# **Texera:** A System for Collaborative and Interactive Data Analytics Using Workflows

Zuozhi Wang, **Yicong Huang**, Shengquan Ni, Avinash Kumar, Sadeem Alsudais, Xiaozhen Liu, Xinyuan Lin, Yunyan Ding, and Chen Li

UCIRVINE

VLDB2024
GUANGZHOU

# The goto Data Science Tool - Notebook
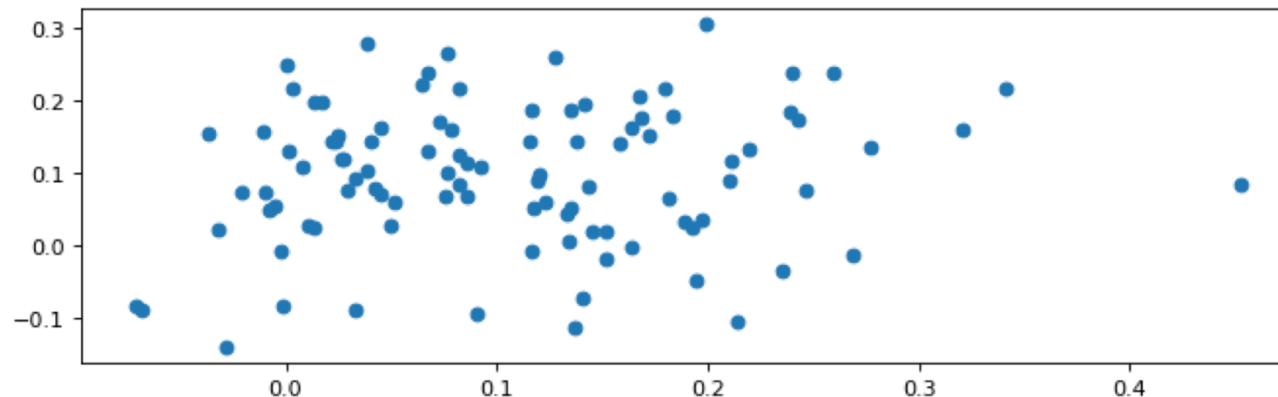
**Step 1: Loading Data**

In [2]:
```python
twenty_train = fetch_20newsgroups(subset='train')
