



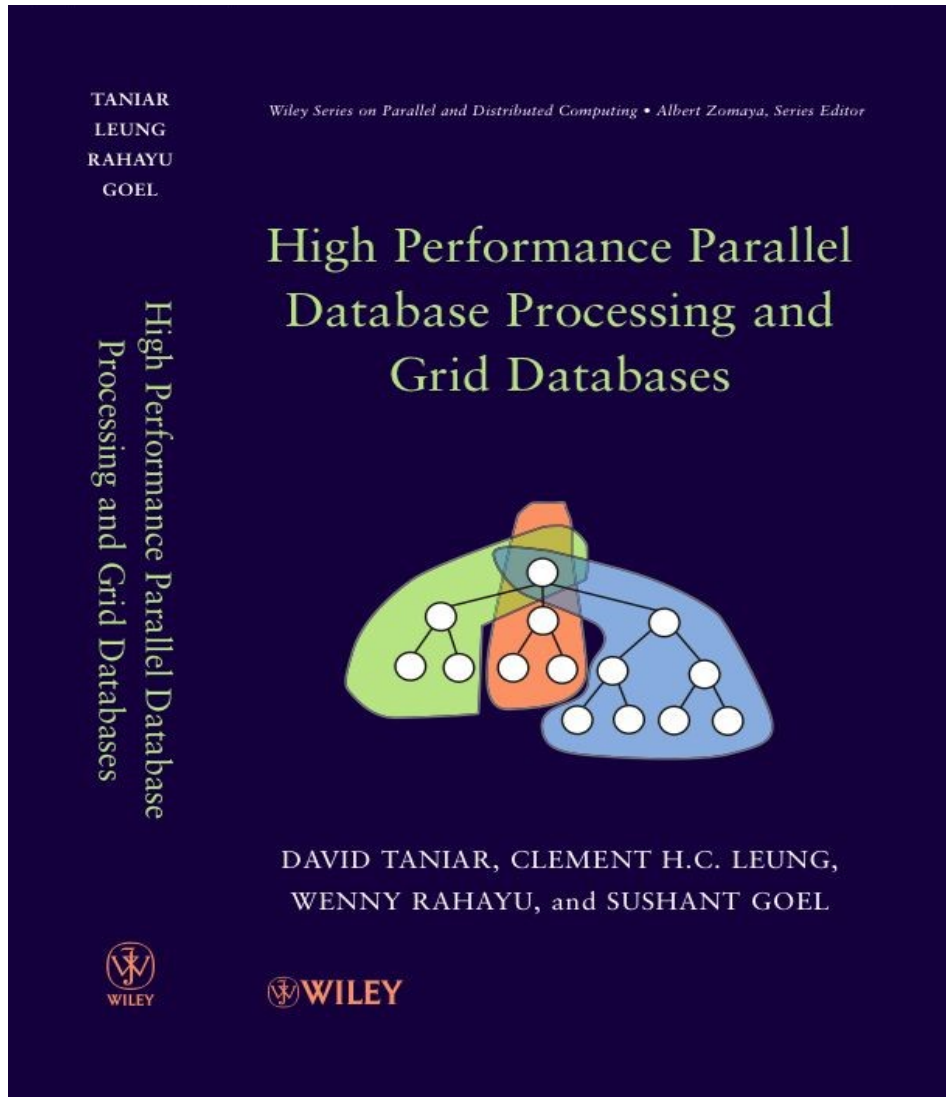
MONASH University

Information Technology

FIT5202 (Volume IV – Sort and Group By)

Week 4b – Parallel Group By

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



Chapter 4

Parallel Sort and GroupBy

- 4.1 Sorting, Duplicate Removal and Aggregate
- 4.2 Serial External Sorting Method
- 4.3 Algorithms for Parallel External Sort
- 4.4 Parallel Algorithms for GroupBy Queries
- 4.5 Cost Models for Parallel Sort
- 4.6 Cost Models for Parallel GroupBy
- 4.7 Summary
- 4.8 Bibliographical Notes
- 4.9 Exercises

