

AutoClean: LLMs Can Prepare Their Training Corpus

Xingyu Shen¹ Shengding Hu¹ Xinrong Zhang¹ Xu Han^{1*}
Xiaojun Meng² Jiansheng Wei² Zhiyuan Liu^{1*} Maosong Sun¹

¹Department of Computer Science and Technology, Institute for Artificial Intelligence,
Beijing National Research Center for Information Science and Technology,
Tsinghua University, Beijing, China.

²Huawei Noah’s Ark Lab

xingyu-c21@mails.tsinghua.edu.cn, {han-xu, liuzy}@tsinghua.edu.cn

Abstract

Recent studies highlight the reliance of Large Language Models (LLMs) on high-quality, diverse data for optimal performance. The data sourced from the Internet is often aggregated into datasets like Common Crawl, which presents significant quality variability and necessitates extensive cleaning. Moreover, specific domain knowledge is usually presented in HTML, but there is a lack of effective methods to clean HTML pages into the pre-training corpus automatically. Traditional cleaning methods involve either labor-intensive human teams that lack scalability or static heuristics that lead to suboptimal outcomes and cannot be applied to specific target domains. In this paper, inspired by the recent progress in employing LLMs as versatile agents for diverse tasks, we take the initiative to explore the potential of these agents in automating data-cleaning methodologies. By configuring LLMs as an agent team that imitates the human data-cleaning team, we can automatically generate cleaning rules that traditionally require the involvement of data-cleaning experts. These rules are developed using a limited number of data samples and can then be applied broadly to substantial portions of raw data from the same domain. We demonstrate the efficiency and effectiveness of AutoClean on both pre-training corpora such as Common Crawl and specific target websites. Both automatic and human evaluations of the quality of the cleaned content highlight the feasibility of using LLMs to prepare their pre-training corpus.

1 Introduction

The rapid advancement of Large Language Models (LLMs) has opened pathways toward AGI, with capabilities like programming (Roziere et al., 2023) and instruction-following (Wei et al., 2021). However, these models face a critical bottleneck: the scarcity of high-quality training data. According

to scaling laws (Kaplan et al., 2020), larger models demand more extensive datasets. Training data for LLMs primarily comes from two sources: (1) vast web-scale datasets like Common Crawl, which contain a huge corpus but also includes substantial noise; and (2) specific domain repositories rich in specialized knowledge. However, both sources lack effective automated methods for extracting clean, high-quality text, particularly from dynamic or noisy web content in HTML format.

Existing data cleaning methods, such as CC-Net (Wenzek et al., 2020), Pile (Gao et al., 2020), and RefinedWeb (Penedo et al., 2023), rely on fixed pipelines that incorporate deduplication, classification, and filtering. While these approaches have achieved some success, they are constrained by static, heuristic-driven processes, which often struggle to handle diverse or complex data. Additionally, these pipelines require an extensive volume of original data to operate effectively.

We introduce AutoClean, a novel approach that leverages LLMs as autonomous agents for intelligent and scalable data cleaning. Specifically, AutoClean assumes that web pages within the same domain share a consistent structure. Therefore, a set of heuristic rules generated by LLMs can effectively handle all web pages in that domain. By analyzing domain-specific web pages, AutoClean generates cleaning rules tailored to each domain and then applies these rules across all pages within it. The entire process comprises several stages, including web page clustering, HTML processing, text processing, and quality inspection. Each step will be explained in detail in the Methods section.

We conduct comprehensive experiments and analyses to show the effectiveness and efficiency of AutoClean. We instantiate AutoClean with GPT-3.5 (OpenAI, 2021) as the agents and apply the rules generated by these agents to raw data from both Common Crawl and certain websites. Both automatic evaluation and human evaluation demon-

*indicates corresponding authors.

