

# Performance Evaluation of Apache Tajo

Jihoon Son / Gruter Inc.



# Goals

---



- Performance comparison with other systems
- Scalability test of Tajo

# Evaluation on Cloud Environment



- Google Cloud Platform
  - Instance type: n1-standard-8
    - 8 core, 30GB RAM

# TPC-DS



- Data
  - 24 tables
    - Plain text format
    - Stored on Google Cloud Storage
- Query
  - Which can be executed on every system without modifications
    - For Hive, 0.12 doesn't support implicit join, so every query had to be changed



# Performance Comparison with Other Systems



# Target Systems



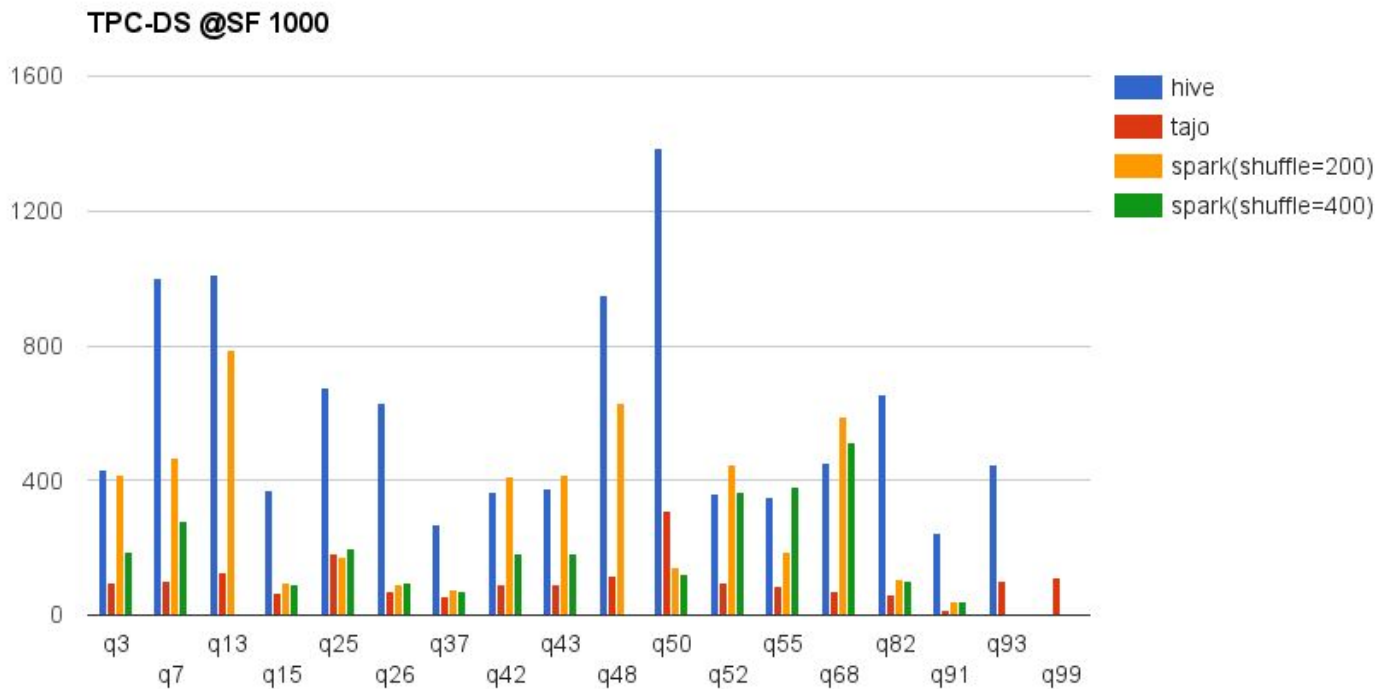
- Tajo (0.11.0)
  - Default configuration provided by GCP
    - Use the whole cpu and memory
- Hive (0.12)
  - Baseline performance
  - Default configuration provided by GCP
    - Use the whole cpu and memory

