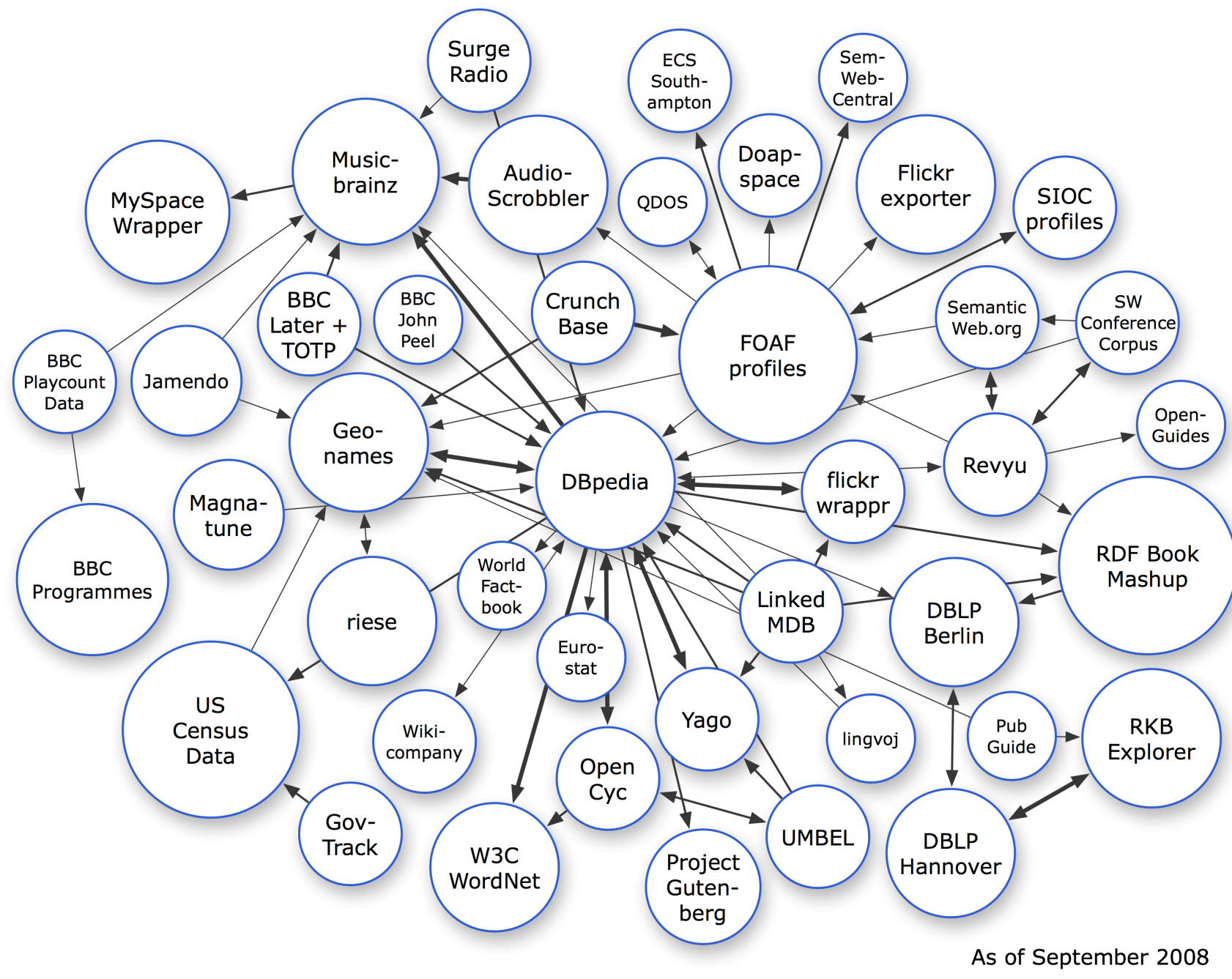
