# ManaTEE: an Open-Source Private Data Analytics Framework with Confidential Computing

Dayeol Lee

dayeol.lee@tiktok.com

# Contents

- **Private Data Analytics**
  - Why is it Hard?
  - What do we Need?
- **Existing Approaches**
- **ManaTEE Project**
- **Tutorial Demo**
- **Project Current & Future**

# Private Data Analytics

# Why is Private Data So Important?

Value Extraction

Public Interest

# Why is Private Data So Important?

Value Extraction

Public Interest

# Sharing Private Data for Public Interest

- **Public Health**: medical data, personal health data
- **Public Safety**: PII (e.g., address, phone number) associated with crimes or illegal activities
- **Education**: academic performance, attendance, and engagement
- **Civic Engagement**: personal beliefs, social activities

… and many more

# Often Requires Cross-Organizational Data

Example: Understanding Illicit Drug Promotion by Using Cross-Platform Data
(Zha et al., CCS'24)

# HDR UK: Trusted Research Environment (TRE)



https://www.hdruk.ac.uk/access-to-health-data/trusted-research-environments/

# Why is it Hard?

# Challenge 1: Data Privacy Risks

- Trust
  - Conflict of interests
  - Risk of abusing data or data fabrication
  - Different trust domain (e.g., company, country, …)

- Compliance
  - Privacy policies such as data retention or purpose limitation
  - Providing raw data might be legally prohibited
  - Changing the geolocation of data can be legally restricted

# Challenge 2: Accountability and Transparency

In the Old Days…

Now it looks like

**Organization**

Data

Compute

**Organization**

Data

**Data Warehouse**

Data

Data

**Cloud Provider**

Compute

Data

# What do we Need?

# We Need A "Standard" Way That Provides...

- **Strong Privacy Protection Mechanisms**
  - Privacy Enhancing Technologies (PETs)
  - Protecting privacy while maximizing the utility of data
- **Technical Enforcement of Policies**
  - Terms and conditions and honor codes are not enough
  - Proactive measures, instead of reactive
- **Accountability and Transparency**
  - Auditability of the full system if necessary
  - Verifiability on the integrity of the results
- **Usability**
  - Provide accurate results
  - Must be easy to deploy, easy to use, and easy to customize

# Existing Approaches

# Existing Industry Solutions for Security & Privacy

**SQL Policy-based Data Clean Room**

Rely on SQL/data platforms provided by 3rd party, who is free of conflict of interest

**Differential Privacy**

Preprocess data or add noise to the aggregated SQL results to limit information leakage

**Trusted Execution Environment**

Use remote attestation to co-verify the code before releasing data; contain data in an isolated environment during execution

# Technical Difficulties of Existing Solutions

| | PET | Technical Enforce | Transparency | Usability | Accuracy |
|---|---|---|---|---|---|
| SQL Policy-based Data Clean Room | No | No | No | Good | Yes |
| Differential Privacy | Yes | Yes | No | Good | No |
| Trusted Execution Environment | Yes | Yes | Yes | Could Be Better | Yes |

# Our Goals

- **Technical Enforcement via PET**: enforce privacy policies such as purpose limitation and data retention via PETs
- **Usability**: provide an interactive tool to utilize the data
- **Accuracy**: provide accurate results on real data, as well as an evidence of execution
- **Transparency and Accountability**: make it auditable and verifiable
- **Deployment**: make it easy to deploy to the cloud
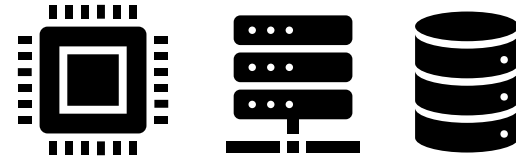
# Different Need at Each Stage

## Programming Stage

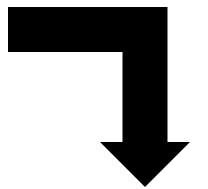| |
|---|
| Smaller Data/Compute |
| Interactive |
| Hard to Control Data |
| Higher Privacy Risk |

