

CatDB: Data-catalog-guided, LLM-based Generation of Data-centric ML Pipelines

Saeed Fathollahzadeh¹ Essam Mansour¹ Matthias Boehm^{2,3}

1. Motivation

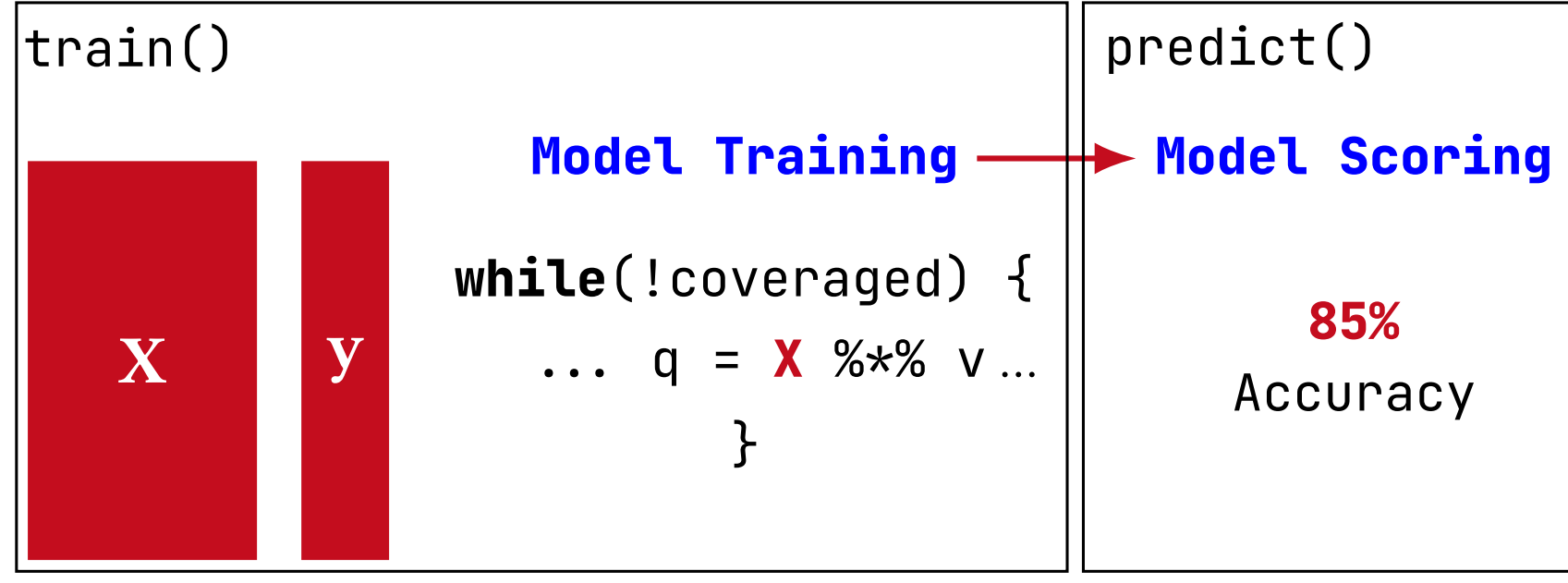
Data Preparation

Data Integration & Data Cleaning

Data Programming & Augmentation

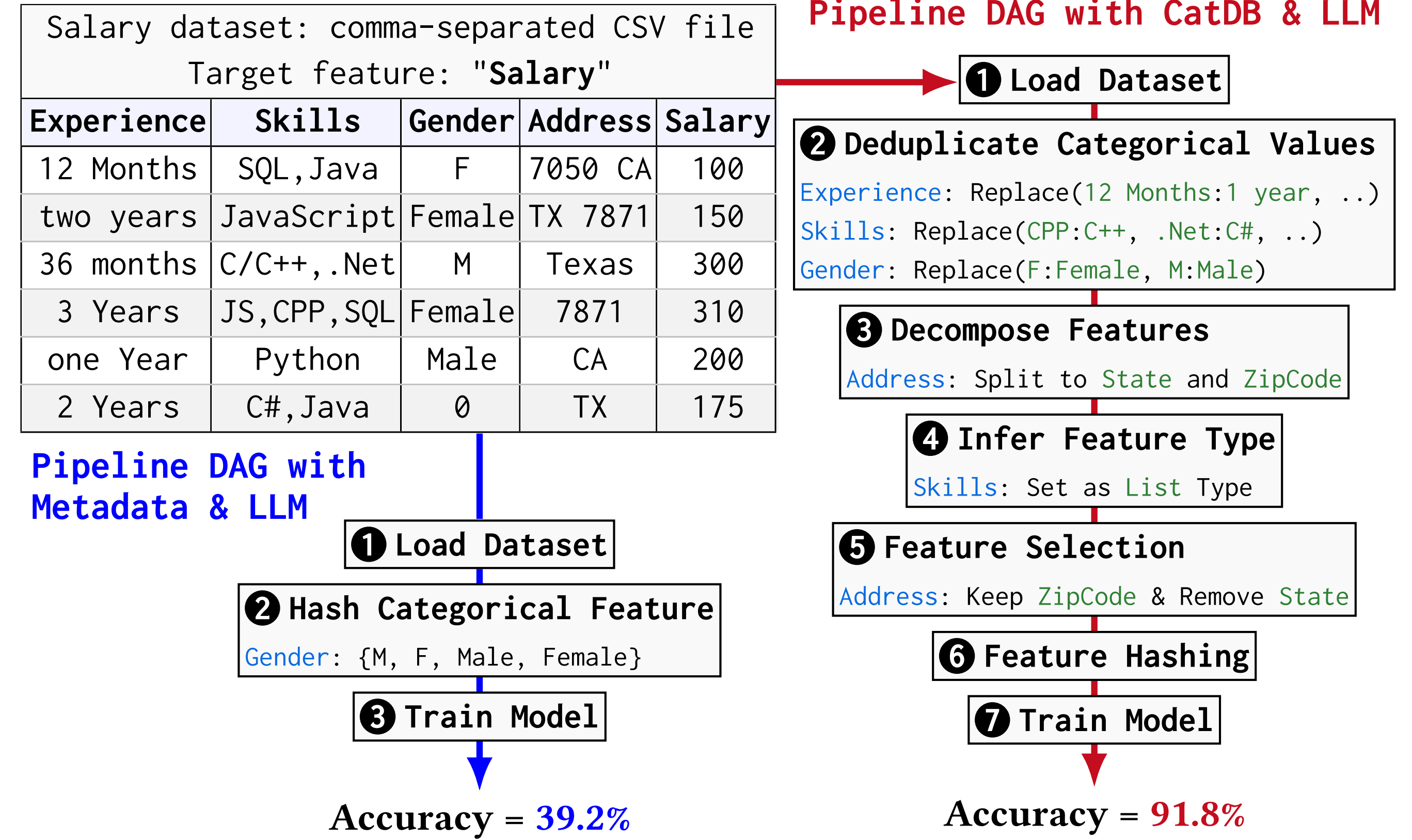
Model and Feature Selection

