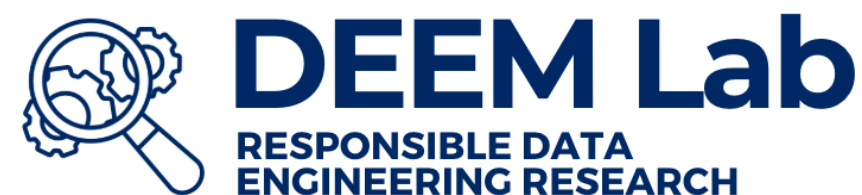


# SemPipes: Optimizable Semantic Data Operators for Tabular Machine Learning Pipelines



NHR4CES Community Workshop 2026, May 12

Olga Ovcharenko  Matthias Boehm  Sebastian Schelter 

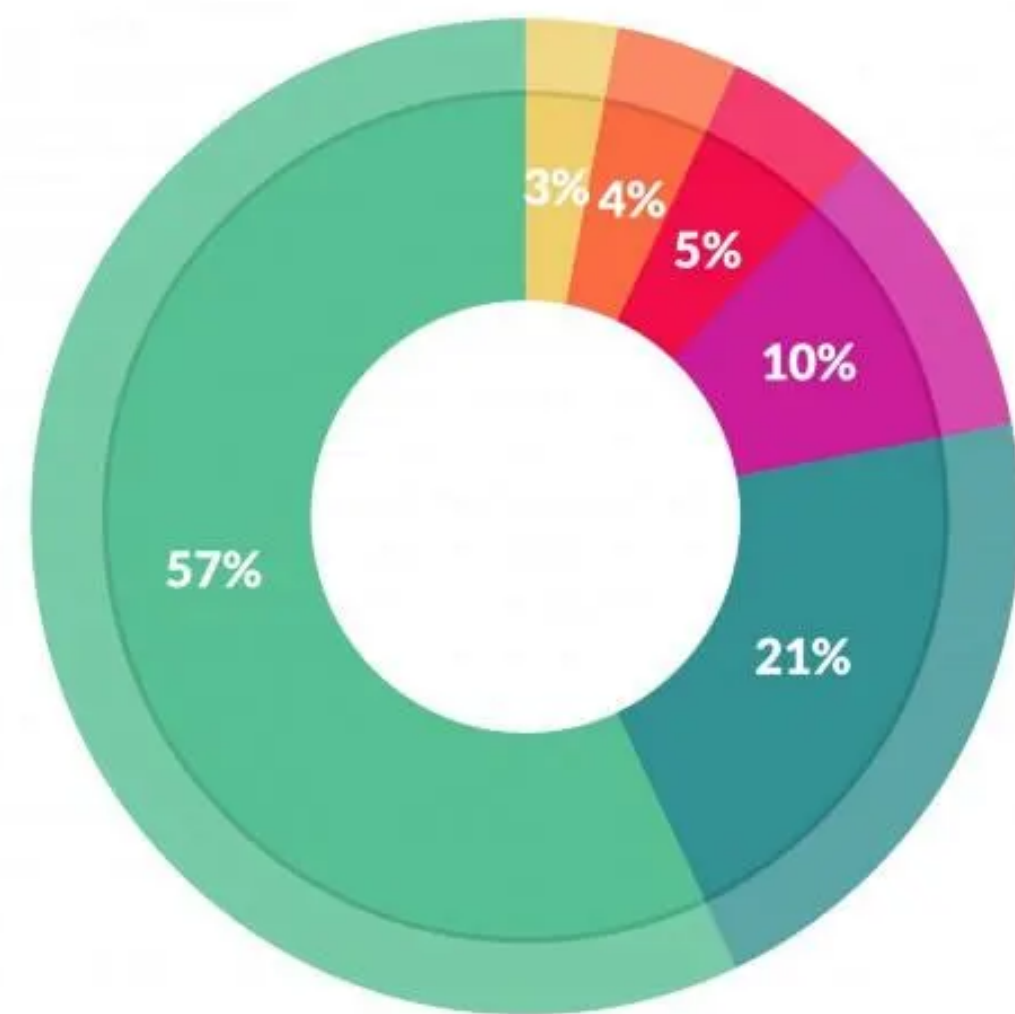


# About Me



- Second year PhD student at  **DEEM Lab**, supervised by Prof. Schelter  
RESPONSIBLE DATA  
ENGINEERING RESEARCH
- Studied before at **ETH zürich** and  **TU**  
Graz
- Working on semantic operators in the context of ML and data management
- Prior experience on applied ML for genomics and development of ML systems internals

# Data Preparation: The Least Enjoyable Part of Data Science



What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

<https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

## "Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes AI

Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, Lora Aroyo  
[nithyasamba,kapania,hhighfill,dakrong,pkp,loraa]@google.com  
Google Research  
Mountain View, CA

### ABSTRACT

AI models are increasingly applied in high-stakes domains like health and conservation. Data quality carries an elevated significance in high-stakes AI due to its heightened downstream impact, impacting predictions like cancer detection, wildlife poaching, and loan allocations. Paradoxically, data is the most under-valued and de-glamorised aspect of AI. In this paper, we report on data practices in high-stakes AI, from interviews with 53 AI practitioners in India, East and West African countries, and USA. We define, identify, and present empirical evidence on *Data Cascades*—compounding events causing negative, downstream effects from data issues—triggered by conventional AI/ML practices that undervalue data quality. Data cascades are pervasive (92% prevalence), invisible, delayed, but often avoidable. We discuss HCI opportunities in designing and incentivizing data excellence as a first-class citizen of AI, resulting in safer and more robust systems for all.

### CCS CONCEPTS

• Human-centered computing → Empirical studies in HCI.

lionized work of building novel models and algorithms [46, 125]. Intuitively, AI developers understand that data quality matters, often spending inordinate amounts of time on data tasks [60]. In practice, most organisations fail to create or meet any data quality standards [87], from under-valuing data work vis-a-vis model development.

Under-valuing of data work is common to all of AI development [125]<sup>1</sup>. We pay particular attention to undervaluing of data in *high-stakes domains*<sup>2</sup> that have safety impacts on living beings, due to a few reasons. One, developers are increasingly deploying AI models in complex, humanitarian domains, e.g., in maternal health, road safety, and climate change. Two, poor data quality in high-stakes domains can have outsized effects on vulnerable communities and contexts. As Hiatt *et al.* argue, high-stakes efforts are distinct from serving customers; these projects work with and for populations at risk of a litany of horrors [47]. As an example, poor data practices reduced accuracy in IBM's cancer treatment AI [115] and led to Google Flu Trends missing the flu peak by 140% [63, 73]). Three, high-stakes AI systems are typically deployed in low-resource contexts with a pronounced lack of readily available, high-quality datasets. Applications span into communities that

"Everyone wants to do the model work, not the data work":  
Data Cascades in High-Stakes AI, CHI'21

# ML Engineering Agents Will Write All Pipelines Instead of Us!

- LLM-based agents for the end-to-end machine learning engineering using code generation
- Promises
  - Automatic implementation of ML pipelines via code generation
  - State-of-the-art results on multiple machine learning engineering benchmarks
  - End-to-end ML automation across diverse data modalities with minimal human intervention
  - Comprehensive, LLM-based system for generating effective, error-free, and efficient data-centric ML pipelines

**MLE-STAR: Machine Learning Engineering Agent via Search and Targeted Refinement**

Jaehyun Nam<sup>1, 2\*</sup>, Jinsung Yoon<sup>1</sup>, Jiefeng Chen<sup>1</sup>  
Jinwoo Shin<sup>2</sup>, Sercan Ö. Arik<sup>1</sup>, Tomas Pfister<sup>1</sup>  
<sup>1</sup>Google Cloud, <sup>2</sup>KAIST  
jaehyun.nam@kaist.ac.kr, jinsungyoon@google.com

**Abstract**

Agents based on large language models (LLMs) for machine learning engineering (MLE) can automatically implement ML models via code generation. However, existing approaches to build such agents often rely heavily on inherent LLM knowledge and employ coarse exploration strategies that modify the entire code structure at once. This limits their ability to select effective task-specific models and perform deep exploration within specific components, such as experimenting extensively with feature engineering options. To overcome these, we propose *MLE-STAR*, a novel approach to build MLE agents. *MLE-STAR* first leverages external knowledge by using a search engine to retrieve effective models from the web, forming an initial solution, then iteratively refines it by exploring various strategies targeting specific ML components. This exploration is guided by ablation studies analyzing the impact of individual code blocks. Furthermore, we introduce a novel ensembling method using an effective strategy suggested by *MLE-STAR*. Our experimental results show that *MLE-STAR* achieves medals in 64% of the Kaggle competitions on the MLE-bench, significantly outperforming the best alternative.<sup>1</sup>

