

Parallelize Join Operations

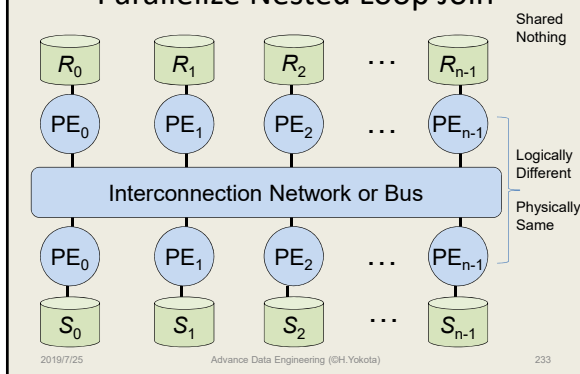
- Remember Join Algorithms
 - Nested Loop Join
 - Sort Merge Join
 - Hash Join
 - Simple Hash Join
 - GRACE Hash Join
 - Hybrid Hash Join
- Here, we consider how to parallelize these algorithms

2019/7/25

Advance Data Engineering (©H.Yokota)

232

Parallelize Nested Loop Join



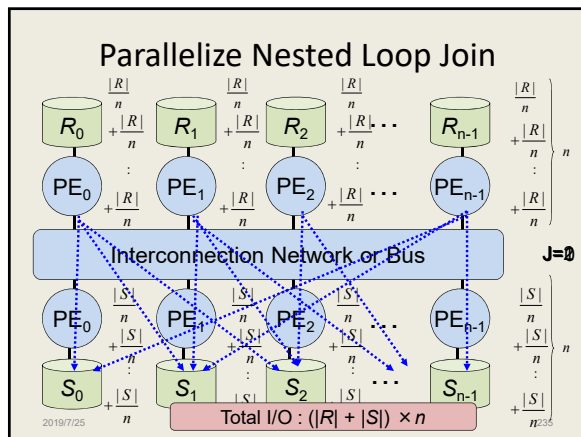
Parallelize Nested Loop Join

- Both relations are horizontally fragmented
 - R_0, R_1, \dots, R_{n-1}
 - S_0, S_1, \dots, S_{n-1}
 - independent of target attribute value
- Place both fragmented relations in each PE
 - S_j and R_i are placed in i -th PE
- Outer Loop: ($0 \leq j \leq n-1$)
 - Send R_j to $\text{mod}(i+j, n)$
 - Read $|R|/n$ page n times in each PE: $|R|$
- Inner Loop:
 - Do Join Operation between S_j and received relation
 - If each PE has enough memory
 - Read $|S|/n$ page n times in each PE: $|S|$
- Total I/O : $(|R| + |S|) \times n$

2019/7/25

Advance Data Engineering

Enlarge size increase total costs
Parallel processing has no effects



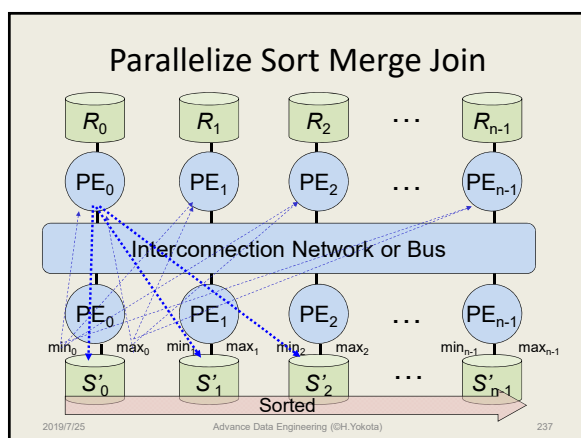
Parallelize Sort Merge Join (1)

- Assumption: Both relations are fragmented
- Type 1
 - Sort one relation in parallel for the target attribute
 - $S'_0, S'_1, \dots, S'_{n-1}$
 - Broadcast maximum and minimum values in S'_i
 - Send each tuple of R_i to a PE correspond to the value
 - Sort all received tuples in the PE
 - Do Sort-Merge Join in each PE

2019/7/25

Advance Data Engineering (©H. Yokota)

236



Parallelize Sort Merge Join (2)

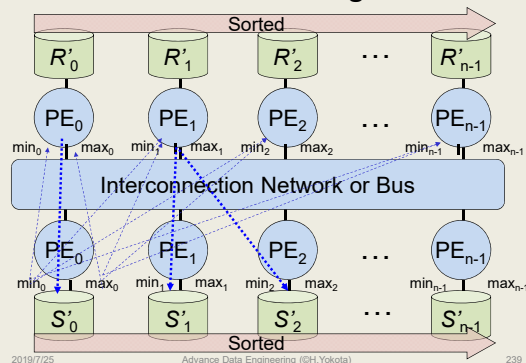
- Type 2
 - Sort both relations in parallel for the target attribute
 - $R'_0, R'_1, \dots, R'_{n-1}$
 - $S'_0, S'_1, \dots, S'_{n-1}$
 - Broadcast maximum and minimum values in S'_i
 - Send each tuple of R'_j to a PE correspond to the value
 - R'_j may be sent to multiple PEs
 - Do Sort-Merge Join in each PE
 - Disk I/O: $|R|/n + |S|/n$, Total Disk I/O $|R| + |S|$
- Parallel Sort Algorithm
 - There are so many parallel sort algorithm

