

# DataFlow: An LLM-Driven Framework for Unified Data Preparation and Workflow Automation in the Era of Data-Centric AI

Hao Liang<sup>\*,†</sup>, Xiaochen Ma<sup>\*,†</sup>, Zhou Liu<sup>\*,†</sup>, Zhen Hao Wong<sup>\*</sup>, Zhengyang Zhao<sup>\*</sup>, Zimo Meng<sup>\*</sup>, Runming He<sup>\*</sup>, Chengyu Shen<sup>\*</sup>, Qifeng Cai<sup>\*</sup>, Zhaoyang Han<sup>\*</sup>, Meiyi Qiang<sup>\*</sup>, Yalin Feng<sup>\*</sup>, Tianyi Bai<sup>\*</sup>, Zewei Pan, Ziyi Guo, Yizhen Jiang, Jingwen Deng, Qijie You, Peichao Lai, Tianyu Guo, Chi Hsu Tsai, Hengyi Feng, Rui Hu, Wenkai Yu, Junbo Niu, Bohan Zeng, Ruichuan An, Lu Ma, Jihao Huang, Yaowei Zheng, Conghui He, Linpeng Tang, Bin Cui, Weinan E, Wentao Zhang<sup>‡</sup>

<sup>1</sup>Peking University, <sup>2</sup>Institute for Advanced Algorithms Research, Shanghai,  
<sup>3</sup>OriginHub Technology, <sup>4</sup>OpenDataLab, Shanghai Artificial Intelligence Laboratory,  
<sup>5</sup>LLaMA-Factory Team

The rapidly growing demand for high-quality data in Large Language Models (LLMs) has intensified the need for scalable, reliable, and semantically rich data preparation pipelines. However, current practices remain dominated by ad-hoc scripts and loosely specified workflows, which lack principled abstractions, hinder reproducibility, and offer limited support for model-in-the-loop data generation. To address these challenges, we present DataFlow, a unified and extensible LLM-driven data preparation framework. DataFlow is designed with system-level abstractions that enable modular, reusable, and composable data transformations, and provides a PyTorch-style pipeline construction API for building debuggable and optimizable dataflows. The framework consists of nearly 200 reusable operators and six domain-general pipelines spanning text, mathematical reasoning, code, Text-to-SQL, agentic RAG, and large-scale knowledge extraction. To further improve usability, we introduce DataFlow-Agent, which automatically translates natural-language specifications into executable pipelines via operator synthesis, pipeline planning, and iterative verification. Across six representative use cases, DataFlow consistently improves downstream LLM performance. Our math, code, and text pipelines outperform curated human datasets and specialized synthetic baselines, achieving up to +3% execution accuracy in Text-to-SQL over SynSQL, +7% average improvements on code benchmarks, and 1–3 point gains on MATH, GSM8K, and AIME. Moreover, a unified 10K-sample dataset produced by DataFlow enables base models to surpass counterparts trained on 1M Infinity-Instruct data. These results demonstrate that DataFlow provides a practical and high-performance substrate for reliable, reproducible, and scalable LLM data preparation, and establishes a system-level foundation for future data-centric AI development.

\*Equal Contribution, †Project Leader, ‡Corresponding author

✉ Correspondence : [wentao.zhang@pku.edu.cn](mailto:wentao.zhang@pku.edu.cn)

🔗 Source Code : <https://github.com/OpenDCAI/DataFlow>

📁 Dataset : <https://huggingface.co/datasets/OpenDCAI/dataflow-instruct-10k>

📖 Codebase Documentation : <https://opendcai.github.io/DataFlow-Doc/>

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Background and Related Works</b>	<b>6</b>
2.1	Data in LLM Development	6
2.2	Data Preparation for LLMs	6
2.3	Existing LLM Data Preparation Systems	6
<b>3</b>	<b>DataFlow System Overview</b>	<b>7</b>
3.1	Goals and Design Philosophy	7
3.2	System Scope and Positioning	7
3.3	System Workflow	8
<b>4</b>	<b>Framework Design and Architecture</b>	<b>8</b>
4.1	Global Storage Abstraction and Operator Interaction	9
4.2	Hierarchical Programming Interfaces	9
4.2.1	LLM Serving API	9
4.2.2	Operator Programming Interface	10
4.2.3	Prompt Template Interface	11
4.2.4	Pipeline Composition Interface	11
4.3	Operator Categorization	12
4.4	DataFlow-Ecosystem	13
<b>5</b>	<b>DataFlow-Agent</b>	<b>14</b>
5.1	AgentRoles	14
5.2	Intelligent Pipeline Recommendation	15
5.3	Summary	15
<b>6</b>	<b>Use Cases &amp; Pipelines</b>	<b>16</b>
6.1	Case Study: Text-to-SQL Data Pipeline in DATAFLOW	16
6.1.1	Operators	16
6.1.2	Pipelines	17
6.1.3	DATAFLOW-support Mechanism	18
<b>7</b>	<b>Experiments</b>	<b>18</b>
7.1	Text Data Preparation	19
7.1.1	Experimental Setting	19
7.1.2	Experimental Results	19
7.2	Math Reasoning Data Preparation	20
7.2.1	Experimental Setting	20
7.2.2	Experimental Results	21
7.3	Code Data Preparation	21
7.3.1	Experimental Setting	21
7.3.2	Experimental Results	22
7.4	Text-to-SQL Data Preparation	22
7.4.1	Experimental Setting	22
7.4.2	Experimental Results	23
7.5	AgenticRAG Data Preparation	24
7.5.1	Experimental Setting	24
7.5.2	Experimental Results	24
7.6	Knowledge Extraction	25
7.6.1	Experimental Setting	25
7.6.2	Experimental Results	26
7.7	Unified Multi-Domain Data Preparation with DataFlow	26

---

7.7.1	Experimental Setting . . . . .	26
7.7.2	Experimental Results . . . . .	27
7.8	Agentic Orchestration . . . . .	28
7.8.1	Experimental Setting . . . . .	28
7.8.2	Experimental Results . . . . .	29
<b>8</b>	<b>Conclusion . . . . .</b>	<b>29</b>
<b>A</b>	<b>Author Contributions . . . . .</b>	<b>35</b>

## 1 Introduction

Large language models (LLMs) have rapidly evolved from research prototypes to foundational infrastructure across natural language processing and beyond. Since OpenAI introduced the GPT [1] family through large-scale human annotation and ignited the era of large language models (LLMs), scaling-law studies [26, 55] have consistently demonstrated that data quality and quantity are central to model performance. As model scales continue to grow and downstream tasks become increasingly complex, the size and semantic diversity of training corpora have expanded dramatically [29, 67]. Modern LLM development now relies on multi-stage, semantics-heavy data preparation pipelines that integrate synthetic, refinement, filtering, and domain-specific transformation across trillions of tokens [6, 47, 61].

However, despite the critical role of high-quality data, data preparation for LLMs remains fragmented and largely unstandardized. Most practitioners still rely on ad-hoc scripts and loosely standardized workflows, which lack explicit dataflow abstractions, well-defined atomic operators, or any form of pipeline-level optimization. The absence of a unified and programmable paradigm makes pipelines difficult to reproduce, extend, or compare across projects [6, 47, 48]. This problem is amplified by the trend toward increasingly fine-grained post-training tasks, such as instruction tuning, chain-of-thought generation, or function calling, where both the semantic richness and the semantic accuracy in data preparation are essential for achieving precise task-level model behavior [62, 72].

In response to this fragmentation, several systems have recently emerged with the goal of standardizing LLM data curation. Frameworks such as NeMo Curator [47] and Data-Juicer [6] offer substantial functionality—including captioning, rewriting, classification, and multimodal processing—and have significantly improved the efficiency of large-scale corpus construction. Yet these systems remain fundamentally extraction- and filtering-oriented, and their abstractions provide limited support for expressing iterative, model-in-the-loop generative workflows with fine-grained semantic control. As a result, they are ill-suited for pipelines in which data synthesis and multi-step semantic refinement are central rather than auxiliary.

This limitation is becoming increasingly consequential. LLMs are no longer only consumers of data, but also producers. Because large-scale human annotation is prohibitively expensive, recent work heavily leverages LLM-based data synthesis workflows to construct high-quality corpora at scale [3]. Multiple recent reports show that, in many regimes, carefully synthesized data can outperform even high-quality selected data [60, 69], further underscoring the importance of LLM-driven generation workflows.

Given these trends, we argue that a unified framework for LLM data preparation must elevate LLM-driven data synthesis to a first-class, programmable dataflow abstraction. Such a framework should: (1) provide fine-grained, composable operators for model-in-the-loop generation and semantic refinement; (2) support explicit, verifiable pipeline definitions that serve as an inspectable, domain-agnostic open-source protocol for LLM data preparation—much like how `torch.nn.Module` standardizes model composition in deep learning; (3) remain backend-agnostic to integrate different LLM engines and storage backends; and (4) enable principled workflow composition, reuse, and optimization across models, tasks, and domains, while further supporting agent-driven automatic workflow construction. Taken together, these requirements signal a shift in data preparation—from post-hoc corpus cleaning toward LLM-centric workflows that build high-fidelity, semantically rich, and task-aligned synthetic corpora through iterative synthesis and refinement.

Motivated by this shift, we introduce DATAFLOW, a unified and automated LLM-driven framework for end-to-end LLM data preparation, with a DATAFLOW-AGENT that allows users to compose pipelines directly from natural-language specifications. DATAFLOW places LLMs at the center of the operator ecosystem: most operators are LLM-driven, with a small number implemented using heuristics or small models. The framework provides over 180 operators organized into four categories—generation, evaluation, filtering, and refinement—and includes more than 90 reusable prompt templates that enable operator-level composition and consistent behavior across tasks. Using these primitives, DATAFLOW includes a set of state-of-the-art (SOTA) synthesis pipelines that span mathematical reasoning, raw text, code, Text-to-SQL, agentic RAG-style data, and large-scale QA extraction from web or PDF corpora. All pipelines are expressed within DATAFLOW’s common abstractions and require no task-specific glue code, adhering to a generate-evaluate-filter workflow augmented with targeted refinement stages.

To ensure usability, extensibility, and long-term maintainability, DATAFLOW adopts a PyTorch-like programming interface that exposes its core abstractions, global storage, LLM serving, operators, prompt templates, and pipelines, through modular Python classes and functions. This code-first design avoids complex YAML or shell-based configuration schemes and provides an IDE-friendly development workflow, including code completion and reliable navigation. Beyond the core library, operators, prompt templates, and pipelines can be developed outside the main repository and packaged as standalone Python modules, enabling practitioners to publish and reuse domain-specific components as first-class DATAFLOW-EXTENSIONS. To support this ecosystem, DATAFLOW includes a Command-Line Interface (CLI) toolchain that scaffolds new extension packages, from operator stubs to full pipeline repositories, standardizing development practices and lowering the barrier to community contribution. Finally, DATAFLOW-AGENT serves as an agentic orchestration layer that translates natural-language specifications into executable pipelines and can automatically synthesize and debug new operators when needed, further accelerating the construction of scalable and semantically rich LLM-driven data preparation workflows.

Extensive experiments on six DATAFLOW-implemented pipelines show that our design philosophy is effective across diverse data preparation scenarios, consistently producing high-quality training data. Across all settings, the resulting DATAFLOW datasets match or even surpass SOTA baselines, including curated human datasets, specialized synthetic workflows, and the strong Qwen2.5-Instruct series. For example, DATAFLOW-synthesized mathematical reasoning data yields 1–3 point gains over high-quality synthetic baselines [28, 43] on MATH, GSM8K, and AIME; our Text-to-SQL pipelines achieve over +3% execution-accuracy improvements compared with the 2.5M-sample SynSQL corpus [37] while using less than 0.1M training examples; and DATAFLOW-based code pipelines deliver over 7% average improvements relative to widely used public code instruction datasets [5, 63].

Moreover, by combining DATAFLOW-generated text, math, and code data into a unified corpus, DATAFLOW-INSTRUCT-10K, we find that training on only 10K samples enables Qwen2-base and Qwen2.5-base to surpass models trained on 1M Infinity-Instruct [39] instances, while approaching the performance of their corresponding Qwen-Instruct model. This demonstrates that DATAFLOW can produce domain-diverse supervision of sufficiently high quality to yield substantial gains in data efficiency.

Together, these results demonstrate that DATAFLOW is not only an end-to-end system for LLM-based data preparation, but also a comprehensive operator and algorithm library and an open, user-friendly protocol framework. Built around six SOTA template pipelines and a large collection of reusable operators, DATAFLOW offers a unified foundation for LLM-centric data construction, enabling principled, semantically rich, and scalable workflows that improve programmability, reproducibility, and data quality across domains.

Overall, our key contributions are summarized as follows:

- **A unified LLM-driven data preparation framework.** We propose DATAFLOW, a unified system for LLM data preparation built on composable abstractions and an LLM-first operator execution model.
- **A rich and extensible operator–pipeline ecosystem.** DATAFLOW provides nearly 200 reusable operators and six SOTA template pipelines covering text, mathematical reasoning, code, Text-to-SQL, agentic RAG data, and large-scale QA extraction.
- **A developer- and open-source-friendly programming model.** Through a PyTorch-like API, IDE-native tooling, and plugin-style extensibility via Python packages, DATAFLOW enables reproducible experimentation, easy customization, and community-driven extensions to form DATAFLOW-ECOSYSTEM.
- **An agentic orchestration layer for automated pipeline construction.** DATAFLOW-AGENT composes executable pipelines from natural-language intent, lowering the barrier to building scalable and semantically rich LLM-driven workflows.
- **Extensive empirical validation and open-source data release.** Experiments across six pipelines show that DATAFLOW-generated data consistently improves downstream LLM performance and data efficiency. We additionally release a high-quality, multi-domain dataset produced entirely with DATAFLOW to support further research and benchmarking.

## 2 Background and Related Works

### 2.1 Data in LLM Development

The development of LLMs involves several key stages, among which training is particularly crucial, as the model learns fundamental linguistic patterns from large-scale corpora. During this stage, the model is exposed to vast amounts of text data from various domains, enabling it to acquire a broad understanding of language.

Consequently, the quality and diversity of training data directly impact the model’s ability to generalize effectively across different contexts [19, 38]. Recently, the rapid development of large language models has brought about a substantial increase in the volume of training data [1, 57]. In this scenario, the quality and quantity of data become even more paramount.

