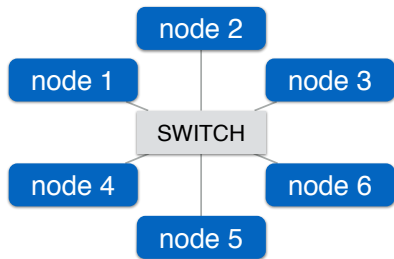# Locality-Sensitive Operators for Parallel Main-Memory Database Clusters

**Wolf Rödiger**, Tobias Mühlbauer, Philipp Unterbrunner*,
Angelika Reiser, Alfons Kemper, Thomas Neumann
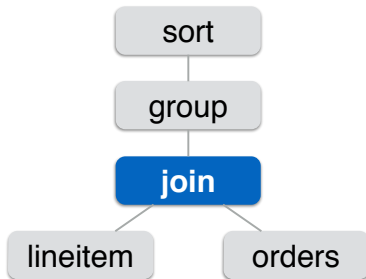
Technische Universität München, *Snowflake Computing, Inc.

# Scale Out

- **HyPer:** High-performance in-memory transaction and query processing system
- **Scale out** to process extremely large inputs
- Aiming at **clusters** with large main memory capacity

# Running Example (1)

- Focus on **analytical** query processing in this talk
- TPC-H query 12 used as **running example**
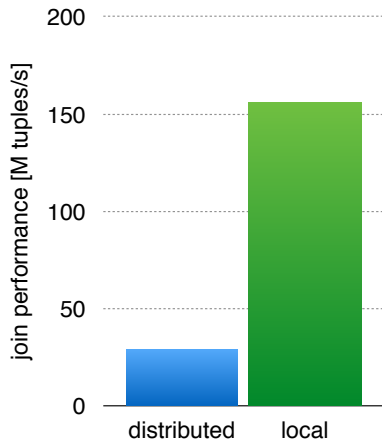- Runtime dominated by **join** orders ⋈ lineitem

# Running Example (2)

▸ Relations are **equally** distributed across nodes

▸ We make **no** other assumptions on data distribution

▸ **Network communication** required as tuples join with tuples on remote nodes

**orders**

| key | priority |
|-----|----------|
| 1 | 1-URGENT |
| 2 | 2-HIGH |
| 3 | 1-URGENT |
| 4 | 5-LOW |
| 5 | 3-MEDIUM |
| 6 | 1-URGENT |
| 7 | 2-HIGH |
| 8 | 1-URGENT |
| 9 | 1-URGENT |
| 10 | 2-HIGH |
| 11 | 3-MEDIUM |
| 12 | 5-LOW |
| 13 | 1-URGENT |
| 14 | 3-MEDIUM |
| 15 | 1-URGENT |
| 16 | 3-MEDIUM |
| 17 | 2-HIGH |
| 18 | 3-MEDIUM |
| 19 | 5-LOW |
| 20 | 1-URGENT |
| 21 | 2-HIGH |

**lineitem**

| key | shipmode |
|-----|----------|
| 1 | MAIL |
| 1 | MAIL |
| 1 | MAIL |
| 2 | SHIP |
| 2 | MAIL |
| 6 | SHIP |
| 6 | SHIP |
| 6 | SHIP |
| 6 | MAIL |
| 10 | SHIP |
| 11 | MAIL |
| 11 | MAIL |
| 13 | MAIL |
| 13 | MAIL |
| 13 | MAIL |
| 13 | SHIP |
| 17 | MAIL |
| 18 | MAIL |
| 18 | MAIL |
| 19 | SHIP |
| 20 | SHIP |

node 1 | node 2 | node 3

# Scale Out: Network is the Bottleneck

- **Single node:** Performance is bound algorithmically
- **Cluster:** Network is bottleneck for query processing
- We propose a novel join algorithm called **Neo-Join**
- **Goal:**
  Increase local processing to close the performance gap