## Execution Stage
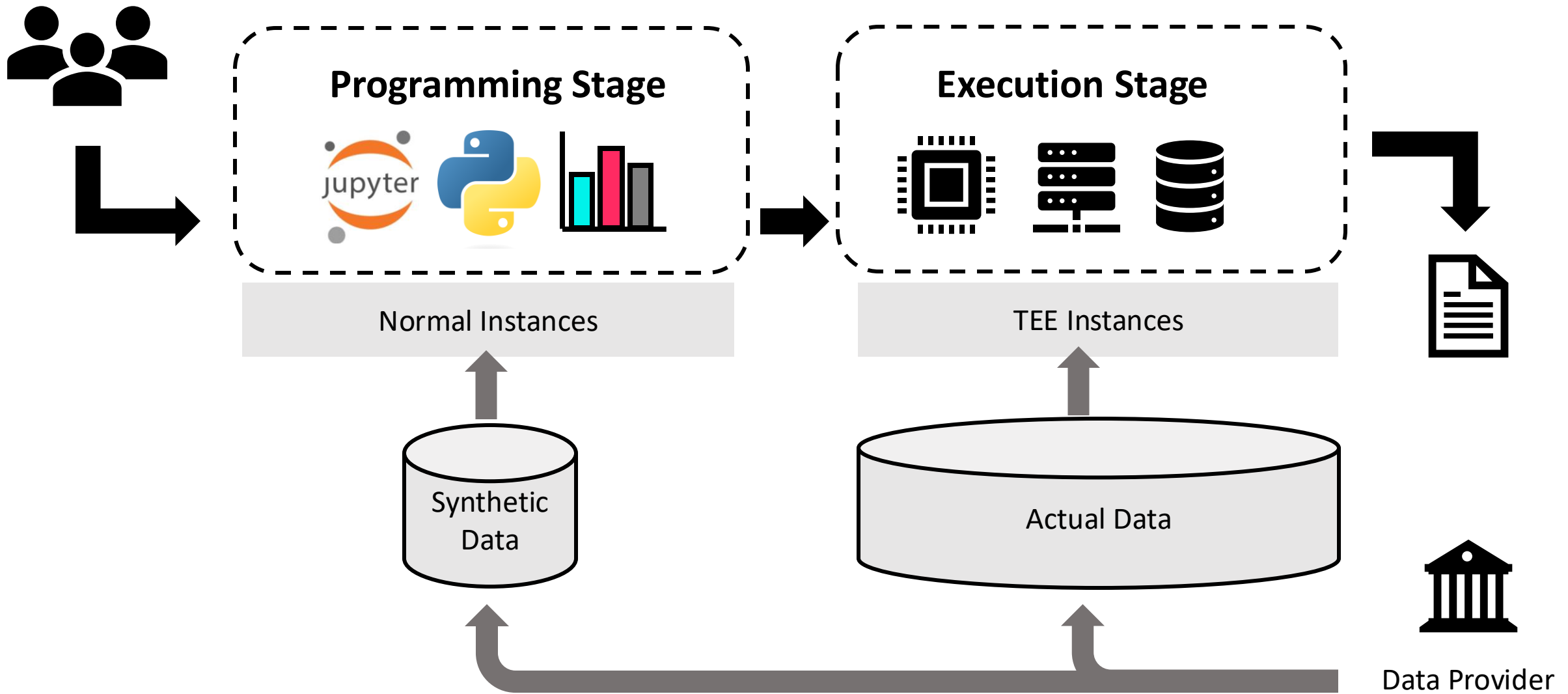
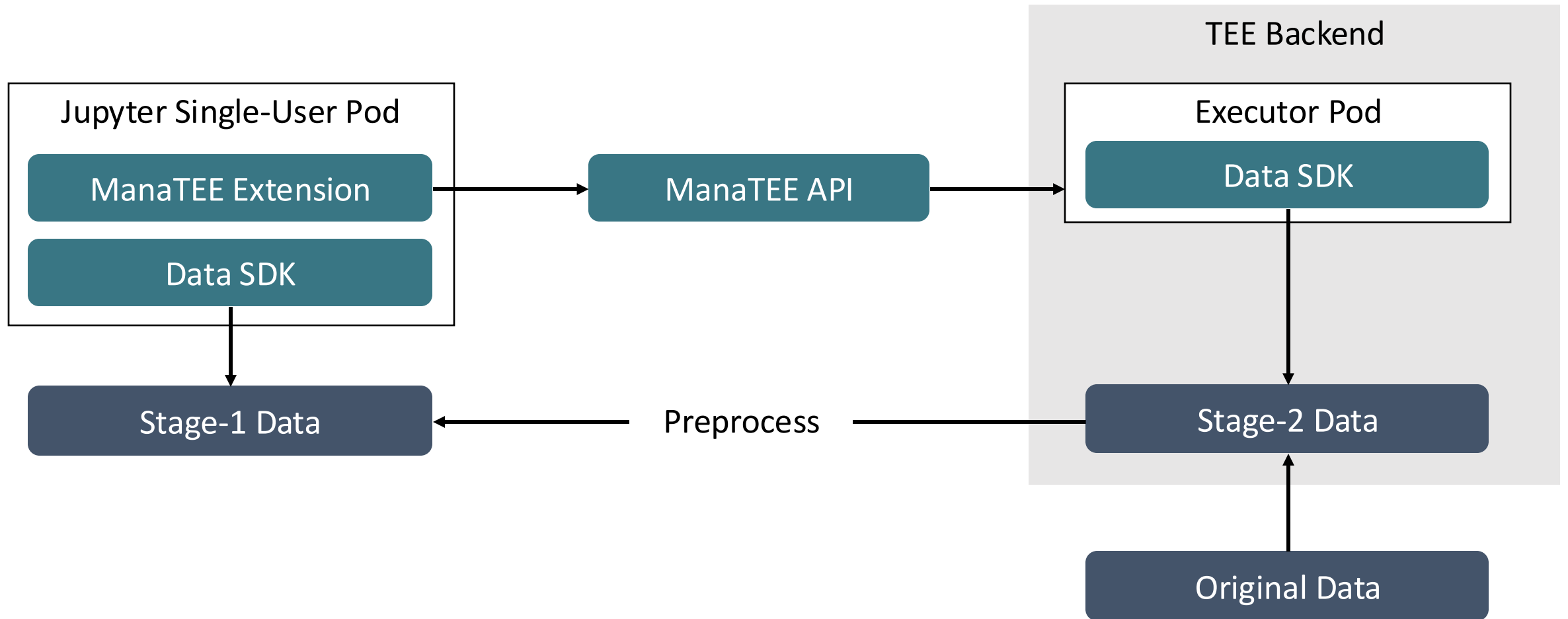| |
|---|
| Larger Data/Compute |
| One-Time Execution |
| Easier to Control Data |
| Lower Privacy Risk |