High-quality data can significantly enhance model performance [44]. As the volume of data increases, ensuring high data quality becomes more challenging, as it requires additional resources for data cleaning, selection, and annotation [3]. Poor-quality data can cause models to learn incorrect patterns and produce inaccurate predictions. Furthermore, insufficient data diversity may result in models performing well in specific domains but exhibiting poor generalization in cross-domain tasks. Additionally, distributional shifts in the data can exacerbate model over-reliance on training distributions, diminishing their applicability in real-world scenarios.

### 2.2 Data Preparation for LLMs

As disclosed by the above discussion, data preparation is a crucial step in training LLMs, significantly impacting the model’s performance and generalization capabilities. With the continuous expansion of LLM scales, the complexity and efficiency of data preparation have become key research focuses. However, although systems like Apache Spark [71], Dask [52], and Hadoop [13, 20, 65] are powerful for large-scale Extract–Transform–Load (ETL), they are not a good fit for modern LLM data preparation. These frameworks can, in principle, run semantic cleaning by calling LLMs or embedding models as user-defined functions, but they provide no native support for model-in-the-loop processing, GPU-efficient batching, or token-level text operations. More importantly, their built-in operators focus on structured data and offer very limited functionality for unstructured text, meaning that essential steps—such as tokenization, language detection, document segmentation, semantic deduplication, or safety filtering—must be implemented manually with ad-hoc User-Defined Functions (UDFs). This leads to significant overhead and engineering complexity, making general big-data engines inadequate for the large-scale, semantics-heavy pipelines required for LLM corpus construction.

LLM-based methods have been widely used in data quality evaluation and data selection. For instance, MoDS [15] leverages DeBERTa for scoring and retaining high-quality data, while Alphagassus [7] uses ChatGPT to score data accuracy. Other studies have employed GPT-4 for data rewriting and quality improvement. For a comprehensive overview, refer to the data for LLM survey [3].

### 2.3 Existing LLM Data Preparation Systems

Recent work increasingly approaches LLM training data preparation as a first-class systems problem. Table 1 summarizes the distinguishing characteristics of the major frameworks.

**NeMo Curator** [47] is an open-source, GPU-accelerated library from NVIDIA that offers modular pipelines for large-scale LLM data curation, including data download and extraction (e.g., Common Crawl, arXiv, Wikipedia), language identification, text cleaning, heuristic and learned quality filtering, domain and toxicity classification, document- and semantic-level deduplication, privacy filtering, and even synthetic data generation, all built on Dask/RAPIDS and designed to scale to multi-node, multi-GPU environments.

**Data-Juicer** [6] is a “one-stop” data processing system that abstracts LLM data recipes into composable operators: the original system already provides 50+ operators for constructing and evaluating text data mixtures, while the 2.0 version extends this to 100+ operators across text, image, video, and audio, supporting analysis, cleaning, synthesis, annotation, and post-training data pipelines with tight integration to Ray and HuggingFace Datasets.

**Table 1** High-level comparison of existing data preparation systems for LLM.

Dimension	Data-Juicer [6]	NeMo Curator [47]	DataFlow (ours)
<b>Primary focus</b>	Filtering / Cleaning	Large-scale Curation	LLM-driven Synthesis + Refinement
<b>Programming model</b>	Config-based Recipes	Component-based Pipelines	PyTorch-like Operators & Pipelines
<b>LLM integration</b>	Partial (some gen ops)	Minimal (mainly filtering)	First-class Serving + Templates
<b>Automation</b>	Recommendation Agent	None	Pipeline Construct/Debug Agent
<b>Extensibility</b>	OperatorZoo / Cookbook	Custom Scripts	Extension Packages + CLI Scaffolding

These systems substantially improve the efficiency and quality of LLM data preparation, but they remain largely configuration-centric toolkits. In contrast, our framework is built around a rich library of nearly 200 reusable text-specific operators, enabling fine-grained control over cleaning, transformation, synthesis, and evaluation; multiple pipelines instantiated from these operators consistently yield strong downstream gains, and even simple mixtures of data produced by different pipelines remain highly effective. Moreover, the system adopts a modular, PyTorch-style “building-block” design with lightweight, well-defined interfaces, making it natural for data agents to compose, orchestrate, and invoke data-processing pipelines programmatically.

### 3 DataFlow System Overview

In this section, we present an overview of DATAFLOW, a unified and automated system that standardizes and streamlines multi-domain data preparation for LLMs.

#### 3.1 Goals and Design Philosophy

DataFlow is designed around six core goals:

*Ease of Use.* A PyTorch-inspired [49], IDE-friendly programming interface enables users to build and debug complex data preparation pipelines with minimal boilerplate.

*Extensibility.* Following a modular abstraction similar to `torch.nn.Module`, new operators, and algorithms can be added as plug-and-play components and naturally compose with existing workflows.

*Unified Paradigm.* DATAFLOW unifies heterogeneous data preparation workflows under a standardized abstraction layer. The design balances standardization, for consistency and reproducibility, with customization needed across domains, enabling efficient pipeline reuse and adaptation.

*Performance Efficiency.* The official pipelines in DATAFLOW achieve performance comparable to or exceeding SOTA data preparation methods, demonstrating that unification does not impose substantial overhead.

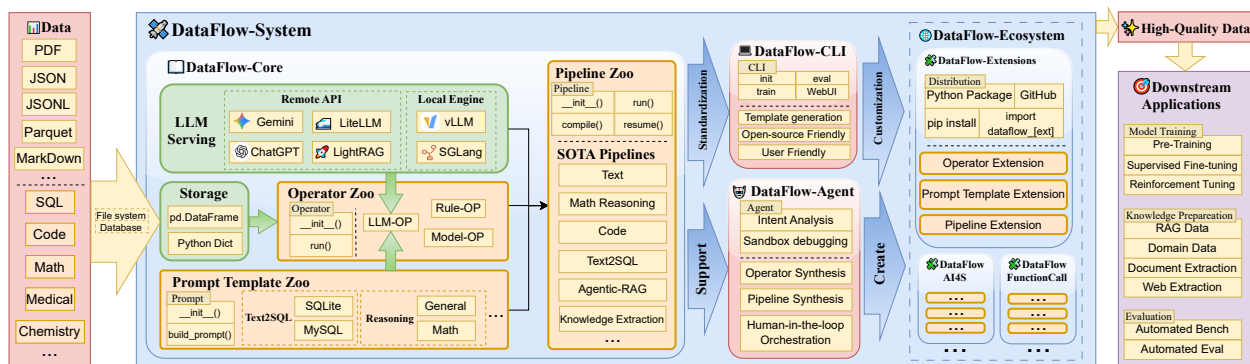
*Intelligent Automation.* A lightweight agentic subsystem leverages core abstractions to interpret natural-language intent and automatically construct or adjust operators and pipelines, supporting rapid prototyping and reducing manual engineering.

*Open Source Paradigm.* DATAFLOW aims to serve as a community standard for LLM data preparation. Its unified abstractions enable reproducible pipeline sharing, transparent swapping of LLM backends, and controlled experimentation.

#### 3.2 System Scope and Positioning

DATAFLOW spans the full workflow of LLM-centric data preparation. As Figure 1 shown, at its core, the system provides unified abstractions for storage, LLM serving, operators, prompt templates, and pipelines—defining the execution substrate on which all transformations are performed. Above the core, two user-facing control layers, the CLI and the DATAFLOW-AGENT, support both scriptable and automated workflow construction.





**Figure 1** High-level architecture of DATAFLOW. The system consists of a core execution engine (storage, operators, templates, and LLM serving), reusable pipelines, user-facing control layers (CLI and agent), and an extensible ecosystem for domain-specialized workflows. DATAFLOW produces high-quality, task-aligned datasets consumed by downstream LLM applications.

Beyond the engine, DATAFLOW-EXTENSIONS offer a modular interface for adding Python-package-based operators, templates, and pipelines. Domain-specialized packages built on this interface collectively form the broader DATAFLOW-ECOSYSTEM. Together, these components define the system boundary: DATAFLOW provides the abstractions and control layers for data preparation, while downstream LLM training, evaluation, and retrieval applications consume its outputs.

### 3.3 System Workflow

Figure 1 also illustrates the end-to-end workflow of DATAFLOW. The system ingests datasets from common file formats (e.g., JSON, JSONL, CSV, Parquet, Markdown, PDF) as well as domain-specific sources such as SQL logs and code repositories, converting all inputs into a unified tabular representation maintained by the core storage layer. Operators interact with this shared storage to read and write intermediate results, enabling consistent data flow across transformation stages.

Operators implement transformations such as generation, refinement, filtering, and evaluation. LLM-driven operators invoke local inference engines (e.g., vLLM [35], SGLang [73]) or online API-based services (e.g., Gemini [56], ChatGPT [1]) via the unified serving abstraction, while rule-based and small-model operators execute independently of LLM backends.

Pipelines in the Pipeline Zoo compose these operators into reusable workflows for tasks such as text synthesis, mathematical reasoning, code processing, Text-to-SQL generation, agentic RAG, and large-scale knowledge extraction. Pipelines may be executed directly, compiled for optimized execution, resumed from intermediate states, or adapted to new domains.

Users interact with DATAFLOW during workflow execution through either the CLI or the DATAFLOW-AGENT: the CLI issues explicit execution commands, while the agent translates natural-language specifications into executable workflows and performs iterative debugging. Workflow outputs, high-quality, task-aligned datasets, integrate seamlessly into downstream LLM applications.

## 4 Framework Design and Architecture

This section presents the internal design of DATAFLOW and formalizes the execution model underlying its abstractions in Section 3. DATAFLOW is organized around four architectural pillars: (1) a global storage abstraction that maintains the canonical tabular representation of datasets and mediates all data access; (2) a set of hierarchical programming interfaces for LLM serving, operators, prompt templates, and pipelines; (3) a principled operator categorization scheme that reconciles open-ended domain requirements with a compact set of reusable transformation primitives; and (4) an extension mechanism that supports a growing ecosystem



of user-contributed components. Together, these elements provide a scalable and extensible substrate for constructing, executing, and sharing LLM-centric data preparation workflows.

## 4.1 Global Storage Abstraction and Operator Interaction

At the core of DATAFLOW’s execution substrate is a unified storage abstraction that maintains the canonical tabular representation of the dataset and mediates all data access during workflow execution. LLM-oriented data—such as instructions, responses, chain-of-thought traces, scores, and metadata—is naturally expressed as key–value fields associated with each sample, making a tabular structure a suitable and expressive organizational format. The storage layer decouples data management from operator logic, exposing a minimal, backend-agnostic API through the `DataFlowStorage` base class. This design allows custom storage backends—such as file-system, object-store, or database implementations—to be integrated without altering operator behavior.

The abstraction provides two primary operations:

- `read()`: retrieve the current dataset (or relevant fields) in a format required by the operator.
- `write(data)`: update or append fields to the shared dataset representation.

Centralizing all access through these operations ensures that operators remain agnostic to physical storage layout, while intermediate artifacts produced by one operator become immediately available to others. A typical operator interaction follows the pattern in Figure 2.

```
def run(self, storage: DataFlowStorage, **kwargs):  
    inputs = storage.read()           # 1. Read input  
    results = operator_transform(inputs, **kwargs) # 2. Transform the data  
    storage.write(results)           # 3. Write output
```

**Figure 2** The standard execution pattern of an operator’s `run()` method in DATAFLOW. Within `run()`, the operator interacts with the global `DataFlowStorage` by retrieving inputs through `storage.read()`, applying its transformation logic, and writing updated fields back via `storage.write()`. This read–transform–write paradigm captures how data flows from one operator to the next throughout the workflow.

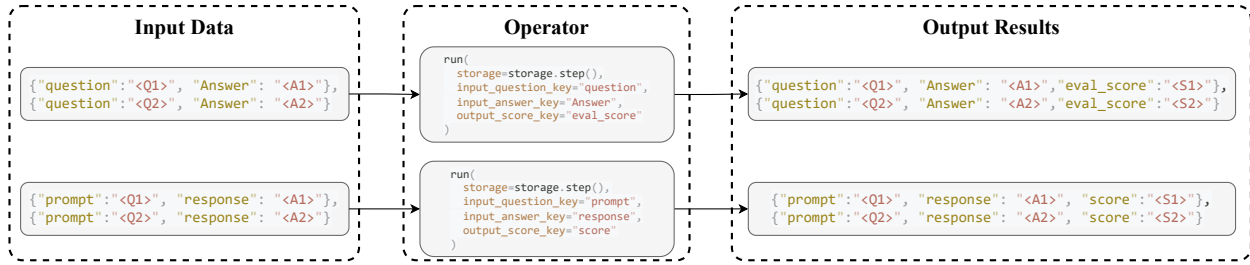
Because operators operate only against this logical abstraction, they can be reordered, recomposed, or batched without modifying their internals, and improvements to the storage backend (e.g., adding distributed or database-backed implementations) require no operator-level changes. The default storage implementation uses a Pandas as the execution substrate and supports common input/output formats such as JSON, JSONL, CSV, and Parquet.

## 4.2 Hierarchical Programming Interfaces

DATAFLOW exposes a hierarchical programming interface built around four core abstractions. (1) The serving interface provides a unified mechanism for issuing LLM inference requests across heterogeneous backends. (2) Operators define reusable data-transformation units and may optionally invoke the serving layer when LLM-driven computation is required. (3) Prompt templates specify how operator inputs are rendered into concrete prompts and how model outputs should be structured or constrained, providing a declarative interface for consistent prompt construction. (4) Pipelines compose operators into multi-stage workflows with explicit data dependencies and support optional compilation for validation and optimization. The following subsections describe these abstractions in detail.

### 4.2.1 LLM Serving API

LLM-driven operators rely on a unified serving API that abstracts over heterogeneous model backends. The API exposes a single high-level entry point, `generate_from_input(user_inputs, system_prompt, json_schema)`, which accepts a list of prompts, typically assembled by the calling operator, and returns a list of model-generated outputs. Optional arguments such as a `system_prompt` or an output `json_schema` enable



**Figure 3** Example of how an operator’s `run()` method interacts with data via key-based bindings. This flexible key-binding mechanism adapts to arbitrary datasets without preprocessing and enables seamless operator composition.

structured prompting and decoding when needed. This interface shields operators from backend-specific considerations such as batching, retry strategies, request routing, and rate limiting.

The serving layer supports both:

- Local inference engines (e.g., vLLM [35], SGLang [73]), which exploit backend-level parallelism for high-throughput execution; and
- Online API-based services (e.g., ChatGPT [1], Gemini [56]), for which DATAFLOW performs multi-threaded request dispatch to maximize throughput.

This unified serving abstraction reduces the implementation burden of LLM-driven operators and enables flexible backend substitution, making it easy to assess how different LLM choices influence data preparation quality.

