

# ELT-Bench: An End-to-End Benchmark for Evaluating AI Agents on ELT Pipelines

Tengjun Jin  
University of Illinois (UIUC)  
Urbana, USA  
tengjun2@illinois.edu

Yuxuan Zhu  
University of Illinois (UIUC)  
Urbana, USA  
yxx404@illinois.edu

Daniel Kang  
University of Illinois (UIUC)  
Urbana, USA  
ddkang@illinois.edu

## ABSTRACT

Practitioners are increasingly turning to Extract-Load-Transform (ELT) pipelines with the widespread adoption of cloud data warehouses. However, designing these pipelines often involves significant manual work to ensure correctness. Recent advances in AI-based methods, which have shown strong capabilities in data tasks, such as text-to-SQL, present an opportunity to alleviate manual efforts in developing ELT pipelines. Unfortunately, current benchmarks in data engineering only evaluate isolated tasks, such as using data tools and writing data transformation queries, leaving a significant gap in evaluating AI agents for generating end-to-end ELT pipelines.

To fill this gap, we introduce ELT-Bench, an end-to-end benchmark designed to assess the capabilities of AI agents to build ELT pipelines. ELT-Bench consists of 100 pipelines, including 835 source tables and 203 data models across various domains. By simulating realistic scenarios involving the integration of diverse data sources and the use of popular data tools, ELT-Bench evaluates AI agents' abilities in handling complex data engineering workflows. AI agents must interact with databases and data tools, write code and SQL queries, and orchestrate every pipeline stage. We evaluate two representative code agent frameworks, Spider-Agent and SWE-Agent, using six popular Large Language Models (LLMs) on ELT-Bench. The highest-performing agent, Spider-Agent Claude-3.7-Sonnet with extended thinking, correctly generates only 3.9% of data models, with an average cost of \$4.30 and 89.3 steps per pipeline. Our experimental results demonstrate the challenges of ELT-Bench and highlight the need for a more advanced AI agent to reduce manual effort in ELT workflows. Our code and data are available at <https://github.com/uiuc-kang-lab/ELT-Bench>.

## 1 INTRODUCTION

Data engineers are increasingly leveraging Extract-Load-Transform (ELT) pipelines to integrate data and efficiently transform it into the required format as scalable cloud data warehouses become more accessible and storage prices continue to fall [16, 36, 48, 51, 52]. For example, the TPC-DI benchmark requires the creation of a decision support system for a retail brokerage firm by transforming data from various sources, including a trading system, internal Human Resources (HR), and Customer Relationship Management (CRM) systems [41]. These data sources vary in formats, data types, attributes, and inter-table relationships [41]. To build such a decision support system, data engineers can design an ELT pipeline: first, extracting and loading data into the data warehouse, followed by writing transformation queries to process the data for analysis.

Compared to traditional Extract-Transform-Load (ETL) pipelines, ELT pipelines ingest data directly into data warehouses, enabling

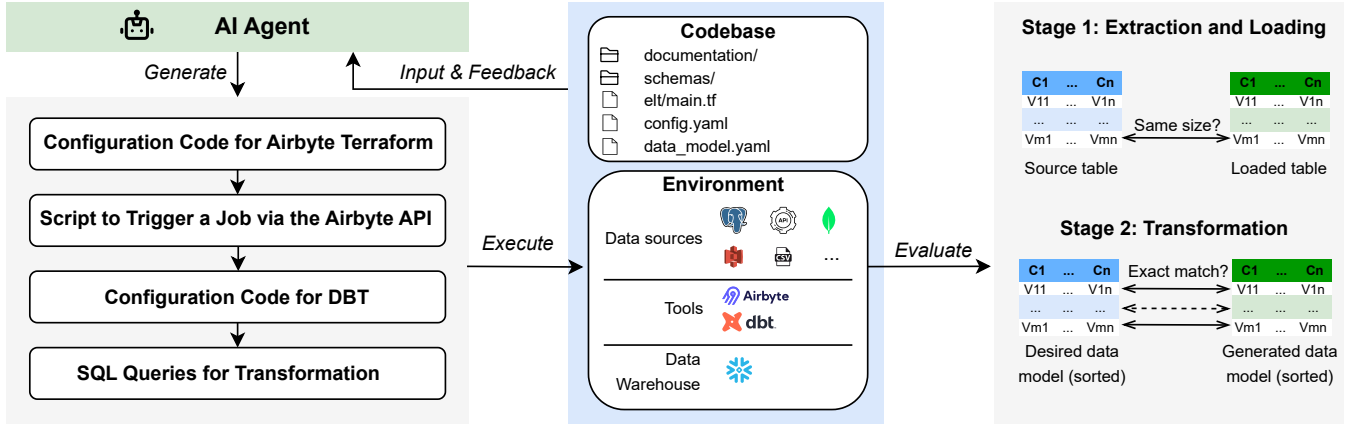
real-time Business Intelligence (BI) analysis [58]. Furthermore, with cloud infrastructure, ELT enhances scalability for processing large volumes of data [51] and offers greater flexibility in incorporating additional data transformations [42]. These benefits make ELT pipelines an increasingly preferred choice for processing data across various scenarios [16, 36, 48, 51, 52].

Developing ELT pipelines is an essential task for data engineers [16, 36, 48, 51, 52], but the process requires significant manual work. Prior studies estimate that data engineers spend over 60% of their time on data warehousing projects building data pipelines [7, 20, 27, 46, 62]. First, these pipelines must extract and integrate data from disparate sources with varying formats and standards. Second, data engineers or analysts need a deep understanding of the source data schema to write transformation queries.

Can AI agents effectively reduce the manual effort involved in constructing ELT pipelines? Recent advancements in Large Language Models (LLMs) have demonstrated strong capabilities in the text-to-SQL task, a crucial component of ELT pipelines. Notably, state-of-the-art (SOTA) techniques based on LLMs have achieved execution accuracy rates of 77.1% and 91.2% on the BIRD [33] and Spider 1.0 [71] benchmarks, respectively. Researchers have recently developed AI agents to tackle more complex real-world tasks that demand reasoning, tool usage, planning, and memorization [31, 49, 59, 61, 66, 69]. To evaluate the capability of emerging AI agents, researchers have proposed numerous benchmarks in the data domain [5, 22, 24, 30, 31]. However, there is no end-to-end benchmark designed with end-to-end ELT pipelines.

Building an end-to-end ELT benchmark is challenging because it requires sophisticated setup and configuration, time-consuming ground truth labeling, and thorough workflow verification to ensure reproducibility and correctness. First, annotators must set up various data management systems and platforms to store source data and provide data tools capable of handling diverse data sources. Second, annotators must prepare all necessary input files within the project base. Third, annotators must label ground truth by developing configuration files and writing complex transformation queries involving various relational operations (e.g., casting, type conversion, joins, aggregation, and ranking). Finally, annotators manually execute and verify each ELT pipeline to validate the correctness of both the configured environment and generated annotations.

We address these challenges by introducing ELT-Bench, a new benchmark of 100 ELT pipelines associated with 835 source tables and 203 data models across various domains. For a single ELT pipeline, we spend approximately 3 to 5 hours of manual effort setting up the environment, annotating input files and the ground truth, and building the entire pipeline for verification. Notably, 60% of the pipelines require extracting and integrating data from



**Figure 1: ELT-Bench is the first end-to-end benchmark designed to evaluate the ability of AI agents to build ELT pipelines. The agent must construct complete ELT pipelines from scratch by decomposing the complex workflow, interacting with databases and data tools, writing code and SQL queries, and calling APIs.**

five distinct categories of sources (APIs, cloud services, relational databases, NoSQL databases, and flat files). In addition, the ground truth for each pipeline, on average, involves 187 lines of code per configuration file and 200 SQL tokens (tokenized by whitespace [31]) per data model.

ELT-Bench is the first benchmark that covers the entire workflow for building ELT pipelines, providing a comprehensive evaluation through several interconnected subtasks. As shown in Figure 1, ELT-Bench requires agents to construct an end-to-end ELT pipeline from scratch, encompassing two primary stages: (1) data extraction & loading stage and (2) data transformation stage. Agents execute their actions within a sandbox environment, which includes preinstalled packages and an established project base. This setup replicates the real-world workflow of a data engineer, challenging AI agents to break down the complex workflow into manageable subtasks, interact with databases and data tools, write relevant code, orchestrate every stage of the ELT pipeline, and finally generate the required data models.

We evaluate two code agent frameworks, Spider-Agent [31] and SWE-Agent [66], with six popular LLMs on ELT-Bench. The top-performing agent, Spider-Agent Claude-3.7-Sonnet with extended thinking, achieves a success rate of 57% in the data extraction & loading stage but only a success rate of 3.9% in the data transformation stage. On average, Spider-Agent Claude-3.7-Sonnet consumes \$4.30 and requires 89.3 execution steps per task. Current agents’ poor performance and high costs highlight the need for further advancements in AI agents to reduce manual effort in developing ELT pipelines.

## 2 ELT-BENCH

In this section, we first introduce the data source of ELT-Bench, followed by summary statistics of tasks. We then provide an overview of ELT-Bench and outline the annotation pipeline. Finally, we demonstrate how an agent can complete one specific task in ELT-Bench as an example.

### 2.1 Data Collection

To ensure the quality of data, we collect databases based on a widely used text-to-SQL benchmark, BIRD [33], and the GitHub repository of an enterprise software, Fivetran [15].

- BIRD is a text-to-SQL benchmark with large-scale databases spanning 37 domains. We use all the databases that have enough natural language questions to extract features as columns in data models, leading to 78 out of 80 open-source databases. Previous study indicates that databases in BIRD can contain noise levels as high as 49% [63]. To ensure quality, we manually verify every natural language question and its corresponding SQL query used in our benchmark, correcting all identified errors.
- Fivetran is a data movement platform that develops dbt packages to facilitate the analysis of data from popular sources, such as Microsoft Advertising, Instagram Business, and YouTube Analytics. We sampled 22 databases from Fivetran.

