



MONASH University

程序代写代做 CS 编程辅导

Information Technology

FIT5202 (Volume III - Join)



Week 3a – Parallel Join

WeChat: cstutorcs

Assignment Project Exam Help

**algorithm distributed systems database**

Email: tutorcs@163.com

systems **computation knowledge management**

QQ: 749389476

**design e-business model data mining intelligent**

**distributed systems database software engineering**

<https://tutorcs.com>

**computation knowledge management analysis**

程序代写代做 CS编程辅导

# Chapter 3 Parallel Search

TANIA  
LEUNG  
RAHAYU  
GOEL

Wiley Series on Parallel and Distributed Computi



High Performance  
Database Processing  
Grid Databases

High Performance Parallel Database  
Processing and Grid Databases

WeChat: cstutorcs

Last Week

Email: tutorcs@163.com

QQ: 749389476

DAVID TANIA, CLEMENT H.C. LEUNG,  
WENNY RAHAYU, and SUSHANT GOEL

<https://tutorcs.com>

- 3.1 Search Queries
- 3.2 Data Partitioning
- 3.3 Search Algorithms
- 3.4 Summary
- 3.5 Bibliographical Notes
- 3.6 Exercises



WILEY



MONASH University

# Revision

程序代写代做 CS编程辅导

## Exercise 1 (FLUX Quiz)

- If a query runs on a multi-core machine(e.g. Windows server with 16 cores), it is called parallel query processing if multiple different queries run at the same time in a multi-core machine?



WeChat: cstutorcs

- A. Intra-query parallelism
- B. Inter-query parallelism
- C. Intra-operation parallelism
- D. Inter-operation parallelism

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Revision

程序代写代做 CS编程辅导

## Exercise 2 (FLUX Quiz)

- If a **hash data partition** is used to store the data, and the query is a **discrete range search**, how many processors need to be used to process such a query efficiently in order to get the query results?

WeChat: cstutorcs

- A. 1 processor
- B. Selected processors
- C. All processors

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



程序代写代做 CS编程辅导

TANIA  
R LEUNG  
RAHAYU  
GOEL

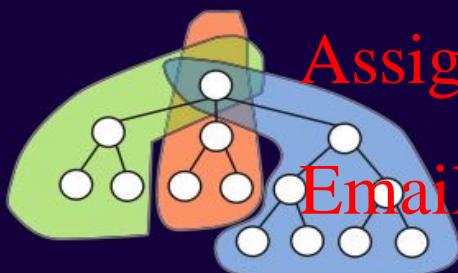
Wiley Series on Parallel and Distributed Computi



High Performance  
Database Processing  
Grid Databases

High Performance Parallel Database  
Processing and Grid Databases

WeChat: cstutorcs



Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

DAVID TANIA, CLEMENT H.C. LEUNG,  
WENNY RAHAYU, and SUSHANT GOEL

<https://tutorcs.com>



WILEY

# Chapter 5

## Parallel Join

- 5.1 Join Operations
- 5.2 Serial Join Algorithms
- 5.3 Parallel Join Algorithms
- 5.4 Cost Models
- 5.5 Parallel Join Optimization
- 5.6 Summary
- 5.7 Bibliographical Notes
- 5.8 Exercises

# 5.1. Join Operations



- Join operations to link tables based on the nominated attributes
  - one from each table

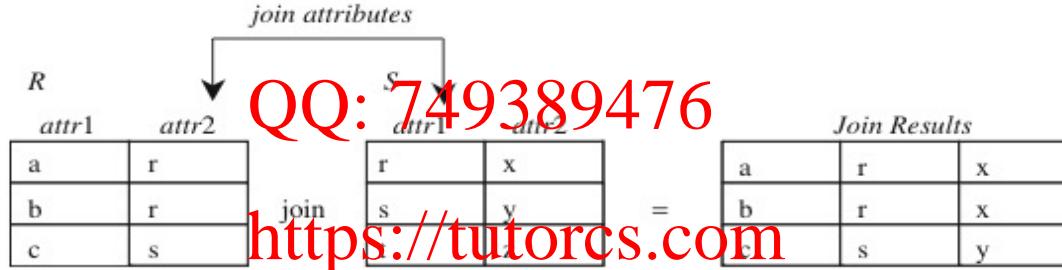
**Query 5.1:**

```
Select *  
From STUDENT S, ENROLMENT E  
Where S.Sid = E.Sid;
```

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com



**Figure 5.1** The join operation

## 5.2. Serial Join Algorithms

- Three serial join algorithms
  - Nested loop join algorithm
  - Sort-merge join algorithm
  - Hash-based join algorithm



WeChat: cstutorcs

Table R	
Adele	8
Bob	22
Clement	16
Dave	23
Ed	11
Fung	22
Goel	3
Harry	17
Irene	14
Joanna	2
Kelly	6
Lim	20
Meng	1
Noor	5
Omar	19

Table S	
Arts	8
Business	15
CompSc	2
Dance	12
Engineering	7
Finance	20
Geology	10
Health	11
IT	18

Join Results		
Adele	8	Arts
Ed	11	Health
Joanna	2	CompSc

Assignment Project Exam Help

QQ: 749389476

<https://tutorcs.com>

Figure 5.2 Sample data

## 5.2. Serial Join Algorithms (cont'd)



### Nested-Loop Join

- For each record of table R, it goes through all records of table S

Table R		Table S		Join Results	
Adele	8	Arts	8	Adele	8
Bob	22	Business	15	Ed	11
Clement	16	Computer	22	Joanna	2
Dave	23	Dance	12		
Ed	11	Engineering	7		
Fung	25	Finance	21		
Goel	3	Geology	10		
Harry	17	Health	11		
Irene	14	IT	18		
Joanna	2				
Kelly	6				
Lim	20				
Meng	1				
Noor	5				
Omar	19				

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com  
QQ: 749389476  
<https://tutorcs.com>

## 5.2. Serial Join Algorithms (cont'd)



### • Sort-Merge Join

- Both tables must be pre-sorted based on the join attribute(s). If not, then both tables must be sorted first
- Then merge the two sorted tables

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

## 5.2. Serial Join Algorithms (cont'd)

程序代写代做CS编程辅导

Table R	
Meng	1
Joanna	2
Goel	3
Noor	5
Kelly	6
Adele	8
Ed	11
Irene	14
Clement	16
Harry	17
Omar	19
Lim	20
Bob	22
Dave	23
Fung	25



Major	Score
Arts	7
Geology	10
Health	11
Dance	12
Business	15
Finance	21

Join Results		
Joanna	2	CompSc
Adele	8	Arts
Ed	11	Health

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Figure 5.4. Sorted tables

## 5.2. Serial Join Algorithms (cont'd)



Table R

Meng	1						
Joanna	2						
Goel	3						
Noor	5						
Kelly	6						
Adele	8						
Ed	11						
Irene	14						
Clement	16						
Harry	17						
Omar	19						
Lim	20						
Bob	22						
Dave	23						
Fung	25						

Join Results

Joanna	2	CompSc
Adele	8	Arts
Ed	11	Health

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Figure 5.4. Sorted tables

## 5.2. Serial Join Algorithms (cont'd)



### • Hash-based Join

- The records of files  $R$  and  $S$  are both hashed to the *same hash file*, using the *same hashing function* on the join attributes  $A$  of  $R$  and  $B$  of  $S$  as hash keys
- A single pass through the file with fewer records (say,  $R$ ) hashes its records to the hash file buckets
- A single pass through the other file ( $S$ ) then hashes each of its records to the appropriate bucket, where the record is combined with all matching records from  $R$

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

## 5.2. Serial Join Algorithms (cont'd)



Table S	
Arts	8
Business	15
CompSc	2
Dance	12
Engineering	7
Finance	21
Geology	10
Health	11
IT	18

x	Entries
1	Geology/10
2	CompSc/2
3	Health/11
4	Dance/12
5	Finance/21
6	Business/15
7	Engineering/7
8	Arts/8
9	IT/18
10	
11	
12	

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Figure 5.6. Hashing Table S

## 5.2. Serial Join Algorithms (cont'd)



**Table R**

Adele	8
Bob	22
Clement	16
Dave	23
Ed	11
Fung	25
Goel	3
Harry	17
Irene	14
Joanna	2
Kelly	6
Lim	20
Meng	1
Noor	5
Omar	19

WeChat: cstutors  
Assignment Project Exam Help  
Email: tutorcs@163.com  
QQ: 749389476

1	Geology/10
2	CompSc/2
3	Dance/12
4	Finance/21
5	
6	Business/15
7	Engineering/7
8	Arts/8
9	IT/18
10	
11	
12	

**Join Results**

Adele	8	Arts
Ed	11	Health
Joanna	2	CompSc

<https://tutorcs.com>

Figure 5.7. Probing Table R

## 5.3. Parallel Join Algorithms



- Parallelism of join is achieved through *data parallelism*, whereby the same task is applied to different parts of the data
- After data partitioning is completed, each processor will have its own data to work with using any serial join algorithm
- Data partitioning for parallel join algorithms:
  - Divide and broadcast
  - Disjoint data partitioning

WeChat: csutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

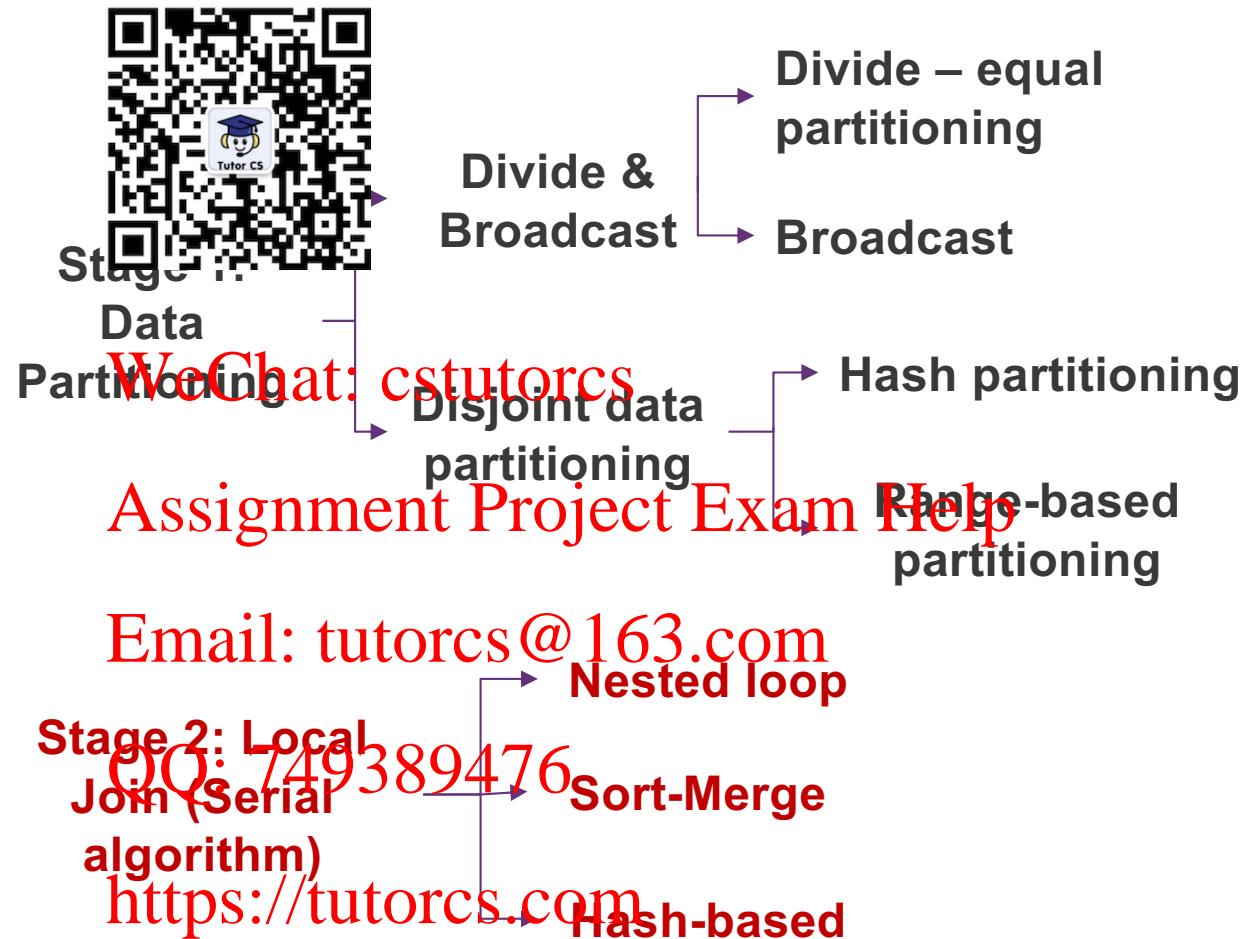
QQ: 749389476

<https://tutorcs.com>

# Overview

程序代写代做 CS编程辅导

Parallel Join  
Algorithms



## 5.3. Parallel Join Algorithms (cont'd)



### • Divide and Broadcast Parallel Join Algorithms

- Two stages: data partitioning using the divide and broadcast method, and a local join
- Divide and Broadcast method: Divide one table into multiple disjoint partitions, where each partition is allocated a processor, and broadcast the other table to all available processors
- Dividing one table can simply use equal division
- Broadcast means replicate the table to all processors
- Hence, choose the smaller table to broadcast and the larger table to divide

QQ: 749389476

<https://tutorcs.com>

### 5.3. Parallel Join Algorithms (cont'd)



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Processor 1		Processor 2		Processor 3	
R1	S1	R2	S2	R3	S3
Adele	8	Arts	8	Fung	25
Bob	22	Dance	15	Gael	3
Clement	16	Celeste	10	Harry	18
Dave	23			Irene	14
Ed	11			Joanna	2