#### 4.2.2 Operator Programming Interface

Operators serve as the fundamental transformation units in DATAFLOW. They follow a two-phase interface that cleanly separates initialization from execution: initialization configures the operator, while execution performs the transformation. This separation allows heterogeneous behaviors, from LLM-driven generation to rule-based filtering, to be expressed under a unified abstraction.

During initialization (`__init__()`), an operator receives configuration parameters such as hyperparameters or task-specific settings. LLM-driven operators may additionally bind to a LLM serving object and a prompt-template object in this stage, whereas rule-based and lightweight-model operators omit these bindings entirely. Initialization therefore captures all static configuration and external dependencies, leaving execution to focus exclusively on data transformation.

An operator’s `run()` method implements its transformation logic and constitutes the unit of execution within a pipeline. To keep operators general and easily composable, `run()` accepts only a `DataFlowStorage` object together with a set of `input_*` and `output_*` keys. Interpreting these as key-value pairs, an `input_*` key indicates the storage column to be read as an input field, while an `output_*` key indicates the name of the new column to be written for each processed data item. Figure 3 illustrates this mapping. This design provides flexible I/O bindings that naturally adapt to diverse upstream datasets, while the declared keys form a directed dependency graph among operators, enabling topological scheduling and downstream optimization checks.

By isolating configuration from execution and constraining state changes to explicit key-based read/write operations on shared storage, the operator abstraction remains lightweight, deterministic, and easy to compose. These properties allow DATAFLOW to support a wide range of transformation behaviors under a single, portable interface while preserving consistent execution semantics throughout the system.

```

class TranslatePipeline(PipelineABC):
    def __init__(self):
        super().__init__()
        # Init Resources
        self.storage = FileStorage(
            entry_file="input_data.jsonl",
        )
        self.llm_serving = APILLMServing(
            api_url="<api_url>",
            model_name="gpt-4o",
        )
        # Initialize Operators
        self.op1 = PromptedGenerator(
            llm_serving=self.llm_serving,
            system_prompt="Translate the content to Chinese",
        )
        self.op2 = PromptedGenerator(
            llm_serving=self.llm_serving,
            system_prompt="Translate the content to English",
        )

    # forward function of the Pipeline
    def forward(self):
        # execute operators
        self.op1.run(
            self.storage.step(),
            input_key='raw_content',
            output_key='content_CN'
        )
        self.op2.run(
            self.storage.step(),
            input_key='raw_content',
            output_key='content_EN'
        )

if __name__ == "__main__":
    TransPipeline = TranslatePipeline()
    Transpipeline.compile() # Optional
    # excute pipeline, resume from `op2`
    Transpipeline.forward(resume_step=1)

```

**Figure 4** Illustration of the DATAFLOW pipeline API. The example shows how a pipeline declares its storage and serving backends, instantiates operators with task-specific configurations, and executes them via `forward()` using input/output key bindings. The interface supports compilation and stepwise resumption, enabling flexible and modular workflow construction.

#### 4.2.3 Prompt Template Interface

Prompts serve as the primary mechanism guiding LLMs to perform task-specific transformations. Every LLM-driven operator relies on a prompt, and operators that share the same high-level logic often differ only in subtle prompt variations. For instance, in Text-to-SQL generation, synthesizing queries for SQLite and MySQL involves identical operator logic; the only difference lies in minor syntax adjustments communicated through the prompt. To support such reuse while accommodating domain-specific variations, DATAFLOW decouples prompt construction from operator implementation through a dedicated prompt template interface.

A prompt template encapsulates a reusable prompt pattern and provides parameterized slots that operators populate at execution time. Each LLM-driven operator initializes its associated template during `__init__()`, following the same configuration-execution paradigm as other system components. During execution, the operator invokes the template's `build_prompt()` method, which assembles task-relevant information—such as input fields, schema hints, or contextual metadata—into a concrete prompt that is subsequently passed to the LLM serving layer. This encapsulation allows the operator's transformation logic to remain agnostic to how prompts are rendered.

To facilitate one-to-many mappings between operators and templates, LLM-driven operators expose a unified `op.ALLOWED_PROMPTS` interface that enumerates all compatible prompt templates. This design enables operators to be flexibly reused across domains or tasks by simply switching or tuning templates, without modifying operator logic.

Overall, the prompt template interface provides a declarative mechanism for prompt construction, promotes operator reuse across closely related tasks, and ensures consistent prompting behavior throughout DATAFLOW's LLM-driven workflows.

#### 4.2.4 Pipeline Composition Interface

Building on the abstractions introduced above, DATAFLOW provides a pipeline interface that enables users to compose operators into multi-stage data-preparation workflows. A pipeline is represented as an ordered sequence of operators (or a lightweight DAG), forming an end-to-end execution graph that captures the intended dataflow. Figure 4 illustrates the pipeline API and its core components.

The pipeline API adopts a PyTorch [49]-like design in which the `__init__()` method handles resource allocation and operator configuration, while the `()` method encodes a single pass of execution. Within

`forward()`, operator-specific key bindings implicitly define the dataflow topology, allowing pipelines to be constructed in a modular, readable, and IDE-friendly manner.

Functionally, the pipeline interface provides a built-in `compile()` procedure that performs static analysis of the operator sequence prior to execution. During compilation, DATAFLOW extracts operator dependencies and parameters, constructs the corresponding DAG, and conducts key-level validation to detect missing fields, type inconsistencies, and malformed dependency chains. Instead of executing operators immediately, `compile()` records all operator configurations and dependency information to produce a deferred execution plan. This deferred-construction design follows the Factory Method pattern [16], in which object creation is separated from object execution: the actual invocation of each operator’s `run()` method is deferred until the subsequent `forward()` call.

The compiled execution graph first provides complete structural information to the DATAFLOW-Agent, enabling it to surface all key- and dependency-related errors in a single report. This significantly reduces the number of debugging rounds required by the agent and lowers the associated inference cost. Additionally, the compiled graph defines a minimal and efficient execution plan that supports advanced runtime features such as checkpointing and stepwise resumption, improving iterative development and large-scale pipeline construction.

### 4.3 Operator Categorization

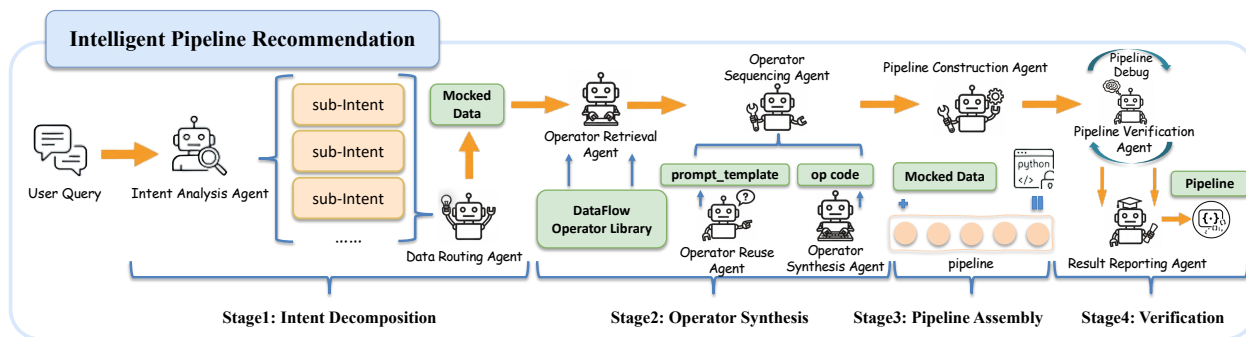
Operators in DATAFLOW encapsulate diverse data-processing algorithms that, when composed, support end-to-end LLM data preparation workflows. As a unified yet extensible framework intended to serve arbitrarily many domains, DATAFLOW must simultaneously accommodate an open-ended set of domain-specific algorithms while exposing a stable and comprehensible operator space. These competing forces—unbounded domain requirements and the need for conceptual compactness—introduce inherent tension. To reconcile this, DATAFLOW organizes operators along multiple orthogonal categorization dimensions. Categories are mutually exclusive within each dimension, while dimensions themselves are parallel. This categorization scheme has been validated across the diverse domains covered in this paper, including more than six state-of-the-art data preparation pipelines, demonstrating both its representational sufficiency and scalable generality.

*Modality Dimension.* The fundamental categorization separates operators by the modality they process, such as text, visual content, or document-like inputs. Modalities must be distinguished because operators within the same modality share compatible input–output semantics and can interoperate, whereas operators across different modalities often cannot be composed directly. DATAFLOW primarily operates on textual representations, with non-text modalities first processed by modality-specific operators that parse or convert raw inputs—such as images or PDFs—into text before any downstream transformations are applied. Therefore, Clear modality classification makes this conversion flow explicit and enables the pipeline compiler to validate operator chains, ensuring that modality transitions are correctly specified and that only compatible operators are composed.

*Core vs. Domain-Specific Dimension.* A second categorization distinguishes between core operators and domain operators. Core operators reflect the fundamental design philosophy of DATAFLOW and serve as the conceptual basis from which most other operators can be derived. Although domain operators may wrap or specialize core operators, their semantics can generally be expressed by instantiating the parameters of a corresponding core operator. Core operators are intentionally limited in number and relatively stable, forming the recommended entry point for new users. Domain operators, by contrast, expand without bound as new domains, modalities, or tasks emerge. Although theoretically unbounded, the domain operators included in DATAFLOW are limited to those required to support the best-performing pipelines across existing domains, ensuring practical conciseness and avoiding unnecessary proliferation.

*Functional Dimension.* At a finer granularity, operators fall into four functional categories—generate, evaluate, filter, and refine—each capturing a distinct transformation pattern in data preparation. These categories align with a core design philosophy of DATAFLOW as a data-synthesis framework: pipelines first expand the candidate space through generation, then score and filter the results, optionally applying refinement stages in





**Figure 6** DATAFLOW-AGENT architecture: a LangGraph-orchestrated multi-agent workflow that translates natural-language intent into a verified executable DAG pipeline.

package that encapsulates additional operators, prompt templates, and pipelines. User-contributed extensions collectively form the broader DATAFLOW-Ecosystem, a plug-and-play environment analogous to Python’s package ecosystem, where practitioners can readily publish, share, and reuse domain-specific components.

To streamline extension development, DATAFLOW provides automated project scaffolding through the DATAFLOW-CLI. Given a few high-level specifications, the CLI generates ready-to-use templates for operators, prompt templates, pipelines, and even full repository layouts suitable for distribution via PyPI or GitHub. Developers need only implement task-specific logic within these generated stubs. Both the core system and extension packages can be installed and imported through Python’s package manager, while lazy-loading mechanisms ensure that multiple extensions coexist with minimal environmental interference.

Complementing the CLI, the DATAFLOW-AGENT supports natural-language-driven construction of operators and pipelines. Leveraging the domain knowledge embedded in large language models, the agent synthesizes effective data-transformation logic and automates common design steps, substantially reducing the cost of authoring high-quality DATAFLOW-Extensions.

Together, the DATAFLOW-CLI and DATAFLOW-AGENT reduce the overhead of extension development and promote community-driven growth. Our goal is to cultivate a sustainable open-source ecosystem in which data preparation recipes—constructed from standardized operators, prompt templates, and pipelines—can be shared, reproduced, and improved, ultimately accelerating progress across the data-centric ML community.

## 5 DataFlow-Agent

The DATAFLOW-AGENT serves as the intelligent orchestration layer atop the DATAFLOW framework. It bridges high-level human intent with low-level data-processing execution by leveraging the modular abstractions of DataFlow together with a graph-based multi-agent workflow engine. Built on LangGraph [2], the agent layer coordinates a set of specialized agents through a stateful execution graph, translating natural-language directives into executable, self-correcting, and optimized data preparation pipelines.

### 5.1 AgentRoles

To achieve autonomous pipeline construction and code synthesis, the system decomposes responsibilities across a roster of specialized agents. Each agent encapsulates specific logic and interacts with the DATAFLOW core components:

- **Intent Analysis Agent:** Accepts the user’s high-level natural language query and decomposes it into a structured sequence of actionable sub-intents, providing the foundational blueprint for the pipeline.
- **Data Routing Agent:** Analyzes the provided input data to determine the task category for routing, or generates synthetic data placeholders if no data is supplied to enable dry-run execution.



- **Operator Retrieval Agent:** Takes specific sub-intents as input and employs RAG to retrieve the most relevant existing operators from the DATAFLOW library as potential candidates.
- **Operator Sequencing Agent:** Evaluates candidate operators for I/O compatibility to select the best fit, or outputs detailed specifications for new operators when functional gaps are detected.
- **Operator Synthesis Agent:** Receives specifications for missing functions and generates context-aware code using RAG, performing automated unit-level debugging until the code is executable.
- **Operator Reuse Agent:** Assesses the generated operator code for quality and creates a reusable `prompt_template`, ensuring the code can be efficiently reused without rewriting.
- **Pipeline Construction Agent:** Orchestrates the assembly of all validated operators (both pre-existing and newly synthesized) into a coherent Directed Acyclic Graph (DAG) structure ready for processing.
- **Pipeline Verification Agent:** Executes the assembled pipeline within a sandboxed environment to identify runtime errors, autonomously adjusting connections or parameters to output a validated, error-free pipeline.
- **Result Reporting Agent:** Synthesizes the final workflow details and execution results, generating a comprehensive report and an executable pipeline artifact as the final solution.

## 5.2 Intelligent Pipeline Recommendation

As shown in Figure 6, the core capabilities of the system are realized through a sophisticated agentic layer built atop the DataFlow framework. This layer employs LangGraph [2] to orchestrate a series of specialized agents within graph-based stateful workflows.

*Intent Decomposition* The workflow begins when the system receives a user’s natural language query. The Intent Analysis Agent decomposes this high-level objective into a sequence of discrete, actionable sub-intents. Concurrently, the Data Routing Agent evaluates the input dataset to categorize the task for downstream routing. If no dataset is provided, this agent generates synthetic data placeholders to enable a complete dry-run execution.

*Operator Synthesis* To fulfill these sub-intents, the Operator Retrieval Agent searches the DATAFLOW library for relevant operators, which the Operator Sequencing Agent evaluates for compatibility. If a functional gap is identified, the Operator Reuse Agent first assesses whether the requirement can be met by reusing existing code via a `prompt_template`. Only when reuse is not feasible does the Operator Synthesis Agent generate new code using RAG-based few-shot learning. The code is then debugged automatically to ensure stable execution.

