # Search module

```

```

1674     convert_json_format.py           # JSON format conversion
1675     convert_json_format2.py          # Alternative JSON format conversion
1676     GetRawResponse.py                # Fetch raw responses from search queries
1677     get_firstpage_links.py           # Retrieve first-page search links
1678     links_eval.py                     # Evaluate and rank retrieved links
1679     main_s.py                         # Main script for Search module
1680     prompt_generation.py             # Generate search prompts for LLM
1681     readme.md                         # Documentation for Search module
1682     __init__.py                       # Init file for the Search module
1683
1684     E/                                # Evaluate module
1685     analyze_ref_pdfs.py              # Analyze reference papers in PDF format
1686     get_dataset_metadata.py          # Extract metadata from dataset sources
1687     get_paper.py                     # Retrieve reference papers for the dataset
1688     get_pdfs.py                      # Download and parse PDFs
1689     get_sorted_ref_papers.py         # Sort and rank reference papers by citations
1690     longtext_api.py                  # Handle long text input/output for LLMs
1691     main_e.py                        # Main script for Evaluate module
1692     sortgs_update.py                 # Update sorting logic for references
1693     __init__.py                     # Init file for the Evaluate module
1694
1695     A/                                # Analyze module
1696     analyze_dataset.py               # Generate analysis and visualizations for datasets
1697     get_download_method.py           # Determine download method for datasets
1698     main_a.py                        # Main script for Analyze module
1699     main_sea.py                      # Integrated script for Search, Evaluate, Analyze
1700     try_download_ideas.py            # Try different download ideas for dataset
1701     zip_files_final.py               # Handle final dataset packaging
1702     __init__.py                     # Init file for the Analyze module

```

The structure is modular, with each module containing its own set of scripts that handle specific steps in the DataSEA workflow. The modules are integrated by the `app.py` file, which orchestrates the end-to-end pipeline.

And below are details of using the code.

A.10 CODE EXPLANATION

This subsection provides detailed explanations for each Python file in the DataSEA system, covering the functionality, logic, and interactions with other modules.

A.10.1 QUICKSTART

For setup, install `requirement.txt`, and make sure the `openai` api key is set in your environment variable.

Then run the `app.py`. It will ask you to input a dataset name and some optional descriptive info, and then you only need to wait for about 5-10 minutes to get a zip file that stores the infos about the dataset!

A.10.2 ADVANCED RUNNING

You can also do the `s,e,a` pipelines separately by calling `s_pipeline`, `e_pipeline`, `a_pipeline` function one by one, just check the `main_s`, `main_e`, `main_a` functions.

A.10.3 SEARCH MODULE (S)

convert_json_format.py

- **Code Usage:**

- 1728 - process_judge_info(data): Processes entries to parse the judge_info field
- 1729 as JSON, if possible.
- 1730 - convert_judge_info_in_file(input_file, output_file): Reads
- 1731 JSON data from input_file, processes it, and saves it to output_file.
- 1732 - eval_pipeline(): Runs the dataset evaluation pipeline from the links_eval
- 1733 module.
- 1734

1735 convert_json_format2.py

- 1736
- 1737 • Code Usage:
- 1738 - process_judge_info(data): Processes and converts the judge_info field to
- 1739 a valid JSON object if possible.
- 1740 - filter_dataset_websites(data): Filters entries where
- 1741 is_dataset_website in judge_info is True.
- 1742 - filter_judge_info_in_file(input_file, output_file): Reads JSON
- 1743 from input_file, processes and filters the entries, and saves the filtered result to
- 1744 output_file.
- 1745

1746 GetRawResponse.py

- 1747
- 1748 • Code Usage:
- 1749 - google_response(query): Simulates a Google search by sending a search
- 1750 query to Google's search engine and saves the raw HTML response to
- 1751 raw_search_response.html.
- 1752

1753 get_firstpage_links.py

- 1754 • Code Usage:
- 1755 - search_google(query): Sends a search query to Google, parses the HTML re-
- 1756 sponse, and returns the search results as HTML.
- 1757 - save_to_file(content, filename): Saves the provided content (HTML or
- 1758 text) to a file with the specified filename.
