



CSSI Element: Empowering Data-driven Discovery with a Provenance Collection, Management, and Analysis Software Infrastructure

PI: Yong Chen, Co-PIs: Dong Dai, Brian Ancell, William Hase

Institutions: Texas Tech University

Objective: To design and develop a provenance collection, management, and analysis software infrastructure for high performance computing (HPC) systems and to train students for the future workforce.

What is provenance?

- Entities, such as:
 - Users, machines/hosts,
 - Jobs, programs/processes/threads,
 - Files, etc., and
- Relationships among them

What are use cases?

- Describing the history of a piece of data
- Identifying sources, parameters, or assumptions behind a given result, how an input transformed to the output
- Auditing data and usage
- Monitoring the system & performing predictive analysis

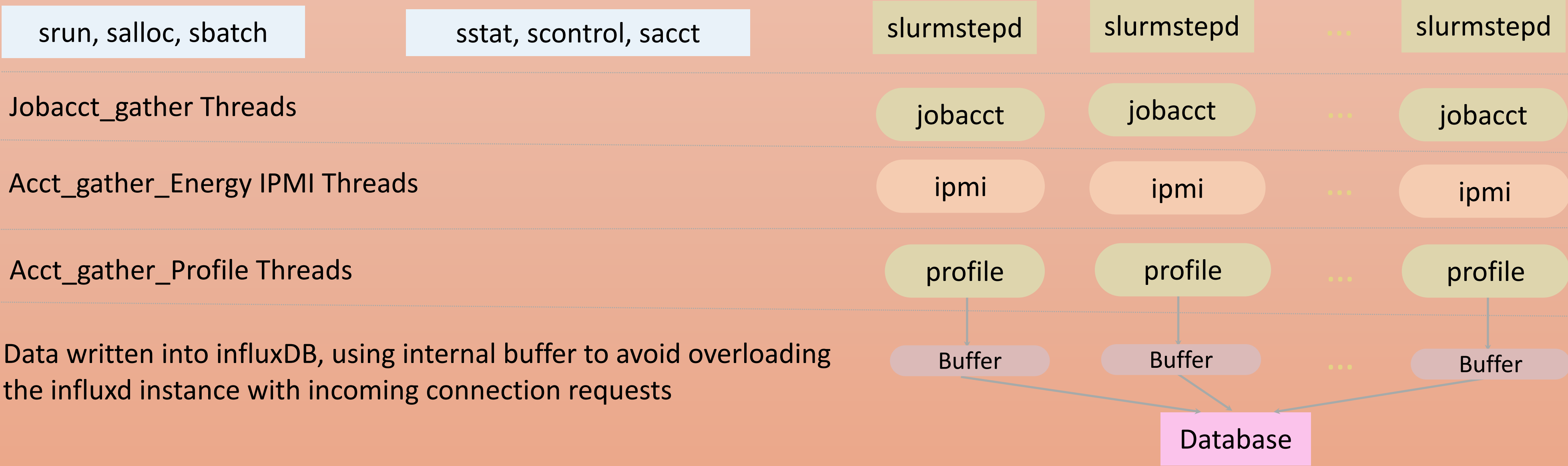
Who are the users?

- Users primarily include:
- HPC system administrators
 - Domain scientists
 - Data scientists such as these who provide visual and predictive analysis

What is the current progress?

