



MONASH University

Information Technology

FIT5202 (Volume III - Join)

Week 3b – Parallel Outer Join

algorithm distributed systems **database**
systems **computation** knowledge ma
design e-business **model** data mining int
distributed systems **database** software
computation knowledge management an

Join Queries

- Two types of Join Queries

- **Inner Join**

```
Select R.x, R.a, S.y, S.b  
From R, S  
Where R.a = S.b;
```

- **Outer Join**

```
Select R.x, R.a, S.y, S.b  
From R left outer join S on R.a = S.b;
```

Join Queries

R	
x	a
0	1
1	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

- **Inner Join**

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

From R, S

Where R.a = S.b;

Join Queries

R	
x	a
0	1
1	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

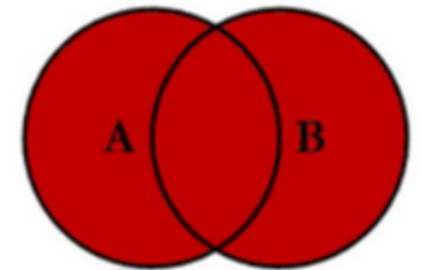
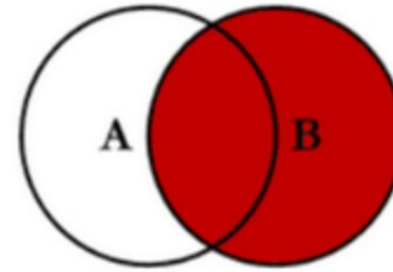
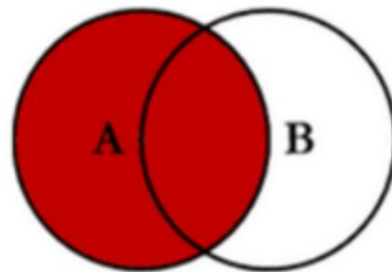
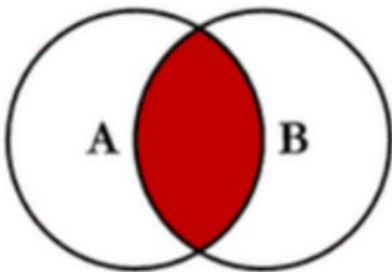
Results			
x	a	y	b
0	1	3	1
0	1	7	1
0	1	2	1
1	2	6	2
1	2	4	2
1	2	7	2
2	3	Null	Null

- Outer Join**

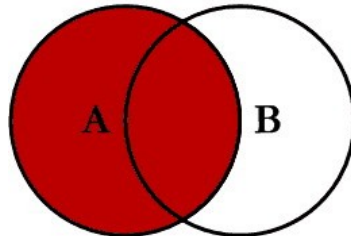
Select R.x, R.a, S.y, S.b
From R left outer join S
On R.a = S.b;

Exercise 1

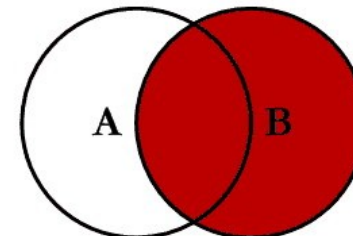
- Identify the LEFT OUTER JOIN?



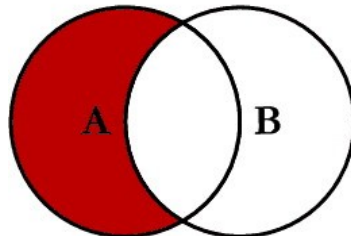
SQL JOINS



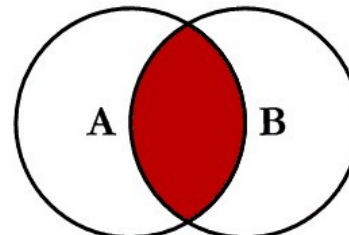
```
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
```



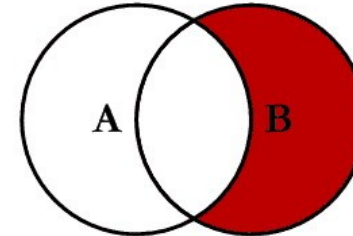
```
SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
```