Figure 5.10 Initial data placement

### 5.3. Parallel Join Algorithms (cont'd)

程序代写代做 CS 编程辅导

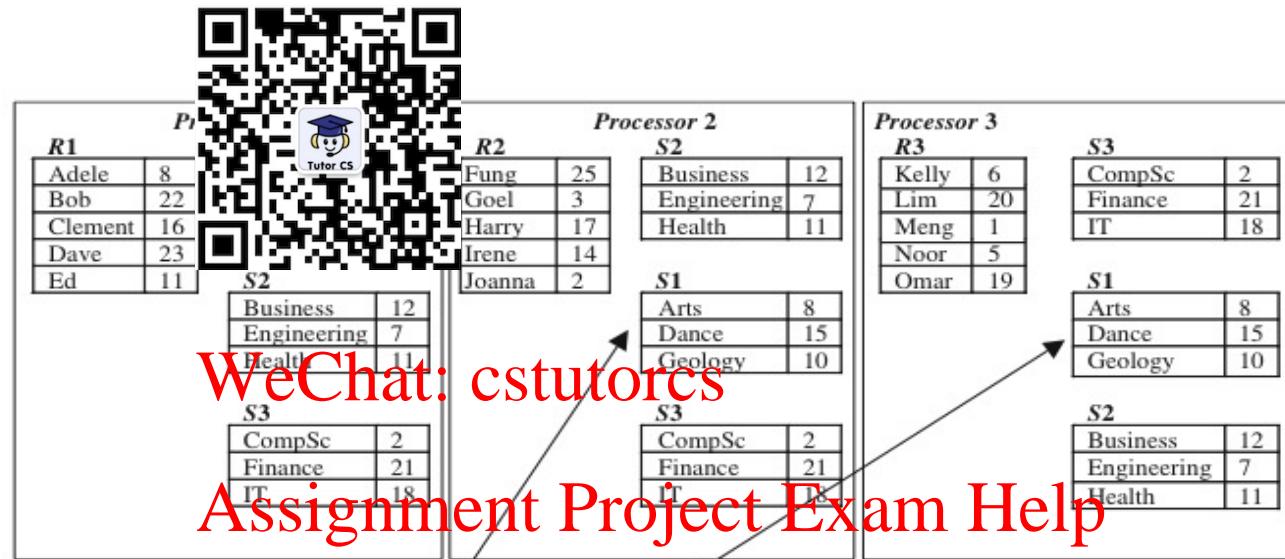


Figure 5.11 Divide and broadcast result

### 5.3. Parallel Join Algorithms (cont'd)

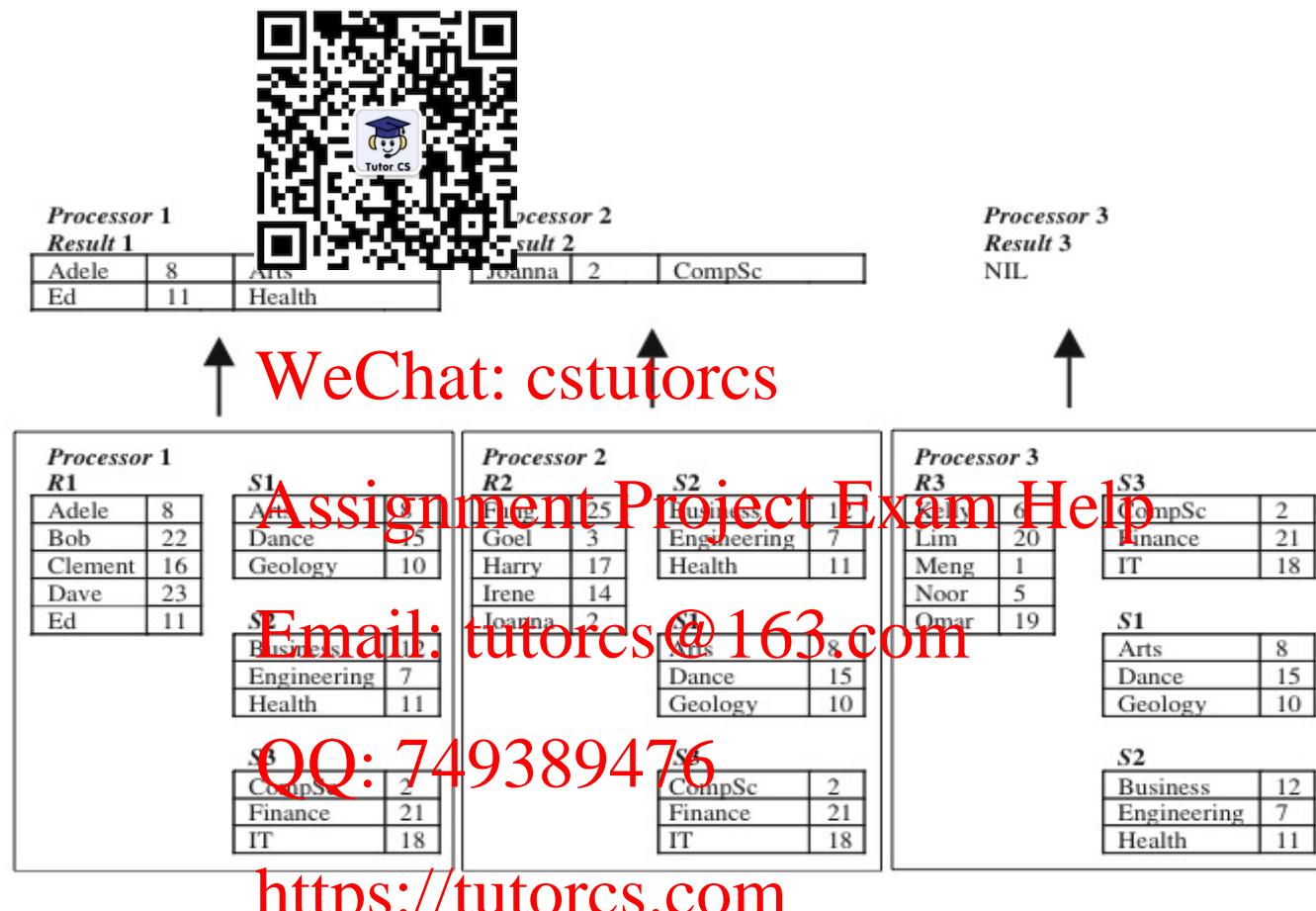


Figure 5.12 Join results based on divide and broadcast

## 5.3. Parallel Join Algorithms (cont'd)



### • Divide and Broadcast Parallel Join Algorithms

- No load imbalance problem, but the broadcasting method is inefficient
- The problem of workload imbalance will occur if the table is already partitioned using random-unequal partitioning
- If shared-memory is used, then there is no replication of the broadcast table. Each processor will access the entire table  $S$  and a portion of table  $R$ . But if each processor does not have enough working space, then the local join might not be able to use a hash-based join.

WeChat: cstutorcs

Assignment Project Exam Help

Email: [tutorcs@163.com](mailto:tutorcs@163.com)

QQ: 749389476

<https://tutorcs.com>

## 5.3. Parallel Join Algorithms (cont'd)



### • Disjoint Partitioning

### Parallel Join Algorithms

- Two stages: data partitioning using a disjoint partitioning, and local join
- Disjoint partitioning: range or hash partitioning
- Local join: any serial local join algorithm

WeChat: cstutorcs  
Assignment Project Exam Help

Email: tutorcs@163.com

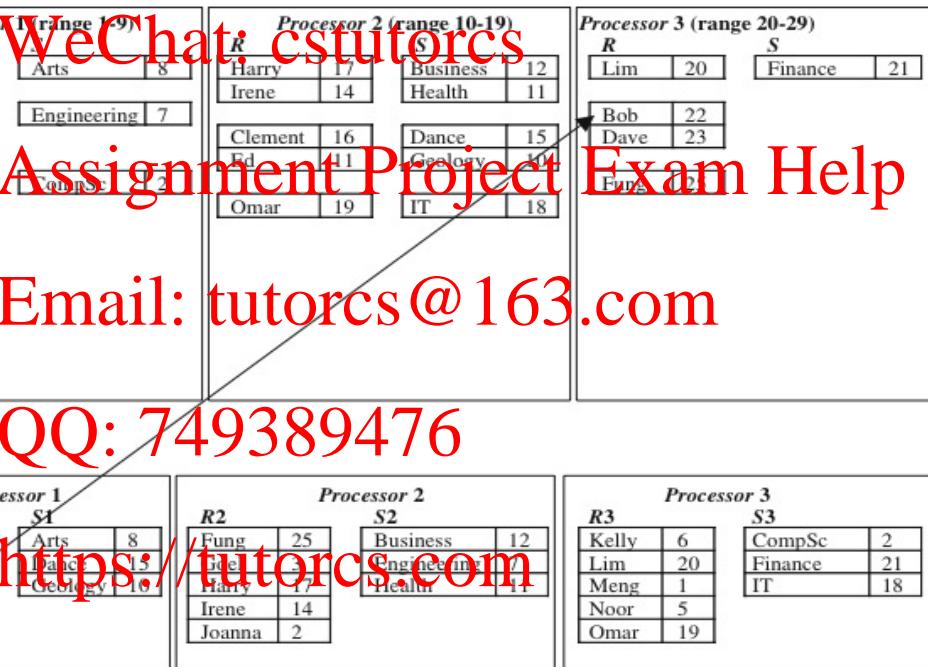
QQ: 749389476

<https://tutorcs.com>

### 5.3. Parallel Join Algorithms (cont'd)



- Example 1: Range Partitioning



The figure illustrates range partitioning for two relations, R and S, across three processors.

**Processor 1 (range 0-9)**

R	
Adele	8
Goel	3
Joanna	2
Kelly	6
Meng	1
Noor	5

S1	
Arts	8
Biology	15
Engineering	7
Geology	10

**Processor 2 (range 10-19)**

R	
Harry	17
Irene	14
Clement	16
Ed	11
Omar	19

S2	
Business	12
Dance	15
Geology	10
IT	18
Health	11

**Processor 3 (range 20-29)**

R	
Lim	20
Bob	22
Dave	23
Fung	21

S3	
Finance	21

**Processor 1 (range 0-9)**

R1	S1
Adele	8
Bob	22
Clement	16
Dave	23
Ed	11

**Processor 2 (range 10-19)**

R2	S2
Fung	25
Lim	27
Harry	17
Irene	14
Joanna	2

**Processor 3 (range 20-29)**

R3	S3
Kelly	6
Lim	20
Meng	1
Noor	5
Omar	19

Figure 5.14 Range partitioning

### 5.3. Parallel Join Algorithms (cont'd)



- Example 1: Range Partitioning

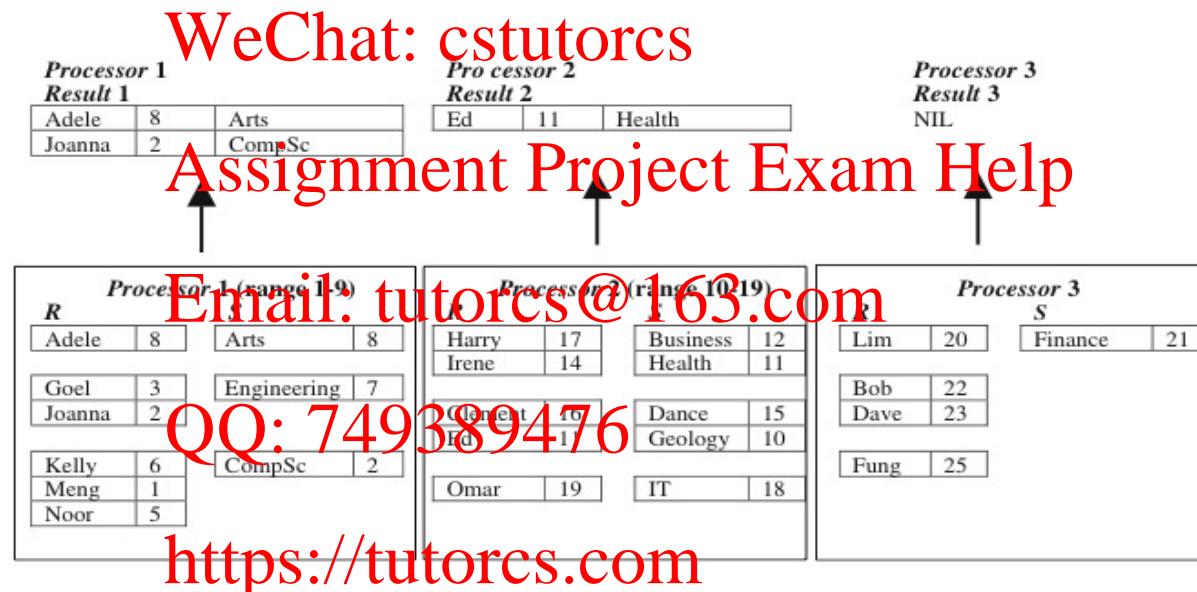


Figure 5.15 Join results based on range partitioning

## 5.3. Parallel Join Algorithms (cont'd)



- Example 2: Hash Partitioning

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com  
QQ: 749389476  
<https://tutorcs.com>

Processor 1 (Hash 1,4,7)		Processor 2 (Hash 2,5,8)		Processor 3 (Hash 3,6,9)	
R	S	R	S	R	S
Bob	22	Geology	10	Harold	7
Clement	16	Irene	14	Irene	11
Fung	25	Joanna	2	Adele	8
Meng	1	Lia	23	Dave	15
Omar	19	Lim	41	Noor	5
				Lim	20
				CompSc	2
				Noor	
					Business
					12

Processor 1		Processor 2		Processor 3	
R1	S1	R2	S2	R3	S3
Adele	8	Arts	8	Kelly	6
Bob	22	Dance	15	Finance	21
Clement	16	Geology	10	IT	18
Dave	23	Health	11		
Ed	11				