# Target Systems



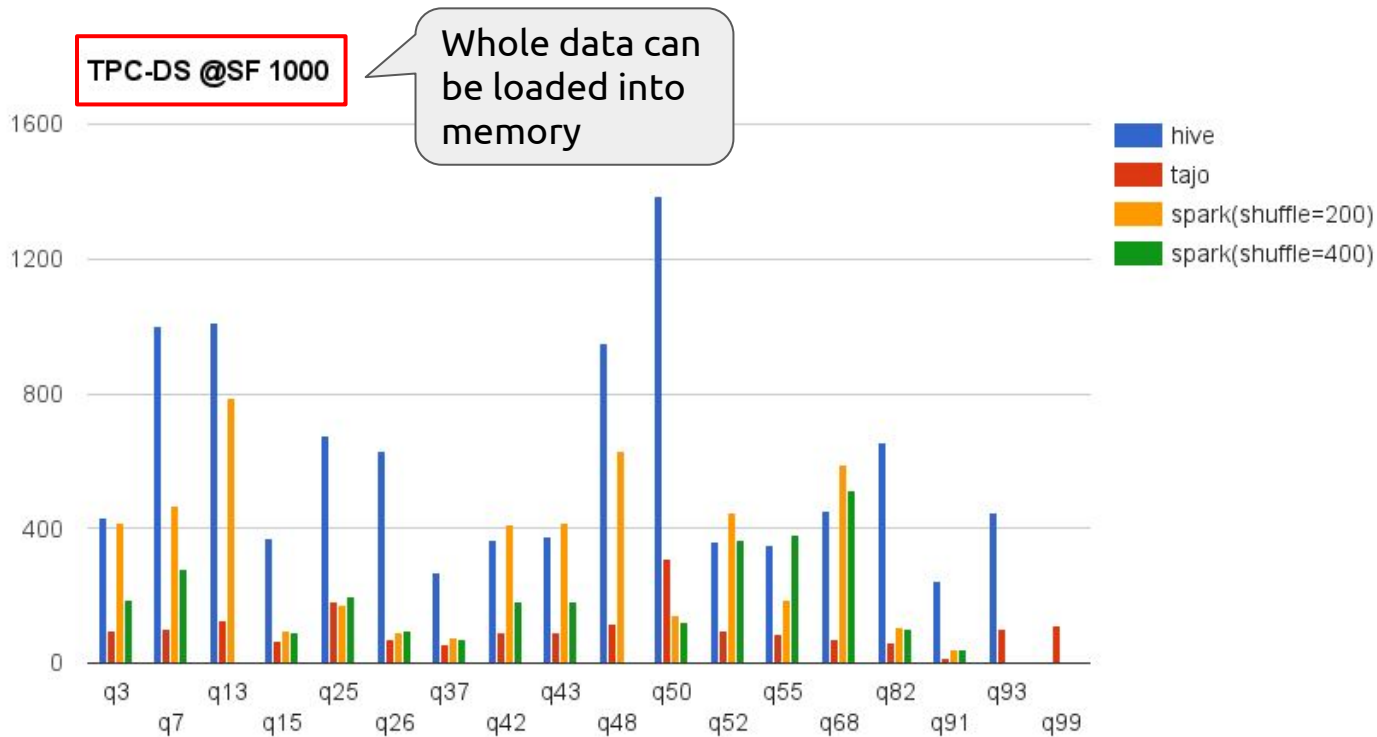
- Spark-SQL (1.5.0)
  - Default configuration provided by GCP
    - Use the whole cpu and memory
    - Tungsten enabled by default
  - `spark.sql.shuffle.partitions` is adjusted for better performance