2019/7/25

Advance Data Engineering (©H.Yokota)

238

Parallelize Sort Merge Join

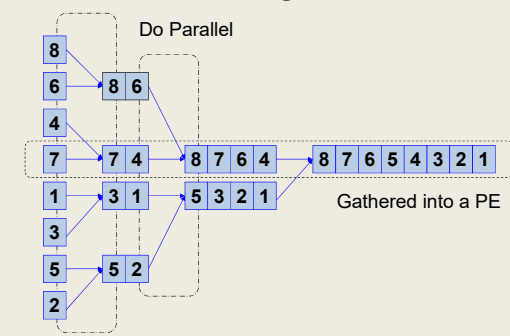


2019/7/25

Advance Data Engineering (©H.Yokota)

239

Parallel Merge Sort (1)



2019/7/25

Advance Data Engineering (©H.Yokota)

240

Parallel Merge Sort (2)

- Construct $\log(\{R\})$ stages for sorting a stream
 - Use $\{R\}/2$ PEs at first
 - Finally, all tuples are gathered into a PE
 - Comparisons in each stage can be done in parallel
- Cost for sorting $\{R\} = 2^m$ tuples
 - Communication paths: $2^m + 2^{m-1} + \dots + 2^1 = 2 \times (\{R\} - 1)$
 - Data transfer: $m \times \{R\} = \{R\} \log_2 \{R\}$
 - Total comparison: $((1+1-1) \times \{R\}/2) + ((2+2-1) \times \{R\}/2^2 + \dots + ((\{R\}/2 + \{R\}/2 - 1) \times \{R\}/2^m) = \{R\} \log_2 \{R\} - (\{R\} - 1)$
 - By parallel processing (time for corresponding): $(2-1) + (2^2-1) + \dots + (2^m-1) = 2(\{R\} - 1) - \log_2 \{R\}$

2019/7/25

Advance Data Engineering (©H.Yokota)

241

Costs for Parallel Block Merge Sort

- For sorting tuples larger than the number of PEs
 - Using $N = 2^n$ PEs, $n \geq 1$ ($N \geq 2$)
- A comparison operation is replaced by a stream merge
- Total Data transfer: $n \times \{R\} = \log_2 N \times \{R\}$
 - By parallel processing (time for corresponding)

$$(1/N + 2/N + \dots + 1/2) \times \{R\}$$

$$= (1 - (1/2)^{n+1}) / (1 - 1/2) \times \{R\}$$

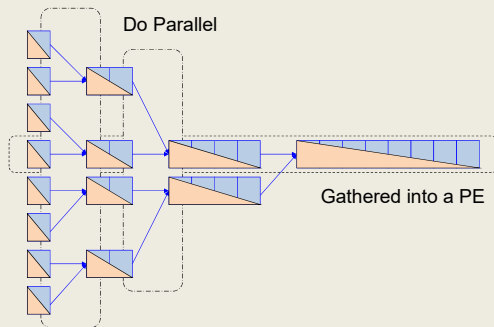
$$= (2 - 1/N) \times \{R\} \quad (\text{how many tuples are transferred})$$

2019/7/25

Advance Data Engineering (©H.Yokota)

242

Parallel Block Merge Sort



2019/7/25

Advance Data Engineering (©H.Yokota)

243

Parallel Block Merge Sort (2)

- If disks are used in each stage
 - Total I/O : $2n \times |R| = 2 \log_2 N \times |R|$
 - the i -th stage: read $|R|/2^{n-(i-1)}$, write $|R|/2^{n-i}$
 - By parallel processing (time for corresponding)

$$\sum_{i=1}^n |R|/2^i + \sum_{i=1}^n |R|/2^{i-1}$$

$$= ((1-(1/2)^n)/(1-1/2)) \times |R| + 2((1-(1/2)^n)/(1-1/2)) \times |R|$$

$$= 6 \times (1 - 1/N) \times |R|$$
 - Expect Pipeline effect

2019/7/25

Advance Data Engineering (©H.Yokota)

244

Bitonic Sort (1)

- Bitonic Sequence
 - There exist $1 \leq j \leq 2n$ which satisfies
 - $a_1 \leq a_2 \leq \dots \leq a_j \geq a_{j+1} \geq \dots \geq a_{2n}$
- For a bitonic sequence a_1, a_2, \dots, a_{2n}
 - Let $d_i = \min(a_i, a_{n+i})$ and $e_i = \max(a_i, a_{n+i})$ where $1 \leq i \leq n$
 - Then d_1, d_2, \dots, d_n and e_1, e_2, \dots, e_n are also bitonic sequences
 - And $\max(d_1, d_2, \dots, d_n) \leq \min(e_1, e_2, \dots, e_n)$

2019/7/25

Advance Data Engineering (©H.Yokota)