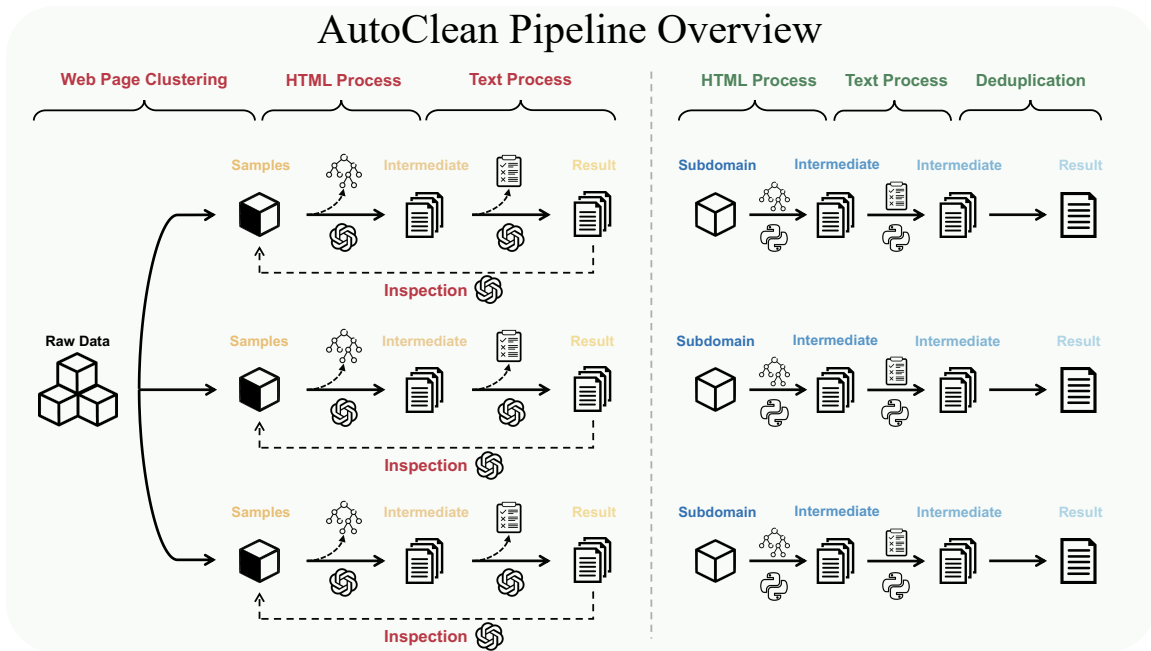


Figure 1: AutoClean consists of two parts: rule generation and data cleaning. The left part shows the cleaning rules generated by AutoClean based on the randomly collected samples, while the right part shows AutoClean cleaning the entire corpus according to the generated rules.

strate superior data cleaning performance compared to previous heuristic methods. We provide a [Demo](#) to show the effectiveness of AutoClean. To summarize, our contributions are as follows:

1. Design a pipeline leveraging LLMs for autonomous corpus cleaning.
2. Demonstrate the effectiveness of AutoClean in processing large-scale corpora and specific website cleaning.
3. Show through both automatic and human evaluations that our method achieves significantly cleaner text compared to heuristic approaches.

2 Related Work

Two lines of work are related to this paper: data cleaning methods, and agent automation.

2.1 Data Cleaning Methods

Prior to LLMs, training datasets were manually curated for training task-specific models (Zhu et al., 2015; Zhang et al., 2015). The scaling of pre-trained models necessitated larger datasets, leading to the adoption of web-crawled data like Common Crawl (Kreutzer et al., 2022; Dodge et al., 2021; Penedo et al., 2023). However, such data contains noise, such as low-quality or unsafe content (Trinh and Le, 2018; Kreutzer et al., 2022), requiring effective cleaning methods.

Cleaning approaches primarily aim to remove low-quality content. Early methods like fast-Text (Grave et al., 2018) used deduplication and language filtering. CCNet (Wenzek et al., 2020) employed PPL scores, while heuristic methods like punctuation counts were explored (Raffel et al., 2020). Pile (Gao et al., 2020) applied selectors trained on OpenWebText2 to filter Common Crawl data. These refined datasets underpin many LLMs (Brown et al., 2020; Raffel et al., 2020). However, these heuristic-heavy methods require significant human effort and lack scalability. In contrast, AutoClean is a flexible and scalable solution for cleaning large-scale web data.

2.2 LLM Agent

LLM agents have emerged as tools for automating complex tasks. We focus on two aspects: tool usage and multi-agent collaboration. Innovations like AutoGPT (Significant Gravititas) and XAgent (XAgent, 2023) enable LLMs to perform multi-step operations via APIs. AutoClean integrates cleaning tools to enhance LLM utility.

In multi-agent collaboration, frameworks like MetaGPT (Hong et al., 2023) and AgentVerse (Chen et al., 2023) simulate human-like interactions, while works like ChatDev (Qian et al., 2023) automate software design. Building on these advancements, AutoClean mimics human data-