- 1759 - extract_links(html_content): Extracts all the links from the provided
- 1760 HTML content and returns them as a list.
- 1761 - get_links(input_text): Performs a Google search for the given input text,
- 1762 extracts the links, and returns them as a list.
- 1763

1764 links_eval.py

- 1765
- 1766 • Code Usage:
- 1767 - LLMapi(input_text): Sends a request to the OpenAI API using the provided
- 1768 input text and returns the LLM's response.
- 1769 - test(dataset_name="", desc="", need_input=True): Retrieves
- 1770 dataset links, generates prompts, and sends them to the LLM API for evaluation,
- 1771 returning the results.
- 1772 - save_array_to_json(array, file_path="draft/evals.json"): Saves an array to a specified JSON file.
- 1773 - eval_pipeline(dataset_name="", dataset_desc="",
- 1774 need_input=True): Runs the evaluation pipeline, gathering and saving the
- 1775 LLM evaluations for a given dataset.
- 1776
- 1777

1778 main.s.py

- 1779
- 1780 • Code Usage:
- 1781 - process_judge_info(data): Processes the judge_info field, converting it
- to JSON if valid.

1782 - `convert_judge_info_in_file(input_file, output_file)`: Reads
 1783 JSON from `input_file`, processes `judge_info`, and writes the result to
 1784 `output_file`.
 1785 - `create_folders(base_folder="draft")`: Deletes contents in the draft
 1786 folder and creates a folder structure for storing documents and metadata.
 1787 - `create_metadata_file(base_folder)`: Creates an empty
 1788 `metadata.json` file with fields for dataset metadata.
 1789 - `s_pipeline(dataset_name="", dataset_desc="",`
 1790 `need_input=True)`: Runs the full search pipeline, including folder creation,
 1791 dataset evaluation, and processing `judge_info` into JSON.
 1792

1793 **prompt_generation.py**

- 1795 • **Code Usage:**

1796 - `fetch_html_from_link(link)`: Fetches the HTML content from a given link.
 1797 Returns the HTML as a string or None if an error occurs.
 1798 - `generate_prompt(link, dataset_name, desc="")`: Generates a
 1799 prompt based on the HTML content of the link and the dataset description. The
 1800 prompt is used to check if the link is a dataset website.
 1801 - `save_prompt_to_file(link, dataset_name,`
 1802 `filename="gen_pro.txt")`: Fetches HTML, generates a prompt, and
 1803 saves it to a file.
 1804 - `clamp_prompt(long_string, char_limit=8000)`: Clamps a string to a
 1805 specified character limit (default: 8000 characters).
 1806 - `prompts_links(dataset_name, desc="")`: Fetches dataset-related links,
 1807 generates prompts, and returns them as a list of dictionaries with `link` and `prompt`.
 1808 - `test()`: Prompts the user for a dataset name, fetches the first link, and saves a
 1809 generated prompt to a file.
 1810 - `test2()`: Prompts the user for a dataset name and a specific link, then saves a
 1811 generated prompt to a file.
 1812
 1813

1814 **__init__.py (S)**

1815 A.10.4 EVALUATE MODULE (E)

1816 The Evaluate module processes and extracts metadata from the dataset links obtained from the
 1817 Search module.

1818 **analyze_ref_pdfs.py**

- 1821 • **Code Usage:**

1822 - `extract_text_from_pdf(pdf_path)`: Extracts text from a PDF file and returns
 1823 it as a string.
 1824 - `analyze_ref_papers()`: Reads research paper links from a JSON file, extracts
 1825 PDF links, downloads PDFs, and runs analysis on them with the dataset.
 1826 - `analyze_pdfs_with_dataset(folder_path, output_file)`: Analyzes
 1827 PDFs in a folder by checking for dataset references and saves the results to a JSON
 1828 file.
 1829 - `generate_instruction_prompt(dataset_name, dataset_info)`:
 1830 Generates a prompt for an LLM to analyze how a research paper uses the given
 1831 dataset.
 1832 - `analyze_pdf_with_dataset(text)`: Sends the extracted text from a research
 1833 paper to the LLM for analysis, checking for dataset references.