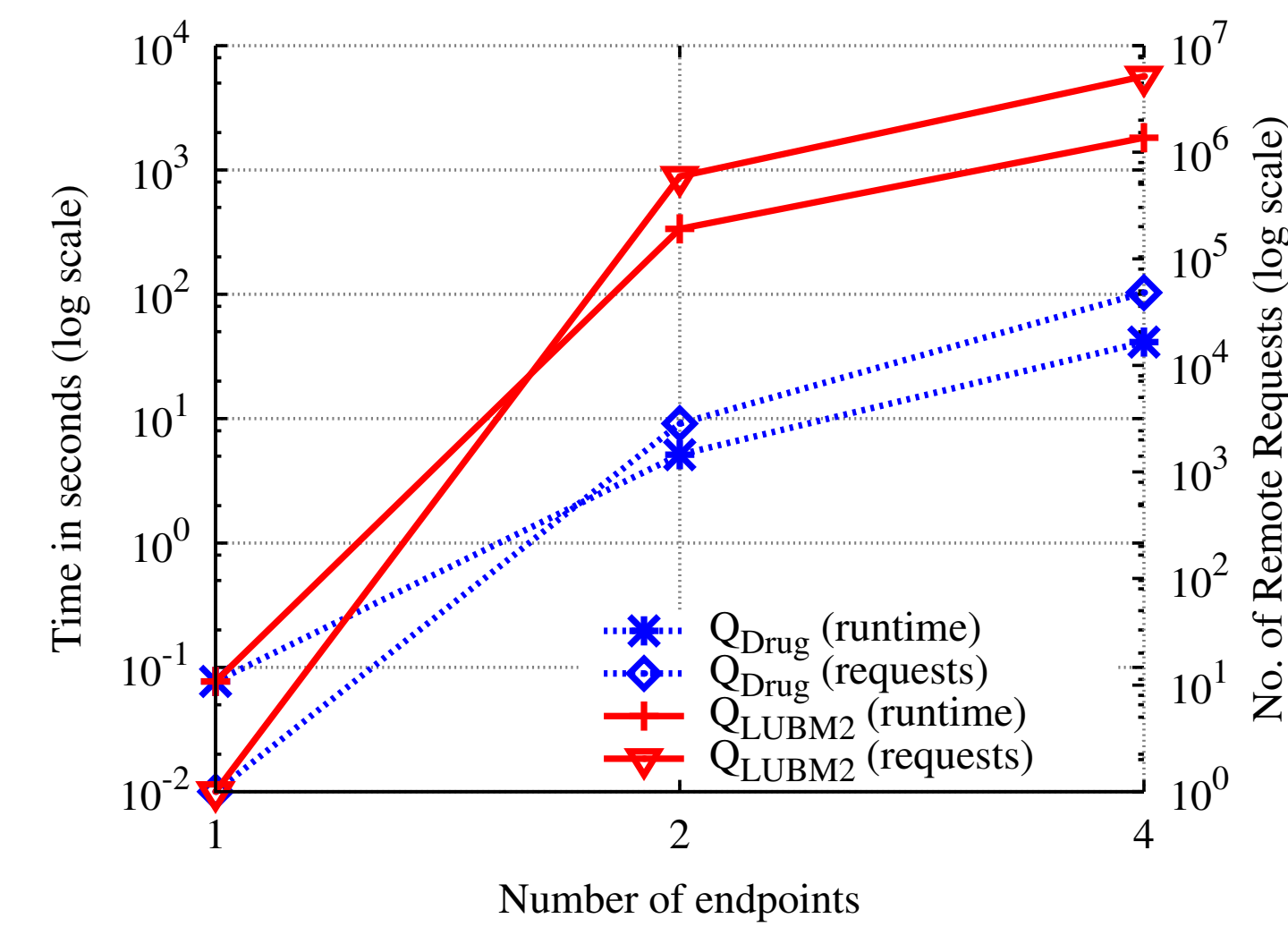