```

**Step 2: Sentiment Analysis Step**

In [4]:
```python
text_clf = Pipeline([CountVectorizer(),TfidfTransformer(),SGDClassifier()])
text_clf.fit(twenty_train.data, twenty_train.target)
predicted = text_clf.predict(docs_test)
```
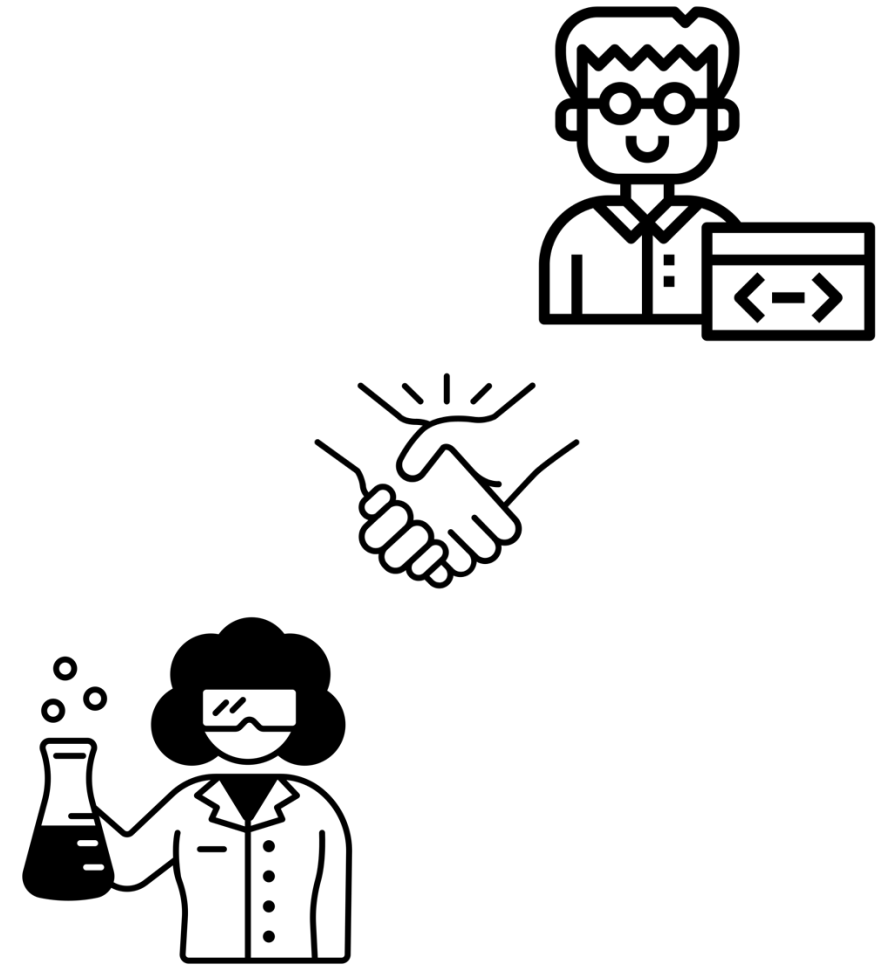
**Step 3: Plotting Data**
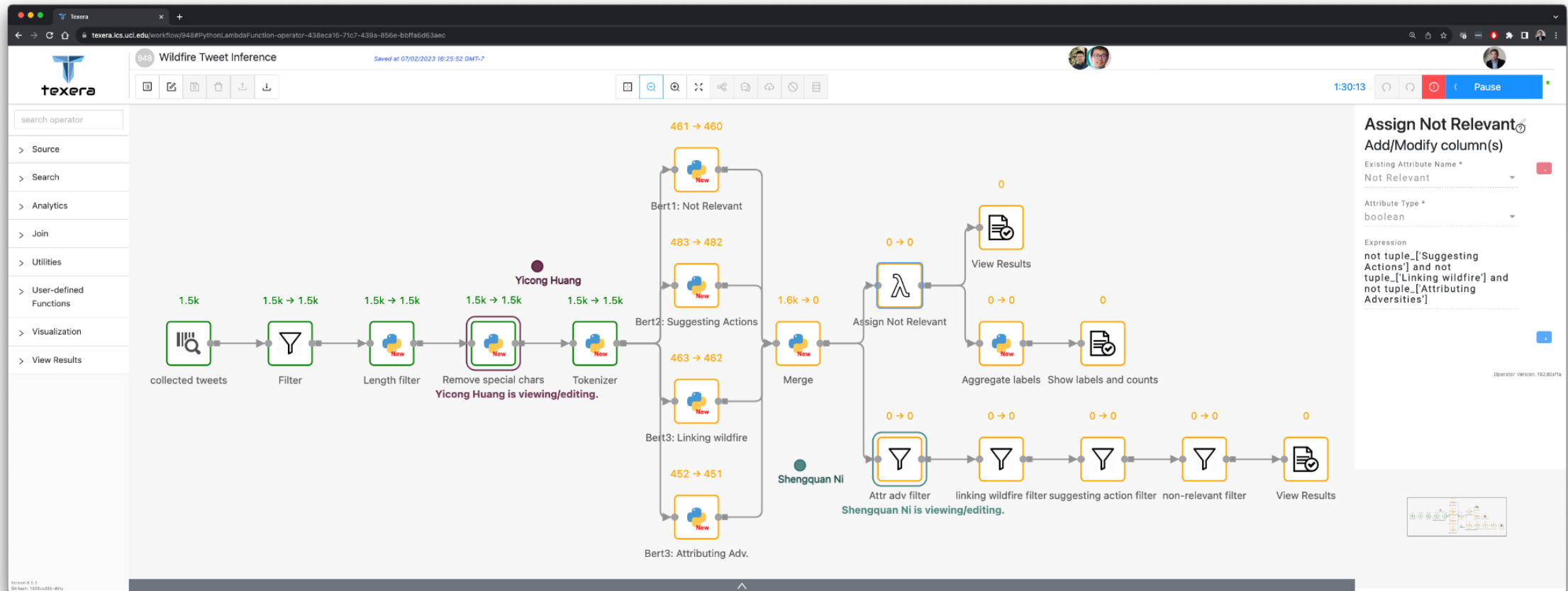
In [36]:
```python
ax = sns.scatterplot(data=predicted)
```