# Neo-Join: Network-optimized Join

1. **Open Shop Scheduling**
   Efficient network communication

2. **Optimal Partition Assignment**
   Increase local processing

3. **Selective Broadcast**
   Handle value skew

# Open Shop Scheduling

Efficient network communication

# Standard Network Model

- **Star topology**
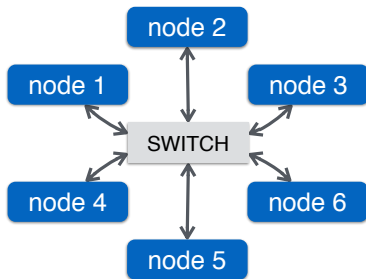  Nodes are connected to a
  central switch
- **Fully switched**
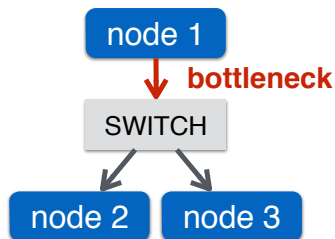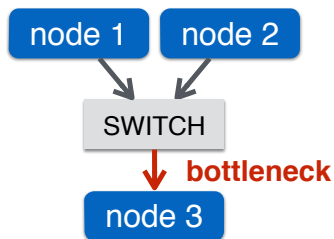  All links can be used
  simultaneously
- **Fully duplex**
  Nodes can both send and
  receive at full speed
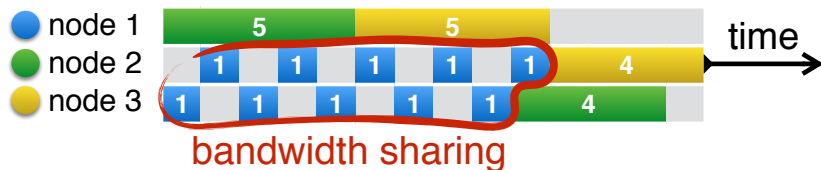
# Bandwidth Sharing



- Simultaneous use of a single link creates a **bottleneck**
- **Reduces bandwidth** by at least a factor of **2**

# Naïve Schedule



- ‣ Node 2 and 3 send to node 1 **at the same time**
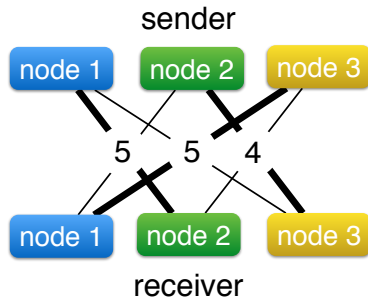- ‣ Bandwidth sharing increases **network duration** significantly

# Open Shop Scheduling (1)

- Avoiding bandwidth sharing **translates directly** to open shop scheduling
- **Network Transfer:** Receivers receive from at most **one** sender, senders send to at most **one** receiver
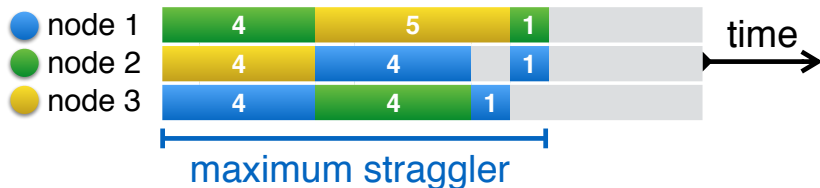- **Open Shop:** Processors perform **one** task at a time, only **one** task of a job is processed at a time

| Open Shop | Network Transfer |
|-----------|------------------|
| task | data transfer |
| processor | receiver |
| job | sender |
| execution time | message size |

# Open Shop Scheduling (2)

- **Bipartite graph** of senders and receivers
- **Edge weights** represent transfer size
- Scheduler repeatedly finds **perfect matchings**
- Each matching specifies one **communication phase**
- Transfers in a phase will **never** share bandwidth