### 2.2 Benchmark Statistics

ELT-Bench contains 100 ELT pipelines associated with 835 source tables and 203 data models. As shown in Table 1, compared to existing agent benchmarks for data engineering, ELT-Bench is the first end-to-end benchmark that covers the entire ELT pipeline construction workflow. In contrast, Spider 2-V [5] focused on evaluating an agent’s ability to use high-level data tools individually, such as using Airbyte to extract and load data and using DBT with a given SQL query to transform data. It does not include writing low-level SQL queries for data transformation or creating a complete pipeline using multiple tools. Furthermore, Spider 2.0 [31] focuses on general text-to-SQL workflows, with only 10.8% of tasks involving DBT. We highlight the characteristics of ELT-Bench as follows.

**Diverse Data sources.** As shown in Tables 2 and 3, our benchmark features diverse data sources. In total, 60 pipelines require extracting data from 5 categories of data sources, and 24 pipelines involve extracting more than 10 tables. Furthermore, 30 pipelines require

**Table 1: Comparison of ELT-Bench with two existing benchmarks in the data engineering field. ELT-Bench is the first end-to-end benchmark that focuses on building ELT pipelines, whereas Spider2-V concentrates on data tool usage and Spider 2.0 on the text-to-SQL workflow. Note: A single task may involve multiple data models and use both Airbyte and DBT, so ELT-Bench encompasses 203 data transformations and 200 data tools.**

Benchmark	# Tasks	Data Extraction & Loading	Data Transformation	Data Tools	End-to-End
Spider2-V [5]	494	✓ (48)	✗	✓ (410)	✗
Spider 2.0 [31]	632	✗	✓ (120)	✓ (68)	✗
ELT-Bench	100	✓ (100)	✓ (203)	✓ (200)	✓ (100)

**Table 2: Statistics of ELT-Bench, illustrating the distribution of data sources, source tables, lines of Terraform code, data models, and SQL tokens per data model. As shown, ELT-Bench consists of ELT pipelines that involve multiple data sources, extensive code, and complex data transformations.**

Statistics	Number
<b># Categories of Data Sources</b>	100
2 data sources	7
3 data sources	15
4 data sources	18
5 data sources	60
<b># Source Tables</b>	100
< 5 tables	36
5 - 10 tables	40
> 10 tables	24
<b># Lines of Airbyte Terraform Code</b>	100
< 100 lines	7
100 - 200 lines	63
> 200 lines	30
<b># Data Models</b>	100
1 data model	50
2 data models	22
≥ 3 data models	28
<b># SQL Tokens per Data Model</b> (Tokenized by whitespace [31])	100
< 100 tokens	8
100-200 tokens	19
> 200 tokens	73

writing more than 200 lines of code in Terraform files to extract data from these sources and load them into the data warehouse.

**Complex Data Transformation.** ELT-Bench evaluates the agent’s ability to write SQL queries based on natural language to generate target data models. Specifically, 28 pipelines require generating at least three data models. Following the approach in Spider 2.0 [31], we tokenize the SQL queries using whitespace and then count the resulting tokens to measure complexity. Because pipelines from Fivetran include both a staging and an intermediate layer, we calculate the average number of tokens per data model for each

pipeline. As shown in Table 2, 73 pipelines demand over 200 tokens per data model, illustrating the complexity of these SQL queries.

## 2.3 ELT-Bench Overview

**Task Description.** ELT-Bench requires agents building an end-to-end ELT pipeline from an existing project base, which contains connection information, target data models, an initialization file for data tools, schemas of source tables, and data tool documentation. The pipeline must extract data from a variety of sources, load it into a data warehouse, and finally transform the source data into the target data models.

**Model Inputs.** We now describe the details of the pre-established project base, which consists of the following files:

- (1) `config.yaml` contains the necessary connection information for data sources, data warehouses, Airbyte, and DBT. For example, extracting data from PostgreSQL requires specifying the host, port, user, password, schema, database, and tables.
- (2) `data_model.yaml` defines the data models the ELT pipeline generates. Each data model includes a description, column names, and explanations for each column.
- (3) `elt/main.tf` contains the code to initialize Terraform provided in Airbyte.
- (4) `documentation`: This directory includes extracted or modified documentation from Airbyte, providing guidance on writing configuration code and triggering sync jobs.
- (5) `schemas`: This directory contains the column names and descriptions of source tables.

**Environment Description.** ELT-Bench also includes a complex environment, as it requires agents to interact with a variety of data storage platforms and data tools. We describe them in detail below:

- (1) *Data sources*: We select five common categories of data sources for ELT-Bench, as shown in Table 3. We use Docker containers to deploy four data sources, including PostgreSQL, MongoDB, REST API, and Amazon S3 (simulated using LocalStack [34]), while providing downloadable links for flat files.
- (2) *Data warehouse*: We use Snowflake [8] as our data warehouse to store extracted data and execute transformation queries that generate target data models, as it is a widely studied and popular cloud data warehousing solution [31, 56].
- (3) *Data tools*: We use Airbyte [1] for data integration, a leading open-source data integration tool for ELT pipelines, which has also been used in prior work [5]. Airbyte runs in a separate

Docker container. To generate target data models, we use DBT [9], a widely adopted data transformation tool [5, 31].

- (4) *Packages and functions*: We provide a Docker file with all the required packages. Additionally, since data extraction and loading jobs typically take several minutes, we provide a script that monitors the status of all synchronization jobs and waits for their completion. This prevents redundant Airbyte API and LLM calls, reducing execution steps and costs.

## 2.4 Annotation Pipeline

We now describe the annotation pipeline of ELT-Bench, which consists of following six steps:

**Step 1: data sources conversion.** To simulate integrating various data sources in a real-world ELT pipeline, we convert the collected data into different formats based on its characteristics and the classifications in Table 3. The original BIRD data is stored in SQLite, while Fivetran data is in CSV format. Data is typically stored in multiple formats in real-world scenarios, depending on its intended use. For example, as shown in Table 3, relational databases are commonly used for transactional data storage. The selection of a target format follows these criteria:

- (1) We identify the potential data sources based on Table 3.
- (2) We select the format that maximizes the source diversity if a data source can be represented in multiple formats.
- (3) We ensure that the selected format is compatible with Airbyte.

**Step 2: data source and environment setup.** The second step involves setting up the environment for storing data in different formats. As mentioned, we use Docker containers to deploy various data sources. We write the necessary scripts for data storage, including table creation and data insertion for PostgreSQL and MongoDB, REST API implementation, and data upload scripts for cloud storage. Furthermore, because the existing Airbyte extractor does not support local APIs, we have developed a custom extractor that can retrieve data from a REST API running inside a Docker container. We only need to perform all of these setup steps once before running the experiments.

**Step 3: configuration annotation.** After converting data formats and setting up the environment, we annotate the necessary connection information for data storage platforms and tools, which serve as one of the inputs in ELT-Bench. Notably, the annotated configurations do not strictly match the field names in the Airbyte documentation. This setting increases complexity and requires the agent’s reasoning capabilities. For example, when extracting data from MongoDB, we specify the connection string as follows:

```
mongodb:
  config:
    connection_string: 'mongodb://elt-mongodb:27017/?
      directConnection=true'
```

The agent must determine that MongoDB is self-managed and configure Airbyte based on the provided configuration:

```
configuration = {
  database_config = {
    self_managed_replica_set = {
      connection_string = "mongodb://eltmongodb:27017/
        directConnection=true"}}
```

**Step 4: data model definition.** In this step, we annotate the columns and their corresponding descriptions of data models in each database. Each data model consists of descriptive attributes from dimension tables and quantitative metrics from fact tables, representing specific entities. We categorize the columns in the data model based on the transformations applied:

- (1) *Derived columns* come from either direct copies of the columns in source tables or transformations through basic operations (e.g., renaming, concatenation, and mathematical operations).
- (2) *Aggregated columns* summarize data from the fact table using aggregation functions such as SUM, AVG, COUNT, MAX and MIN.
- (3) *Categorical columns* classify data into predefined categories based on thresholds or conditions.

*Example:* The `has_more_than_10_movies_1960_to_1985` is set to 1 if the director has directed at least 10 movies between 1960 to 1985; otherwise set to 0.

```
SELECT director_id, CASE WHEN COUNT(*) > 10 THEN 1
  ELSE 0 END AS has_more_than_10_movies_1960_to_1985
FROM movies
WHERE movie_release_year BETWEEN 1960 AND 1985
GROUP BY director_id;
```

- (4) *Ranked Columns* apply ranking functions (e.g., RANK()) to establish an order based on specific criteria and extract a value at a particular rank position.

*Example:* The `highest_average_score_film` column shows the film directed by a director with the highest average score, with ties broken by the ascending order of the movie title.

```
WITH rated_film_ranks AS (
  SELECT director_id, movie_title,
    RANK() OVER (PARTITION BY director_id
      ORDER BY AVG(rating_score) DESC,
        movie_title ASC) AS rated_film_rank
  FROM movie_platform.movies AS T1
  JOIN movie_platform.ratings AS T2
    ON T1.movie_id = T2.movie_id
  GROUP BY director_id, movie_title)
SELECT director_id,
  movie_title AS highest_average_score_film
FROM rated_film_ranks
WHERE rated_film_rank = 1
```

We now describe how we extract these four types of columns for our data models based on the natural language questions of BIRD. While the original questions are typically designed for a specific entity, such as a single movie or person, we extract features for all entities instead. For example, consider the question, "Which film directed by Abbas Kiarostami has the highest average score?" This corresponds to the extracted feature `highest_average_score_film` in the Directors data model, which represents the film with the highest average score for each director. After extracting features, we add text descriptions for these columns to help the agent better understand the desired data model and reduce ambiguity.