# The Need for Collaborative Data Analytics

- **IT experts**
  - Limited domain knowledge
  - Strong coding skills

- **Domain experts**
  - Rich domain knowledge
  - Limited IT/coding skills

# Introducing **texera**
# The **GUI-based Workflow** System for Data Analytics



| Cloud Service | Collaboration | Domain Friendly | AI/ML Access |

# More Features...

**Rich Built-in Operators**
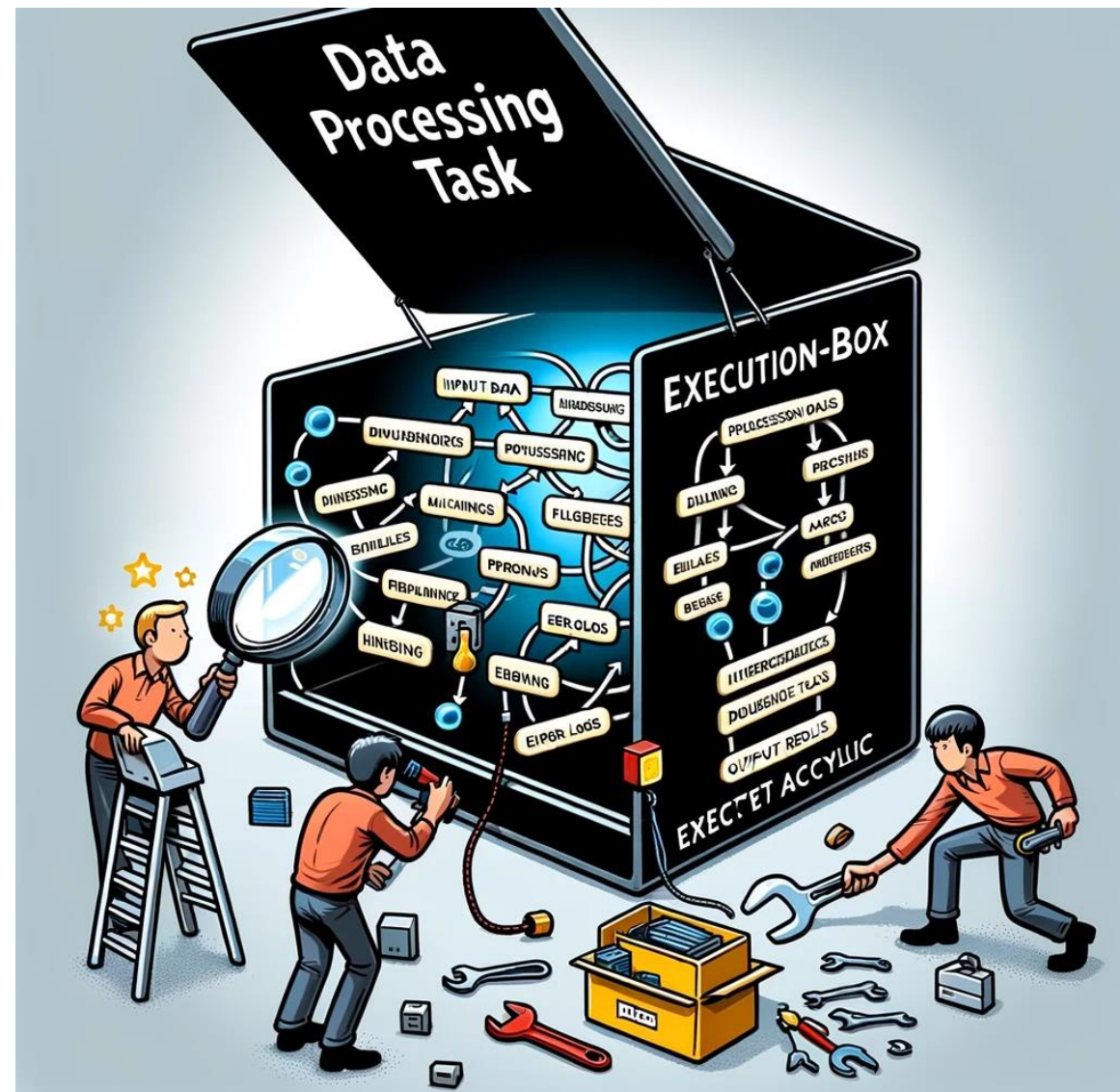