sender

node 1   node 2   node 3

5   5   4
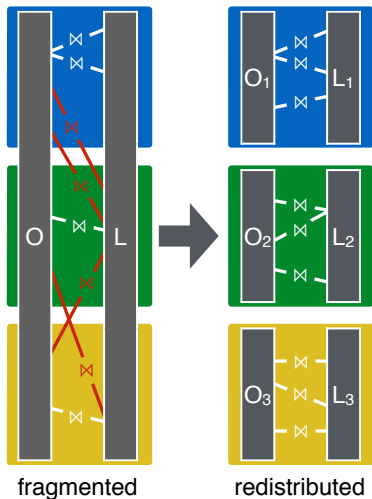
node 1   node 2   node 3

receiver

# Optimal Schedule



- Open shop schedule achieves minimal **network duration**
- Schedule duration determined by **maximum straggler**
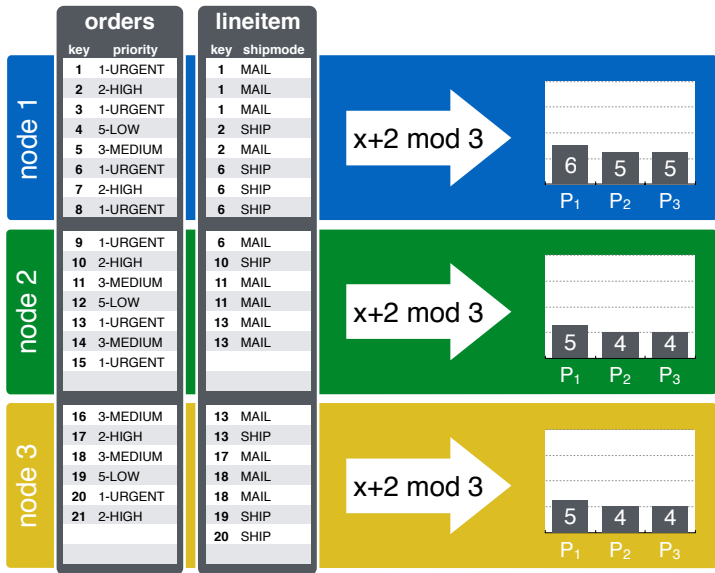
# Optimal Partition Assignment

Minimize network duration for distributed joins

# Distributed Join

- Tuples may join with tuples on **remote nodes**
- Repartition and redistribute **both relations** for local join
- Tuples will join only with the **corresponding partition**
- Using hash, range, radix, or other **partitioning** scheme
- **In any case:** Decide how to **assign** partitions to nodes



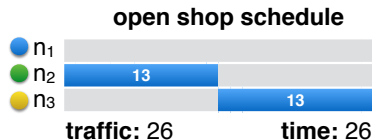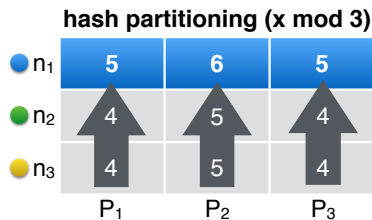fragmented          redistributed

# Running Example: Hash Partitioning

# Assign Partitions to Nodes (1)

**Option 1:** Minimize network traffic

- Assign partition to node that owns its **largest part**
- Only the **small fragments** of a partition sent over the network
- Schedule with minimal network traffic may have **high duration**

**hash partitioning (x mod 3)**

|  | | | |
|---|---|---|---|
| n$_1$ | 5 | 6 | 5 |
| n$_2$ | 4 | 5 | 4 |
| n$_3$ | 4 | 5 | 4 |
|  | P$_1$ | P$_2$ | P$_3$ |

**open shop schedule**

n$_1$

n$_2$ 13

n$_3$ 13

**traffic:** 26          **time:** 26

# Assign Partitions to Nodes (2)

**Option 2:** Minimize response time:

- **Query response time** is time from request to result
- Query response time dominated by **network duration**
- To minimize network duration, minimize **maximum straggler**

**hash partitioning (x mod 3)**



| | $P_1$ | $P_2$ | $P_3$ |
|---|---|---|---|
| $n_1$ | 5 | 6 | 5 |
| $n_2$ | 4 | 5 | 4 |
| $n_3$ | 4 | 5 | 4 |

**open shop schedule**



$n_1$ | 4 | 5 | 1
$n_2$ | 4 | 4 | 1
$n_3$ | 4 | 4 | 1

**traffic:** 28          **time:** 10

# Minimize Maximum Straggler

- Formalized as mixed-integer **linear program**
- Objective function minimizes **maximum straggler**
- Shown to be **NP-hard** (see paper for proof sketch)
- In practice **fast enough** using CPLEX or Gurobi (< 0.5 % overhead for 32 nodes, 200 M tuples each)
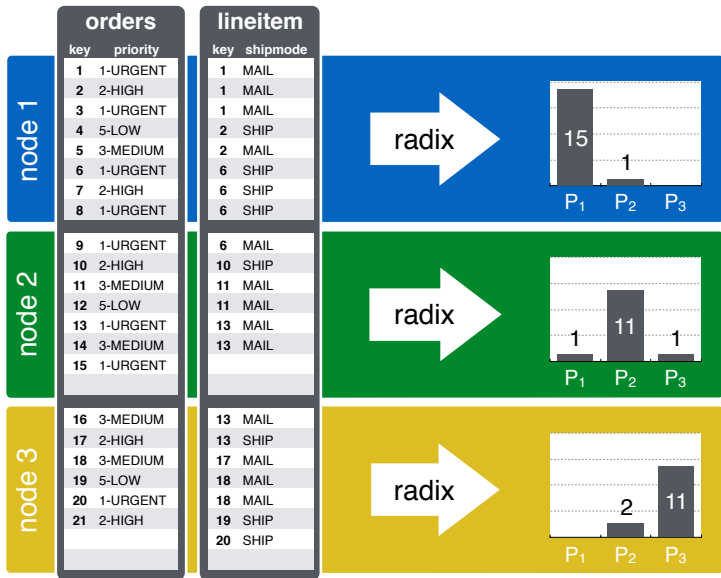- Partition assignment can optimize **any partitioning**