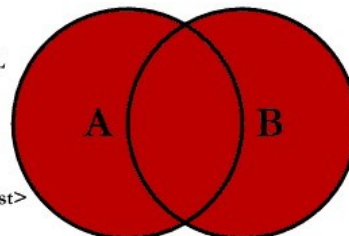
```
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
WHERE B.Key IS NULL
```



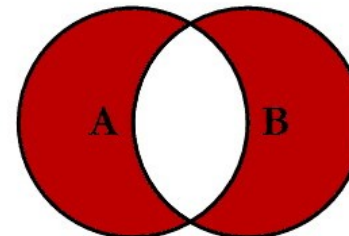
```
SELECT <select_list>
FROM TableA A
INNER JOIN TableB B
ON A.Key = B.Key
```



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



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



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

© C.L. Moffatt, 2008

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

Parallel Join Query Processing

- Parallel Inner Join components
 - **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**

Parallel Join Query Processing

- Parallel Outer Join processing methods
 - **ROJA** (Redistribution Outer Join Algorithm)
 - **DOJA** (Duplication Outer Join Algorithm)
 - **DER** (Duplication & Efficient Redistribution)
- Load Balancing
 - **OJSO** (Outer Join Skew Optimization)

1. ROJA

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

Step 2: Each processor performs local outer Join.

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

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

$h(i) = 1$

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

Processor 1

$h(i) = 2$

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

Processor 2

$h(i) = 3$

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

Processor 3

2. DOJA

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

Step 1: Replicate small table.

Step 2: Local Inner Join

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

<i>R</i>		<i>S</i>	
x	a	y	b
0	1	0	6
1	2	0	4
2	3	3	1
		6	2

Processor 1

<i>R</i>		<i>S</i>	
x	a	y	b
1	2	1	6
0	1	4	2
4	3	7	2
		7	1

Processor 2

<i>R</i>		<i>S</i>	
x	a	y	b
2	3	2	1
1	2	5	5
0	1	5	6
		8	9

Processor 3

2. DOJA

Step 4: Local outer join

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

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

Processor 1

R		J			
x	a	x	a	y	b
1	2	1	2	6	2
		1	2	4	2
		1	2	7	2

Processor 2

R		J			
x	a	x	a	y	b
2	3			N	N

Processor 3

3. DER

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

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	1	6
1	1	2	4	2
2	2	3	7	2
			7	1

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

Parallel Join Query Processing

- Parallel Outer Join processing methods
 - **ROJA** (Redistribution Outer Join Algorithm)
 - **DOJA** (Duplication Outer Join Algorithm)
 - **DER** (Duplication & Efficient Redistribution)
- Load Balancing
 - **OJSO** (Outer Join Skew Optimization)

	ROJA	DOJA	DER
Steps	<p>Step 1: Distribute or reshuffle the data based on the join attribute.</p> <p>Step 2: Each processor performs the Local outer Join.</p>	<p>Step 1: Replication. We duplicate the small table.</p> <p>Step 2: Local Inner Join</p> <p>Step 3: Hash redistribute the inner join result based on attribute X.</p> <p>Step 4: Local outer join</p>	<p>Step 1: Replication. We broadcast the left table.</p> <p>Step 2: Local Inner Join</p> <p>Step 3: Select the ROW ID of left table with no matches.</p> <p>Step 4: Redistribute the ROW ID.</p> <p>Step 5: Store the ROW ID that appears as many times as the number of processors.</p> <p>Step 6: Inner join</p>
Pros	fast performance, only two steps	None. ROJA is faster than DOJA.	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>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.

Another example...