 1834
 1835

1836 **get_dataset_metadata.py**

1837

1838 • **Code Usage:**

1839

1840 – `extract_links_from_file(file_path)`: Extracts links from a JSON file,
1841 looking for `link` and `download_link_dataset` fields.

1842 – `extract_all_links2(file_path)`: Extracts links from another JSON struc-
1843 ture, including nested fields like `download_link_paper` and metadata URLs.

1844 – `download_files_dataset()`: Combines all extracted links from
1845 `extract_links_from_file` and `extract_all_links2`, then processes
1846 these links for downloading.

1847 – `download_link_content(url)`: Downloads content from a URL if the file size
1848 is less than 10MB.

1849 – `save_content_to_file(content, url, content_type)`: Saves the
1850 downloaded content to a file, naming it based on the URL.

1851 – `process_links(all_links)`: Processes a list of links by downloading content
1852 for each and saving it to the appropriate folder.

1853 – `extract_text_from_file(file_path)`: Extracts text from various file types
1854 (PDF, HTML, CSV, TXT) and returns the content.

1855 – `process_folder(input_folder, output_folder)`: Extracts and pro-
1856 cesses text from all files in a folder and saves the cleaned text to the output folder.

1857 – `generate_instruction_prompt()`: Generates a prompt for LLMs to extract
1858 dataset information from concatenated text.

1859 – `process_folder_and_generate_prompt(folder_path)`: Concatenates
1860 text from multiple files, generates an LLM prompt, and processes the results.

1861 – `merge_jsons(generated_data, file_path)`: Merges generated LLM re-
1862 sults with an existing JSON metadata file.

1863 – `whole_pipeline_get_metadata_and_txt_info()`: Runs the entire pro-
1864 cess—downloads dataset files, processes text, generates a prompt, and merges results
1865 with metadata.
1866

1867

1868 **get_paper.py**

1869

1870 • **Code Usage:**

1871 – `prompts_links(dataset_name, desc="")`: Retrieves links for potential
1872 dataset papers, generates prompts for each link, and returns a list of links with as-
1873 sociated prompts.

1874 – `generate_prompt_paper(link, dataset_name, desc="")`: Generates
1875 a prompt to determine if the given link corresponds to the original paper of the dataset.

1876 – `get_json_evals()`: Retrieves dataset and paper-related links, generates prompts,
1877 and evaluates them using LLM.

1878 – `save_json_prompts()`: Retrieves evaluations from LLM for dataset and paper
1879 links and saves them in JSON format.

1880 – `dataset_link_prompts(dataset_name, desc="")`: Retrieves and gener-
1881 ates prompts for dataset-related links from the `dataset_res.json` file.

1882 – `getValidLinks(json_path)`: Filters valid links from a JSON file based on cer-
1883 tain criteria like `is_dataset_website`, `download_link_dataset_exists`,
1884 and `is_direct_data`.

1885 – `merge_link_prompts(link_prompts, dataset_link_prompts_array)`:
1886 Merges two arrays of link prompts, counting the occurrences of links and adding a
1887 number property.

1888 – `get_possible_papers()`: Runs the full process to retrieve, evaluate, and convert
1889 potential paper links into a JSON file.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

get_pdfs.py

- **Code Usage:**

- `filter_json_data(json_file, callback=None)`: Filters and returns relevant data from a JSON file based on certain paper-related attributes (`is_dataset_paper_website`, `download_link_paper_exists`, `is_direct_paper`).
- `extract_links_and_paper_links()`: Extracts both dataset and paper download links from filtered JSON data, evaluates them, and saves them in a separate JSON file.
- `get_potential_pdf_link(link, dataset_name, desc="")`: Fetches the HTML content of a link and generates a prompt to find direct download links for the original paper of the dataset.
- `save_download_links_to_json(download_links_array, file_path)`: Saves the extracted download links to a specified JSON file.
- `get_pdf_links_from_single_link(link)`: Extracts PDF links from a given URL by generating a prompt using the dataset name and metadata.
- `download_file(link, file_path)`: Downloads the content from a URL and saves it in a specified folder. Supports formats like PDF, TXT, and CSV.