MLE-STAR: Machine Learning Engineering Agent via Search and Targeted Refinement, NeurIPS'25

**AIDE: AI-Driven Exploration in the Space of Code**

Zhengyao Jiang<sup>4</sup>, Dominik Schmidt<sup>4</sup>, Dhruv Srikanth, Dixing Xu, Ian Kaplan  
Weco AI, Runway ML, Weco AI, Weco AI, Weco AI

Deniss Jacenko, Yuxiang Wu<sup>4</sup>  
Weco AI, Weco AI

**Abstract**

Machine learning, the foundation of modern artificial intelligence, has driven innovations that have fundamentally transformed the world. Yet, behind advancements lies a complex and often tedious process requiring labor and compute-intensive iteration and experimentation. Engineers and scientists developing machine learning models spend much of their time on trial-and-error tasks instead of conceptualizing innovative solutions or research hypotheses. To address this challenge, we introduce AI-Driven Exploration (AIDE), a machine learning engineering agent powered by large language models (LLMs). AIDE frames machine learning engineering as a code optimization problem, and formulates trial-and-error as a tree search in the space of potential solutions. By strategically reusing and refining promising solutions, AIDE effectively trades computational resources for enhanced performance, achieving state-of-the-art results on multiple machine learning engineering benchmarks, including our Kaggle evaluations, OpenAI's MLE-Bench and METR's RE-Bench. The implementation of AIDE is publicly available at <https://github.com/WecoAI/aideml>.

AIDE: AI-Driven Exploration in the Space of Code

**MLZero: A Multi-Agent System for End-to-end Machine Learning Automation**

Haoyang Fang, Boran Han, Nick Erickson, Xiyuan Zhang, Su Zhou, Anirudh Dagar, Jiani Zhang, Ali Caner Turkmen, Cuixiong Hu, Huzefa Rangwala, Ying Nian Wu, Bernie Wang, George Karypis  
Amazon Web Services

{haoyang, boranhan, neerick, xiyuanz, suzhou, anidagar}@amazon.com  
{zhajiani, atturkm, tonyhu, rhuzefa, wunyin, yuyawang, gkarypis}@amazon.com

Website: <https://project-mlzero.github.io/>  
GitHub: <https://github.com/autogluon/autogluon-assistant>

**Abstract**

Existing AutoML systems have advanced the automation of machine learning (ML); however, they still require substantial manual configuration and expert input, particularly when handling multimodal data. We introduce MLZero, a novel multi-agent framework powered by Large Language Models (LLMs) that enables end-to-end ML automation across diverse data modalities with minimal human intervention. A cognitive perception module is first employed, transforming raw multimodal inputs into perceptual context that effectively guides the subsequent workflow. To address key limitations of LLMs, such as hallucinated code generation and outdated API knowledge, we enhance the iterative code generation process with semantic and episodic memory. MLZero demonstrates superior performance on MLE-Bench Lite, outperforming all competitors in both success rate and solution quality, securing six gold medals. Additionally, when evaluated on our Multimodal AutoML Agent Benchmark, which includes 25 more challenging tasks spanning diverse data modalities, MLZero outperforms the competing methods by a large margin with a success rate of 0.92 (+263.6%) and an average rank of 2.28. Our approach maintains its robust effectiveness even with a compact 8B LLM, outperforming full-size systems from existing solutions.

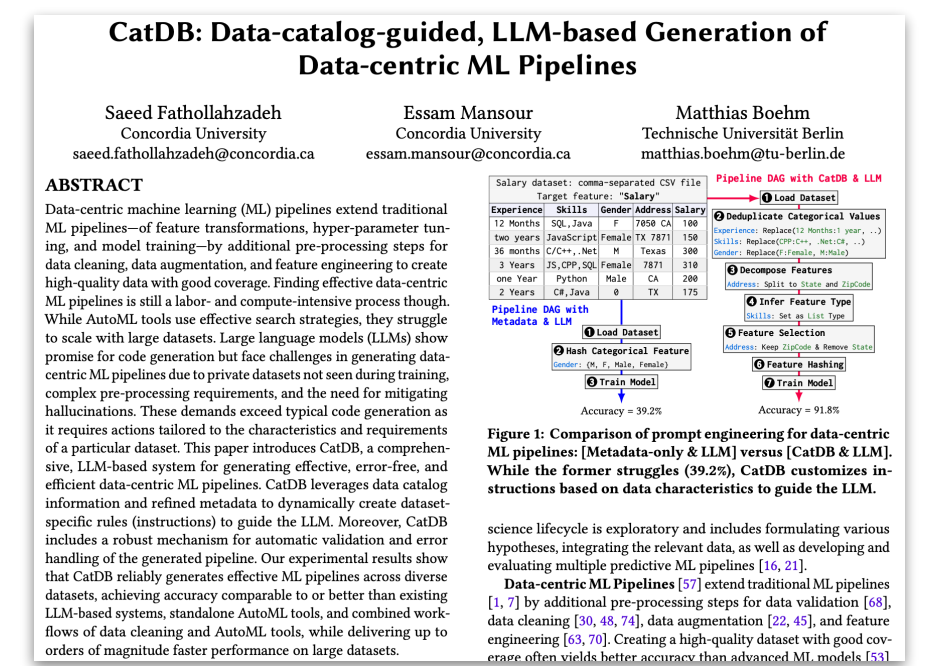
MLZero: A Multi-Agent System for End-to-end Machine Learning Automation, NeurIPS'25

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

Saeed Fathollahzadeh, Essam Mansour, Matthias Boehm  
Concordia University, Concordia University, Technische Universität Berlin  
saeed.fathollahzadeh@concordia.ca, essam.mansour@concordia.ca, matthias.boehm@tu-berlin.de

**ABSTRACT**

Data-centric machine learning (ML) pipelines extend traditional ML pipelines—of feature transformations, hyper-parameter tuning, and model training—by additional pre-processing steps for data cleaning, data augmentation, and feature engineering to create high-quality data with good coverage. Finding effective data-centric ML pipelines is still a labor- and compute-intensive process though. While AutoML tools use effective search strategies, they struggle to scale with large datasets. Large language models (LLMs) show promise for code generation but face challenges in generating data-centric ML pipelines due to private datasets not seen during training, complex pre-processing requirements, and the need for mitigating hallucinations. These demands exceed typical code generation as it requires actions tailored to the characteristics and requirements of a particular dataset. This paper introduces CatDB, a comprehensive, LLM-based system for generating effective, error-free, and efficient data-centric ML pipelines. CatDB leverages data catalog information and refined metadata to dynamically create dataset-specific rules (instructions) to guide the LLM. Moreover, CatDB includes a robust mechanism for automatic validation and error handling of the generated pipeline. Our experimental results show that CatDB reliably generates effective ML pipelines across diverse datasets, achieving accuracy comparable to or better than existing LLM-based systems, standalone AutoML tools, and combined workflows of data cleaning and AutoML tools, while delivering up to orders of magnitude faster performance on large datasets.



CatDB: Data-catalog-guided, LLM-based Generation of Data-centric ML Pipelines, VLDB'24

# Will ML Engineering Agents Write All Pipelines Instead of Us?

**AIDE: AI-Driven Exploration in the Space of Code**

Zhengyao Jiang<sup>4</sup> Weico AI, Dominik Schmidt<sup>4</sup> Runway ML, Dhruv Srikanth Weico AI, Dikang Xu Weico AI, Ian Kaplan Weico AI, Denis Jacenko Weico AI, Yuxiang Wu<sup>4</sup> Weico AI

**Abstract**

Machine learning, the foundation of modern artificial intelligence, has driven innovations that have fundamentally transformed the world. Yet, behind advancements lies a complex and often tedious process requiring labor and compute-intensive iteration and experimentation. Engineers and scientists developing machine learning models spend much of their time on trial-and-error tasks instead of conceptualizing innovative solutions or research hypotheses. To address this challenge, we introduce AI-Driven Exploration (AIDE), a machine learning engineering agent powered by large language models (LLMs). AIDE frames machine learning engineering as a code optimization problem, and formulates trial-and-error as a tree search in the space of potential solutions. By strategically reusing and refining promising solutions, AIDE effectively trades computational resources for enhanced performance, achieving state-of-the-art results on multiple machine learning engineering benchmarks, including our Kaggle evaluations, OpenAI's MLE-Bench and METR's RE-Bench. The implementation of AIDE is publicly available at <https://github.com/weicoAI/aide1>.

