

Parallelize GRACE Hash Join

- Parallelize Bucket Decomposition in Phase 1
 - Each relation is partitioned in advance
 - The average number of tuples in each disk $\{R\}/N$ and $\{S\}/N$
 - Selection operations are also executed in advance in parallel
 - Virtually divide connected disk into Read Disk and Write Disk
 - Assign each bucket to each processor
 - Send each tuple by its hash value
- Parallelize Join Operation in Phase 2
 - There is no communication within Phase 2

2019/7/29

Advance Data Engineering (H. Yokota)

258

Pseudo Code for Phase1

- Each PE has two threads:
 - Thread 1:


```
for (j = 1; j ≤ {R/N}; j++) {
  read j-th tuple t and attribute value v in t;
  x = h0(v); /* e.g. h(v) = mod(v, N) */
  send t to PEx
}
```
 - Thread 2

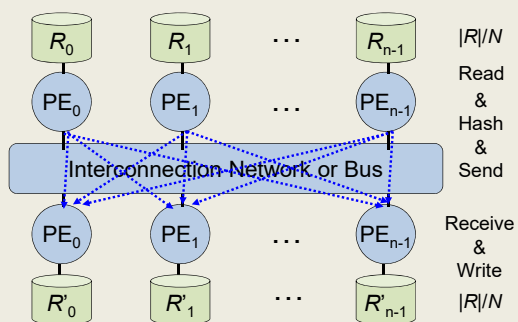

```
for (;;) {
  receive t and attribute value v in t;
  y = h1(v);
  write t into a file for y
}
```
- It should be combined with the phase switch of all-to-all communication

2019/7/29

Advance Data Engineering (H. Yokota)

259

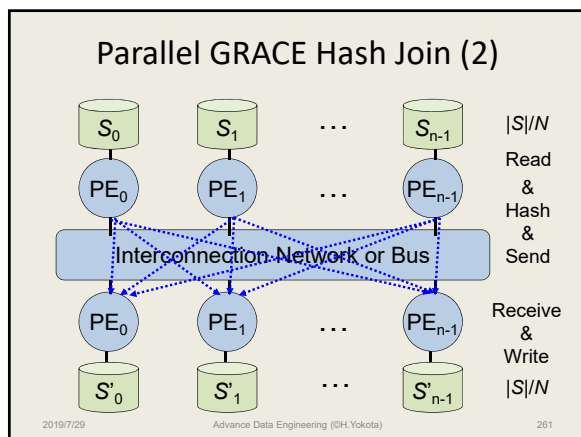
Parallel GRACE Hash Join (1)

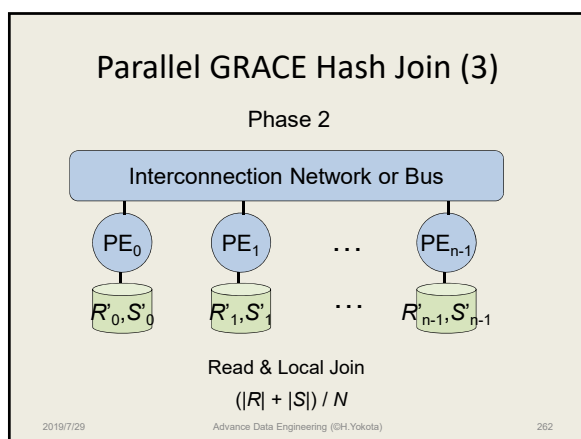


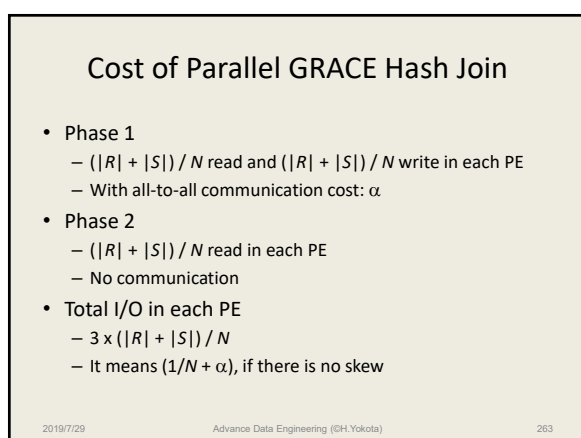
2019/7/29

Advance Data Engineering (H. Yokota)

260







Estimate α (1/3)