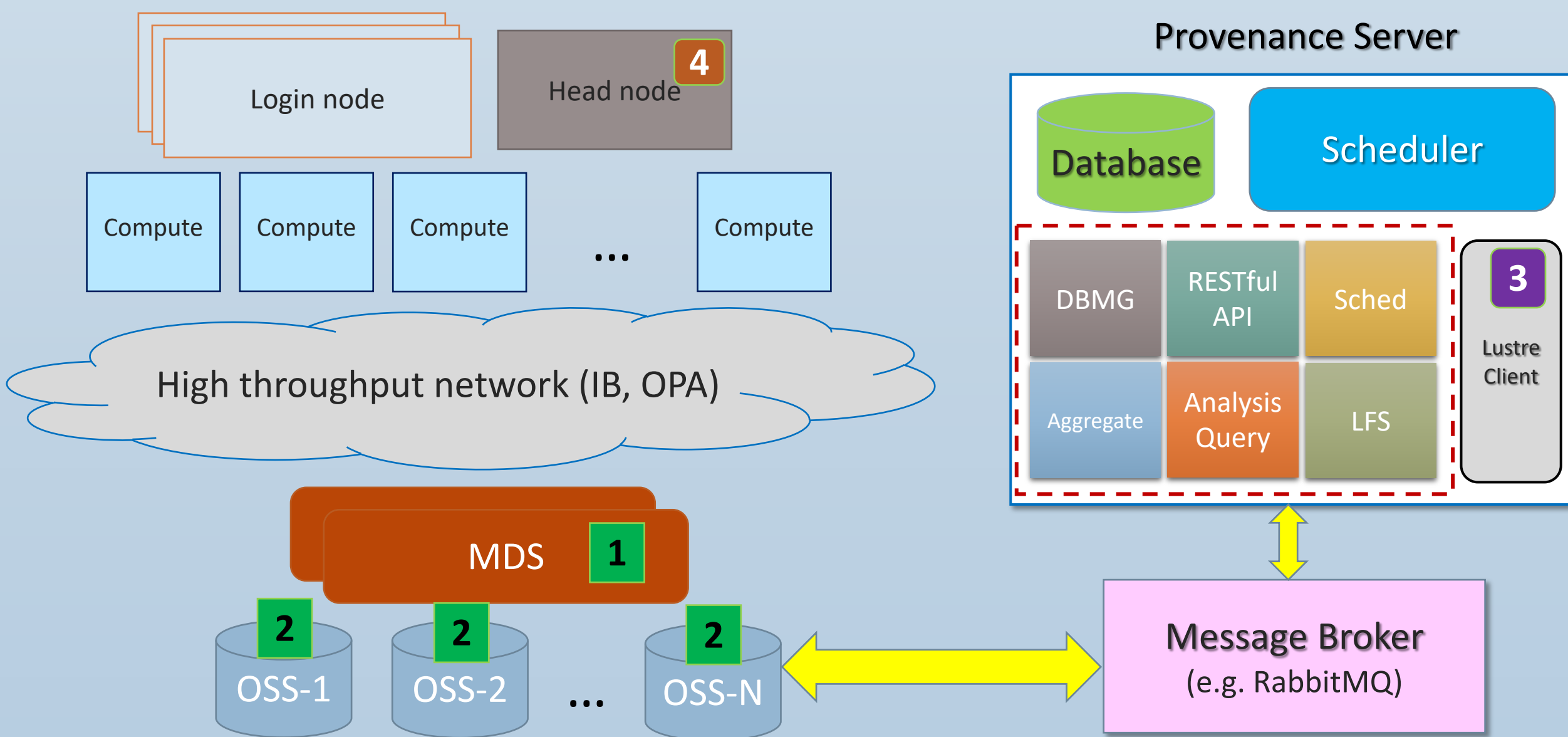
- Developed a component to collect users/jobs from UGE job scheduler (working on Slurm scheduler)
- Developed a component to leverage jobstats and changelogs to gather file system provenance
- Used on a 467-node Quanah cluster for preliminary visual analytics

Collect users and jobs information from Slurm job scheduler



File system provenance

- RabbitMQ as the message broker
- Pika library in Python to communicate with RabbitMQ
- Compatible with all Linux-based cluster
- Expects to be scalable for:
 - 100+ MDSs & OSSs
 - 1,000+ jobs
 - 10,000+ nodes



File system provenance

1. Lustre JobStats on MDS servers
 2. Lustre JobStats on OSS servers
 3. Lustre ChangeLogs on a Lustre client
 4. Job scheduler data (reporting/accounting)
- Collecting JobStats data from OSSs and MDSs filtered by HPC job scheduler IDs
 - Collecting file operations data (Lustre Changelog) from a Lustre client
 - Collecting job scheduler information
 - Aggregate all data into a usable format
 - Store generated provenance into database
 - Providing a RESTful API and data analysis tools