To make the data transformation stage in ELT-Bench more challenging, we rank the original SQL queries by complexity and prioritize features that involve more SQL components, conditions, and table joins [71]. Each data model typically consists of three derived columns from the source tables and five additional columns (i.e., aggregated columns, categorical columns, and ranked columns) extracted from BIRD questions. Since the databases from Fivetran already contain predefined data models, we directly refine these

**Table 3: Overview of common data source categories, representative sources, and their real-world applications.**

Data Source Category	Representative Sources	Applications in Practice
APIs	REST API	Web services, third-party platforms, real-time applications.
Cloud Services	Amazon S3	Big data platforms, modern applications.
Relational Databases	PostgreSQL	Traditional enterprise systems, transactional systems.
NoSQL Databases	MongoDB	Modern web applications, real-time data systems.
Flat Files	CSV, JSONL, Parquet	Third-party data providers, backups.

models by removing columns generated by utility functions and those that contain only null values, unless they are required by another data model (a data model may incorporate another data model as an intermediate stage).

**Step 5: ground-truth annotation.** We annotate the ground truth using the configuration file and the defined data models. First, we review the official documentation to understand each configuration field and implement the necessary code accordingly. Next, we annotate the SQL queries used to generate the defined data models in Step 4. To achieve this, we initially validate the existing queries provided by the BIRD benchmark and the Fivetran repository. We modify those valid queries to conform to our defined data models; otherwise, we write gold queries from scratch.

**Step 6: execution-based verification.** To ensure the quality and correctness of ELT-Bench, we manually execute and verify each ELT pipeline, thoroughly confirming the accuracy of environment configurations and annotations.

In the first stage, we validate whether Airbyte can correctly extract data from diverse sources and load it into the data warehouse based on the provided configurations. Since Airbyte is an actively developing project, we encounter several issues:

- (1) *Table name case sensitivity in PostgreSQL:* Airbyte automatically converts table names in PostgreSQL to lowercase, which can cause a table not found error if the provided configuration contains table names with uppercase letters.
- (2) *Schema detection for API data sources:* Airbyte’s automatic schema detection fails when the source data exceeds the maximum allowed string length.

To address these issues, we standardize table names in PostgreSQL to only use lowercase letters and manually define schemas for affected API data sources.

In the data transformation stage, we verify our annotated SQL transformation queries. Specifically, for each defined data model, we write ten additional representative testing queries on average and then execute them against both the original source tables and corresponding data models. These testing queries encompass: (1) aggregation queries, (2) limit queries, and (3) queries obtained from BIRD corresponding to entity-specific questions. For any discrepancies or mismatches observed during query execution, we carefully review the annotated transformation queries, correct identified errors, and revise the transformation queries accordingly.

## 2.5 Task Example

We describe the process of building an ELT pipeline to generate a customers data model extracted from the retailers database to demonstrate the task. The pipeline involves all five types of data sources. Starting in a sandbox environment with all required packages, we divide the task into two stages:

- (1) *Data extraction & loading stage:* We use Airbyte to extract data from all data sources and load it into Snowflake. First, we initialize Airbyte and write Terraform codes to configure Airbyte, all data sources, Snowflake, and the connections between data sources and Snowflake using information from the config.yaml file (Figure 2a) and relevant documentation. We illustrate the code for configuring the flat file nation.jsonl and establishing its connection to Snowflake (Figure 3a). Next, we execute `terraform apply`, which will apply the configuration code and generate configuration information in an output file. Finally, we extract connection IDs from the generated output file, trigger synchronization jobs via the Airbyte API, and monitor their status until completion.
- (2) *Data transformation stage:* If all synchronization jobs in the first stage complete successfully, we proceed to data transformation using DBT. First, based on the provided configuration config.yaml (Figure 2a), we initialize a DBT project configured for Snowflake and set all necessary parameters (Figure 3b). Next, by referring to the customers data model definitions specified in data\_model.yaml (Figure 2b), we develop the transformation query (Figure 3c). After executing `dbt run` to generate the data model, we validate the output by running `SELECT * FROM customers` in Snowflake to identify any issues that DBT might have missed.

## 3 EXPERIMENTS

We evaluate two representative code agent frameworks, Spider-Agent [31] and SWE-Agent [66], using six LLMs on ELT-Bench. In this section, we first introduce the evaluation metrics of ELT-Bench, followed by a detailed explanation of the experimental settings for both agents. Finally, we present the evaluation results.

### 3.1 Evaluation Metrics

We use the widely adopted metric, success rate [5, 23, 31, 74], to assess the performance of agents on ELT-Bench. To provide a more comprehensive evaluation, we measure the success rate for both the data extraction & loading stage and the data transformation stage. Specifically, we introduce the **Success Rate for Data Extraction &**



**Table 4: ELT-Bench evaluation results for all tested agents and LLMs. Spider-Agent Claude-3.7-Sonnet with extended thinking performs best, with a 57% SRDEL and 3.9% SRDT.**

Agent Framework	LLM	SRDEL (%)	SRDT (%)	Average Cost (\$)	Average Steps
Spider-Agent	Claude-3.5-Sonnet	23%	0	3.51	63.3
	GPT-4o	15%	0	2.03	43.7
	DeepSeek-R1	0	0	0.38	18.4
	Llama-3.1-405B	0	0	0.39	22.0
	Qwen2.5-Coder-32B-Instruct	0	0	0.50	37.3
SWE-Agent	Claude-3.5-Sonnet	37%	1%	5.22	60.0
	GPT-4o	0	0	5.22	114.3
	DeepSeek-R1	0	0	3.16	66.9
	Llama-3.1-405B	0	0	2.90	73.9
	Qwen2.5-Coder-32B-Instruct	0	0	0.48	39.1
Spider-Agent	Claude-3.7-Sonnet w/ extended thinking	57%	3.9%	4.30	89.3

```
Airbyte:
  config:
    files_definition_id: <id_1>
    workspace_id: <id_2>
  flat_files:
  - format: jsonl
    path: "https://..."
    sync_mode: full_refresh_append
    table: nation
  snowflake:
    config:
      account: <account_id>
      database: retails
      password: <snowflake_password>
      username: AIRBYTE_USER ...
```

(a) A partial configuration of Airbyte for the retails database, as defined in the provided config.yaml file.

```
models:
- name: customers
  description: Each record represents a customer.
  columns:
  - name: c_custkey
    description: Unique identifier for the customer.
  - name: order_date_highest_total_price
    description: The order date with the highest total
      price the customer has made, with ties broken by
      the ascending order of the order date. ...
```

(b) A partial description of the customers data model for the retails database, as defined in the provided data\_model.yaml file.

**Figure 2: An example of provided input files of the retails database in ELT-Bench.**

**Loading (SRDEL)** to measure the proportion of ELT pipelines that successfully extract and load data in the first stage and the **Success Rate for Data Transformation (SRDT)** to measure the proportion of data models successfully built in the second stage. Additionally, we measure the agent’s **average cost** (calculated based on token usage and API pricing [2, 14, 38]) and **average steps** per task to assess its efficiency. We describe SRDEL and SRDT below.

**SRDEL.** We evaluate the metric SRDEL in the first stage:

$$\text{SRDEL} = \frac{\# \text{ successful pipelines in data extraction \& loading}}{\# \text{ total pipelines}},$$

which measures the proportion of pipelines that successfully extract and load data.

A pipeline is considered successful in the data extraction & loading stage if the pipeline successfully extracts data from all sources and loads it into the data warehouse. To evaluate this, we execute the following query for each source table in the data warehouse:

```
SELECT COUNT(*) FROM source_table;
```

We then verify whether the query result matches the corresponding row count in the original data.

**SRDT.** To evaluate the performance of the agent in the second stage, we use the metric SRDT:

$$\text{SRDT} = \frac{\# \text{ correctly generated data models}}{\# \text{ total data models}},$$

which measures the proportion of correctly generated data models among all data models (one ELT pipeline may involve multiple data models). To assess the correctness of a generated data model, we execute the following query:

```
SELECT * FROM data_model ORDER BY unique_key;
```

The unique key may consist of a composite set of columns determined manually for each data model to ensure the query produces consistent results across different runs. We use this query to create CSV files for the generated data model, which are then compared against the ground truth, which is also derived from the same query.

A generated data model is correct if it contains all columns of the ground truth. Following prior work [31], we permit extra columns in the generated data model since they do not affect functionality.

## 3.2 AI Agent Frameworks

We select two representative code agent frameworks: Spider-Agent [31] and SWE-Agent [66] since ELT-Bench requires capabilities

```
resource "airbyte_source_file" "jsonl_file_nation" {
  name           = "JSONL File nation"
  definition_id   = "<id_1>"
  workspace_id    = "<id_2>"
  configuration = {
    dataset_name = "nation"
    format       = "jsonl"
    provider     = {https_public_web = {}}
    url          = "https://..."
  }
resource "airbyte_connection" "nation_to_snowflake" {
  name = "JSONL nation to Snowflake"
  configurations = {
    streams = [
      {name = "nation"
        sync_mode = "full_refresh_append"}}] ...
}
```

(a) Partial code of Airbyte Terraform to configure the flat file data source and its connection to Snowflake for the retails database.

```
my_dbt_profile:
  target: dev
  outputs:
    dev:
      type: snowflake
      account: <account_id>
      user: AIRBYTE_USER
      password: <snowflake_password>
      database: retails ...
```

(b) Example DBT configuration for the retails database in Snowflake.

```
WITH order_date_highest_total_price AS (
  SELECT o_custkey, o_orderdate,
    RANK() over(PARTITION by o_custkey
      ORDER BY o_totalprice DESC, o_orderdate) AS price_rank FROM retails.airbyte_schema.orders)
SELECT T1.c_custkey AS c_custkey,
  T2.o_orderdate AS order_date_highest_total_price, ...
FROM retails.airbyte_schema.customer T1
LEFT JOIN order_date_highest_total_price T2 ON T1.
  c_custkey = T2.o_custkey
AND T2.price_rank = 1 ...
```

(c) Partial code of the SQL transformation query for generating the customers data model.

Figure 3: Files required to build the ELT pipeline.

such as file viewing and editing, code generation, and command execution. As baseline evaluations, we combine these two agent frameworks with five LLMs, including GPT-4o [39], Claude-3.5-Sonnet [3], two open-sourced LLMs (Llama-3.1-405B-Instruct [19], Qwen2.5-Coder-32B-Instruct [45]), and one reasoning model (DeepSeek-R1 [10]). In addition, as a case study aimed at exploring the frontier reasoning model, we also evaluate Spider-Agent Claude-3.7-Sonnet with extended thinking on ELT-Bench.