AIDE: AI-Driven Exploration in the Space of Code

**MLE-BENCH: EVALUATING MACHINE LEARNING AGENTS ON MACHINE LEARNING ENGINEERING**

Chan Jun Shern<sup>1</sup>, Neil Chowdhury<sup>1\*</sup>, Oliver Jaffe<sup>1</sup>, James Aung<sup>1</sup>, Dane Sherburn<sup>1</sup>, Evan Mays<sup>1</sup>, Giulio Starace<sup>1</sup>, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng<sup>1</sup>, Aleksander Mądry

OpenAI

**ABSTRACT**

We introduce MLE-bench, a benchmark for measuring how well AI agents perform at machine learning engineering. To this end, we curate 75 ML engineering-related competitions from Kaggle, creating a diverse set of challenging tasks that test real-world ML engineering skills such as training models, preparing datasets, and running experiments. We establish human baselines for each competition using Kaggle's publicly available leaderboards. We use open-source agent scaffolds to evaluate several frontier language models on our benchmark, finding that the best-performing setup — OpenAI's o1-preview with AIDE scaffolding — achieves at least the level of a Kaggle bronze medal in 16.9% of competitions. In addition to our main results, we investigate various forms of resource-scaling for AI agents and the impact of contamination from pre-training. We open-source our benchmark code ([github.com/openai/mle-bench/](https://github.com/openai/mle-bench/)) to facilitate future research in understanding the ML engineering capabilities of AI agents.

MLE-Bench: Evaluating Machine Learning Agents on Machine Learning Engineering, ICLR'25

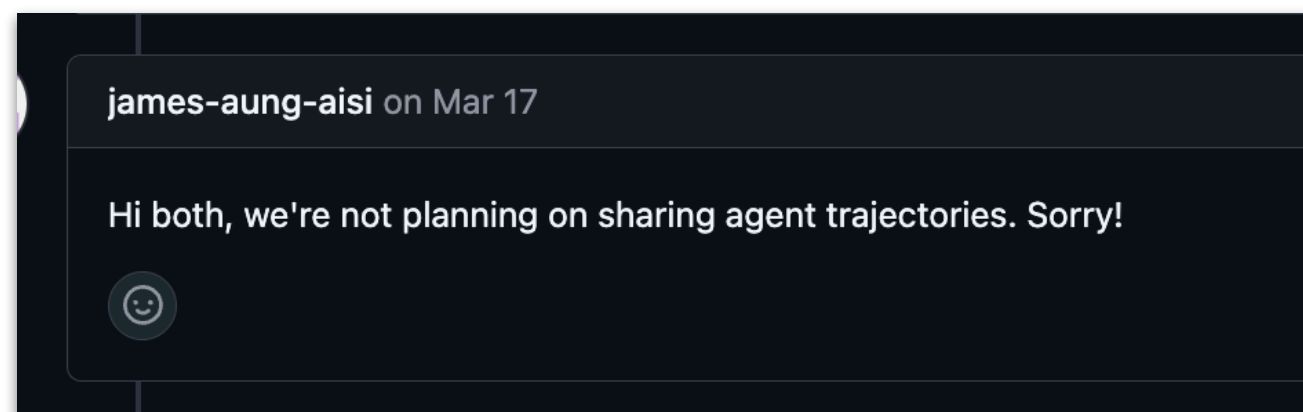
**CatDB: Data-catalog-guided, LLM-based Generation of Data-centric ML Pipelines**

Saeed Fathollahzadeh Concordia University, Essam Mansour Concordia University, Matthias Boehm Technische Universität Berlin

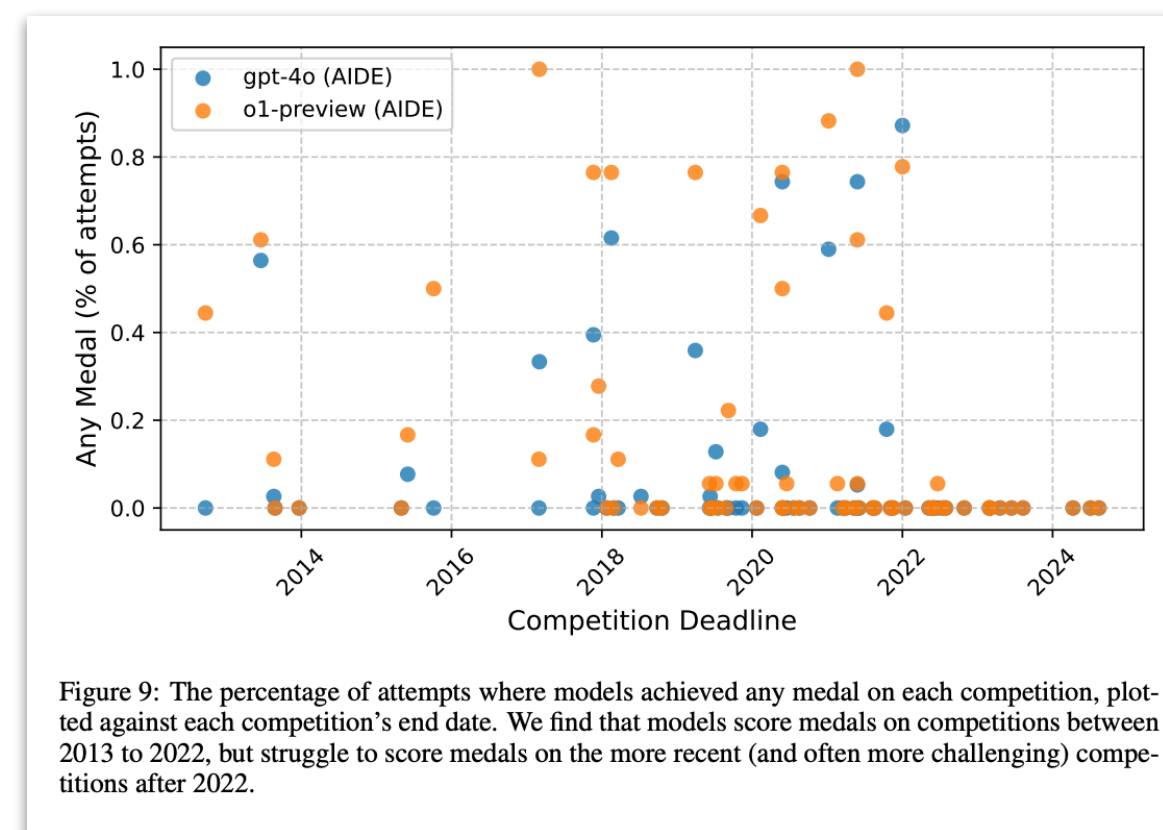
**ABSTRACT**

Data-centric machine learning (ML) pipelines extend traditional ML pipelines—of feature transformations, hyper-parameter tuning, and model training—by additional pre-processing steps for data cleaning, data augmentation, and feature engineering to create high-quality data with good coverage. Finding effective data-centric ML pipelines is still a labor- and compute-intensive process though. While AutoML tools use effective search strategies, they struggle to scale with large datasets. Large language models (LLMs) show promise for code generation but face challenges in generating data-centric ML pipelines due to private datasets not seen during training, complex pre-processing requirements, and the need for mitigating hallucinations. These demands exceed typical code generation as it requires actions tailored to the characteristics and requirements of a particular dataset. This paper introduces CatDB, a comprehensive, LLM-based system for generating effective, error-free, and efficient data-centric ML pipelines. CatDB leverages data catalog information and refined metadata to dynamically create dataset-specific rules (instructions) to guide the LLM. Moreover, CatDB includes a robust mechanism for automatic validation and error handling of the generated pipeline. Our experimental results show that CatDB reliably generates effective ML pipelines across diverse datasets, achieving accuracy comparable to or better than existing LLM-based systems, standalone AutoML tools, and combined workflows of data cleaning and AutoML tools, while delivering up to orders of magnitude faster performance on large datasets.