**Version & Restore**

# Collaboration only during Editing is not Enough for Data Analytics…

```
▷     data_path = "/kaggle/input/CORD-19-research-challenge/metadata.csv"
      source_data = preprocess_data(data_path, 100)
```
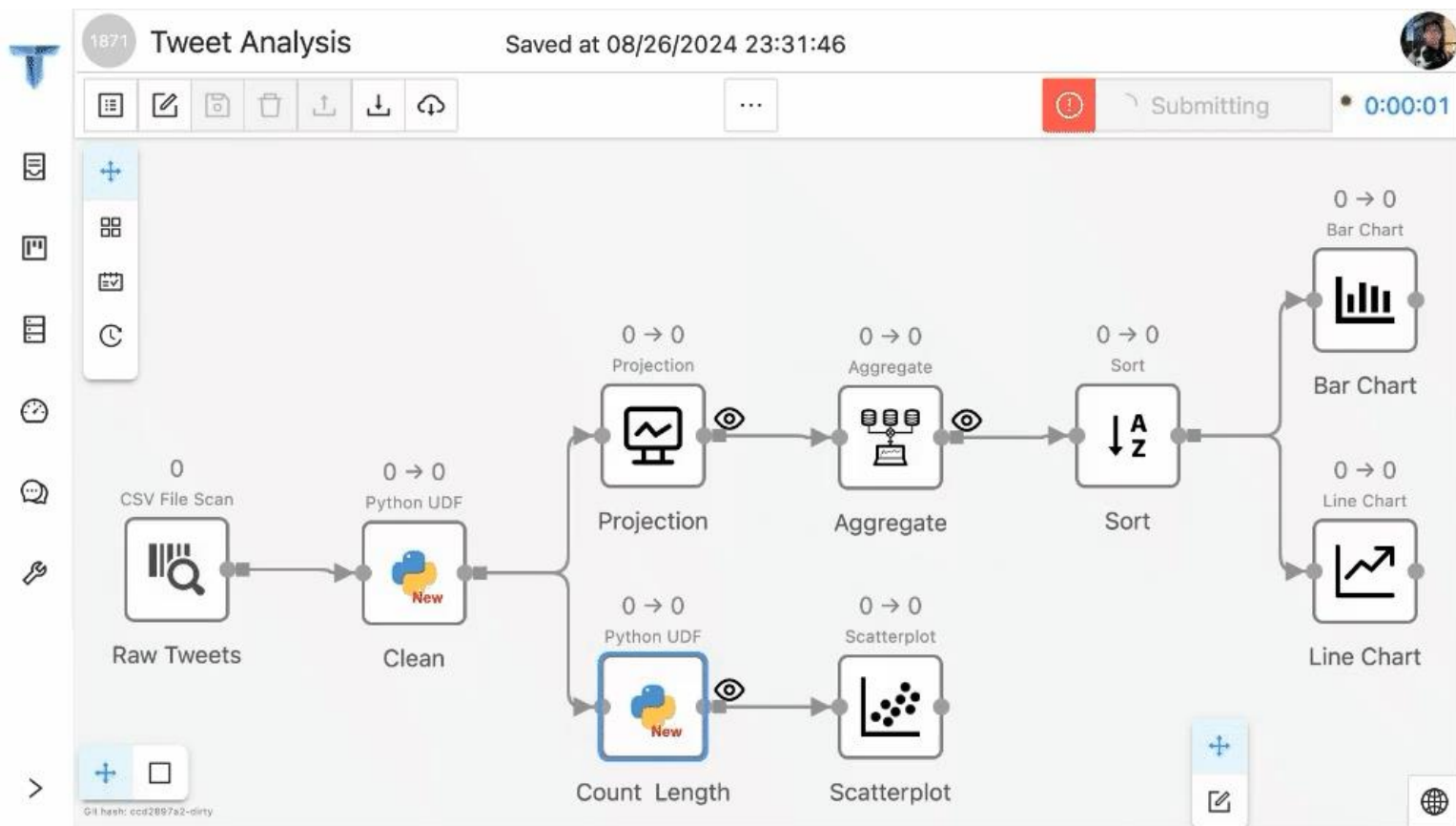
Unlike Google Doc or Overleaf, data analytics requires an (extensive) execution phase.