Proportion of Characters Cleaned in Each Step

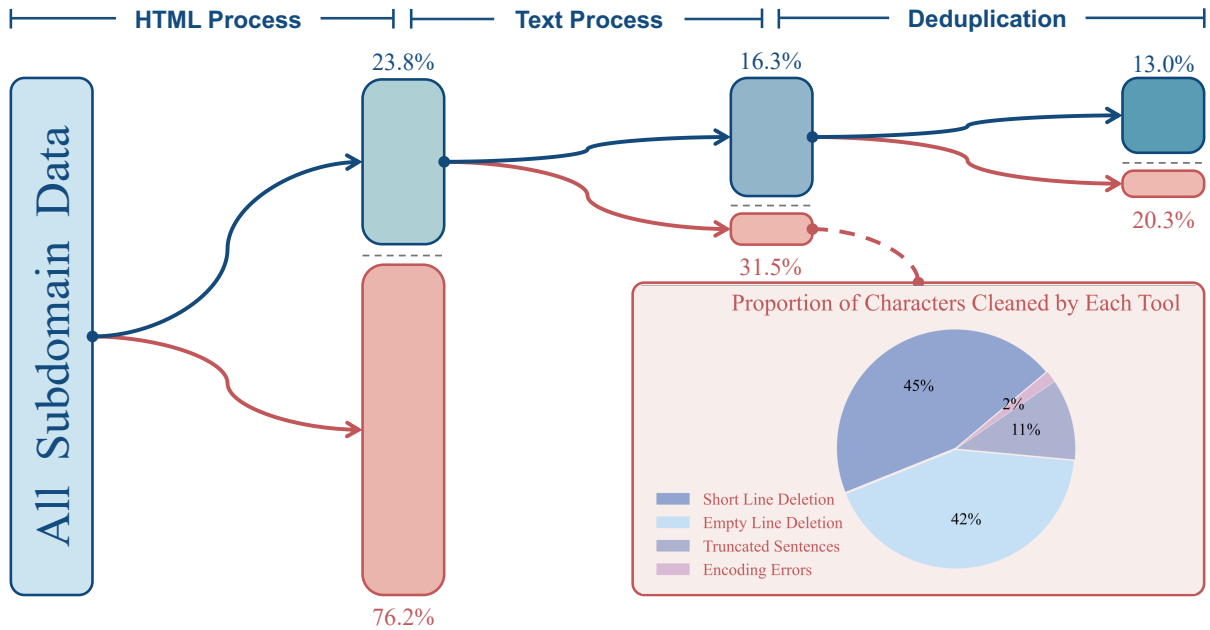


Figure 2: This figure shows the proportion of characters removed at each cleaning step. The blue numbers indicate the proportion of remaining characters relative to the original data. The red numbers indicate the proportion of discarded characters relative to the data before the current processing step. The pie chart illustrates the proportion of characters removed by each cleaning tool in the text process step.

cleaning workflows, achieving results comparable to human engineers.

3 Method

The AutoClean method starts by generating cleaning rules based on a small sample of web pages. These rules are then applied to clean the entire dataset. Importantly, the rules are created from a few random web pages but can be applied to all pages within the same domain. This allows for efficient and cost-effective cleaning of large-scale datasets without requiring agents. Below, we outline the stages of rule generation.

3.1 Web Page Clustering

AutoClean primarily leverages the similarity between different web pages under the same subdomain to clean the web pages. Hence the first step is to partition all web pages in a domain into subdomains. The desired outcome is that the web pages within each subdomain are highly similar. The similarity of a subdomain is defined by the similarity of randomly selected pairs of web pages within the subdomain. For the similarity of web pages, they are deemed similar if the HTML nodes with a depth of less than a fixed threshold are identical. We use a recursive method to divide a large domain

into highly similar subdomains, initially checking similarity, dividing by next-level domain names if needed, and merging smaller subsets during the process.

3.2 HTML Process

In this step, we utilize the Observer Agent to observe web pages, identifying nodes in the HTML structure tree that we wish to retain or delete. Based on a large amount of such data, we derive the Xpath paths where high-quality and low-quality texts are located for web pages in this subdomain.

3.2.1 HTML Observation

Node Quality Identification. 🤖 We randomly sample some web pages for the Observer Agent to select nodes with high-quality and low-quality content. Specifically, leaf nodes are defined as all nodes whose HTML tags belong to a whitelist that includes all tags used to display text. And <div> nodes containing text directly are also leaf nodes. Then the Observer Agent will select high-quality nodes from all leaf nodes, while the unselected leaf nodes are considered low-quality nodes.

Xpath Generation. We use two distinct strategies to identify high-quality and low-quality Xpath

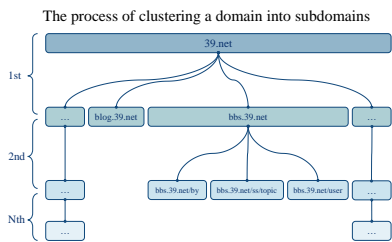


Figure 3: The process of web page clustering applying on 39.net. The nodes with insufficient similarity are divided into child nodes. Leaf nodes represent the resulting subdomains.

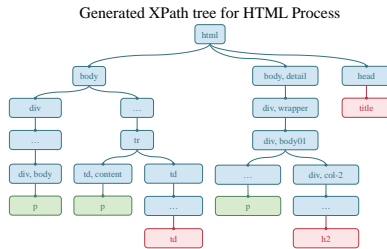


Figure 4: An example of the XPath tree for a subdomain. The green/red nodes represent high-quality/low-quality Xpaths.

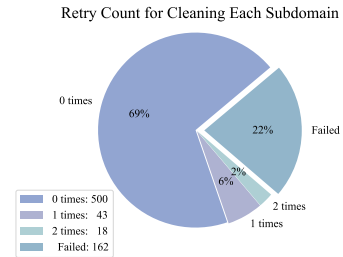


Figure 5: This pie chart describes the number of cleaning retries for all subdomains. The legend shows the subdomain counts in each category in the pie chart.

paths. Paths that lead from the root to all high-quality nodes are termed H-paths, while those leading to low-quality nodes are called L-paths. A path is considered valid if the number of H-paths using it as a prefix surpasses a certain threshold. For a path to be deemed high-quality, it must be valid and must not be a strict prefix of any other valid path. Conversely, a path is classified as low-quality if its occurrence among L-paths exceeds another threshold. For any web page within this subdomain, we start by removing all content associated with low-quality XPath paths. We extract the portions under high-quality XPath paths from the remaining content to complete the HTML processing for that web page.