Figure 5.16 Hash partitioning

## 5.3. Parallel Join Algorithms (cont'd)



- Example 2: Hash Joining

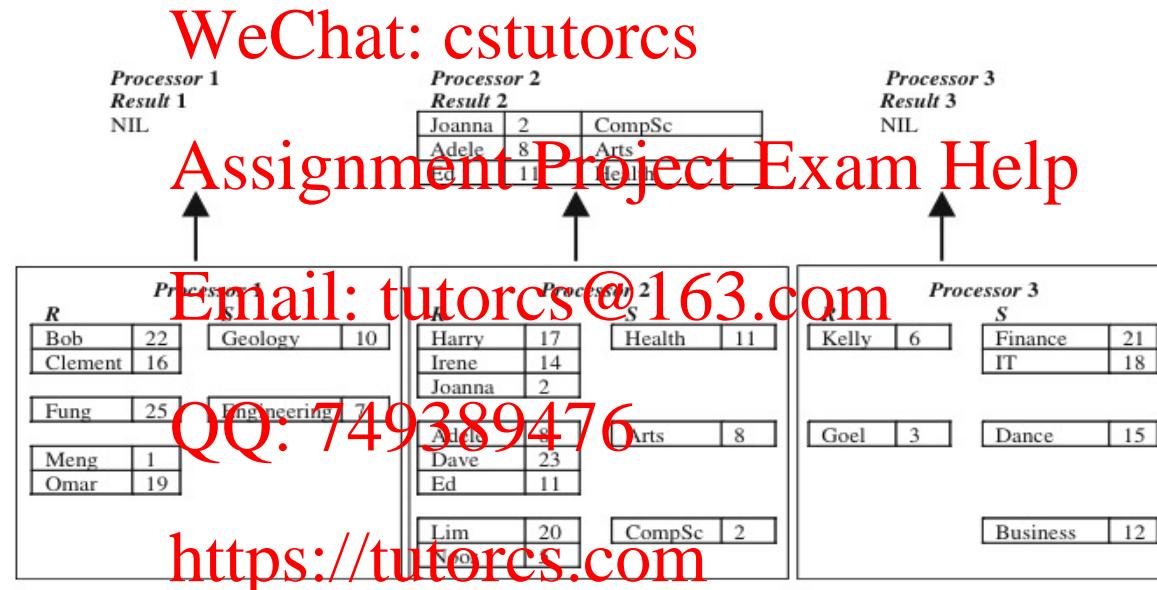


Figure 5.17 Join results based on hash partitioning

## 5.4. Cost Models for Parallel Join

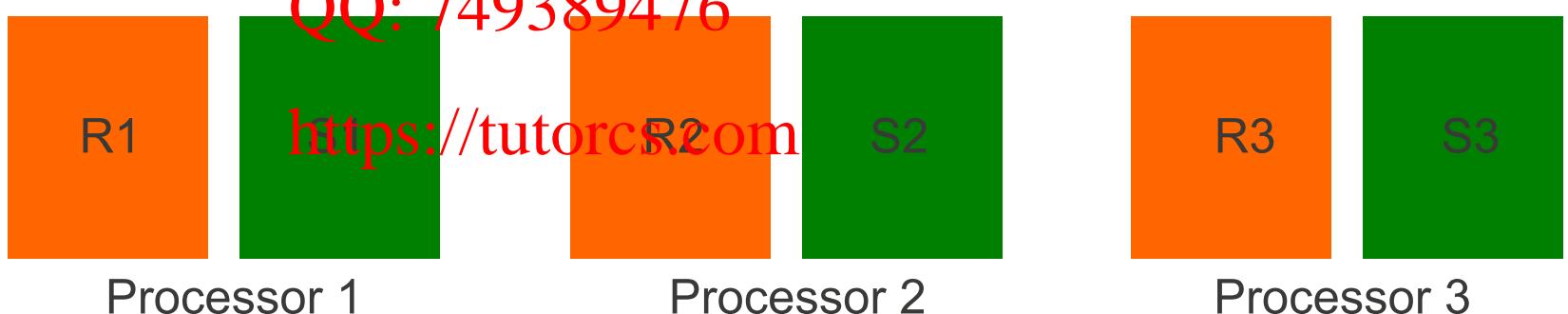


- **Divide and Broadcast**

- Join two tables (**table R** and **table S**)
- The two tables have been partitioned and stored in 3 processors
- The tables have been partitioned using the random-equal data partitioning method  
**WeChat: cstutorcs**
- The table fragments are called R1, R2, R3, and S1, S2, S3 (in general, each fragment is called **R<sub>H</sub>** or **S<sub>H</sub>**, where **H** is the processor number)

Email: tutorcs@163.com

Initial data placement  
(random-equal partitioning)

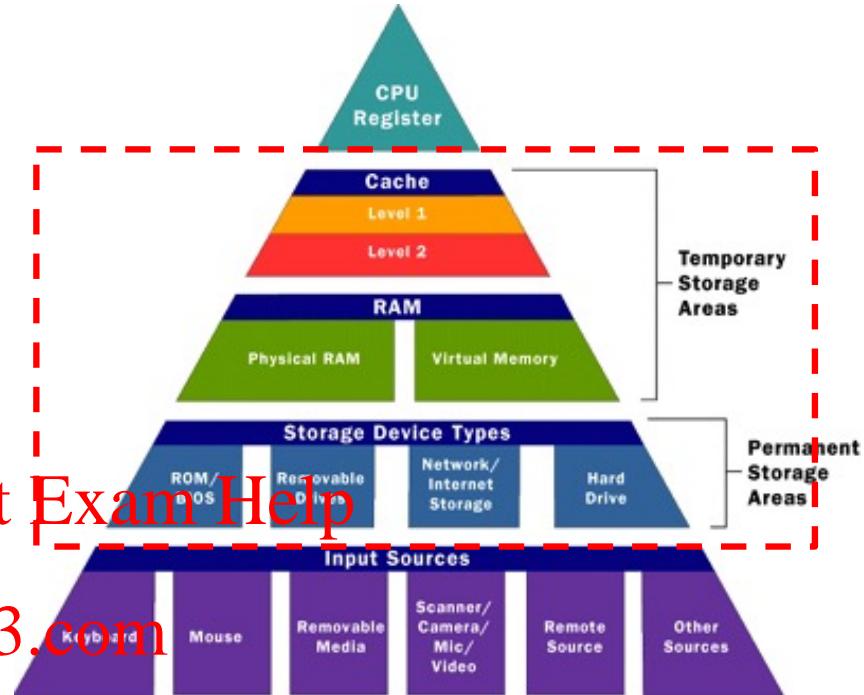


## 5.4. Cost Models for Parallel Join (cont'd)

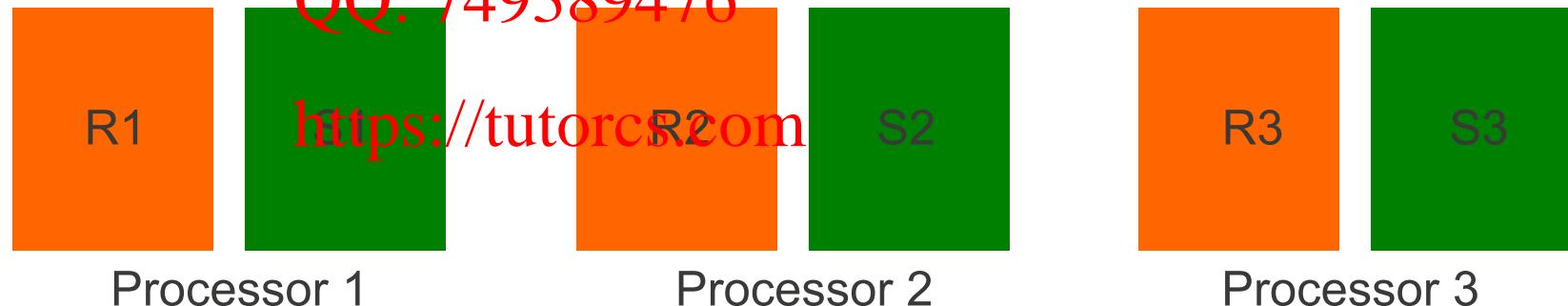


### Divide and Broadcast

- Phase 1: Data Load
  - Come the table to be broadcast.
  - Read all records from each processor
  - If there are 30,000 records in table S, how long does it take to complete the reading before the process can continue to the next phase?
  - Is it after reading 30,000 records, or after reading 10,000 records?



Initial data placement  
(random-equal partitioning)



QQ: 749389476

<https://tutorcs.com>

## 5.4. Cost Models for Parallel Join (cont'd)

- Divide and Broadcast

- $|S| = 30,000$  records
- $N = 3$  processors
- $|S_i| = 10,000$  records ( $|S_i| = |S|/N$ )



WeChat: cstutorcs

- Each record has a fixed length, in bytes
- The size of table S is denoted as S (in bytes)

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Initial data placement  
(random-equal partitioning)



Processor 1



Processor 2



Processor 3

Table 2.1 Cost notations

Symbol	Description
<b>Data parameters</b>	
$R$	Size of table in bytes
$R_i$	Size of table fragment in bytes on processor $i$
$ R $	Number of records in table $R$
$ R_i $	Number of records in table $R$ on processor $i$
<b>Systems parameters</b>	
$N$	Number of processors
$P$	Page size
$H$	Hash table size
<b>Query parameters</b>	
$\pi$	Projectivity ratio
$\sigma$	Selectivity ratio
<b>Time unit cost</b>	
$t_O$	Effective time to read a page from disk
$t_r$	Time to read a record in the main memory
$t_w$	Time to write a record to the main memory
$t_d$	Time to compute destination
<b>Communication cost</b>	
$m_p$	Message protocol cost per page
$m_l$	Message latency for one page

# 程序代写代做 CS编程辅导

## • Exercise 3 (FLUX Quiz)

- $|S| = 600$  records, each record has a length of 100 bytes, and  $N=3$ .
- A.  $S_i = 200$  records
- B.  $S_i = 20,000$  bytes
- C.  $|S_i| = 20,000$  bytes
- D.  $|S_i| = 200$  records
- E. A and C
- F. B and D
- G. All of the above



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

## 5.4. Cost Models for Parallel Join (cont'd)

### Divide and Broadcast

- When a table fragment is read from disk, it is based on the size of the table (in bytes), not how many records.
- Hence, it uses  $S_i$ , and not  $|S_i|$



WeChat: cstutorcs

- When the disk reads a table fragment ( $S_i$ ), it reads a disk block at a time. The size of a disk block is  $P$  (or page size),
- The loading time or  $\text{Scan cost} = (S_i/P) \times IO$ , where  $IO$  is the time taken to load 1 page from disk to main memory



## 5.4. Cost Models for Parallel Join (cont'd)

程序代写代做CS编程辅导

### Divide and Broadcast

- Once the loading is complete, the records are not yet ready in memory. The records must be read from memory and be written to the main memory data page.
- All processing in the processor is based on the number of records, and not the byte size



WeChat: cstutors

Assignment Project Exam Help

Email: tutorcs@163.com

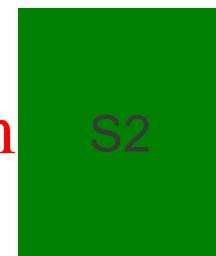
QQ: 749389476

<https://tutorcs.com>

Initial data placement  
(random-equal partitioning)



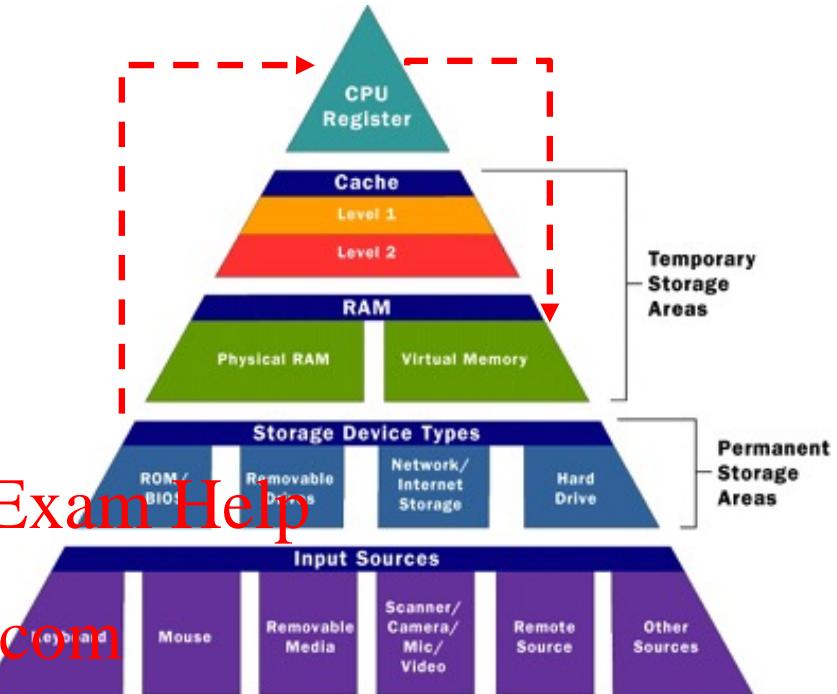
Processor 1



Processor 2



Processor 3



## 5.4. Cost Models for Parallel Join (cont'd)



### Cost Models for Broadcast

- Phase 1: data loading consists of the *scan costs* and the *select costs*

WeChat: cstutorcs

- Scan cost for loading data from local disk in each processor is:

Assignment Project Exam Help

- Select cost for getting record out of data page is:

$$|S_i| \times (tr + tw)$$

Email: tutorcs@163.com  
QQ: 749389476

<https://tutorcs.com>

# 程序代写代做 CS编程辅导

## • Exercise 4 (FLUX Quiz)

- The cost (the time taken) from disk is called...
- A. Scan Cost
- B. Select Cost
- C. Both A and B



WeChat: cstutorcs

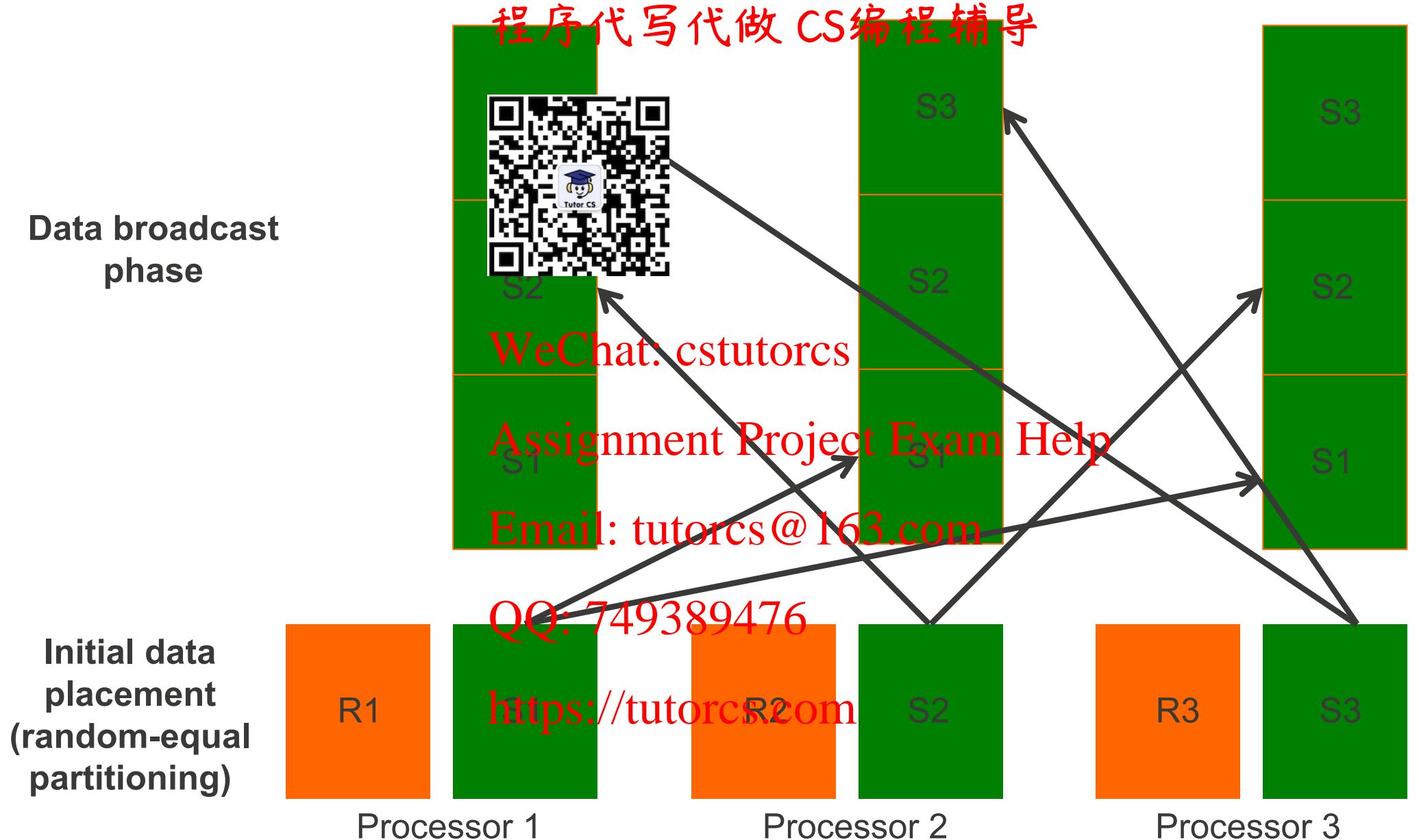
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