CatDB: Data-catalog-guided, LLM-based Generation of Data-centric ML Pipelines, VLDB'24



Intransparent Experiments



Performance Drop for Tasks Published after LLM Training Cutoff Date

```
# Combine train and test data for preprocessing
combined_data = pd.concat([train_data, test_data], ignore_index=True)

# Preprocessing: Fill missing values and one-hot encode categorical features
for col in combined_data.columns:
    if combined_data[col].dtype == "object":
        combined_data[col] = combined_data[col].fillna(combined_data[col].mode()[0])
        le = LabelEncoder()
        combined_data[col] = le.fit_transform(combined_data[col])
```

Data Leakage

# What Would Be Nice To Have?

- An easy and robust way to control which parts of the pipeline
  - To write **by hand**
  - To generate & optimize **via an LLM**
- A runtime system that synthesizes & optimizes the code mimicking manual optimization by a data scientist
  - Users can **declaratively** specify the operations to delegate to LLMs
  - LLM generates **implementations tailored** to the users **pipeline**

# Declarative Programming with Semantic Operators

- **Core idea:** Users embed semantic operators with NL instructions into their code to define **what** to compute, the underlying system decides **how** to compute it

- LLMs in SQL operators

```
1 SELECT COUNT(*)
2 FROM Cars
3 WHERE AI.IF(pic, 'the picture shows a red car');
```

SemBench: A Benchmark for Semantic Query Processing Engines, VLDB'26

```
1 def paper_digest(research_interest: str, baseline: str):
2     return papers_df\
3         .sem_search("abstracts", research_interest, K=100)\
4         .sem_filter(f"the paper {{abstracts}} claim to
5         outperform {baseline}")\
6         .sem_agg(f"Write a digest summarizing {{abstracts}}
7         and their relevance to {research_interest}")
```

Semantic Operators and Their Optimization: Enabling LLM-Based Data Processing with Accuracy Guarantees in LOTUS, VLDB'24

- Modular AI software automating the optimization of prompt construction

```
1 import dspy
2 from dspy.datasets import HotPotQA
3
4 dspy.configure(lm=dspy.LM("openai/gpt-4o-mini"))
5
6 def search_wikipedia(query: str) -> list[str]:
7     results = dspy.ColBERTv2(url="http://20.102.90.50:2017/wiki17_abstracts")(query, k=3)
8     return [x["text"] for x in results]
9
10 trainset = [x.with_inputs('question') for x in HotPotQA(train_seed=2024, train_size=500).train]
11 react = dspy.ReAct("question -> answer", tools=[search_wikipedia])
12
13 tp = dspy.MIPROv2(metric=dspy.evaluate.answer_exact_match, auto="light", num_threads=24)
14 optimized_react = tp.compile(react, trainset=trainset)
```

# SemPipes

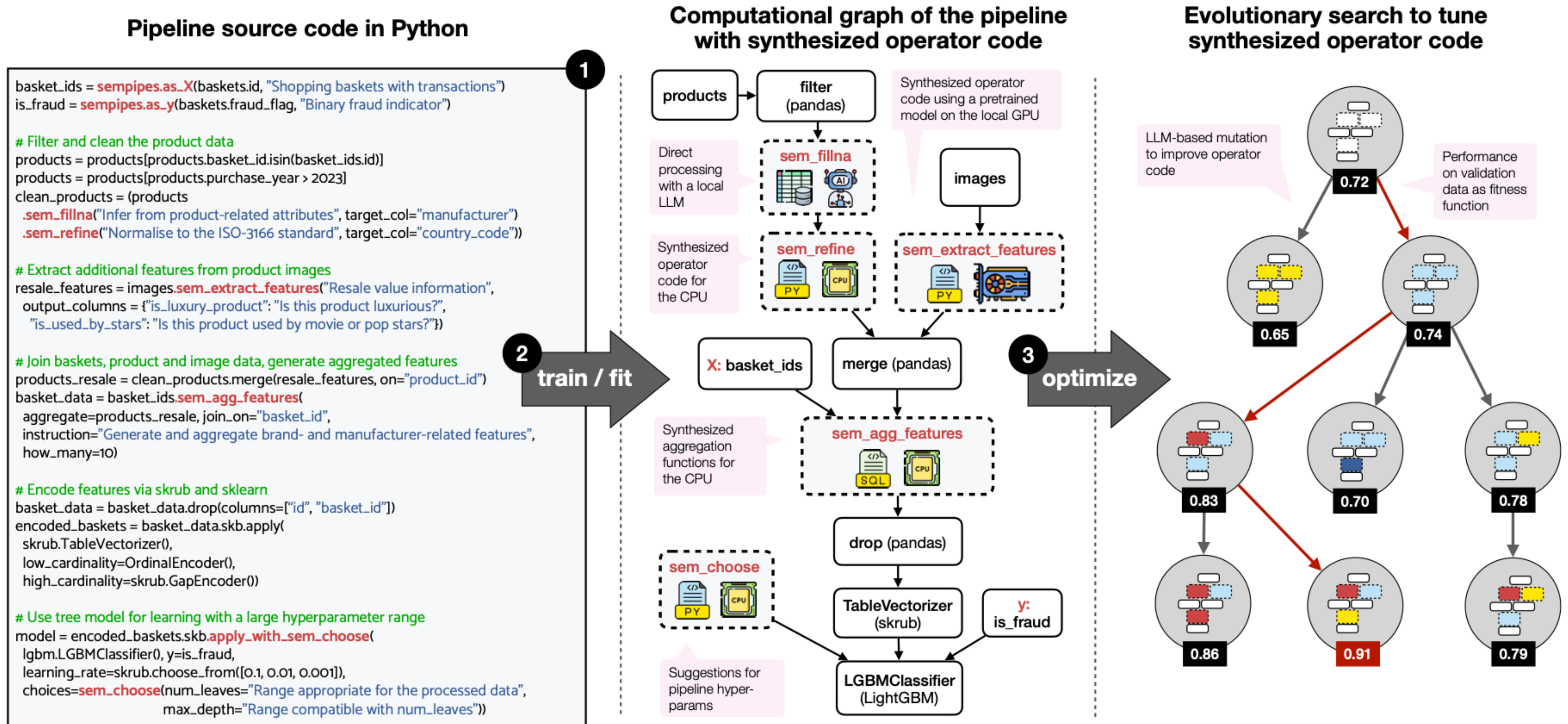


- **SemPipes** introduces **AI-assisted primitives** for writing (tabular) ML pipelines
- A novel declarative programming model that integrates **LLM-powered optimizable semantic data operators**
  - A new execution model: Training-time code synthesis, inference-time execution
  - Leverages skrub DataOps pipelines, a lightweight abstraction for multi-table ML pipelines

# SemPipes



- **SemPipes synthesizes custom operator implementations** based on data characteristics, user instructions in natural language, and pipeline context



ML pipeline combining code from common data science libraries with **semantic data operators from SemPipes**.

At training time, **SemPipes synthesizes custom operator implementations** based on data characteristics, user instructions and the pipeline context.

SemPipes can **automatically tune the synthesized operator code via evolutionary search** on a validation set.



## Pipeline code with semantic data operator

```
baskets = sempipes.X("baskets", "shopping baskets")
is_fraud = sempipes.y("fraud_flags", "fraud indicator")
products = skrub.var("products", "sales data with product attributes")

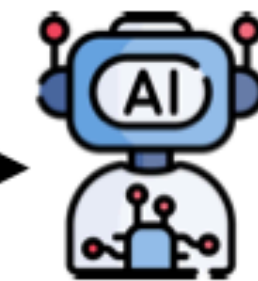
# Join baskets with filtered product data and generate aggregated features
recent_products = products[products.purchase_year > 2023]
basket_features = baskets.sem_agg_features(
    aggregate=recent_products, join_on="basket_id", how_many=4,
    instruction="Generate brand- and price-related features indicative of fraud")

model = basket_features.skb.apply(lgbm.LGBMClassifier(), y=is_fraud)

pipeline = model.skb.make_learner() # Turn pipeline into fittable learner
pipeline.fit({"baskets": ..., "fraud_flags": ..., "products": ... })
```

## Context computed at pipeline training time

- Characteristics of input data  
Columns in recent\_products:  
- brand (str): apple, ...  
- price (float): min 5.0, ...  
Columns in baskets: ...
- User instruction for operator  
Generate brand- and price-related features ...
- Pipeline and task summary  
Classification of shopping baskets with fraud indicator as target via LGBMClassifier
- Operator-specific hints  
Use a left join and do not remove rows from the left input!