4.1. GroupBy, and Serial GroupBy

Select Suburb, Count(*)
From Student
Group By Suburb;

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Chris	Doncaster
Daniel	Caulfield
Eric	Kew
Fred	Richmond
Garry	Hawthorn
Harold	Elwood
Irene	Clayton
Jessica	Caulfield
Katie	Malvern
Leonard	Balwyn
Mary	Hawthorn

4.1. Serial GroupBy Processing (cont'd)

Select Suburb, Count(*)
From Student
Group By Suburb;

Processing Steps:

1. Read the first student record, and hash the suburb to the hash table

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Chris	Doncaster
Daniel	Caulfield
Eric	Kew
Fred	Richmond
Garry	Hawthorn
Harold	Elwood
Irene	Clayton
Jessica	Caulfield
Katie	Malvern
Leonard	Balwyn
Mary	Hawthorn

Hash
the
record
using a
certain
hash
function

Hash Table

1		
2		
3		
4		
5		
6		
7		
8	Clayton	1
9		

4.1. Serial GroupBy Processing (cont'd)

Select Suburb, Count(*)
From Student
Group By Suburb;

Processing Steps:

1. Read the first student record, and hash the suburb to the hash table
2. Read the second record and hash it

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Chris	Doncaster
Daniel	Caulfield
Eric	Kew
Fred	Richmond
Garry	Hawthorn
Harold	Elwood
Irene	Clayton
Jessica	Caulfield
Katie	Malvern
Leonard	Balwyn
Mary	Hawthorn

Hash Table

1	Hawthorn	1
2		
3		
4		
5		
6		
7		
8	Clayton	1
9		

4.1. Serial GroupBy Processing (cont'd)

Select Suburb, Count(*)
From Student
Group By Suburb;

Processing Steps:

1. Read the first student record, and hash the suburb to the hash table
2. Read the second record and hash it
3. Read the subsequent records one-by-one and hash them

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Chris	Doncaster
Daniel	Caulfield
Eric	Kew
Fred	Richmond
Garry	Hawthorn
Harold	Elwood
Irene	Clayton
Jessica	Caulfield
Katie	Malvern
Leonard	Balwyn
Mary	Hawthorn

Hash Table

1	Hawthorn	3
2	Caulfield	2
3	Malvern	1
4	Balwyn	1
5	Kew	1
6	Richmond	1
7	Elwood	1
8	Clayton	2
9	Doncaster	1

4.1. Serial GroupBy Processing (cont'd)

Select Suburb, Count(*)
From Student
Group By Suburb;

Processing Steps:

1. Read the first student record, and hash the suburb to the hash table
2. Read the second record and hash it
3. Read the subsequent records one-by-one and hash them
4. Read the Hash Table, and store this in disk as the query results

Query Results in Disk

Hawthorn	3
Caulfield	2
Malvern	1
Balwyn	1
Kew	1
Richmond	1
Elwood	1
Clayton	2
Doncaster	1

Hash Table in Main-Memory

1	Hawthorn	3
2	Caulfield	2
3	Malvern	1
4	Balwyn	1
5	Kew	1
6	Richmond	1
7	Elwood	1
8	Clayton	2
9	Doncaster	1

4.1. Serial GroupBy Processing (cont'd)

Select Suburb, Count(*)
From Student
Group By Suburb;

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Clayton	Clayton

This will work, if we assume that the main-memory can hold the entire Hash Table.

How about if the Hash Table is so big that it cannot fit into the main-memory.

For example, how about if the main-memory can only hold 4 hash records at a time? How does the Group By processing work?