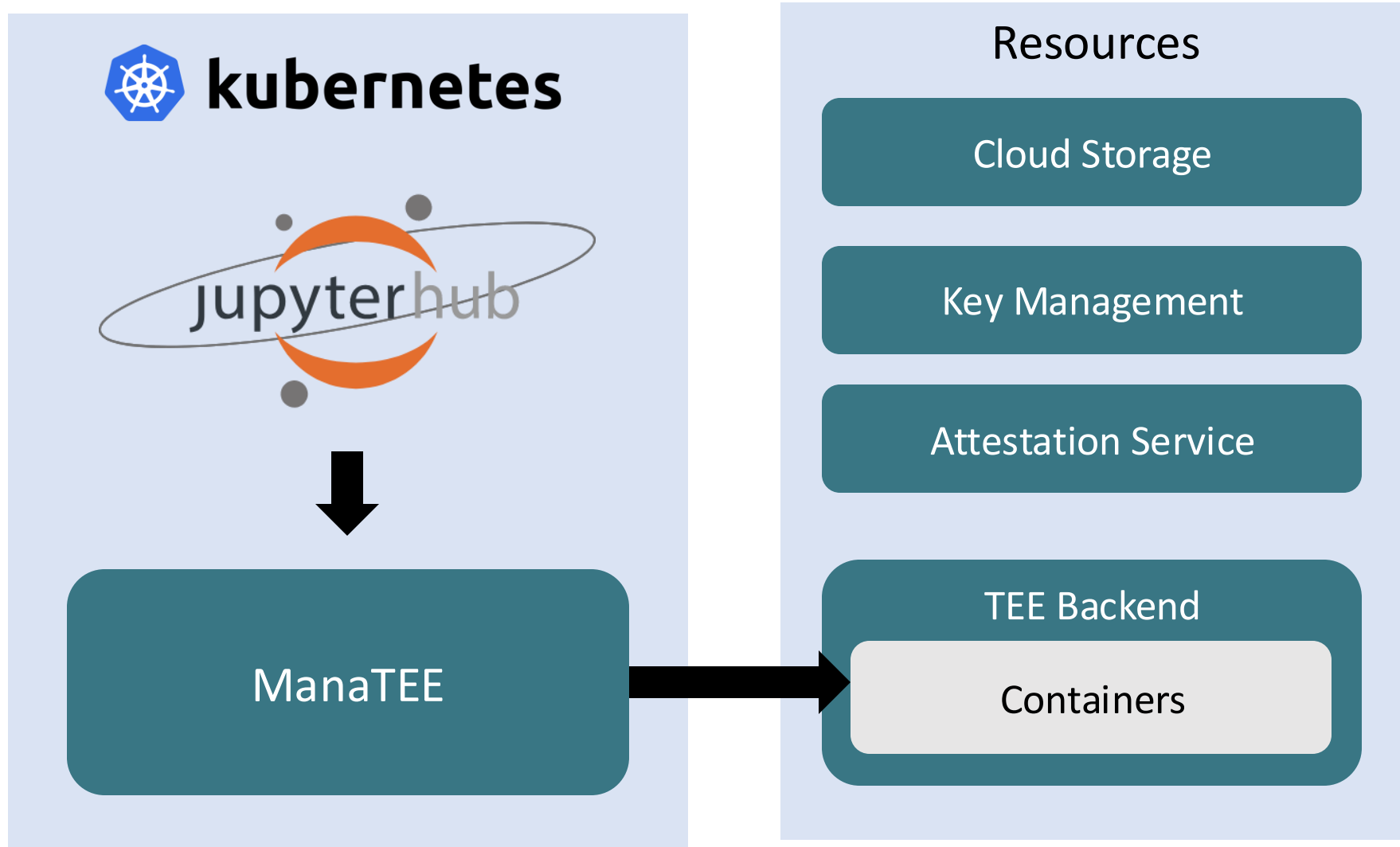# Our Approach: Two-Stage Data Clean Room

