



# Parallel Joins in Main-Memory Data(base) Processing

Alfons Kemper, et al.

Fakultät für Informatik, TUM

# Joining two Relations $R$ and $S$ : $R \bowtie S$

$R$			$S$		
$A$	$B$	$C$	$C$	$D$	$E$
$a_1$	$b_1$	$c_1$	$c_1$	$d_1$	$e_1$
$a_2$	$b_2$	$c_2$	$c_3$	$d_2$	$e_2$
$a_3$	$b_3$	$c_1$	$c_4$	$d_3$	$e_3$
$a_4$	$b_4$	$c_2$	$c_5$	$d_4$	$e_4$
$a_5$	$b_5$	$c_3$	$c_7$	$d_5$	$e_5$
$a_6$	$b_6$	$c_2$	$c_8$	$d_6$	$e_6$
$a_7$	$b_7$	$c_6$	$c_5$	$d_7$	$e_7$

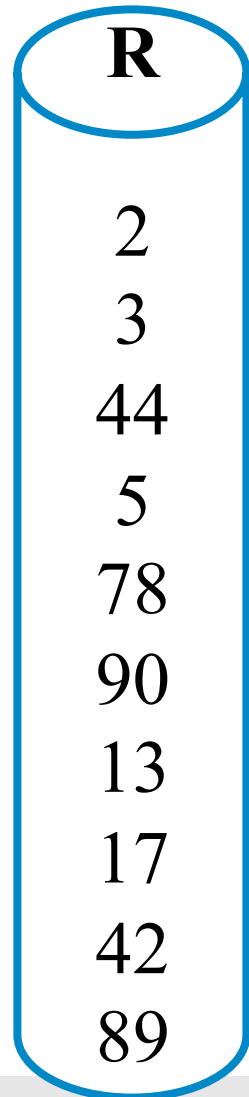
$\bowtie$

$R \bowtie S$				
$A$	$B$	$C$	$D$	$E$
$a_1$	$b_1$	$c_1$	$d_1$	$e_1$
$a_3$	$b_3$	$c_1$	$d_1$	$e_1$
$a_5$	$b_5$	$c_3$	$d_2$	$e_2$

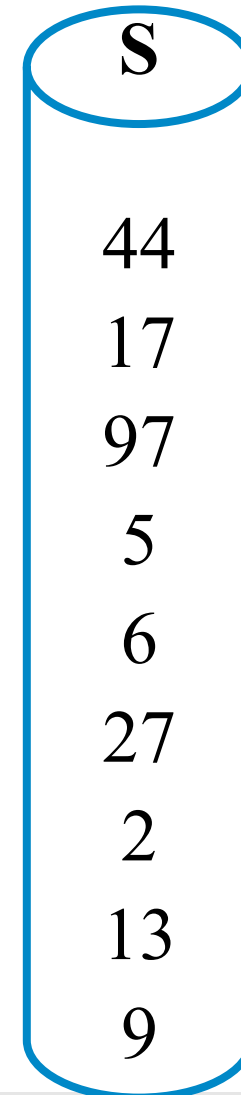
$=$



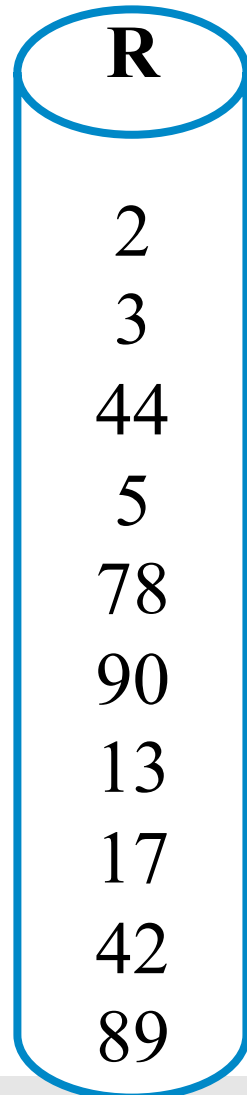
## Algorithmics:



$$R \cap S$$

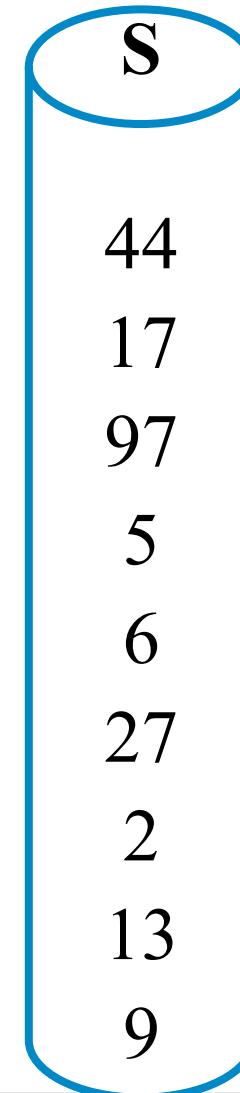


## Algorithmics:



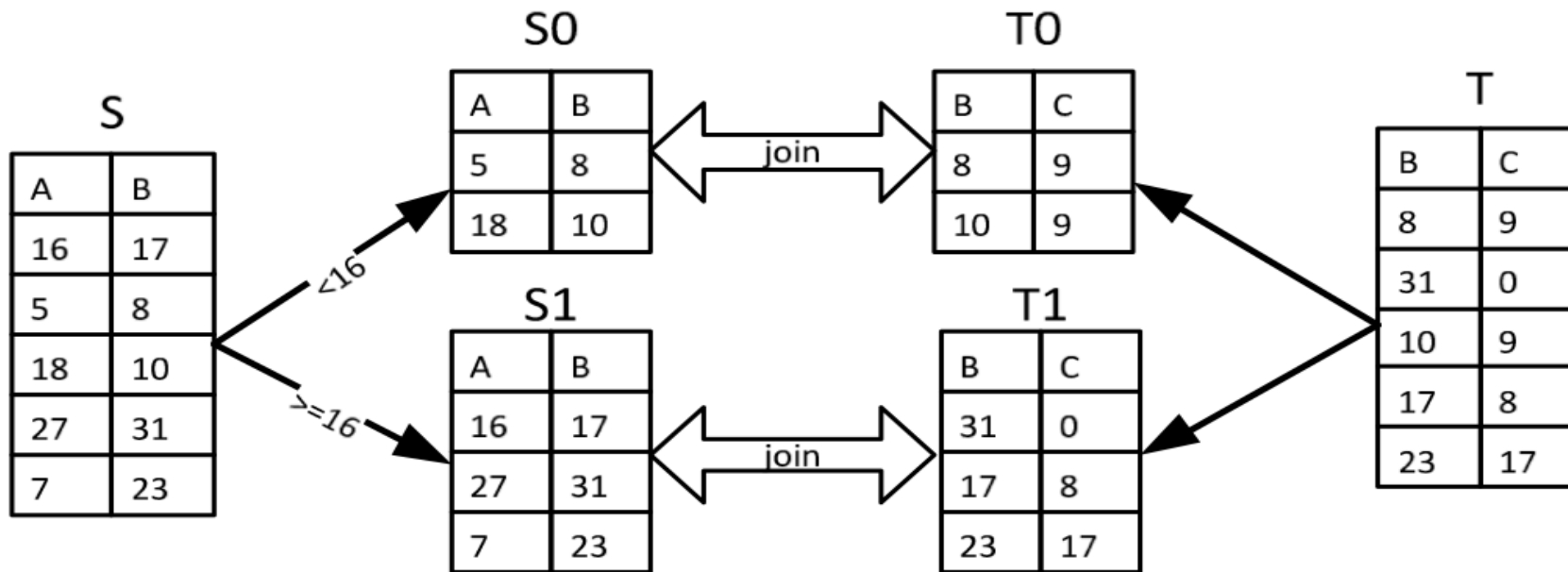
$R \cap S$

- Nested Loop:  $O(N^2)$
- Sort-Based:  $O(N \log N)$
- Partitioning and Hashing

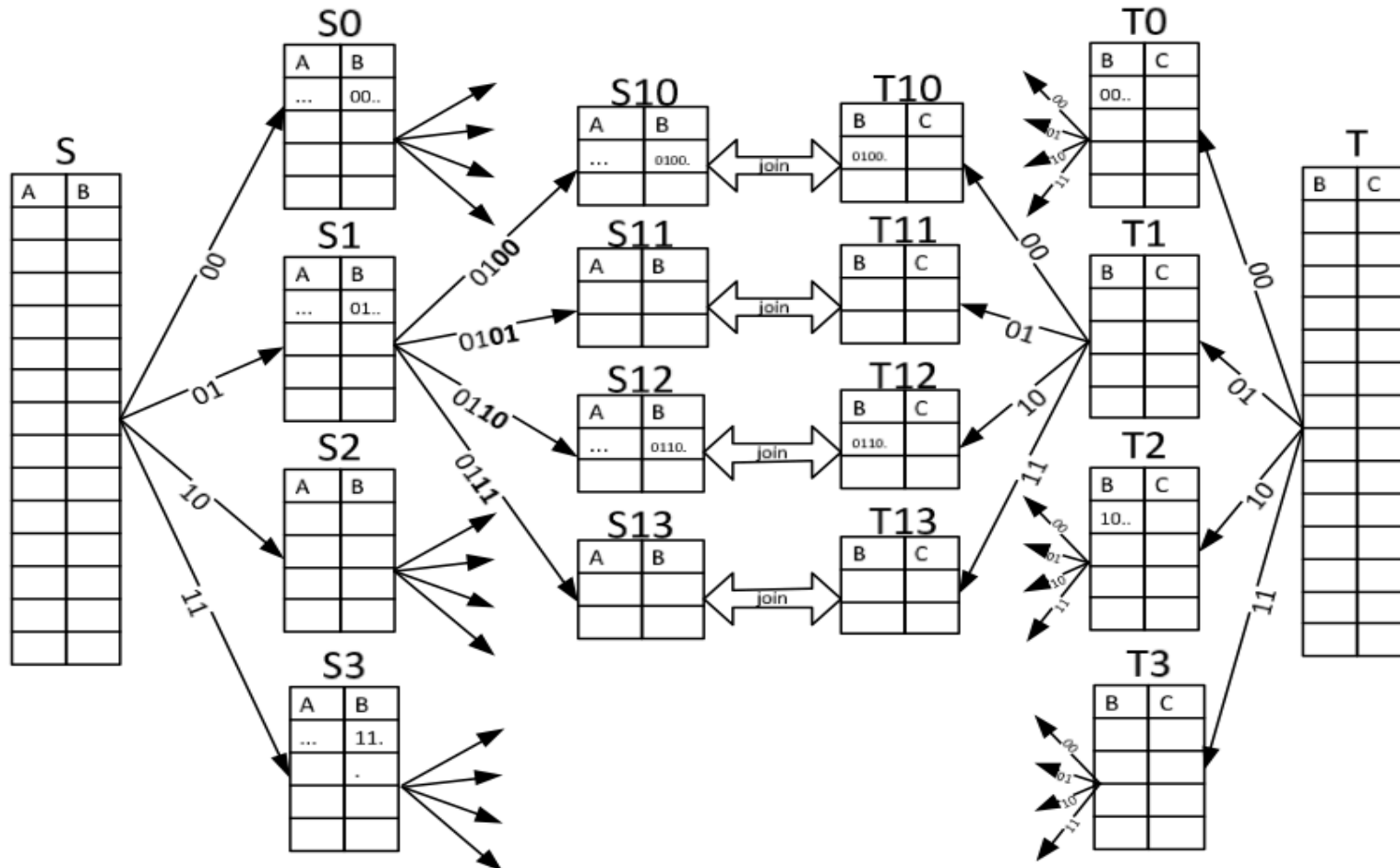




# Parallel Radix-Join

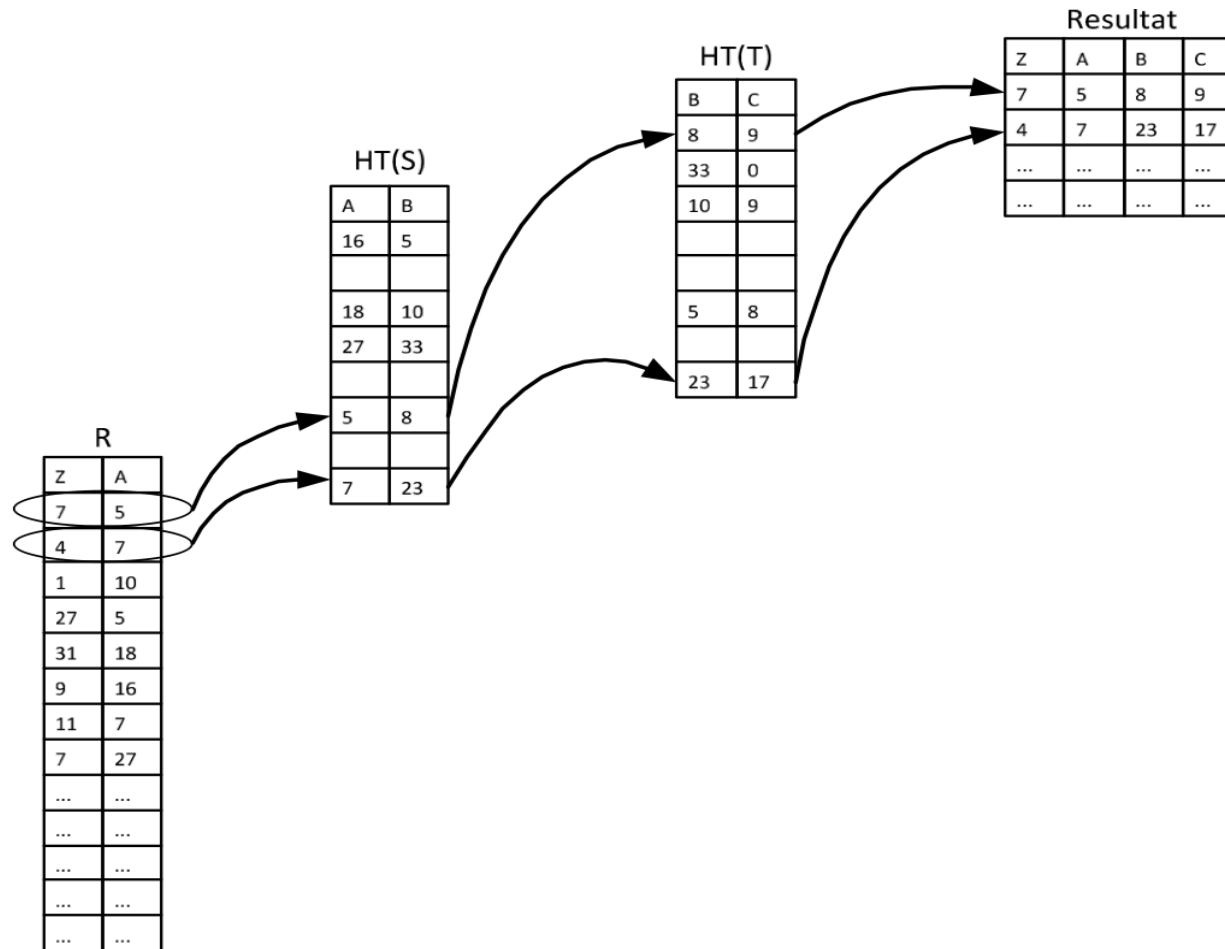


# Multiple Phase Radix-Joins: Cache-Locality





# Hash-Join-Teams: Global Hash Table



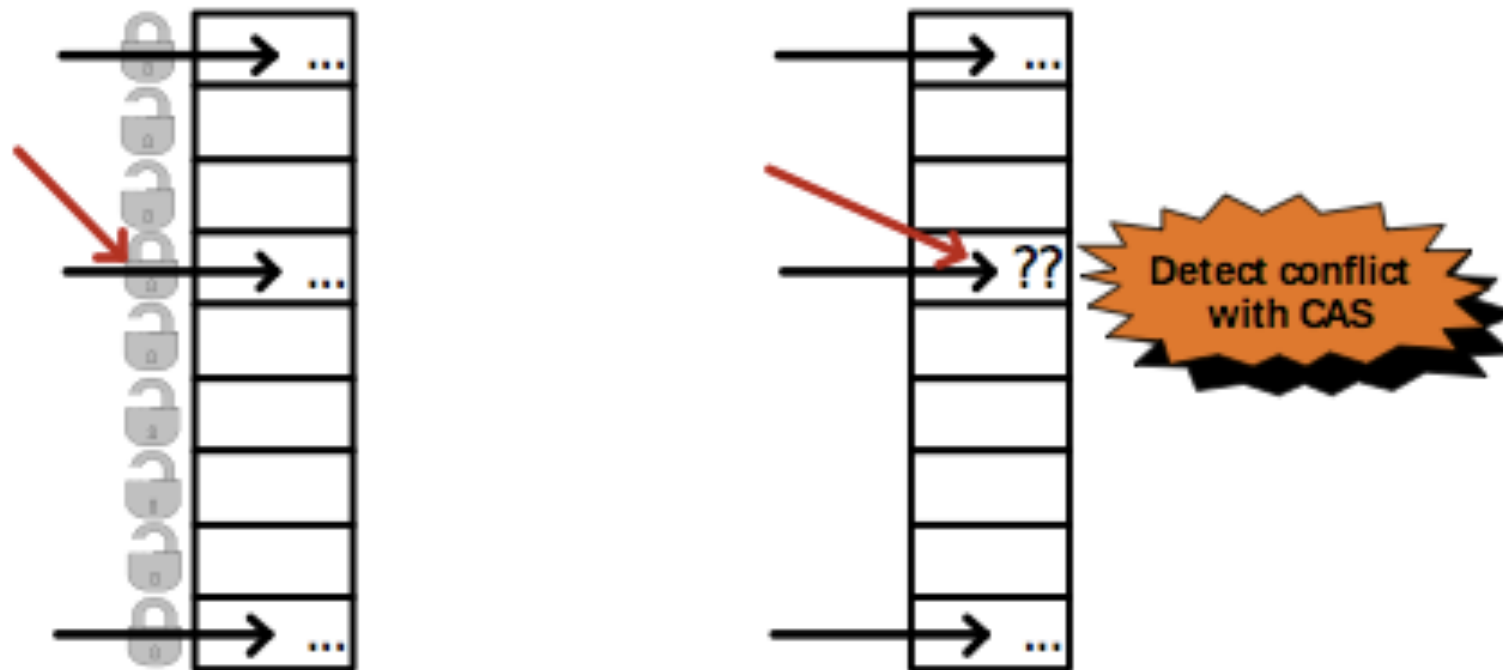


Fig. 1: Pessimistic vs. optimistic write access to a hash table



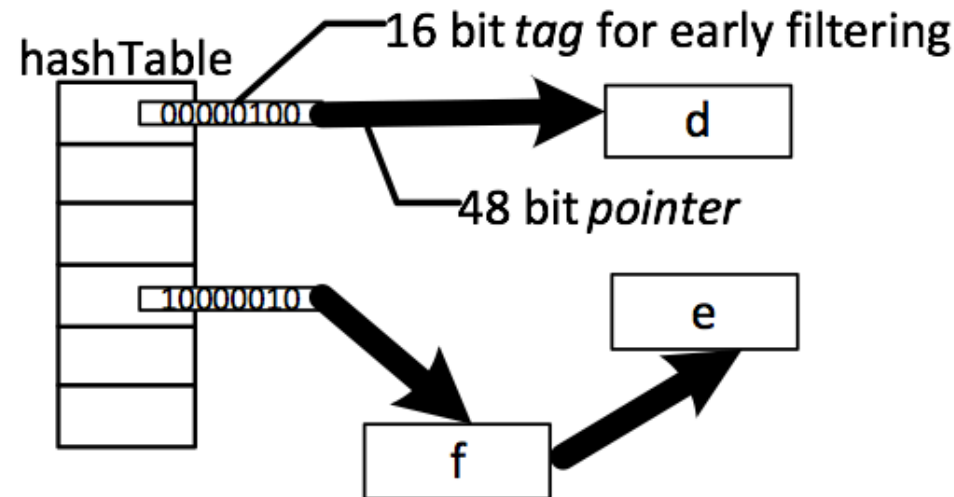
### Listing 1.1: Atomic insert function

---

```
1 insertAtomic(uint64_t key, uint64_t value) {
2     uint64_t hash = hashFunction(key);
3     uint64_t pos = hash & mask;
4     while (table[pos].h != 0
5           || (! CAS(&table[pos].h, 0, hash))) {
6         pos = (pos + 1) & mask;
7     }
8     table[pos].k = key;
9     table[pos].v = value;
10 }
```