Leonard	Balwyn
Mary	Hawthorn

Hash Table

1	Hawthorn	3
2	Caulfield	2
3	Malvern	1
4	Balwyn	1
5	Kew	1
6	Richmond	1
7	Elwood	1
8	Clayton	2
9	Doncaster	1

4.1. Serial GroupBy Processing (cont'd)

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Chris	Doncaster
Daniel	Caulfield
Eric	Kew
Fred	Richmond
Garry	Hawthorn
Harold	Elwood
Irene	Clayton
Jessica	Caulfield
Katie	Malvern
Leonard	Balwyn
Mary	Hawthorn

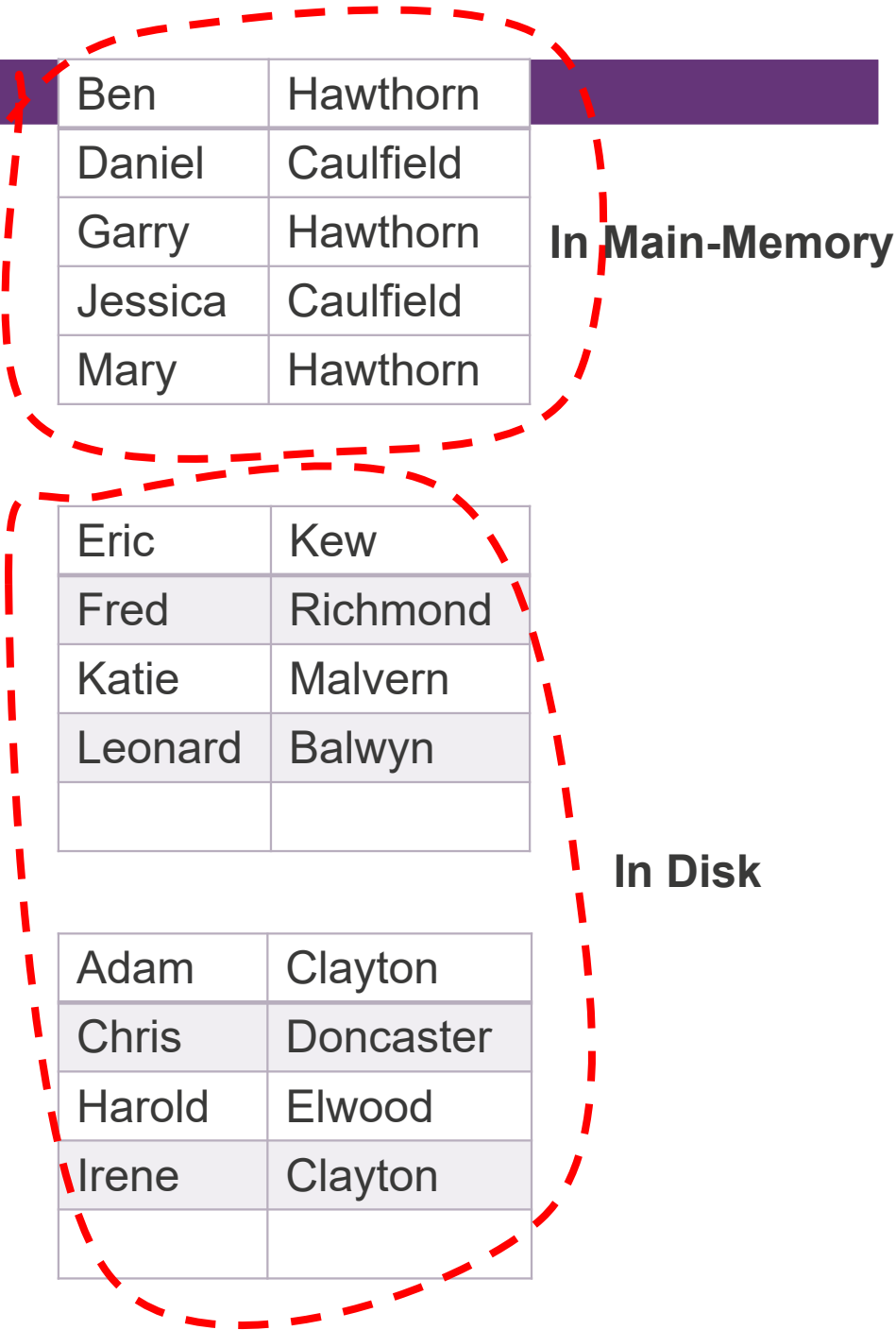
Hash Table

Assume that the main-memory can hold 4 records in the hash table.