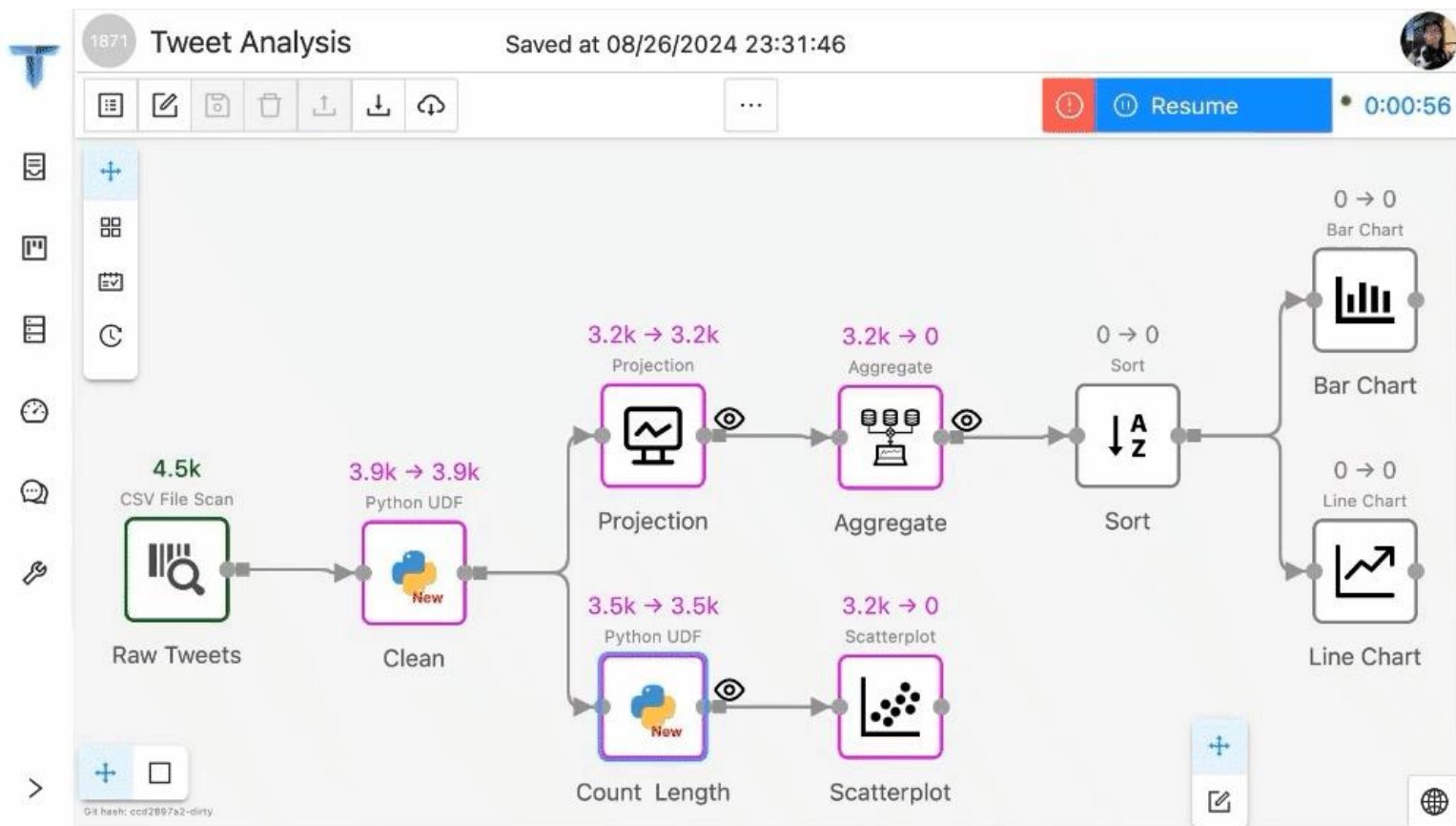# Interactions during Collaborative Execution
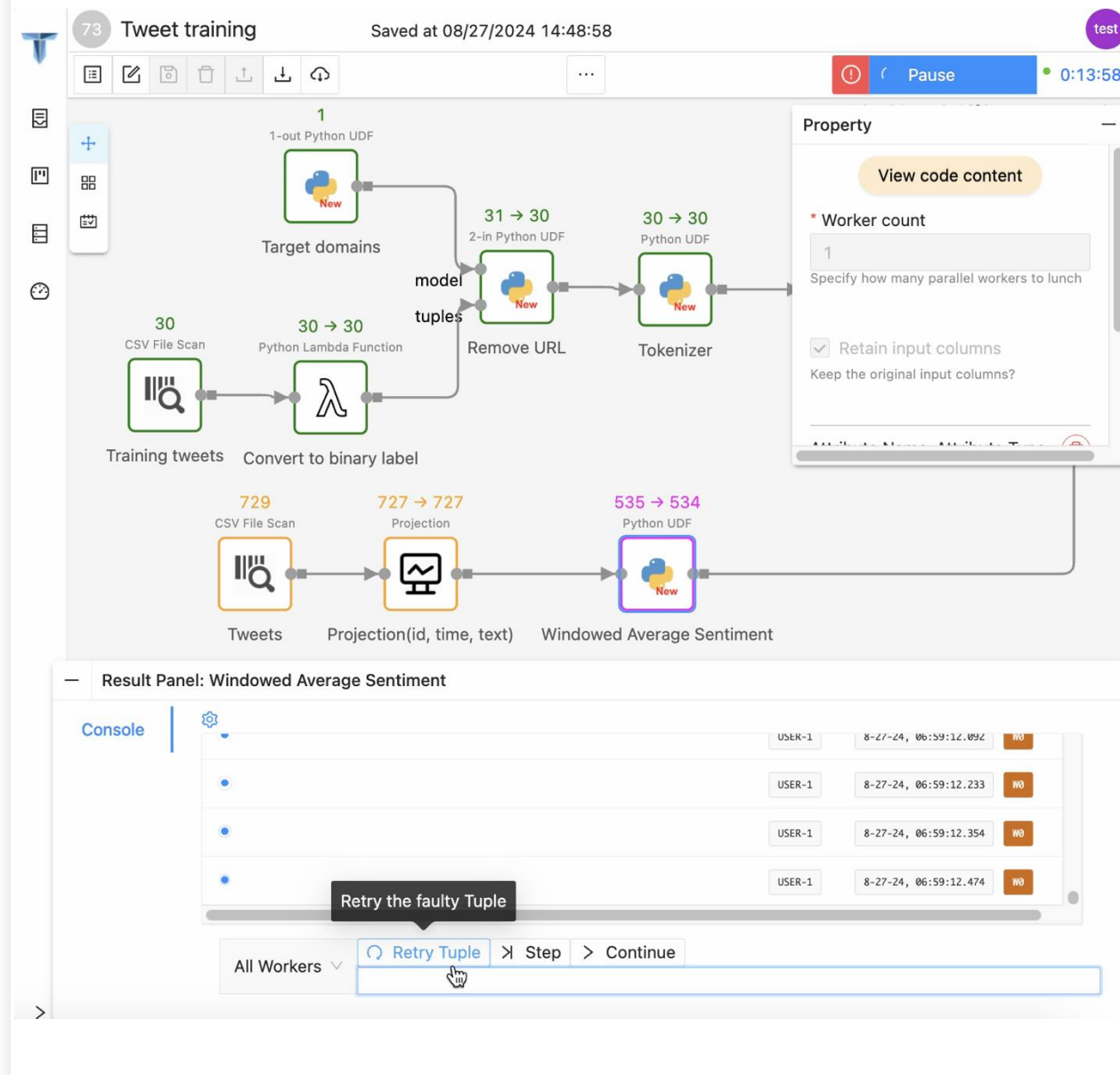
# What Kind of Interactions?