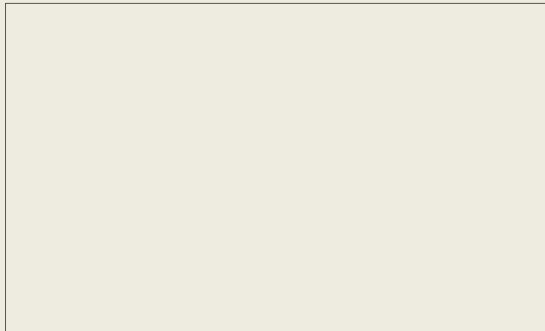
- Bandwidth of each connection of network: 10MB/s
- Network setup time for each connection: 50 μ s
- $|R|$ and $|S|$: 64MB each
- The number of processors (N): 8
- Consider the cost for communication, assuming each processing element has enough large buffer space to keep $|R|/N$ or $|S|/N$
- Also consider the cost when each processing element has memory for two pages (8KB)

2019/7/29

Advance Data Engineering (H.Yokota)

264

Estimate α (2/3)

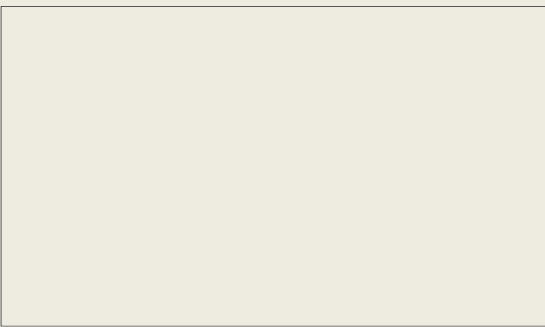


2019/7/29

Advance Data Engineering (H.Yokota)

265

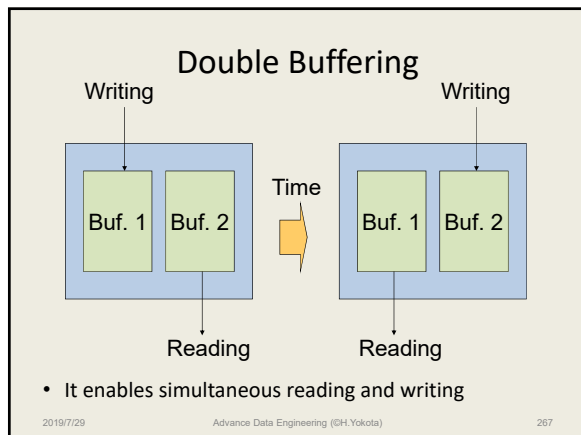
Estimate α (3/3)

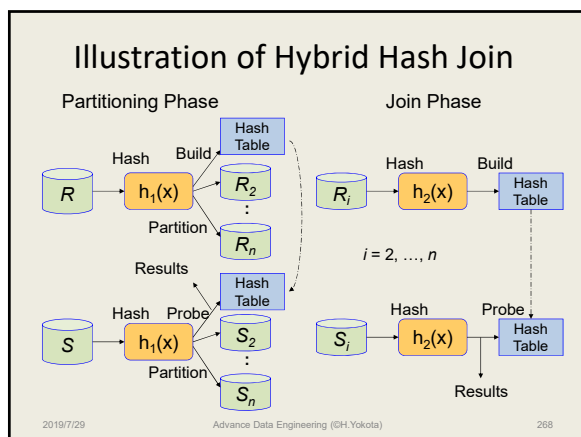


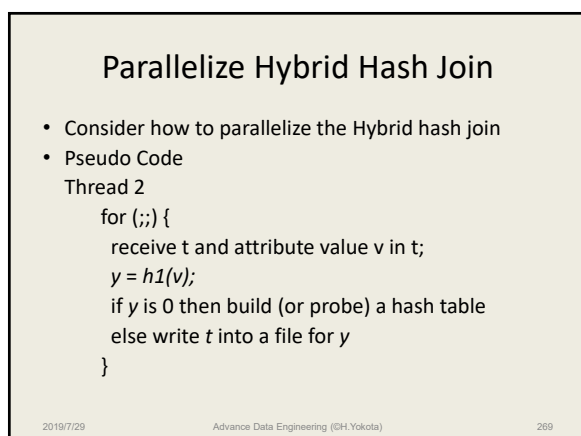
2019/7/29

Advance Data Engineering (H.Yokota)

266







Multiple Joins

- There tend to be a number of join operations in a query
 - The query construct a query tree
 - The configuration of the tree deeply influence the performance of query processing
 - especially under parallel environment
- Parallel executions in a query tree
 - Independent executions
 - Pipeline executions
- Here, we assume hash join algorithm