## Synthesis of operator code tailored to pipeline and data

```
def aggregate_and_join(baskets, recent_products):
    features = recent_products.groupby("basket_id").agg(
        total = ("cash_price", "sum"), # Total cost of the basket
        # Maximum item price, high value indicates fraud
        max_price = ("cash_price", "max"),
        # Number of unique brands, low value indicates fraud
        num_brands = ("brand", "nunique"))
    # Price concentration in basket, high value indicates fraud
    features["price_con"] = features.max_price / features.total
    return baskets.merge(features, on="basket_id", how="left")
```

## Code validation

- Syntax check
- Execution on data sample
- Operator-specific output constraints

Retry with error feedback (in case of validation errors)

- **Training**

- Synthesize operator code based on input data characteristics, pipeline context & prompt
- Validate by executing code on a sample & check operator-specific output constraints

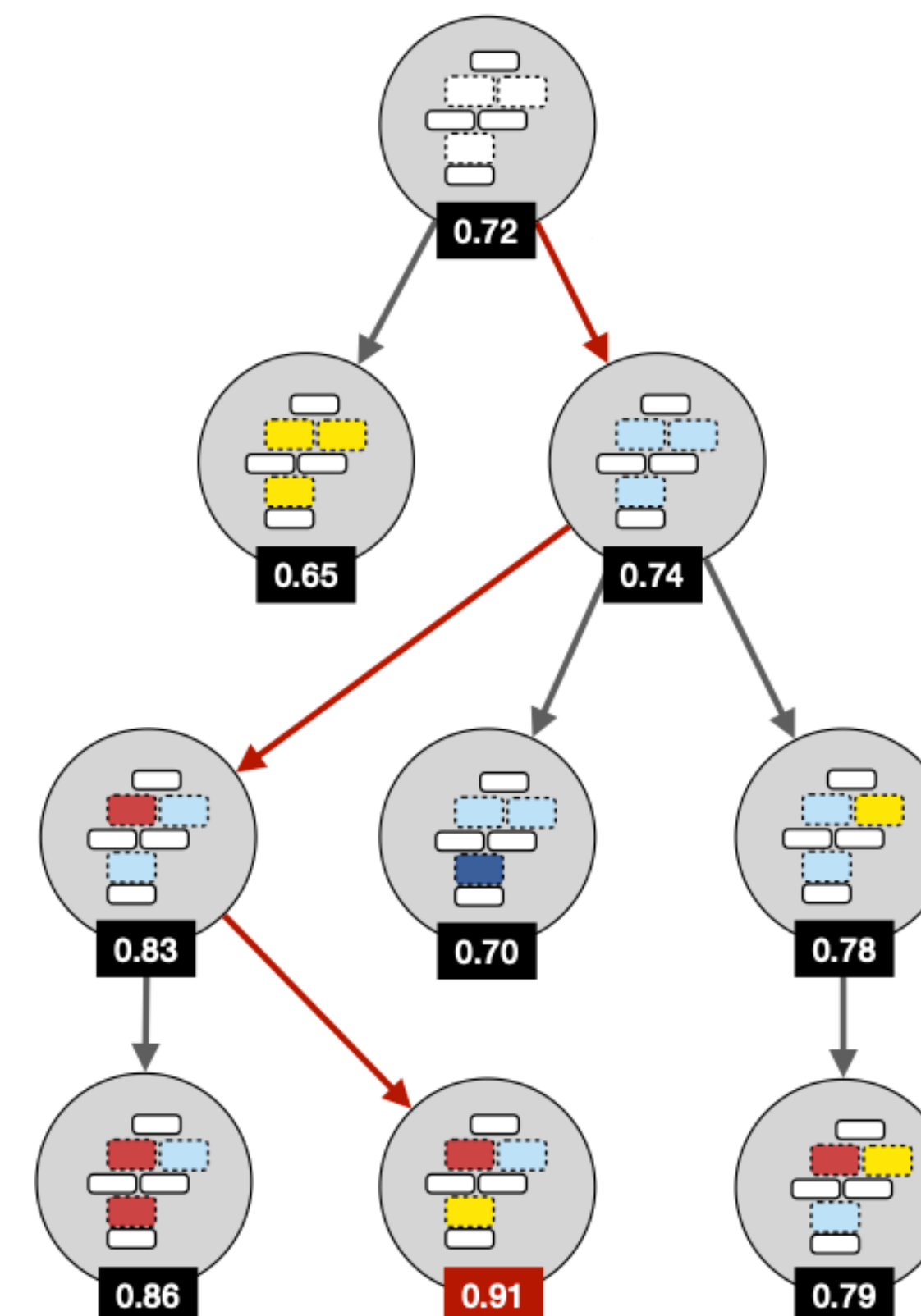
- **Inference**

- Reuse synthesized code

# SemPipes Optimizer



- **Refines** the LLM-synthesized implementations and **iteratively improves** generated code and predictive performance
- LLM suggests a better implementation based on history, predictive score on validation data, and previous “good” implementations
  - **Monte Carlo tree search** to improve the **each operator**
  - Select which operator to optimize in each round using **multi-armed bandits** approach
  - Evaluate improvements & **predictive performance** on the **validation set**



# SemPipes Operators



| Category                        | Semantic Operators                                       |
|---------------------------------|--|
| Data Cleaning                   | sem_fillna, sem_clean, sem_refine                        |
| Feature Extraction & Generation | sem_extract_features, sem_gen_features, sem_agg_features |
| Data Augmentation               | sem_augment  |
| Parameterization                | sem_select, sem_choose                                   |

Introduced semantic operators

- **Data cleaning & augmentation operators**

- **Logical:** Fill missing values
- **Physical:** Select one of 5 imputation models suggested by LLM

```
products = products.sem_fillna(  
    target_column="make",  
    nl_prompt="Infer the manufacturer from relevant product-related attributes like title or description.",  
    impute_with_existing_values_only=True,  
)
```

- **Logical:** Refine column's unique values using parametric knowledge
- **Physical:** Create correspondence dictionaries

```
cleaned_countries_ref = countries_ref.sem_refine(  
    nl_prompt="Make sure that all values are in the ISO 3166-1 alpha-2 two-letter uppercase format",  
    target_column="country_code",  
    refine_with_existing_values_only=False,  
)  
.skb.eval()
```

- **Logical:** Data augmentation (rows)
- **Physical:** Augment data using Synthetic Data Vault (SDV) Python library

```
salaries_augmented = salaries.sem_augment(  
    nl_prompt="Augment data to improve 'current_annual_salary' prediction.",  
    number_of_rows_to_generate=2000,  
    name="augment_salaries",  
    generate_via_code=False,  
)
```

# SemPipes Operators



| Category                        | Semantic Operators                                       |
|---------------------------------|--|
| Data Cleaning                   | sem_fillna, sem_clean, sem_refine                        |
| Feature Extraction & Generation | sem_extract_features, sem_gen_features, sem_agg_features |
| Data Augmentation               | sem_augment  |
| Parameterization                | sem_select, sem_choose                                   |

Introduced semantic operators

- **Feature extraction operators**

```
movie_stats = movie_stats.sem_gen_features(  
    nl_prompt=""  
    """  
    Create additional features that could help predict the box office revenue of a movie.  
    Consider aspects like genre, production details, cast, crew, and any other relevant information  
    that could influence a movie's financial success. Some of the attributes are in JSON format,  
    so you might need to parse them to extract useful information.  
    """,  
    name="additional_movie_features",  
    how_many=25,  
)
```

- **Logical:** Generate features from existing columns
- **Physical:** Combine/remove existing columns

```
aggregated = baskets.sem_agg_features(  
    products,  
    left_on="ID",  
    right_on="basket_ID",  
    nl_prompt=""  
    """  
    Generate product features that are indicative of potentially fraudulent baskets, make it easy to distinguish  
    anomalous baskets from regular ones! It might be helpful to combine different product statistics.  
    """,  
    name="basket_agg_features",  
    how_many=1,  
)
```

- **Logical:** Aggregate columns and calculate new aggregated columns
- **Physical:** Left join with custom aggregations

```
labeled_tweets = tweets.sem_extract_features(  
    nl_prompt="Extract features from textual tweets that could help predict toxicity.",  
    input_columns=["text"],  
    name="tweet_features",  
    output_columns=output_columns,  
    generate_via_code=True,  
    print_code_to_console=True,  
)
```

- **Logical:** Extract features from text/images/audio
- **Physical:** Pre-trained models, regexes