# Benefits of the Two-Stage Execution Model

- Separate **Data Policy** and **Code Policy**
  - Flexible data policy on programming stage – LDP perturbation, sampled data, or DP synthetic data
  - Code policy enforced only at the execution stage
- **Accurate** Results in Execution Stage
  - Full data access is securely enabled via confidential computing
- Why **Confidential Computing**?
  - Provides transition of trust, making it work with various trust model (Cross-organizational data providers)
  - Integrity of the execution
  - Proof of execution (attestation report)
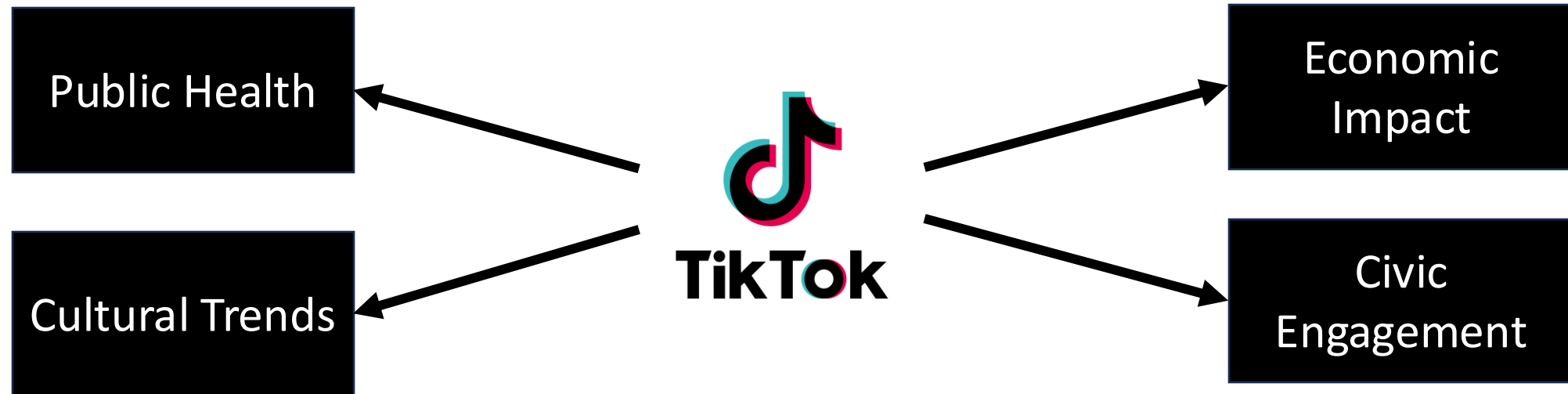
# ManaTEE Data and Code Pipeline
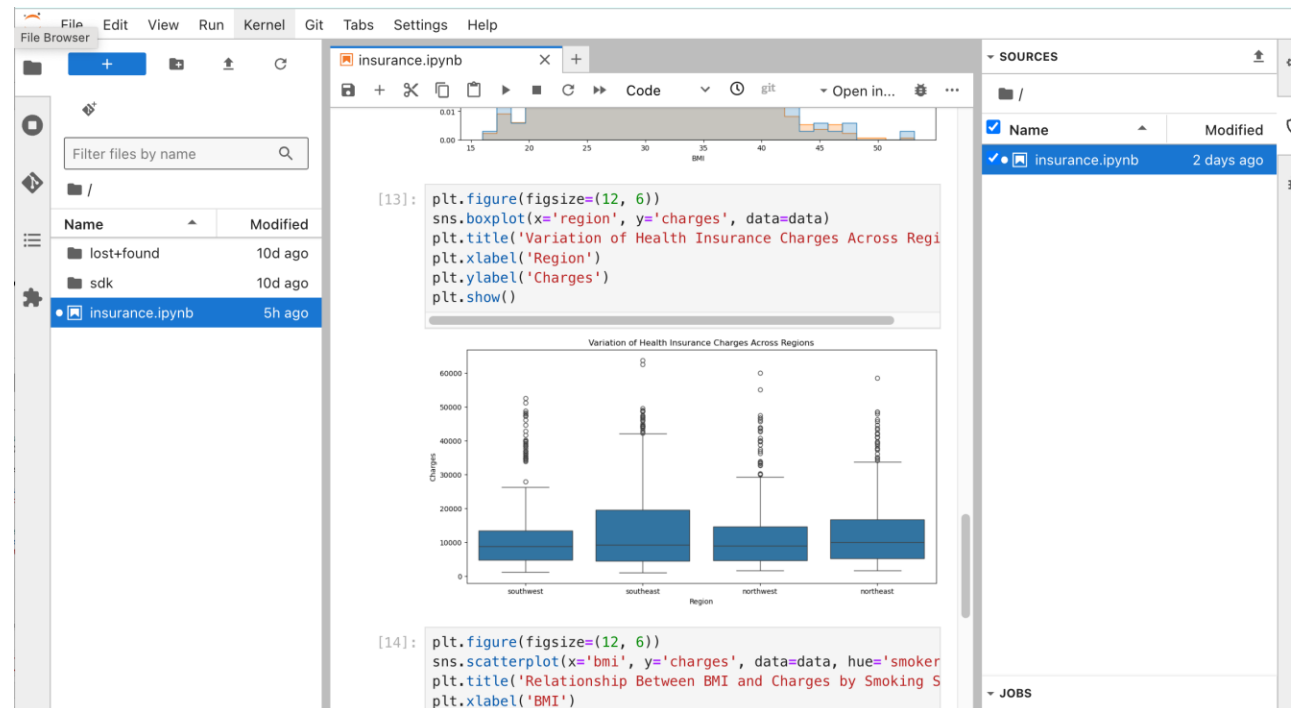
# Easy Cloud Deployment via Terraform