Hyperparameter Tuning + CV

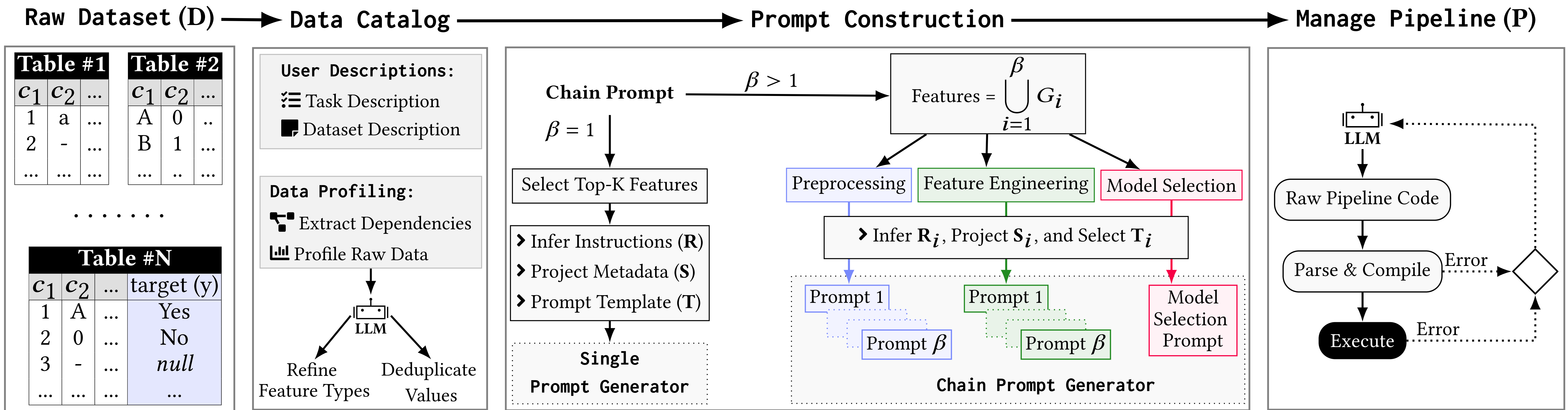


- ❑ Data-centric ML pipelines are crucial → **labor-intensive**.
- ❑ AutoML systems → **struggle with large datasets**.
- ❑ LLMs demonstrate strong capabilities in coding → **struggle on unseen data**.
- ❑ LLM-based pipeline generation → **lacks tailored dataset context**.