# Alternative Hash Table with Overflow

## Buc

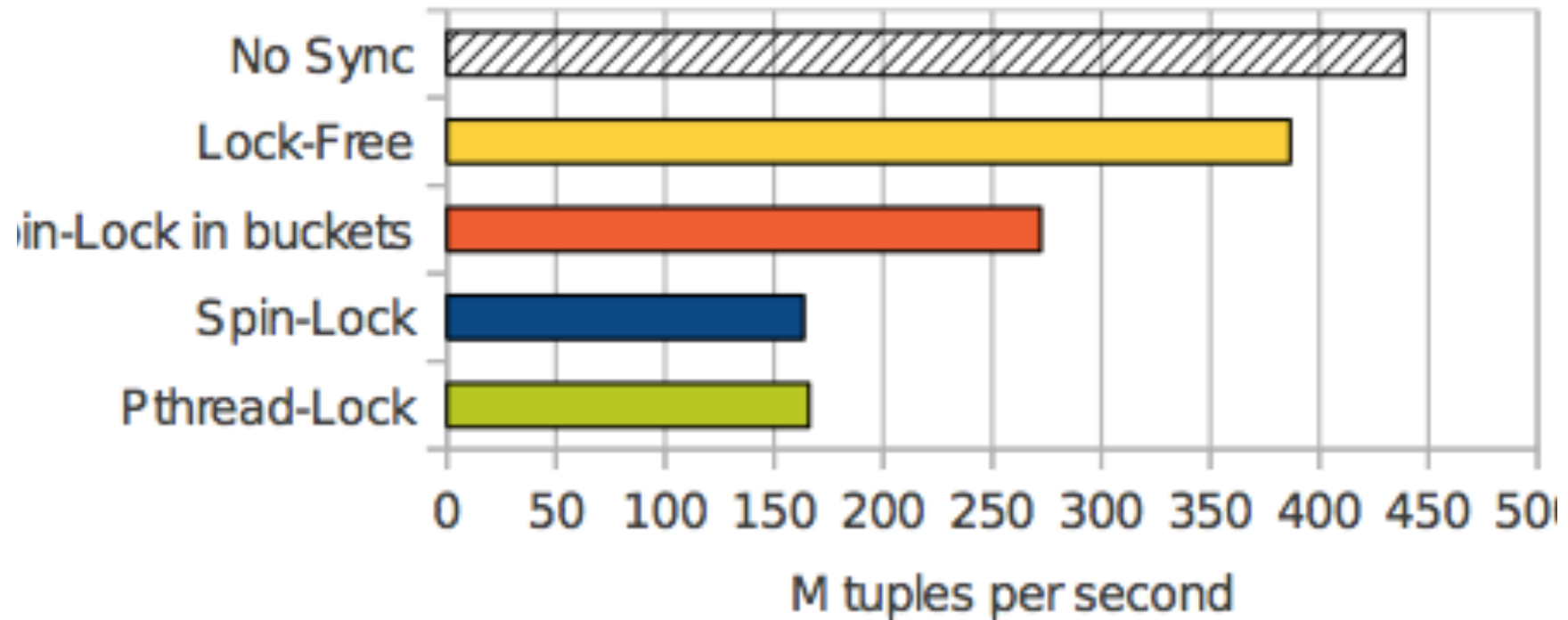



---

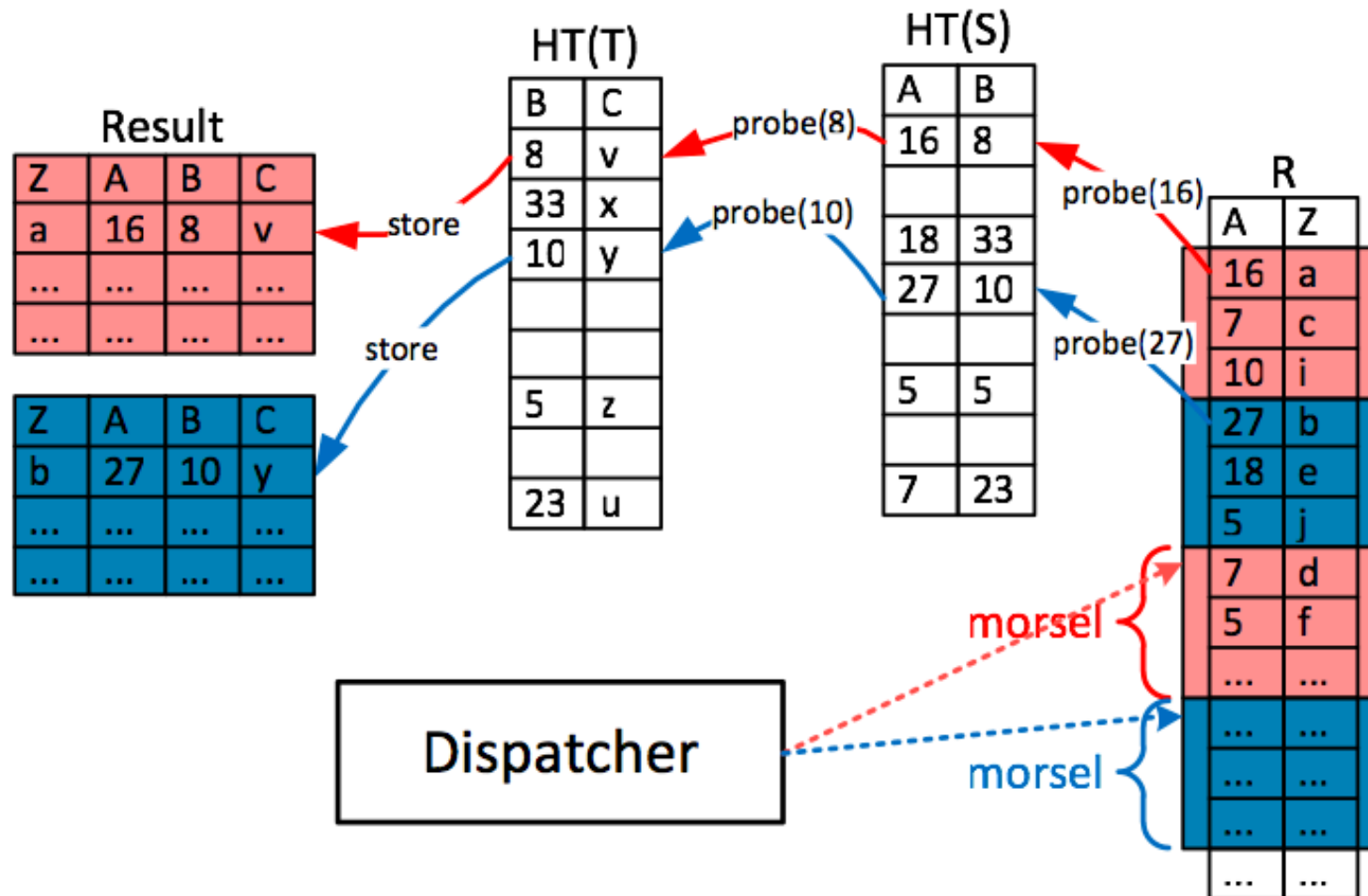
```

1 insert(entry) {
2     // determine slot in hash table
3     slot = entry->hash >> hashTableShift
4     do {
5         old = hashTable[slot]
6         // set next to old entry without tag
7         entry->next = removeTag(old)
8         // add old and new tag
9         new = entry | (old&tagMask) | tag(entry->hash)
10        // try to set new value, repeat on failure
11    } while (!CAS(hashTable[slot], old, new))
12 }
    
```

---



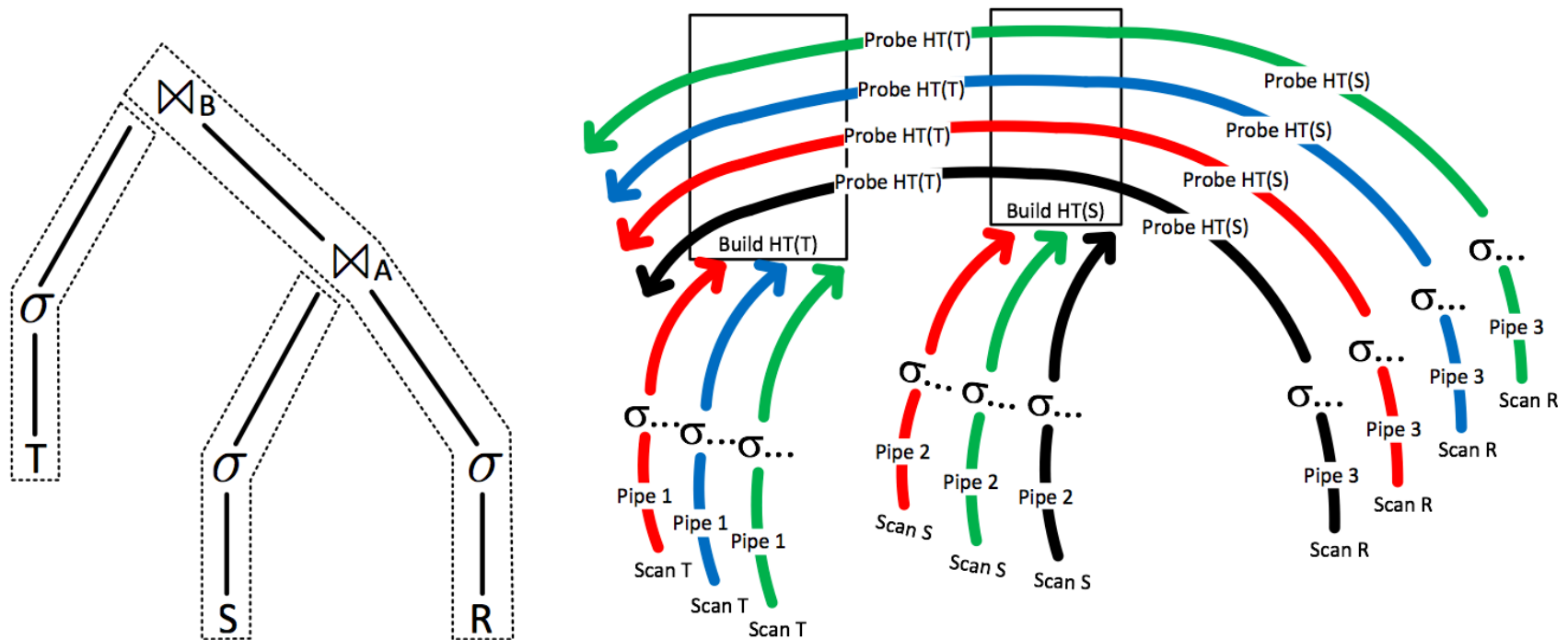
Build performance using different synchronization mechanisms