# SF 1000, 50 instances

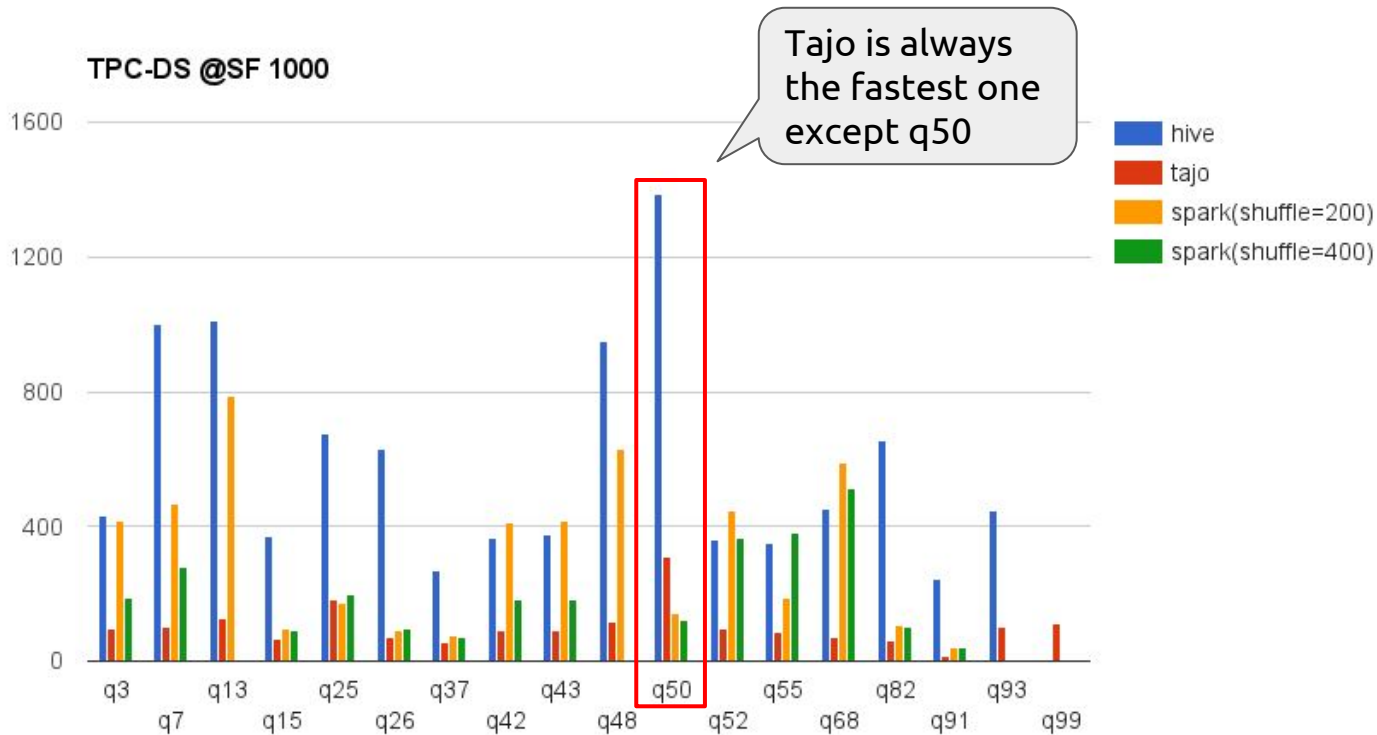




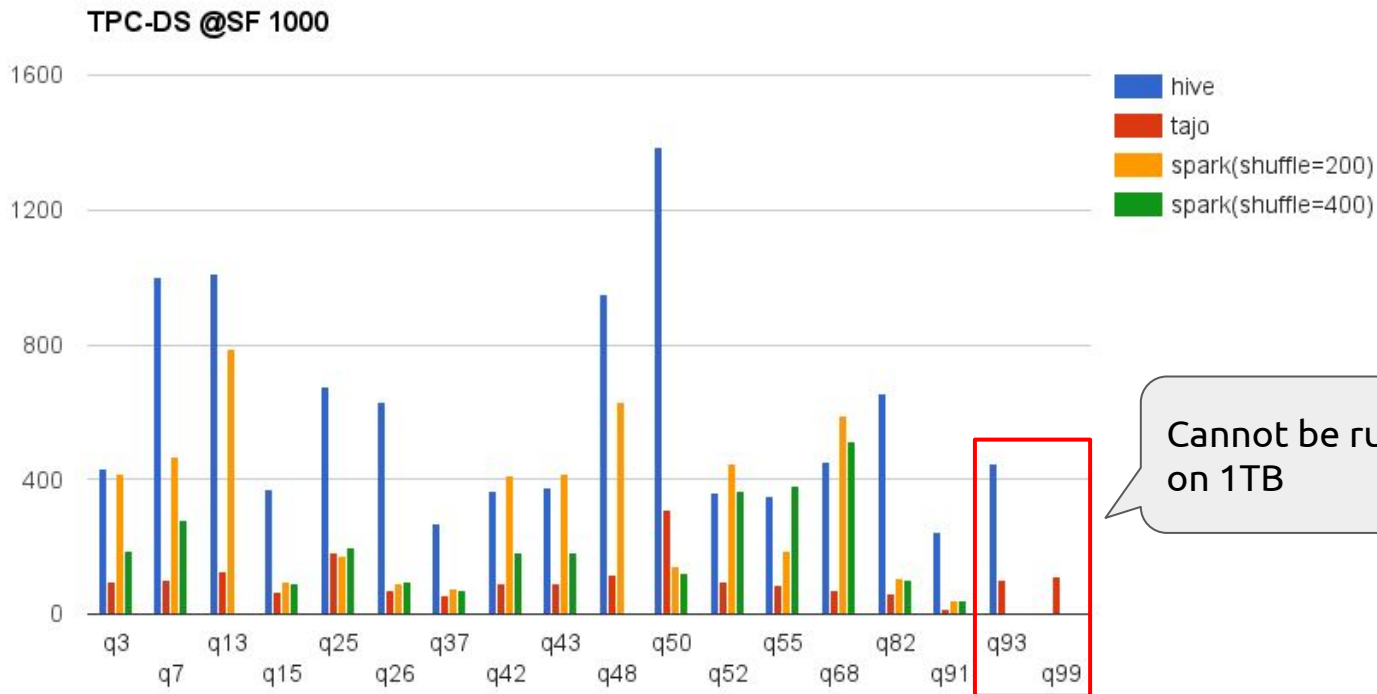
# SF 1000, 50 instances



# SF 1000, 50 instances

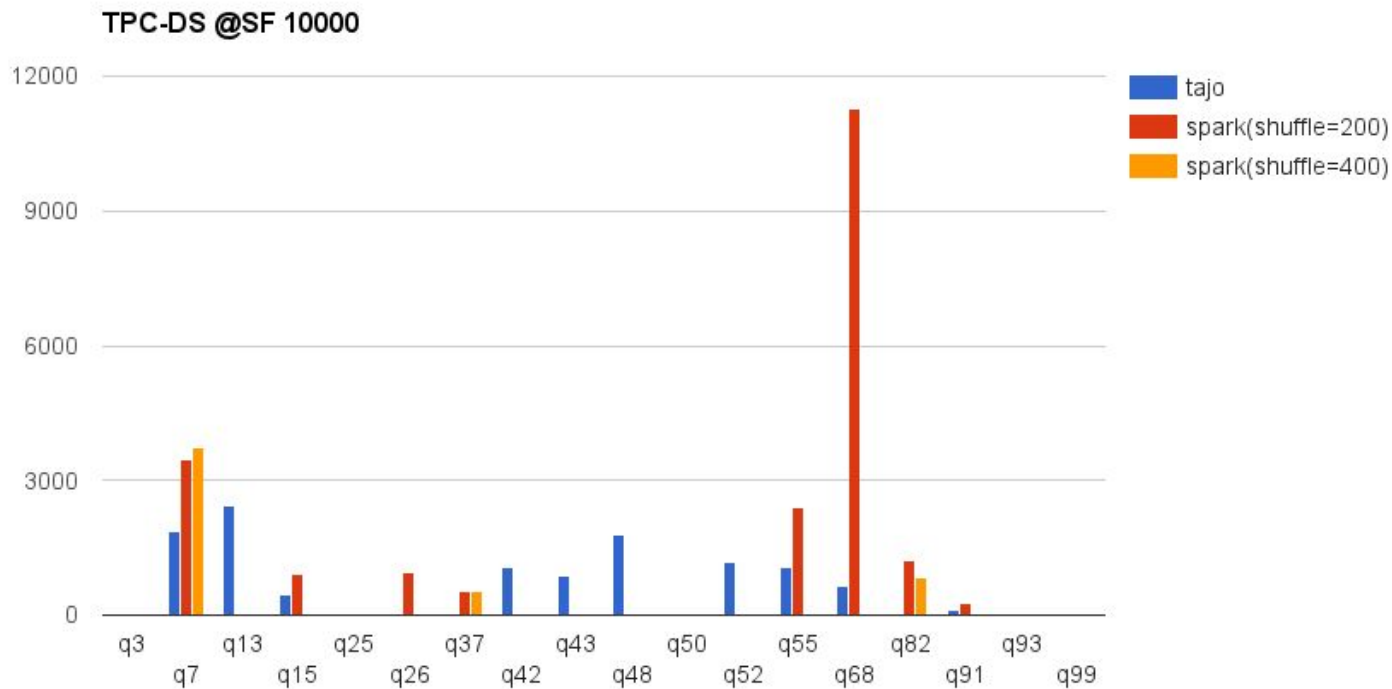