## MOTIVATION

- Linked Open Data (LOD) contains >1,000 datasets with 130 billion triples and more than 500M interlinks.
- Emerging Applications in life sciences, decentralized social networks, and Internet of Things need to integrate data from multiple RDF datasets via SPARQL queries.
- Existing RDF systems have been proposed to support a small number of data sources by utilizing schema information.
- Objective: decompose queries to avoid unnecessary data communication and inefficient computation.



As of September 2008



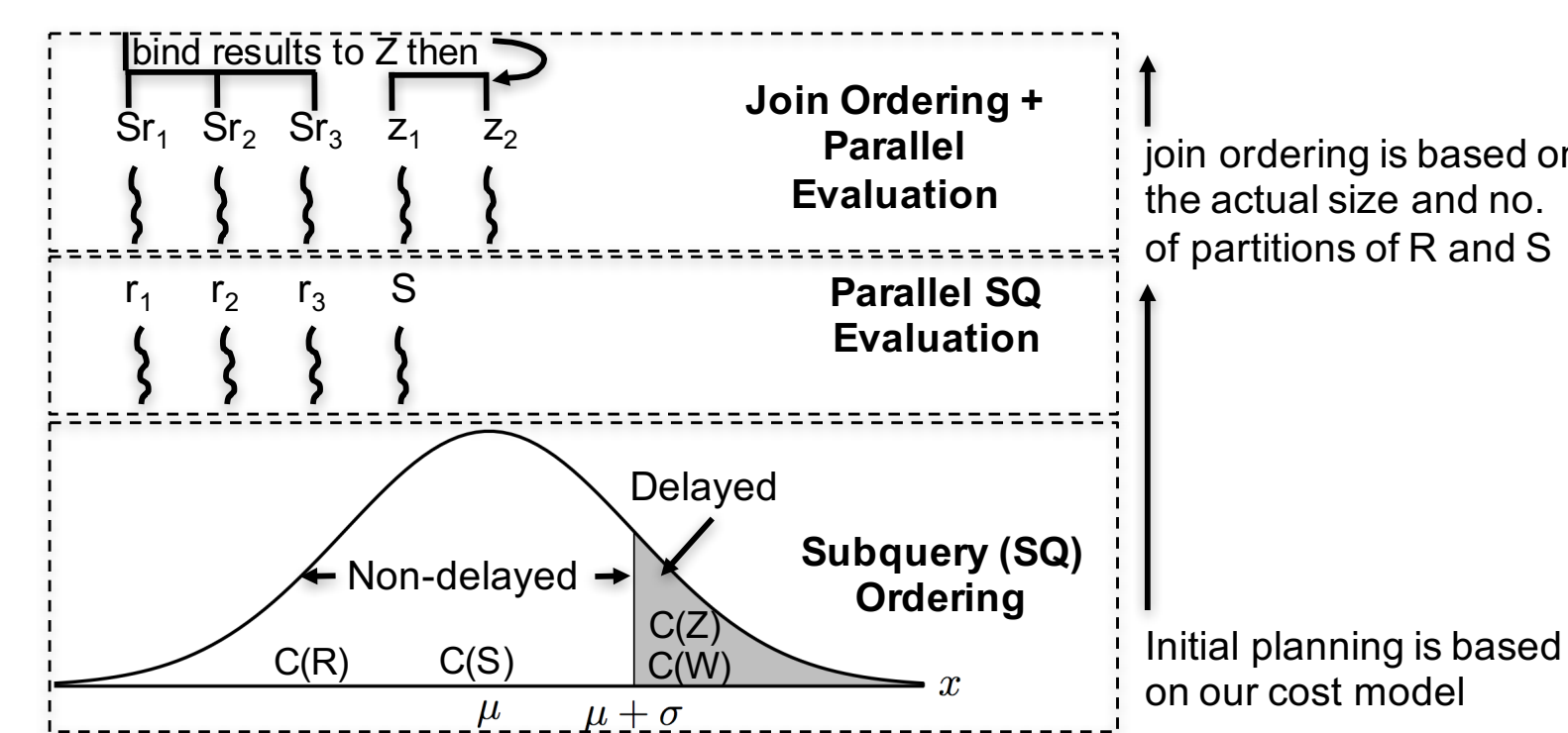
Scalability limitations of existing systems for querying linked data using FedX, the state-of-the-art system.