3.3 Text Process

In this phase, we adopt an observation-cleaning cycle to apply various cleaning tools [2]. This phase is completed collaboratively by two agents. The Observer Agent first samples the dataset and generates an observation report detailing the current data quality issues. The Programmer Agent then reads the observation report and intelligently selects and applies some or all of the tools from the provided tool library to clean the data.

Observation Agent. 🤖 This agent samples text from the last stage’s result and generates an observation report based on a set of data quality criteria. If a text is too long, it will be split into multiple chunks, with each chunk summarized individually. The final summaries of all chunks from samples are then condensed into a single observation report, which is passed on to the Programmer Agent.

Programmer Agent. 🤖 This agent reviews each cleaning tool by reading both the observation report and the tool usage instructions. Then this agent

will determine whether each tool is applicable. Ultimately, all applicable tools will be added to the rules of this subdomain. Hence these tools will be applied to all web pages within this subdomain.

3.4 Quality Inspection

In this stage, we will inspect the results from the previous step by resampling a portion of web pages and applying the cleaning rules developed earlier. The Inspector Agent will evaluate the outcomes, splitting lengthy articles into chunks based on a fixed threshold. Each chunk will be assessed for quality, determining if it aligns more with high-quality content or spam.

If good chunks exceed a certain proportion of total characters, the subdomain and its cleaning rules are deemed valid. Otherwise, the agent will request re-cleaning from the HTML process stage. If retries reach three or more, the domain will be classified as uncleanable and abandoned.

3.5 Deduplication

The web pages within a subdomain are highly similar, so we apply a line deduplication operation. Specifically, in each subdomain, we retain only the first occurrence of any completely identical text line, removing all subsequent duplicates.

Finally, we obtain a series of subdomains and their corresponding cleaning rules. By matching a web page’s URL with the subdomain, we can apply the appropriate cleaning rules to extract high-quality textual data from the web page.

4 Experiments

In this section, we present the experiments. We conducted an experiment using AutoClean on Common Crawl and compared it with traditional

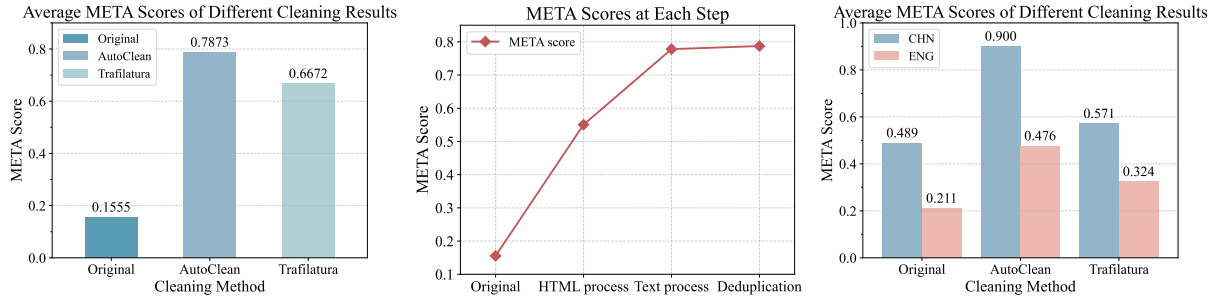


Figure 6: Average META Scores of Different Cleaning Results.

Figure 7: META Scores for intermediate results after different stages.

Figure 8: The META scores on www.gushiwen.cn (CHN) and hug-gingface.co (ENG).

pipelines to demonstrate the advantages of AutoClean. Specifically, we run AutoClean on a 1TB¹ Common Crawl corpus containing 20 domains. All the steps described above are performed, resulting in a dataset approximately 14GB in size, containing 6,000,000 documents.

4.1 The Result of Web Page Clustering

723 subdomains are obtained through the web page clustering. Figure 3 shows the process of dividing one of the 20 domains into several subdomains. First, the domain 39.net is divided into several subsets according to the next level of domain name, including (blog.)39.net and (bbs.)39.net. The subset blog.39.net has already achieved sufficient similarity, so it stops splitting and becomes a subdomain. And bbs.39.net continues to be divided into three subsets. The leaf nodes represent the resulting subdomains derived from the domain.

4.2 The Result of HTML Process

In Figure 4, We demonstrated an Xpath tree to visualize high/low-quality paths generated in the HTML process stage. These paths are rules used to extract high-quality content from web pages. The extraction will be operated on the HTML structure tree. Firstly, everything under the low-quality paths represented by red nodes will be removed. Then all contents under high-quality paths represented by green nodes will be extracted as the result.

4.3 Quality Inspection

Figure 5 shows the different outcomes of the 723 subdomains. Most subdomains met the quality requirements within 3 cleaning attempts, forming a rule set. 69% of the subdomains passed on the first attempt, while a small portion of subdomains

passed within 2 retries. 22% of the subdomains did not pass after 3 attempts and were abandoned.

4.4 Dataset Quality Assessment by Human Experts

Trafilatura (Barbaresi, 2021) is a traditional method that extracts high-quality text directly from HTML. It obtains high-quality content by using a large number of empirical Xpath and regular expressions. To assess the quality of the dataset produced by our pipeline, 1000 web pages are randomly sampled from the 14GB dataset. We compare the cleaning result of AutoClean and Trafilatura (Barbaresi, 2021) on these web pages.