**Figure 1: Idea of morsel-driven parallelism:  $R \bowtie_A S \bowtie_B T$**



# Massively parallel processing



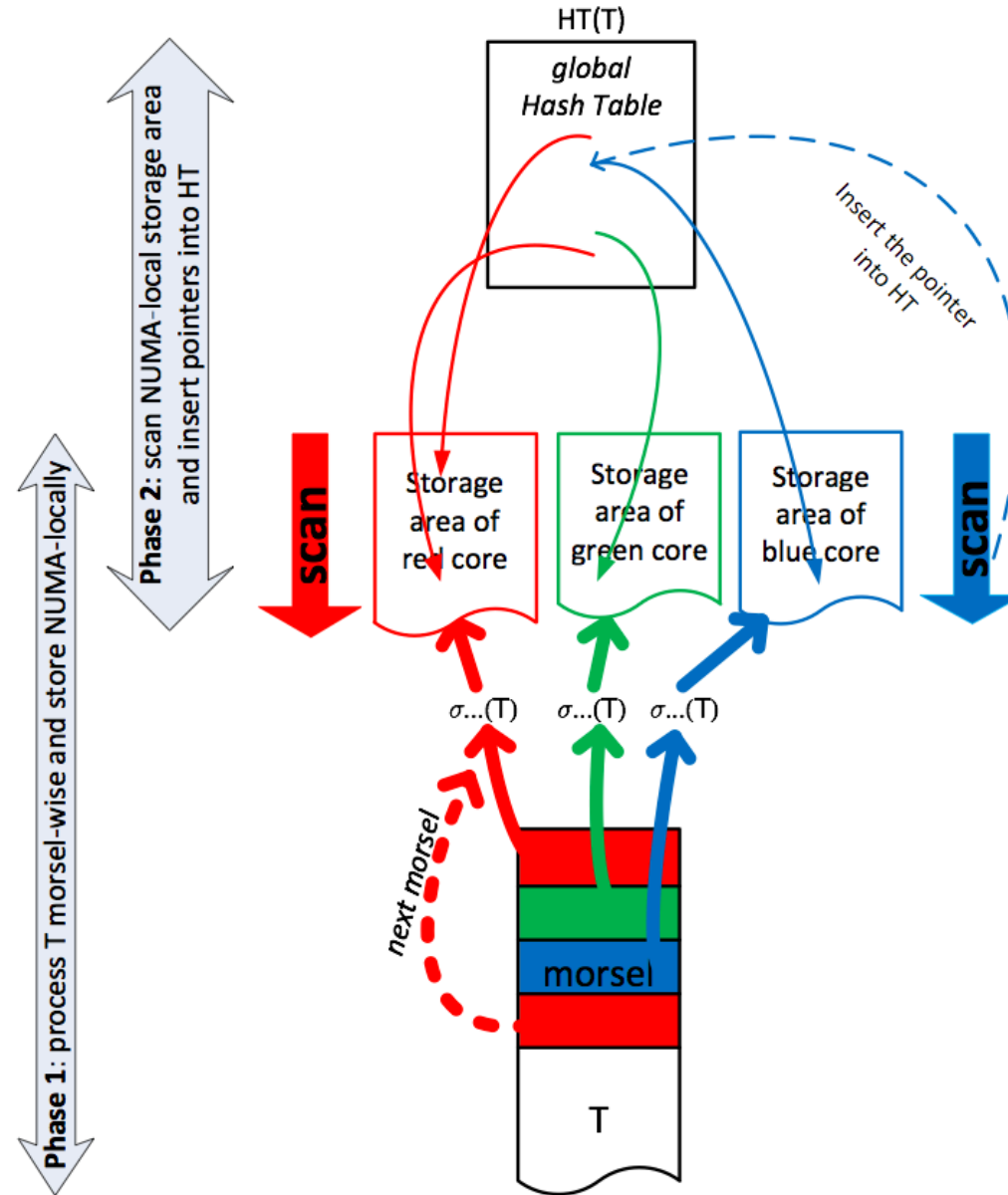
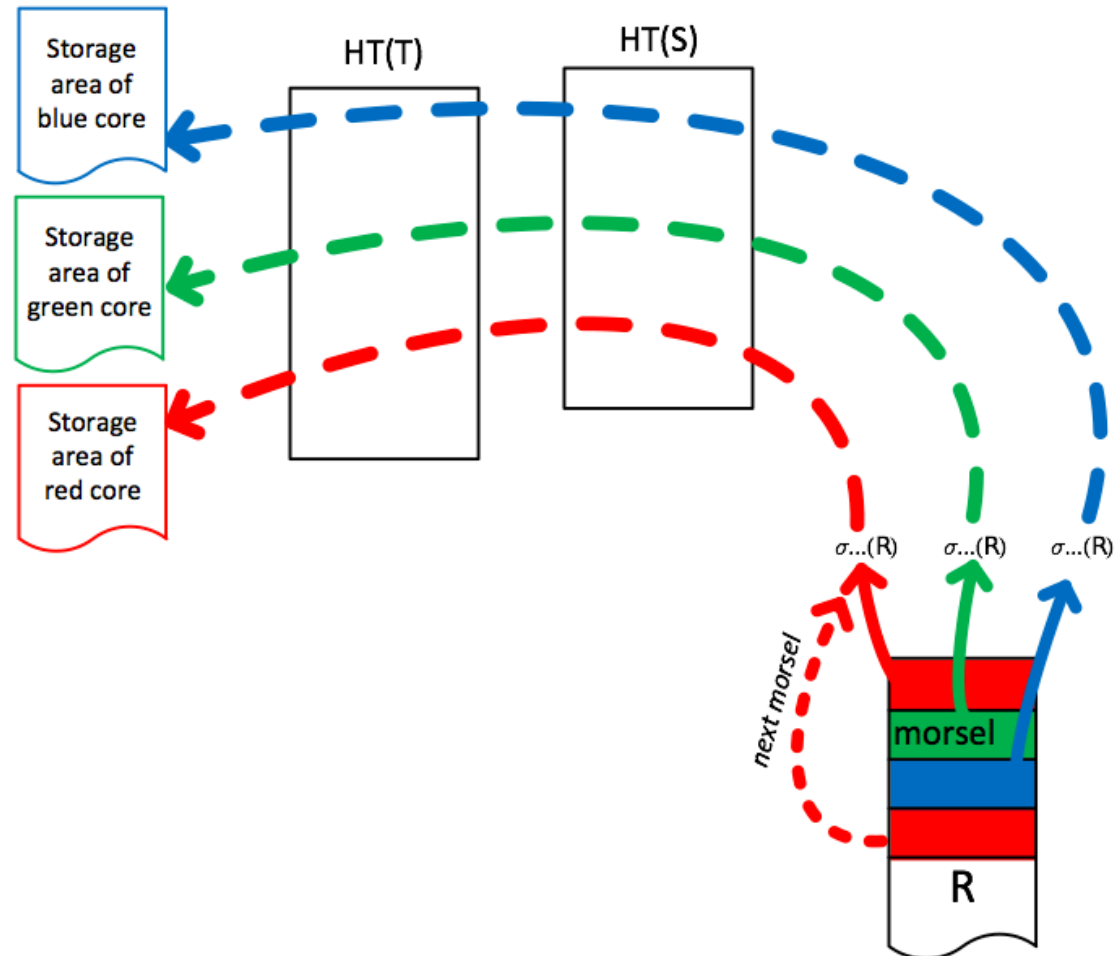