## LOCALITY AWARE DECOMPOSITION

- Maximizes the computations at the endpoints and minimize intermediate results.
- Utilizes the knowledge of the locations of the actual RDF triple instances matching a query variable.
- Determines triple patterns that can be sent together to an endpoint.
- Decomposes the input query into a set of independent subqueries.

## SELECTIVITY-AWARE EXECUTION

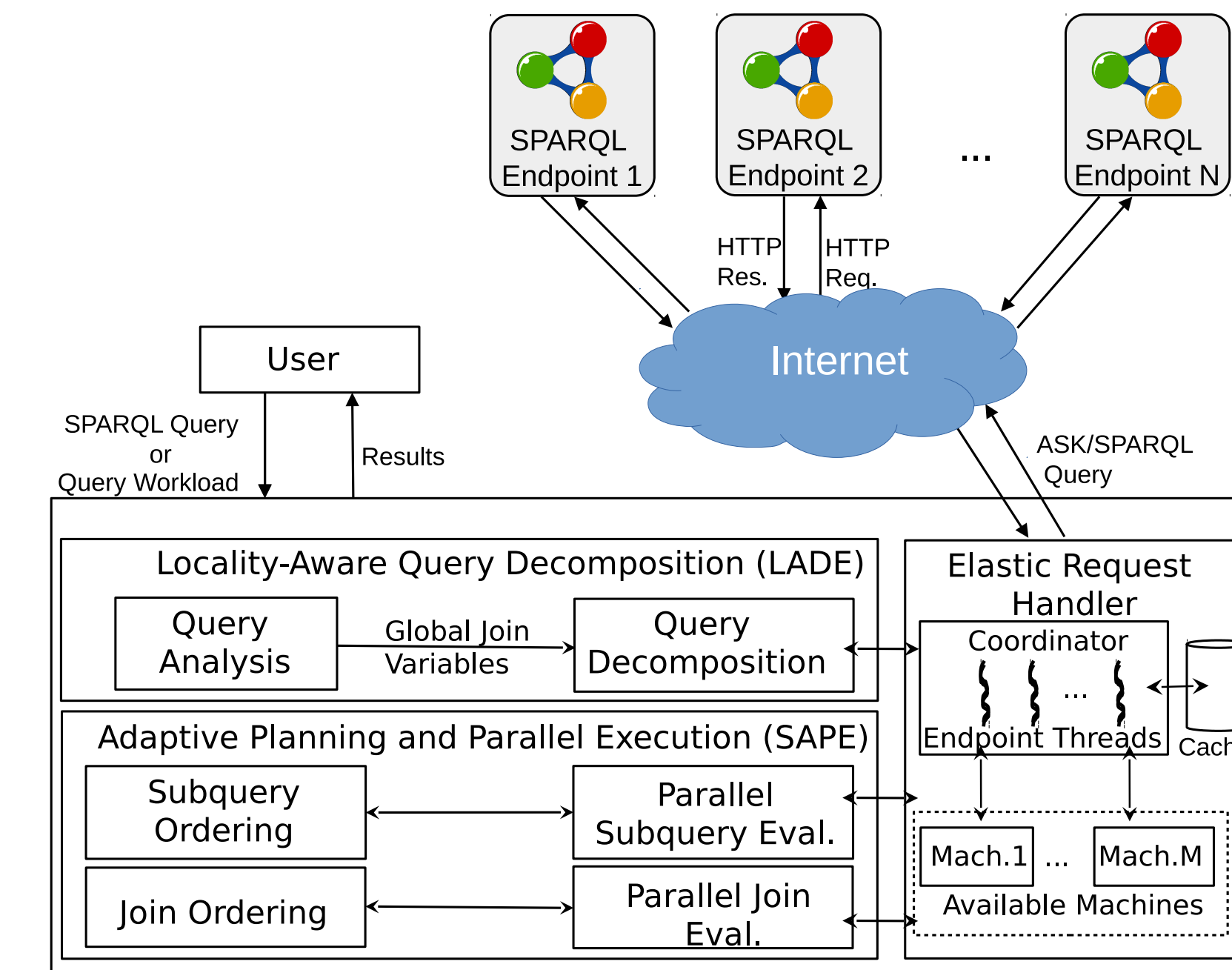
- Cost model delays subqueries expected to return large results.
- Achieve a high degree of parallelism while minimizing the communication cost.