Web Extraction Issues	AutoClean	Trafilatura
Navigation Information	6.00%	30.33%
Irrelevant Information	3.00%	18.33%
Pagination Information	1.67%	2.67%
Top Navigation Bar	0.00%	4.67%
HTML Tags/Codes	0.67%	0.00%

Table 1: The percentages indicate the proportion of samples in the dataset that exhibit the issues described in each row.

Table 1 shows the results of comparing the cleaning effect on 300 web pages out of 1000 samples. The comparison method is human annotation. Annotation jobs are completed by 4 professional data annotators from large language model companies. The judgment criteria are based on a data cleaning standards manual which includes 26 rules, and approximately 2000 words, with over 60 illustrative examples. The relevant parts of the document are summarized in the appendix B.

Table 1 shows that our method significantly improves the removal of navigation bars and irrelevant information from web pages compared to the Trafilatura (Barbaresi, 2021) method.

¹Disk space, including the HTML scripts.

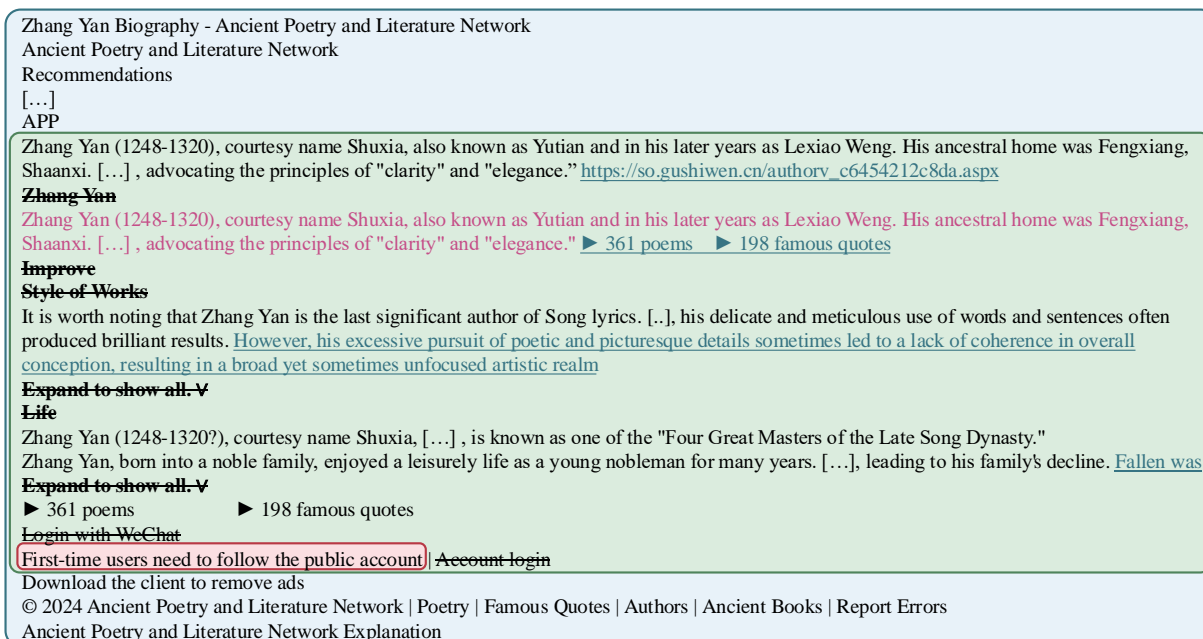


Figure 9: The selected page from www.gushiwen.cn for demonstrating in detail how AutoClean converts it into clean corpus. Different colored boxes and different font styles represent the content removed from each part during the cleaning process. The specific details of each color/style corresponding to the cleaning steps are described in Section 5. The omitted parts indicated by the ellipsis are similar to the adjacent ones. In high-quality/low-quality content, the ellipsis also represents high-quality/low-quality content.

4.5 Quantitative Dataset Quality Assessment

We also adopted the META method (Sharma et al., 2024) for a more comprehensive evaluation of our data cleaning effectiveness. The META (Sharma et al., 2024) method classifies high-quality corpus by scoring each corpus. It will first set a series of heuristic rules, and then calculate the weight of these heuristic rules based on the text perplexity changes caused by these rules. The weights are then used to score the corpus. We applied the META method (Detailed settings are in Appendix C) to evaluate 1000 samples. Text directly extracted by the HTML parser, as well as the cleaning results from AutoClean and Trafilatura are scored. The results in Figure 6 show that Trafilatura (Sharma et al., 2024) optimizes the quality, while AutoClean provides a better performance.

4.6 Analysis of Various Stages During Cleaning Process

As mentioned in Section 4.5, we also conducted a quantitative evaluation of the quality of the intermediate results in AutoClean. It demonstrates the role of each stage in improving data quality. In Figure 7, it can be observed that after the HTML process, the data quality improves significantly, while the text process and deduplication also contribute to some

improvement in data quality.