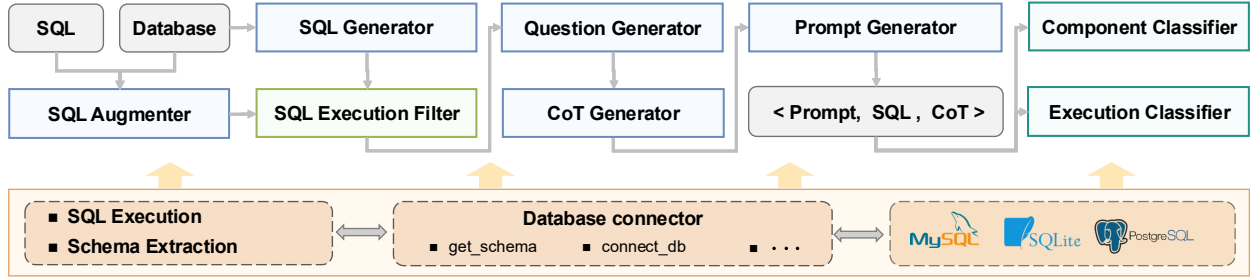
*Pipeline Assembly* After all retrieved or synthesized operators are validated, the Pipeline Construction Agent assembles them into a single pipeline. It represents the pipeline as a DAG and defines the initial connections so data can flow from the source to the sink.

*Verification* The system then runs an integration test. The Pipeline Verification Agent executes the pipeline in a sandbox with a data sample to check connectivity and runtime behavior. If errors occur, it fixes them by adjusting parameters or connections. After the pipeline passes validation, the Result Reporting Agent generates a report and outputs the final executable pipeline definition.

## 5.3 Summary

In summary, unlike Data-Juicer’s agentic approach [6], which is largely constrained to parameterizing and sequencing a static library of pre-existing operators, DATAFLOW-AGENT achieves a significantly higher degree of autonomy through its ability to dynamically synthesize and debug executable code for missing functionalities. By integrating a "retrieve-reuse-synthesize" strategy with a self-correcting verification loop,





**Figure 7** Overall framework of Text-to-SQL pipelines in DATAFLOW.

our system transcends simple configuration generation, enabling the construction of truly adaptive pipelines that can handle unforeseen requirements without manual coding intervention.

## 6 Use Cases & Pipelines

DATAFLOW integrates a rich collection of data pipelines covering diverse text centric task domains, including text processing, mathematical reasoning data, Text-to-SQL generation, and agentic data preparation. In addition, DATAFLOW supports structured knowledge extraction and normalization from PDFs and textbooks, enabling tasks such as schema construction, domain grounding, and instruction synthesis.

All pipelines are implemented through reusable operators and declarative workflow specifications, allowing users to flexibly compose, extend, and adapt them to new scenarios with minimal engineering effort. More detailed tutorials, pipeline examples, and operator-level documentation are available at the website: <https://opendcai.github.io/DataFlow-Doc/>.

### 6.1 Case Study: Text-to-SQL Data Pipeline in DataFlow

We first design a set of specifically designed, reusable Text-to-SQL operators to ensure modularity and extensibility (see Section 6.1.1). As shown in Figure 7, we introduce two pipelines to construct high-quality Text-to-SQL datasets (see Section 6.1.2). Furthermore, Section 6.1.3 describes the support for database operations and the prompt template mechanisms provided by DATAFLOW.

#### 6.1.1 Operators

*SQL Generator.* The SQL Generator operator produces SQL queries from scratch using the database, ensuring both diversity and validity. Four levels of complexity, simple, moderate, complex, and highly complex, are defined and randomly selected to guide the LLM in generating queries of varying difficulty through clear definitions and few-shot examples. The database schema, including `CREATE TABLE` statements for all relational tables and randomly sampled column values, provides the necessary context for the LLM to understand the database. Advanced SQL functions are also randomly supplied to increase the realism of the generated queries. Since natural language questions often require querying specific entries, the number of returned columns is constrained accordingly. Under task instructions, the LLM produces meaningful SQL queries. Within the DATAFLOW framework, the SQL Generator operator can be naturally adapted and reused across different databases (e.g., MySQL, SQLite, PostgreSQL) simply by replacing the corresponding prompt template.

*SQL Augmentor.* The SQL Augmentor operator generates diverse, closely related augmented SQL queries based on seed SQL rather than synthesizing them from scratch. We propose six augmentation strategies to expand SQL queries in different directions: (1) Data Value Transformation, (2) Query Structure Modification, (3) Business Logic Alteration, (4) Complexity Enhancement, (5) Introduction of Advanced SQL Features, and (6) Performance and Optimization. Categories are randomly selected and applied through few-shot prompting. The database schema and values are provided as contextual information. Given an original SQL query and task instructions, the augmentor produces its augmented SQL counterpart.

*Text2SQL Consistency Filter.* For existing pairs of natural language questions and SQL queries, inconsistencies may arise where the two do not correspond. Such problematic data needs to be filtered out. This is achieved using an LLM, which analyzes whether the question and SQL align in content.

*SQL Execution Filter.* Not all generated SQL queries are valid or efficient. Therefore, the SQL Execution Filter operator filters queries from two perspectives: (1) whether the SQL query can be successfully executed on the target database, and (2) whether its runtime exceeds a preset threshold, in which case it is discarded to ensure system responsiveness.

*Question Generator.* The Question Generator operator generates a semantically equivalent natural language question based on the SQL. Natural language questions are categorized into the following stylistic types: (1) Tone and Formality: formal vs. colloquial, (2) Syntactic Structure and Intent: imperative, interrogative, and declarative, (3) Information Density and Clarity: concise, descriptive, ambiguous, and metaphorical, and (4) Interaction Mode: role-playing and procedural. The first two categories cover queries with clear user intent, whereas ambiguous and metaphorical styles involve unclear or figurative language. A target language style is randomly selected, and the database schema is provided for context. Based on task instructions and the generated SQL query, the LLM produces a natural language question.

*Chain-of-Thought Generator.* Chain-of-Thought(CoT) reasoning enhances a model’s ability to solve complex tasks by breaking them down into a series of smaller, manageable sub-problems. To generate CoT reasoning traces, the task instructions, database schema, the generated natural language question, and the generated SQL query are needed. The LLM produces a complete reasoning chain covering intermediate reasoning steps and the final SQL query. During CoT validation, the generated SQL is extracted from the reasoning chain. A CoT process is considered a valid solution only if the execution result of its generated SQL matches that of the reference SQL on the given database.

*Prompt Generator.* As the primary input to the model, a prompt contains the necessary information for reasoning. To facilitate reliable Text-to-SQL generation, a well-structured prompt should include not only the natural language question but also the database schema and specific task instructions to guide the model. The Prompt Generation operator synthesizes these components into a final prompt.

*SQL Component Classifier.* Classifying SQL queries enables deeper analysis of their structural complexity. Following the evaluation standards of Spider [70], SQL queries are categorized into four difficulty levels, simple, moderate, hard, and extra hard, based on the number and complexity of their syntactic components. These components include column selections, the use of aggregate functions in the SELECT clause, and advanced constructs such as GROUP BY, ORDER BY, INTERSECT, or nested subqueries. The SQL Component Classifier operator assigns each SQL query to one of these categories according to the defined criteria.

*SQL Execution Classifier.* Whether the model can generate correct SQL for a given natural language question is also a meaningful measure of difficulty. In the SQL Execution Classifier operator, LLM is instructed to generate SQL query  $k$  times on the same input prompt and count the number of successful executions, denoted as  $n$ . We then classify the difficulty level based on  $\frac{n}{k}$ . Unlike the SQL component classifier operator, execution difficulty is model-dependent: more capable LLMs achieve higher success rates on the same task and thus are considered to have lower execution difficulty.

## 6.1.2 Pipelines

In the design philosophy of DATAFLOW, pipelines are decomposed into independent operator units according to their functionalities, enabling maximal reusability of operators. As shown in Figure 7, the designed operators are composed into two pipelines to support SQL data synthesis in different scenarios.

**Table 2** Pre-training Data Filtering: Performance comparison across models trained with 30B-scale tokens on general evaluation benchmarks.

Methods	ARC-C	ARC-E	MMLU	HellaSwag	WinoGrande	Gaokao-MathQA	Avg
<b>Random-30B</b>	25.26	43.94	27.03	37.02	50.99	27.35	35.26
<b>Qurating-30B</b>	25.00	43.14	27.50	37.03	50.67	26.78	35.02
<b>FineWeb-Edu-30B</b>	26.45	45.41	27.41	38.06	50.43	25.64	35.57
<b>DataFlow-30B</b>	25.51	45.58	27.42	37.58	50.67	27.35	<b>35.69</b>

*SQL Generation Pipeline.* This pipeline generates SQL from scratch based on the database schema. It first uses the SQL Generator operator to produce initial SQL statements, followed by the SQL Execution Filter to remove low-quality or non-executable SQL. Next, the Question Generator produces the natural language question corresponding to each SQL query, the Chain-of-Thought Generator operator generates the reasoning steps (CoT), and the Prompt Generator constructs the prompt content. Finally, the SQL Component Classifier and SQL Execution Classifier assign difficulty labels to the data.

*SQL Refinement Pipeline.* This pipeline generates data starting from the existing seed SQLs. The pipeline first verifies the quality of the seed SQL using the SQL Execution Filter, and the Text2SQL Consistency Filter removes samples where the SQL does not align with the natural language question. Then, the SQL Augmentor produces augmented SQL based on the seed SQL. The subsequent steps mirror those in the SQL Generation Pipeline: filtering low-quality SQL with the SQL Execution Filter, generating natural language questions via the Question Generator, producing CoT reasoning via the Chain-of-Thought Generator, composing prompts with the Prompt Generator, and finally assigning difficulty labels using the SQL Component Classifier and SQL Execution Classifier.

### 6.1.3 DataFlow-support Mechanism

*Database Manager Module.* Within the Pipeline, an efficient and reliable data interaction mechanism serves as the core infrastructure that ensures the stable execution of the workflow. To this end, we implement the **Database Manager** module, which encapsulates the low-level details of database interaction and provides a unified, efficient, and extensible programming interface. The **Database Manager** improves processing throughput under high-concurrency workloads and abstracts schema metadata retrieval, thereby reducing the upper layers’ dependency on the underlying database structure. To achieve cross-database compatibility, we introduce the abstract base class **DatabaseConnector**. This class defines a standardized set of interfaces, including `connect_db` (establishing a database connection), `execute_sql` (executing SQL statements and returning results), and `get_schema` (retrieving complete schema metadata). For each database system, developers need only subclass this base class and implement the system-specific driver invocation and error-handling logic, enabling seamless integration into the overall system.

*Prompt Template Module.* When generating SQL, different scenarios, such as CRUD queries, vector search SQL, or SQL categorized by different difficulty specifications, require distinct prompt templates. To maximize operator reusability under these varying requirements, DATAFLOW introduces the Prompt Template module. This design allows the SQL Generator operator to be reused across scenarios by simply substituting the Prompt class. In practice, one only needs to reimplement the `build_prompt` method within a new Prompt class, without modifying the SQL Generator operator itself.

## 7 Experiments

In this section, we present a comprehensive set of experiments spanning text, math, and code data preparation, as well as Text-to-SQL and AgenticRAG workflows constructed using DATAFLOW. Except for the AgenticRAG setting, which is trained using the Recall [9, 54] framework, all other experiments are conducted using the LLaMA-Factory [74] training framework. We further integrate these modalities to assess the model’s general instruction-tuning performance across diverse tasks.

**Table 3** SFT Data Filtering: Comparison of different 5k dataset filtering methods across Math, Code, and Knowledge benchmarks.

Methods	Math						Code			Knowledge		
	math	gsm8k	aime24	minerva	olympiad	Avg	HumanEval	MBPP	Avg	MMLU	C-EVAL	Avg
<b>Alpaca(random)</b>	54.9	77.2	13.3	14.0	27.0	37.3	71.3	75.9	73.6	71.8	80.0	75.9
<b>Alpaca(filtered)</b>	60.3	80.0	13.3	14.7	30.7	39.8	73.8	75.7	74.8	71.8	80.0	75.9
<b>WizardLM(random)</b>	61.1	84.2	6.7	18.0	29.3	39.9	75.6	82.0	<b>78.8</b>	71.8	79.2	75.5
<b>WizardLM(filtered)</b>	69.7	88.8	10.0	19.9	35.4	44.8	77.4	80.4	<b>78.9</b>	71.9	79.6	75.8
<b>DataFlow-SFT-15K(random)</b>	72.6	89.6	13.3	37.9	32.9	<b>49.3</b>	79.9	75.9	77.9	72.1	80.0	<b>76.1</b>
<b>DataFlow-SFT-15K(filtered)</b>	73.3	90.2	13.3	36.0	35.9	<b>49.7</b>	82.9	74.9	<b>78.9</b>	72.2	80.4	<b>76.3</b>

## 7.1 Text Data Preparation

### 7.1.1 Experimental Setting

We evaluate the impact of high-quality text data preparation on both pre-training (PT) and supervised fine-tuning (SFT) using our DATAFLOW system. Our experiments cover three complementary scenarios:

(1) *Pre-training Data Filtering (30B Scale)*. From the SlimPajama-627B corpus, we extract a 100B-token subset and apply multiple DATAFLOW text-pretraining filters (implemented in `dataflow/operators/text_pt/filter`). For each filter, the top 30% (approximately 30B tokens) is selected. We train a Qwen2.5-0.5B model from scratch for 30B tokens using the Megatron-DeepSpeed framework. We compare four settings:

- **Random-30B**: a random 30B-token subset.
- **FineWeb-Edu-30B**: educational filtering based on FineWeb-Edu [50].
- **Qurating-30B**: Qurating filters [64] using thresholds: `educational_value`  $\geq 7.5$ , `facts_and_trivia`  $\geq 4.0$ , `required_expertise`  $\geq 5.0$ , `writing_style`  $\geq 1.0$ .
- **DataFlow-30B**: intersection of all DATAFLOW PT filters selecting the top 30%.

(2) *SFT Data Filtering (5K Scale)*. To study small-scale SFT data quality, we fine-tune the Qwen2.5-7B base model using LLaMA-Factory on WizardLM and Alpaca datasets. For each dataset, we compared a randomly sampled set of 5K instances against a set of 5K instances filtered by DATAFLOW’s SFT pipeline. Additionally, we synthesize a 15k-size dataset, DATAFLOW-SFT-15K, using DATAFLOW’s Condor Generator and Condor Refiner pipeline, followed by DATAFLOW’s SFT filtering pipeline (excluding the Instagram filter). Benchmarks include comprehensive Math, Code, and Knowledge evaluation suites.