## REFERENCES

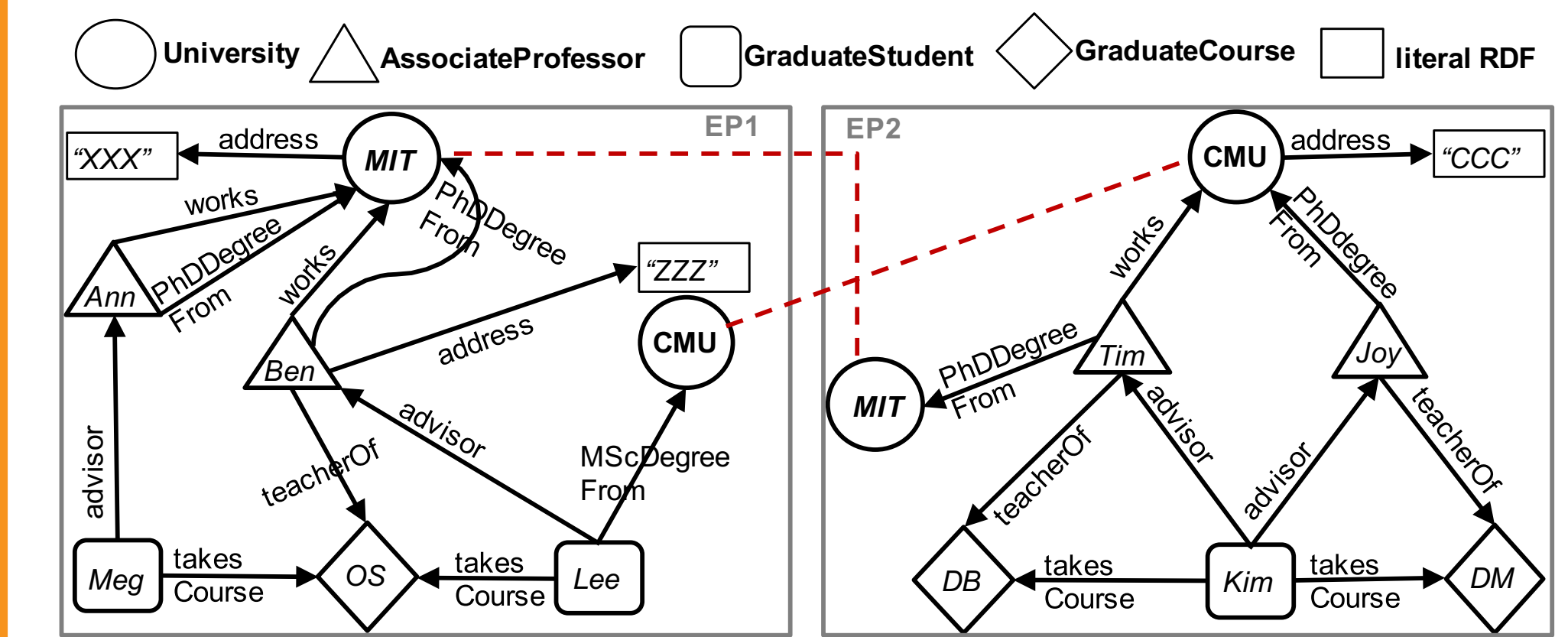
- I. Abdelaziz, E. Mansour, M. Ouzzani, A. Aboulnaga, P. Kalnis, "Lusail: A System for Querying Linked Data at Scale", PVLDB 11(4), 2018
- E. Mansour, I. Abdelaziz, M. Ouzzani, A. Aboulnaga, P. Kalnis, "A Demonstration of Lusail – Querying Linked Data at Scale", SIGMOD 2017 (Demo).

## LUSAIL



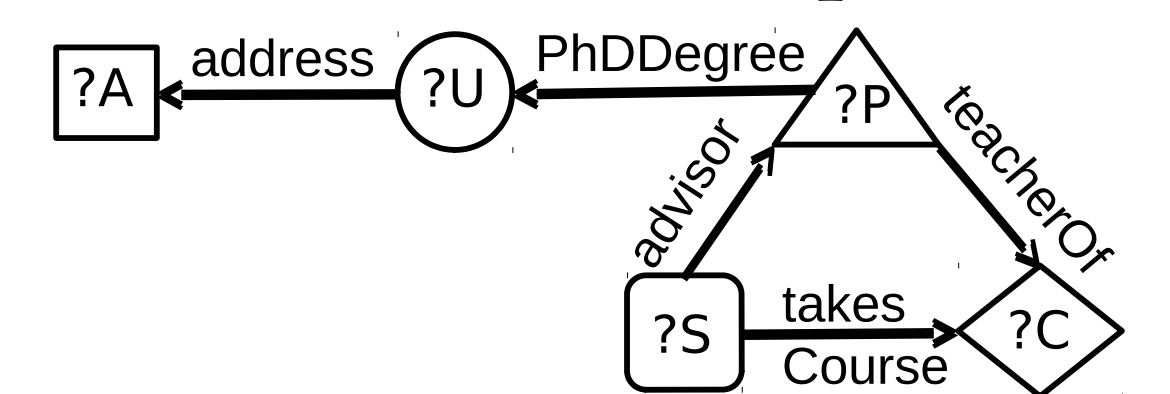
Lusail: a scalable and efficient federated RDF engine for querying linked data at scale.