# SF 1000, 50 instances

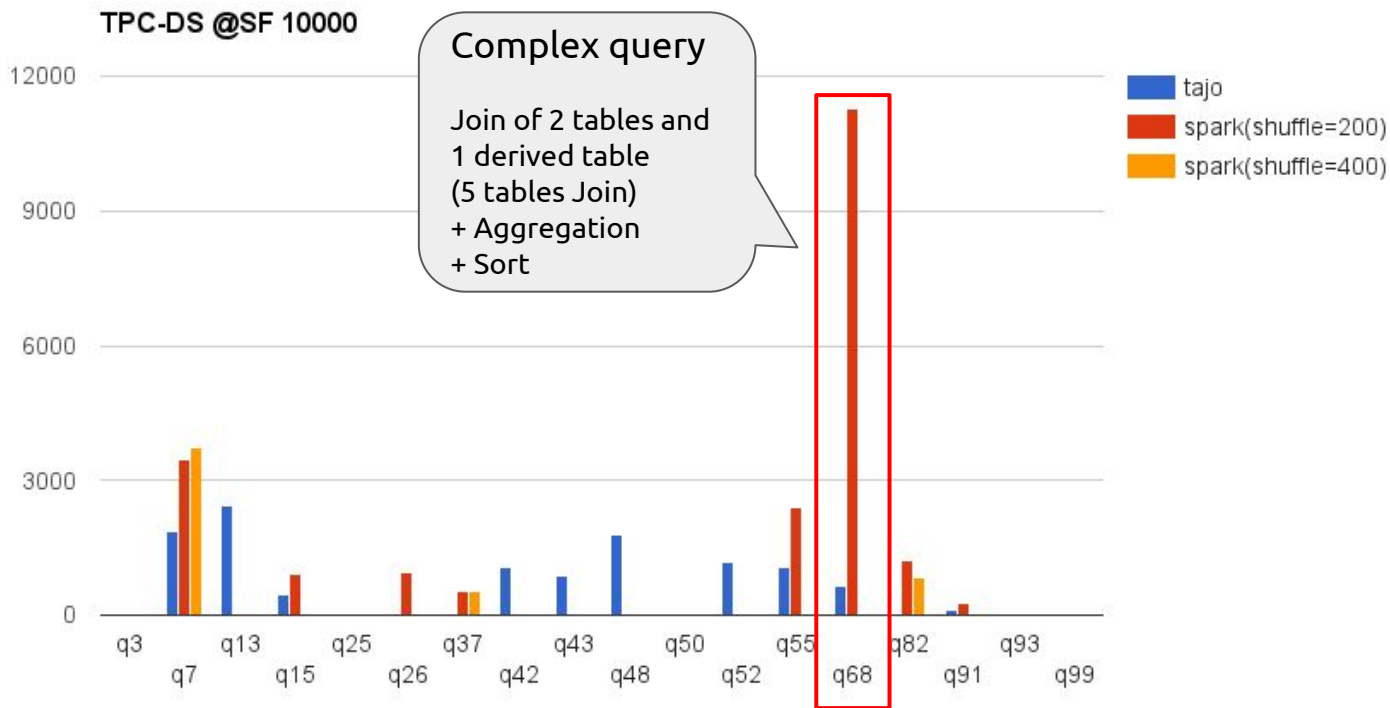


Cannot be run  
on 1TB

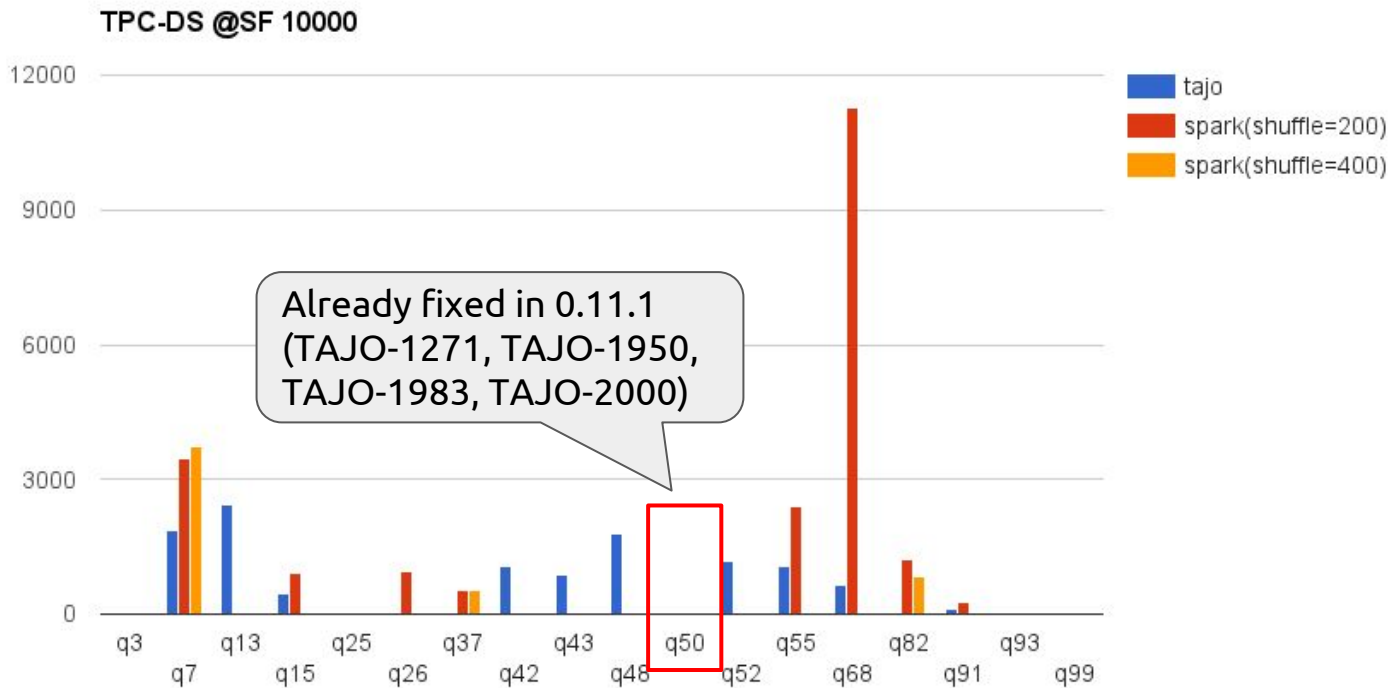
# SF 10000, 50 instances



# SF 10000, 50 