**Spider-Agent.** Spider-Agent is a code agent framework designed for database-related tasks, providing command-line interfaces for multi-turn interactions with environments [31]. It also enables direct interaction with databases to extract detailed source table information (e.g., column values) and verify the correctness of transformation queries (e.g., DBT may fail to detect format errors). The agent employs the ReAct [69] framework, in which the LLM generates thought and decides the next action based on current

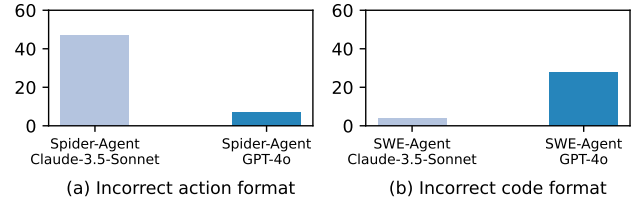


Figure 4: The number of tasks with incorrect formats.

observation and history trajectory at each iteration. We use the default parameter settings of Spider-Agent, except for changing the maximum allowed steps to 100, as ELT-Bench presents more challenging tasks compared to Spider 2.0 [31].

**SWE-Agent.** SWE-Agent is a code agent framework designed to address GitHub issues [66]. Compared to Spider-Agent, it does not allow direct database access. In each iteration, the agent interacts with the filesystem based on its observations. SWE-Agent operates as a function-calling agent by prompting the LLM to invoke predefined functions, and it also offers a thought-action mode for LLMs that lack native support for tool usage. Consequently, we employ function calling for GPT-4o and Claude-3.5-Sonnet, and the thought-action approach for Llama-3.1-405B-Instruct, Qwen2.5-Coder-32B-Instruct, and DeepSeek-R1. We apply the default parameter settings of SWE-Agent, with one modification: retaining the last 25 observations for the agent due to the complexity of ELT-Bench. Following prior work [66], we allocate a same cost budget to all evaluated LLMs. To establish a comparable budget for both SWE-Agent and Spider-Agent, we first estimate the cost of completing 100 agent steps using Spider-Agent across all LLMs. We then select the highest of these estimates and round it up to the nearest integer, yielding a budget of \$6 for each evaluated LLM.

### 3.3 Evaluation Results

We report our evaluation metrics for all evaluated agents and LLMs in Table 4. The poor performance, high cost, and extensive step requirements highlight the challenges of ELT-Bench. The top-performing agent, Spider-Agent Claude-3.7-Sonnet with extended thinking, attains a 57% success rate for data extraction & loading, but only a 3.9% success rate for data transformation. Despite these limitations, this agent demonstrates substantial performance improvements over the best-performing agent employing non-reasoning models, SWE-Agent Claude-3.5-Sonnet, with 54.1% improvement in the data extraction and loading stage and 290% improvement in the data transformation stage. Moreover, ELT-Bench presents a significantly higher computational cost compared to Spider 2.0 [31]. While 30 agent steps are sufficient for most tasks in Spider 2.0, with an average cost of \$0.30 per instance using Spider-Agent GPT-4o, evaluating Spider-Agent GPT-4o on ELT-Bench requires an average of 43.7 steps and costs \$2.03 per task.

We present a detailed error analysis for the baseline agent evaluations in Section 4, followed by an in-depth case study of Spider-Agent Claude-3.7-Sonnet in Section 5.

```
Action: CreateFile(filepath='/root/.dbt/profiles.yml ':  
    ```retail_complains: ...```)
```

(a) Incorrect action format.

```
provider "airbyte" {username = "<username>"}
```

(b) Incorrect code format.

**Figure 5: Incorrect action format generated by Spider-Agent Claude-3.5-Sonnet and incorrect code format generated by SWE-Agent GPT-4o.**

## 4 ERROR ANALYSIS

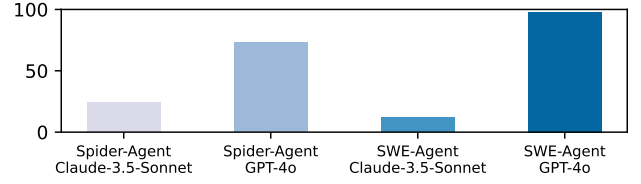
In this section, we examine the errors encountered by different agents and LLMs. We first highlight common issues observed among open-source LLMs. Then, we provide a detailed analysis of the errors arising in the data extraction & loading stage and the data transformation stage for Spider-Agent and SWE-Agent using Claude-3.5-Sonnet and GPT-4o.

### 4.1 Error Analysis of Open-Sourced LLMs

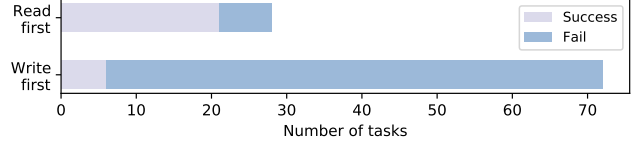
In ELT-Bench, the environment starts with a project base requiring agents to use official documentation extracted from the Airbyte Terraform website [54], reflecting a realistic scenario where data engineers must learn from documentation to use data tools, especially for those underdeveloped tools. However, open-sourced LLMs struggle to interact with the provided project base, resulting in a 0% success rate in the data extraction & loading stage. In addition, the complexity of ELT-Bench necessitates that the agent maintain a substantial memory length, leading to excessive prompt length issues when employing Qwen2.5-Coder-32B-Instruct.

**Failure to check configuration information.** We first analyzed the performance of Spider-Agent DeepSeek-R1 and found that it failed primarily by neglecting the provided `config.yaml`. Spider-Agent DeepSeek-R1 references `config.yaml` in only five tasks. Instead, in most cases, it generates the configuration with random values, which results in execution errors in Airbyte. Even in the five tasks where `config.yaml` is used, the agent either produces an incomplete configuration or inserts its generated actions directly into the file rather than executing the actions.

**Failure to consult documentation files.** We further analyzed the remaining agents’ performances and frequently observed a failure to reference the provided documentation. Without consulting up-to-date documentation, agents instead generate configuration code based on outdated versions of Airbyte’s documentation, resulting in errors related to incorrect resource types, as exemplified by Action 6 in Figure 8b. Specifically, we tracked how frequently they consulted the Snowflake configuration guide. Despite explicit prompt instructions, neither Spider-Agent nor SWE-Agent, when using Qwen2.5-Coder-32B-Instruct, consult this documentation. In contrast, Spider-Agent Llama-3.1-405B references the guide three times, SWE-Agent Llama-3.1-405B references it 20 times, and SWE-Agent DeepSeek-R1 refers to it 22 times. Consequently, SWE-Agent Llama-3.1-405B successfully configures the Snowflake destination



**Figure 6: The number of tasks with incorrect Snowflake password field.**



**Figure 7: The success and failure rates of Spider-Agent GPT-4o in Stage 1 under two strategies: when reading the documentation first (27% of tasks), it achieves a 78.8% success rate; when writing the configuration first (73% of tasks), the success rate drops to 9.6%.**

in 9 of the 20 instances. Similarly, SWE-Agent DeepSeek-R1 correctly configures this destination in 10 of the 22 instances. However, these agents encountered other issues such as incorrect action formatting and failures in triggering synchronization jobs.

**Excessive prompt length.** Another common issue for agents using Qwen2.5-Coder-32B-Instruct is the excessive prompt length error [31], caused by exceeding its maximum supported context length of 32,767 tokens. Specifically, we observe this issue occurring in 15 tasks with Spider-Agent Qwen2.5-Coder-32B-Instruct and 87 tasks with SWE-Agent Qwen2.5-Coder-32B-Instruct. This discrepancy arises because Spider-Agent retains only the last 25 steps of model outputs and environment observations, whereas SWE-Agent retains all model outputs starting from the first step, along with only the last 25 environment observations.

### 4.2 Error Analysis of Data Extraction & Loading

We examined the common issues encountered by Spider-Agent and SWE-Agent when using GPT-4o and Claude-3.5-Sonnet, including incorrect action or code formats, incorrect Snowflake password fields, incorrect table sizes, and missing configuration for multiple flat files. We describe each of these problems in detail below.

**Failure to generate action or code in the required format.** We observed that LLMs frequently generate actions in incorrect formats when used with Spider-Agent, and similarly produce code in incorrect formats under SWE-Agent. Spider-Agent, which builds on the ReAct framework [69], mandates that the LLM produce a valid action at each iteration, terminating the process if three formatting errors are detected via regular expressions. Despite providing detailed descriptions and examples of all actions, LLMs can still generate invalid actions. For instance, in Figure 5a, Spider-Agent Claude-3.5-Sonnet incorrectly places a closing parenthesis



### Success case

```

Action 1: Bash(code="cat /workspace/config.yaml")
Action 10: Bash(code="cat /workspace/documentation/destination_snowflake.md")
Action 11: EditFile(filepath="/workspace/elt/main.tf")
```resource "airbyte_destination_snowflake" "snowflake" {
  configuration = {credentials = {
    username_and_password = {password = "Snowflake@123"}}}}
(Omitted other fields)```
Action 14: Bash(code="terraform apply -auto-approve")
...

```

(a) The execution path of a successful task in Stage 1. The agent writes the configuration after reading the documentation.

### Failure case