## EXAMPLE DATA AND QUERY

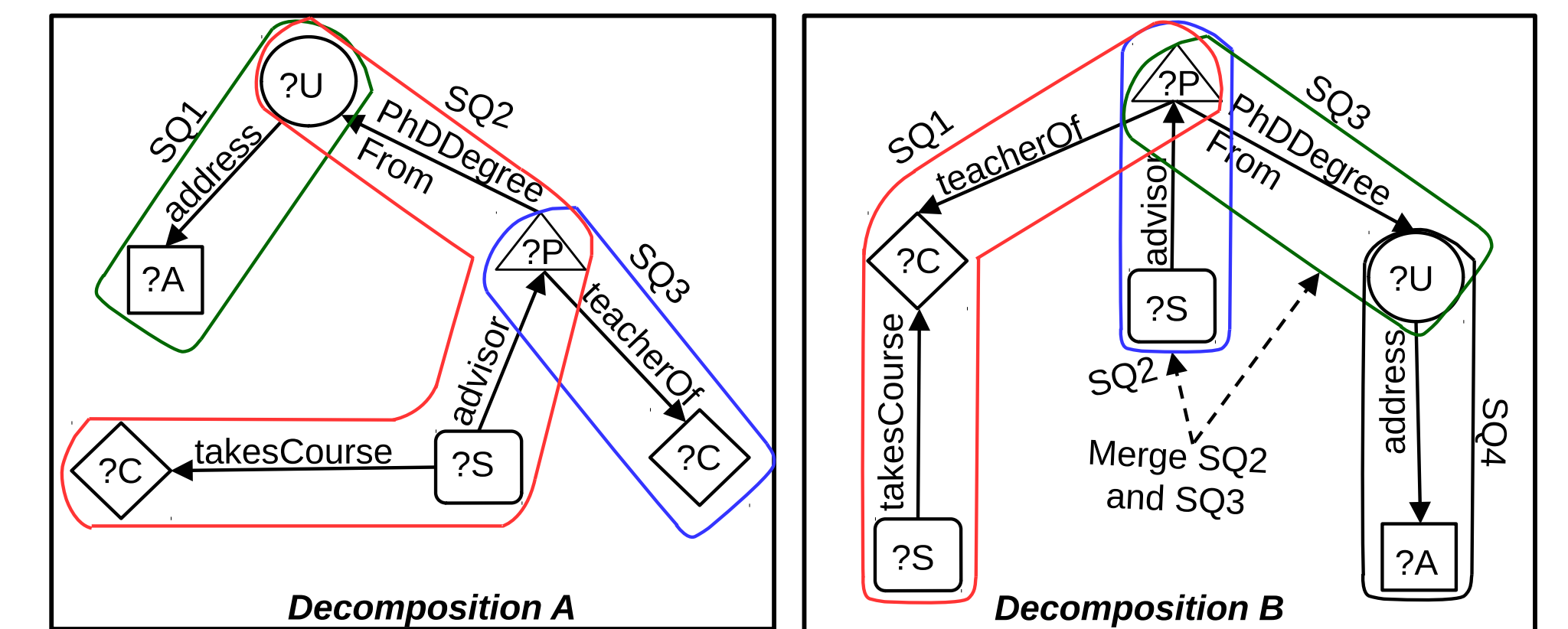
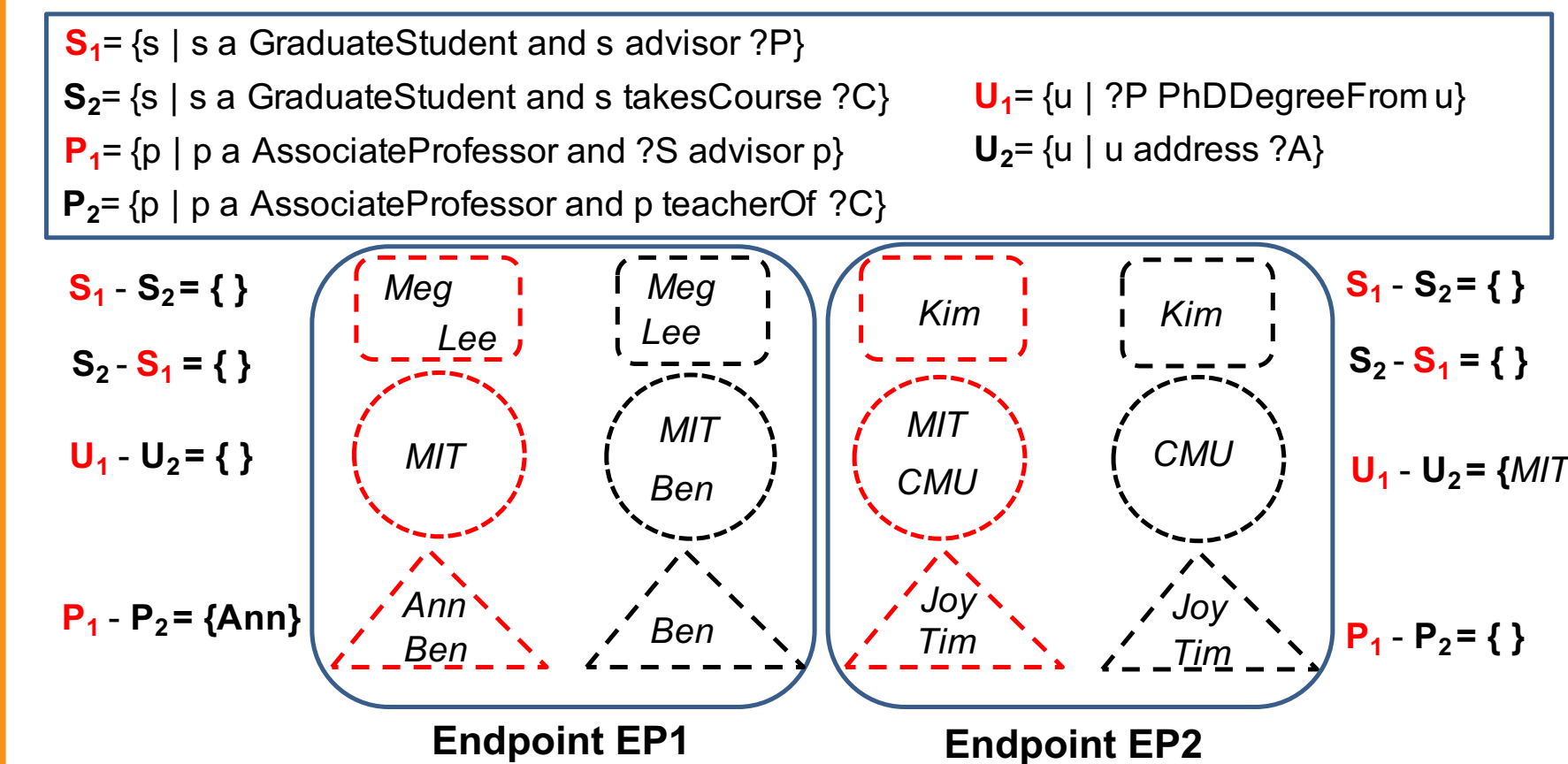