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

- Initial Data Placement

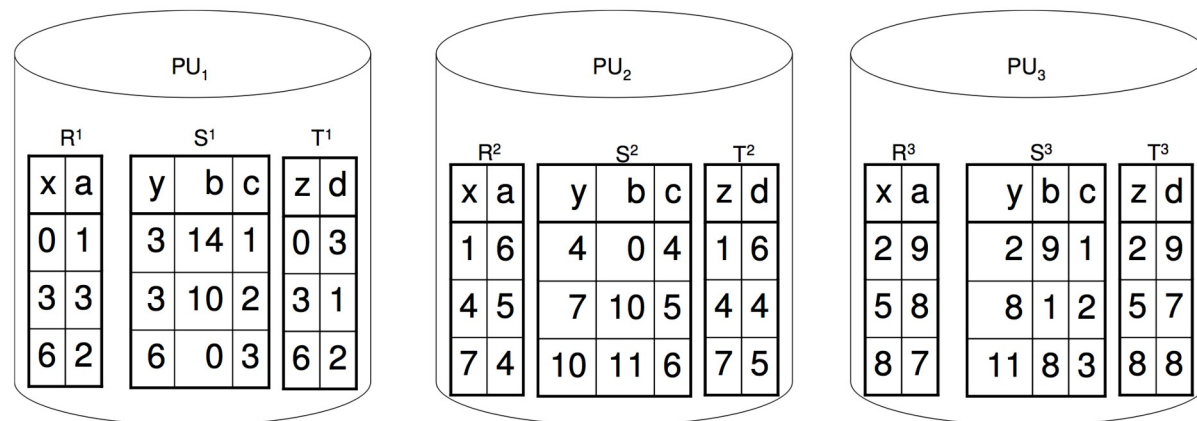


Figure 1: Three relations R , S and T are hash partitioned on a three parallel-unit system. The partitioning columns are $R.x$, $S.y$ and $T.z$ respectively. The hash function, $h(i) = i \bmod 3 + 1$, places a tuple with value i in the partitioning column on the $h(i)$ -th PU.

Another example...

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

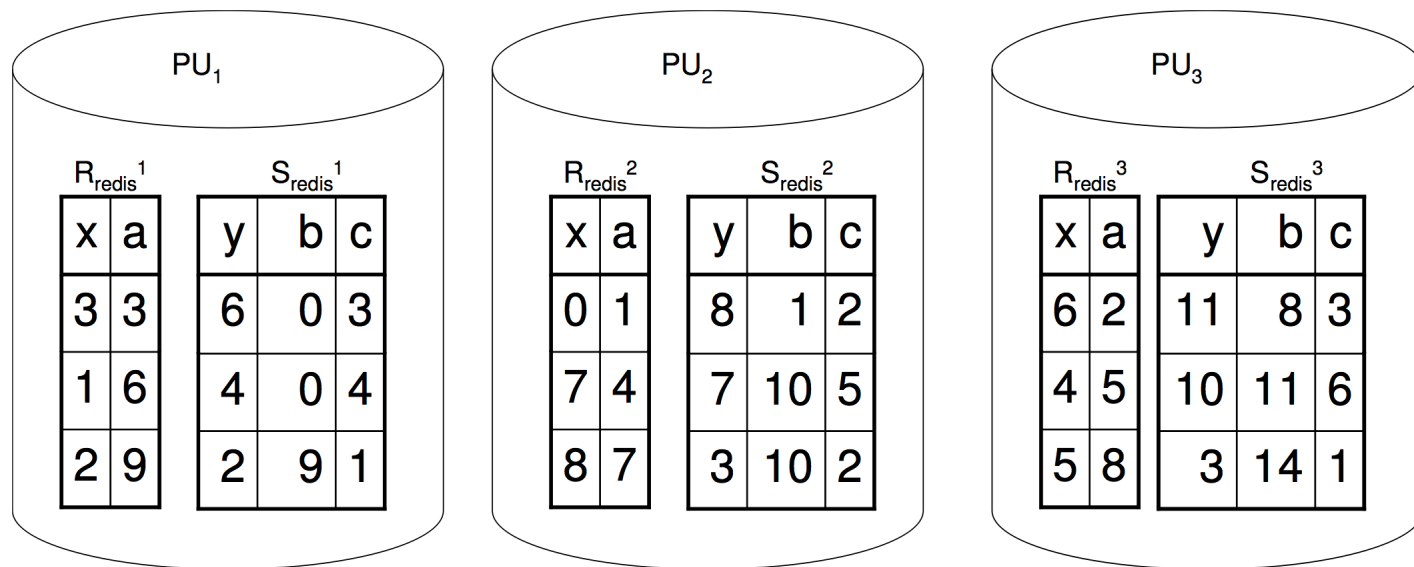


Figure 2: The result of hash redistributing R and S on their join attributes ($R.a$ and $S.b$) to two temporary tables R_{redis} and S_{redis} .

Another example...