instances



# SF 10000, 50 instances

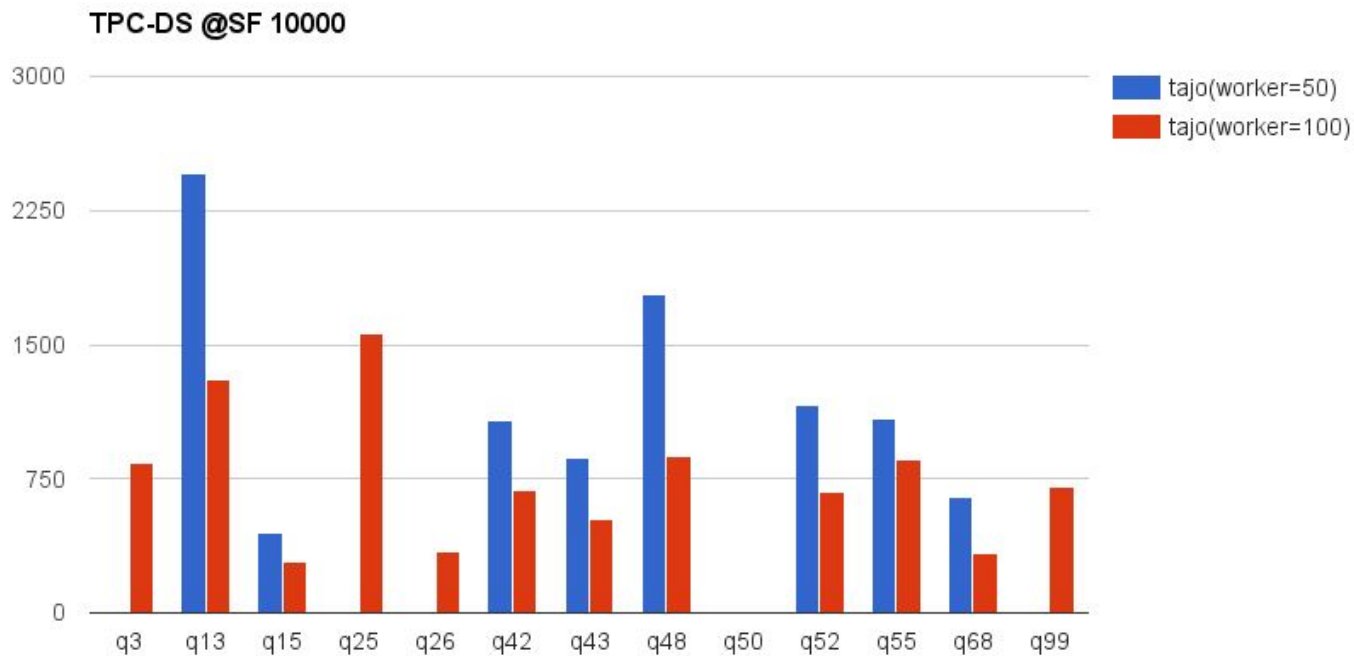




# Scalability Test

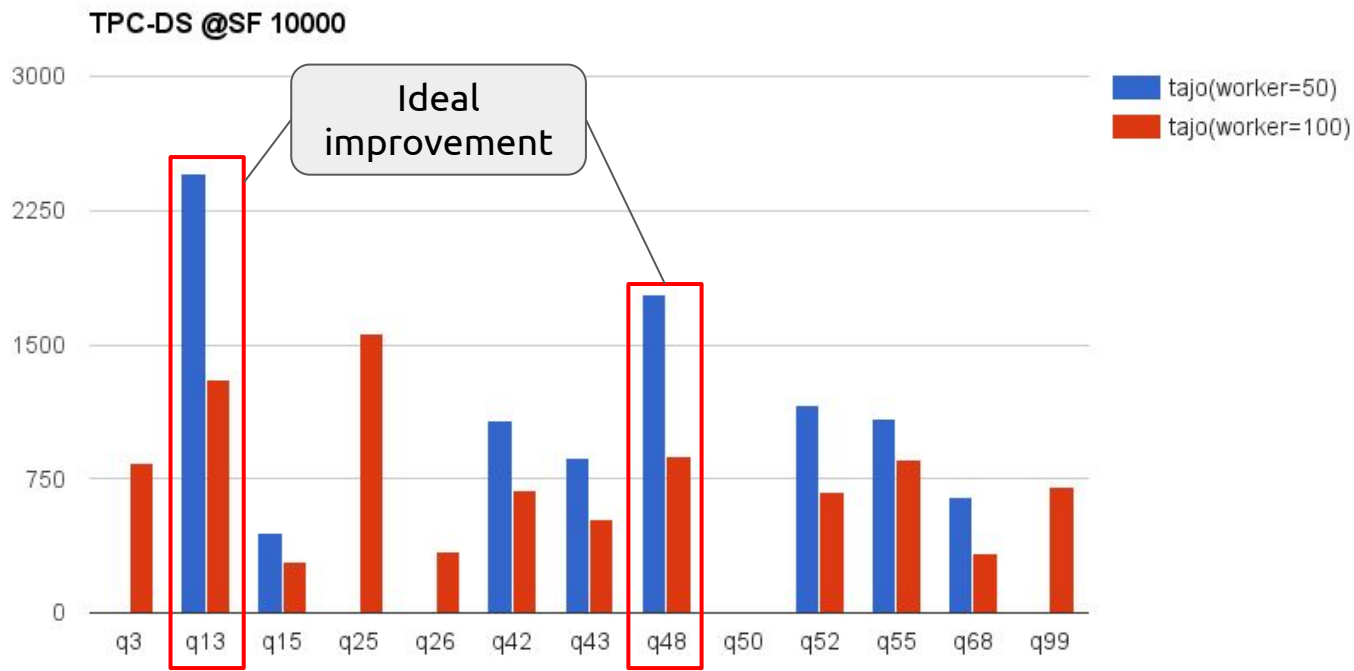


# SF 10000

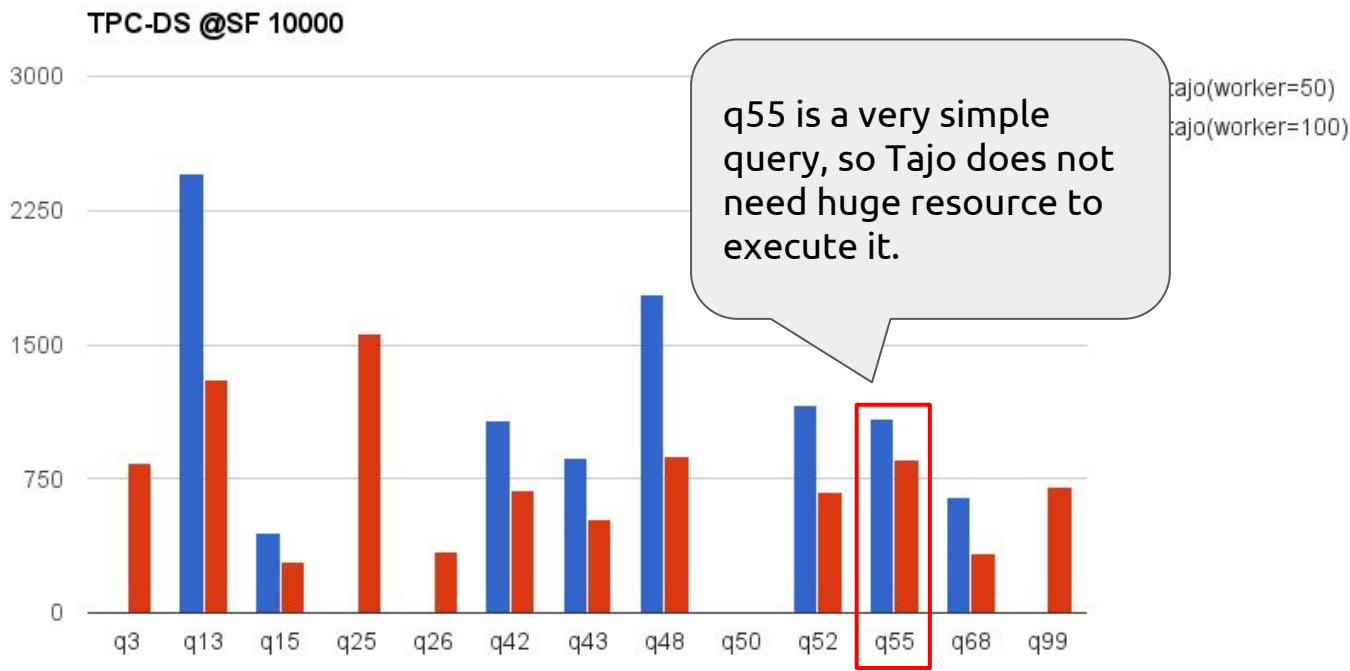




# SF 10000



# SF 10000



# Get Involved!



- We are recruiting contributors!
- General
  - <http://tajo.apache.org/>
- Getting Started
  - [http://tajo.apache.org/docs/current/getting\\_started.html](http://tajo.apache.org/docs/current/getting_started.html)
- Downloads
  - <http://tajo.apache.org/downloads.html>
- Issue tracker
  - <http://issues.apache.org/jira/browse/TAJO>
- Join the mailing list
  - [dev-subscribe@tajo.apache.org](mailto:dev-subscribe@tajo.apache.org)
  - [issues-subscribe@tajo.apache.org](mailto:issues-subscribe@tajo.apache.org)

# Q & A