**minimize $w$, subject to**

$$w \geq \sum_{j=0}^{p-1} h_{ij}(1 - x_{ij}) \qquad 0 \leq i < n$$

$$w \geq \sum_{j=0}^{p-1} \left( x_{ij} \sum_{k=0, i \neq k}^{n-1} h_{kj} \right) \qquad 0 \leq i < n$$
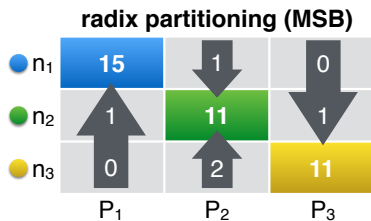
$$1 = \sum_{i=0}^{n-1} x_{ij} \qquad 0 \leq j < p$$
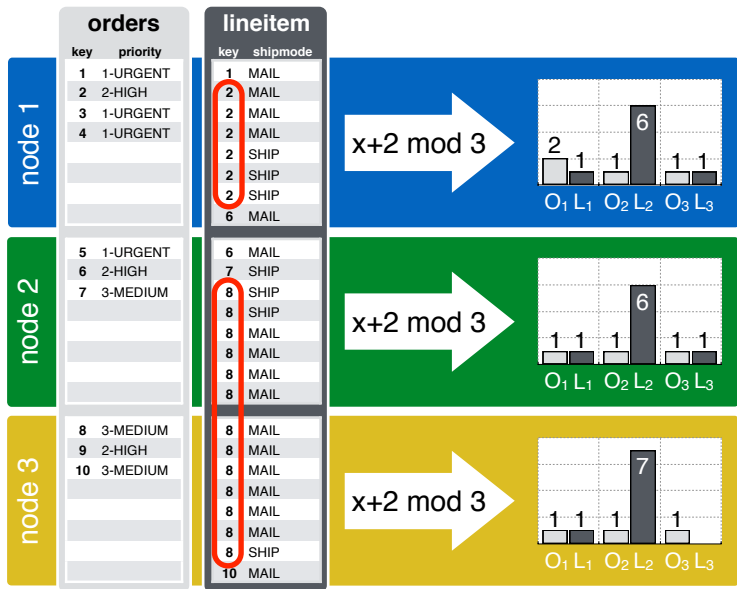
# Running Example: Locality

# Locality

- Running example exhibits **time-of-creation** clustering
- **Radix repartitioning** on most significant bits retains locality
- Partition assignment can **exploit locality**
- Significantly reduces **query response time**

**radix partitioning (MSB)**

| | $P_1$ | $P_2$ | $P_3$ |
|---|---|---|---|
| $n_1$ | 15 | 1 | 0 |
| $n_2$ | 1 | 11 | 1 |
| $n_3$ | 0 | 2 | 11 |

**open shop schedule**

$n_1$ 1
$n_2$ 1 1
$n_3$ 2

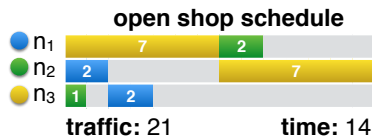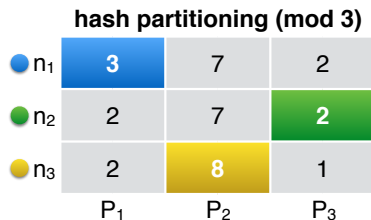**traffic:** 5          **time:** 3

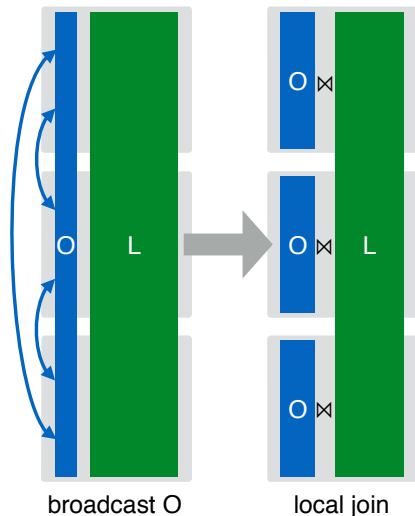# Selective Broadcast

Handle value skew

# Running Example: Skew

# Skew

- **Value skew** can lead to some very large partitions
- Assignment of these partitions increases **network duration**
- One may try to balance skewed partitions by partitioning the input into **more partitions**
- **High skew** is still a problem