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

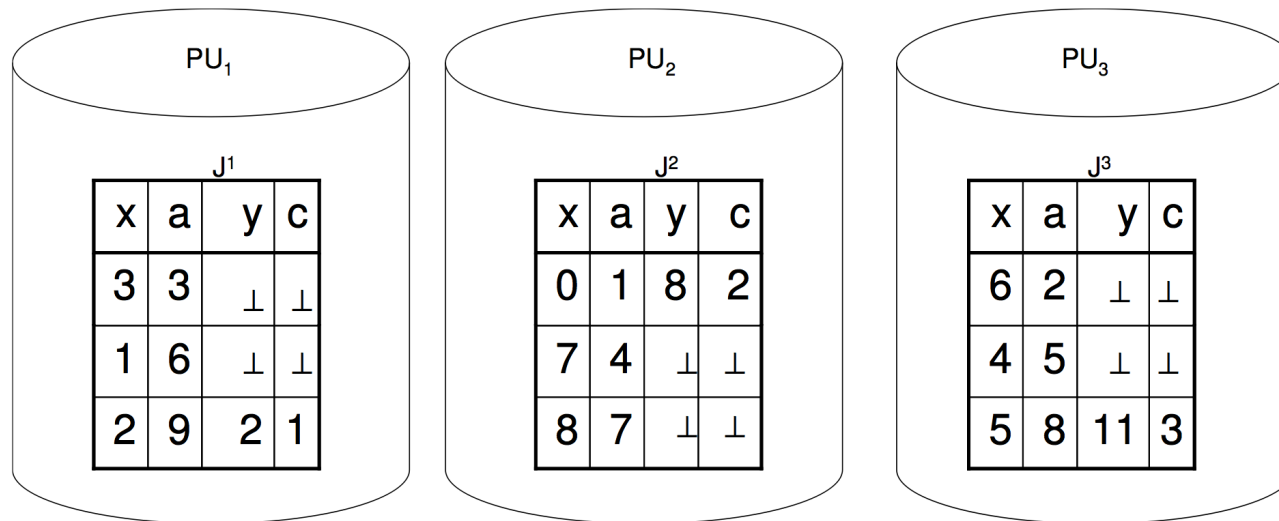


Figure 3: The results of left outer joining R_{redis} and S_{redis} (R_{redis} and S_{redis} are shown in Figure 2) are stored in a temporary table J .

Another example...