- `download_pdfs_from_links(links, file_path)`: Downloads PDF files from a list of links and saves them to the specified folder.
- `download_all_pdfs()`: Runs the complete process of extracting, filtering, and downloading dataset-related PDFs from the provided links.
- `delete_all_files_in_folder(folder_path)`: Deletes all files in a specified folder.
- `delete_all_contents_in_folder(folder_path)`: Deletes all files and subfolders within a specified folder.

get_sorted_ref_papers.py

- **Code Usage:**

- `evaluate_paper(obj)`: Placeholder function for evaluating a paper. No functionality implemented yet.
- `get_gs_rank_res()`: Reads the dataset name and calls the `sortgs_main()` function to rank results based on the dataset name.
- `csv_to_json(csv_file, json_file)`: Converts a CSV file to a JSON format, saving the result in the specified JSON file.
- `get_gs_papers()`: Retrieves Google Scholar ranking results for the dataset and converts them from CSV to JSON format.

longtext_api.py

- **Code Usage:**

- `split_into_chunks(text, max_char_len=8888)`: Splits a long text into smaller chunks based on a character length limit.
- `call_llm_with_chunks(instruction, text, max_tokens_per_chunk=8888, max_chunk_number=50, model="gpt-4o-mini")`: Processes text in chunks using an LLM, based on the provided instruction and model.
- `generate_chunk_prompt(instruction, chunk, number)`: Creates a prompt for an LLM to process a specific chunk of text based on the provided instruction.
- `generate_combination_prompt(instruction, chunk_responses)`: Generates a prompt to combine multiple LLM responses from different chunks into a single cohesive output.

1944 - `LLM.long_api(instruction, input_text, max_chunk=100,`
 1945 `model="gpt-4o-mini")`: Processes a long text using an LLM by splitting
 1946 it into chunks, generating responses for each, and then combining the results.
 1947

1948 `main_e.py`

1949 • **Code Usage:**
 1950 - `get_final_metadata()`: Combines information from various sources like dataset
 1951 websites, original papers, and reference papers into the `metadata.json` file.
 1952 - `prune_metadata()`: Refines the metadata by pruning and enhancing fields like
 1953 description, size, scale, author, and usage based on evaluation data and papers. Saves
 1954 the pruned metadata to `metadata_pruned.json` and updates `metadata.json`.
 1955 - `get_prune_metadata()`: Runs both `get_final_metadata()` and
 1956 `prune_metadata()` to generate and refine the metadata.
 1957 - `e_pipeline()`: Runs the complete pipeline for retrieving papers, downloading
 1958 PDFs, processing Google Scholar papers, generating metadata, analyzing reference
 1959 papers, and pruning metadata.
 1960
 1961

1962 `main_es.py`

1963 • **Code Usage:**
 1964 - `se_pipeline()`: Combines two pipelines, `s_pipeline()` and `e_pipeline()`,
 1965 running them sequentially to process both the "S" and "E" workflows.
 1966
 1967

1968 **sortgs_update.py** NOTE: This code has source <https://github.com/WittmannF/sort-google-scholar>, and I update it for convenience.
 1969
 1970

1971 • **Code Usage:**
 1972 - `get_command_line_args()`: Parses command-line arguments for keyword, num-
 1973 ber of results, output path, sorting criteria, language filter, and other options related to
 1974 Google Scholar scraping.
 1975 - `get_citations(content)`: Extracts the number of citations from the provided
 1976 HTML content.
 1977 - `get_year(content)`: Extracts the publication year from the provided HTML con-
 1978 tent.
 1979 - `setup_driver()`: Sets up and returns a Selenium WebDriver instance to handle
 1980 Google Scholar requests.
 1981 - `get_author(content)`: Extracts the author information from the HTML content.
 1982 - `get_element(driver, xpath)`: Safely retrieves an element from the webpage
 1983 using an XPath expression with multiple attempts.
 1984 - `get_content_with_selenium(url)`: Uses Selenium to retrieve the page con-
 1985 tent from a URL, handling CAPTCHA challenges when required.