4.7 The Cost of AutoClean

In this experiment, an average of approximately 1.71×10^6 GPT-3.5-Turbo tokens are used in each subdomain, and we handle a total of 723 generated subdomains from 20 domains. This means that generating a suitable rule for a domain costs around \$40. According to our experiment, the 20 domains can obtain about 1.3% of high-quality data from the same domain corpus in Common Crawl, especially considering the large volume of data in Common Crawl and the high quality of our method’s cleaning results.

5 Domain Specific Data Acquisition

AutoClean offers unique advantages in acquiring specific types of data from targeted websites on any scale. Unlike traditional methods such as CC-Net (Wenzek et al., 2020), which depend on large-scale corpus and rely on selecting segments with the lowest PPL, AutoClean can effectively clean a domain at any scale, including smaller datasets where low-PPL segments may lack reliability.

For demonstration, we applied AutoClean on www.gushiwen.com and huggingface.co. A visual example is shown in Figure 9, where AutoClean

extracts high-quality paths (green) while removing low-quality paths (red) during the HTML process stage. Additional text cleaning tools—short line deletion (bold strikethrough) and truncated sentence removal (blue underline)—were selected by the Programmer Agent. Identical lines (magenta) were deduplicated, effectively removing navigation bars and web components.

Figure 11 illustrates the high-quality/low-quality paths generated for www.gushiwen.cn, and Figure 10 highlights the cleaning contribution of each stage, showing that the HTML process handles most junk content, followed by text processing tools. Only 14.2% of the original corpus was retained after cleaning.

We evaluated the cleaned data using the method from the previous section. As shown in Figure 8, AutoClean achieved the highest quality scores on both www.gushiwen.cn and huggingface.co, significantly outperforming Trafilaturo (Barbatesi, 2021), even at the scale of 100 web pages. A Demo is provided to demonstrate AutoClean on any domain.

6 Conclusion

In this paper, we present AutoClean, a framework designed for automatic data cleaning by utilizing LLMs as agents. AutoClean generates a comprehensive set of cleaning rules using agents for each domain, thereby ensuring scalability, flexibility, and effectiveness. Future research directions include augmenting the AutoClean intelligence level to support more sophisticated data cleaning processes and distilling the capability of LLMs into smaller models to ensure effectiveness.

Limitations

There are several limitations of our work. The flexibility of the pipeline is somewhat constrained. Even though AutoClean generates rules intelligently, the pipeline adheres to a predetermined workflow, mirroring a typical data cleaning team’s process. This could potentially limit the adaptability of the approach in diverse scenarios. The scale of the experiment presents another limitation. Owing to resource constraints, AutoClean has not been tested extensively with the raw data corpus of Terabytes. The demonstration of its effectiveness is, therefore, restricted to a limited number of domains. We currently do not provide a direct comparison of model performances trained with corpus cleaned by AutoClean and other methods.

However, we anticipate that a cleaner corpus will bring substantial performance improvement.

Ethical Considerations

In this paper, we present AutoClean, a novel data-cleaning workflow empowered by LLM agents. In the cleaning process of data, we currently do not include the step of screening for unsafe content, such as material exhibiting political bias. While there remains a possibility that the cleaned corpus may still contain such content, it is crucial to note that AutoClean’s output comprises a set of rules for refining raw data, rather than content generated by the LLM itself. Consequently, AutoClean inherently avoids the introduction of additional unsafe content into the cleaned corpus. The data source utilized in this study is open-source. During the comparative analysis between AutoClean and human data cleaning, the engineer tasked with refining four domains of the corpus was formally employed under our supporting affiliations, ensuring legal compliance. We used GPT-4 as a tool for grammar correction in our paper writing.

Acknowledgements

Thanks for the kind suggestions and support from AI Data Technology Laboratory of Huawei Noah’s Ark Lab.

References

- Adrien Barbaresi. 2021. *Trafilatura: A Web Scraping Library and Command-Line Tool for Text Discovery and Extraction*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. volume 33, pages 1877–1901.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chen Qian, Chi-Min Chan, Yujia Qin, Yaxi Lu, Ruobing Xie, et al. 2023. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors in agents. *arXiv preprint arXiv:2308.10848*.
- Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. 2021. Documenting large webtext corpora: A case study on the colossal clean crawled corpus. *arXiv preprint arXiv:2104.08758*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Édouard Grave, Piotr Bojanowski, Prakhara Gupta, Armand Joulin, and Tomáš Mikolov. 2018. Learning word vectors for 157 languages. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Julia Kreutzer, Isaac Caswell, Lisa Wang, Ahsan Wahab, Daan van Esch, Nasanbayar Ulzii-Orshikh, Allahsera Tapo, Nishant Subramani, Artem Sokolov, Claytone Sikasote, et al. 2022. Quality at a glance: An audit of web-crawled multilingual datasets. *Transactions of the Association for Computational Linguistics*, 10:50–72.
- OpenAI. 2021. [Introducing chatgpt](#).
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.
- Vasu Sharma, Karthik Padthe, Newsha Ardalani, Kushal Tirumala, Russell Howes, Hu Xu, Po-Yao Huang, Shang-Wen Li, Armen Aghajanyan, and Gargi Ghosh. 2024. Text quality-based pruning for efficient training of language models. *arXiv preprint arXiv:2405.01582*.
- Significant Gravitas. [AutoGPT](#).
- Trieu H Trinh and Quoc V Le. 2018. A simple method for commonsense reasoning. *arXiv preprint arXiv:1806.02847*.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Édouard Grave. 2020. Ccnet: Extracting high quality monolingual datasets from web crawl data. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4003–4012.
- XAgent. 2023. Xagent: An autonomous agent for complex task solving.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *NIPS*.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27.