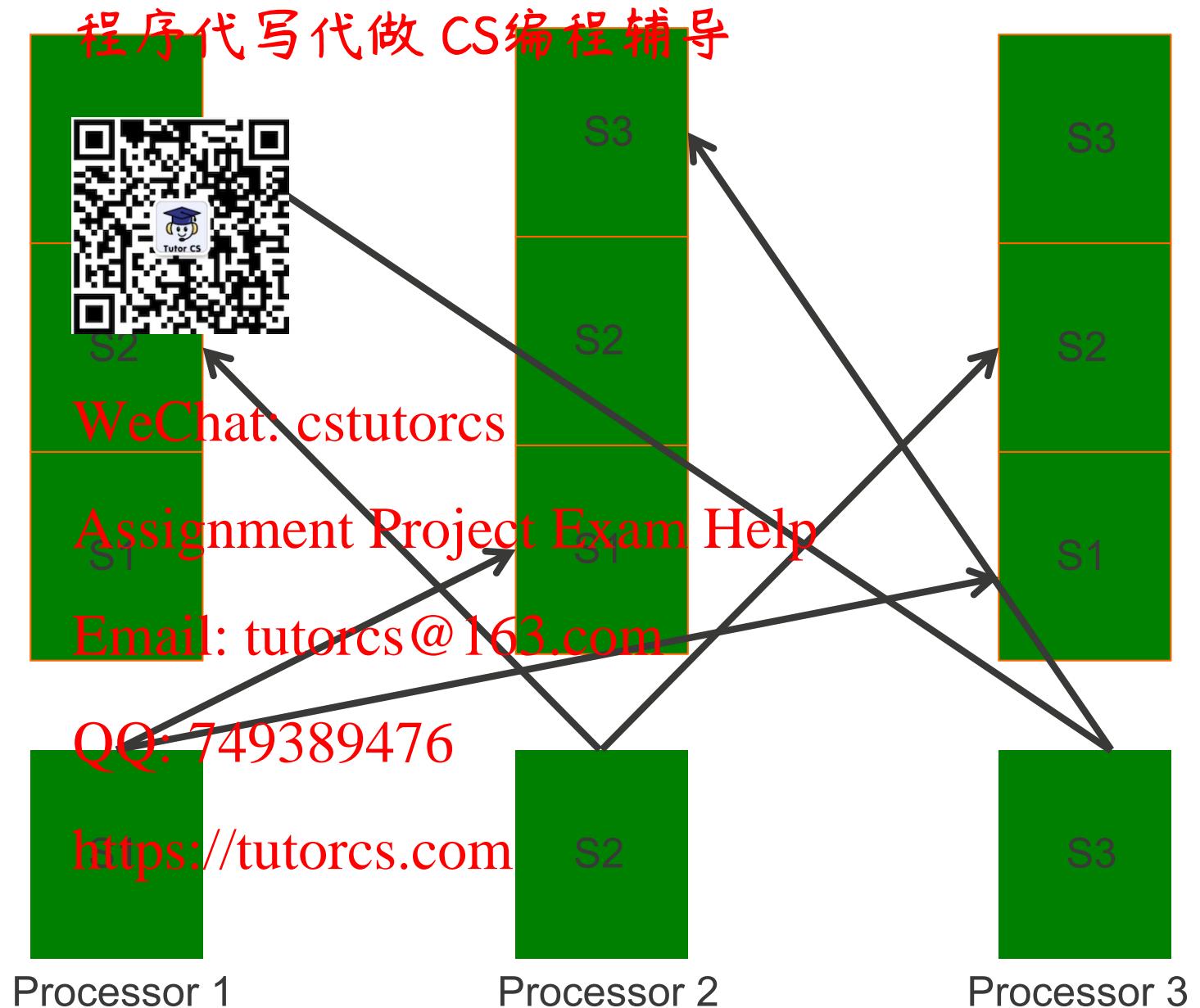
<https://tutorcs.com>

**Phase 2: Data Broadcasting.** Table fragment S1 must be broadcasted (copied) to processors 2 and 3. Table fragment S2 to processors 1 and 3, etc...



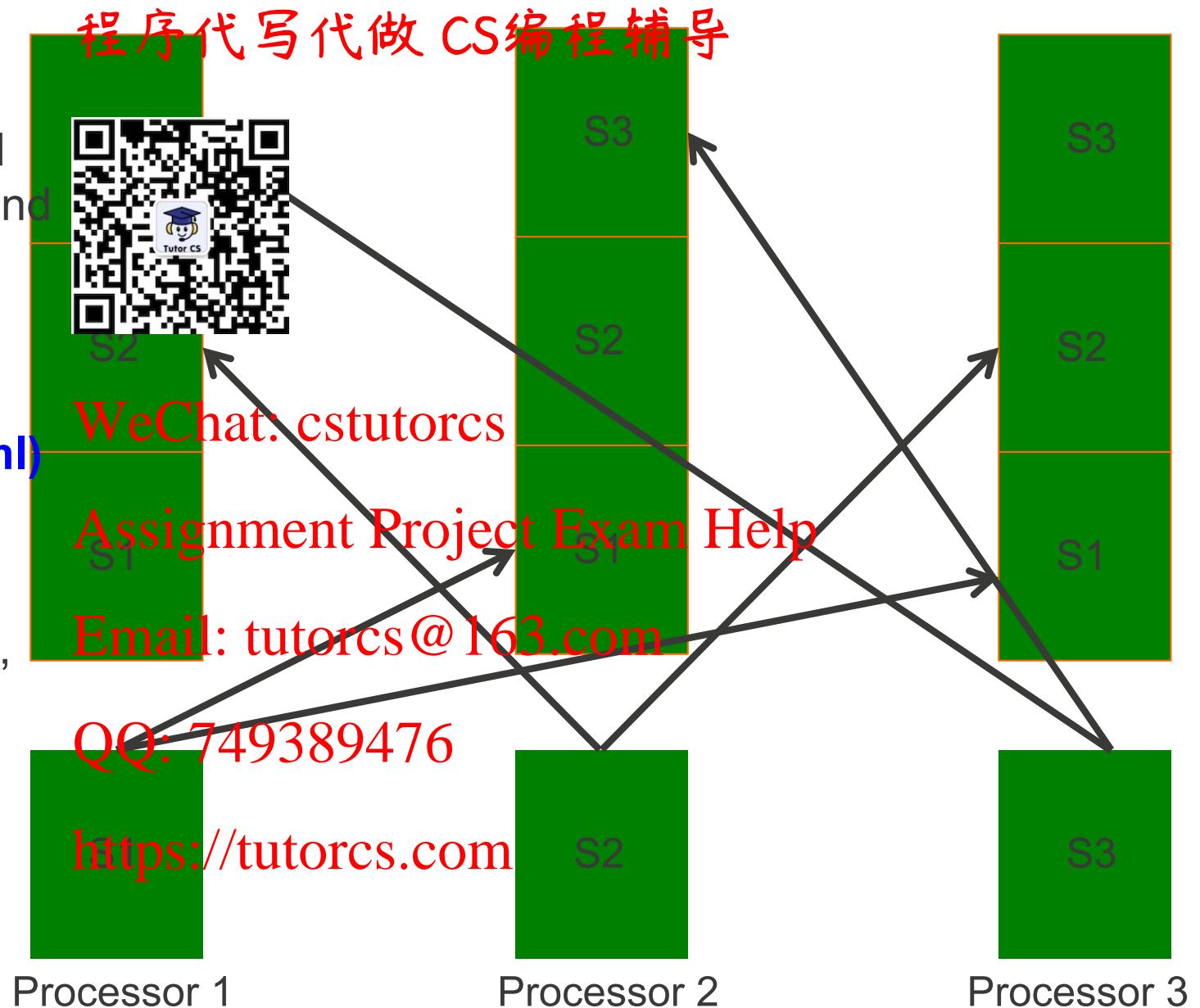
**Phase 2: Data Broadcasting.** Table fragment S1 must be broadcasted (copied) to processors 2 and 3. Table fragment S2 to processors 1 and 3, etc...

- Data broadcast is done through network (so it is a network data transfer)
- The transfer cost is based on **how many bytes** of data being transferred
- Hence, we use **Si**, instead of  $|S_i|$
- Data transfer is also done page-per-page (**P**)
- Hence, **Si/P**



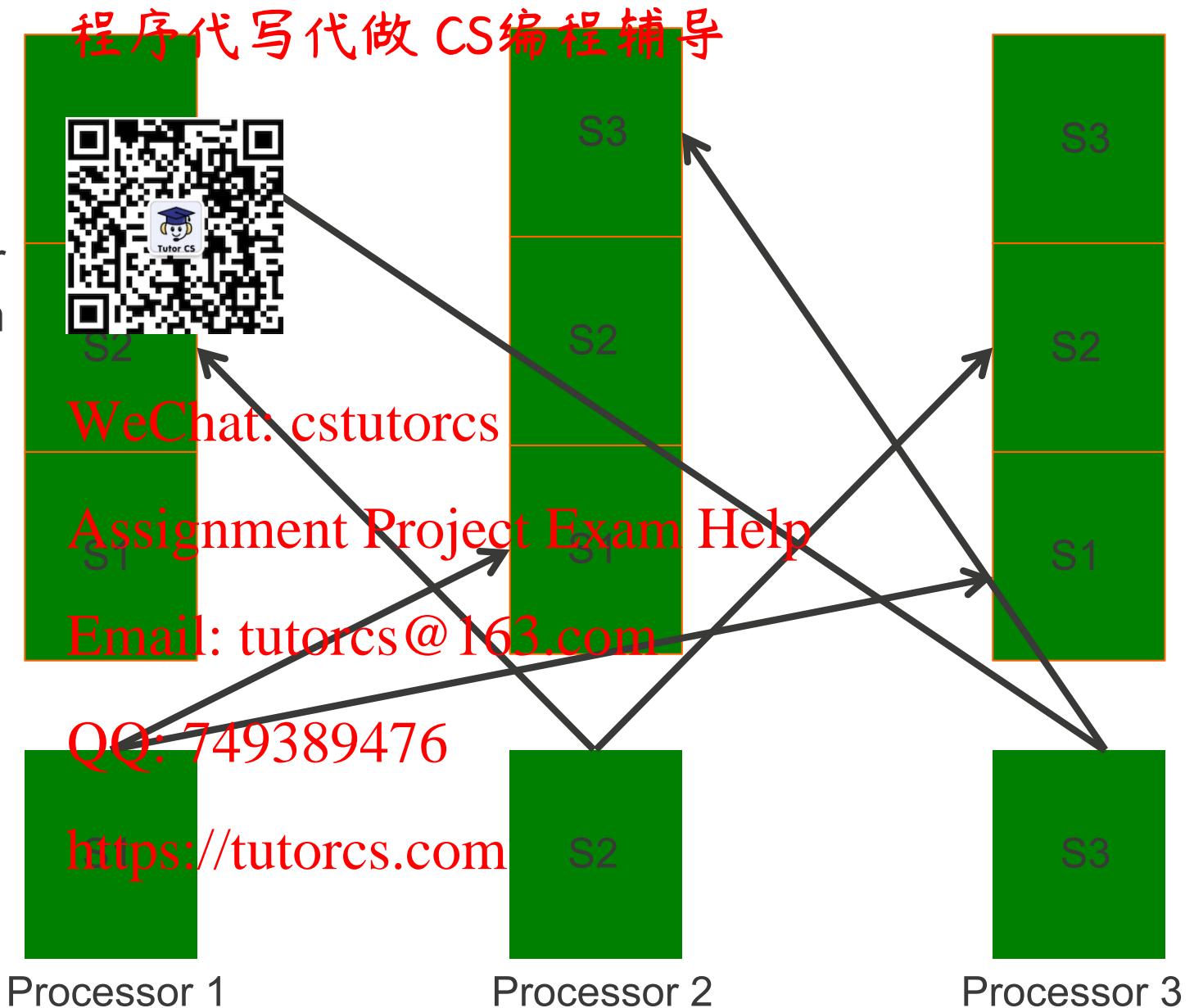
**Phase 2: Data Broadcasting.** Table fragment S1 must be broadcasted (copied) to processors 2 and 3. Table fragment S2 to processors 1 and 3, etc...

- S1 must be transferred twice (to processor 2 and to processor 3)
- **Transfer cost =**  
$$(S_i/P) \times (N-1) \times (mp+ml)$$
- Where **mp** is message protocol per data page, and **ml** is message latency per data page.



**Phase 2: Data Broadcasting.** Table fragment S1 must be broadcasted (copied) to processors 2 and 3. Table fragment S2 to processors 1 and 3, etc...

- Processor 2 must receive S1 from processor 1; Processor 3 must receive S1 from processor 1
- **Receiving cost =  $(S/P - Si/P) \times (mp)$**
- Why  $(S/P - Si/P)$ ?
- Why  $(mp)$  only?



## 5.4. Cost Models for Parallel Join (cont'd)



### Cost Models for Record Broadcast

- Phase 2: The broadcast cost by each processor broadcasting its fragment to all other processors

WeChat: cstutorcs

- Data transfer cost is:  $(S_i / P) \times (N - 1) \times (m_p + m_l)$
- The  $(N-1)$  indicates that each processor must broadcast to all other processors. Note that broadcasting from one processor to the others has to be done one processor at a time, although all processors send the broadcast in parallel. The above cost equation would be the same as  $(S/P - S_i/P) \times (m_p + m_l)$ , where  $(S/P - S_i/P)$  is the size of other fragments.