 1986 - `sortgs_main()`: Scrapes Google Scholar for papers related to a dataset, extract-
 1987 ing metadata like citations, authors, and years. Saves the results in a CSV file and
 1988 optionally plots the number of citations vs. rank.
 1989

1990 `__init__.py` (E)

1992 A.10.5 ANALYZE MODULE (A)

1993 The Analyze module is responsible for downloading, organizing, and visualizing the dataset.
 1994
 1995

1996 `analyze_dataset.py`

1997 • **Code Usage:**

- 1998 - `delete_py_files_in_folder(folder_path)`: Recursively deletes all Python
- 1999 (.py) files in the specified folder.
- 2000 - `delete_log_json_files_in_folder(folder_path)`: Recursively deletes
- 2001 all JSON log files ending in `_log.json` in the specified folder.
- 2002 - `get_analyze_code_for_all()`: Cleans up the dataset folder and generates
- 2003 Python code to analyze dataset files, extracting the first 10 samples and visualizing
- 2004 them.
- 2005 - `get_file_info_list(dataset_folder, n=500)`: Reads the first 500 char-
- 2006 acters from each file in the specified folder, returning a list of dictionaries with file-
- 2007 names and file content.
- 2008 - `generate_code_for_analyzing(files_info, path, error_info)`:
- 2009 Generates Python code for analyzing dataset files based on file content, dataset
- 2010 metadata, and past error logs.
- 2011 - `generate_instruction_prompt(files_info, path, error_info)`:
- 2012 Generates a prompt for an LLM to create Python code for loading, analyzing, and
- 2013 visualizing a dataset.
- 2014 - `analyze_and_run_code()`: Generates and runs Python code to analyze all dataset
- 2015 files.
- 2016 - `analyze_and_run_code_with_self_repair()`: Attempts to run generated
- 2017 Python files up to three times with self-repair functionality if an error occurs.
- 2018 - `regenerate_idea(file_path, e)`: Regenerates Python code for a given file
- 2019 if an error occurs during execution.

2020 **get_download_method.py**

- 2021 • **Code Usage:**

- 2022 - `delete_py_files_in_folder(folder_path)`: Recursively deletes all Python
- 2023 (.py) files in the specified folder.
- 2024 - `delete_log_json_files_in_folder(folder_path)`: Recursively deletes
- 2025 all JSON log files ending in `_log.json` in the specified folder.
- 2026 - `get_analyze_code_for_all()`: Cleans up the dataset folder and generates
- 2027 Python code to analyze dataset files, extracting the first 10 samples and visualizing
- 2028 them.
- 2029 - `get_file_info_list(dataset_folder, n=500)`: Reads the first 500 char-
- 2030 acters from each file in the specified folder, returning a list of dictionaries with file-
- 2031 names and file content.
- 2032 - `generate_code_for_analyzing(files_info, path, error_info)`:
- 2033 Generates Python code for analyzing dataset files based on file content, dataset
- 2034 metadata, and past error logs.
- 2035 - `generate_instruction_prompt(files_info, path, error_info)`:
- 2036 Generates a prompt for an LLM to create Python code for loading, analyzing, and
- 2037 visualizing a dataset.
- 2038 - `analyze_and_run_code()`: Generates and runs Python code to analyze all dataset
- 2039 files.
- 2040 - `analyze_and_run_code_with_self_repair()`: Attempts to run generated
- 2041 Python files up to three times with self-repair functionality if an error occurs.
- 2042 - `regenerate_idea(file_path, e)`: Regenerates Python code for a given file
- 2043 if an error occurs during execution.

2044 **main_a.py**

- 2045 • `a_pipeline()`: A pipeline that automates the process of:
- 2046 - `get_download_ideas()`: Retrieves ideas for how to download datasets.
- 2047 - `try_ideas_and_run_code()`: Attempts various download methods and runs the
- 2048 corresponding code.
- 2049
- 2050
- 2051

- 2052 - `analyze_and_run_code()`: Analyzes the dataset and runs the generated analysis
- 2053 code.
- 2054 - `zip_folder_with_uuid()`: Zips the dataset folder with a unique identifier.