It needs a bigger hash table, but it doesn't have.

Student	Suburb
Adam	Clayton
Ben	Hawthorn
Chris	Doncaster
Daniel	Caulfield
Eric	Kew
Fred	Richmond
Garry	Hawthorn
Harold	Elwood
Irene	Clayton
Jessica	Caulfield
Katie	Malvern
Leonard	Balwyn
Mary	Hawthorn

Hash Data
Partitioning
based on the
Suburb



Hash Processing

In Main-Memory

Ben	Hawthorn
Daniel	Caulfield
Garry	Hawthorn
Jessica	Caulfield
Mary	Hawthorn

Still in Disk

Eric	Kew
Fred	Richmond
Katie	Malvern
Leonard	Balwyn

Adam	Clayton
Chris	Doncaster
Harold	Elwood
Irene	Clayton

Hash Table in Main-Memory

Hawthorn	3
Caulfield	2

Query Results in Disk

Hawthorn	3
Caulfield	2

Hash Processing

Hash Table in
Main-Memory

Ben	Hawthorn
Daniel	Caulfield
Garry	Hawthorn
Jessica	Caulfield
Mary	Hawthorn

Flush to Disk

Eric	Kew
Fred	Richmond
Katie	Malvern
Leonard	Balwyn

Load to Main-
Memory

Adam	Clayton
Chris	Doncaster
Harold	Elwood
Irene	Clayton

Still in Disk

Kew	1
Malvern	1
Richmond	1
Balwyn	1

Query Results
in Disk

Hawthorn	3
Caulfield	2
Kew	1
Caulfield	1
Richmond	1
Balwyn	1

Hash Processing

Hash Table in
Main-Memory

Ben	Hawthorn
Daniel	Caulfield
Garry	Hawthorn
Jessica	Caulfield
Mary	Hawthorn

Flush to Disk

Eric	Kew
Fred	Richmond
Katie	Malvern
Leonard	Balwyn

Flush to Disk

Adam	Clayton
Chris	Doncaster
Harold	Elwood
Irene	Clayton

Load to Main-
Memory

Clayton	2
Doncaster	1
Elwood	1

Hawthorn	3
Caulfield	2
Kew	1
Caulfield	1
Richmond	1
Balwyn	1
Clayton	2
Doncaster	1
Elwood	1

Query Results
in Disk

4.4. Parallel GroupBy

- Traditional methods (Merge-All and Hierarchical Merging)
- Two-phase method
- Redistribution method

Without data
redistribution

With data
redistribution

4.4. Parallel GroupBy (cont'd)

Traditional Methods

- Step 1: local aggregate in each processor
- Step 2: global aggregation
- May use a Merge-All or Hierarchical method
- Need to pay a special attention to some aggregate functions (AVG) when performing a local aggregate process