Assignment Project Exam Help

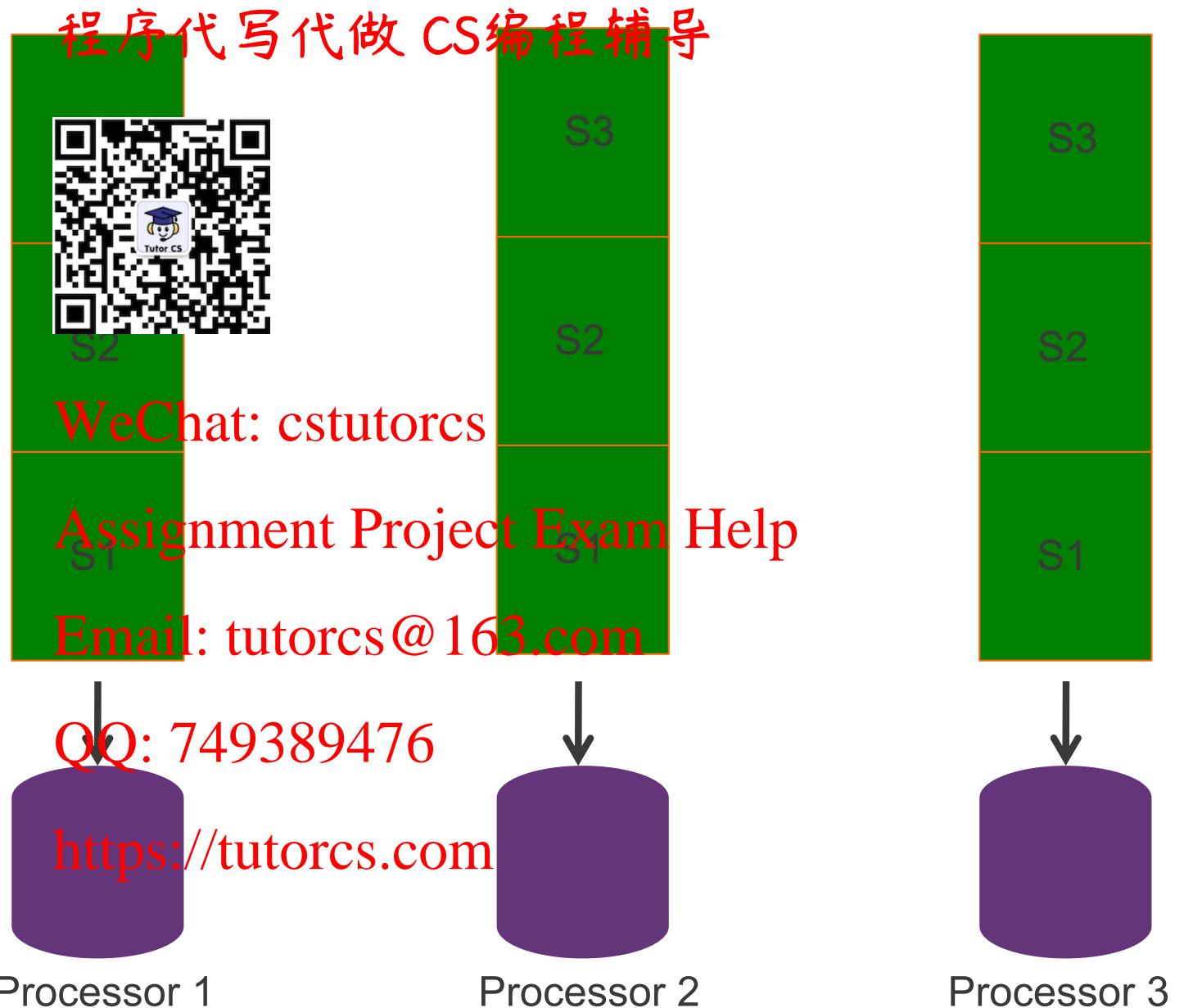
Email: [tutors@163.com](mailto:tutors@163.com)

QQ: 749389476

- Receiving records cost is:  $(S/P - S_i/P) \times (m_p)$

**Phase 3: Data Storing.** Each fragment in each processor must be stored/written in the local disks

- **Storing cost =  $(S/P - Si/P) \times (IO)$**
- Why  $(S/P - Si/P)$ ?



## 5.4. Cost Models for Parallel Join (cont'd)



### Cost Models for Broadcast

- Phase 3: Each processor after receiving all other fragments of table S, needs to be stored on local disk.

WeChat: cstutorcs

- Disk cost for storing the table is:  $(S/P - Si/P) \times IO$

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

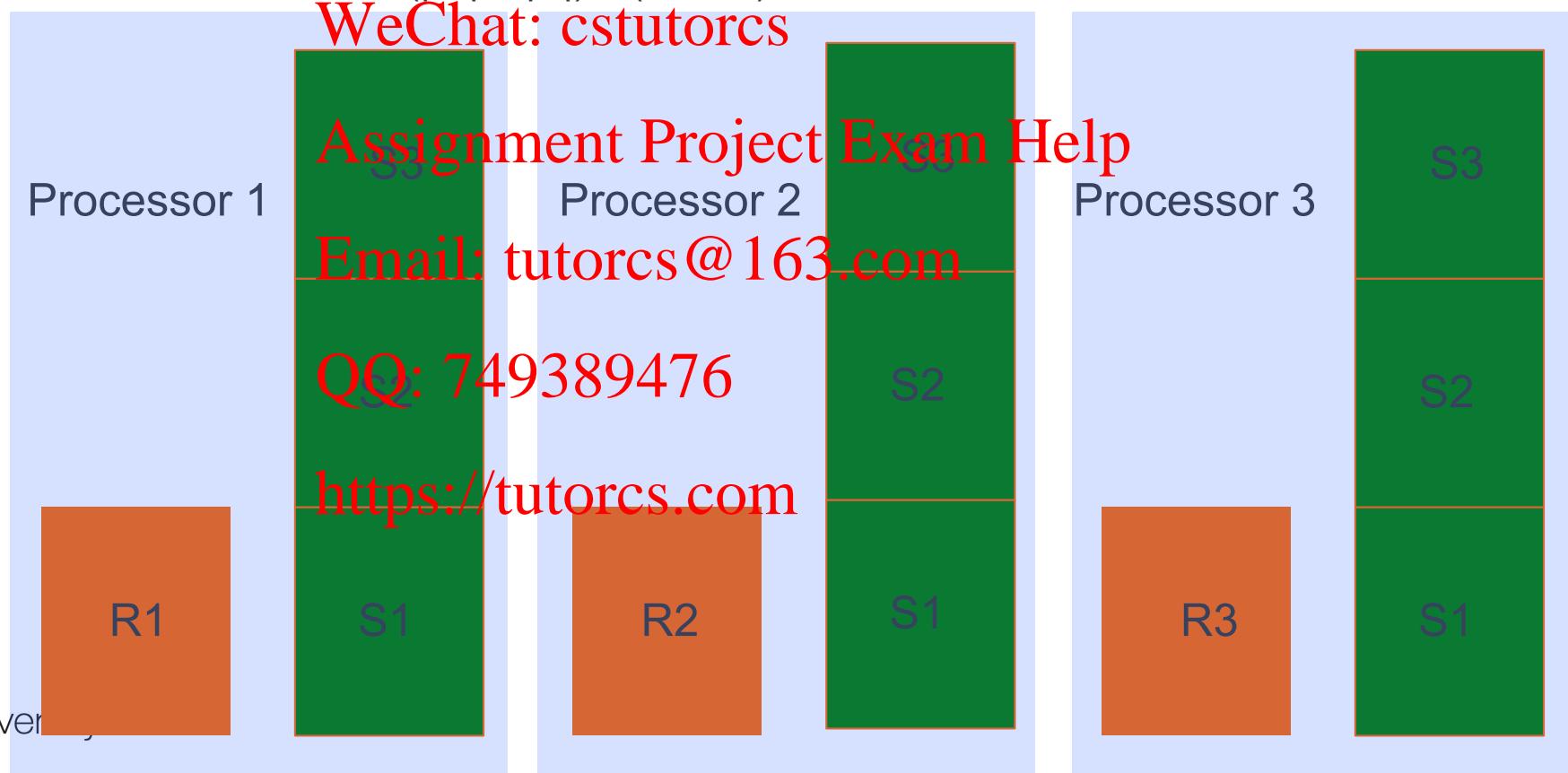
## 5.4. Cost Models for Parallel Join (cont'd)



### Local Join

- Each processor performs a local join (using a Hash Join Algorithm)
- Phase 1: Loading

$$\text{Scan cost} = ((R_i / P) + (S / P)) \times \text{IO}$$
$$\text{Select cost} = (|R_i| + |S|) \times (tr + tw)$$



## 5.4. Cost Models for Parallel Join (cont'd)



### Cost Models for Parallel Join

- Assume to use hash-based join
- Three main phases: data loading from each processor, the joining process (hashing and probing), and result storing in each processor.

Assignment Project Exam Help

- Phase 1: The data loading consists of scan costs and select costs
- Scan cost =  $((R_i / P) + (S / P)) \times IO$
- Select cost =  $(|R_i| + |S|) \times (tr + tw)$

QQ: 749389476

<https://tutorcs.com>

# Query Parameters



- **Projectivity ratio  $\pi$  :**
  - Ratio between projected attribute size and original record length
- **Selectivity ratio  $\sigma$  :**
  - Ratio between number of records in the query result and original total number of records

Example: Join selectivity ratio

If the query operation involves two tables (like in a join operation), a selectivity ratio can be written as  $\sigma_j$ , for example. The value of  $\sigma_j$  indicates the ratio between the number of records produced by a join operation and the number of records of the Cartesian product of the two tables to be joined. For example,  $|R_i| = 1000$  records and  $|S_i| = 500$  records; if the join produces 5 records only, then the join selectivity ratio  $\sigma_j$  is  $5/(1,000 \times 500) = 0.00001$ .

## 5.4. Cost Models for Parallel Join (cont'd)

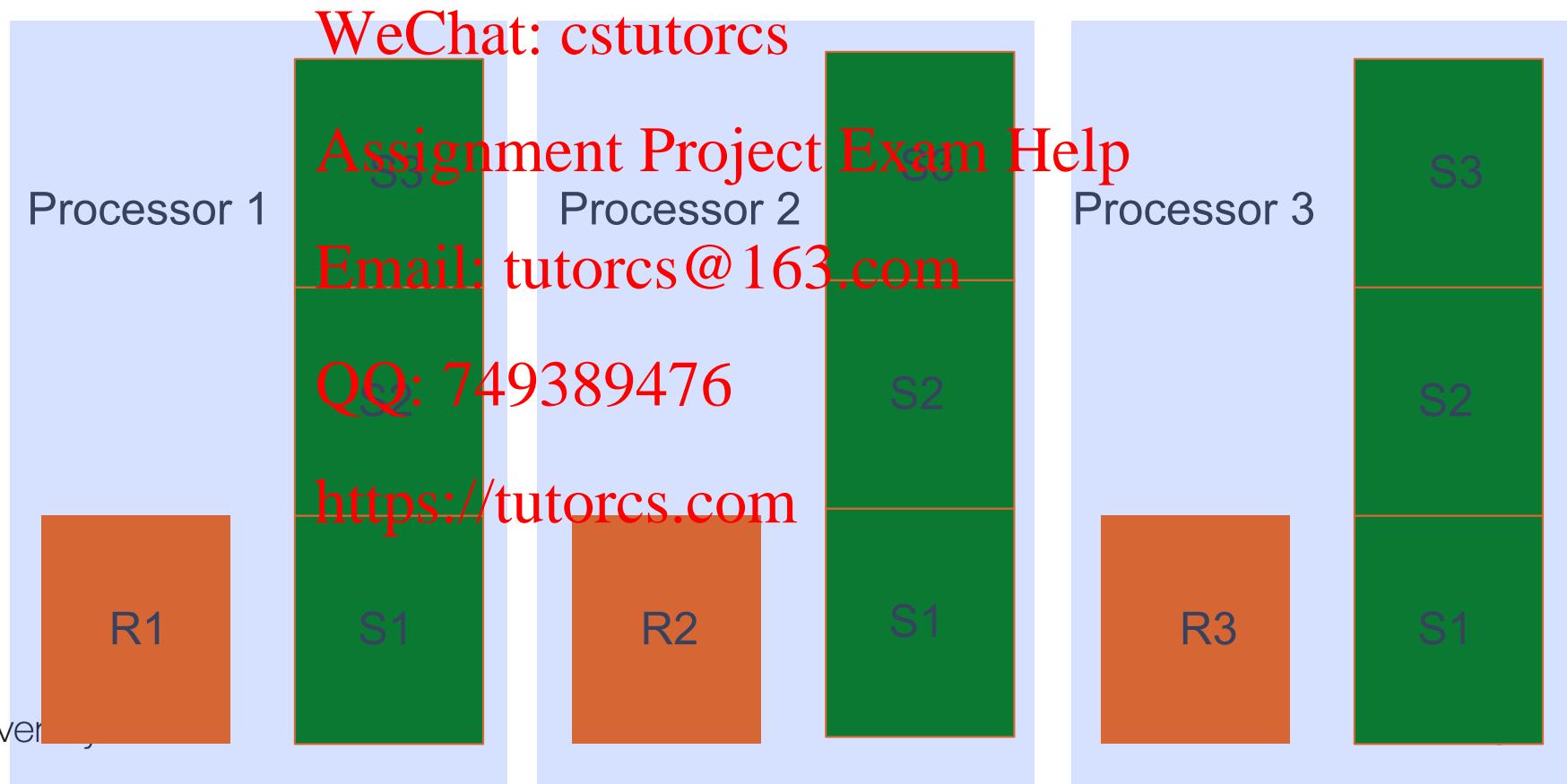


### Local Join

- Phase 2: Join cost (using Hash Join Algorithm)

$$(|R_i| \times (tr + th + tj))$$

Where  $tr$  is reading cost,  $th$  is hashing cost, and  $tj$  is joining cost



## 5.4. Cost Models for Parallel Join (cont'd)



### Cost Models for Join

- Phase 2: The join process is the hashing and probing costs
- Join costs involve reading, hashing, and probing:  
$$(|R_i| \times (tr + th) + (|S| \times (tr + th + t_j)))$$

WeChat: cstutorcs

### Assignment Project Exam Help

- If the memory size is smaller than the hash table size, we normally partition the hash table into multiple buckets whereby each bucket can perfectly fit into main memory. All but the first bucket is spilled to disk.
- Reading/Writing of overflow buckets cost is the I/O cost associated with the limited ability of main memory to accommodate the entire hash table.

QQ: 749389476

$$\frac{H}{|R_i|} \times \left( \frac{R_i}{P} \times 2 \times IO \right)$$

## 5.4. Cost Models for Parallel Join (cont'd)



### Cost Models for Parallel Join

- Reading/Writing of overflow buckets cost is the I/O cost associated with the limited ability of main memory to accommodate the entire hash table.

WeChat: cstutorcs

$$\left(1 - \min\left(\frac{H}{|R_i|}, 1\right)\right) \times \left(\frac{R_i}{P} \times 2 \times IO\right)$$

Assignment Project Exam Help

- For example, the Hash table can only occupy 10 records at a time from table  $R_i$ . Assume  $|R_i|=50$  records. That means that there will be 5 buckets. Because the main memory can take one bucket only, it means the 4 buckets must be stored on disk.
- $(1-\min(H/|R_i|,1)) = 1-\min(0.2,1) = 1-0.2 = 0.8$
- If  $|R_i|=10$  or less, then  $(1-\min(H/|R_i|,1)) = 1-1 = 0$ ; That means there is no overflow bucket cost.
- $(R_i/P \times 2 \times IO) \rightarrow$  the constant 2 means two input/output accesses: one for spooling, and the other for reading it back from the disk

## 5.4. Cost Models for Parallel Join (cont'd)



- **Cost Models for Parallel Join**

- Phase 3: query results storing cost, consisting of generating result cost and disk cost.