- Step 2: (b) Redistribute J and T

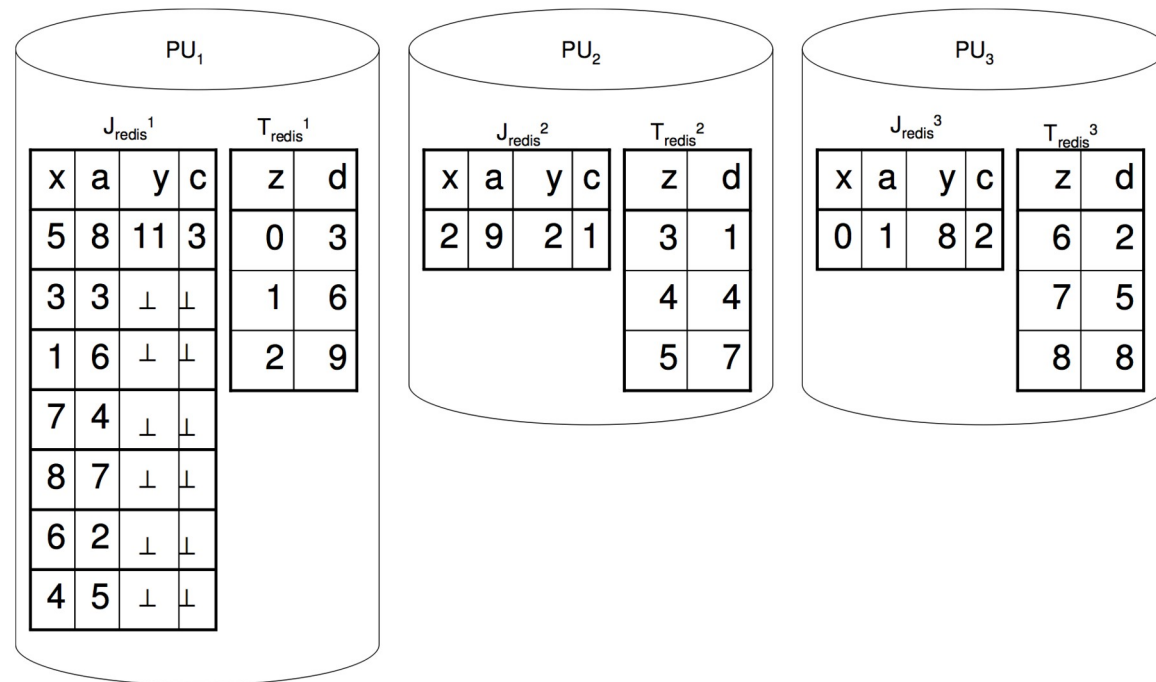


Figure 4: The result of hash redistributing J (shown in Figure 3) and T (shown in Figure 1) on their join attributes ($J.c$ and $T.d$) to two temporary tables J_{redis} and T_{redis} .

Another example...

- Step 3: Outer Join J and T → Final Results

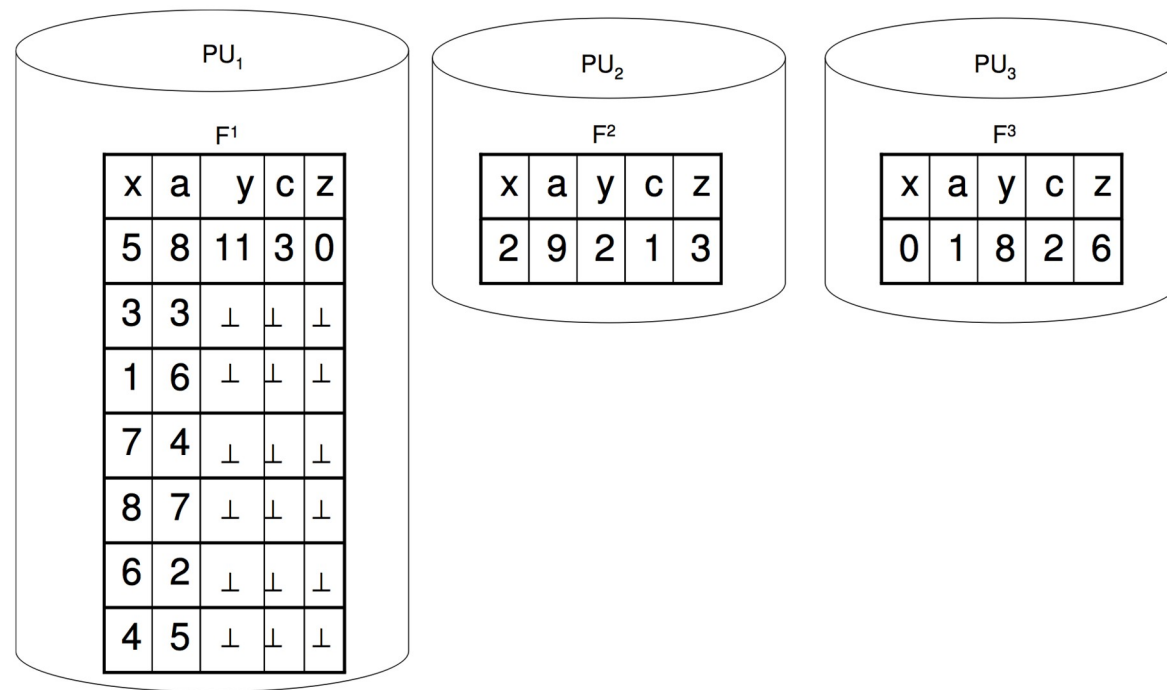


Figure 5: The final results of the two outer joins in Query 1 are stored in a temporary table F .

Conclusion...

- Skew can easily happen easily in Outer Join queries

• **Exercise 1 (FLUX Quiz)**

- In the previous example (R outer join S outer join T), which parallel outer join method was used?
- A. ROJA
- B. DOJA
- C. DER
- D. None of the above

A better solution... OJSO

- Step 1: Redistribute R and S (same as the previous example)

Figure 1 shows three parallel units (PU₁, PU₂, PU₃) each containing three tables: R¹, S¹, T¹ on PU₁; R², S², T² on PU₂; and R³, S³, T³ on PU₃. The tables are hash-partitioned based on their first attribute (x, y, z).

R ¹	S ¹	T ¹
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

R ²	S ²	T ²
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

R ³	S ³	T ³
x a	y b c	z d
2 9	2 9 1	2 9
5 8	8 1 2	5 7
8 7	11 8 3	8 8

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.