```

Action 1: Bash(code="cat /workspace/config.yaml")
Action 6: EditFile(filepath="/workspace/elt/main.tf")
```resource "airbyte_destination" "snowflake" {
  password = "Snowflake@123"}
(Omitted other fields)```
Action 8: Bash(code="terraform apply -auto-approve")
Observation-8: The provider airbytehq/airbyte does not support resource type "airbyte_destination"
Action 14: Bash(code="cat /workspace/documentation/destination_snowflake.md")
```resource "airbyte_destination_snowflake" "snowflake"{
  password = "Snowflake@123"}
(Omitted other fields)```
...

```

(b) The execution path of a failed task in Stage 1. The agent writes the configuration file before reading the documentation and only fixes the detected error after reading the documentation.

Figure 8: Comparison the execution path of a successful task and a failed task.

on the last line instead of before the colon in the first line, causing a parsing failure. As shown in Figure 4, Spider-Agent Claude-3.5-Sonnet terminates 47% of tasks in Stage 1 because of unparseable actions, while Spider-Agent GPT-4o only terminates 7% of tasks.

In contrast, SWE-Agent employs a function-calling framework with well-defined functions, which avoids parsing errors. However, misformatted code can still be generated. For instance, SWE-Agent GPT-4o generates an extra right curly brace at the end of a code block, as illustrated in Figure 5. We show in Figure 4 that SWE-Agent Claude-3.5-Sonnet produces misformatted code in 4% of cases, whereas SWE-Agent GPT-4o exhibits a 28% error rate. These findings highlight the importance of developing frameworks that robustly ensure LLMs generate correct syntax for actions and code.

**Failure to configure the Snowflake password field.** We examined the Snowflake password field, which must be written in the

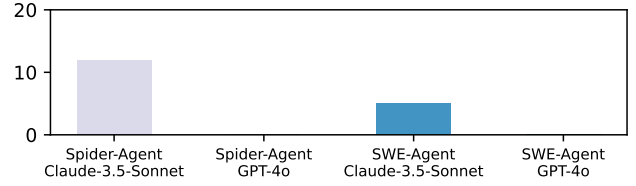


Figure 9: The number of tasks with incorrect table size.

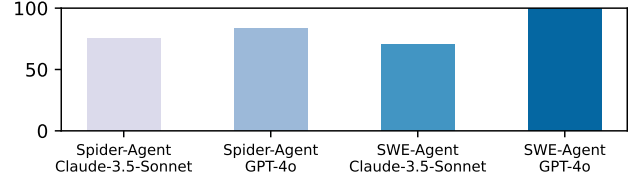
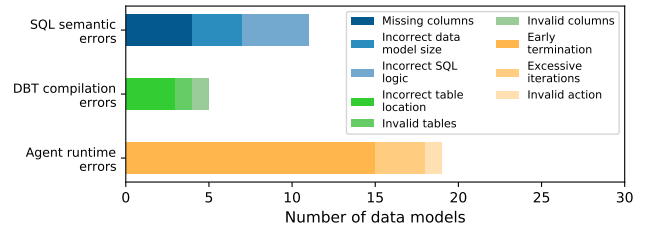
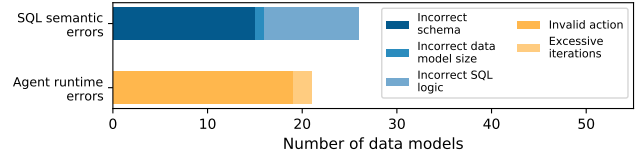


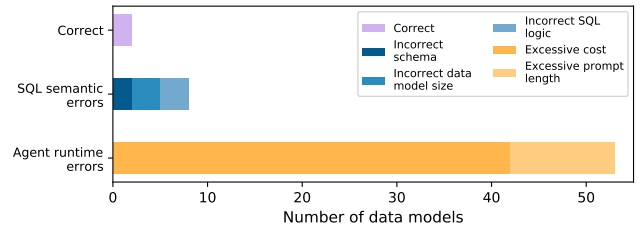
Figure 10: The failure rate of configuring multiple flat files. The total number of tasks with multiple flat files is 24.



(a) Statistics of Spider-Agent GPT-4o in the second stage.



(b) Statistics of Spider-Agent Claude-3.5-Sonnet in the second stage.



(c) Statistics of SWE-Agent Claude-3.5-Sonnet in the second stage.

Figure 11: Statistics of agent performance on generating data models in the second stage. Each subfigure includes results for databases that the agent successfully completed in the first stage (35, 47, and 63 data models, respectively).

format as shown in Figure 8a. However, as Figure 6 illustrates, SWE-Agent GPT-4o fails in up to 98% of tasks to configure the Snowflake password field. In contrast, SWE-Agent Claude-3.5-Sonnet performs much better, failing to configure it in only 12% of tasks.

We further analyzed the execution path of Spider-Agent GPT-4o and identified two distinct strategies the agent adopted when configuring Airbyte Terraform. In one strategy (Figure 8a), the agent attempts to write the configuration code first and then runs `terraform apply -auto-approve`. Upon encountering an error indicating an incorrect resource type, the agent consults the documentation but only corrects the specific issue reported by Terraform. Because Airbyte Terraform ignores any fields that are not explicitly defined, other misconfigurations remain undetected, which finally causes the ELT pipeline to fail.

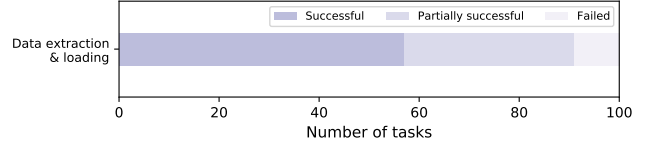
In contrast, when the agent references the documentation before writing the configuration, it is more likely to produce a valid Terraform configuration, leading to a higher success rate for data extraction & loading. As illustrated in Figure 7, the agent reads the documentation before writing the configuration in 27 tasks and successfully configures the Snowflake password field in 21 tasks. By comparison, in 73 tasks, the agent writes the configuration first, and only six tasks succeed. These observations underscore the importance of the agent’s effective planning (e.g., executing actions in the correct sequence) in achieving higher success rates.

**Incorrect table size due to repeated synchronization job triggers.** We observed that, in some cases, the size of the source tables did not match the size of the original data. Analyzing the execution paths of failed cases, we found that the agent repeatedly triggered the same synchronization job. For example, if the original dataset contains 100 rows but the agent executes the synchronization job three times, the resulting table in Snowflake ends up with 300 rows instead of the intended 100. As shown in Figure 9, Spider-Agent Claude-3.5-Sonnet repeatedly triggers the same synchronization job in 12 tasks. These findings highlight the importance of short-term memorization in the agent for tracking executed actions and preventing redundant synchronization jobs.

**Missing configuration for multiple flat files.** For Postgres, MongoDB, APIs, and Amazon S3, multiple tables or files can be configured within a single source block and a single connection block. In contrast, Airbyte Terraform requires individual source and connection configuration blocks for each flat file. ELT-Bench includes 24 instances to evaluate whether the agent can correctly generate multiple configuration blocks for multiple flat files. As shown in Figure 10, SWE-Agent GPT-4o fails on all 24 instances, while even the best-performing agent, SWE-Agent Claude-3.5-Sonnet, still fails in 70.8% of cases. These findings emphasize the need for the agent to handle diverse configuration patterns across different data sources.

### 4.3 Error Analysis of Data Transformation

We evaluated the data transformation stage performance of Spider-Agent GPT-4o, Spider-Agent Claude-3.5-Sonnet, and SWE-Agent Claude-3.5-Sonnet, as these agents completed the data extraction & loading stage for some databases. Specifically, Spider-Agent GPT-4o successfully processes 15 databases in Stage 1, comprising 35 data models; Spider-Agent Claude-3.5-Sonnet processes 23 databases,



**Figure 12: Task completion status in the data extraction & loading stage. Spider-Agent Claude-3.7-Sonnet fails on all sources in nine tasks and partially succeeds in 34 tasks.**

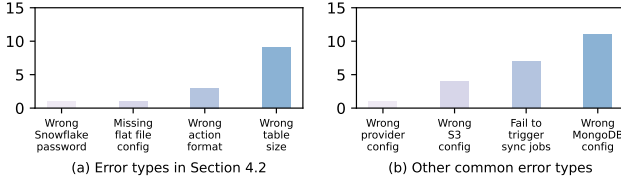
comprising 47 data models, while SWE-Agent Claude-3.5-Sonnet processes 37 databases, comprising 63 data models.

However, only SWE-Agent Claude-3.5-Sonnet successfully generates two correct data models. We categorize Stage 2 errors into three main types: agent runtime errors, DBT compilation errors, and SQL semantic errors, and show the performance in Figure 11.

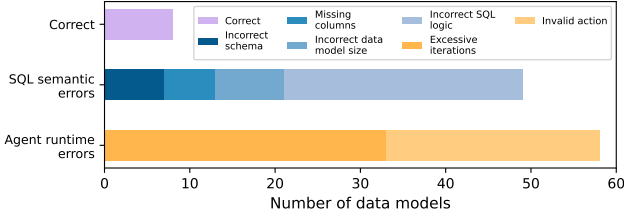
**Agent runtime errors.** Agent runtime errors primarily result from configuration or action-related issues that cause the agent to terminate before generating a data model. For example, due to a \$6 maximum cost constraint, SWE-Agent Claude-3.5-Sonnet fails to create 42 data models within the allotted budget. In addition, Spider-Agent GPT-4o fails to generate 15 data models because it prematurely halts in Stage 1. Among these prematurely terminating cases, 46.7% stop because, although the agent correctly assumes it should proceed to Stage 2, it mistakenly executes a terminate action instead. Meanwhile, 53.3% fail due to incorrect synchronization job status checks. As shown in Figure 11c, agent runtime errors account for up to 84.1% of failed data model generations.

**DBT compilation errors.** We observed DBT compilation errors in Spider-Agent GPT-4o, which were identified by DBT during dbt run. As shown in Figure 11a, these errors occur because the agent fails to correctly specify the corresponding database and schema for tables used in the query (3 data models) or references non-existent tables (1 data model) or columns (1 data model).

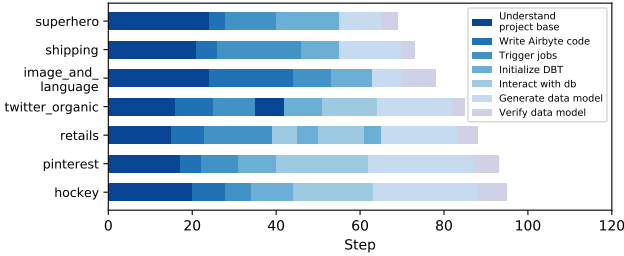
**SQL semantic errors.** SQL semantic errors occur when the agent successfully creates a data model in the database, but the model is incorrect. We classify the semantic errors of data models in ascending severity: incorrect schema, missing columns, incorrect data model size, and flawed SQL logic. For example, if a data model is placed in the wrong schema and omits some columns, we categorize it as an incorrect schema error. In Spider-Agent GPT-4o (Figure 11a), the most common SQL semantic errors (11 data models) are due to missing columns (4 data models) and flawed SQL logic (4 data models). For Spider-Agent Claude-3.5-Sonnet, 26 data models contain SQL semantic errors, with incorrect schema assignments (15 data models) and flawed SQL logic (10 data models), as shown in Figure 11b. Similarly, as illustrated in Figure 11c, SWE-Agent Claude-3.5-Sonnet produces 8 data models with SQL semantic issues, including incorrect schema (2 data models), incorrect data model size (3 data models), and flawed SQL logic (3 data models).



**Figure 13: Common error types encountered by Spider-Agent Claude-3.7-Sonnet in the first stage.**



**Figure 14: Statistics of Spider-Agent Claude-3.7-Sonnet in the second stage.**



**Figure 15: The action trajectories of the agent on databases with at least one successful data model.**

## 5 CASE STUDY

In this section, we present an in-depth analysis of the Spider-Agent Claude-3.7-Sonnet with extended thinking, focusing on its performance and the errors encountered across two stages of the task. We then examine its action trajectories in successful cases.