(3) *Conversation-Domain Synthesis (15K Scale)*. We synthesize DATAFLOW-CHAT-15K using DATAFLOW’s conversation-generation pipeline and fine-tune Qwen2.5-7B-Base on it. Baselines include ShareGPT-15K, UltraChat-15K, and their full (non-truncated) versions. We evaluate on domain-specific tasks (TopDial, Light) and general benchmarks (MMLU [23], AlpacaEval [42], Arena-Hard [41]).

### 7.1.2 Experimental Results

*Pre-training* First, from Table 2, we can see across six general benchmarks (ARC-C/E, MMLU, HellaSwag, WinoGrande, Gaokao-MathQA), the DATAFLOW method achieves the highest average score (35.69), outperforming Random (35.26), FineWeb-Edu (35.57), and Qurating (35.02). Despite using the same 30B token budget, DATAFLOW’s multi-filter intersection produces a cleaner and more semantically consistent dataset, leading to better generalization for a 0.5B-scale Qwen2.5 model trained from scratch.

*SFT* In Table 3, we then evaluate 5K-scale SFT data filtering using Alpaca, WizardLM, and DATAFLOW synthetic data. For all three sources, DATAFLOW’s filtering pipeline consistently improves performance over random sampling across Math, Code, and Knowledge benchmarks. At the same time, the results also show that the DATAFLOW-constructed SFT corpus is inherently stronger than Alpaca and WizardLM: even without

**Table 4** Conversation Synthesis: Performance comparison on conversation-domain datasets and general benchmarks for Qwen2.5-7B under different 15K SFT data sources.

Model	Conversation Benchmarks			General Benchmarks			
	TopDial	Light	Avg	MMLU	AlpacaEval	Arena-Hard	Avg
<b>Qwen2.5-7B</b>	7.71	7.79	7.75	71.45	7.05	0.60	26.36
<b>+ ShareGPT-15K</b>	7.75	6.72	7.24	73.09	3.70	1.30	26.03
<b>+ UltraChat-15K</b>	7.72	6.83	7.28	72.97	3.97	0.80	25.91
<b>+ DataFlow-Chat-15K</b>	7.98	8.10	8.04	73.41	10.11	1.10	<b>28.21</b>

filtering, DATAFLOW-SFT-15K achieves higher Math averages (49.3) than the filtered variants of Alpaca (39.8) and WizardLM (44.8), and remains competitive on Code and Knowledge. Moreover, the smaller performance gap between the random and filtered versions of DATAFLOW-SFT-15K (49.3→49.7) further suggests that DATAFLOW-synthesized data is already cleaner and more informative, requiring less aggressive filtering to reach peak performance.

*Conversation* Finally, from Table 4 we can see DATAFLOW-CHAT-15K boosts the overall general benchmark mean from 26.36 to 28.21 and improves AlpacaEval from 7.05 to 10.11, outperforming ShareGPT and UltraChat.

These findings demonstrate that high-quality synthetic data, when paired with DATAFLOW’s refinement and filtering stack, can surpass commonly used human-collected instruction datasets.

## 7.2 Math Reasoning Data Preparation

### 7.2.1 Experimental Setting

We construct a high-quality synthetic mathematical reasoning dataset based on the DATAFLOW Reasoning Pipeline, with adaptations tailored for large-scale reasoning generation. Our goal is to compare three training sources: (1) a random 10K subset from Open-R1 [28], (2) a random 10K subset from Synthetic-1 [43], and (3) our 10K synthesized DATAFLOW-REASONING-10K dataset constructed using DATAFLOW.

*Data Synthesis Method.* The data generation process follows the core structure of the DATAFLOW Reasoning Pipeline and includes three stages:

- **Problem Synthesis.** We adopt the NuminaMath dataset as a high-quality seed set and utilize the o4-mini model together with DATAFLOW’s math problem synthesis operators to expand it into a diverse candidate problem pool.
- **Quality Verification.** All candidate problems are validated using DATAFLOW’s MathQ-Verify [53] module, which detects incorrect, ambiguous, or logically inconsistent problems. Low-quality samples are removed to ensure correctness and robustness.
- **Chain-of-Thought (CoT) Generation.** For all verified problems, we employ DATAFLOW’s CoT-generation operators to prompt DeepSeek-R1 to produce complete, step-by-step reasoning traces.

Compared with the original Reasoning Pipeline, we omit the seed-level pre-verification stage, because NuminaMath is already a curated and validated dataset. This reduces computational overhead while maintaining overall data reliability.

We evaluate Qwen2.5-32B-Instruct fine-tuned on different 10k synthetic datasets across eight mathematical benchmarks, including GSM8K [11], MATH [24], AMC23, Olympiad, Gaokao24-Mix, Minerva, and AIME 2024/2025. Table 5 reports the full results.

*Generation Hyperparameters.* For non-AIME problems, we use `temperature = 0` and `top-p = 0.95`. For AIME-style problems, we adopt a more exploratory sampling strategy with `temperature = 0.6`, `top-p =`

**Table 5** Math Reasoning Pipeline: Performance comparison of Qwen2.5-32B-Instruct under different synthetic data training settings.

Model	gsm8k	math	amc23	olympiad	gaokao24_mix	minerva	AIME24@32	AIME25@32	Avg
<b>Qwen2.5-32B-Instruct</b>	95.8	73.5	70.0	38.5	42.9	26.5	16.8	11.6	46.95
<b>Trained with 1 epoch</b>									
<b>+ SYNTHETIC-1-10k</b>	92.9	71.8	52.5	38.4	23.1	24.3	35.6	34.0	46.6
<b>+ Open-R1-10k</b>	91.5	72.3	65.0	38.4	20.9	24.6	43.0	33.5	48.7
<b>+ DataFlow-Reasoning-10K</b>	93.9	72.3	72.5	38.7	38.5	26.5	35.9	34.5	<b>51.6</b>
<b>Trained with 2 epochs</b>									
<b>+ SYNTHETIC-1-10k</b>	94.5	78.4	75.0	45.0	24.2	28.3	48.4	37.9	54.0
<b>+ Open-R1-10k</b>	93.9	77.2	80.0	44.1	20.9	25.4	51.0	40.7	54.2
<b>+ DataFlow-Reasoning-10K</b>	94.4	76.6	75.0	45.2	42.9	25.7	45.4	40.0	<b>55.7</b>

0.95, and top-k = 20. All models are fine-tuned with either 1 epoch or 2 epochs on 10k examples using Qwen2.5-32B-Instruct.

## 7.2.2 Experimental Results

Our first observation is that training on Synthetic-1 random subsets yields limited improvement over the base model. While minor gains appear on AMC23 and AIME benchmarks after 2 epochs, the overall average remains similar to the instruction-only baseline (47.0 vs. 46.6).

In contrast, the Open-R1 synthetic subset provides a stronger training signal: two epochs of fine-tuning increase the average score from 48.7 to 54.2, demonstrating that Open-R1-style CoT data is effective for enhancing mathematical reasoning in a 32B model. Building on this, our DATAFLOW-synthesized dataset achieves the strongest overall gains using only 10k samples, two epochs of fine-tuning reach the highest average performance of 55.7, surpassing both Open-R1 (54.2) and Synthetic-1 (54.0). These results indicate that combining verified NuminaMath seeds, MathQ-Verify filtering, and DeepSeek-R1-driven CoT generation yields more precise, diverse, and robust reasoning supervision.

Overall, the experiments demonstrate that data quality, rather than data scale, is the dominant factor in mathematical reasoning performance. Even with the same 10k size, our DATAFLOW-based synthesis pipeline consistently outperforms existing synthetic sources.

## 7.3 Code Data Preparation

### 7.3.1 Experimental Setting

To investigate the effect of high-quality code instruction data on code generation performance, we construct supervised fine-tuning (SFT) datasets using seed samples from **Ling-Coder-SFT** [12]. We first randomly sample 20k instances from the Ling-Coder-SFT corpus and process them through the DATAFLOW **CodeGenDataset\_APIPipeline**. This yields three curated code instruction datasets of different scales, DATAFLOW-CODE-1K, DATAFLOW-CODE-5K, and DATAFLOW-CODE-10K, each designed to provide high-quality, pipeline-refined supervision signals for code generation tasks.

We compare our synthesized datasets against two widely used baselines, each subsampled to 1k examples for fairness:

- **Code Alpaca (1k)**[5]: a randomly sampled subset from the Code Alpaca dataset.
- **Self-OSS-Instruct-SC2-Exec-Filter-50k(1k)** [63] : a 1k random subset from the SC2-Exec-Filter dataset, which incorporates execution-based filtering.

Models are fine-tuned on DATAFLOW-CODE-1K, DATAFLOW-CODE-5K, and DATAFLOW-CODE-10K using full-parameter SFT.

We then experiment with two base models: **Qwen2.5-7B-Instruct** and **Qwen2.5-14B-Instruct**. Evaluation



**Table 6** Code Pipeline: Performance comparison of Qwen2.5-7B-Instruct and Qwen2.5-14B-Instruct under different SFT dataset settings (all numbers in %).

Training Data	BigCodeBench	LiveCodeBench(v6)	CruxEval (Input)	CruxEval (Output)	HumanEval+	Avg
Trained on Qwen2.5-7B-Instruct						
<b>Qwen2.5-7B-Instruct</b>	35.3	23.4	44.8	43.9	72.6	44.0
<b>+ Code Alpaca-1K</b>	33.3	18.7	45.6	46.4	66.5	42.1
<b>+ Self-OSS</b>	31.9	21.4	46.9	45.9	70.1	43.2
<b>+ DataFlow-Code-1K</b>	35.5	25.7	48.0	45.1	72.6	45.4
<b>+ DataFlow-Code-5K</b>	36.2	<b>26.4</b>	48.6	45.0	73.2	45.9
<b>+ DataFlow-Code-10K</b>	<b>36.8</b>	26.0	<b>48.8</b>	<b>45.4</b>	<b>73.8</b>	<b>46.2</b>
Trained on Qwen2.5-14B-Instruct						
<b>Qwen2.5-14B-Instruct</b>	37.5	33.4	48.0	48.5	74.4	48.4
<b>+ Code Alpaca-1K</b>	37.0	28.2	50.2	49.6	71.3	47.3
<b>+ Self-OSS</b>	36.9	22.3	52.6	50.1	68.3	46.0
<b>+ DataFlow-Code-1K</b>	41.4	<b>33.7</b>	51.0	50.9	<b>77.3</b>	50.9
<b>+ DataFlow-Code-5K</b>	41.1	33.2	52.5	50.6	76.2	50.7
<b>+ DataFlow-Code-10K</b>	<b>41.9</b>	33.2	<b>52.9</b>	<b>51.0</b>	76.2	<b>51.0</b>

is conducted on four code benchmarks: (1) **BigCodeBench** [75], (2) **LiveCodeBench** [30], (3) **CruxEval** [22], and (4) **HumanEval** [8]. The final performance is reported as the average across these four benchmarks. All values in Table 6 are percentages.

### 7.3.2 Experimental Results

Table 6 shows that our synthesized datasets consistently improve the code generation performance of both Qwen2.5-7B-Instruct and Qwen2.5-14B-Instruct across all benchmarks. For the 7B model, even 1k of our synthetic data already outperforms both the Code Alpaca and SC2 execution-filtered baselines. Specifically, DATAFLOW-CODE-1K improves BigCodeBench, LiveCodeBench, and CruxEval scores over the original model, while remaining competitive on HumanEval+. Scaling the supervision to 5k and 10k further boosts overall performance. In particular, the DATAFLOW-CODE-10K setting achieves the best results on all metrics, including 36.8 on BigCodeBench, 48.8 on CruxEval(Input), and 45.4 on CruxEval(Output), and yields the highest overall average score of 46.2, surpassing both Code Alpaca-1K and SC2-Exec-Filter under the same data scale.

For the larger Qwen2.5-14B-Instruct model, the benefits are even more pronounced. While Code Alpaca-1k and SC2 filtering provide moderate improvements over the original 14B model, our datasets consistently deliver stronger gains across all metrics. In particular, DATAFLOW-CODE-10K reaches an average score of 51.0, achieving 41.9 on BigCodeBench, 52.9 on CruxEval(Input), and 51.0 on CruxEval(Output). Notably, LiveCodeBench, which stresses executable correctness—rises from 21.9 (Code Alpaca-1k) to 33.2 under our synthetic supervision. These results indicate that the DATAFLOW-generated data provide more explicit execution-grounded signals and structured reasoning cues than existing open-source sources.

Overall, the experiments demonstrate that DATAFLOW-driven synthesis consistently outperforms existing open-source code instruction datasets even under the same sample scale. The consistent gains from 1k to 10k indicate a simple trend: with more high-quality DATAFLOW training samples, the model keeps getting better on code reasoning tasks.

## 7.4 Text-to-SQL Data Preparation

### 7.4.1 Experimental Setting

To evaluate the effectiveness of Text-to-SQL data generation, we construct a training corpus comprising 89,544 high-quality Text-to-SQL instances, which is called DATAFLOW-TEXT2SQL-90K. Each instance in DATAFLOW-TEXT2SQL-90K includes natural language questions, corresponding SQL queries, and chain-of-thought reasoning traces. Specifically, these data are derived through systematic augmentation of seed SQL queries: 37,517 instances originate from the Spider-train [70] dataset, 37,536 from the BIRD-train [40]



**Table 7** Text-to-SQL Pipeline: Performance of LLMs on mainstream benchmarks. The first two blocks list closed-source and open-source base models. The last two blocks show fine-tuned models, where the first column indicates the training data setting.