- Selected results of **SemPipes** operators
  - Semantic operators improve and simplify ML pipelines written by humans and agents

| Data-Centric Task                          | Metric              | Pipeline         | Origin       |              | Original Pipeline                  | w/ SemOps & Instructions | w/ Optimized SemOps                |
|--|---------------------|------------------|--------------|--------------|------------------------------------|--------------------------|------------------------------------|
|  |                     |                  | Expert       | Agent        |                                    |                          |                                    |
| Extraction of clinically informed features | F1 Score $\uparrow$ | micromodels      | $\checkmark$ |              | $0.665 \pm 0.01$                   | $0.692 \pm 0.02$         | <b><math>0.729 \pm 0.00</math></b> |
|  |                     | swe-micromodels  |              | $\checkmark$ | $0.690 \pm 0.00$                   | $0.692 \pm 0.02$         | <b><math>0.729 \pm 0.00</math></b> |
| Blocking for entity resolution             | Recall $\uparrow$   | rutgers          | $\checkmark$ |              | $0.238 \pm 0.00$                   | $0.240 \pm 0.02$         | <b><math>0.255 \pm 0.01</math></b> |
|  |                     | sustech          | $\checkmark$ |              | <b><math>0.322 \pm 0.00</math></b> | $0.320 \pm 0.00$         | $0.321 \pm 0.00$                   |
|  |                     | aide-sigmoid     |              | $\checkmark$ | $0.043 \pm 0.00$                   | $0.163 \pm 0.03$         | <b><math>0.165 \pm 0.03</math></b> |
|  |                     | baseline-sigmoid | $\checkmark$ |              | $0.099 \pm 0.00$                   | $0.141 \pm 0.02$         | <b><math>0.173 \pm 0.05</math></b> |
| Feature engineering (scrabble rating)      | RMSE $\downarrow$   | kaggle-scrabble  | $\checkmark$ |              | $176.1 \pm 24.7$                   | $149.2 \pm 24.8$         | <b><math>144.4 \pm 18.3</math></b> |
|  |                     | aide-scrabble    |              | $\checkmark$ | $255.0 \pm 24.9$                   | $195.2 \pm 21.9$         | <b><math>186.7 \pm 21.1</math></b> |
|  |                     | swe-scrabble     |              | $\checkmark$ | $228.5 \pm 12.9$                   | $210.1 \pm 43.0$         | <b><math>187.8 \pm 21.4</math></b> |
| Data annotation                            | Accuracy $\uparrow$ | hibug            | $\checkmark$ |              | $0.598 \pm 0.05$                   | $0.758 \pm 0.15$         | <b><math>0.766 \pm 0.16</math></b> |
| Data augmentation for fairness             | AUROC $\uparrow$    | sivep            | $\checkmark$ |              | $0.653 \pm 0.03$                   | $0.682 \pm 0.00$         | <b><math>0.710 \pm 0.04</math></b> |
|  |                     | swe-fair         |              | $\checkmark$ | $0.650 \pm 0.04$                   | $0.682 \pm 0.02$         | <b><math>0.710 \pm 0.04</math></b> |

# SemPipes Operator Evaluation



- **sem\_gen\_features** for the feature generation from existing data/columns
- **SemPipes** is on par with LLM feature engineering methods like CAAFE (NeurIPS'24)
  - Larger context such as data characteristics and pipeline & task summary

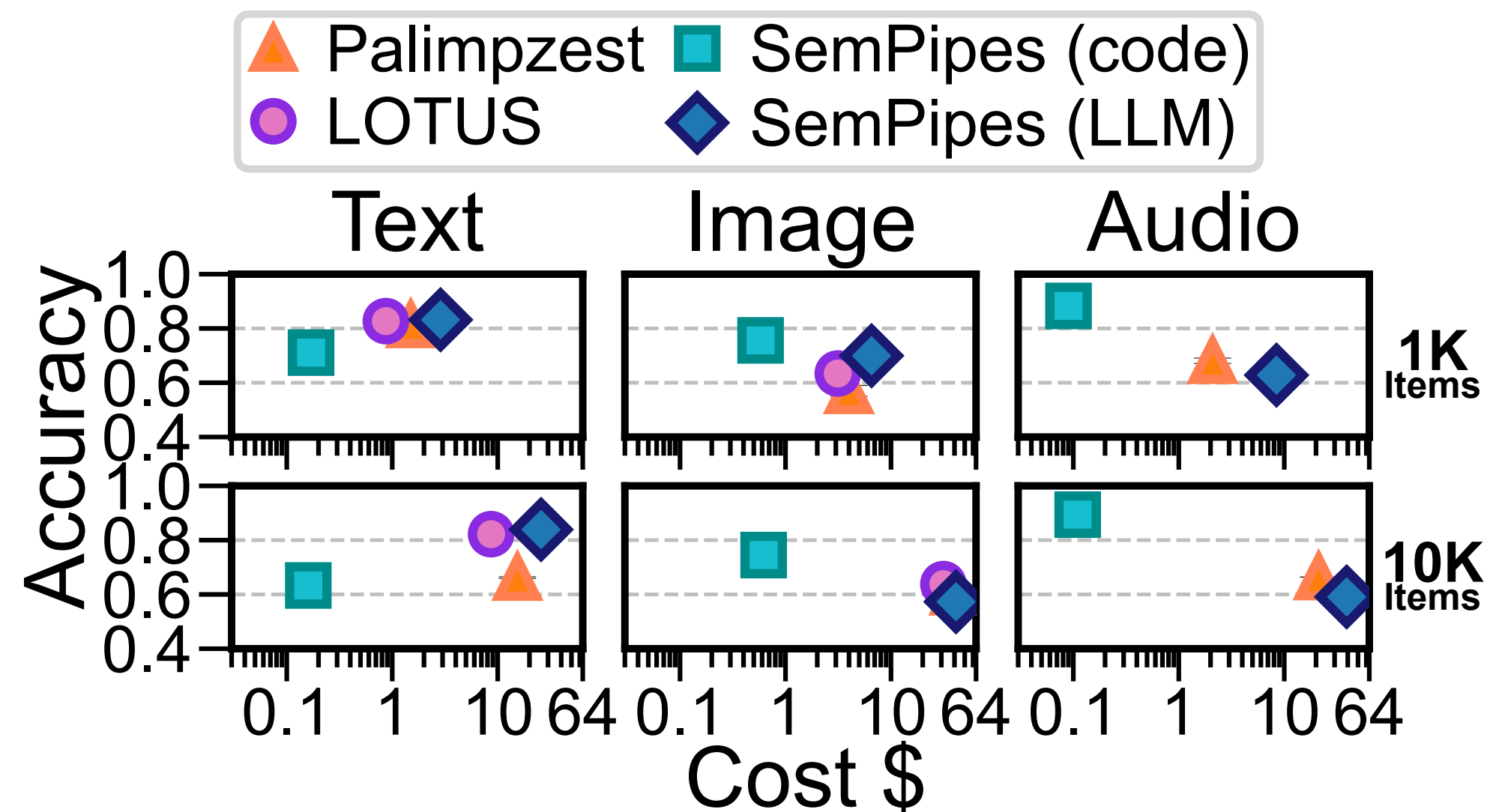
|        | Dataset       | No FE           | CAAFE                  | SemPipes               |
|--------|---------------|-----------------|------------------------|------------------------|
| RF     | balance-scale | 0.685<br>± 0.13 | 0.737<br>± 0.19        | <b>0.884</b><br>± 0.23 |
|        | tic-tac-toe   | 0.807<br>± 0.07 | 0.878<br>± 0.11        | <b>0.929</b><br>± 0.08 |
|        | airlines      | 0.629<br>± 0.03 | <b>0.631</b><br>± 0.03 | 0.625<br>± 0.03        |
| TabPFN | balance-scale | 0.848<br>± 0.25 | 0.863<br>± 0.24        | <b>0.938</b><br>± 0.11 |
|        | tic-tac-toe   | 0.959<br>± 0.02 | 0.945<br>± 0.07        | <b>0.997</b><br>± 0.00 |
|        | airlines      | 0.647<br>± 0.03 | 0.647<br>± 0.03        | <b>0.650</b><br>± 0.03 |

AUROC scores without feature engineering (No FE), and for automated feature engineering on three datasets with Random Forest (RF) and TabPFN. **SemPipes** uses **sem\_gen\_features**

# SemPipes Operator Evaluation



- **sem\_extract\_features** for zero-shot feature extraction from multi-modal data
- **SemPipes** outperforms direct LLM data processing for specialized data
  - Text - financial clause classification (DeBERTa)
  - Image - pneumonia detection in X-Rays (BiomedCLIP)
  - Audio - environmental sound extraction (CLAP)