WeChat: cstutorcs

- Generating result records cost is:  $|R_i| \times \sigma_j \times |S| \times tw$

Assignment Project Exam Help

- Disk cost for storing the final result is:  $(\pi_R \times R_i \times \sigma_j \times \pi_S \times S / P) \times IO$

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

## 5.4. Cost Models for Parallel Join (cont'd)



- Home Work

- Disjoint data partitioning based parallel join algorithm
- Cost Model for Disjoint data partitioning based parallel join algorithm

WeChat: **ostutors**  
**Assignment Project Exam Help**

<https://onlinelibrary-wiley-com.ezproxy.lib.monash.edu.au/doi/pdf/10.1002/9780470391365.ch5>  
Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

## 5.5. Parallel Join Optimization



- The aim of query processing optimization in general is to speed up the query processing time
- In terms of parallelism, the reduction in the query elapsed time is achieved by having each processor finish its execution as early as possible and as evenly as possible
- In the disjoint partitioning, after the data is distributed to the designated processors, the data has to be stored on disk. Then in the local join, the data has to be loaded from the disk again → **managing main memory issue**

**Email: tutorcs@163.com**

**QQ: 749389476**

**<https://tutorcs.com>**

## 5.5. Parallel Join Optimization (cont'd)



### Optimizing Main

- Disk access is the most expensive operations, so need to reduce disk access as much as possible
- If it is possible, one single scan of data should be done. If not, then minimize the number of scan
- If main memory size is unlimited, single disk scan is possible
- However, main memory size is not unlimited, hence optimizing main memory is critical
- Problem: In the distribution, when the data arrives at a processor, it is stored in disk. In the local join, the data needs to be reloaded from disk
- This is inefficient. When the data arrives after being distributed from other processor, the data should be left in main memory, so that the data remain available in the local join process
- The data left in the main memory can be as big as the allocated size for data in the main memory

## 5.5. Parallel Join Optimization (cont'd)



### Optimizing Main

- Assuming that the size of main memory for data is  $M$  (in bytes), the disk cost for storing data distribution with a disjoint partitioning is:

WeChat: cstutorcs  
 $((R_i / P) + (S_i / P) - (M/P)) \times IO$

- And the local join scan cost is then reduced by  $M$  as well:

Email: tutorcs@163.com

- When the data from this main memory block is processed, it can be swapped with a new block. Therefore, the saving is really achieved by not having to load/scan the disk for one main memory block

QQ: 749389476  
<https://tutorcs.com>

## 5.5. Parallel Join Optimization (cont'd)



### Load Balancing

- Load imbalance problem in parallel query processing. It is normally caused by data skew or processing skew
- No load imbalance in hash-based and broadcast-based parallel join. But this kind of parallel join is unattractive, due to the heavy broadcasting
- In disjoint-based parallel join algorithms, processing skew is common
- To solve this skew problem, create more fragments than the available processors, and then rearrange the placement of the fragments

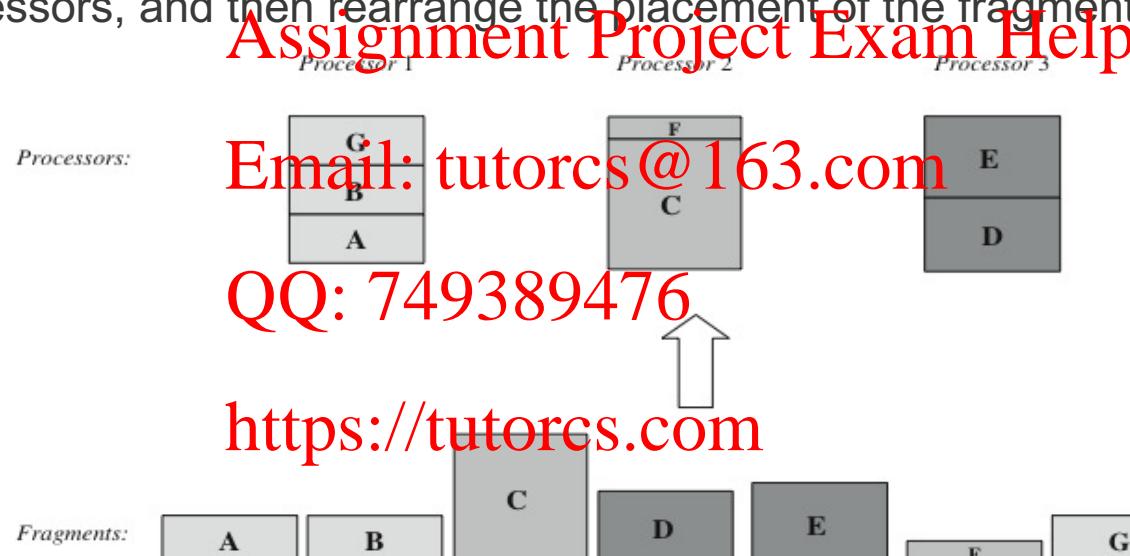


Figure 5.19 Load balancing

## 5.6. Summary

程序代写代做 CS编程辅导

- Parallel join is one of the important operations in parallel database systems
- Parallel join algorithms have two stages
  - Data partitioning
  - Local join
- Two types of data partitioning
  - Divide and broadcast
  - Disjoint partitioning
- Three types of local join
  - Nested-loop join
  - Sort-merge join
  - Hash-based join



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



MONASH University

程序代写代做 CS 编程辅导

Information Technology

FIT5202 (Week III - Join)



Week 3b – Parallel Outer Join

WeChat: cstutorcs

Assignment Project Exam Help

**algorithm distributed systems database**

Email: tutorcs@163.com

systems **computation knowledge management**

QQ: 749389476

**design e-business model data mining intelligent**

**distributed systems database software engineering**

<https://tutorcs.com>

**computation knowledge management analysis**

# Join Queries

程序代写代做 CS编程辅导

- Two types of Join Qu

- Inner Join



Select R.x, R.a, S.y, S.b

WeChat: cstutorcs  
From R, S

Where R.a = S.b;  
Assignment Project Exam Help

- Outer Join

Email: tutorcs@163.com

Select R.x, R.a, S.y, S.b  
QQ: 749389476

From R left outer join S on R.a = S.b;

<https://tutorcs.com>

# Join Queries

程序代写代做 CS 编程辅导

R	a
1	1
2	2
2	3

S	y	b
0	6	
0	4	
3	1	
6	2	
1	6	
4	2	
7	2	
7	1	
2	1	
5	5	
5	6	
8	9	

Results	x	a	y	b
0	1	3	1	
0	1	7	1	
0	1	2	1	
1	2	6	2	
1	2	4	2	
1	2	7	2	

WeChat: cstutorcs

Assignment Project Exam Help

## - Inner Join

Select R.x, R.a, S.y, S.b

From R, S

Where R.a = S.b;

QQ: 749389476

<https://tutorcs.com>

# Join Queries

程序代写代做 CS 编程辅导

R

	a
1	1
2	2
2	3

S

y	b
0	6
0	4
3	1
6	2
1	6
4	2
7	2
7	1
2	1
5	5
5	6
8	9

Results

x	a	y	b
0	1	3	1
0	1	7	1
0	1	2	1
1	2	6	2
1	2	4	2
1	2	7	2
2	3	Null	Null

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

From R left outer join S

QQ: 749389476

On R.a = S.b;

<https://tutorcs.com>

# 程序代写代做 CS编程辅导

## Exercise 1

- Identify the LEFT COIN?

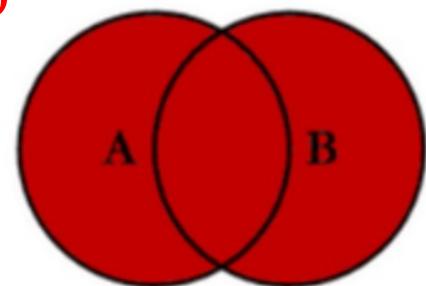
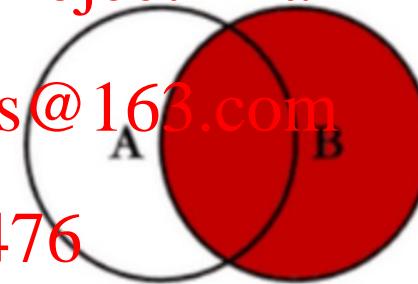
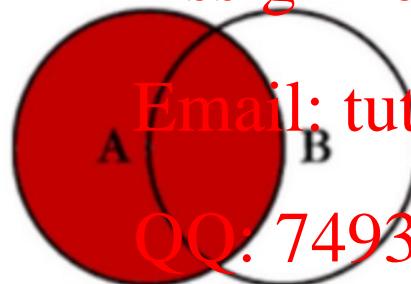
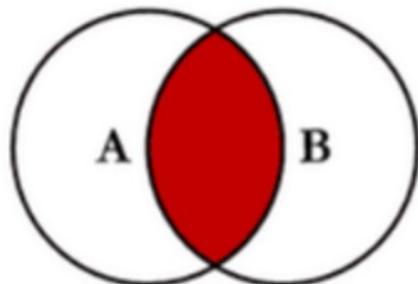


WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

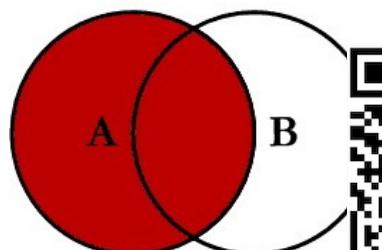
QQ: 749389476



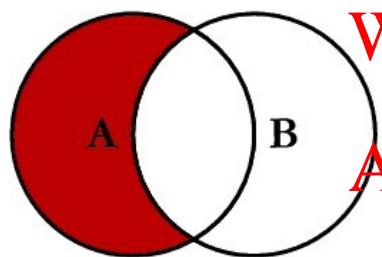
<https://tutorcs.com>

程序代写代做 CS编程辅导

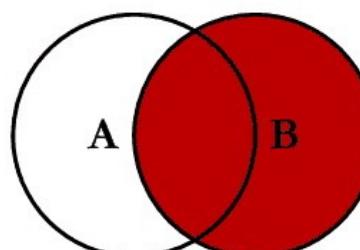
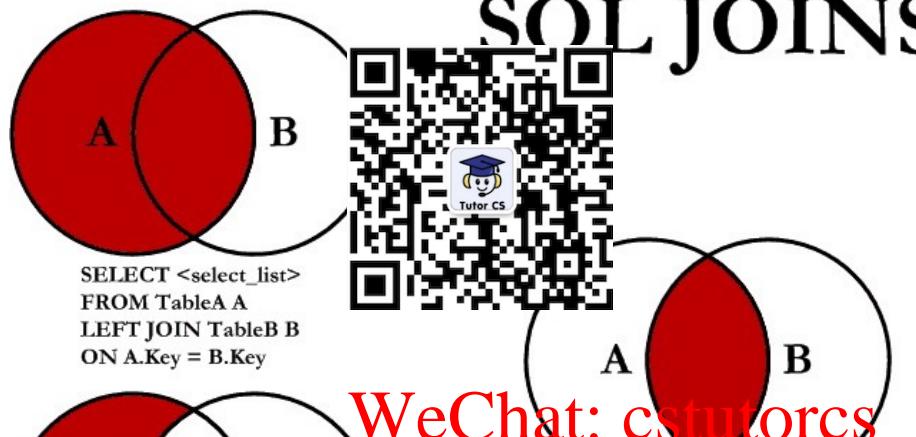
## SQL JOINS



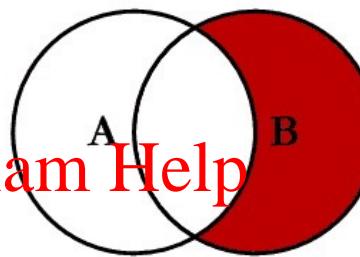
```
SELECT <select_list>  
FROM TableA A  
LEFT JOIN TableB B  
ON A.Key = B.Key
```



```
SELECT <select_list>  
FROM TableA A  
LEFT JOIN TableB B  
ON A.Key = B.Key  
WHERE B.Key IS NULL
```



```
SELECT <select_list>  
FROM TableA A  
RIGHT JOIN TableB B  
ON A.Key = B.Key
```



```
SELECT <select_list>  
FROM TableA A  
RIGHT JOIN TableB B  
ON A.Key = B.Key  
WHERE A.Key IS NULL
```

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

© C.L. Moffatt, 2008

<https://www.codeproject.com/Articles/33052/Visual-Representation-of-SQL-Joins>

# Parallel Join Query Processing



- Parallel Inner Join covers:
  - **Data Partitioning**
    - Divide and Broadcast
    - Disjoint Partitioning
  - **Local Join**
    - Nested-Loop Join
    - Sort-Merge Join
    - Hash Join
- Example of a Parallel Inner Join Algorithm
  - **Divide and Broadcast, plus Hash Join**

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

# Parallel Join Query Processing



- Parallel Outer Join processing methods
  - ROJA (Redistribution Outer Join Algorithm)
  - DOJA (Duplication Outer Join Algorithm)
  - DER (Duplication & Efficient Redistribution)

Assignment Project Exam Help

- Load Balancing
  - OJSO (Outer Join Skew Optimization)

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

WeChat: cstutorcs

Email: tutorcs@163.com

# 1. ROJA

```
SELECT R.x, R.a, S.y, S.b  
FROM R left outer join S on R.a = S.y
```

$h(i)=1$

R		S	
x	a	y	b
0	1	2	3
0	4	4	2
3	1	7	2
6	2	7	1

Processor 1

 MONASH University

程序代写代做CS编程辅导

Step 1: Distribute or reshuffle data based on join attribute.

Step 2: Each processor performs local outer Join.



Eg. Using hash func  $h(i) = i \bmod 3 + 1$

$h(i)=2$

WeChat: cstutorcs

Assignment Project Exam Help

R      S  
x    a      y    b

Email: 1tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

$h(i)=3$

R		S	
x	a	y	b
2	3	2	1
5	5		
5	6		
8	9		

Processor 2

## 2. DOJA

SELECT R.x, R.a, S.y, S.b  
FROM R left outer join S on R.a

程序代写代做CS编程辅导

Step 1: Replicate small table.

Step 2: Local Inner Join



Step 3: Hash redistribute inner join result based  
on attribute x.

