

Projekt – Ein Query-Optimizer für Cloudsysteme

Projektleiter: Foo Bar
Team: Autor 1, Autor 2

Datenbanksysteme und Informationsmanagement
Technische Universität Berlin



14. Oktober 2009

Agenda



Projektziele

Untersuchte Joinverfahren

Projektziele



Effiziente Implementierung relationaler Joins im Hadoop-Framework

- ▶ 2 Verfahren
 - Realisierung ohne Veränderung des Hadoop-Frameworks
 - Erweiterung des Hadoop-Frameworks um einen Merge-Operator
- ▶ Tests der Verfahren gegeneinander

Agenda



Projektziele

Untersuchte Joinverfahren



Repartition Join

- ▶ Equi-Join über 2 Relationen
- ▶ Mapper: wird jedes Tupel einmal weitergegeben oder ausgefiltert
- ▶ Partitioner: Modulo-Division des Schlüssel-HashWerts

```
public void map(...) throws IOException
{
    RelationTuple tuple = new RelationTuple(value);
    ...
    if (!filter.eval(tuple)) output.collect(outputKey, tuple);
}
```

```
public int getPartition (RepartitionKey k, RelationTuple v, int ptns)
{
    return Math.abs(k.getKey().hashCode() % ptns);
}
```

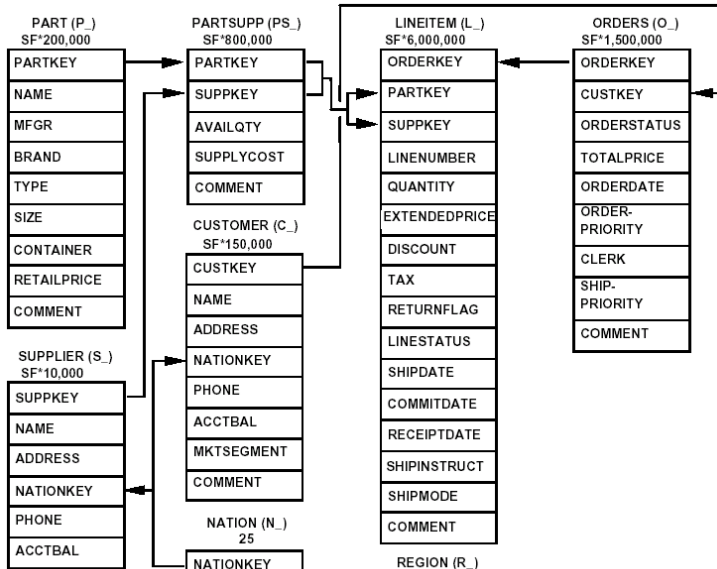
Repartition Join



- ▶ Sort:
 - primär nach Key
 - sekundär nach Relation (L Relation kommt zuerst)
- ▶ Reduce:
 - alle L -Tupel zum aktuellen Schlüssel sammeln
 - mit allen passenden R -Tupel kombinieren (kommen gleich danach)



TPCH Schema



Vielen Dank für Ihre Aufmerksamkeit



Fragen?