Figure 2 shows 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} . The redistribution is done across the three parallel units.

R _{redis} ¹	S _{redis} ¹
x a	y b c
3 3	6 0 3
1 6	4 0 4
2 9	2 9 1

R _{redis} ²	S _{redis} ²
x a	y b c
0 1	8 1 2
7 4	7 10 5
8 7	3 10 2

R _{redis} ³	S _{redis} ³
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 R and S, but the results are divided into J_{2redis} and J_{local}

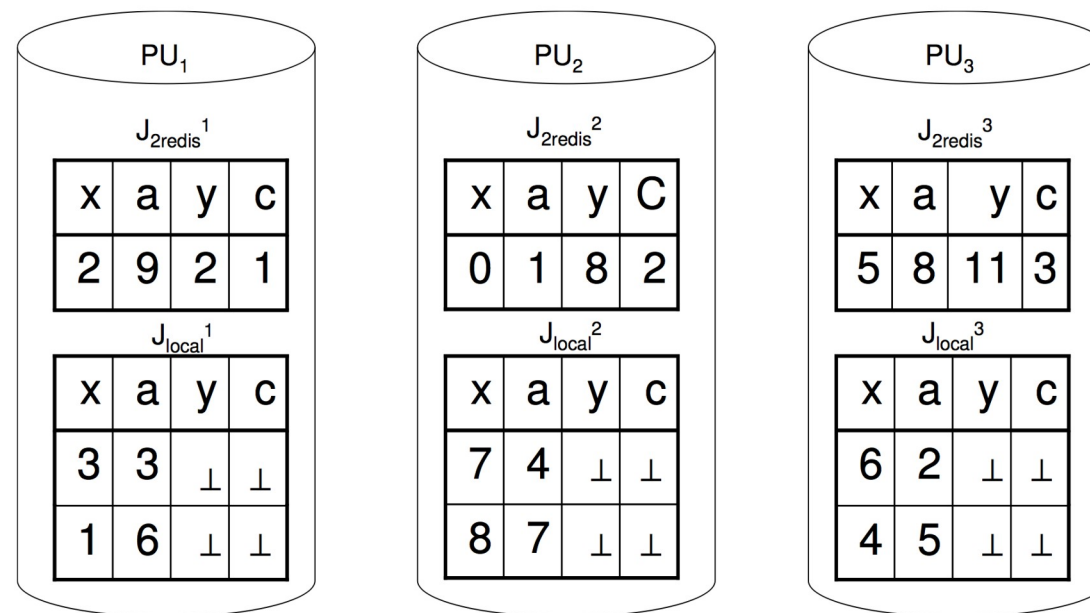


Figure 9: The results of left outer joining R_{redis} and S_{redis} are split into two temporary tables J_{2redis} and J_{local} .