R		S	
x	a	y	b
0	1	0	6
1	2	0	4
2	3	3	1
	6	2	

Processor 1

 MONASH University

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: 1tutorcs@163.com  
QQ: 749389476  
<https://tutorcs.com>

R		S	
x	a	y	b
0	1	4	2
2	3	7	2
4	2	1	1

Processor 2

R		S	
x	a	y	b
2	3	2	1
1	2	5	5
0	1	5	6
	8	7	9

Processor 3

## 2. DOJA

SELECT R.x, R.a, S.y, S.b  
FROM R left outer join S on R.a = S.y



程序代写代做 CS 编程辅导  
Step 4: Local outer join

R		J			
x	a	x	a	y	b
0	1	0	1	3	1
0	1	7	1		
0	1	2	1		

Processor 1

 MONASH University

R		J			
x	a	x	a	y	b
1	2	1	2	6	2
1	2	7	2		
1	2	2	7	2	

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: [tutorcs@163.com](mailto:tutorcs@163.com)  
QQ: 749389476  
<https://tutorcs.com>

Processor 2

R		J			
x	a	x	a	y	b
2	3				

Processor 3

### 3. DER

SELECT R.x, R.a, S.y, S.b  
FROM R left outer join S on R.a = S.y

程序代写代做 CS 编程辅导

Step 1: Replicate small table (left)

Step 2: Local Inner Join



Step 3: Select ROW ID of left table with no matches.

Step 4: Redistribute the ROW ID.

Step 5: Store the ROW ID that appears as many times as the number of processors

Row ID	R		S	
	x	a	y	b
0	0	1	0	6
1	1	2	0	4
2	2	3	3	1
			6	2

Processor 1

Row ID	R		S	
	x	a	y	b
0	0	1	0	6
1	1	2	0	4
2	2	3	3	1
			6	2

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: tutorcs@163.com  
QQ: 749389476  
<https://tutorcs.com>

Processor 2

Row ID	R		S	
	x	a	y	b
0	0	1	2	1
1	1	2	5	5
2	2	3	5	6
			8	9

Processor 3

### 3. DER

```
SELECT R.x, R.a, S.y, S.b  
FROM R left outer join S on R.x = S.a
```



程序代写代做 CS编程辅导

Step 6: Inner join

R		Row ID		
x	a	y	b	
0	1	1	3	1
1	2	6	2	

Processor 1

 MONASH University

R		Row ID		
x	a	y	b	
1	2	6	2	
0	1	7	1	
1	2	7	2	

WeChat: cstutorcs  
Assignment Project Exam Help  
Email: [tutorcs@163.com](mailto:tutorcs@163.com)  
QQ: 749389476  
<https://tutorcs.com>

Processor 2

R		Row ID		
x	a	y	b	
2	3	0	1	2
1	2	7	2	
2	3	N	N	

Processor 3

# Parallel Join Query Processing



- Parallel Outer Join processing methods
  - ROJA (Redistribution Outer Join Algorithm)
  - DOJA (Duplication Outer Join Algorithm)
  - DER (Duplication & Efficient Redistribution)

Assignment Project Exam Help

- Load Balancing
  - OJSO (Outer Join Skew Optimization)

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

WeChat: cstutorcs

Email: tutorcs@163.com

	ROJA	DOJA	DER
Steps	<p><b>Step 1:</b> Distribute or reshuffle the data based on the join attribute.</p> <p><b>Step 2:</b> Each processor performs the Local outer Join.</p>	<p><b>Step 1:</b> Replication. We duplicate the small table.</p> <p><b>Step 2:</b> Local Inner Join</p>  <p><b>Step 3:</b> Redistribute the inner join result based on attribute X.</p> <p><b>Step 4:</b> Local outer join</p> <p>WeChat: cstutorcs</p> <p>Assignment Project Exam Help</p>	<p><b>Step 1:</b> Replication. We broadcast the left table.</p> <p><b>Step 2:</b> Local Inner Join</p> <p><b>Step 3:</b> Select the ROW ID of left table with no matches.</p> <p><b>Step 4:</b> Redistribute the ROW ID.</p> <p><b>Step 5:</b> Store the ROW ID that appears as many times as the number of processors.</p> <p><b>Step 6:</b> Inner join</p>
Pros	fast performance, only two steps	<p>Note: ROJA is faster than DOJA.</p> <p>Email: <a href="mailto:tutorcs@163.com">tutorcs@163.com</a></p>	Redistributes dangling row IDs instead of actual records.
Cons	redistribution of data -> data skew, communication cost	<p>In the replication step, if the table is large, the replication cost is expensive.</p> <p><a href="https://tutorcs.com">https://tutorcs.com</a></p> <p>In the distribution step, data skew and communication cost similar to ROJA</p>	In the replication step, if the table is large, the replication cost is expensive.

程序代写代做 CS编程辅导

## Another example...

Select x, y, z, a  
From R left outer join S on R.a=S.b  
left outer join T on S.c=T.d;



- Initial Data Placement

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Figure 1: Three relations  $R$ ,  $S$  and  $T$  are hash partitioned on a three parallel-unit system. The partitioning columns are  $R.x$ ,  $S.y$  and  $T.z$  respectively. The hash function,  $h(i) = i \bmod 3 + 1$ , places a tuple with value  $i$  in the partitioning column on the  $h(i)$ -th PU.

程序代写代做 CS编程辅导

## Another example...

- Step 1: Redistribution of  $R$  and  $S$  (why do we need to redistribute?)



Figure 2: The result of hash redistributing  $R$  and  $S$  on their join attributes ( $R.a$  and  $S.b$ ) to two temporary tables  $R_{redis}$  and  $S_{redis}$ .

## Another example...

- Step 2: (a) Outer Join and store in J

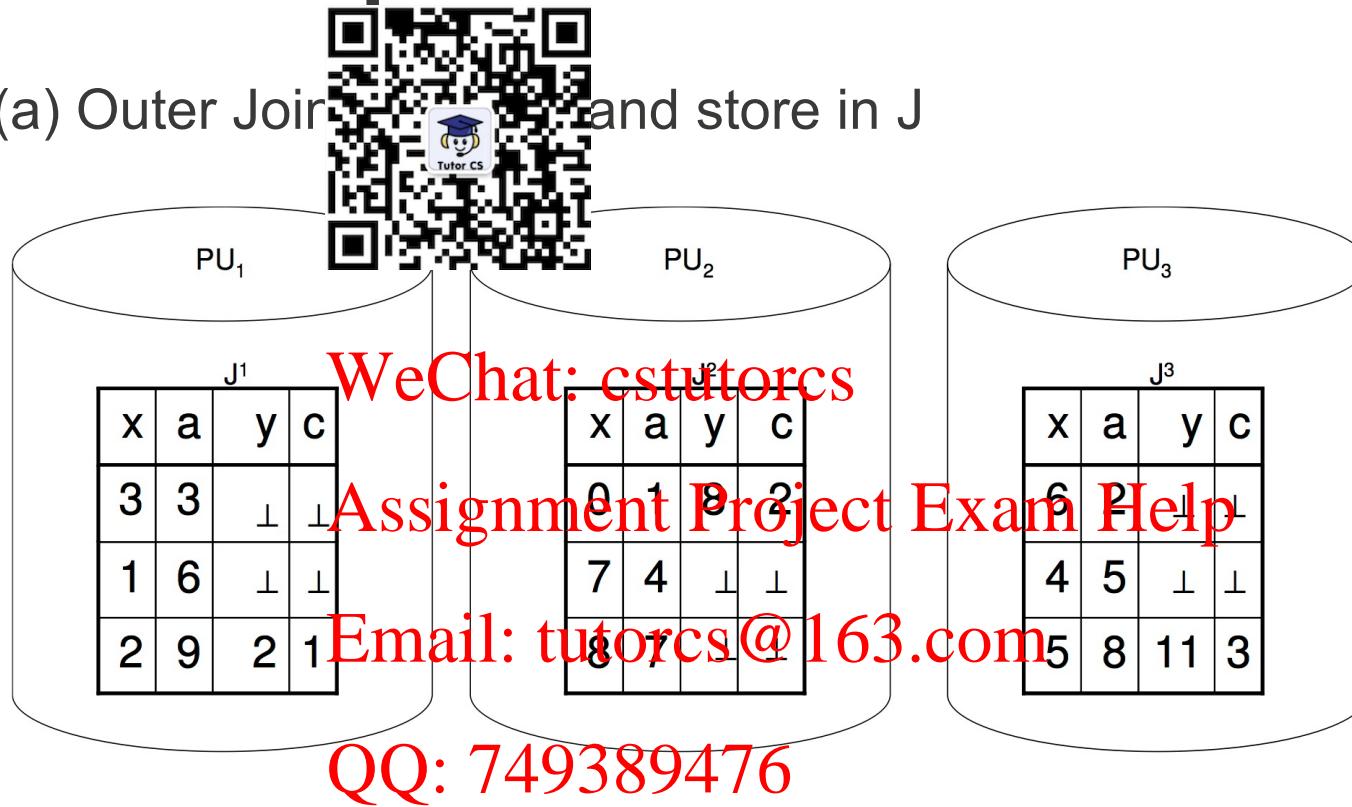


Figure 3: The results of left outer joining  $R_{redis}$  and  $S_{redis}$  ( $R_{redis}$  and  $S_{redis}$  are shown in Figure 2) are stored in a temporary table  $J$ .

## Another example...

- Step 2: (b) Redistribu

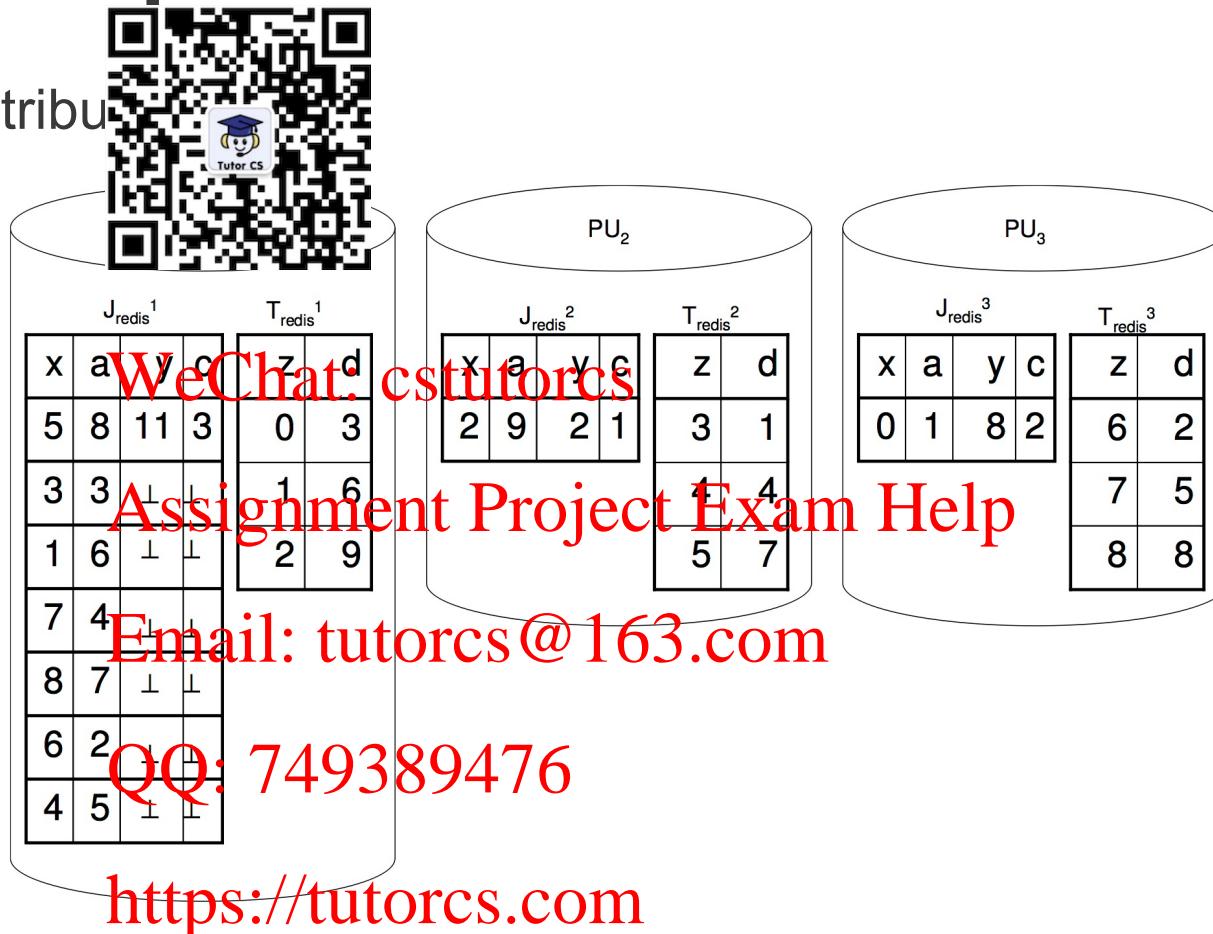


Figure 4: The result of hash redistributing  $J$  (shown in Figure 3) and  $T$  (shown in Figure 1) on their join attributes ( $J.c$  and  $T.d$ ) to two temporary tables  $J_{redis}$  and  $T_{redis}$ .

## Another example...

- Step 3: Outer Join J<sub>2</sub> and J<sub>3</sub> to get final Results

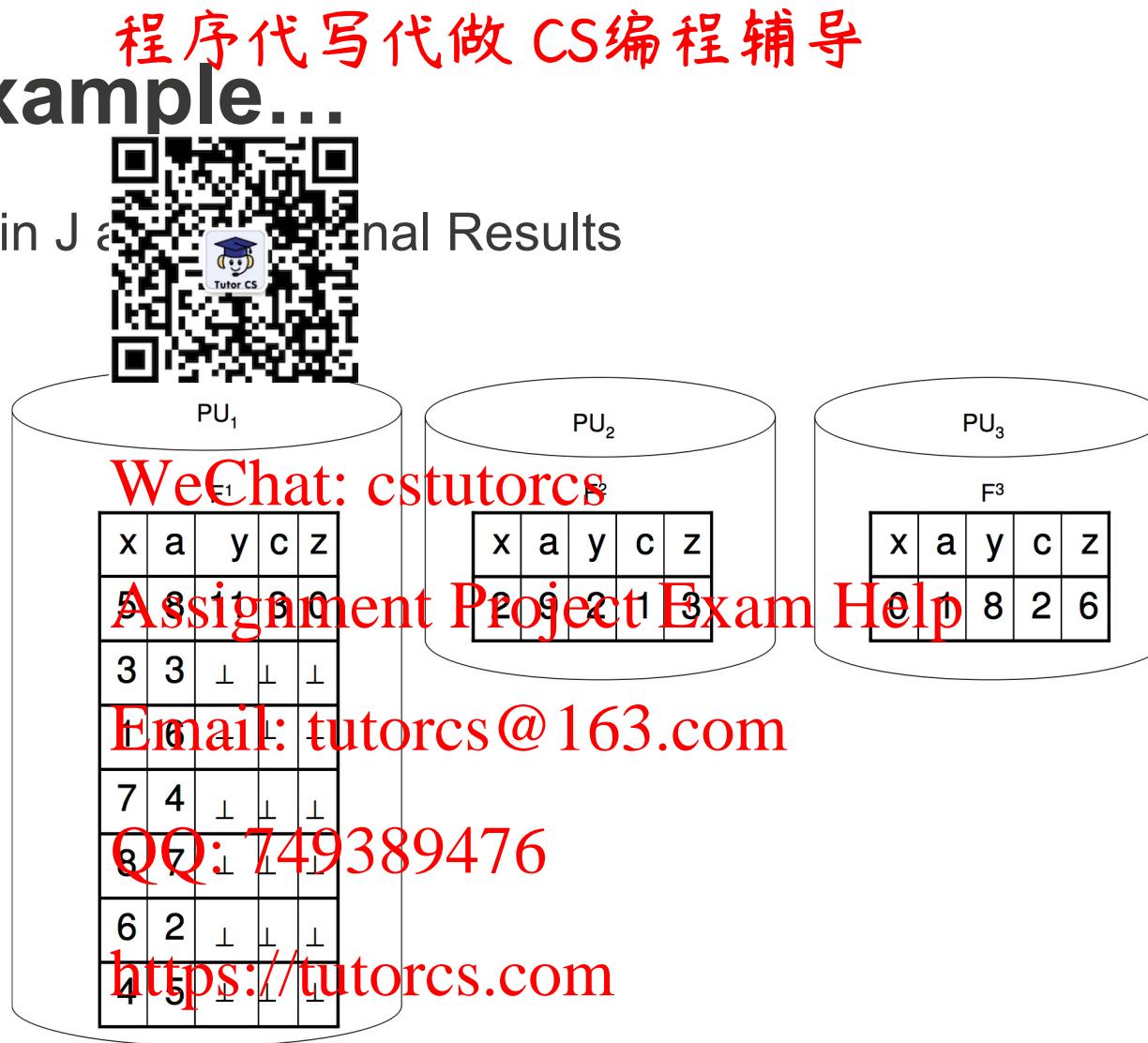


Figure 5: The final results of the two outer joins in Query 1 are stored in a temporary table  $F$ .