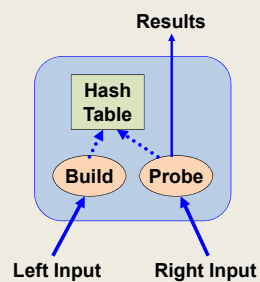
2019/7/29

Advance Data Engineering (H.Yokota)

270

Notation of a Hash Join Node

- Let left input is for **Build** and right for **Probe**
 - Hash Join cannot start Probe until the Build Phase is finished

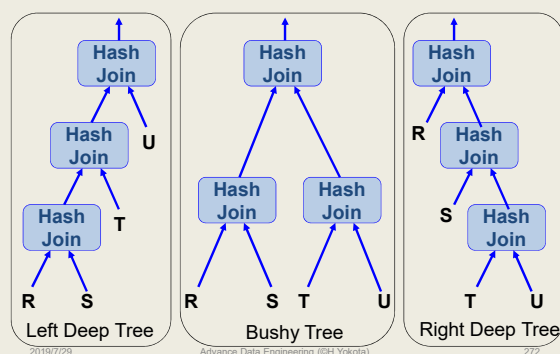


2019/7/29

Advance Data Engineering (H.Yokota)

271

Query Tree Configurations



2019/7/29

Advance Data Engineering (H.Yokota)

272

Parallel Multiple Join

- Parallel execution of multiple join is depend on the structure of the query tree
 - Left Deep Tree:
 - Sequential
 - Right Deep Tree:
 - Pipeline Execution (Parallel Build)
 - Bushy Tree:
 - Child node can be executed in parallel
 - Pipeline Execution can also be done

2019/7/29

Advance Data Engineering (H.Yokota)

273

Parallel Aggregation

- Based on sequential aggregation algorithm using hash
 - Hashing for a Group-By operation
 - Applying aggregation functions to each group
 - Count, Sum, Average, Max, Min
- Three algorithms for parallel execution
 - Centralized Two Phase Algorithm (C-2P)
 - Two-Phase Algorithm (2P)
 - Repartitioning Algorithm (Rep)

2019/7/29

Advance Data Engineering (H.Yokota)

274

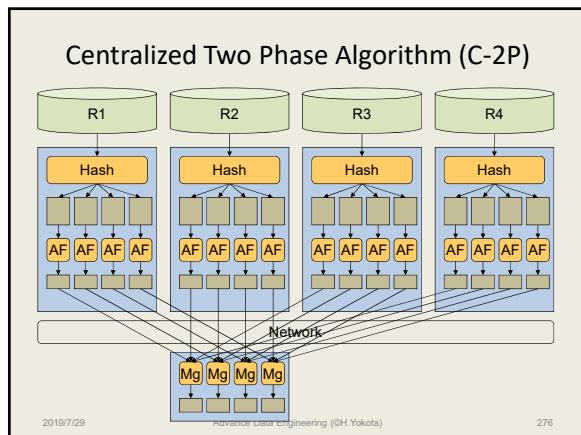
Centralized Two Phase Algorithm (C-2P)

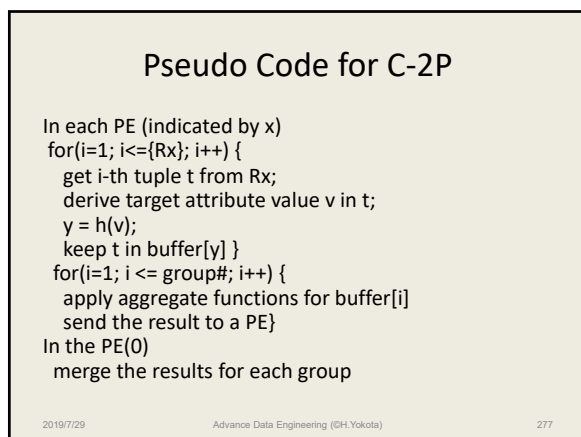
1. Read tuples from each local disk, apply a hash function, and execute aggregate function for hash buckets in each PE
 2. Send the results of aggregation to a **node** to merge them
- The centralized node will be a bottleneck when the number of PE increases

