



# CatDB: LLM-based Generation of Data-centric ML Pipelines

Saeed Fathollahzadeh<sup>1</sup>

Essam Mansour<sup>1</sup>

Matthias Boehm<sup>2</sup>



## 1. Motivation

- Data-centric ML pipelines are crucial → **labor-intensive**.
  - AutoML systems → **struggle with large datasets**.
  - LLMs demonstrate strong capabilities in coding → **struggle on unseen data**.
  - LLM-based pipeline generation → **lacks tailored dataset context**.
- ➔ **Need for efficient, scalable pipeline generation for diverse data.**

### Raw Dataset

Table #1	Table #2
c <sub>1</sub> c <sub>2</sub> ...	c <sub>1</sub> c <sub>2</sub> ...
1 a ...	A 0 ...
2 - ...	B 1 ...
...	...

Table #N
c <sub>1</sub> c <sub>2</sub> ... target (y)
1 A ... Yes
2 0 ... No
...

### Pipeline Generation via CatDB

```
from catdb import config, generate_pipeline
from dataprofilng import build_catalog

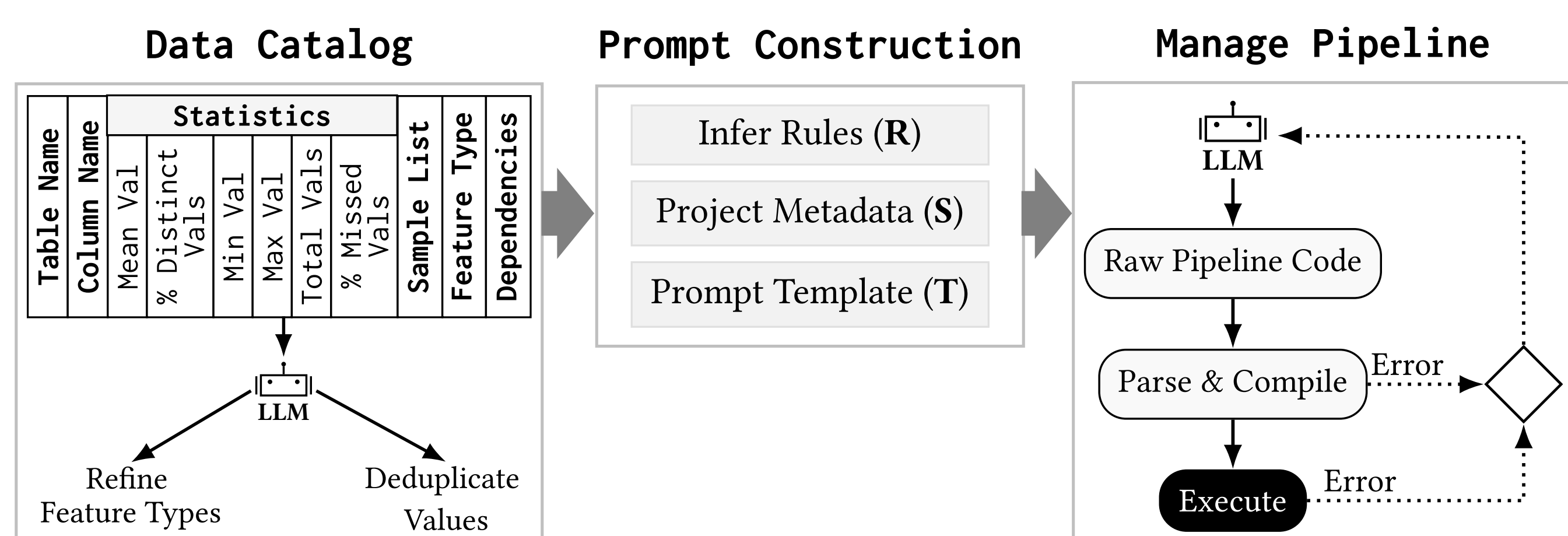
cfg = config(llm_model='gpt-4o')
cat = build_catalog(path='raw_dataset')
p = generate_pipeline(catalog=cat, config=cfg)
```

