

Cost-efficiency and Performance Robustness in Serverless Data Exchange

ACM SIGMOD SRC | June 12 - 17, 2022 | Philadelphia, PA, US

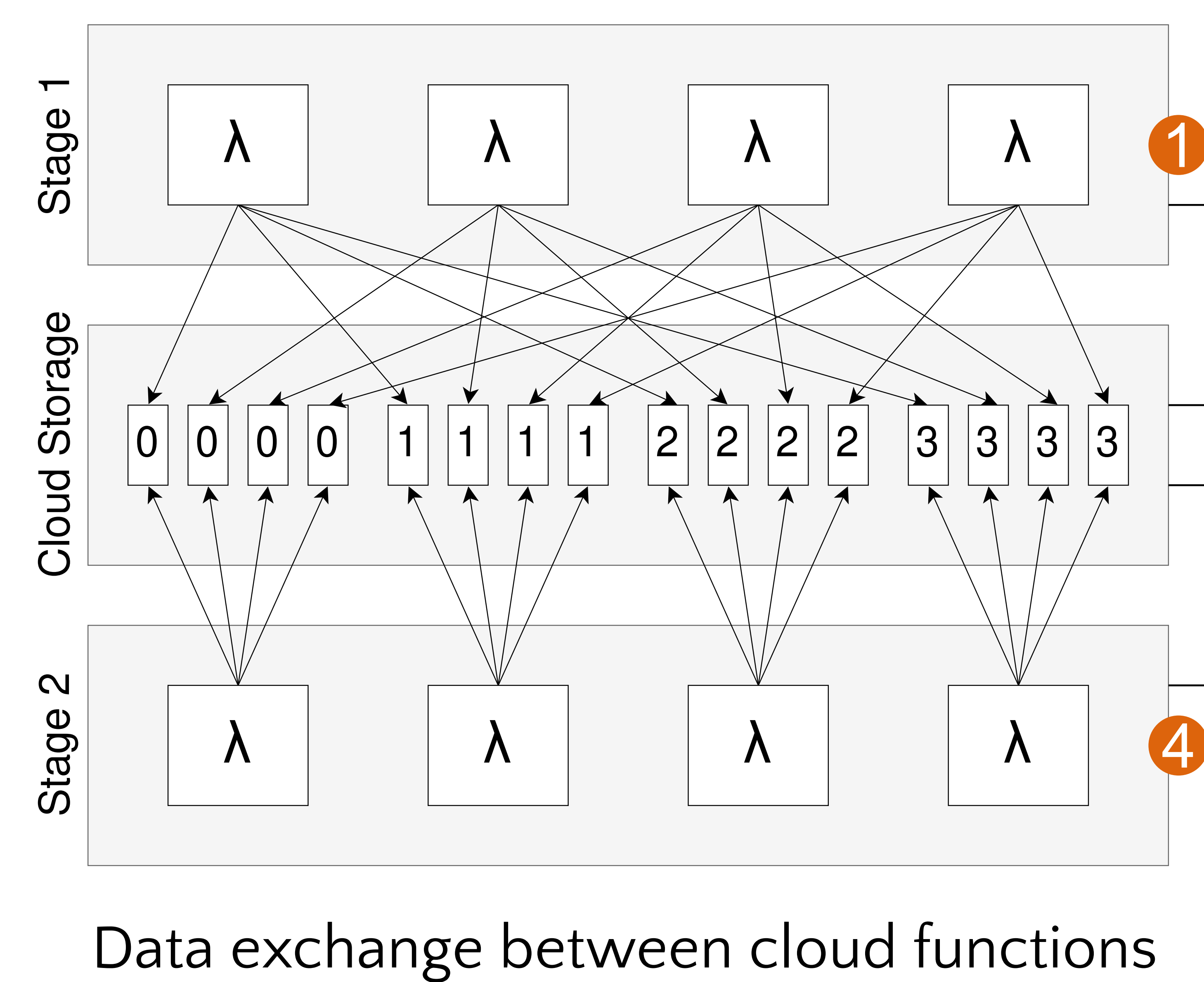
David Justen

Contribution

A *cost-efficient* data exchange strategy with *robust performance* for query processing on *cloud functions*, considering these problems:

- ⚡ Worker-to-worker communication only possible via cloud storage
- ⚡ High storage request prices
- ⚡ Tight request rate limits
- ⚡ Input load imbalance from data skew

General Approach



- 1) Partition input data in each cloud function
- 2) Export partitions to shared cloud storage
- 3) Import and combine partitions
- 4) Execute next pipeline on distinct partitions

Strategy

Storage Request Management

Multi-partition Writes

- ❑ Write all partitions into one single output file
- ❑ Append partition row ranges to file footer
- ❑ Read multi-partitioned files partially in consumer stage