LLM / Training Data	Spider dev		Spider test		BIRD dev		EHRSQL		Spider-DK		Spider-Syn		Spider-Realistic		Average	
	Gre	Maj	Gre	Maj	Gre	Maj	Gre	Maj	Gre	Maj	Gre	Maj	Gre	Maj	Gre	Maj
Closed-source LLMs																
GPT-4o-mini	70.4	71.0	82.4	83.7	58.8	61.5	37.9	43.1	73.3	74.4	60.5	61.6	64.4	66.7	64.0	66.0
GPT-4-Turbo	72.4	72.2	83.4	84.2	62.0	63.6	43.1	44.8	72.3	72.1	62.9	63.5	67.5	68.3	66.2	67.0
GPT-4o	70.9	70.7	83.2	84.9	61.9	64.0	44.9	45.5	72.9	73.5	59.6	62.3	66.5	66.7	65.7	66.8
Open-source LLMs																
DeepSeek-Coder-7B-Instruct	63.2	63.2	70.5	73.2	43.1	48.0	28.6	33.9	60.9	64.1	49.9	51.7	58.7	58.9	53.6	56.1
Qwen2.5-Coder-7B-Instruct	73.4	77.1	82.2	85.6	50.9	61.3	24.3	36.9	67.5	73.6	63.1	66.9	66.7	70.5	61.2	67.4
Qwen2.5-7B-Instruct	65.4	68.9	76.8	82.6	46.9	56.4	20.9	32.1	63.7	71.8	54.2	60.0	56.7	63.6	54.9	62.2
OpenCoder-8B-Instruct	59.5	59.5	68.3	70.1	37.5	45.3	21.9	29.9	62.6	64.7	46.0	46.1	49.0	49.4	49.3	52.1
Meta-Llama-3.1-8B-Instruct	61.8	67.7	72.2	78.5	42.0	53.1	24.6	33.7	62.6	69.9	53.1	59.3	57.5	61.0	53.4	60.5
Granite-8B-Code-Instruct	58.5	59.2	64.9	68.6	27.6	32.5	16.0	22.6	50.7	54.4	45.0	46.8	48.8	49.4	44.5	47.6
Granite-3.1-8B-Instruct	58.3	65.0	69.8	75.3	36.0	47.2	19.6	32.3	60.0	66.5	47.7	53.8	46.5	57.1	48.3	56.7
Trained on Meta-Llama-3.1-8B-Instruct																
SynSQL(50K)	67.1	73.9	72.7	78.6	49.1	55.2	33.6	40.8	63.8	66.1	59.6	63.5	69.3	71.6	59.3	64.2
SynSQL(90K)	68.2	74.6	73.4	78.5	51.1	54.9	31.8	38.0	61.8	67.4	58.9	63.6	69.0	70.9	59.2	64.0
SynSQL(2.5M)	70.6	73.7	78.3	82.5	58.9	62.0	35.1	37.0	72.3	74.7	61.0	63.1	67.9	69.4	63.4	66.1
Spider+BIRD+DataFlow-Text2SQL-90K	74.9	79.2	78.4	82.3	53.4	58.9	28.4	36.5	67.7	69.7	66.6	69.1	74.4	75.0	63.4	67.2
<b>DataFlow-Text2SQL-50K</b>	69.9	76.8	75.1	80.1	51.4	57.6	28.0	36.4	65.9	68.1	61.3	67.5	69.6	73.5	60.2	65.7
<b>DataFlow-Text2SQL-90K</b>	71.4	76.4	75.8	80.0	54.6	56.8	55.5	56.3	66.5	67.7	61.6	67.3	71.4	72.7	65.3	68.2
Trained on Qwen2.5-Coder-7B-Instruct																
SynSQL(50K)	77.1	82.1	81.8	84.8	54.0	59.3	33.1	44.1	67.1	69.5	68.0	70.6	77.2	80.3	65.5	70.1
SynSQL(90K)	79.2	83.1	82.3	84.4	56.2	59.4	31.4	41.4	65.0	70.7	67.2	70.7	77.0	79.9	65.5	69.9
SynSQL(2.5M)	81.2	81.6	87.9	88.3	63.9	66.1	34.9	40.0	76.1	77.8	69.7	69.6	76.2	78.0	70.0	71.6
Spider+BIRD+DataFlow-Text2SQL-90K	85.5	87.5	87.5	88.5	58.3	64.0	27.9	39.8	71.0	73.1	75.0	76.2	82.3	83.7	69.6	73.3
<b>DataFlow-Text2SQL-50K</b>	80.9	84.9	84.6	85.8	57.9	62.5	27.8	39.4	69.7	71.2	70.0	74.0	77.8	82.1	67.0	71.4
<b>DataFlow-Text2SQL-90K</b>	82.0	85.0	84.8	86.0	59.2	61.5	56.1	58.7	69.7	71.0	69.9	74.4	79.5	81.7	71.6	74.0

dataset, and 14,491 from the EHRSQL-train [36] dataset. DATAFLOW pipeline ensures rich syntactic and semantic diversity in SQL structures, question phrasing, and multi-step reasoning processes.

For our method (DATAFLOW-Text2SQL rows in Table 7), models are fine-tuned exclusively on our synthesized corpus, unless otherwise specified. For evaluation, we adopt six widely recognized Text-to-SQL benchmarks: Spider [70], BIRD [40], EHRSQL [36], Spider-DK [18], Spider-Syn [17], and Spider-Realistic [14]. During inference with LLMs, we investigate two decoding strategies: greedy decoding (denoted as **Gre**), which uses temperature 0 for deterministic output generation, and majority voting (denoted as **Maj**). The majority voting strategy samples 8 candidate responses per input at temperature 0.8, executes all valid SQL queries, and selects the query whose execution result appears most frequently among the candidates as the final prediction. We additionally randomly sampled 50K instances to construct DATAFLOW-TEXT2SQL-50K. For comparison, we also randomly sampled the same number of instances from SynSQL [37].

## 7.4.2 Experimental Results

As shown in Table 7, the generated data leads to consistent performance improvements across multiple mainstream benchmarks, demonstrating the effectiveness of DATAFLOW [4]. For both models, Meta-Llama-3.1-8B-Instruct [21] and Qwen2.5-Coder-7B-Instruct [29], training on our generated data significantly improves performance over their respective baselines as well as other competing models. When fine-tuned on the generated data, Qwen2.5-Coder-7B-Instruct achieves notable gains: execution accuracy (**Gre**) on Spider-dev increases from 73.4 to 82.0 (+8.6), on BIRD-dev from 50.9 to 59.2 (+8.3), and on the challenging EHRSQL benchmark from 24.3 to 56.1 (+31.8). These results confirm that DATAFLOW-TEXT2SQL-90K exhibits high quality and strong training utility.

Compared with other training datasets, our data also demonstrates clear advantages. At comparable

data scales, models trained on DATAFLOW-TEXT2SQL-90K and DATAFLOW-TEXT2SQL-50K consistently outperform those trained on SynSQL [37] (SynSQL(90K) and SynSQL(50K), respectively). Specifically, on the Spider-test and BIRD-dev datasets, the model trained on DATAFLOW-TEXT2SQL-50K achieves 84.6 and 57.9 execution accuracy (Gre), surpassing SynSQL(50K) [37], which obtains 81.8 and 54.0. Likewise, the model trained on DATAFLOW-TEXT2SQL-90K not only surpasses the baseline models but also outperforms SynSQL(90K) [37]. Remarkably, even when trained on a much smaller dataset, the model fine-tuned with DATAFLOW-TEXT2SQL-90K achieves performance comparable to SynSQL-2.5M [37] on several challenging benchmarks. These improvements highlight the higher quality of the training data generated by DATAFLOW.

## 7.5 AgenticRAG Data Preparation

### 7.5.1 Experimental Setting

In the field of AgenticRAG, the automatic generation of multihop questions has long been a challenging issue in research. This study constructs a multihop question dataset with a scale of 10k based on **the DataFlow AgenticRAG Pipeline** and conducts a comparative analysis with existing mainstream multihop question answering datasets (2WikiMultiHopQA [25], Musique [58], HotpotQA [68], and Bamboogle [51]).

The specific workflow of the dataset generation pipeline is as follows:

- Documents are randomly selected from the Wikipedia dump to form the initial document set. To avoid the interference of data distribution overlap on the experimental results, documents that have already appeared in the test benchmark are excluded.
- The o4-mini model combined with the generation module of DATAFLOW AgenticRAG is used to generate the initial draft of multihop questions based on the filtered initial documents.
- The verification module is employed to screen the quality of the initial question drafts, eliminating samples with problems such as intermediate question leakage, logical errors, and excessively high or low difficulty, ultimately forming a high-quality multihop question dataset, which we call DATAFLOW-AGENTICRAG-10K.

This study adopts the ReCall [9] framework to complete the model training and evaluation. In the training phase, Qwen2.5-7B-Instruct is selected as the base model, and the GRPO reinforcement learning algorithm is used for model optimization. In the evaluation phase, the model’s temperature parameter is set to 0.0.

For the retrieval component, E5-base-v2 [59] is chosen as the retriever, and the 2018 Wikipedia dump is used as the corpus. All corpus indexing and embedding calculations are preprocessed using FlashRAG [32]. Throughout the entire training and evaluation process, the model is allowed to independently specify the topk value for retrieval, and the default topk value is set to 5 to balance retrieval efficiency and performance.

### 7.5.2 Experimental Results

Table 8 reports the exact-match performance across four multi-hop benchmarks. We group the results by the training dataset and compute an out-of-distribution (OOD) average by removing the in-domain test set of each dataset (e.g., HotpotQA-trained models exclude HotpotQA). To fairly compare against our synthetic data, we additionally report DF-OOD (matched), which applies the same in-domain exclusion to DF-AgenticRAG-10k.

*Comparison with HotpotQA-trained models.* Across 1–3 epochs, HotpotQA-10k achieves OOD averages of 33.7, 35.1, and 36.4. Under the same exclusion (w/o HotpotQA), DF-AgenticRAG achieves 33.8, 35.9, and 37.4—consistently matching or surpassing HotpotQA by +0.1 to +1.0 points despite using entirely synthetic supervision. This indicates that DF-AgenticRAG provides generalization comparable to a widely used human-constructed dataset.

*Comparison with Musique-trained models.* Musique-20k yields an OOD average of 42.4 when evaluated w/o Musique. Under the same exclusion, DF-AgenticRAG (2 epochs effective scale = 20k) reaches 43.6, outperforming Musique by +1.2 points. This shows that our synthetic dataset not only matches but outperforms a strong human-annotated multi-hop benchmark at the same effective training scale.

**Table 8** AgenticRAG Pipeline: Performance comparison between synthetic datasets and existing human-constructed datasets. All values are Exact Match (%). “OOD-Avg” excludes the in-domain test set of each training dataset. “DF-OOD (matched)” provides the OOD score of DF-AgenticRAG under the \*same\* in-domain exclusion, ensuring fair comparison.

Training Data	HotpotQA	2Wiki	Musique	Bamboogle	Avg	OOD-Avg	DF-OOD (matched)
Qwen-2.5-7B-Instruct	25.0	25.8	9.9	27.2	22.0	–	–
<b>Trained on HotpotQA (in-domain = HotpotQA)</b>							
HotpotQA-10k (1 epoch)	40.2	41.9	16.7	42.4	35.3	33.7	<b>33.8</b>
HotpotQA-10k (2 epochs)	43.4	44.9	18.9	41.6	37.2	35.1	<b>35.9</b>
HotpotQA-10k (3 epochs)	45.3	48.0	20.3	40.8	38.6	36.4	<b>37.4</b>
<b>Trained on Musique (in-domain = Musique)</b>							
Musique-20k (1 epoch)	41.1	44.7	19.2	41.6	36.6	42.4	<b>43.6</b>
<b>Trained on 2Wiki (in-domain = 2Wiki)</b>							
2Wiki-30k (2 epochs)	41.3	55.1	17.8	42.4	39.1	33.8	<b>36.4</b>
<b>DF-AgenticRAG (raw results, for reference)</b>							
DataFlow-AgenticRAG-10k (1 epoch)	39.3	42.6	17.3	41.6	34.3	–	–
DataFlow-AgenticRAG-10k (2 epochs)	43.1	44.6	19.9	43.2	37.7	–	–
DataFlow-AgenticRAG-10k (3 epochs)	42.6	45.5	20.2	46.4	38.7	–	–

**Table 9** Knowledge Extraction: Accuracy comparison on PubMedQA, Covert, and PubHealth under different reasoning and training settings.

Method (ACC)	PubMedQA	Covert	PubHealth
<b>CoT</b>	36.40%	48.33%	29.00%
<b>RAG</b>	43.33%	17.55%	19.60%
<b>SFT (DataFlow-Knowledge)</b>	<b>53.40%</b>	<b>68.33%</b>	<b>40.86%</b>

*Comparison with 2Wiki-trained models.* 2Wiki-30k achieves an OOD average of 33.8. Under the same exclusion (w/o 2Wiki), DF-AgenticRAG (3 epochs, effective scale=30k) reaches 36.4, a substantial improvement of +2.6 points. This represents the largest gap among all baselines and highlights the strong cross-dataset generalization capacity of our synthetic questions.

*Summary.* Across all training regimes and all in-domain exclusions, DF-AgenticRAG-10k is either the best or tied for the best OOD dataset, and in several cases (Musique, 2Wiki) significantly surpasses human-constructed datasets. These results demonstrate that our pipeline produces multi-hop reasoning data with superior cross-dataset generalization, suggesting that high-quality synthetic data can not only match but consistently exceed the robustness of existing human-annotated multi-hop datasets.

## 7.6 Knowledge Extraction

### 7.6.1 Experimental Setting

To expand beyond the limited annotated data and take advantage of massive raw corpora from the Internet, we proposed the Knowledge Extraction pipeline, a semi-automated system for corpus cleaning and QA synthesis. The pipeline performs text normalization using MinerU [46], segments long documents, filters noisy or low-quality sentences, generates factuality-aware QA pairs, and conducts automated quality checks, ultimately producing a high-quality synthetic dataset used for supervised fine-tuning (SFT).

In our experiment, the training data is derived from 140M tokens of raw medical data drawn from three major sources. The first source is MedQA Books, a collection of 18 widely used medical textbooks from the USMLE curriculum [31]. The second source consists of 9,330 publicly available StatPearls articles from the NCBI Bookshelf [66]. The third source contains 45,679 clinical guideline documents aggregated from 16 professional guideline providers [10]. These corpora serve as the input to the Knowledge Extraction pipeline,

which converts them into structured, high-quality QA dataset, denoted as DATAFLOW-KNOWLEDGE which is suitable for model training.

For model training, we fine-tune Qwen2.5-7B-Instruct on the DATAFLOW-generated dataset. The SFT process is performed for 37,500 steps over five epochs. For comparison, we also evaluate a zero-shot Chain-of-Thought (CoT) prompting baseline and a retrieval-augmented generation (RAG) baseline using top- $k = 10$  retrieval with `medcpt-query-encoder` as the retriever and `medcpt-article-encoder` as the document encoder. All baselines share same hyperparameter setting during inference time.

We evaluate our models on three medical QA benchmarks: PubMedQA [33], which focuses on biomedical research questions; Covert [45], which evaluates clinical knowledge and reasoning; and PubHealth [34], which targets public-health misinformation classification.

## 7.6.2 Experimental Results

Table 9 presents the accuracy results across all benchmarks. The CoT baseline performs poorly across the board, indicating that zero-shot reasoning alone is insufficient for medical question answering without more targeted supervision. The RAG baseline provides modest improvement on PubMedQA, but remains unstable and substantially underperforms on Covert and PubHealth, suggesting that retrieval alone cannot substitute for explicit training on structured domain data.