Interaction: Pause a Workflow

Interaction: Resume a Workflow

# Interaction: Read a Workflow's State

Interaction: Modify a Workflow's Logic

# Advanced Interaction – Debug

# How to Support Interactions?

# texera
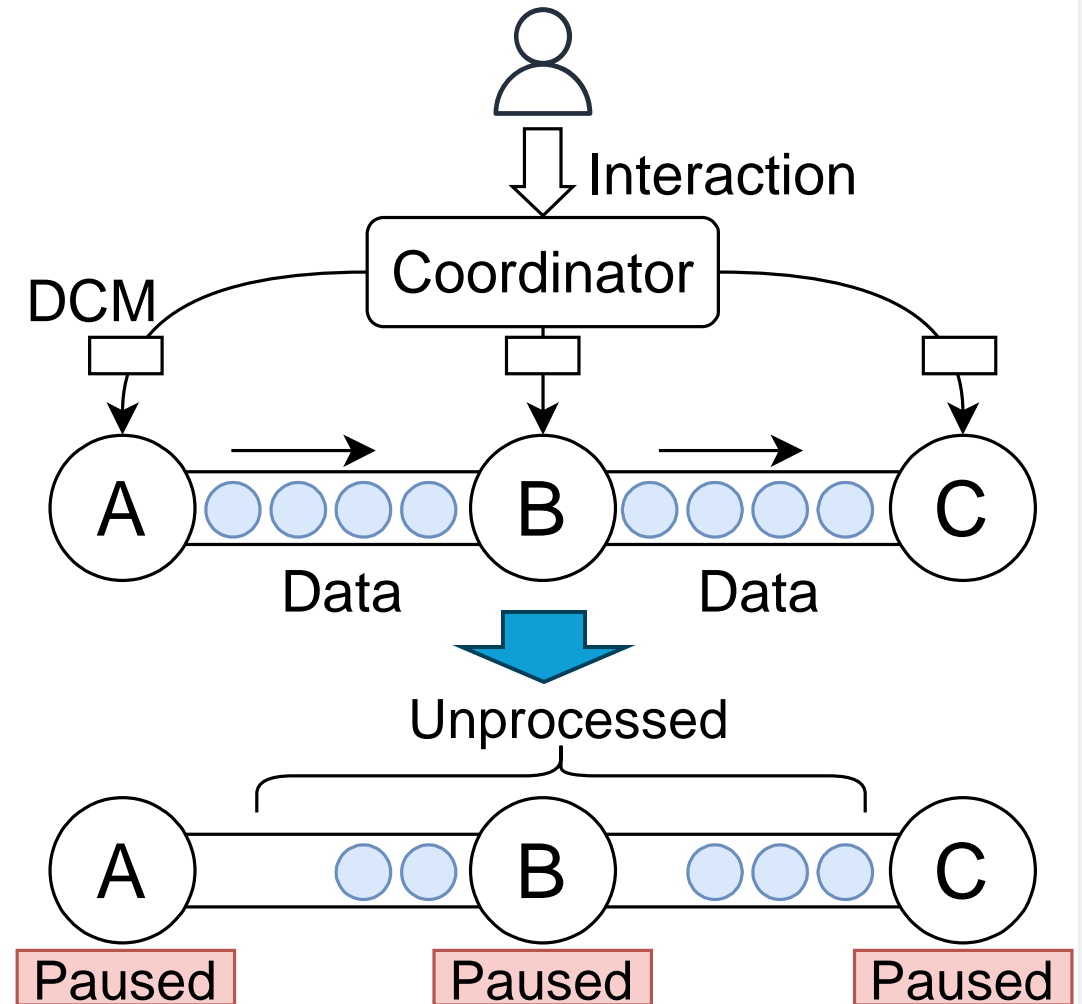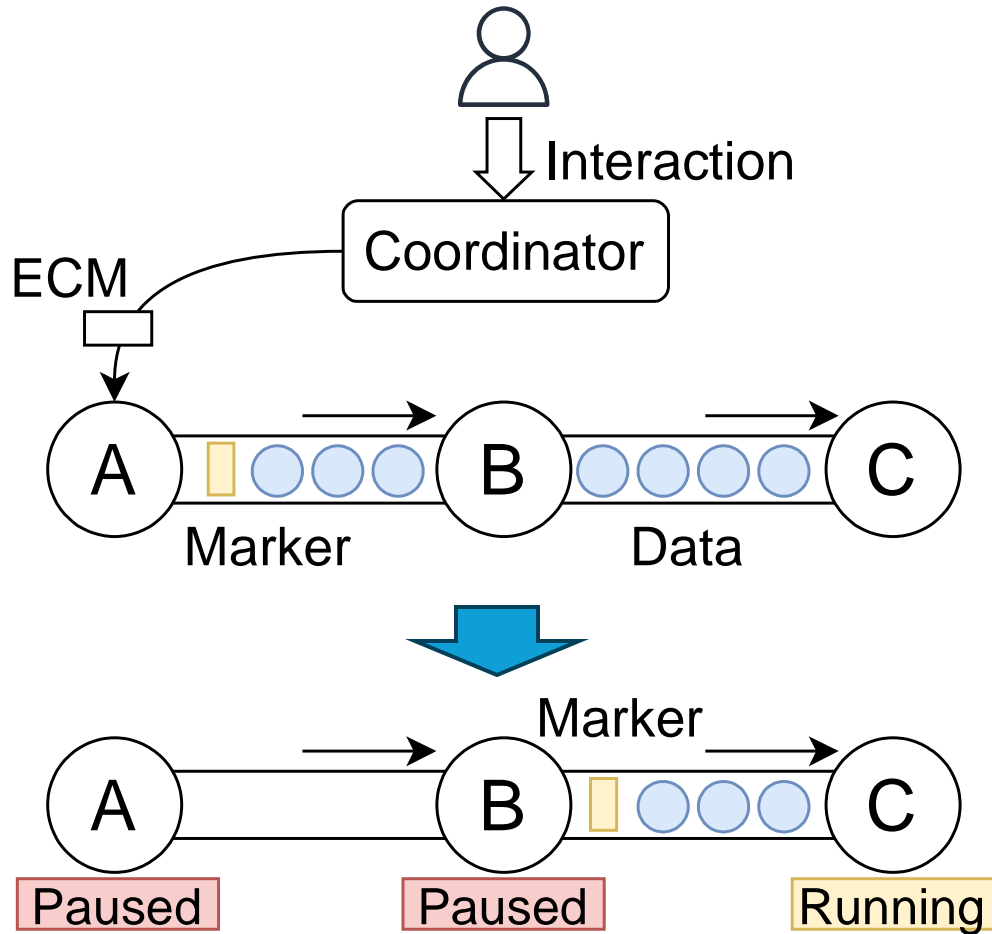# System Architecture