Tool Name	Usage
Encoding Errors	Identify and correct all characters with encoding errors.
Short Line Deletion	Delete all lines containing fewer than 20 characters.
Empty Line Deletion	Delete all lines that contain only spaces.
Adjacent Deduplicate	Delete adjacent lines that are completely identical.
Full-width to Half-width	Convert all full-width characters to half-width characters.
Truncated Sentences	Delete the last sentence if it does not end with a Chinese or English period to address the issue of text truncation.

Table 2: All tools provided in text process. The cleaning effects of each tool are described in the Usage column.

A Cleaning Tools Provided in Text Process

Table 2 presents all the tools provided to the Programmer Agent during the text processing stage. Each tool consists of a Python function and a textual instruction describing its usage.

B Data Cleaning Standards Manual

In this section, we explain the specific criteria for determining each web extraction issue listed in Table 1.

Navigation Information. There exists a list containing a series of hyperlinks. The texts in hyperlinks should include ellipses or be truncated. It’s also suitable when the hyperlink texts include dates or comment numbers.

Irrelevant Information. There exist entire lines on the web page that do not relate to the main content. Irrelevant content inserted in lines related to the main content is not included in this category.

Pagination Information. "Previous page", "next page" and page numbers information used for navigating web pages appears in the cleaning result.

Top Navigation Bar. There is a navigation bar appearing at the top of a web page containing numerous categories. It usually includes many hyperlinked phrases for navigating between major categories.

HTML Tags/Codes. Information such as HTML tags and codes, like `[tag][/tag]` and `>`; appears in the cleaning result.

C Heuristic Rules Used in META Method

We set up new heuristic rules to address different languages (META (Sharma et al., 2024) only provides heuristic rules for English). Table 3 and

Characters Cleaned by Each Step

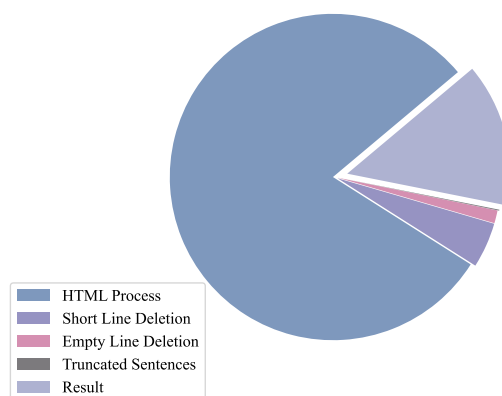


Figure 10: Characters cleaned by each part. The result represents characters retained in the end. Short line deletion, empty line deletion, and truncated sentences are tools selected in the text process.

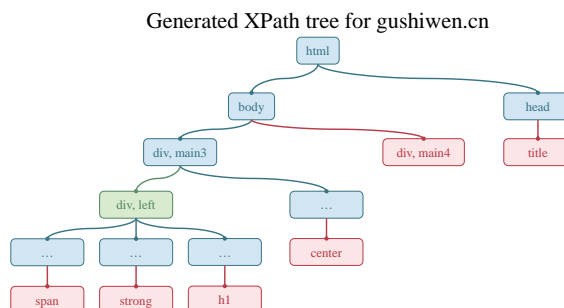


Figure 11: The high-quality/low-quality paths for www.gushiwen.cn with the same meaning as Figure 4. A node is valid if and only if it has a green ancestor and does not have any red ancestors.

Table 4 show the heuristic rules we set for Chinese and English, respectively.

We use the corpus extracted directly by the HTML parser to calculate the weights of all heuristic rules. Then, we use these weights to score the results cleaned by AutoClean and Trafilatura (Barbresi, 2021).

Filter Name	Description
token_count_ge_3	Check if the token count is > 3.
word_count_3_256	Check if line word count is > 3 and < 256.
stop_word_match_2	Check if the line contains at least 2 stop words.
no_special_characters	Check if there is any '{' or '}'.
has_personal_pronoun	Check if there is any personal pronoun in the line.
terminal_punctuation	Check if the lines end with one of the Chinese punctuation marks.
digit_punctuation_ratio_0_25	Identify lines with a ratio of digits/punctuation to words in a line is > 0.25.

Table 3: Heuristic rules for the META method applying to Chinese corpus.

Filter Name	Description
has_noun	Check if the line has a noun.
has_determiner	Check if the line has a determiner.
token_count_ge_3	Check if the token count is > 3.
word_count_3_256	Check if line word count is > 3 and < 256.
stop_word_match_2	Check if the line contains at least 2 stop words.
no_special_characters	Check if there is any '{' or '}'.
has_object	Check if there is any object identified by the parser in this line.
terminal_punctuation	Check if the lines end with one of the English punctuation marks.
digit_punctuation_ratio_0_25	Identify lines with a ratio of digits/punctuation to words in a line is > 0.25.