### ML Pipeline

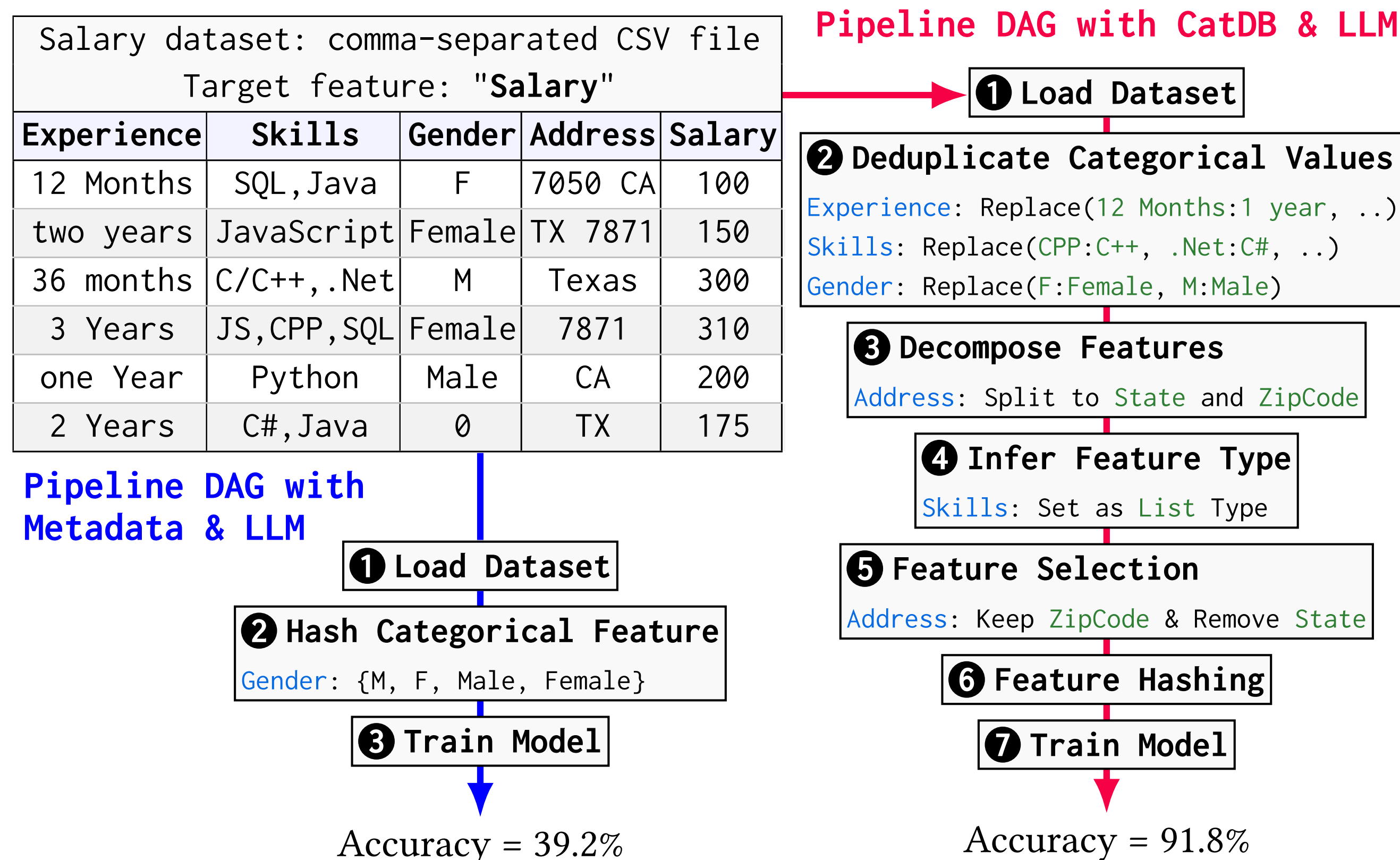
```
1: import pandas as pd
2: import SimpleImputer
3: import OneHotEncoder
4: import ColumnTransformer
5: import Pipeline
6: import RandomForestRegressor
7: import r2_score

8: train = pd.read_csv("train.csv")
9: test = pd.read_csv("test.csv")
10: cat = ["Experience", "Gender"]
11: cat = Pipeline(steps=[
12:     ("imputer", SimpleImputer(...)),
13:     ("encoder", OneHotEncoder(...))
14: ])
15: p = Pipeline(steps=[
16:     ("preprocessor", ColumnTransformer(
17:         transformers=[("cat", cat)]
18:     )),
19:     ("model", RandomForestRegressor(...))
20: ])
21: p.fit(X_train, y_train)
22: y_test_pred = p.predict(X_test)
```