# How to Support Pause?

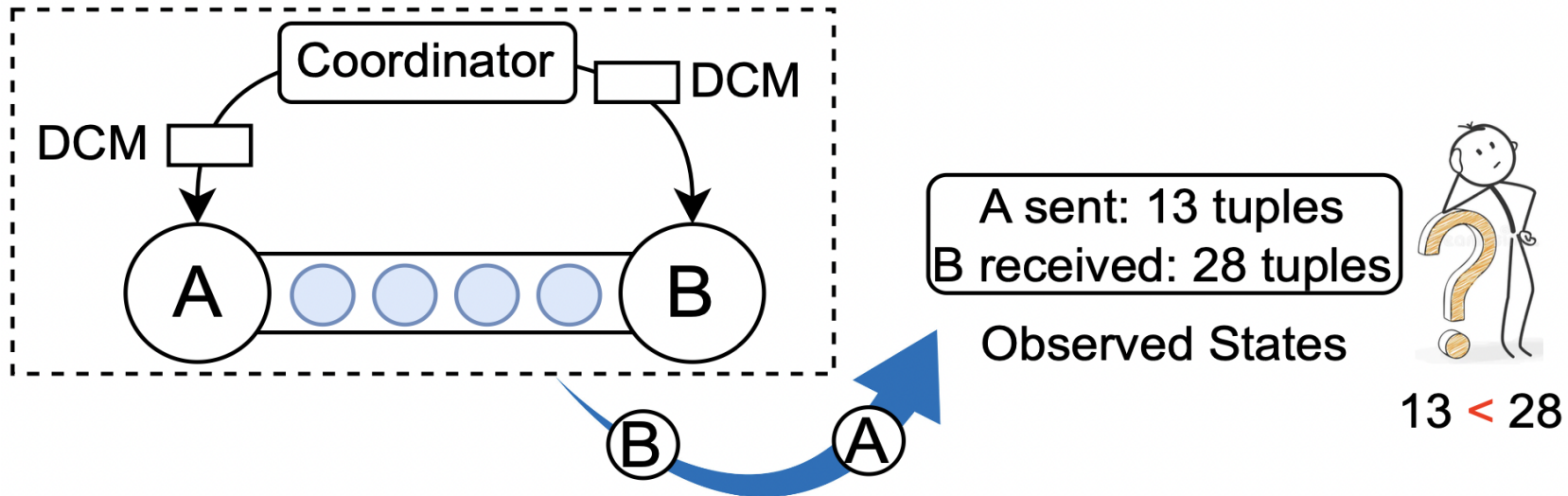Method 1: Pause a Workflow with **Direct** Control Message (DCM)
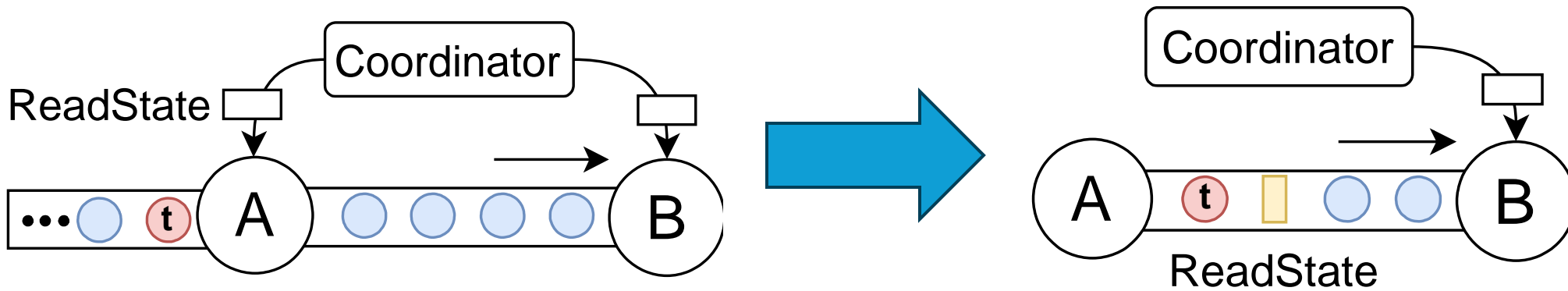
# How to Support Pause?

Method 2: Pause a Workflow with **Embedded** Control Message (ECM)
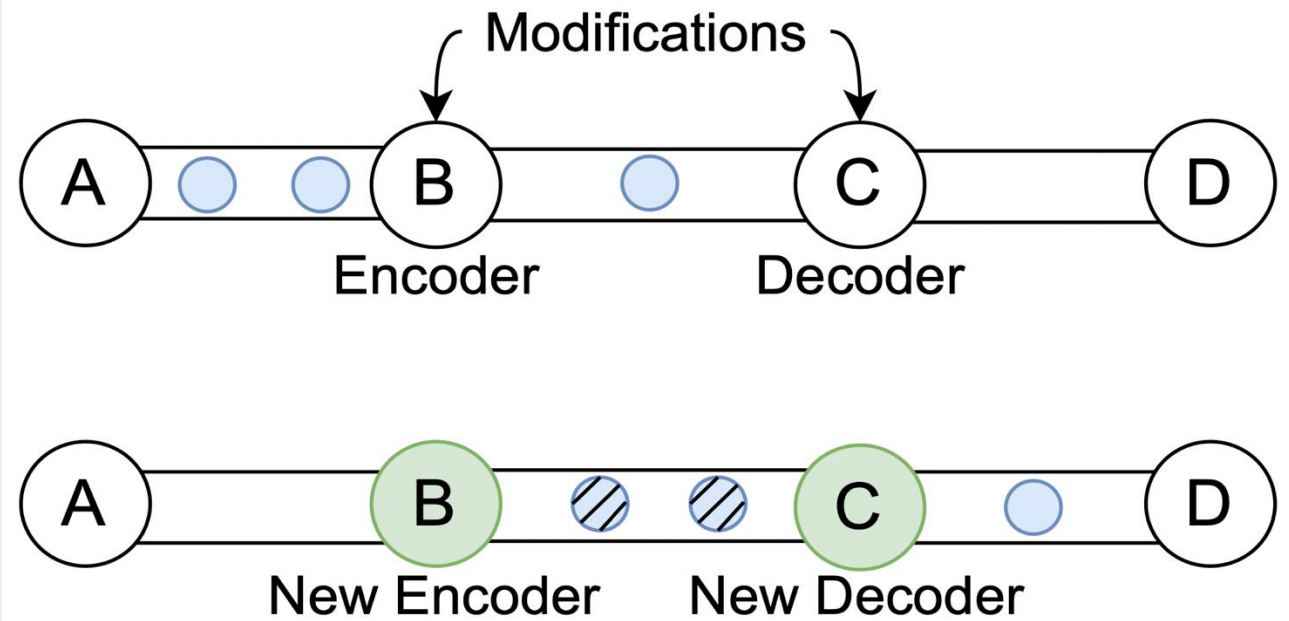
# How to Support Read States?