Table 4: Heuristic rules for the META method applying to English corpus.

D Prompts

quality nodes.

In this section, we list the prompts we used in each agent.

D.1 Observer Agent’s Prompts

<TASK>:
Imagine you are a data cleansing engineer and now you are given a web page with some paragraphs and their HTML tags and asked to weed out low-quality content such as advertisements, buttons, page components, related recommendations, page sidebars, etc., and select semantically rich and coherent body paragraphs. Output their numbers, one number per line. If there are no semantic paragraphs, output "NONE".
{INPUT}

<TASK>:
Imagine you are a data cleansing engineer and now you are given a web page with some paragraphs and their HTML tags and asked to weed out low-quality content such as advertisements, buttons, page components, related recommendations, page sidebars, etc., and select semantically rich and coherent body paragraphs. Output their numbers, one number per line. If there are no semantic paragraphs, output "NONE".
{INPUT}

<TASK>: Imagine you are a data engineer. Could you please check whether this text, which is the training corpus for a large model, contains low-quality content, incorrect punctuation, garbage short lines, and a host of other issues that degrade the quality of the corpus?
These issues specifically include but are not limited to:
1. the text contains redundant Markdown characters or has extra-long markdown reference paragraphs.
2. truncation problems in the text and semantic disjunctions at the end of the data.
3. extra line breaks, blank characters, wrong indentation, and other formatting problems.
4. incorrect use of punctuation, mixed use of full and half-width symbols, a large number of abnormal continuous symbols
5. irrelevant content in the paragraph, usually inserted advertisements or page components.
6. low-quality short lines, a large number of low-quality lines of short length in the article.
7. other problems.
Please output the problems you found in this text.
<TEXT>: {INPUT}

The two prompts above are used in section 3.2.1 to motivate the Observer Agent to select high-

This prompt provides a standard for the Observer Agent to summarize the problem in one document.

Zhang Yan (1248-1320), courtesy name Shuxia, also known as Yutian and in his later years as Lexiao Weng. His ancestral home was Fengxiang, Shaanxi. [...], advocating the principles of "clarity" and "elegance."

It is worth noting that Zhang Yan is the last significant author of Song lyrics. [...]. Due to his expertise in music theory, his delicate and meticulous use of words and sentences often produced brilliant results.

Zhang Yan (1248-1320?), courtesy name Shuxia, also known as Yutian and in his later years as Lexiao Weng. [...]. In literary history, he and another famous lyricist, Jiang Kui, are collectively known as "Jiang and Zhang." He, along with famous lyricists of the late Song Dynasty, Jiang Jie, Wang Yisun, and Zhou Mi, is known as one of the "Four Great Masters of the Late Song Dynasty."

Zhang Yan, born into a noble family, enjoyed a leisurely life as a young nobleman for many years. In 1276, when Yuan soldiers captured Lin'an, Zhang Yan's grandfather, Zhang Ru, was killed by the Yuan forces, and their family wealth was confiscated, leading to his family's decline.

► 361 poems ► 198 famous quotes

Figure 12: The cleaning result of Figure 9.

```
<TASK>: Please summarize the following multiple reports on the issue of training corpus quality into a report, you only need to output the summary.  
<REPORTS>: {INPUT}
```

This prompt is used to combine all summarized problems into a single report.

D.2 Programmer Agent's Prompts

```
<TASK>: Below you will be given a description of what a data cleansing tool does and a report of a problem with the existing data and asked to determine if the data should be cleansed using this tool, if yes please output YES, otherwise output NO.  
<TOOL DESCRIPTION>: {INPUT}  
<REPORT>: {INPUT}
```

The prompt above hints to the Programmer Agent to determine whether a tool is suitable for this subdomain.

D.3 Inspector Agent's Prompts

```
Imagine you're a data cleansing engineer. You're given a paragraph and asked to determine whether it's more like a semantically rich and more coherent piece of text or more like the grossly incoherent garbage phrase content generated by page components, buttons, recommendations, sidebars, and other machinery. If it's more like normal text, output "MAIN"; if it's more like spammy phrases, output "SPAM". Note that you only need to output "MAIN" or "SPAM".  
<PARAGRAPH>: {INPUT}
```

This prompt is used to enable the Inspector Agent to determine whether each paragraph is qualified. The percentage of qualified characters will be used to determine whether the subdomain is

qualified.

E Case of High PPL with High-quality

In this section, we present an example showing that the cleaned corpus has a higher PPL per token than the raw corpus with low quality. The raw corpus is all the content within the blue rectangle in Figure 9, while the cleaned corpus is shown in Figure 12.

We employ the model from CCNet (Wenzek et al., 2020) to compute the perplexity of the texts. The raw corpus has 2039 tokens with a total perplexity of 255.8, resulting in a PPL per token of 0.1255. While the cleaned corpus has only 1303 tokens and a total perplexity of 217.0, leading to a PPL per token of 0.1665.

It can be observed that the PPL per token of the cleaned corpus is higher than that of the raw corpus. However, by comparing Figure 9 and Figure 12, it is evident that cleaned corpus in Figure 12 contains less low-quality content such as navigation bars in Figure 9.