# Conclusion...

- Skew can easily happen in Outer Join queries



程序代写代做 CS编程辅导  
WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# 程序代写代做 CS编程辅导

## • Exercise 1 (FLUX Quiz)

- In the previous example, which outer join S outer join T), which parallel outer join method



WeChat: cstutorcs

- A. ROJA
- B. DOJA
- C. DER
- D. None of the above

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

# A better solution... OJSO

- Step 1: Redistribute Relations (Same as the previous example)



PU <sub>1</sub>		
R <sup>1</sup>	S <sup>1</sup>	T <sup>1</sup>
x a	y b c	z d
0 1	3 14 1	0 3
3 3	3 10 2	3 1
6 2	6 0 3	6 2

PU <sub>2</sub>		
R <sup>2</sup>	S <sup>2</sup>	T <sup>2</sup>
x a	y b c	z d
1 6	4 0 4	1 6
4 5	7 10 5	4 4
7 4	10 11 6	7 5

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

Figure 1: Three relations  $R$ ,  $S$  and  $T$  are hash partitioned on a three parallel-unit system. The partitioning columns are  $R.x$ ,  $S.y$  and  $T.z$  respectively. The hash function,  $h(i) = i \bmod 3 + 1$ , places a tuple with value  $i$  in the partitioning column on the  $h(i)$ -th PU.

<https://tutorcs.com>

PU <sub>1</sub>		
R <sub>redis</sub> <sup>1</sup>	S <sub>redis</sub> <sup>1</sup>	
x a	y b c	
3 3	6 0 3	
1 6	4 0 4	
2 9	2 9 1	
5 8	3 1 2	
8 7	11 8 3	
	8 8	

PU <sub>2</sub>		
R <sub>redis</sub> <sup>2</sup>	S <sub>redis</sub> <sup>2</sup>	
x a	y b c	
0 1	8 1 2	
7 4	7 10 5	
8 7	3 10 2	

PU <sub>3</sub>		
R <sub>redis</sub> <sup>3</sup>	S <sub>redis</sub> <sup>3</sup>	
x a	y b c	
6 2	11 8 3	
4 5	10 11 6	
5 8	3 14 1	

Figure 2: The result of hash redistributing  $R$  and  $S$  on their join attributes ( $R.a$  and  $S.b$ ) to two temporary tables  $R_{redis}$  and  $S_{redis}$ .

# A better solution... OJSO

- Step 2: (a) Outer Join

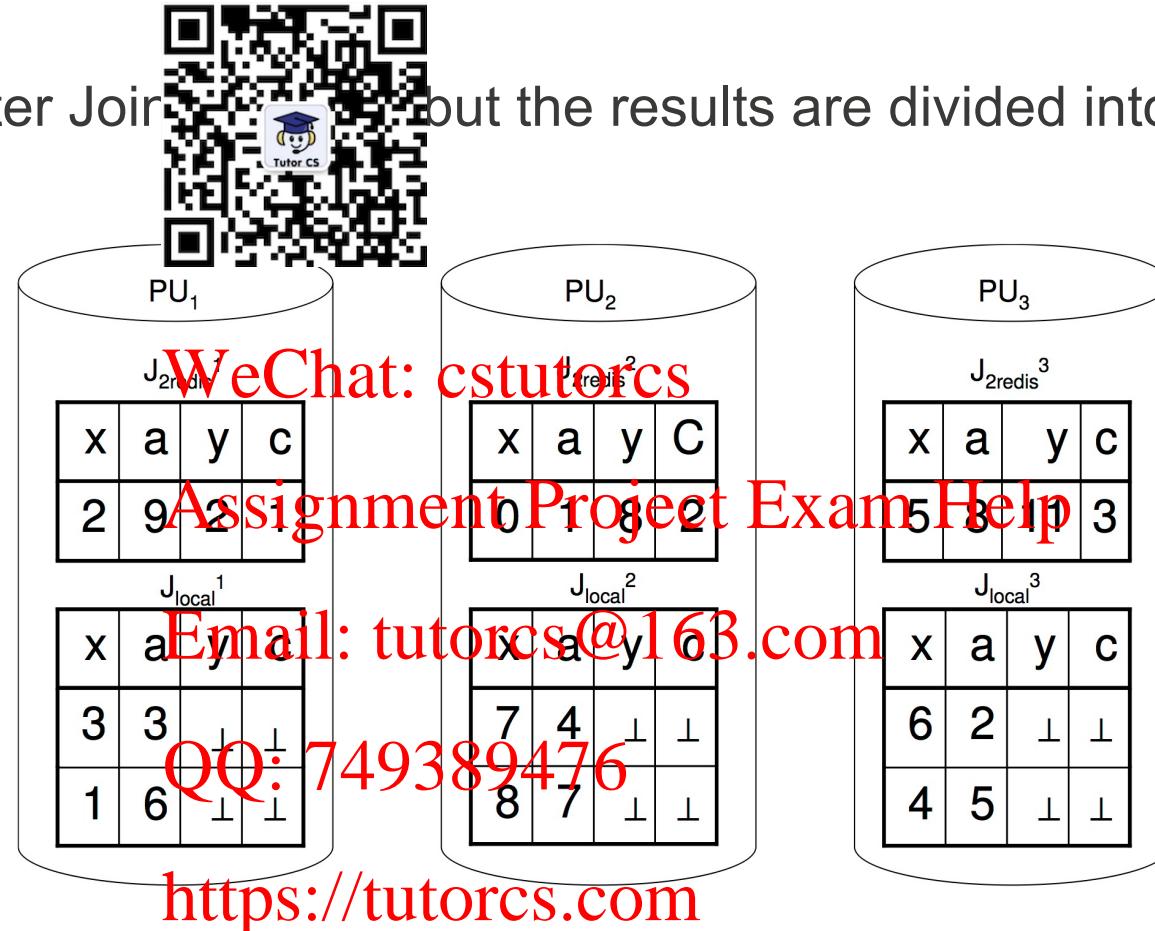
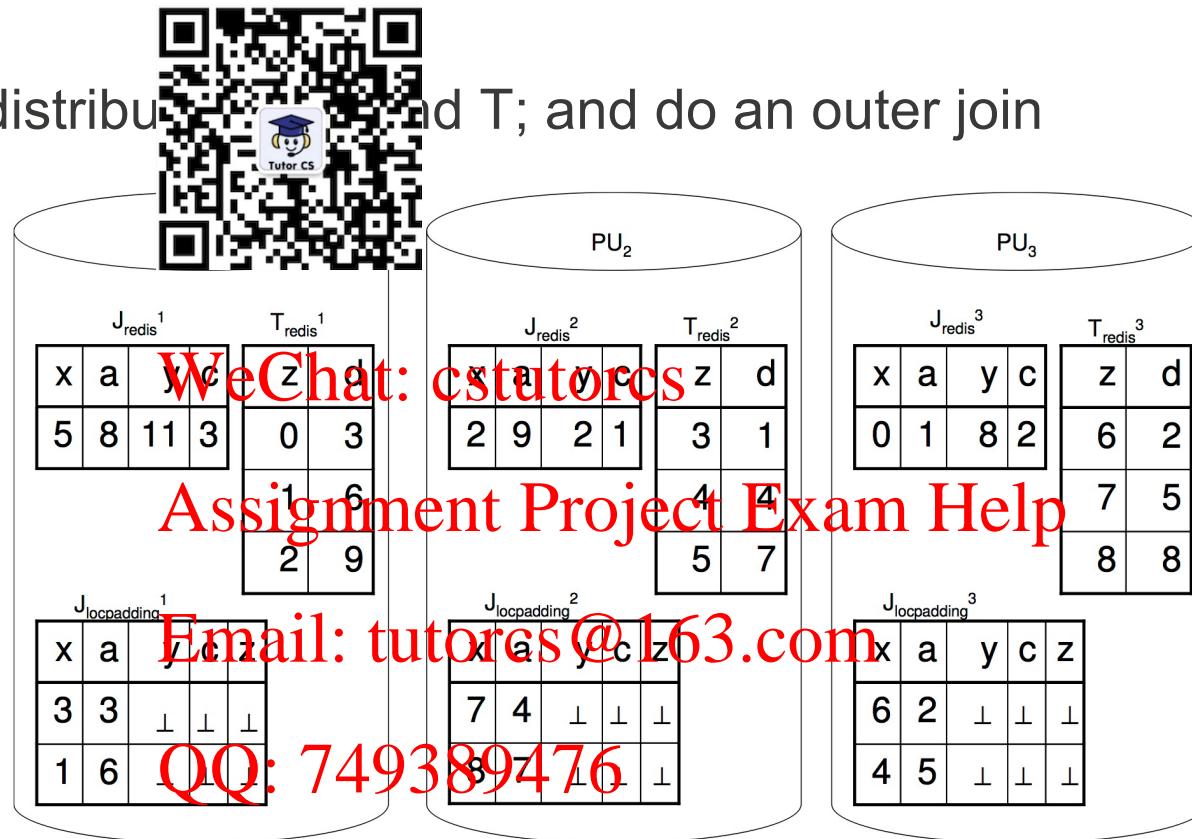


Figure 9: The results of left outer joining  $R_{redis}$  and  $S_{redis}$  are split into two temporary tables  $J_{2redis}$  and  $J_{local}$ .

# A better solution... OJSO

- Step 2: (b) Redistribute  $J_{2redis}$  and  $T$ ; and do an outer join



<https://tutorcs.com>

Figure 10: The result of hash redistributing  $J_{2redis}$  (shown in Figure 9) and  $T$  (shown in Figure 1) on their join attributes to two temporary tables  $J_{redis}$  and  $T_{redis}$ .  $J_{locpadding}$  is created from  $J_{local}$  (shown in Figure 9) with padded nulls.

# A better solution... OJSO

程序代写代做 CS编程辅导

- Step 3: Union the final results from each processor

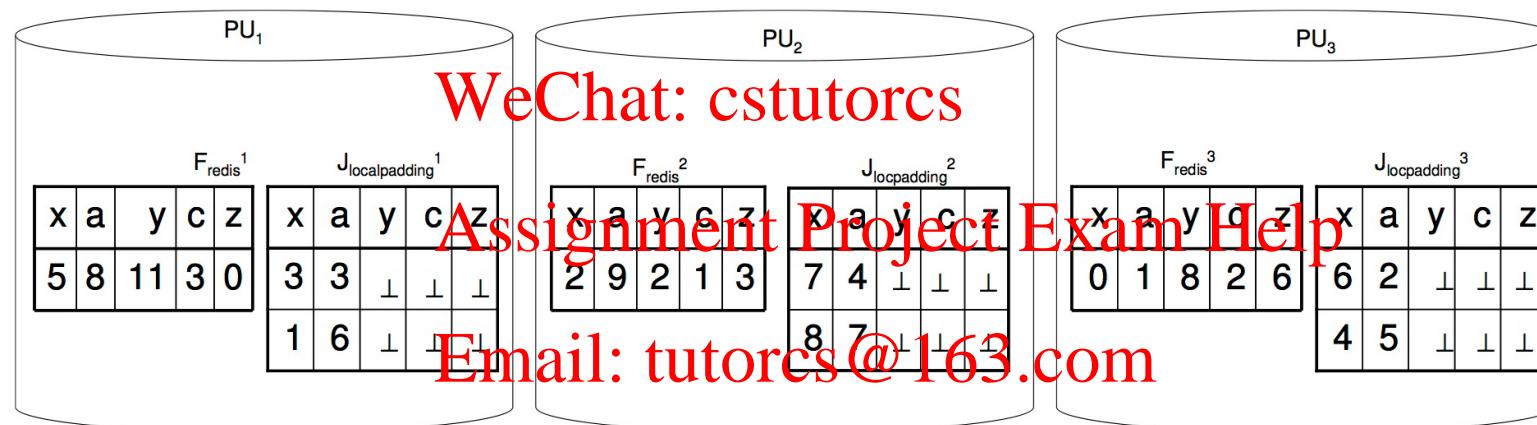


Figure 11: The results of the second outer join in Query 1 are stored in a temporary table  $F$ . The final result for Query 1 is the union of  $F$  and  $J_{locpadding}$ .

# OJSO Conclusion...

程序代写代做 CS编程辅导

- Do not redistribute the records from the previous outer join



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

# Summary...

程序代写代做 CS编程辅导



- Parallel Outer Join problem methods
  - **ROJA** (Redistribution Outer Join Algorithm)
  - **DOJA** (Duplication Outer Join Algorithm)
  - **DER** (Duplication & Efficient Redistribution)

Assignment Project Exam Help

- Load Balancing
  - **OJSO** (Outer Join Skew Optimization)

Email: tutorcs@163.com  
QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导

## References

- Xu, Y. & Kostamaa, P. (2010). A new algorithm for small-large table outer joins in parallel DBMS. In *Proceedings of the 26<sup>th</sup> Intl Conference on Data Engineering (ICDE'2010)* (pp. 1013-1024), IEEE Comp Society Press.
- Xu, Y. & Kostamaa, P. (2009). Efficient Outer Join Data Skew Handling in Parallel DBMS. In *Proceedings of the 35<sup>th</sup> International Conference on Very Large Data Bases (VLDB'2009)* (pp. 1390-1396), VLDB Endowment.



WeChat: cstutors

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>