→ n writes instead of $n*m$

n : producer stage worker count
 m : consumer stage worker count

Multi-stage Exchange

- ❑ Re-shuffle data with an additional stage
- ❑ Group workers that process the same input objects and pre-combine them

→ $nc + mg$ reads instead of $n*m$

g : group count in combiner stage
 c : worker count per group

Data Skew Mitigation

Over-partitioning

- ❑ Create p times as many partitions than workers in the next stage
- ❑ Aggregate partition sizes
- ❑ Assign each worker in the consumer stage $1..m*p$ partitions to even out partition size differences

Partitioned Cartesian Join

- ❑ For skewed joins, divide heavy hitter partitions horizontally into r (right) and s (left) sub-partitions
- ❑ Invoke $r*s$ additional workers to join each partition on the right side with every partition from the left side

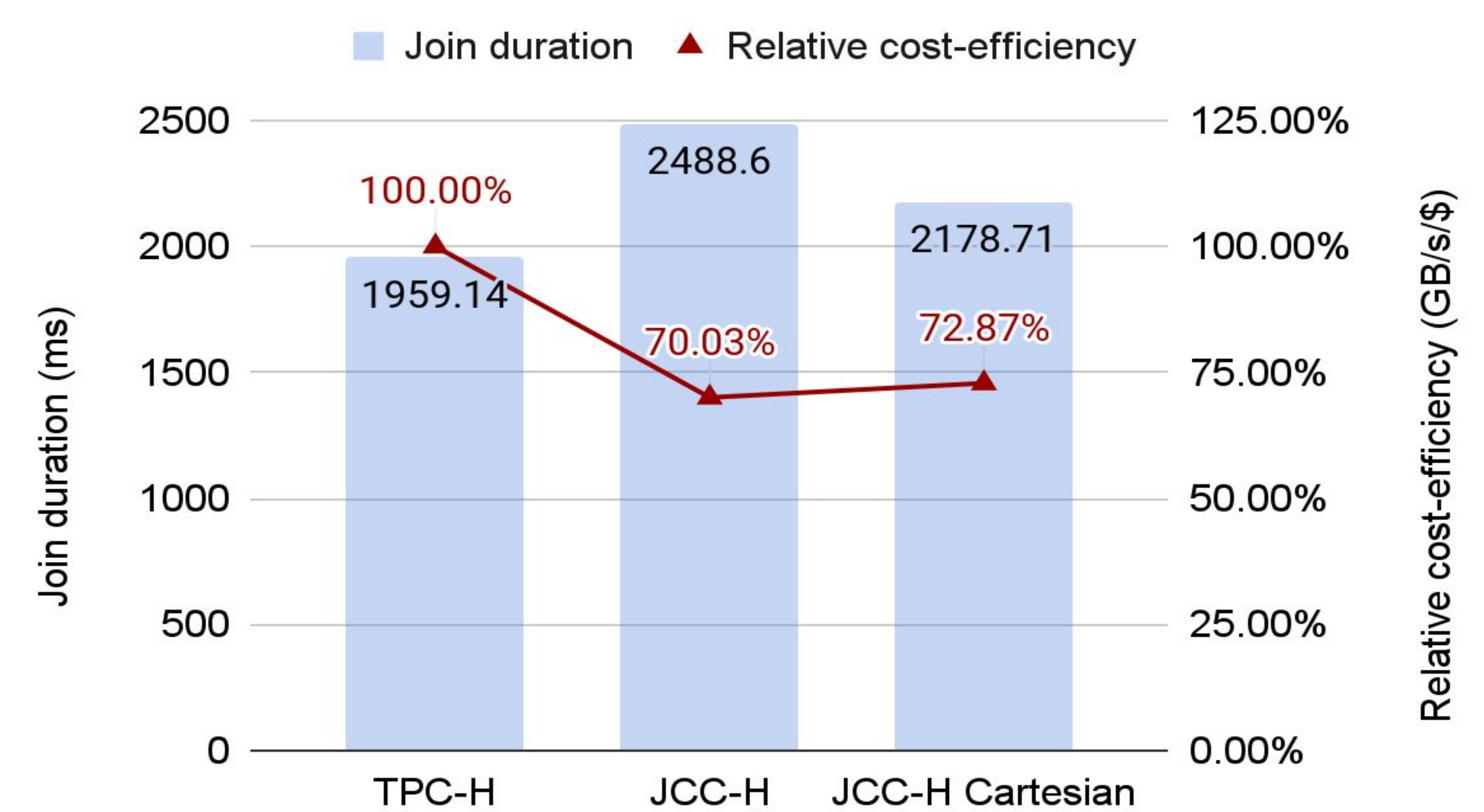
Evaluation

Single-stage vs. 2-stage Data Exchange



Shuffle TPC-H SF100 *lineitem* table on *L_orderkey* column

Data Skew Unmitigated vs. Partitioned Cartesian Join



Join SF10 *lineitem* table with *order* table from TPC-H and JCC-H