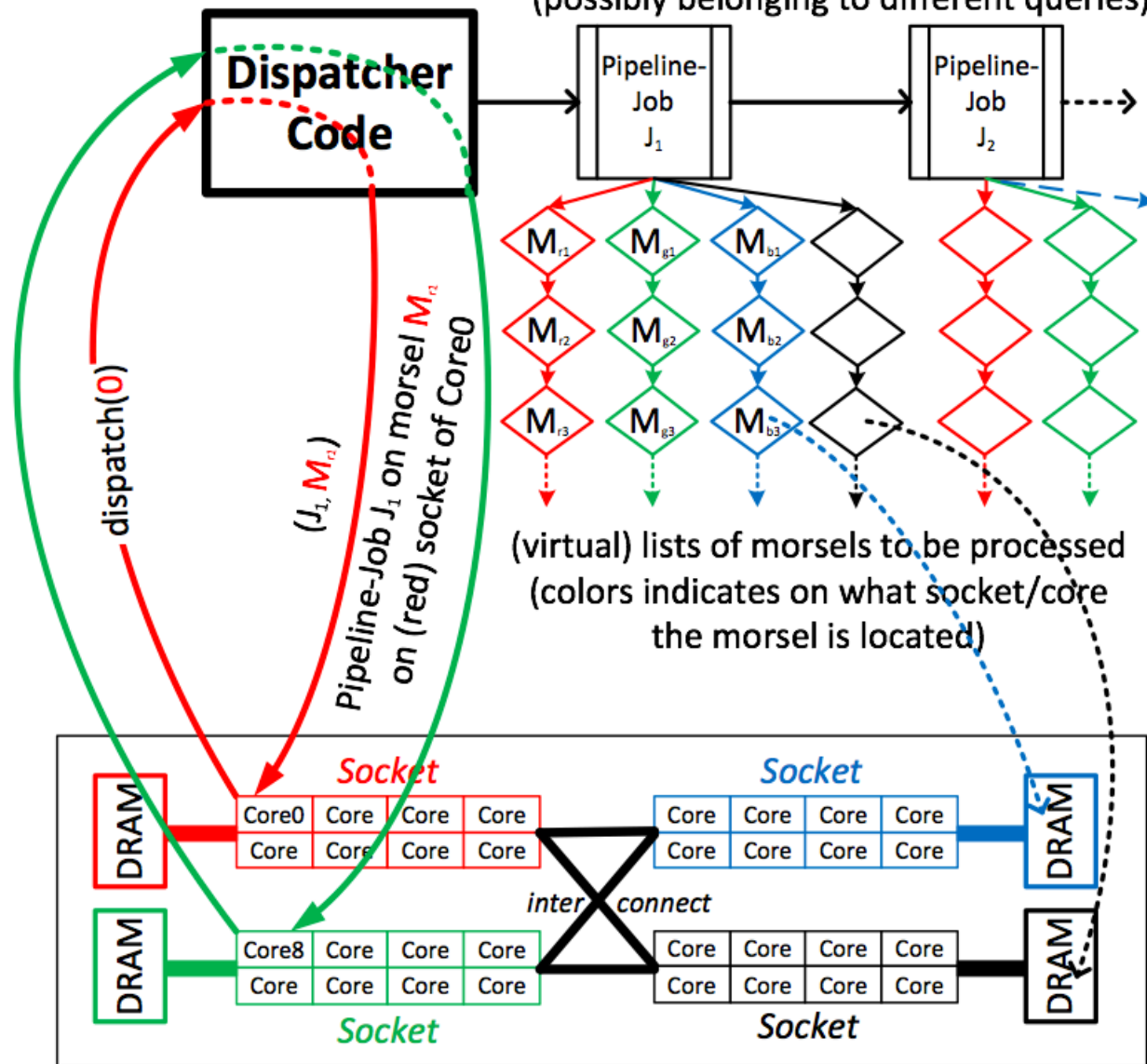
Figure 3: NUMA-aware processing of the build-phase



**Figure 4: Morsel-wise processing of the probe phase**

### Lock-free Data Structures of Dispatcher

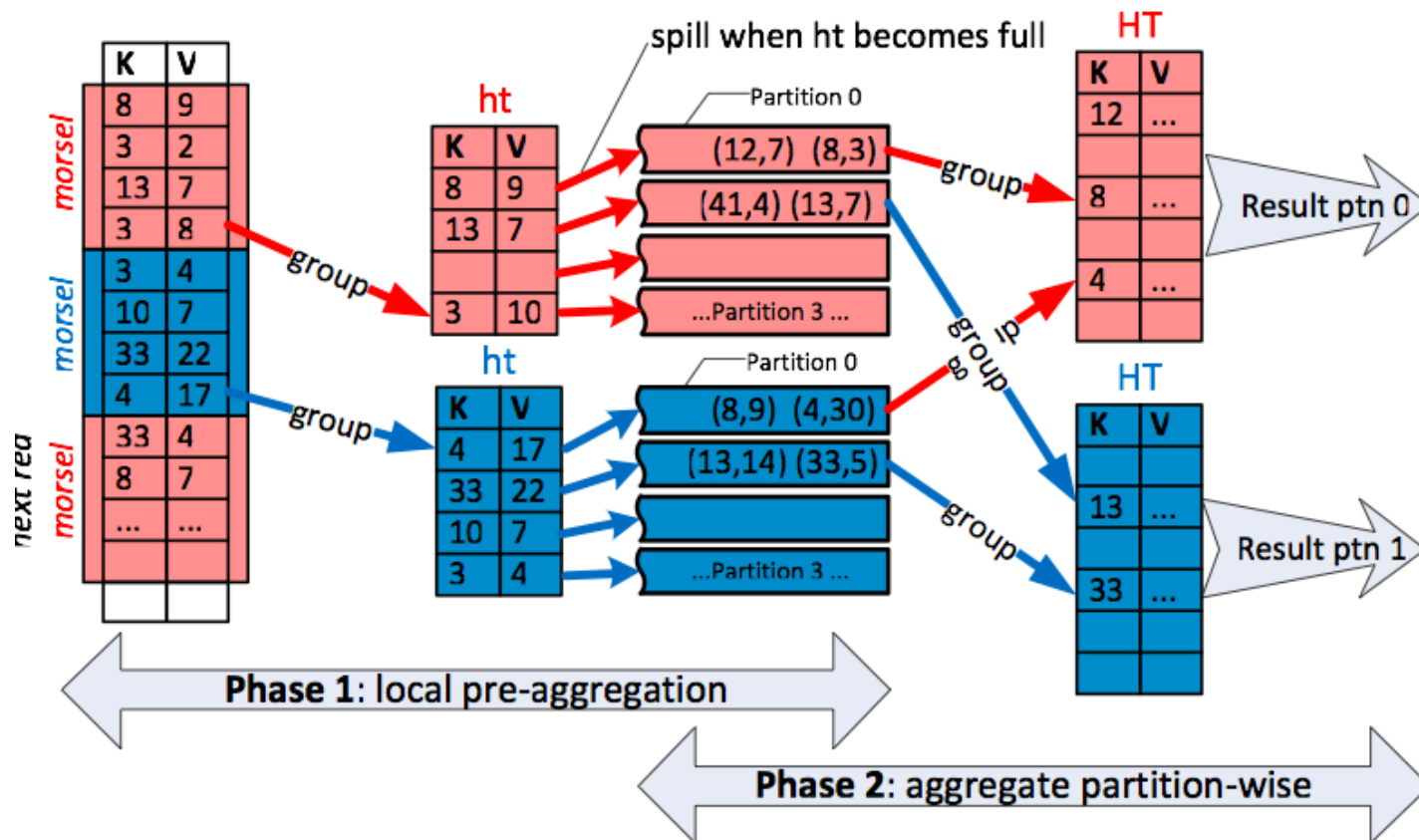
List of pending pipeline-jobs  
(possibly belonging to different queries)