**hash partitioning (mod 3)**

|     | $P_1$ | $P_2$ | $P_3$ |
|-----|-------|-------|-------|
| $n_1$ | **3** | 7 | 2 |
| $n_2$ | 2 | 7 | **2** |
| $n_3$ | 2 | **8** | 1 |

**open shop schedule**

| $n_1$ | 7 | 2 | |
| $n_2$ | 2 | | 7 |
| $n_3$ | 1 | 2 | |

**traffic:** 21        **time:** 14

# Broadcast

- **Alternative** data redistribution scheme
- **Replicate** the smaller relation between all nodes
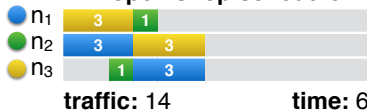- Larger relation **remains fragmented** across nodes



broadcast O          local join

# Selective Broadcast

- Decide **per partition** whether to assign or broadcast
- **Broadcast** partitions with large relation size difference
- **Assign** the other partitions taking locality into account
- **Role reversal possible:** Broadcast different partitions by different relations

**hash partitioning (mod 3)**

| | | | | | | |
|---|---|---|---|---|---|---|
| $n_1$ | 2 | 1 | 1 | 6 | 1 | 1 |
| $n_2$ | 1 | 1 | 1 | 6 | 1 | 1 |
| $n_3$ | 1 | 1 | 1 | 5 | 2 | 2 |
| | $O_1$ | $L_1$ | $O_2$ | $L_2$ | $O_3$ | $L_3$ |

**open shop schedule**

| | |
|---|---|
| $n_1$ | 3 1 |
| $n_2$ | 3 3 |
| $n_3$ | 1 3 |

**traffic:** 14          **time:** 6
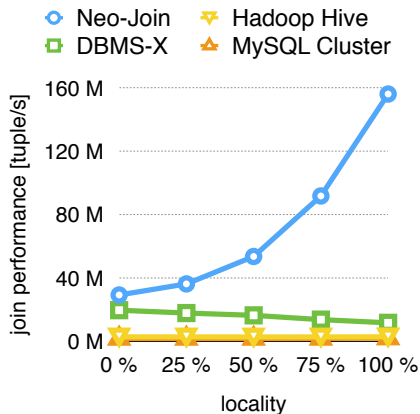
# Evaluation

# Locality

- Vary **locality** from **0 %** (uniform distribution) to **100 %** (range partitioning)
- Neo-Join improves **join performance** from 29 M to 156 M tuples/s (> 500 %)
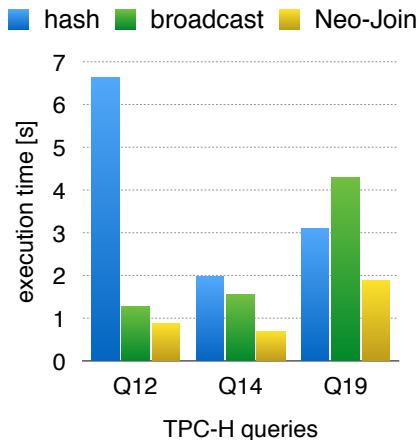- 3 nodes (Core i7, 4 cores, 3.4 GHz, 32 GB RAM), 600 M tuples (64 bit key, 64 bit payload)

# Skew

- **Zipfian distribution** models realistic data skew
- Using **more partitions** alleviates the problem
- Selective broadcast actually **improves** performance for skewed inputs
- 4 nodes, 400 M tuples

| partitions | Zipf factor s | | | | |
|---|---|---|---|---|---|
| | 0.00 | 0.25 | 0.50 | 0.75 | 1.00 |
| 16 | 27 s | 24 s | 23 s | **29 s** | **44 s** |
| 512 | 23 s | 23 s | 23 s | 23 s | **33 s** |
| 16 (SB) | 24 s | 24 s | 23 s | **20 s** | **10 s** |

# TPC-H Results (scale factor 100)

- Results for three selected **TPC-H** queries
- **Broadcast** outperforms **hash** for large relation size differences
- Neo-Join always performs better due to **selective broadcast** and **locality**
- 4 nodes, scale factor 100

# Summary

Motivation:

- Network is the **bottleneck** for distributed query processing
- Increase **local processing** to close the performance gap

Contributions:

- **Open Shop Scheduling** avoids bandwidth sharing
- **Optimal Partition Assignment** minimizes query response time and can exploit locality in the data distribution
- **Selective Broadcast** combines repartitioning and broadcast to improve the performance for skewed inputs