# Use Cases

# Providing Transparency to Researchers

# TikTok Research Tools

- Virtual Compute Environment (VCE)
  https://developers.tiktok.com/doc/vce-getting-started

# Another Potential Use Cases

- Ads & Marketing
  - Lookalike segment analysis
  - Measurement and conversion tracking
- Machine Learning
  - Inferencing & training with private dataset
  - Inferencing & fine-tuning private model
  - AI model evaluation (e.g., fairness) on private models

# Tutorial Demo

https://manatee-project.github.io/manatee/getting-started/tutorials/

# Demo Scenario

- Dataset
  - Insurance charge dataset from Kaggle
- Task
  - Train a model predicting the insurance charge
  - XGBoost Regression
- Privacy Protection
  - Differentially-private synthetic data in the first stage
  - MST (2018 NIST synthetic data challenge winner), McKenna et al.

# Launcher +

## Notebook


Python 3
(ipykernel)

## Console


Python 3
(ipykernel)

## Other


Text File


Markdown File


Python File


Show
Contextual Help

ManaTEE: an Open-Source Private Data Analytics Framework

# Project Current & Future

| | Current | Future |
|---|---|---|
| **Users** | One-Way Collaboration (Singler source of data) | Multi-Way Collaboration (Cross-organizational data) |
| **Backend** | Single Backend | Multiple Backend |
| **Data Provisioning** | Manual | Automated |
| **Policy and Attestation** | Manual | Automated |
| **Compute** | CPU | CPU/GPU |

# Project Timeline

- [2024/5] TikTok launched VCE
- [2024/6] TikTok Open sourced PrivacyGo Data Clean Room
- [2024/10] Project renamed to ManaTEE and donated to Confidential Computing Consortium
- [2025/1] ManaTEE community version released
- [Current] Forming Technical Steering Committee

# The First Community Release for Open Collaboration

- Tutorial
  - https://manatee-project.github.io/manatee/getting-started/tutorials/
- Local Deployment (Minikube)
  - https://manatee-project.github.io/manatee/getting-started/minikube/
- Announcement
  - https://manatee-project.github.io/manatee/blog/2025/01/07/first-community-release-of-manatee
- Release Note
  - https://github.com/manatee-project/manatee/releases/tag/0.1.0

# Collaborators

**Please Join Us!**

Google Groups: https://groups.google.com/u/1/g/manatee-project
Github: https://github.com/manatee-project

# Q&A

dayeol.lee@tiktok.com