Decentralized graphs for two universities managed by geo-distributed SPARQL endpoints.

Example Query Q



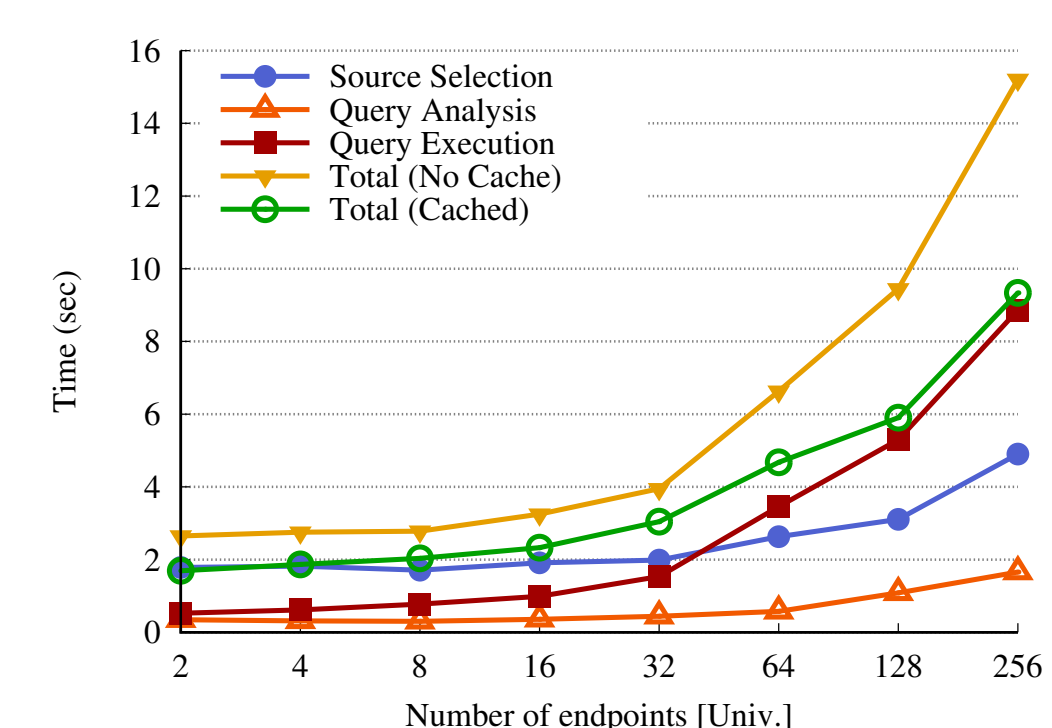
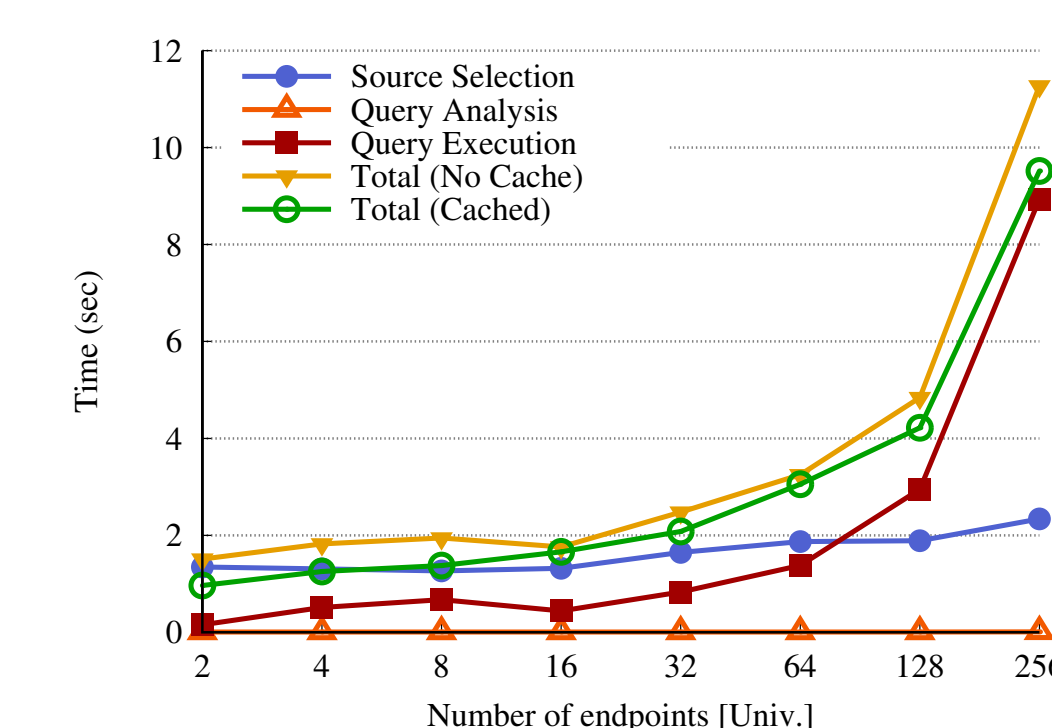
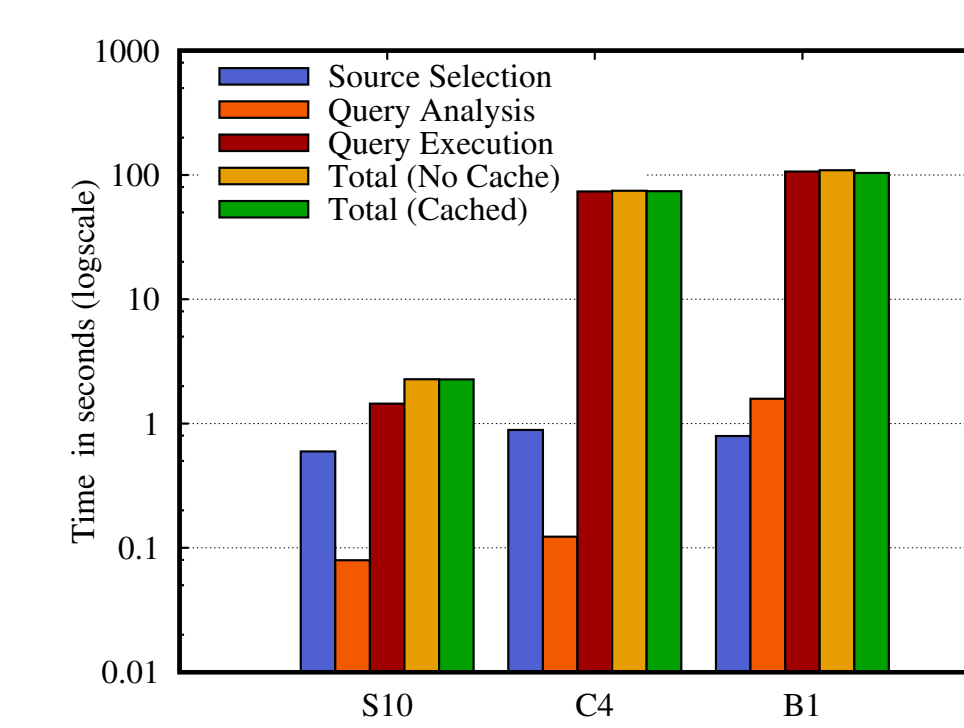
## RUNNING EXAMPLE