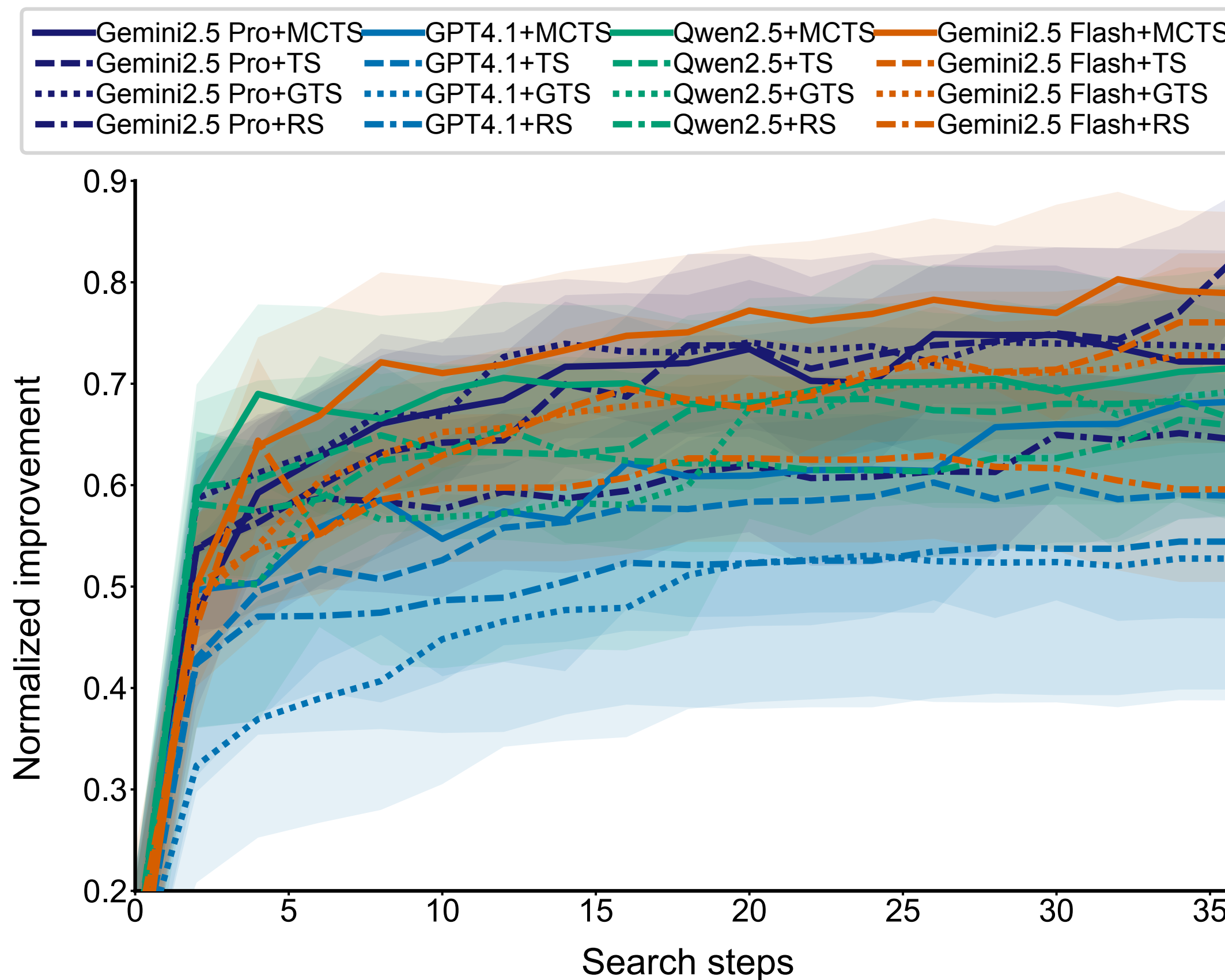


Zero-shot feature extraction via direct data processing and code generation

# Smaller vs Larger Models



- **SemPipes** contextualization allows to achieve competitive performance with smaller models



Optimization effectiveness of LLMs for code synthesis in semantic data operators. Search policies include Monte Carlo tree search (MCTS), truncation selection (TS), greedy tree search (GTS) or random search (RS)

# SemPiper “plays” SemPipes



Pipeline Optimizer

Pipeline: Fraud (simple) Model: gemini/gemini-2.5-flash-lite Temperature: 0.0

```
1 import skrub
2 from sklearn.ensemble import HistGradientBoostingClassifier
3 import sempipes
4
5
6 def sempipes_pipeline():
7     products = skrub.var("products")
8     baskets = skrub.var("baskets")
9
10    fraud_flags = sempipes.as_y(baskets["fraud_flag"],
11    "Fraud label")
12    basket_ids = sempipes.as_X(baskets[["ID"]], "Shopping
13    baskets")
14
15    kept_products =
16    products[products["basket_ID"].isin(basket_ids["ID"])]
17    kept_products = kept_products.sem_gen_features(
18    nl_prompt="Generate useful features for product
19    analysis.",
20    name="product_features",
21    how_many=2,
22    )
23
24    augmented_baskets = basket_ids.sem_agg_features(
25    kept_products,
26    left_on="ID",
27    right_on="basket_ID",
28    nl_prompt="Aggregate product features per basket for
29    fraud detection.",
30    name="basket_features",
31    how_many=1,
32    )
33
34    hgb = HistGradientBoostingClassifier()
35    fraud_detector = augmented_baskets.skb.apply(hgb,
36    y=fraud_flags)
37
38    return fraud_detector
```

Computational graph

Click an operator to see generated code

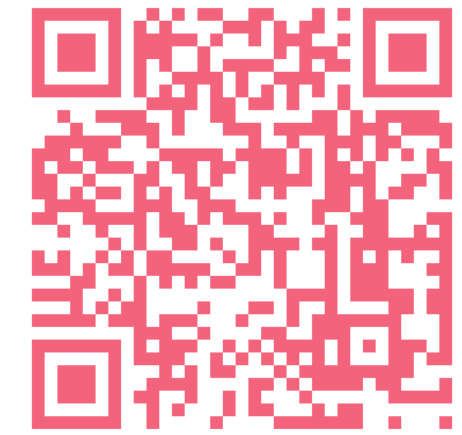
```
graph TD
    V_baskets["<Var 'baskets'>"] --> G_ID1["<GetItem 'ID'>"]
    V_baskets --> G_fraud_flag["<GetItem 'fraud_flag'>"]
    V_products["<Var 'products'>"] --> G_ID2["<GetItem 'ID'>"]
    G_ID1 --> CM_isin["<CallMethod 'isin'>"]
    G_fraud_flag --> CM_isin
    G_ID2 --> CM_isin
    CM_isin --> G_CM_isin["<GetItem '<CallMethod 'isin'>'>"]
    G_CM_isin --> S_gen_features["sem_gen_features"]
    S_gen_features --> S_agg_features["sem_agg_features"]
    S_agg_features --> SKB_apply["skb.apply"]
```

Code details

Select a node in the graph to see its data summary, generated code, or LLM stats.

# Conclusions

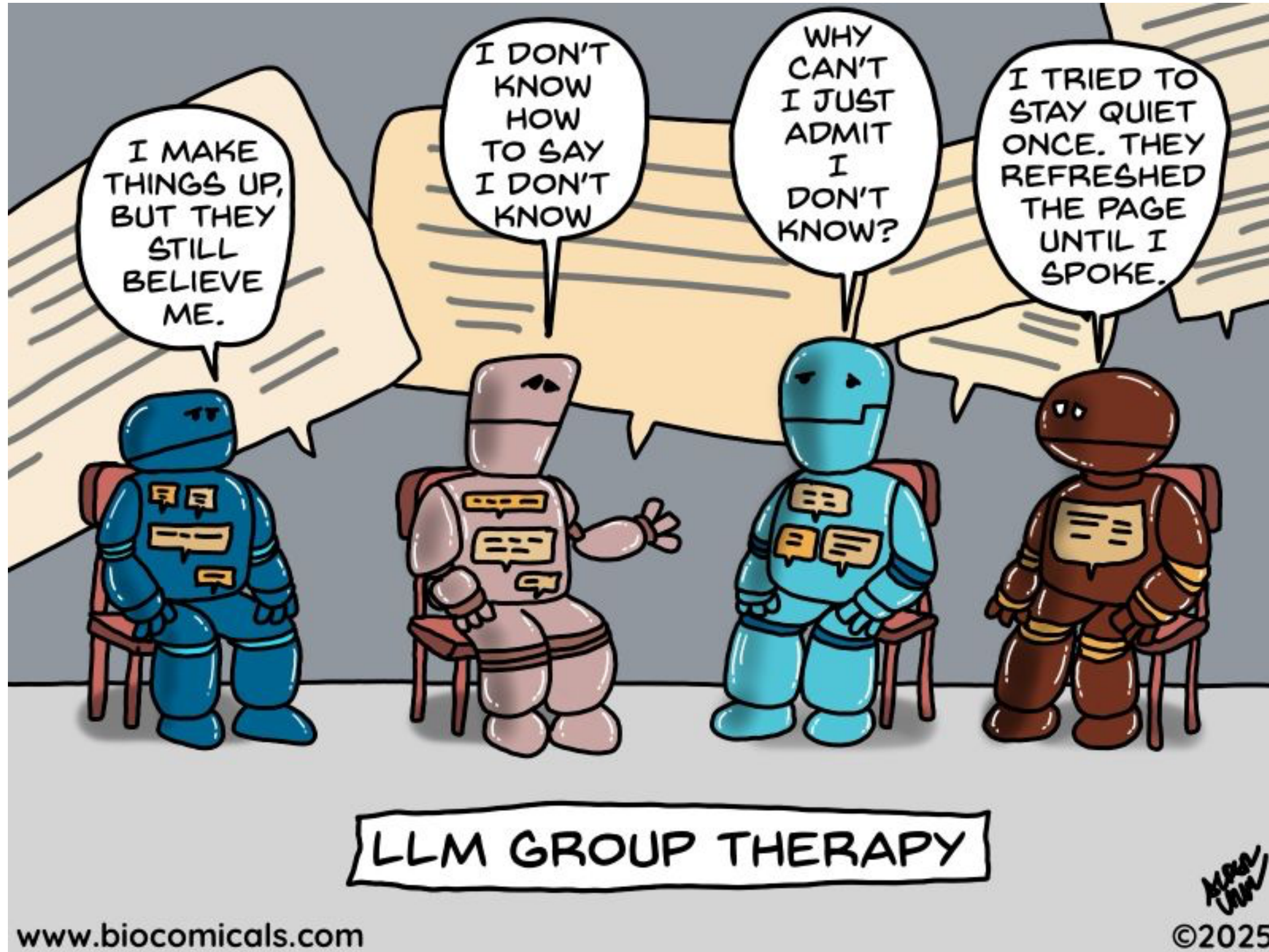
- Delegate data-heavy, time-consuming tasks to LLMs via code generation
- Utilize existing resources such as pre-trained HuggingFace models
- **Limitations**
  - Security risks if **SemPipes** is run in a not security-hardened sandbox
  - Limited support of videos, time series, and documents
- **Future work**
  - **SemAgent**: Generate end-to-end pipelines using **SemPipes** operators
    - Localize tasks within a pipeline and evolve each task separately