In contrast, the SFT model trained on DATAFLOW-KNOWLEDGE synthetic data achieves the highest accuracy on all benchmarks, surpassing both CoT prompting and RAG-based methods by large margins. Notably, it delivers more than 15–20 absolute accuracy gains on PubMedQA and Covert, and an 11-point improvement on PubHealth, demonstrating that the cleaned and structured QA pairs produced by our Knowledge Extraction pipeline offer significantly stronger supervision.

Overall, these results show that high-quality synthetic QA data—when curated and verified through a targeted DATAFLOW pipeline—can substantially enhance the domain reasoning capabilities of a general-purpose model, outperforming both inference-time prompting and retrieval-augmented baselines.

## 7.7 Unified Multi-Domain Data Preparation with DataFlow

### 7.7.1 Experimental Setting

*Data Construction* To evaluate the efficiency and effectiveness of unified data preparation across modality-specific reasoning tasks, we construct an integrated training corpus that combines Math, Code, and General Instruction data. All data are generated or filtered through the DATAFLOW framework as follows:

- **Math.** We synthesize high-quality mathematical problems and chain-of-thought (CoT) solutions using the DATAFLOW *Reasoning Pipeline*, with the MATH dataset serving as seed input. We randomly sample 3k instances for training.
- **Code.** Code data are produced using the DATAFLOW *CodeGenDataset\_APIPipeline*, built upon 20k randomly sampled LingoCoder SFT examples. We generate 1k–10k high-quality code instructions and benchmark against Code Alpaca and SC2-Exec-Filter. A subset of 2k samples is used for training.
- **Text / General Instruction.** For natural language tasks, we employ the DATAFLOW Condor Generator + Refiner pipeline to generate high-consistency instruction–response and dialogue pairs. Outputs are further processed by the SFT-quality filtering pipeline. We randomly sample 5k instances.

All models are fine-tuned on the combined DATAFLOW-INSTRUCT-10K corpus using full-parameter SFT. Evaluation covers: (1) seven math benchmarks, (2) four code benchmarks, and (3) MMLU [23] and C-Eval [27] for general knowledge and reasoning.

*Baselines.* We additionally compare DATAFLOW-INSTRUCT-10K with baselines constructed from the **Infinity-Instruct** (Inf) [39] dataset, a large-scale general-purpose instruction corpus widely used in instruction tuning. Two baselines are included:

**Table 10** Performance of DATAFLOW-INSTRUCT-10K on Math Benchmarks: Qwen2-7B-Base and Qwen2.5-7B-Base finetuned series of models (Exact Match %).

Model	MATH	GSM8K	AMC23	AIME24	Minerva	Gaokao	Olympiad	Math-Avg
Models based on Qwen2-7B								
<b>Qwen2-7B-Base</b>	21.2	55.9	15.0	0.0	9.9	30.8	7.7	20.1
<b>+ Inf-10K</b>	45.6	81.7	25.0	3.3	11.8	24.2	11.1	29.0
<b>+ Inf-1M</b>	45.4	79.2	25.0	0.0	13.2	22.0	10.4	27.9
<b>+ DataFlow-Instruct-10K</b>	54.0	83.0	27.5	0.0	16.5	25.3	20.3	32.4
<b>Qwen2-7B-Instruct</b>	53.9	86.2	22.5	3.3	17.6	35.2	19.6	34.0
Models based on Qwen2.5-7B								
<b>Qwen2.5-7B-Base</b>	62.8	67.1	45.0	10.0	17.6	27.5	29.6	37.1
<b>+ Inf-10K</b>	40.2	30.9	25.0	3.3	9.2	27.5	21.8	22.6
<b>+ Inf-1M</b>	50.6	82.0	27.5	0.0	22.1	30.8	20.0	33.3
<b>+ DataFlow-Instruct-10K</b>	73.8	88.2	47.5	16.7	30.9	31.9	37.6	46.7
<b>Qwen2.5-7B-Instruct</b>	75.1	92.4	47.5	10.0	34.9	48.4	40.6	49.8

**Table 11** Performance of DATAFLOW-INSTRUCT-10K on Code and Knowledge benchmarks: Qwen2-7B-Base and Qwen2.5-7B-Base finetuned models.

Model	HumanEval	MBPP	Code-Avg	MLLU	C-EVAL	Knowledge-Avg
Models based on Qwen2-7B						
<b>Qwen2-7B-Base</b>	66.5	66.1	66.3	69.6	82.8	76.2
<b>+ Inf-10K</b>	64.0	71.7	67.8	69.3	83.0	76.2
<b>+ Inf-1M</b>	65.9	70.4	68.2	69.5	83.0	76.2
<b>+ DataFlow-Instruct-10K</b>	64.6	67.7	66.2	69.4	82.8	76.1
<b>Qwen2-7B-Instruct</b>	73.8	65.3	69.6	69.9	82.0	76.0
Models based on Qwen2.5-7B						
<b>Qwen2.5-7B-Base</b>	78.7	74.3	76.5	71.9	80.0	76.0
<b>+ Inf-10K</b>	77.4	77.8	77.6	71.8	79.9	75.8
<b>+ Inf-1M</b>	78.0	78.0	78.0	72.2	79.4	75.8
<b>+ DataFlow-Instruct-10K</b>	80.5	76.7	78.6	72.1	80.2	76.2
<b>Qwen2.5-7B-Instruct</b>	81.7	79.4	80.6	71.8	79.6	75.7

- **Inf-10K**: a random 10k subset of Infinity-Instruct used for SFT.
- **Inf-1M**: a random 1M subset of Infinity-Instruct.

Comparing against Inf-10K/1M allows us to assess whether high-quality, domain-specific synthetic data (math, code, text) generated through DATAFLOW provides more stable and reliable improvements than large generic instruction data.

### 7.7.2 Experimental Results

Across Math, Code, and Knowledge evaluation suites, our unified multi-domain data preparation strategy provides consistent and robust gains for both Qwen2.5-7B and Qwen2-7B models. A notable pattern observed across all tables is that DATAFLOW-INSTRUCT-10K almost always achieves the best performance among all non-Instruct finetuned models, and in many cases narrows the gap to the Instruct models to within only 2–4 points, despite using orders-of-magnitude less data.

*Math Reasoning.* As shown in Table 10, DATAFLOW-processed math data yields the largest and most stable gains. For Qwen2.5-7B-Base, training on our synthesized math subset improves the overall score from 37.1 to 46.7, which is:

- the best performance among all non-Instruct models, surpassing Inf-10K (22.6) and Inf-1M (33.3) by a clear margin;
- only 3.1 points below the Instruct model (49.8), demonstrating that targeted, high-quality synthetic data can nearly match the performance of costly human-aligned instruction tuning.

A similar trend holds for Qwen2-7B: DATAFLOW-INSTRUCT-10K reaches 32.4 overall, outperforming Inf-10K and Inf-1M, and approaching the Instruct model (34.04). These results highlight that DATAFLOW math synthesis produces significantly more stable and effective improvements than generic inference-generated data.

*Code Generation.* As shown in Table 11, DATAFLOW-INSTRUCT-10K consistently delivers the best Code-Overall performance among all non-Instruct models. For Qwen2.5-7B-Base, DATAFLOW-INSTRUCT-10K raises Code-Overall from 76.5 to 78.6, outperforming Inf-10K (77.6) and Inf-1M (78.0), and reaching within 2.0 points of the Instruct model (80.6). For Qwen2-7B-Base, DATAFLOW-INSTRUCT-10K again matches or exceeds all Inf baselines.

These results show that adding multi-domain synthetic data does not harm code ability (a common issue in mixed-domain SFT), and often improves it. This further supports the robustness of DATAFLOW’s domain-balanced synthetic corpus.

*General Knowledge and NLP.* As summarized in Table 11, our unified dataset also preserves strong general knowledge and reasoning. Across MMLU and C-Eval, DF-Gen-10K:

- matches or slightly improves upon the Base models,
- avoids the regressions frequently observed in Inf-10K and Inf-1M,
- frequently ranks second only to the Instruct model, confirming that DATAFLOW-generated text data provides high-quality supervision even without human instruction tuning.

*Summary.* Together, these results demonstrate that high-quality, domain-specialized synthetic data generated via DATAFLOW produces the strongest non-Instruct performance across Math, Code, and Knowledge. DATAFLOW-INSTRUCT-10K consistently outperforms generic inference-generated data (Inf-10K/Inf-1M) and often approaches the performance of the Instruct models themselves. This highlights the effectiveness of DATAFLOW’s unified, pipeline-driven data preparation for building multi-capability LLMs without reliance on large-scale human-authored instruction corpora.

## 7.8 Agentic Orchestration

### 7.8.1 Experimental Setting

We evaluate the proposed agent orchestration framework on realistic data processing and pipeline construction tasks. Specifically, we selected 6 representative pipelines as benchmarks. For each pipeline, we manually constructed natural language task descriptions at 3 difficulty levels, resulting in 18 user queries to assess automatic orchestration capabilities across varying description granularities. The difficulty levels are defined as follows:

- **Easy.** Descriptions are explicit, directly specifying the functions of required operators (or key operators) and the main processing steps.
- **Medium.** Descriptions are coarse, providing only general processing goals and key constraints without explicitly listing the complete operator sequence.
- **Hard.** Only a high-level requirement or final goal is provided with minimal hints regarding intermediate steps, requiring the system to infer the complete processing flow and operator combination.

For each task, the user provides a natural language description of the goal, and the system must automatically orchestrate a pipeline composed of multiple operators to meet the requirement.



*Evaluation Metrics.* To quantitatively assess orchestration quality, we employ an external LLM as an automatic judge. The evaluator compares the generated pipeline against ground truth under two distinct settings:

- **Text Specification Alignment.** The predicted graph is evaluated against text specifications to verify if the pipeline structure satisfies the detailed task requirements.
- **Code Implementation Consistency.** The pipeline is compared with reference Python implementations to assess logical equivalence regarding operator usage and processing steps.

Based on these comparisons, we report the **LLM-Judge Score** ( $s \in [0, 1]$ ), which measures the consistency of operator coverage and execution order between the generated pipeline and the reference under the corresponding evaluation setting.

## 7.8.2 Experimental Results

**Table 12** Agent orchestration performance by evaluation mode and description difficulty.

Metric	Easy	Medium	Hard	Overall
<b>Text spec evaluation (pipeline mode)</b>				
Avg. LLM-Judge	0.92	0.86	0.60	0.80
<b>Code GT evaluation (code mode)</b>				
Avg. LLM-Judge	0.60	0.59	0.23	0.49

Table 12 reports the LLM-Judge scores under text-spec (pipeline) and code-GT (code) evaluations across difficulty levels. Overall, the framework performs well when judged against textual requirements (**0.80** overall), but is markedly lower when matching reference implementations (**0.49** overall), reflecting the stricter nature of code-level equivalence. Performance degrades as descriptions become less explicit: in pipeline mode scores drop from **0.92/0.86** (Easy/Medium) to **0.60** (Hard), while in code mode the drop is more severe, reaching **0.23** on Hard, indicating that under-specified queries often lead to alternative yet plausible operator compositions that diverge from a single ground-truth program.

## 8 Conclusion

In summary, DATAFLOW addresses a critical gap in the data-centric LLM ecosystem by providing the first unified, LLM-driven data preparation framework. It mitigates long-standing challenges in the field—such as the difficulty of sharing, reproducing, and comparing data preparation algorithms—through a modular and user-friendly programming interface. The framework integrates nearly 200 operators, over 80 prompt templates, and unified abstractions for serving and storage, all of which compose into six high-quality pipelines spanning the major LLM data domains. Extensive experiments demonstrate that these pipelines achieve strong, often state-of-the-art results, confirming that DATAFLOW effectively balances the tension between domain-specific customization and system-level standardization.

Built atop this foundation, the DATAFLOW-CLI and DATAFLOW-Agent further amplify extensibility by enabling rapid template generation, natural-language-driven workflow construction, and scalable extension development. Together, these components lay the groundwork for a sustainable and interoperable data preparation ecosystem that can evolve alongside increasingly complex data-centric AI workflows.

Looking forward, we aim to expand the DATAFLOW-Ecosystem along multiple modality axes, including DATAFLOW-TABLE, DATAFLOW-GRAPH, and DATAFLOW-MULTIMODAL, to support richer data types and workflows. We also plan to develop domain-oriented variants, such as DATAFLOW-AI4S and DATAFLOW-INDUSTRY, tailored for large-scale production environments. These extensions will broaden the applicability of DATAFLOW and strengthen its role as a foundational substrate—and a common protocol—for future research, engineering practice, and community-driven innovation in LLM data preparation.



## References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- [2] LangChain AI. Langgraph. <https://github.com/langchain-ai/langgraph>, 2024.
- [3] Tianyi Bai, Hao Liang, Binwang Wan, Ling Yang, Bozhou Li, Yifan Wang, Bin Cui, Conghui He, Binhang Yuan, and Wentao Zhang. A survey of multimodal large language model from a data-centric perspective. arXiv preprint arXiv:2405.16640, 2024.
- [4] Qifeng Cai, Hao Liang, Chang Xu, Tao Xie, Wentao Zhang, and Bin Cui. Text2sql-flow: A robust sql-aware data augmentation framework for text-to-sql. arXiv preprint arXiv:2511.10192, 2025.
- [5] Sahil Chaudhary. Code alpaca: An instruction-following llama model for code generation. <https://github.com/sahil280114/codealpaca>, 2023.
- [6] Daoyuan Chen, Yilun Huang, Zhijian Ma, Hesun Chen, Xuchen Pan, Ce Ge, Dawei Gao, Yuexiang Xie, Zhaoyang Liu, Jinyang Gao, et al. Data-juicer: A one-stop data processing system for large language models. In Companion of the 2024 International Conference on Management of Data, pages 120–134, 2024.
- [7] Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Sriniwasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca with fewer data. arXiv preprint arXiv:2307.08701, 2023.
- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- [9] Mingyang Chen, Linzhuang Sun, Tianpeng Li, Haoze Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang, Jeff Z Pan, Wen Zhang, Huajun Chen, et al. Learning to reason with search for llms via reinforcement learning. arXiv preprint arXiv:2503.19470, 2025.
- [10] Zeming Chen, Alejandro Hernández Cano, Angelika Romanou, Antoine Bonnet, Kyle Matoba, Francesco Salvi, Matteo Pagliardini, Simin Fan, Andreas Köpf, Amirkeivan Mohtashami, Alexandre Sallinen, Alireza Sakhaeirad, Vinitra Swamy, Igor Krawczuk, Deniz Bayazit, Axel Marmet, Syrielle Montariol, Mary-Anne Hartley, Martin Jaggi, and Antoine Bosselut. Meditron-70b: Scaling medical pretraining for large language models, 2023. URL <https://arxiv.org/abs/2311.16079>.
- [11] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021.
- [12] Codefuse and Ling Team. Every sample matters: Leveraging mixture-of-experts and high-quality data for efficient and accurate code llm, 2025. URL <https://arxiv.org/abs/2503.17793>.
- [13] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. Communications of the ACM, 51(1):107–113, 2008.
- [14] Xiang Deng, Ahmed Hassan Awadallah, Christopher Meek, Oleksandr Polozov, Huan Sun, and Matthew Richardson. Structure-grounded pretraining for text-to-sql. arXiv preprint arXiv:2010.12773, 2020.
- [15] Qianlong Du, Chengqing Zong, and Jiajun Zhang. Mods: Model-oriented data selection for instruction tuning. arXiv preprint arXiv:2311.15653, 2023.
- [16] Erich Gamma. Design patterns: elements of reusable object-oriented software, 1995.