Example NUMA Multi-Core Server with 4 Sockets and 32 Cores



# Massively Parallel Hash Aggregation



**Figure 8: Parallel aggregation**

# Sorting in Parallel

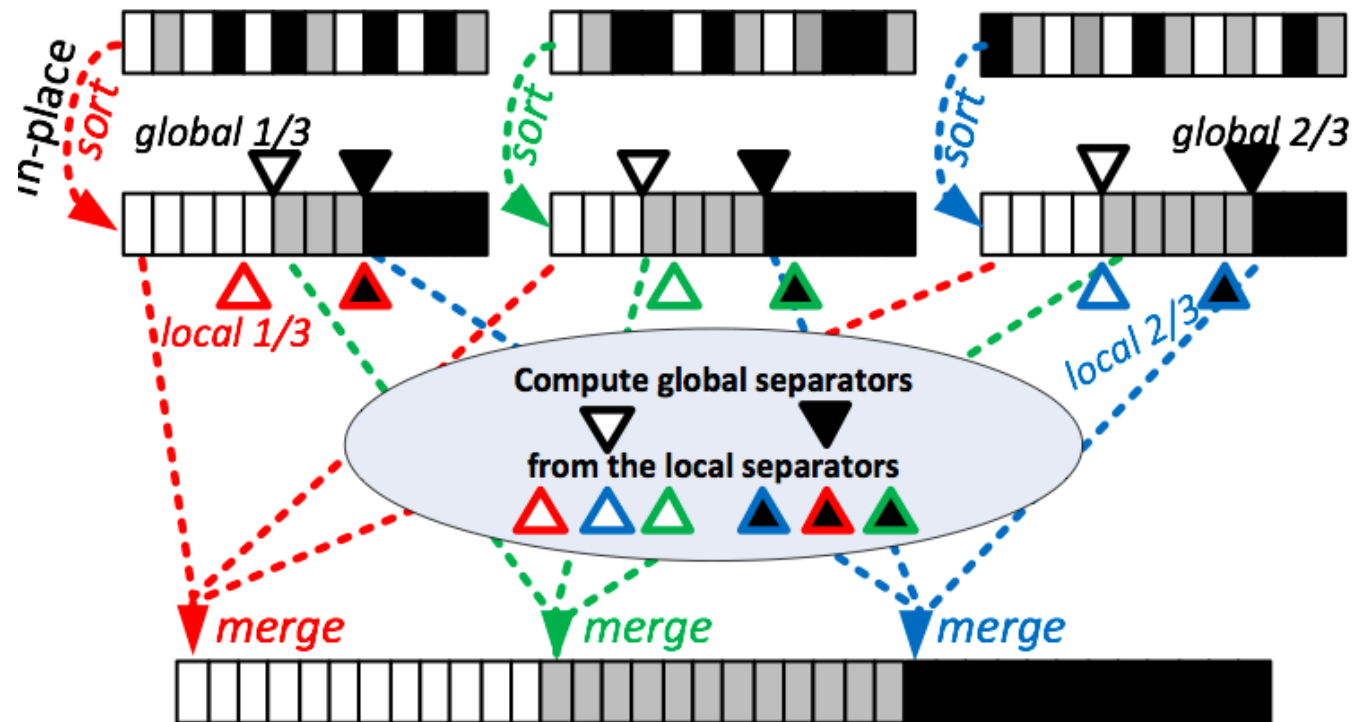


Figure 9: Parallel merge sort



# Massively Parallel Sort-Merge Joins (MPSM) in Main Memory Multi-Core Database Systems

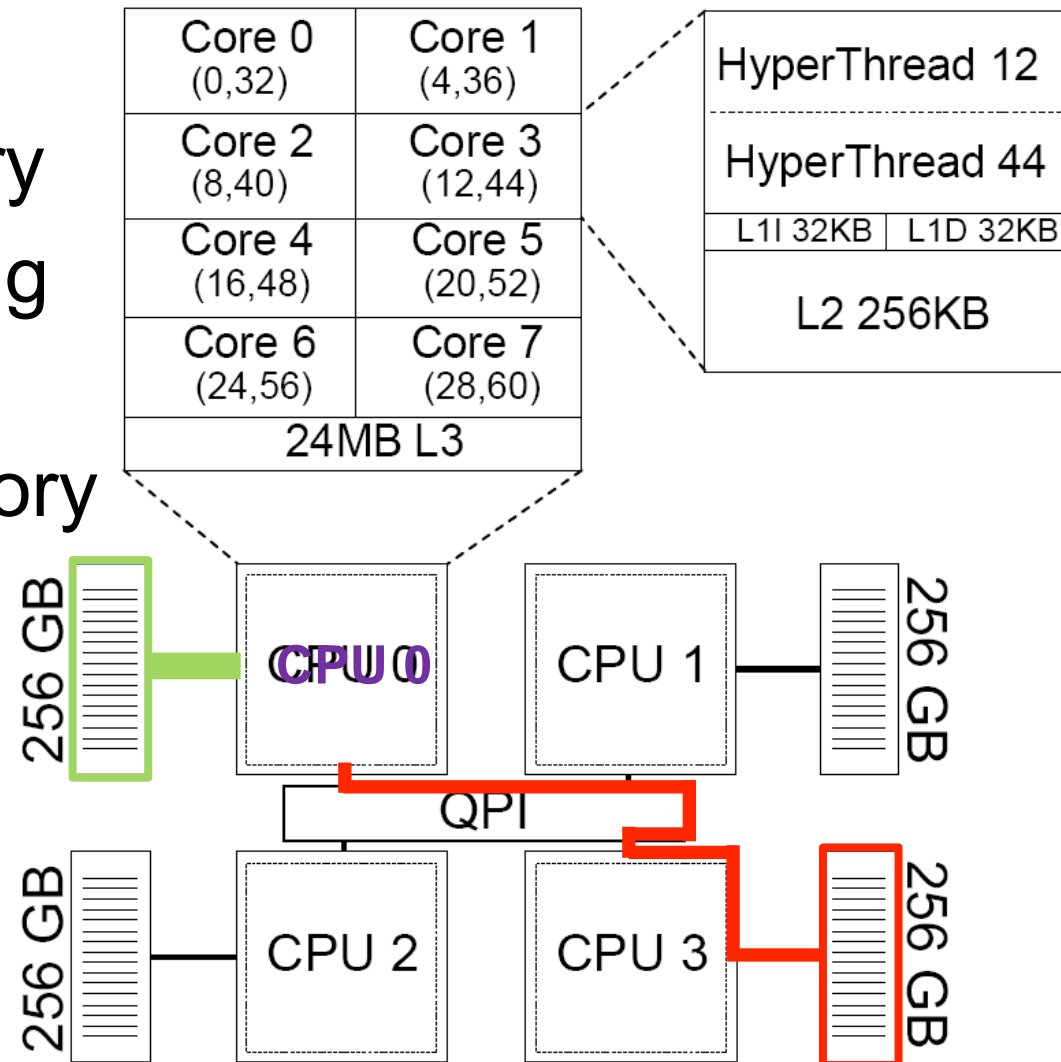
---

Martina Albutiu, Alfons Kemper, and  
Thomas Neumann

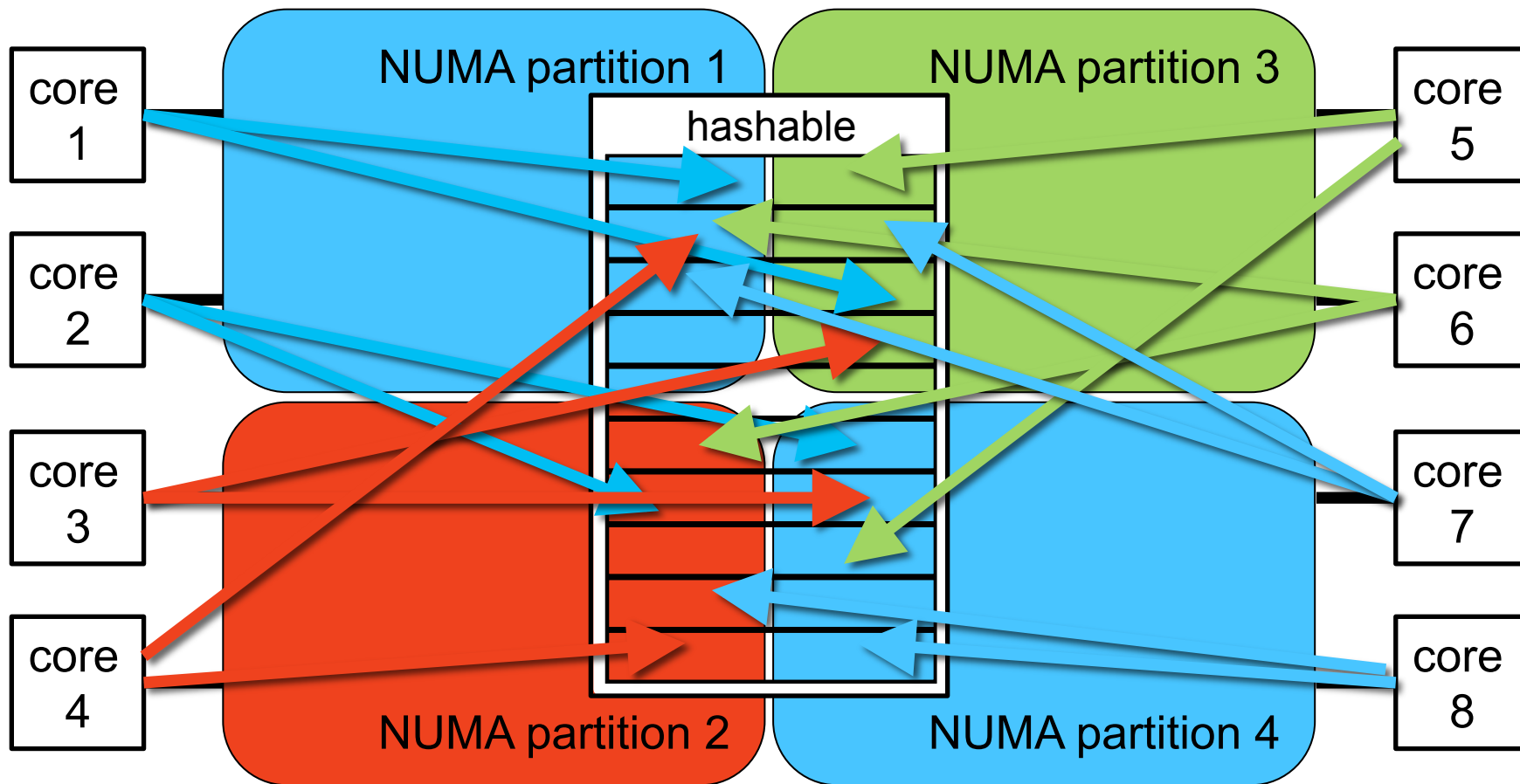
Technische Universität München

## Hardware trends ...

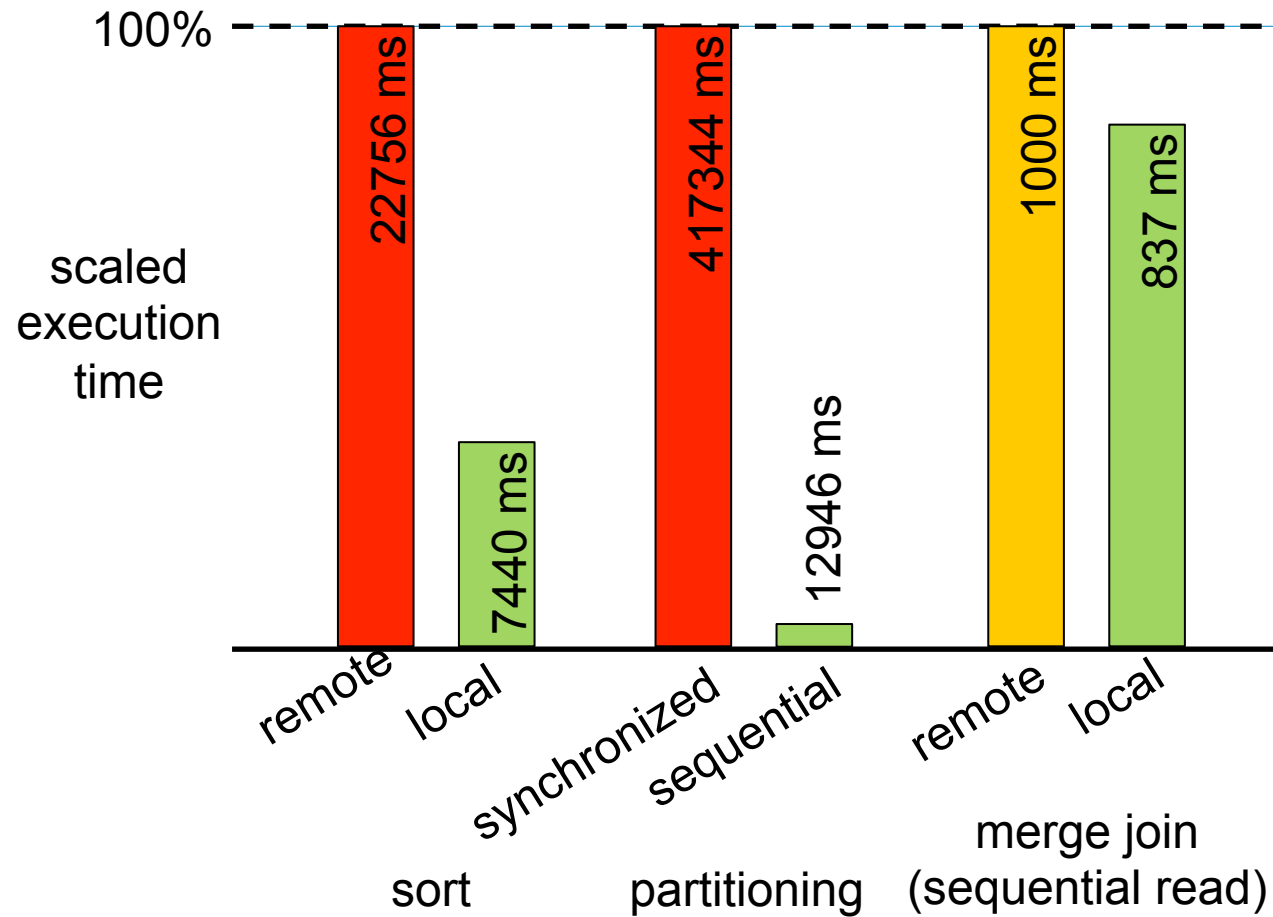
- Huge main memory
- Massive processing parallelism
- Non-uniform Memory Access (NUMA)
- Our server:
  - 4 CPUs
  - 32 cores
  - 1 TB RAM
  - 4 NUMA partitions