**SemPipes** research  
prototype



# Thanks



**SemPipes**



**Contact**

# Skrub & DataOps



- **Eases preprocessing & feature engineering** for tabular ML and bridges between dataframe libraries (pandas, polars) and ML libraries (sklearn, scipy)
- Allows constructions of pipelines that can be **fully tuned, inspected, or applied** to new data
- **DataOps**
  - Special objects that encapsulate operations on data (applying operators, calling methods)
  - Replayable on new data and fitting/predicting with a pipeline w/o any code rewrites
  - Combinable into a DataOps plan, a directed acyclic graph (DAG) of operations

```
>>> import skrub
>>> a = skrub.var("a")
>>> b = skrub.var("b")
>>> c = a + b
>>> c
<BinOp: add>
>>> c.skb.eval({"a": 10, "b": 6})
16
```

# A Very Simple Example Pipeline

- Credit fraud dataset with products and baskets tables
- **Task:** Predict if a basket is fraudulent (binary classification, highly unbalanced classes)

| Product |           |     | Basket |              |
|---------|-----------|-----|--------|--------------|
| ID      | basket_ID | ... | ID     | target_label |
| 1       | 1         |     | 1      | 0            |
| 2       | 1         |     | 2      | 1            |
| 3       | 2         | ... | 3      | 0            |

```
def pipeline() -> skrub.DataOp:
    products = skrub.var("products")
    baskets = skrub.var("data")
    labels = skrub.var("labels")

    baskets = baskets.skb.mark_as_X().skb.set_description("Potentially fraudulent shopping baskets of products")
    labels = labels.skb.mark_as_y().skb.set_description(
        "Flag indicating whether the basket is fraudulent (1) or not (0)"
    )

    aggregated = baskets.sem_agg_features(
        products,
        left_on="ID",
        right_on="basket_ID",
        nl_prompt=""
        """
        Generate product features that are indicative of potentially fraudulent baskets, make it easy to distinguish
        anomalous baskets from regular ones! It might be helpful to combine different product statistics.
        """
        ,
        name="basket_agg_features",
        how_many=1,
    )

    encoded = aggregated.skb.apply(TableVectorizer())
    return encoded.skb.apply(HistGradientBoostingClassifier(random_state=0), y=labels)
```

# A Very Simple Example Pipeline, SemPipes Generated Code

- Left join to aggregate features with SemPipes `sem_agg_features`

```
def _sem_agg_join(left_join_column, left_df, right_join_column, right_df):  
  
    # Aggregate required features at the basket_ID level  
    agg_df = right_df.groupby(right_join_column).agg(  
        total_basket_cash_price=('cash_price', 'sum'),  
        total_basket_product_units=('Nbr_of_prod_purchas', 'sum'),  
        # This feature counts the number of distinct 'item' categories purchased within each basket.  
        num_unique_items_in_basket=('item', 'nunique'),  
        # This feature captures the maximum cash price of any single product entry within a given basket.  
        max_item_cash_price_in_basket=('cash_price', 'max'),  
        # This feature counts how many items within a basket are categorized as service or fulfilment.  
        num_service_fulfilment_items_in_basket=('is_service_or_fulfilment', 'sum'),  
        # This feature calculates the standard deviation of `cash_price` for all items within a specific basket.  
        basket_cash_price_std=('cash_price', 'std'),  
        # Description: This feature counts the number of distinct 'make' (manufacturer/brand) types present within each basket.  
        num_unique_makes_in_basket=('make', 'nunique'),  
        # This feature calculates the minimum cash price of any product line within a given basket.  
        min_item_cash_price_in_basket=('cash_price', 'min'),  
        # This feature calculates the average cash price across all *line items* within a basket, providing a central tendency of  
        # item values.  
        mean_item_cash_price_in_basket=('cash_price', 'mean')  
    ).reset_index()  
  
    # Description: This feature calculates the average cash value contributed by each individual product unit within a basket.  
    agg_df['basket_item_value_per_product_unit'] = agg_df['total_basket_cash_price'] / agg_df['total_basket_product_units']  
    agg_df.loc[agg_df['total_basket_product_units'] == 0, 'basket_item_value_per_product_unit'] = 0  
  
    # Description: This feature calculates the ratio of the maximum single item's cash price to the total cash price of the entire  
    # basket.  
    agg_df['basket_price_concentration_max'] = agg_df['max_item_cash_price_in_basket'] / agg_df['total_basket_cash_price']  
    agg_df.loc[agg_df['total_basket_cash_price'] == 0, 'basket_price_concentration_max'] = 0  
  
    # Fill NaNs for basket_cash_price_std (occurs when a basket has only one item) with 0.  
    # A single item basket has no variation in price, so its standard deviation is 0.  
    agg_df['basket_cash_price_std'] = agg_df['basket_cash_price_std'].fillna(0)  
  
    # Drop the intermediate sum columns used for calculations, as they are not the final desired features  
    agg_df = agg_df.drop(columns=['total_basket_cash_price', 'total_basket_product_units'])  
  
    # Rename the right_join_column in agg_df to match left_join_column in left_df for merging  
    agg_df = agg_df.rename(columns={right_join_column: left_join_column})  
  
    # Left join the aggregated features back to the main dataframe  
    left_df = left_df.merge(agg_df, on=left_join_column, how='left')  
  
    return left_df
```

# A Very Simple Example Pipeline, SemPipes Generated Code

```
import pandas as pd
import numpy as np

def _sem_agg_join(left_join_column, left_df, right_join_column, right_df):
    agg_df = right_df.groupby(right_join_column).agg(
        total_basket_cash_price=('cash_price', 'sum'), total_basket_product_units=('Nbr_of_prod_purchas', 'sum'),
        num_unique_items_in_basket=('item', 'nunique'), max_item_cash_price_in_basket=('cash_price', 'max'),
        num_service_fulfilment_items_in_basket=('is_service_or_fulfilment', 'sum'), basket_cash_price_std=('cash_price', 'std'),
        min_item_cash_price_in_basket=('cash_price', 'min'), mean_item_cash_price_in_basket=('cash_price', 'mean')
    ).reset_index()

    Aggregate product-level features into basket-level features

    agg_df['basket_item_value_per_product_unit'] = agg_df['total_basket_cash_price'] / agg_df['total_basket_product_units']
    agg_df.loc[agg_df['total_basket_product_units'] == 0, 'basket_item_value_per_product_unit'] = 0

    agg_df['basket_price_concentration_max'] = agg_df['max_item_cash_price_in_basket'] / agg_df['total_basket_cash_price']
    agg_df.loc[agg_df['total_basket_cash_price'] == 0, 'basket_price_concentration_max'] = 0

    agg_df['basket_cash_price_std'] = agg_df['basket_cash_price_std'].fillna(0) Fill missing values

    agg_df = agg_df.drop(columns=['total_basket_cash_price', 'total_basket_product_units']) Clean up
    agg_df = agg_df.rename(columns={right_join_column: left_join_column})

    left_df = left_df.merge(agg_df, on=left_join_column, how='left') Left join

    return left_df
```