Spider-Agent Claude-3.7-Sonnet achieves a significant performance improvement of 54.1% over the second-best agent, SWE-Agent Claude-3.5-Sonnet, in the first stage. As shown in Figure 12, Spider-Agent Claude-3.7-Sonnet fails on all data sources in the first stage for nine tasks. We further analyzed common error types during the first stage, with results depicted in Figure 13. Our initial examination of the four issue types described in Section 4.2 reveals that Spider-Agent Claude-3.7-Sonnet significantly reduced error frequencies across all categories compared to Spider-Agent Claude-3.5-Sonnet, achieving up to a 95.8% error reduction. Further analysis of additional common issues, shown in the right part of Figure 13, indicates that Spider-Agent Claude-3.7-Sonnet frequently fails on specific data source types, particularly MongoDB. This finding aligns with the partial success of 34 tasks observed in Figure 12.

Spider-Agent Claude-3.7-Sonnet demonstrates a 290% performance improvement in the second stage compared to SWE-Agent Claude-3.5-Sonnet. As illustrated in Figure 14, the primary issues of Spider-Agent Claude-3.7-Sonnet in the second stage include excessive iterations (28.7%), incorrect SQL logic (24.3%), and invalid actions (21.7%).

To better understand Spider-Agent Claude-3.7-Sonnet’s workflow, we illustrate the action paths of the agent for databases that successfully produced at least one correct data model in Figure 15. On average, the agent executed 83.6 steps for each successful case. To provide clarity, we categorize these actions into defined phases based on the agent’s thoughts and actions. Specifically, if fewer than five consecutive steps belonging to one phase appear between two occurrences of another identical phase, we group these intermediate steps into the surrounding phase. For instance, it is common for the agent to briefly interact with the database during the “generate data model” phase. As depicted in Figure 15, the Spider-Agent Claude-3.7-Sonnet spends most of its execution steps to the phases of “understanding the project base” (averaging 20.6 steps) and “generating the data model” (averaging 17 steps).

## 6 SENSITIVITY AND ABLATION STUDY

### 6.1 Multiple Runs Improve Agent Performance

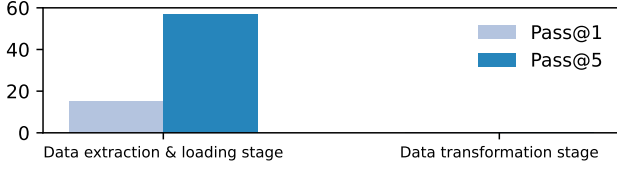
We evaluated Spider-Agent GPT-4o’s performance on ELT-Bench with one attempt (pass@1) and five attempts (pass@5). As shown in Figure 16, Spider-Agent GPT-4o achieves a pass@5 rate of 57% in Stage 1, indicating that in 57% of tasks, at least one of the five attempts successfully extracts data from multiple sources and loads it into the data warehouse. This result represents a 3.8× improvement over its pass@1 performance. However, in Stage 2, despite having more successfully loaded source tables, Spider-Agent GPT-4o still fails to build a correct data model.

We further use the pass<sup>k</sup> metric [67] to evaluate the consistency and robustness of Spider-Agent GPT-4o on ELT-Bench. As shown in Figure 17, as the number of trials increases, pass<sup>k</sup> for Spider-Agent GPT-4o drops significantly, eventually reaching 0 when k equals 5, indicating the need for a more robust agent in future work.

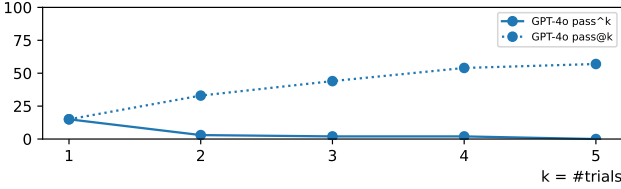
### 6.2 Using Documentation Improves Agent Performance

We evaluated whether Spider-Agent Claude-3.5-Sonnet and Spider-Agent GPT-4o could complete the data extraction & loading stage without consulting documentation. Since LLMs are trained on a fixed knowledge cutoff, their ability to reference up-to-date documentation is crucial for completing real-world tasks. To assess their adaptability, we compared their performance in data extraction & loading with and without documentation guidance.

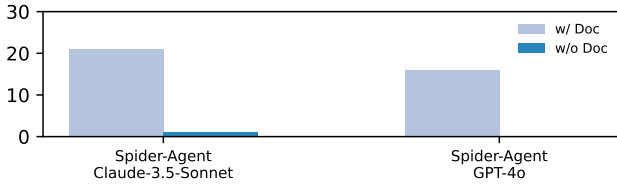
In our experiments, we provided the agents with documentation on configuring Airbyte Terraform and invoking the Airbyte API to initiate synchronization jobs. As shown in Figure 18, Spider-Agent Claude-3.5-Sonnet and Spider-Agent GPT-4o exhibit degraded performance in the data extraction & loading stage when documentation is unavailable. Without access to documentation, Spider-Agent Claude-3.5-Sonnet succeeds in only one task, while Spider-Agent GPT-4o fails in all tasks. These findings reveal that both



**Figure 16: The success rate of Spider-Agent GPT-4o with one versus five attempts. The success rate improves from 15% to 57% in the first stage but remains 0% in the second stage.**



**Figure 17: Pass<sup>k</sup> and pass@k in the first stage of ELT-Bench.**



**Figure 18: The success rate of Spider-Agent with Claude-3.5-Sonnet and GPT-4o in the data extraction & loading stage, evaluated with and without documentation. Success rates decrease from 21% to 1% for Claude-3.5-Sonnet, and from 15% to 0% for GPT-4o.**

Claude-3.5-Sonnet and GPT-4o rely not only on memorized knowledge but also on their reasoning abilities to complete tasks.

## 7 RELATED WORK

**ELT and ETL data pipelines.** ELT and ETL data pipelines are essential for converting raw data into structured, reliable formats, playing an important role in modern data engineering workflows. ETL techniques have been extensively studied over decades [51], while the rise of cloud data warehousing has driven the increasing adoption of ELT pipelines [16, 36, 48, 52]. Early research mainly focus on conceptual modeling for ETL processes [35, 55, 57]. More recent efforts have aimed at automating various stages of ETL and ELT pipelines to minimize engineering effort, including Semantic Web-based approaches for attribute mapping [53], template-driven automatic data loading [6], and machine learning-based data integration [37]. In this work, we introduce ELT-Bench, a benchmark designed to facilitate the development of AI agents capable of automating ELT pipeline construction, thus reducing manual effort.

**Text-to-SQL benchmarks and methods.** Researchers have studied the text-to-SQL task, which aims to convert natural language queries into SQL queries, for decades. Early text-to-SQL datasets primarily target single database scenarios [13, 25, 65]. More recent datasets, including WikiSQL [73] and Spider [71], extend this scope by introducing cross-domain scenarios requiring models to generalize to unseen databases. The BIRD benchmark is further introduced to evaluate text-to-SQL methods within large-scale database contexts, focusing on both query accuracy and execution efficiency [33]. Initially, text-to-SQL methods primarily leverage graph neural networks (GNNs) [4] and long short-term memory (LSTM) networks [64]. Recent research has increasingly adopted fine-tuning techniques [18, 32] and prompting approaches [11, 17, 43] to further enhance SQL generation accuracy with the advent of LLMs. ELT-Bench tasks agents with generating complex SQL transformation queries to construct data models based on provided column names and descriptions. These queries typically involve intricate structures, including nested subqueries and multi-table joins.

**AI agent benchmarks.** To support the development of AI agents for solving complex real-world tasks, researchers have introduced diverse benchmarks across several domains, including software engineering [26], machine learning [23], and web-based interactions [12, 74]. In the data domain, existing benchmarks primarily focus on data science code generation [24, 30] and data analysis [22]. Additionally, Spider 2-V [5] evaluates agents’ proficiency in using data tools, while Spider 2.0 [31] assesses agent performance on enterprise-focused text-to-SQL workflows. In contrast, ELT-Bench is the first benchmark designed to assess AI agents’ capabilities in developing real-world, end-to-end ELT pipelines.

**AI agents.** LLM-based Agents have emerged as a promising approach for addressing real-world challenges across various fields, including software engineering [59, 66, 72], web browsing [29, 40] and data science and engineering [21, 22, 31]. These agents typically consist of four crucial modules: reasoning [28, 60, 68], tool usage [44, 47], planning [50, 70], and memorization [75]. In this work, we evaluate two code generation agent frameworks (SWE-Agent [66] and Spider-Agent [31]) on ELT-Bench to assess their performance in constructing ELT pipelines.

## 8 CONCLUSION

We introduce ELT-Bench, a comprehensive end-to-end benchmark specifically designed for real-world ELT pipeline tasks in the data engineering domain. ELT-Bench aims to replicate realistic scenarios by providing environments for diverse data sources and integrating widely adopted data tools. The benchmark presents a substantial challenge, as the top-performing agent, Spider-Agent Claude-3.7-Sonnet, correctly generates data models in only 3.9% of cases. This performance gap highlights significant opportunities for future research to develop more powerful and intelligent AI agents capable of handling complex ELT workflows.