- 2055

2056 **main_sea.py**

- 2057 • `sea_pipeline()`: A combined pipeline that runs both the S+E and A pipelines:
- 2058 - `se_pipeline()`: Runs both the S and E workflows sequentially.
- 2059 - `a_pipeline()`: Runs the dataset download, analysis, and packaging pipeline.
- 2060
- 2061

2062 **try_download_ideas.py**

- 2063 • **Code Usage:**
- 2064 - `try_ideas()`: Sets up directories, clears previous data, and iterates over dataset
- 2065 download ideas, attempting to generate Python scripts to download datasets based on
- 2066 provided ideas.
- 2067 - `generate_instruction(uid, idea)`: Generates an instruction prompt for
- 2068 the LLM to create Python code for downloading the dataset, handling errors, and
- 2069 saving the file in a specified directory.
- 2070 - `clean_code_block(code_str)`: Cleans up LLM-generated code by removing
- 2071 any surrounding markdown formatting (like ```python).
- 2072 - `evaluate_idea(idea)`: Uses an LLM to generate Python code for a dataset
- 2073 download based on the provided idea, and saves both the code and the status of the
- 2074 evaluation.
- 2075 - `run_all_python_files_in_folder(folder_path)`: Recursively finds and
- 2076 runs all Python files in a given folder and its subfolders, handling errors and logging
- 2077 results.
- 2078 - `try_ideas_and_run_code()`: Combines `try_ideas()` and
- 2079 `run_all_python_files_in_folder()` to first attempt dataset download
- 2080 ideas and then run the generated Python scripts.
- 2081

2082 **zip_files_final.py**

- 2083 • `zip_folder_with_uuid(folder_path="draft", use_uuid=False)`:
- 2084 - This function zips the contents of a specified folder and saves it as a '.zip' file. The
- 2085 zip file is named using the dataset's name, and if `use_uuid` is set to True, a UUID
- 2086 is appended to the filename.
- 2087 - The zip file is saved in the `experiment_results` folder. The function ensures this
- 2088 folder is created if it does not exist.
- 2089 - By default, the `draft` folder is zipped, but you can specify a different folder by
- 2090 passing the `folder_path` argument.
- 2091
- 2092

2093 **__init__.py (A)**

2095 A.10.6 MAIN SYSTEM COORDINATION

2096 **app.py**

- 2097 • `sea_pipeline_without_input(dataset_name, dataset_desc)`:
- 2098 - This function executes the SEA pipeline (Search, Evaluate, Analyze) without requir-
- 2099 ing user input. It accepts a dataset name and description, passing them to the respective
- 2100 pipeline functions `s_pipeline`, `e_pipeline`, and `a_pipeline`.
- 2101 - `batch_get_experiment_res(arr)`:
- 2102 - This function takes a list of dataset names and runs the
- 2103 `sea_pipeline_without_input` for each dataset in the list, automating the
- 2104 execution of the full pipeline for multiple datasets.
- 2105

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

utils.py

- `change_dataset_name(name):`
 - This function updates the dataset name in the `metadata.json` file.
- `read_dataset_name():`
 - Reads the `dataset_name` from the `metadata.json` file.
- `LLMApi(input_text, max_length=8888, model="gpt-4o-mini"):`
 - Sends an API request to an LLM (GPT model) with the given input text, truncating it if it exceeds the character limit.
- `fetch_html_from_link(link):`
 - Fetches raw HTML content from a given URL.
- `fetch_html_from_link_no_script(link):`
 - Fetches HTML content from a given URL, removing any `<script>` tags from the content.
- `clamp_prompt(long_string, char_limit=8888):`
 - Truncates a string if it exceeds a specified character limit.
- `read_metadata(file_path='draft/metadata.json'):`
 - Reads metadata from the specified `metadata.json` file and returns the dataset name and the dataset info as a string.
- `read_metadata_dataset_websites(file_path='draft/metadata.json'):`
 - Reads the `dataset_websites` field from the metadata file.
- `clean_llm_json_res(res):`
 - Cleans and decodes the JSON response from an LLM, removing code block formatting.
- `get_py_files_length(folder_path):`
 - Calculates the total number of lines in all Python files in the specified folder and its subfolders.