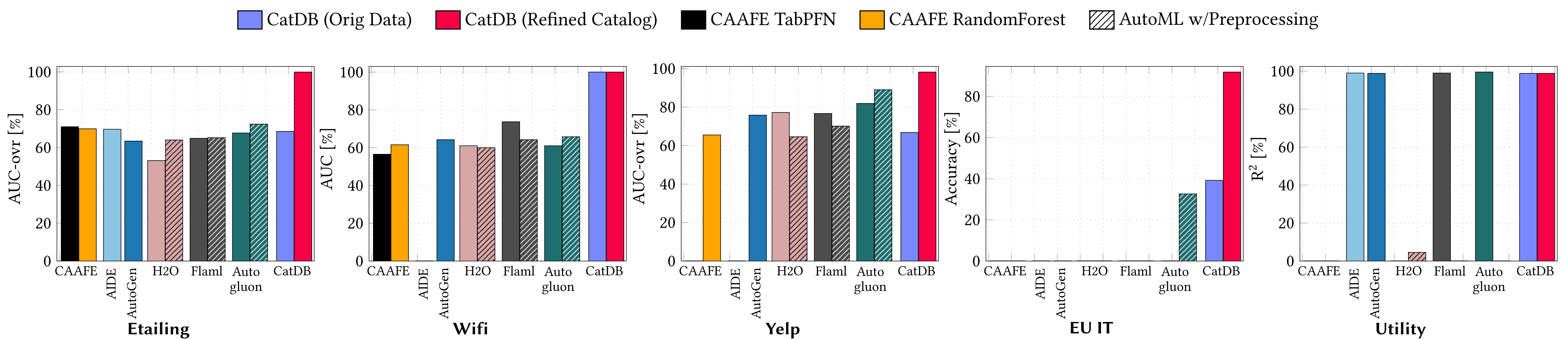
2. Data-centric ML pipelines w/ [Metadata-only vs. CatDB] & LLM



3. Data-centric ML pipeline generation in CatDB



4. Performance Comparison (LLM = Gemini-1.5)



Categorical Features:

- Fixes formatting, transforms complexity.

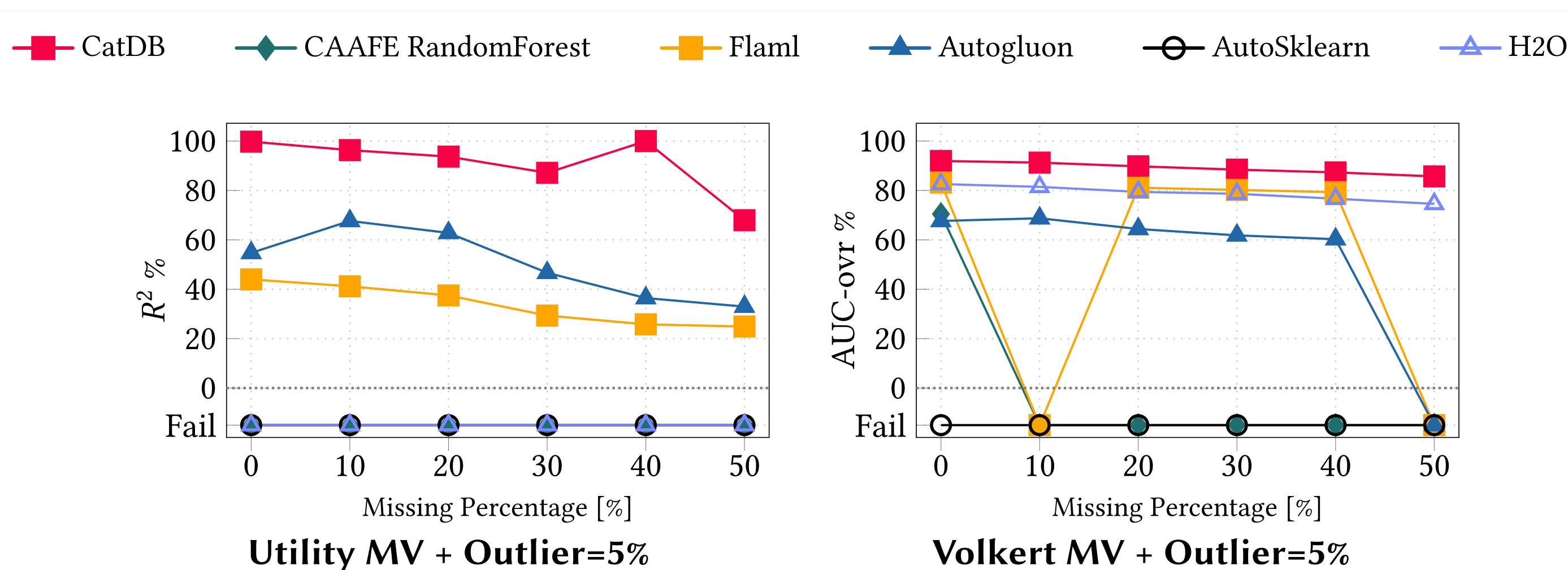
Deduplication:

- Removes duplicates, balances labels.

Feature Refinement:

- Drops constant/misread features, preserves distribution.

5. Outlier and Missing Value Injection (Gemini-1.5)



- ❑ **Robustness** → CatDB maintains high performance even with increasing missing values and 5% outliers.

6. Conclusions

- ❑ **Data Catalog Integration** → Use metadata & rules for tailored pipelines.
- ❑ **Catalog Refinements** → Enhance catalogs to guide ML pipeline creation.
- ❑ **Prompt Chaining** → Sequence prompts to optimize generation.
- ❑ **Error Handling** → Validate, fix with knowledge base for reliable pipelines.



Paper



Reproducibility