245

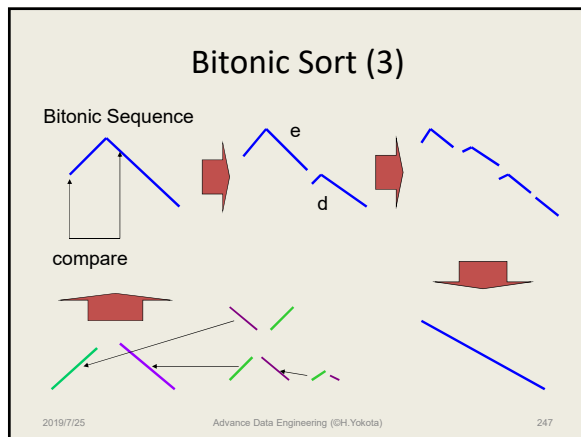
Bitonic Sort (2)

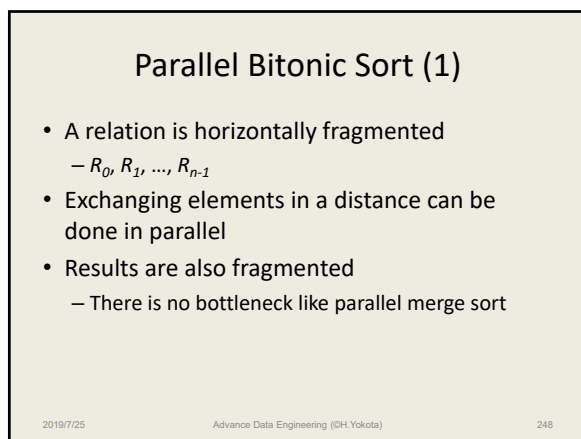
- If there is a bitonic sequence, ordered small (half size) bitonic sequence can be generated by exchanging elements in a distance
 - Finally, the small bitonic sequence becomes an element
 - That is the result of the sort operation
- A Problem: How to generate the first bitonic sequence
 - Answer: Concatenation of two half size sorted sequence
 - Recursively continue until an element

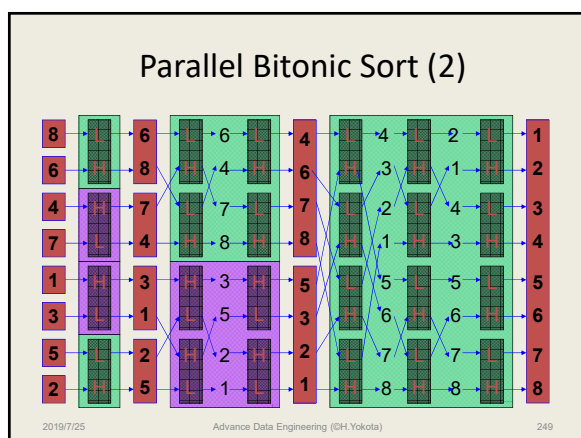
2019/7/25

Advance Data Engineering (©H.Yokota)

246







Cost for Parallel Bitonic Sort

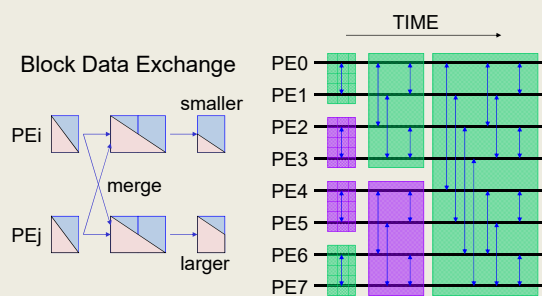
- The number of stages for generating bitonic sequence:
 - $\log_2 \{R\} = m$
- The number of stages for 2^i length bitonic sequence:
 - $\log_2 2^i = i$
- Total stages:
 - $\sum_{i=1}^m \log_2 2^i = \sum_{i=1}^m i = m(m+1)/2 = (\log_2 \{R\}(\log_2 \{R\} + 1)) / 2$
- Total data transfer:
 - $\{R\} \times ((\log_2 \{R\}(\log_2 \{R\} + 1)) / 2)$
- Total data exchanges:
 - $(1/2) \times \{R\} \times ((\log_2 \{R\}(\log_2 \{R\} + 1)) / 2)$
- By parallel processing (time for corresponding):
 - $(\log_2 \{R\}(\log_2 \{R\} + 1)) / 2$

2019/7/25

Advance Data Engineering (©H.Yokota)

250

Parallel Block Bitonic Sort



2019/7/25

Advance Data Engineering (©H.Yokota)

251

Cost for Parallel Block Bitonic Sort

- For sorting tuples larger than the number of PEs ($N = 2^n$)
- Merge two streams and leave smaller (larger) part
- Total stages: $\sum_{i=1}^n i = n(n+1) / 2$
- Total data transfer: $\{R\} \times n(n+1) / 2$
- By parallel processing (time for corresponding):
 - $\{R\} \times n(n+1) / 2N$
- Total I/O (read and write for each stage) :
 - $|R| \times n(n+1)$
- By parallel processing (time for corresponding):
 - $|R| / N \times n(n+1) = (\log_2 N(\log_2 N + 1)) / N \times |R|$

2019/7/25

Advance Data Engineering (©H.Yokota)

252