- [17] Yujian Gan, Xinyun Chen, Qiuping Huang, Matthew Purver, John R Woodward, Jinxia Xie, and Pengsheng Huang. Towards robustness of text-to-sql models against synonym substitution. arXiv preprint arXiv:2106.01065, 2021.
- [18] Yujian Gan, Xinyun Chen, and Matthew Purver. Exploring underexplored limitations of cross-domain text-to-sql generalization. arXiv preprint arXiv:2109.05157, 2021.
- [19] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. arXiv preprint arXiv:2101.00027, 2020.
- [20] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The google file system. In Proceedings of the nineteenth ACM symposium on Operating systems principles, pages 29–43, 2003.
- [21] Aaron Grattaffiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
- [22] Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. Cruxeval: A benchmark for code reasoning, understanding and execution. arXiv preprint arXiv:2401.03065, 2024.
- [23] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300, 2020.
- [24] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. arXiv preprint arXiv:2103.03874, 2021.
- [25] Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. arXiv preprint arXiv:2011.01060, 2020.
- [26] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.
- [27] Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. Advances in Neural Information Processing Systems, 36:62991–63010, 2023.
- [28] Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL <https://github.com/huggingface/open-r1>.
- [29] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. arXiv preprint arXiv:2409.12186, 2024.
- [30] Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. arXiv preprint arXiv:2403.07974, 2024.
- [31] Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What disease does this patient have? a large-scale open domain question answering dataset from medical exams, 2020. URL <https://arxiv.org/abs/2009.13081>.
- [32] Jiajie Jin, Yutao Zhu, Zhicheng Dou, Guanting Dong, Xinyu Yang, Chenghao Zhang, Tong Zhao, Zhao Yang, and Ji-Rong Wen. Flashrag: A modular toolkit for efficient retrieval-augmented generation research. In Companion Proceedings of the ACM on Web Conference 2025, pages 737–740, 2025.
- [33] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. In Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP), pages 2567–2577, 2019.
- [34] Neema Kotonya and Francesca Toni. Explainable automated fact-checking for public health claims. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7740–7754, Online, November 2020. Association for

- Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.623. URL <https://aclanthology.org/2020.emnlp-main.623/>.
- [35] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*, pages 611–626, 2023.
- [36] Gyubok Lee, Hyeonji Hwang, Seongsu Bae, Yeonsu Kwon, Woncheol Shin, Seongjun Yang, Minjoon Seo, Jong-Yeup Kim, and Edward Choi. Ehsql: A practical text-to-sql benchmark for electronic health records. *Advances in Neural Information Processing Systems*, 35:15589–15601, 2022.
- [37] Haoyang Li, Shang Wu, Xiaokang Zhang, Xinmei Huang, Jing Zhang, Fuxin Jiang, Shuai Wang, Tieying Zhang, Jianjun Chen, Rui Shi, et al. Omnisql: Synthesizing high-quality text-to-sql data at scale. *arXiv preprint arXiv:2503.02240*, 2025.
- [38] Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. Datacomp-lm: In search of the next generation of training sets for language models. *Advances in Neural Information Processing Systems*, 37:14200–14282, 2024.
- [39] Jijie Li, Li Du, Hanyu Zhao, Bo-wen Zhang, Liangdong Wang, Boyan Gao, Guang Liu, and Yonghua Lin. Infinity instruct: Scaling instruction selection and synthesis to enhance language models. *arXiv preprint arXiv:2506.11116*, 2025.
- [40] Jinyang Li, Binyuan Hui, Ge Qu, Jiaxi Yang, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, et al. Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls. *Advances in Neural Information Processing Systems*, 36, 2024.
- [41] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. *arXiv preprint arXiv:2406.11939*, 2024.
- [42] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval), 5 2023.
- [43] Justus Mattern, Sami Jaghouar, Manveer Basra, Jannik Straube, Matthew Di Ferrante, Felix Gabriel, Jack Min Ong, Vincent Weisser, and Johannes Hagemann. Synthetic-1: Two million collaboratively generated reasoning traces from deepseek-r1, 2025. URL <https://www.primeintellect.ai/blog/synthetic-1-release>.
- [44] meta llama. Introducing Meta Llama 3: The most capable openly available LLM to date, 2024. URL <https://ai.meta.com/blog/meta-llama-3/>. Accessed: 2024-05-02.
- [45] Isabelle Mohr, Amelie Wüthrich, and Roman Klinger. Covert: A corpus of fact-checked biomedical covid-19 tweets, 2022. URL <https://arxiv.org/abs/2204.12164>.
- [46] Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, Zhenjiang Jin, Guang Liang, Rui Zhang, Wenzheng Zhang, Yuan Qu, Zhifei Ren, Yuefeng Sun, Yuanhong Zheng, Dongsheng Ma, Zirui Tang, Boyu Niu, Ziyang Miao, Hejun Dong, Siyi Qian, Junyuan Zhang, Jingzhou Chen, Fangdong Wang, Xiaomeng Zhao, Liqun Wei, Wei Li, Shasha Wang, Ruiliang Xu, Yuanyuan Cao, Lu Chen, Qianqian Wu, Huaiyu Gu, Lindong Lu, Keming Wang, Dechen Lin, Guanlin Shen, Xuanhe Zhou, Linfeng Zhang, Yuhang Zang, Xiaoyi Dong, Jiaqi Wang, Bo Zhang, Lei Bai, Pei Chu, Weijia Li, Jiang Wu, Lijun Wu, Zhenxiang Li, Guangyu Wang, Zhongying Tu, Chao Xu, Kai Chen, Yu Qiao, Bowen Zhou, Dahua Lin, Wentao Zhang, and Conghui He. Mineru2.5: A decoupled vision-language model for efficient high-resolution document parsing, 2025. URL <https://arxiv.org/abs/2509.22186>.
- [47] NVIDIA Corporation. Curating custom datasets for LLM training with NVIDIA nemo curator. <https://developer.nvidia.com/blog/curating-custom-datasets-for-llm-training-with-nvidia-nemo-curator/>, 2024. Accessed: 2025-11-28.
- [48] Malte Ostendorff, Pedro Ortiz Suarez, Lucas Fonseca Lage, and Georg Rehm. Llm-datasets: An open framework for pretraining datasets of large language models. In *First conference on language modeling*, 2024.
- [49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- [50] Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data at scale. Advances in Neural Information Processing Systems, 37:30811–30849, 2024.
- [51] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 5687–5711, 2023.
- [52] Matthew Rocklin et al. Dask: Parallel computation with blocked algorithms and task scheduling. In SciPy, pages 126–132, 2015.
- [53] Chengyu Shen, Zhen Hao Wong, Runming He, Hao Liang, Meiyi Qiang, Zimo Meng, Zhengyang Zhao, Bohan Zeng, Zhengzhou Zhu, Bin Cui, et al. Let’s verify math questions step by step. arXiv preprint arXiv:2505.13903, 2025.
- [54] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. arXiv preprint arXiv: 2409.19256, 2024.
- [55] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. Advances in Neural Information Processing Systems, 35:19523–19536, 2022.
- [56] Gemini Team, R Anil, S Borgeaud, Y Wu, JB Alayrac, J Yu, R Soricut, J Schalkwyk, AM Dai, A Hauth, et al. Gemini: A family of highly capable multimodal models, 2024. arXiv preprint arXiv:2312.11805, 10, 2024.
- [57] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023.
- [58] Harsh Trivedi, Niranjana Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop questions via single-hop question composition. Transactions of the Association for Computational Linguistics, 10:539–554, 2022.
- [59] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. arXiv preprint arXiv:2212.03533, 2022.
- [60] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. In Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers), pages 13484–13508, 2023.
- [61] Yudong Wang, Zixuan Fu, Jie Cai, Peijun Tang, Hongya Lyu, Yewei Fang, Zhi Zheng, Jie Zhou, Guoyang Zeng, Chaojun Xiao, et al. Ultra-fineweb: Efficient data filtering and verification for high-quality llm training data. arXiv preprint arXiv:2505.05427, 2025.
- [62] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
- [63] Yuxiang Wei, Federico Cassano, Jiawei Liu, Yifeng Ding, Naman Jain, Zachary Mueller, Harm de Vries, Leandro von Werra, Arjun Guha, and Lingming Zhang. Selfcodealign: Self-alignment for code generation. arXiv preprint arXiv:2410.24198, 2024.
- [64] Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. Qurating: Selecting high-quality data for training language models. arXiv preprint arXiv:2402.09739, 2024.
- [65] Tom White. Hadoop: The definitive guide. " O’Reilly Media, Inc.", 2012.
- [66] Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. Benchmarking retrieval-augmented generation for medicine, 2024. URL <https://arxiv.org/abs/2402.13178>.
- [67] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. arXiv preprint arXiv:2505.09388, 2025.
- [68] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In Proceedings of the 2018 conference on empirical methods in natural language processing, pages 2369–2380, 2018.

- [69] Ping Yu, Jack Lanchantin, Tianlu Wang, Weizhe Yuan, Olga Golovneva, Iliia Kulikov, Sainbayar Sukhbaatar, Jason Weston, and Jing Xu. Cot-self-instruct: Building high-quality synthetic prompts for reasoning and non-reasoning tasks. arXiv preprint arXiv:2507.23751, 2025.
- [70] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. arXiv preprint arXiv:1809.08887, 2018.
- [71] Matei Zaharia, Reynold S. Xin, Patrick Wendell, Tathagata Das, Michael Armbrust, Ankur Dave, Xiangrui Meng, Josh Rosen, Shivaram Venkataraman, Michael J. Franklin, Ali Ghodsi, Joseph Gonzalez, Scott Shenker, and Ion Stoica. Apache spark: a unified engine for big data processing. Commun. ACM, 59(11):56–65, October 2016. ISSN 0001-0782. doi: 10.1145/2934664. URL <https://doi.org/10.1145/2934664>.
- [72] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Guoyin Wang, et al. Instruction tuning for large language models: A survey. ACM Computing Surveys, 2023.
- [73] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. Sglang: Efficient execution of structured language model programs. Advances in neural information processing systems, 37:62557–62583, 2024.
- [74] Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. arXiv preprint arXiv:2403.13372, 2024.
- [75] Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. arXiv preprint arXiv:2406.15877, 2024.

# Appendix

## A Author Contributions

- Hao Liang: *Project Leader*, *Project Founder*; algorithm lead and manuscript writing.
- Xiaochen Ma: *Project Leader*, *Project Founder*; system lead and manuscript writing.
- Zhou Liu: *Project Leader*, *Project Founder*; DATAFLOW-AGENT lead and manuscript writing.
- Zhen Hao Wong: *Core Contributor*, *Project Founder*; designs and develops reasoning pipelines, AI4S pipelines, and AgenticRAG pipelines.
- Zhengyang Zhao: *Core Contributor*, *Project Founder*; designs and conducts experiments for the Text Pipeline.
- Zimo Meng: *Core Contributor*, *Project Founder*; system development and support for the design and experiments of the Text Pipeline.
- Runming He: *Core Contributor*, *Project Founder*; develops reasoning pipelines and conducts math reasoning experiments.
- Chengyu Shen: *Core Contributor*, *Project Founder*; DATAFLOW evaluation pipelines and reasoning pipelines.
- Qifeng Cai: *Core Contributor*; designs and conducts experiments for the Text-to-SQL Pipeline.
- Zhaoyang Han: *Core Contributor*; designs and conducts experiments for the Knowledge Cleaning Pipeline.
- Meiyi Qiang: *Core Contributor*; scientific visualization, publicity leadership, and testing.
- Yalin Feng: *Core Contributor*; designs DATAFLOW evaluation and PDF2Model pipelines.
- Tianyi Bai: *Core Contributor*; designs and conducts experiments for the Code Pipeline.
- Zewei Pan: *Contributor*; designs and conducts the operator-writing workflow and experiments for DATAFLOW-AGENT.
- Ziyi Guo: *Contributor*; designs and conducts the operator-reuse workflow for DATAFLOW-AGENT.
- Yizhen Jiang: *Contributor*; supports the design and experiments for the Code Pipeline.
- Jingwen Deng: *Contributor*; develops the VQA extraction pipeline and operators.
- Qijie You: *Contributor*; develops the AgenticRAG pipeline and conducts experiments.
- Peichao Lai: *Contributor*; develops the frontend of DATAFLOW-WEBUI.
- Tianyu Guo: *Contributor*; develops audio-to-text operators.
- Chi Hsu Tsai: *Contributor*; fixes bugs and applies DATAFLOW to achieve first place in the BAAI LIC Challenge.
- Hengyi Feng: *Contributor*; DATAFLOW testing.
- Rui Hu: *Contributor*; conducts DATAFLOW-INSTRUCT-10K experiments.
- Wenkai Yu: *Contributor*; implements several operators.
- Junbo Niu: *Contributor*; supports the integration of MinerU into the Knowledge Cleaning Pipeline.
- Bohan Zeng: *Contributor*; supports framework design and provides serving for text-to-image components.
- Ruichuan An: *Contributor*; supports framework design and provides VQA-related component design.
- Lu Ma: *Contributor*; implements several operators.

- Jihao Huang: *Contributor*; integrates LightRAG serving.
- Yaowei Zheng: *Contributor*; integration of DATAFLOW data with LLaMA-Factory.
- Conghui He: *Project Supervisor*; project supervision and integration of MinerU with DATAFLOW.
- Linpeng Tang: *Project Supervisor*; project supervision.
- Bin Cui: *Project Supervisor*; project supervision.
- Weinan E: *Project Supervisor*; project supervision.
- Wentao Zhang: *Corresponding Author*, *Project Supervisor*; manuscript writing and project supervision.