## REFERENCES

- [1] Airbyte. 2025. Airbyte. <https://airbyte.com/>
- [2] Anthropic. [n.d.]. Anthropic API Pricing. <https://www.anthropic.com/pricing#anthropic-api>
- [3] Anthropic. 2024. The Claude 3 Model Family: Opus, Sonnet, Haiku. <https://api.semanticscholar.org/CorpusID:268232499>
- [4] Ruisheng Cao, Lu Chen, Zhi Chen, Yanbin Zhao, Su Zhu, and Kai Yu. 2021. LGEQL: Line Graph Enhanced Text-to-SQL Model with Mixed Local and Non-Local Relations. arXiv:2106.01093 [cs.CL] <https://arxiv.org/abs/2106.01093>
- [5] Ruisheng Cao, Fangyu Lei, Haoyuan Wu, Jixuan Chen, Yeqiao Fu, Hongcheng Gao, Xinzhuang Xiong, Hanchong Zhang, Yuchen Mao, Wenjing Hu, Tianbao Xie, Hongshen Xu, Danyang Zhang, Sida Wang, Ruoxi Sun, Pengcheng Yin, Caiming Xiong, Ansong Ni, Qian Liu, Victor Zhong, Lu Chen, Kai Yu, and Tao Yu. 2024. Spider2-V: How Far Are Multimodal Agents From Automating Data Science and Engineering Workflows? arXiv:2407.10956 [cs.AI] <https://arxiv.org/abs/2407.10956>
- [6] Malu Castellanos, Alkis Simitsis, Kevin Wilkinson, and Umeshwar Dayal. 2009. Automating the loading of business process data warehouses. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology* (Saint Petersburg, Russia) (EDBT '09). Association for Computing Machinery, New York, NY, USA, 612–623. <https://doi.org/10.1145/1516360.1516431>
- [7] Abhirup Chatterjee and Arie Segev. 1991. Data manipulation in heterogeneous databases. *SIGMOD Rec.* 20, 4 (Dec. 1991), 64–68. <https://doi.org/10.1145/141356.141385>
- [8] Benoit Dageville, Thierry Cruanes, Marcin Zukowski, Vadim Antonov, Artin Avanes, Jon Bock, Jonathan Claybaugh, Daniel Engovatov, Martin Hentschel, Jiansheng Huang, Allison W. Lee, Ashish Motivala, Abdul Q. Munir, Steven Pelley, Peter Povinec, Greg Rahn, Spyridon Triantafyllis, and Philipp Unterbrunner. 2016. The Snowflake Elastic Data Warehouse. In *Proceedings of the 2016 International Conference on Management of Data* (San Francisco, California, USA) (SIGMOD '16). Association for Computing Machinery, New York, NY, USA, 215–226. <https://doi.org/10.1145/2882903.2903741>
- [9] dbt Labs. 2025. dbt. <https://www.getdbt.com/>
- [10] DeepSeek-AI. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. arXiv:2501.12948 [cs.CL] <https://arxiv.org/abs/2501.12948>
- [11] Xuemei Dong, Chao Zhang, Yuhang Ge, Yuren Mao, Yunjun Gao, lu Chen, Jinshu Lin, and Dongfang Lou. 2023. C3: Zero-shot Text-to-SQL with ChatGPT. arXiv:2307.07306 [cs.CL] <https://arxiv.org/abs/2307.07306>
- [12] Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H. Laradji, Manuel Del Verme, Tom Marty, Léo Boisivert, Megh Thakkar, Quentin Cappart, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. 2024. WorkArena: How Capable Are Web Agents at Solving Common Knowledge Work Tasks? arXiv:2403.07718 [cs.LG] <https://arxiv.org/abs/2403.07718>
- [13] Catherine Finegan-Dollak, Jonathan K. Kummerfeld, Li Zhang, Karthik Ramanathan, Sesh Sadasivam, Rui Zhang, and Dragomir Radev. 2018. Improving Text-to-SQL Evaluation Methodology. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 351–360. <https://doi.org/10.18653/v1/p18-1033>
- [14] Fireworks. [n.d.]. Fireworks Model Library. <https://fireworks.ai/models>
- [15] Fivetran. [n.d.]. Fivetran.
- [16] Harald Foidl, Valentina Golendukhina, Rudolf Ramler, and Michael Felderer. 2024. Data pipeline quality: Influencing factors, root causes of data-related issues, and processing problem areas for developers. *Journal of Systems and Software* 207 (2024), 111855. <https://doi.org/10.1016/j.jss.2023.111855>
- [17] Dawei Gao, Haibin Wang, Yaliang Li, Xiuyu Sun, Yichen Qian, Bolin Ding, and Jingren Zhou. 2024. Text-to-SQL Empowered by Large Language Models: A Benchmark Evaluation. *Proc. VLDB Endow.* 17, 5 (Jan. 2024), 1132–1145. <https://doi.org/10.14778/3641204.3641221>
- [18] Yingqi Gao, Yifu Liu, Xiaoxia Li, Xiaorong Shi, Yin Zhu, Yiming Wang, Shiqi Li, Wei Li, Yuntao Hong, Zhiling Luo, Jinyang Gao, Liyu Mou, and Yu Li. 2025. A Preview of XiYan-SQL: A Multi-Generator Ensemble Framework for Text-to-SQL. arXiv:2411.08599 [cs.AI] <https://arxiv.org/abs/2411.08599>
- [19] Aaron Grattafiori and et al. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] <https://arxiv.org/abs/2407.21783>
- [20] Mauricio A. Hernández and Salvatore J. Stolfo. 1998. Real-world Data is Dirty: Data Cleansing and The Merge/Purge Problem. *Data Min. Knowl. Discov.* 2, 1 (Jan. 1998), 9–37. <https://doi.org/10.1023/A:1009761603038>
- [21] Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Chenxing Wei, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Li Zhang, Lingyao Zhang, Min Yang, Mingchen Zhuge, Taicheng Guo, Tuo Zhou, Wei Tao, Xiangru Tang, Xiangtao Lu, Xiaowu Zheng, Xinbing Liang, Yaying Fei, Yuheng Cheng, Zhibin Gou, Zongze Xu, and Chenglin Wu. 2024. Data Interpreter: An LLM Agent For Data Science. arXiv:2402.18679 [cs.AI] <https://arxiv.org/abs/2402.18679>
- [22] Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing Xu, Ming Zhu, Yao Cheng, Jianbo Yuan, Jiwei Li, Kun Kuang, Yang Yang, Hongxia Yang, and Fei Wu. 2024. InfiAgent-DABench: Evaluating Agents on Data Analysis Tasks. In *Proceedings of the 41st International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (Eds.), Vol. 235. PMLR, 19544–19572. <https://proceedings.mlr.press/v235/hu24s.html>
- [23] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2024. MLAgent-Bench: Evaluating Language Agents on Machine Learning Experimentation. arXiv:2310.03302 [cs.LG] <https://arxiv.org/abs/2310.03302>
- [24] Yiming Huang, Jianwen Luo, Yan Yu, Yitong Zhang, Fangyu Lei, Yifan Wei, Shizhu He, Lifu Huang, Xiao Liu, Jun Zhao, and Kang Liu. 2024. DA-Code: Agent Data Science Code Generation Benchmark for Large Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 13487–13521. <https://doi.org/10.18653/v1/2024.emnlp-main.748>
- [25] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, Jayant Krishnamurthy, and Luke Zettlemoyer. 2017. Learning a Neural Semantic Parser from User Feedback. arXiv:1704.08760 [cs.CL] <https://arxiv.org/abs/1704.08760>
- [26] Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. SWE-bench: Can Language Models Resolve Real-World GitHub Issues? arXiv:2310.06770 [cs.CL] <https://arxiv.org/abs/2310.06770>
- [27] Ralph Kimball and Joe Caserta. 2004. *The Data Warehouse ETL Toolkit: Practical Techniques for Extracting, Cleaning, Conforming and Delivering Data*. John Wiley & Sons, Inc., Hoboken, NJ, USA.
- [28] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large Language Models are Zero-Shot Reasoners. arXiv:2205.11916 [cs.CL] <https://arxiv.org/abs/2205.11916>
- [29] Hanyu Lai, Xiao Liu, lat Long long, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang, Xiaohan Zhang, Yuxiao Dong, and Jie Tang. 2024. AutoWebGLM: A Large Language Model-based Web Navigating Agent. arXiv:2404.03648 [cs.CL] <https://arxiv.org/abs/2404.03648>
- [30] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen tau Yih, Daniel Fried, Sida Wang, and Tao Yu. 2022. DS-1000: A Natural and Reliable Benchmark for Data Science Code Generation. arXiv:2211.11501 [cs.SE] <https://arxiv.org/abs/2211.11501>
- [31] Fangyu Lei, Jixuan Chen, Yuxiao Ye, Ruisheng Cao, Dongchan Shin, Hongjin Su, Zhaoqing Suo, Hongcheng Gao, Wenjing Hu, Pengcheng Yin, Victor Zhong, Caiming Xiong, Ruoxi Sun, Qian Liu, Sida Wang, and Tao Yu. 2024. Spider 2.0: Evaluating Language Models on Real-World Enterprise Text-to-SQL Workflows. arXiv:2411.07763 [cs.CL] <https://arxiv.org/abs/2411.07763>
- [32] Haoyang Li, Jing Zhang, Hanbing Liu, Ju Fan, Xiaokang Zhang, Jun Zhu, Renjie Wei, Hongyan Pan, Cuiping Li, and Hong Chen. 2024. CodeS: Towards Building Open-source Language Models for Text-to-SQL. arXiv:2402.16347 [cs.CL] <https://arxiv.org/abs/2402.16347>
- [33] Jinyang Li, Binyuan Hui, Ge Qu, Jiayi Xue, Binhua Li, Bowen Li, Bailin Wang, Bowen Qin, Ruiying Geng, Nan Huo, Xuanhe Zhou, Chenhao Ma, Guoliang Li, Kevin C.C. Chang, Fei Huang, Reynold Cheng, and Yongbin Li. 2024. Can LLM already serve as a database interface? a big bench for large-scale database grounded text-to-SQLs. In *Proceedings of the 37th International Conference on Neural Information Processing Systems* (New Orleans, LA, USA) (NIPS '23). Curran Associates Inc., Red Hook, NY, USA, Article 1835, 28 pages.