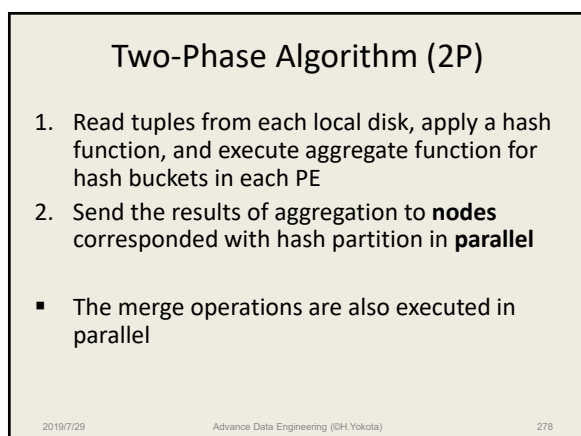
2019/7/29

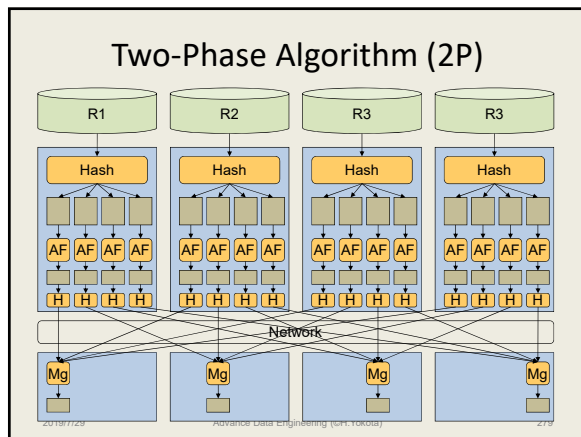
Advance Data Engineering (H.Yokota)

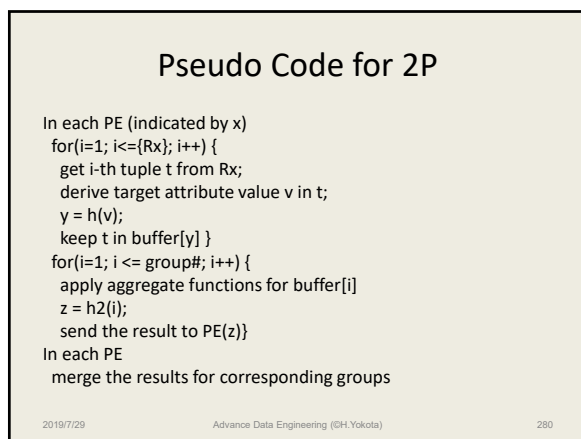
275

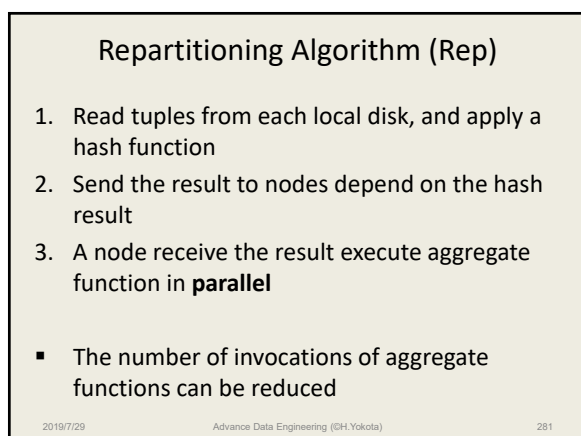


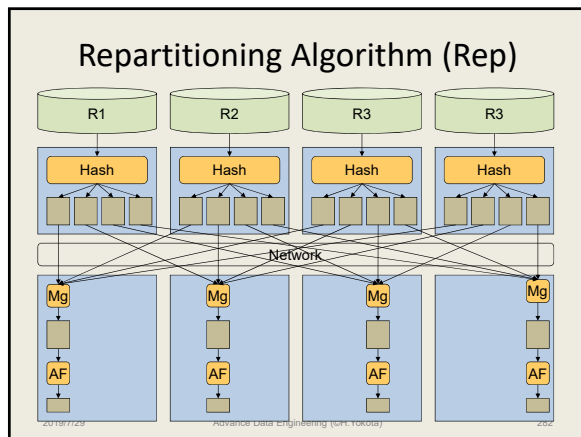












Pseudo Code for Rep (1)

```

In each PE (indicated by x)
for(i=1; i<={Rx}; i++) {
    get i-th tuple t from Rx;
    derive target attribute value v in t;
    y = h2(v);
    send t to PE(y) }

```

2019/7/29

Advanced Data Engineering (H.Yokota)

283

Pseudo Code for Rep (2)

```

In each PE
merge tuples;
n = count tuples;
for(i = 1, i < n ; i++) {
    get i-th tuple from buffer and derive target attribute
    value v in t;
    z = h(v);
    store t in buffer[z];
}
for( j= 1; j < number of buffer in the PE; j++) {
    apply aggregate functions for buffer[j]
}

```

2019/7/29

Advanced Data Engineering (H.Yokota)

284