A better solution... OJSO

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

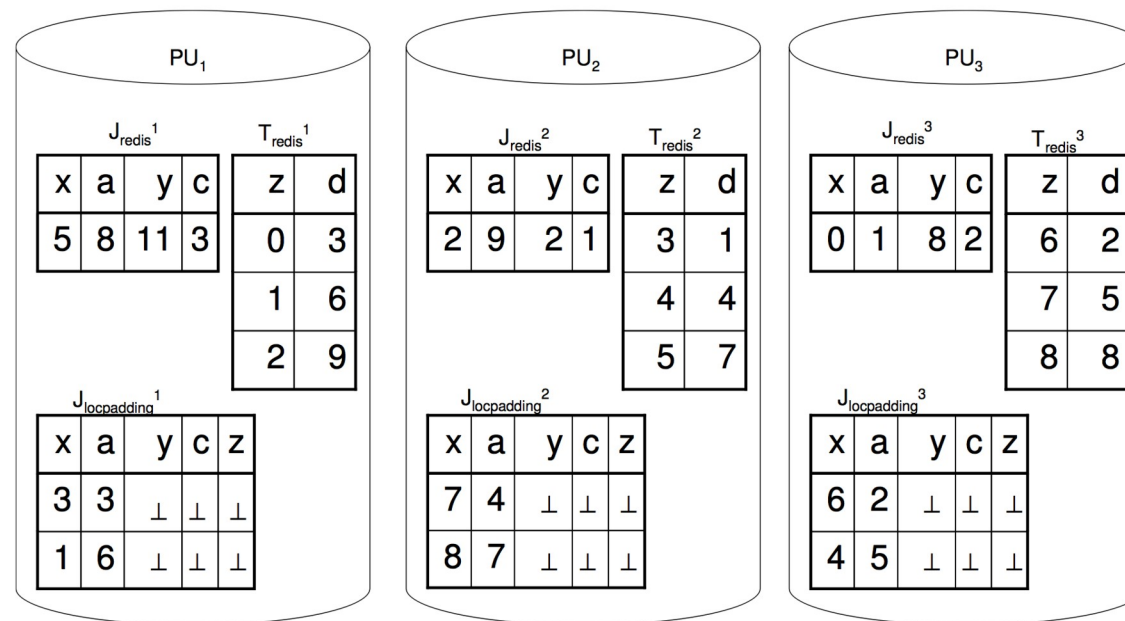


Figure 10: The result of hash redistributing J_{2redis} (shown in Figure 9) and T (shown in Figure 1) on their join attributes to two temporary tables J_{redis} and T_{redis} . $J_{locpadding}$ is created from J_{local} (shown in Figure 9) with padded nulls.

A better solution... OJSO

- Step 3: Union the final results in each processor

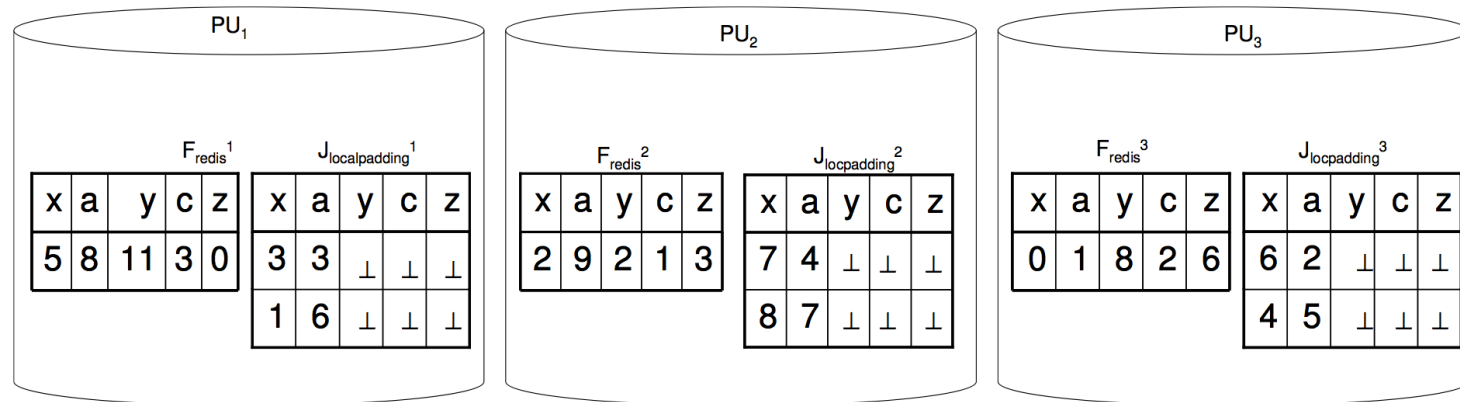


Figure 11: The results of the second outer join in Query 1 are stored in a temporary table F . The final result for Query 1 is the union of F and $J_{localpadding}$.

OJSO Conclusion...

- Do not redistribute the dangling records from the previous outer join

Summary...

- Parallel Outer Join processing methods
 - **ROJA** (Redistribution Outer Join Algorithm)
 - **DOJA** (Duplication Outer Join Algorithm)
 - **DER** (Duplication & Efficient Redistribution)
- Load Balancing
 - **OJSO** (Outer Join Skew Optimization)

References

- Xu, Y. & Kostamaa, P. (2010). A new algorithm for small-large table outer joins in parallel DBMS. In *Proceedings of the 26th Intl Conference on Data Engineering (ICDE'2010)* (pp. 1018-1024), IEEE Comp Society Press.
- Xu, Y. & Kostamaa, P. (2009). Efficient Outer Join Data Skew Handling in Parallel DBMS. In *Proceedings of the 35th International Conference on Very Large Data Bases (VLDB'2009)* (pp. 1390-1396), VLDB Endowment.