## 2. CatDB Overview



## 3. Data-centric ML pipelines w/ [Metadata-only & LLM] vs. [CatDB & LLM]



## 4. An Example of a Prompt (Instructions and Meta- data)

### Task: Generate a data science pipeline in Python 3.10.

### Input: A dataset in CSV format, a schema that describes the columns and data types of the dataset, and a data profiling info that summarizes the statistics and quality of the dataset.

### Output: A pipeline code that performs the following steps:

- #1: Import the necessary libraries and modules.
- #2: Load the training and test datasets. For the training data, utilize the variable "train\_data.csv", and for the test data, employ the variable "test\_data.csv".
- #3: The user will provide the Schema, and Data Profiling Info of the dataset with columns appropriately named as attributes.
- #4: Perform missing value imputation for features 'Address' and 'Zip'.
- #5: Perform feature extraction (dataset contains categorical values).
- #6: Perform feature selection (redundant columns e.g., 'Address' and 'Zip').
- #7: Select an appropriate ML algorithm.
- #8: Assign a default value to a particular hyperparameter.
- #9: Evaluate the model.

### Dataset Description: A The dataset was obtained from multiple sources, including surveys, job posting sites, and other publicly available sources. A total of 6704 data points were collected. The dataset included five variables: age, experience, job role, and education level and salary.

### Schema, and Data Profiling Info:

- # Experience (string), distinct [ 60% ], missing [ 0% ], categorical-vals [ 1 year, 2 years, 3 years ]
- # Gender (string), distinct [ 40% ], missing [ 0% ], categorical-vals [ Male, Female ]
- # Address (string), distinct [ 40% ], missing [ 20% ], categorical-vals [ CA, TX ]
- # Zip (string), distinct [ 40% ], missing [ 40% ], categorical-vals [ 7050, 7871 ]
- # Salary (int, **target feature**), min-max vals [ 100, 310 ], total-vals [ 5 ]