# Ignoring NUMA



# How much difference does NUMA make?





## The three NUMA commandments

C1

*Thou shalt not write thy neighbor's memory randomly*  
-- chunk the data, redistribute, and then sort/

C2

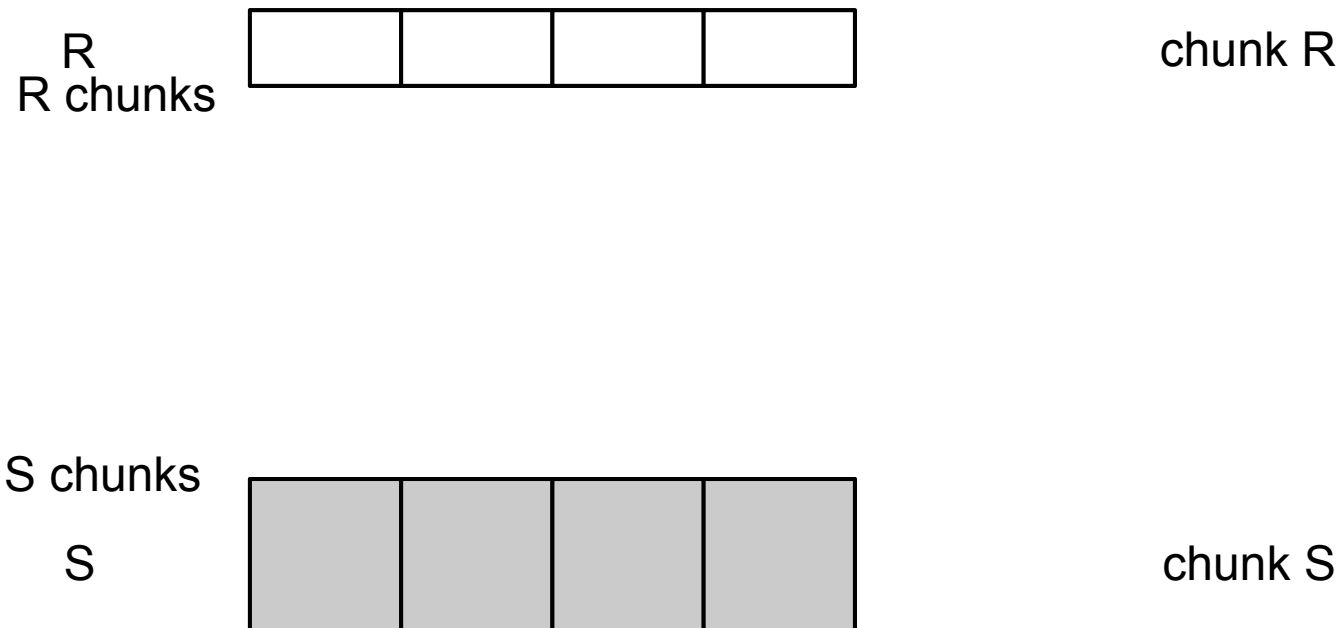
*Thou shalt read thy neighbor's memory only sequentially*  
-- let the prefetcher hide the remote access latency.

C3

*Thou shalt not wait for thy neighbors*  
-- don't use fine-grained latching or locking and avoid synchronization points of parallel threads.



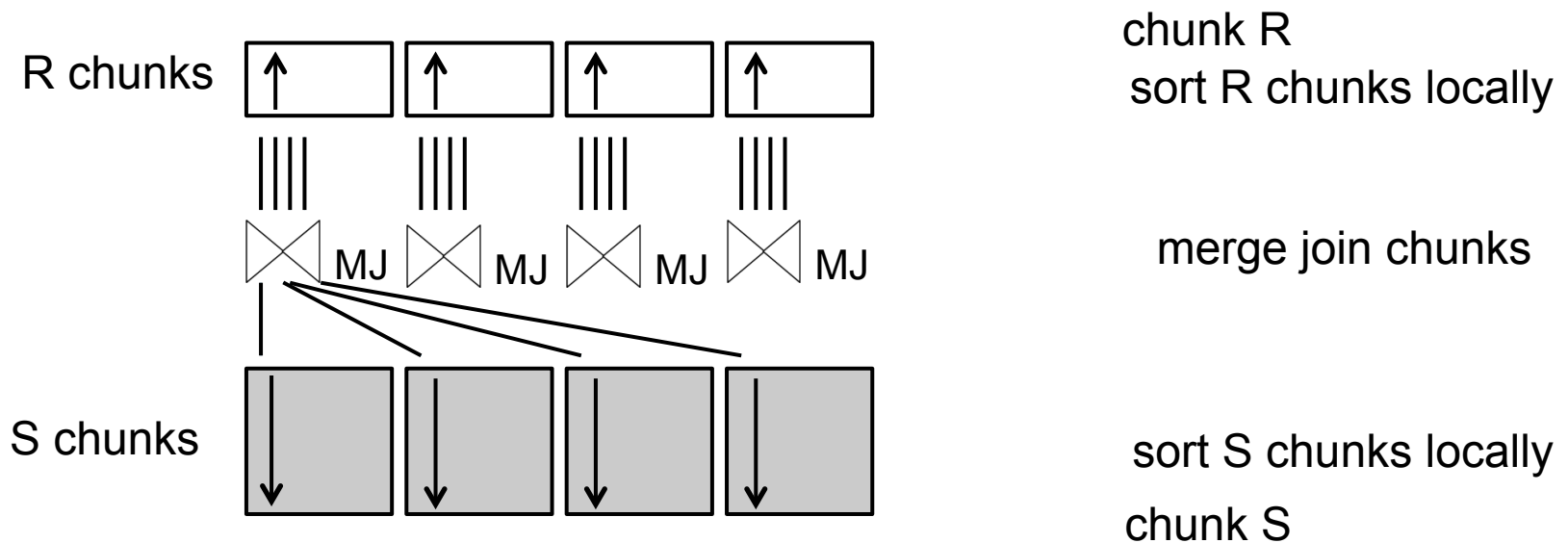
# Basic idea of MPSM





## Basic idea of MPSM

- C1: Work locally: sort
- C3: Work independently: sort and merge join
- C2: Access neighbor's data only sequentially



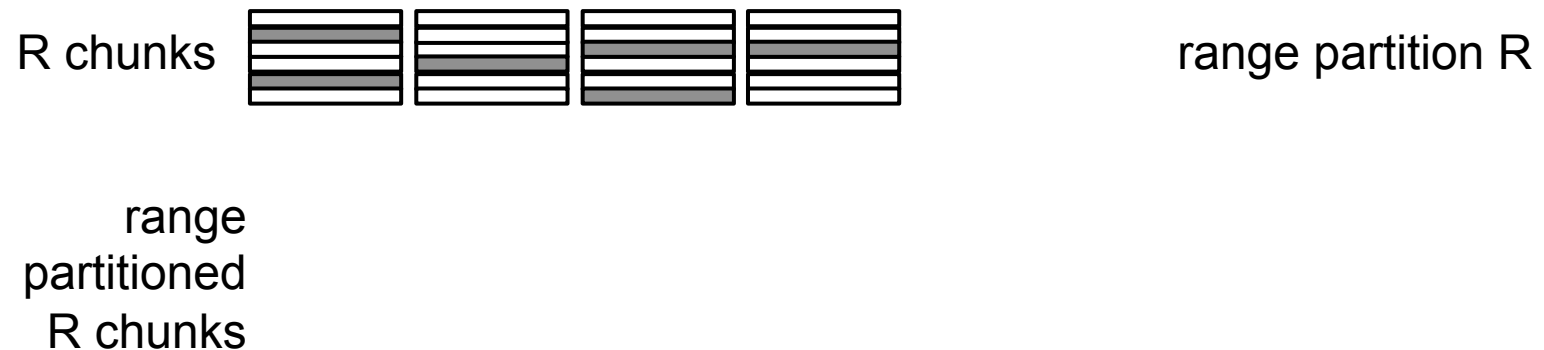


## Range partitioning of private input R

- To constrain merge join work
- To provide scalability in the number of parallel workers

## Range partitioning of private input R

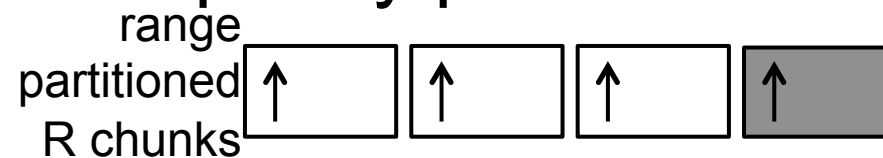
- To constrain merge join work
- To provide scalability in the number of parallel workers



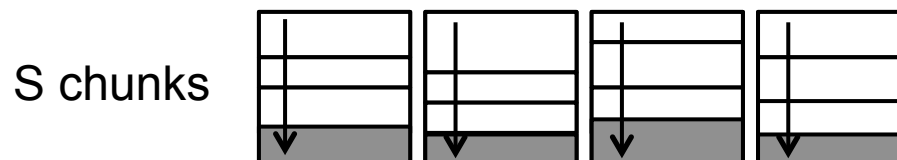
## Range partitioning of private input R

- To constrain merge join work
- To provide scalability in the number of parallel workers