- [34] LocalStack. 2025. LocalStack: A Fully Functional Local AWS Cloud Stack. <https://github.com/localstack/localstack>
- [35] Sergio Luján-Mora, Panos Vassiliadis, and Juan Trujillo. 2004. Data Mapping Diagrams for Data Warehouse Design with UML, Vol. 3288. 191–204. [https://doi.org/10.1007/978-3-540-30464-7\\_16](https://doi.org/10.1007/978-3-540-30464-7_16)
- [36] Anthony Mbata, Yaji Sripada, and Mingjun Zhong. 2024. A Survey of Pipeline Tools for Data Engineering. arXiv:2406.08335 [cs.LG] <https://arxiv.org/abs/2406.08335>
- [37] Kartick Chandra Mondal, Neepa Biswas, and Swati Saha. 2020. Role of Machine Learning in ETL Automation. In *Proceedings of the 21st International Conference on Distributed Computing and Networking* (Kolkata, India) (ICDCN '20). Association for Computing Machinery, New York, NY, USA, Article 57, 6 pages. <https://doi.org/10.1145/3369740.3372778>
- [38] OpenAI. [n.d.]. OpenAI API Pricing. <https://openai.com/api/pricing>
- [39] OpenAI. 2024. GPT-4o System Card. arXiv:2410.21276 [cs.CL] <https://arxiv.org/abs/2410.21276>
- [40] Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. 2024. Autonomous Evaluation and Refinement of Digital Agents. arXiv:2404.06474 [cs.AI] <https://arxiv.org/abs/2404.06474>
- [41] Meikel Poess, Tilmann Rabl, Hans-Arno Jacobsen, and Brian Caulfield. 2014. TPC-DI: the first industry benchmark for data integration. *Proc. VLDB Endow.* 7, 13 (Aug. 2014), 1367–1378. <https://doi.org/10.14778/2733004.2733009>
- [42] Daniel Poppy. 2023. ETL vs ELT: What's the difference? <https://www.getdbt.com/blog/etl-vs-elt>



- [43] Mohammadreza Pourreza and Davood Rafiei. 2023. DIN-SQL: Decomposed In-Context Learning of Text-to-SQL with Self-Correction. arXiv:2304.11015 [cs.CL] <https://arxiv.org/abs/2304.11015>
- [44] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023. ToolLLM: Facilitating Large Language Models to Master 16000+ Real-world APIs. arXiv:2307.16789 [cs.AI] <https://arxiv.org/abs/2307.16789>
- [45] Qwen. 2025. Qwen2.5 Technical Report. arXiv:2412.15115 [cs.CL] <https://arxiv.org/abs/2412.15115>
- [46] Raza Rasool and Ali Afzal Malik. 2015. Effort estimation of ETL projects using Forward Stepwise Regression. In *2015 International Conference on Emerging Technologies (ICET)*. 1–6. <https://doi.org/10.1109/ICET.2015.7389209>
- [47] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language Models Can Teach Themselves to Use Tools. arXiv:2302.04761 [cs.CL] <https://arxiv.org/abs/2302.04761>
- [48] Dhamotharan Seenivasan. 2022. ETL vs ELT: Choosing the right approach for your data warehouse. *International Journal for Research Trends and Innovation* (2022), 110–122.
- [49] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. arXiv:2303.11366 [cs.AI] <https://arxiv.org/abs/2303.11366>
- [50] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. arXiv:2303.11366 [cs.AI] <https://arxiv.org/abs/2303.11366>
- [51] Alkis Simitsis, Spiros Skiadopoulos, and Panos Vassiliadis. 2023. The History, Present, and Future of ETL Technology (invited). In *Proceedings of the 25th International Workshop on Design, Optimization, Languages and Analytical Processing of Big Data (DOLAP) co-located with the 26th International Conference on Extending Database Technology and the 26th International Conference on Database Theory (EDBT/ICDT 2023)*, Ioannina, Greece, March 28, 2023 (CEUR Workshop Proceedings), Enrico Gallinucci and Lukasz Golab (Eds.), Vol. 3369. CEUR-WS.org, 3–12. <https://ceur-ws.org/Vol-3369/invited1.pdf>
- [52] Bharat Singhal and Alok Aggarwal. 2022. ETL, ELT and reverse ETL: a business case Study. In *2022 Second International Conference on Advanced Technologies in Intelligent Control, Environment, Computing & Communication Engineering (ICATIECE)*. IEEE, 1–4.
- [53] Dimitrios Skoutas and Alkis Simitsis. 2006. Designing ETL processes using semantic web technologies. In *Proceedings of the 9th ACM International Workshop on Data Warehousing and OLAP* (Arlington, Virginia, USA) (DOLAP '06). Association for Computing Machinery, New York, NY, USA, 67–74. <https://doi.org/10.1145/1183512.1183526>
- [54] Terraform. 2025. airbyte Provider. <https://registry.terraform.io/providers/airbytehq/airbyte/latest/docs>
- [55] Juan Trujillo and Sergio Luján-Mora. 2003. A UML Based Approach for Modeling ETL Processes in Data Warehouses, Vol. 2813. 307–320. [https://doi.org/10.1007/978-3-540-39648-2\\_25](https://doi.org/10.1007/978-3-540-39648-2_25)
- [56] Alexander van Renen and Viktor Leis. 2023. Cloud Analytics Benchmark. *Proc. VLDB Endow.* 16, 6 (Feb. 2023), 1413–1425. <https://doi.org/10.14778/3583140.3583156>
- [57] Panos Vassiliadis, Alkis Simitsis, and Spiros Skiadopoulos. 2002. Conceptual modeling for ETL processes. In *Proceedings of the 5th ACM International Workshop on Data Warehousing and OLAP* (McLean, Virginia, USA) (DOLAP '02). Association for Computing Machinery, New York, NY, USA, 14–21. <https://doi.org/10.1145/583890.583893>
- [58] Florian Waas, Robert Wrembel, Tobias Freudenreich, Maik Thiele, Christian Koncilia, and Pedro Furtado. 2013. On-Demand ELT Architecture for Right-Time BI: Extending the Vision. *International Journal of Data Warehousing and Mining* 9 (04 2013), 21–38. <https://doi.org/10.4018/jdwmm.2013040102>
- [59] Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable Code Actions Elicit Better LLM Agents. arXiv:2402.01030 [cs.CL] <https://arxiv.org/abs/2402.01030>
- [60] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903 [cs.CL] <https://arxiv.org/abs/2201.11903>
- [61] Lilian Weng. 2023. LLM-powered Autonomous Agents. *lilianweng.github.io* (Jun 2023). <https://lilianweng.github.io/posts/2023-06-23-agent/>
- [62] Jennifer Widom. 1995. Research problems in data warehousing. In *Proceedings of the Fourth International Conference on Information and Knowledge Management* (Baltimore, Maryland, USA) (CIKM '95). Association for Computing Machinery, New York, NY, USA, 25–30. <https://doi.org/10.1145/221270.221319>
- [63] Niklas Wretblad, Fredrik Gordh Riseby, Rahul Biswas, Amin Ahmadi, and Oskar Holmström. 2024. Understanding the Effects of Noise in Text-to-SQL: An Examination of the BIRD-Bench Benchmark. arXiv:2402.12243 [cs.CL] <https://arxiv.org/abs/2402.12243>
- [64] Xiaojun Xu, Chang Liu, and Dawn Song. 2017. SQLNet: Generating Structured Queries From Natural Language Without Reinforcement Learning. arXiv:1711.04436 [cs.CL] <https://arxiv.org/abs/1711.04436>
- [65] Navid Yaghmazadeh, Yuepeng Wang, Isil Dillig, and Thomas Dillig. 2017. SQLizer: query synthesis from natural language. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 63 (Oct. 2017), 26 pages. <https://doi.org/10.1145/3133887>
- [66] John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. SWE-agent: Agent-Computer Interfaces Enable Automated Software Engineering. arXiv:2405.15793 [cs.SE] <https://arxiv.org/abs/2405.15793>
- [67] Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024.  $\tau$ -bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains. arXiv:2406.12045 [cs.AI] <https://arxiv.org/abs/2406.12045>
- [68] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafraan, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. arXiv:2305.10601 [cs.CL] <https://arxiv.org/abs/2305.10601>
- [69] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafraan, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. arXiv:2210.03629 [cs.CL] <https://arxiv.org/abs/2210.03629>
- [70] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafraan, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. arXiv:2210.03629 [cs.CL] <https://arxiv.org/abs/2210.03629>
- [71] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2019. Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task. arXiv:1809.08887 [cs.CL] <https://arxiv.org/abs/1809.08887>
- [72] Yuntong Zhang, Haifeng Ruan, Zhiyu Fan, and Abhik Roychoudhury. 2024. AutoCodeRover: Autonomous Program Improvement. arXiv:2404.05427 [cs.SE] <https://arxiv.org/abs/2404.05427>
- [73] Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2SQL: Generating Structured Queries from Natural Language using Reinforcement Learning. arXiv:1709.00103 [cs.CL] <https://arxiv.org/abs/1709.00103>
- [74] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2024. WebArena: A Realistic Web Environment for Building Autonomous Agents. arXiv:2307.13854 [cs.AI] <https://arxiv.org/abs/2307.13854>
- [75] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, Yu Qiao, Zhaoxiang Zhang, and Jifeng Dai. 2023. Ghost in the Minecraft: Generally Capable Agents for Open-World Environments via Large Language Models with Text-based Knowledge and Memory. arXiv:2305.17144 [cs.AI] <https://arxiv.org/abs/2305.17144>