### Categorical Features: Experience, Gender, Address, Zip

## 5. CatDB Data-Centric ML Pipeline Generation APIs

```
1 # Install Data Profiling
!pip install -U git+https://github.com/CoDS-GCS/CatDB.git@dataprofilng --quiet

# Install CatDB
!pip install -U git+https://github.com/CoDS-GCS/CatDB.git --quiet

2 # Import CatDB APIs
from catdb import config, prepare_dataset, create_report, generate_pipeline

# Import Data Profiling API
from dataprofilng import build_catalog

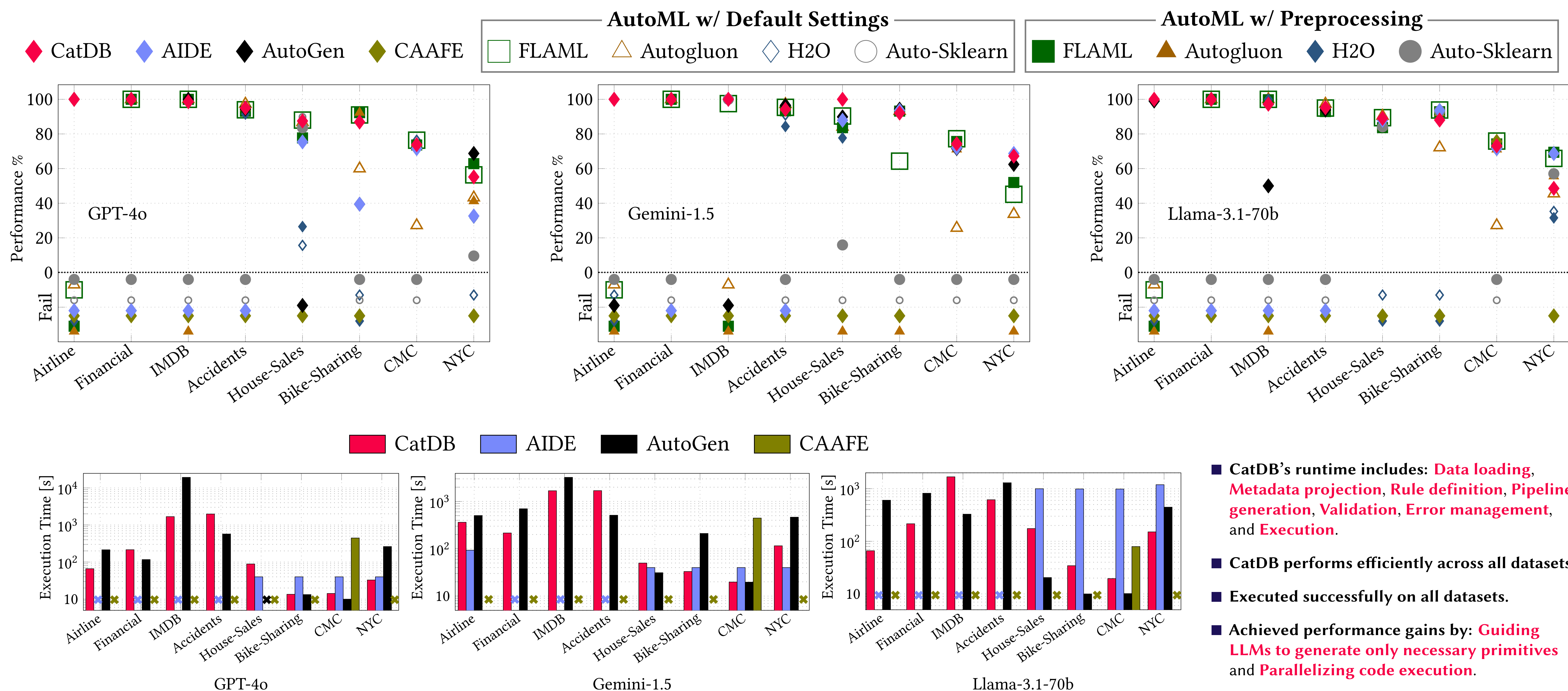
3 cfg = config(model="gpt-4o", API_key="KEY", iteration=7)

4 data = prepare_dataset(path="adult.csv", task_type="binary",
    target_attribute="income")

5 catalog = build_catalog(data=data)
create_report(catalog)

6 pipeline = generate_pipeline(catalog, cfg)
create_report(pipeline)
```

## 6. Performance Comparison of 8 Datasets



- CatDB's runtime includes: **Data loading, Metadata projection, Rule definition, Pipeline generation, Validation, Error management, and Execution.**
- CatDB performs efficiently across all datasets.
- Executed successfully on all datasets.
- Achieved performance gains by: **Guiding LLMs to generate only necessary primitives and Parallelizing code execution.**