→ S is implicitly partitioned



sort R chunks

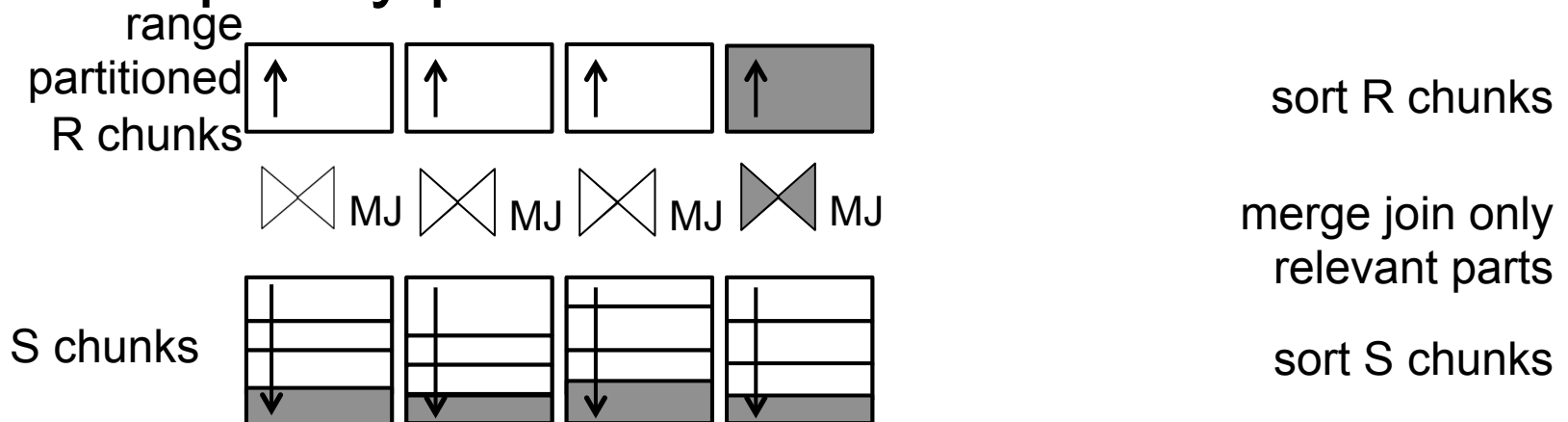


sort S chunks

## Range partitioning of private input R

- To constrain merge join work
- To provide scalability in the number of parallel workers

→ S is implicitly partitioned

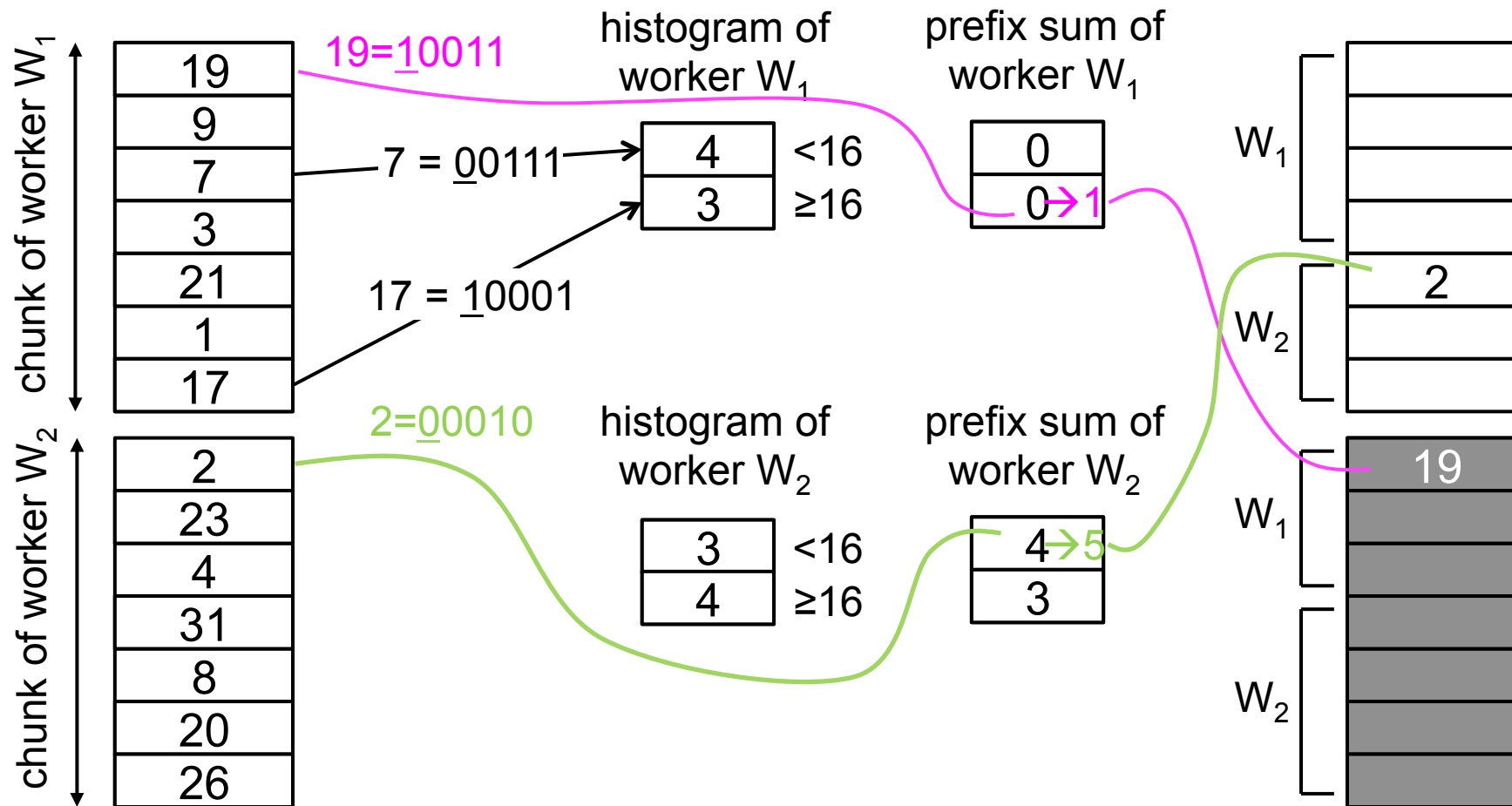


## Range partitioning of private input

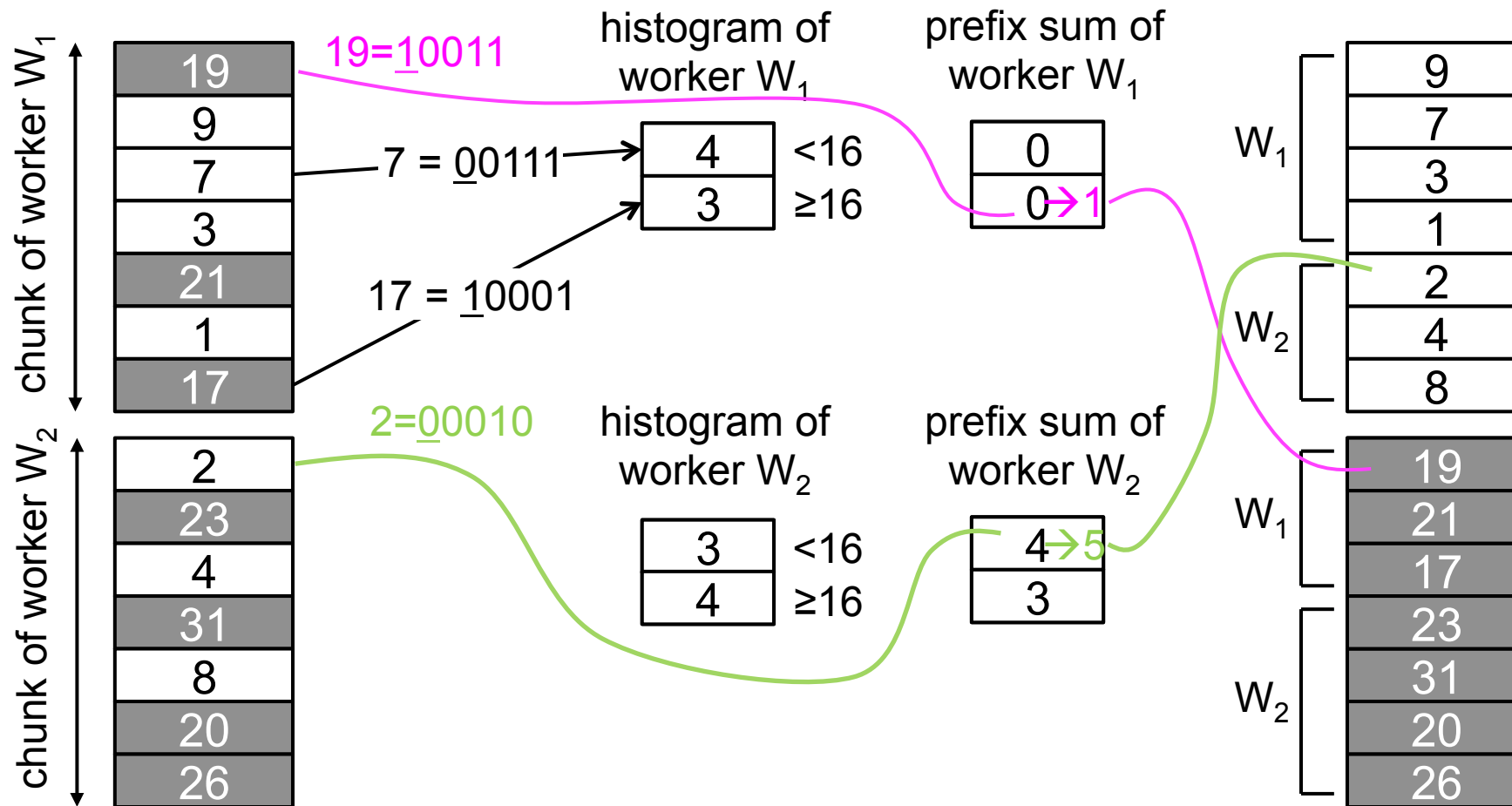
- Time efficient
  - branch-free
  - comparison-free
  - synchronization-freeand
- Space efficient
  - densely packed
  - in-place

by using radix-clustering and precomputed target partitions to scatter data to

# Range partitioning of private input



# Range partitioning of private input







# Real C hacker at work ...

$$ps_i[j] = \&R_j \left[ \left( \sum_{k=1}^{i-1} h_k[j] \right) \right]$$

```
memcpy(ps_i[sp[t.key] >> (64 - B)]++, t, t.size)
```