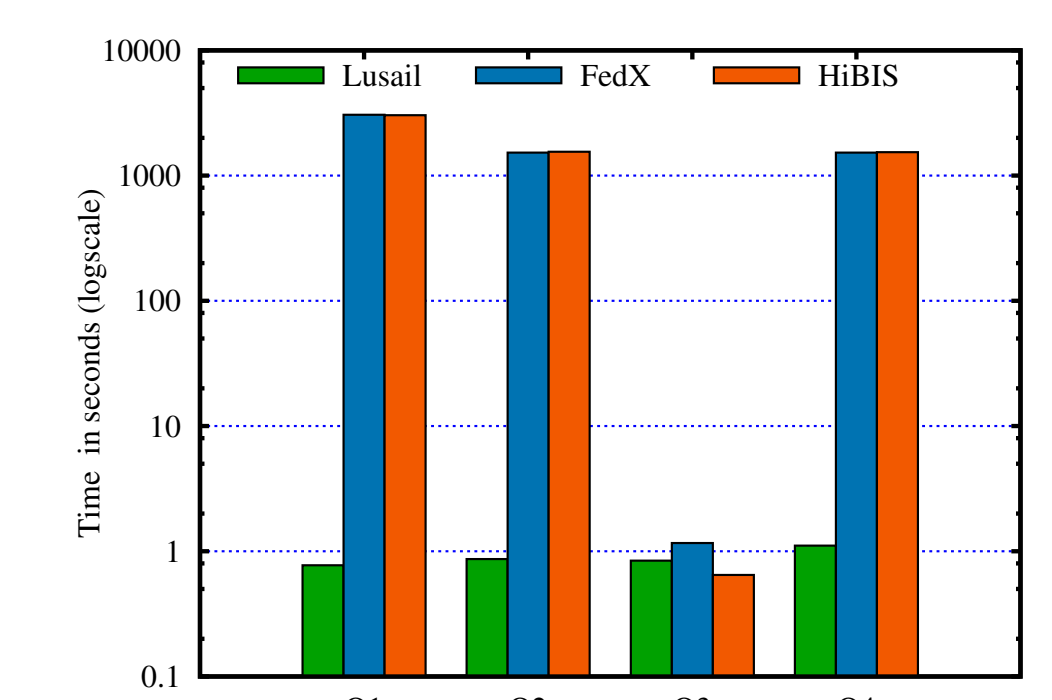
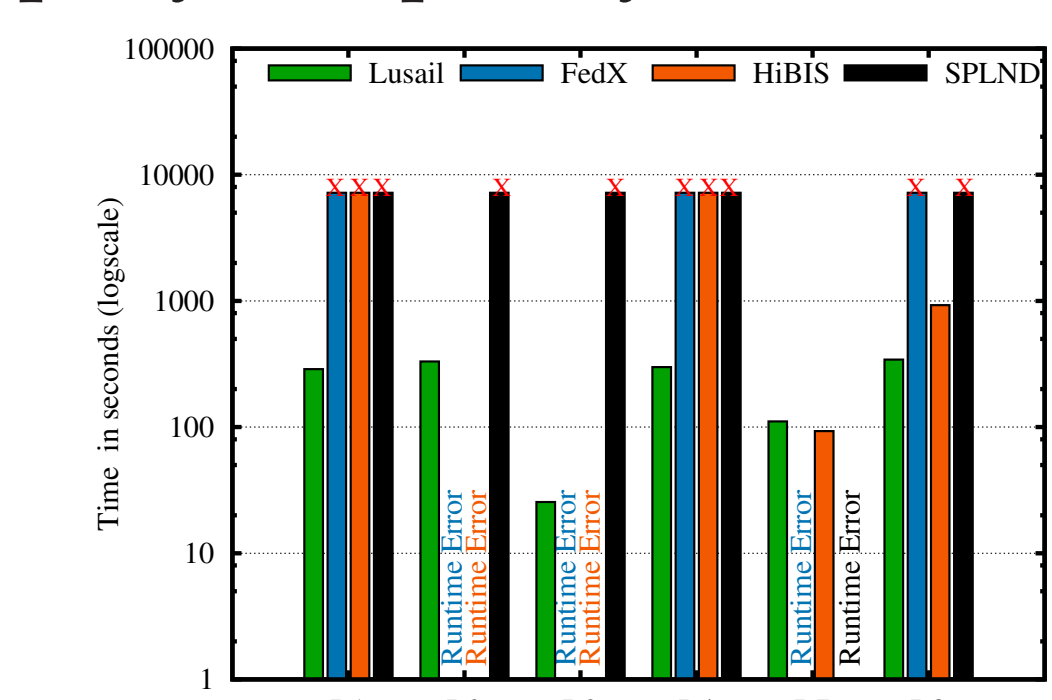
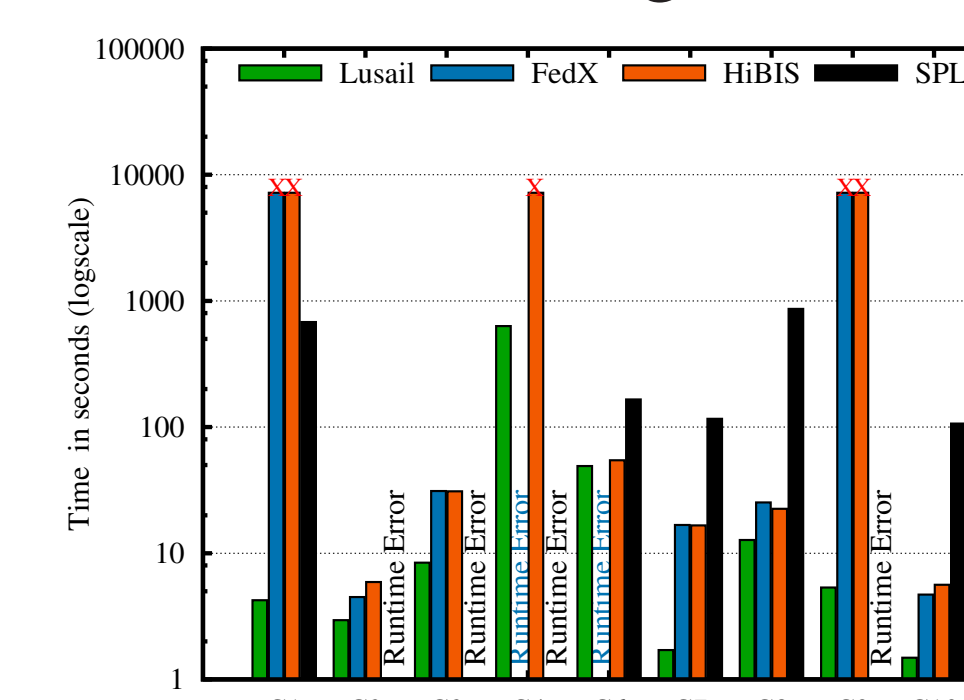
Locality analysis of data instances matching ?S, ?U, and ?P in in Q

Two possible decompositions for Q using detected global join variables ?U and ?P

## EVALUATION



Profiling Lusail by varying the query complexity, the number of endpoints, and the data size.



Geo-distributed federation: endpoints deployed in 7 different regions of the Azure cloud.