**Incorrect global state introduced by DCM**

# Combining ECMs and DCMs to Read States

# How to Support Modifications?

## Strict Consistency



Modifications

A — B (Encoder) — C (Decoder) — D

A — B (New Encoder) — C (New Decoder) — D

A tuple should be processed by the **same** version.

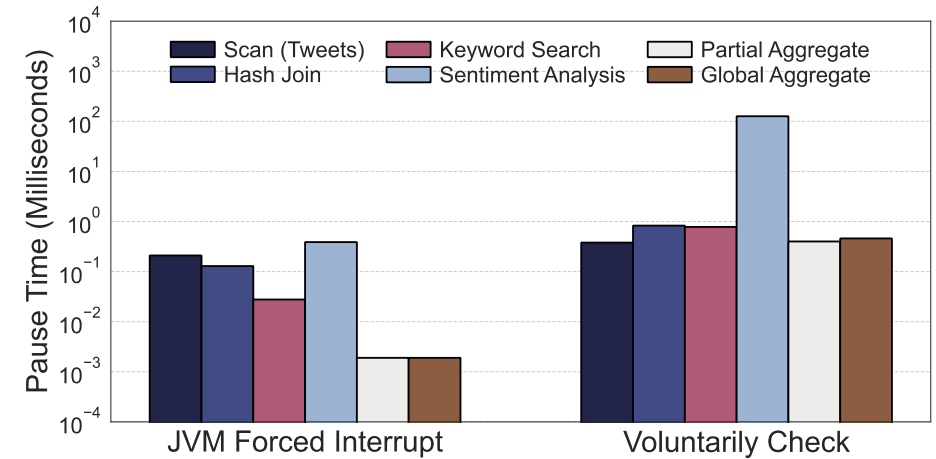# Experiments



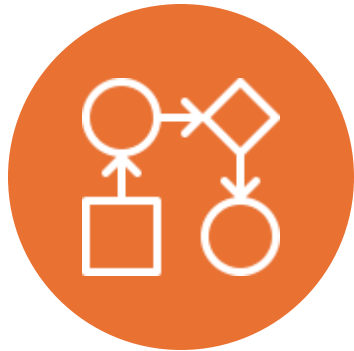- **Low Interaction Latency**
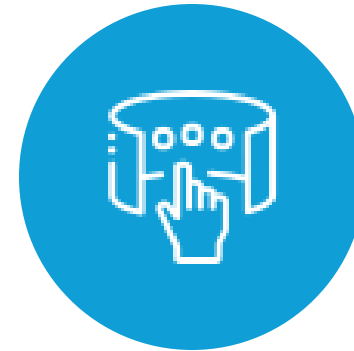
- **Low Runtime Overhead**

- **High Scalability**

# Summary of texera

WORKFLOW          COLLABORATION          INTERACTIONS

Texera Live Service

Texera GitHub Repo

Apache-2.0 License

# Open Source

# Project and Usage Metrics

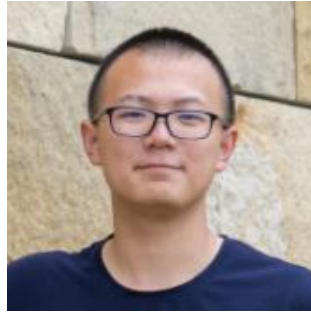| # of user accounts | 332 | # of projects | 86 |
|---|---|---|---|
| # of workflows | 2,481 | # of executions | 50,950 |
| # of workflow versions | 357,000 | # of publications | 23 |
| # of deployed servers | 7 | # of CPU cores in the largest deployment | 400 |
| # of files on GitHub | 1,291 | # of lines of code on GitHub | 101,690 |
| # of pull requests on GitHub | 2,252 | # of current PhD students | 7 |
| # of collaborating professors | 17 | # of involved undergraduates | 80+ |
| # of completed PhD theses | 3 | # of development years | 8 |

# **Texera:** A System for Collaborative and Interactive Data Analytics Using Workflows
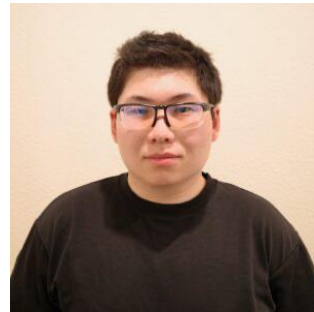


Zuozhi Wang

Yicong Huang

Shengquan Ni

Avinash Kumar

Sadeem Alsudais

Xiaozhen Liu

Xinyuan Lin

Yunyan Ding

Chen Li

Texera at https://texera.io

Texera Live Service

Texera GitHub Repo

University of California, Irvine