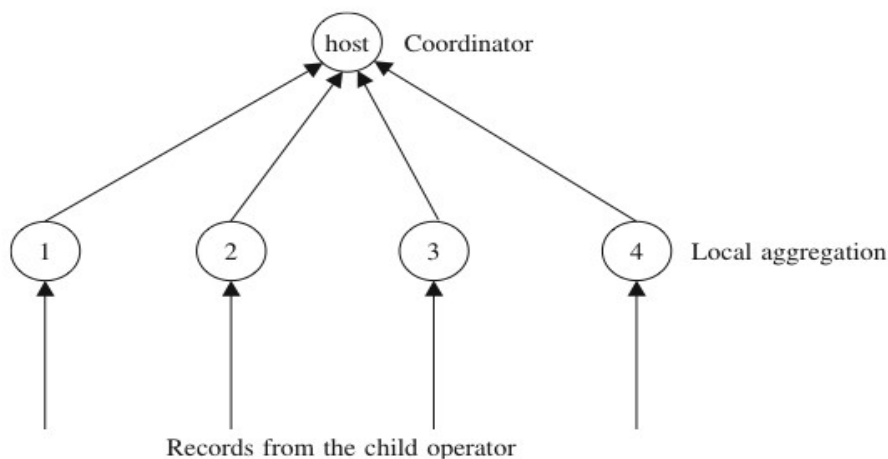


Figure 4.10 Traditional method

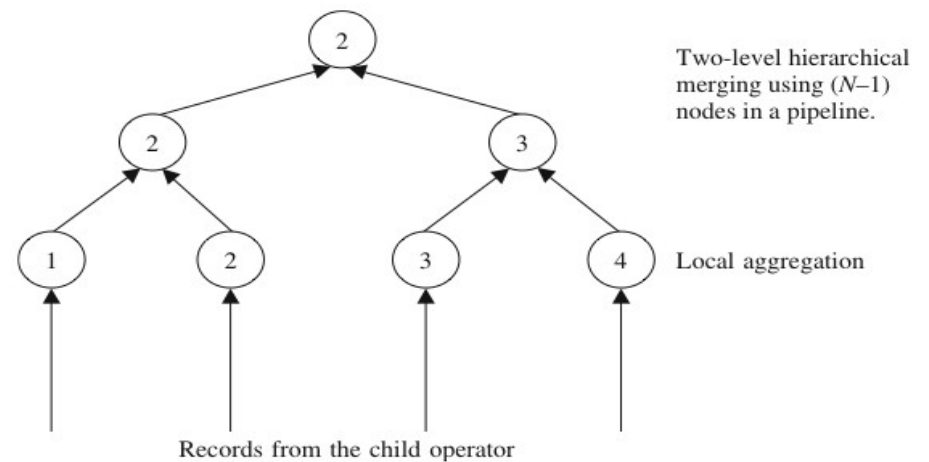
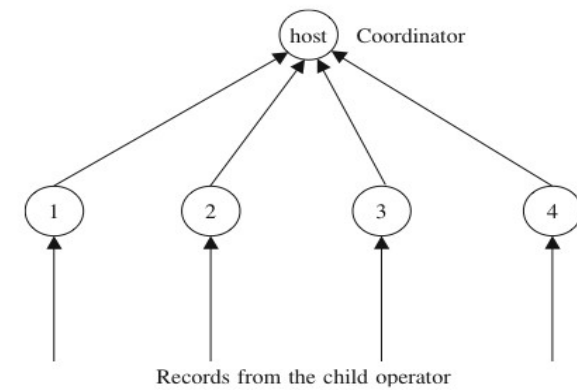


Figure 4.11 Hierarchical merging method

4.4. Parallel GroupBy (cont'd)

- Traditional Method: Merge All



Initial Data Placement

Processor 1

A dam	Clayton
B en	Clayton
C hris	Caulfield
D ennis	Malvern
E ric	Vermont

Processor 2

F red	Hawthorn
G eorge	Richmond
H arold	Elwood
I rene	Malvern
J essica	Kew

Processor 3

K elly	Balwyn
L esley	Hawthorn
M egan	Kew
N aomi	Richmond
O scar	Vermont

Processor 4

P eter	Elwood
Q uin	Kew
R oger	Balwyn
S arah	Malvern
T racy	Clayton

4.4. Parallel GroupBy (cont'd)

- Traditional Method: Merge All



Clayton	2
Caulfield	1
Malvern	1
Vermont	1

Hawthorn	1
Richmond	1
Elwood	1
Malvern	1
Kew	1

Balwyn	1
Hawthorn	1
Kew	1
Richmond	1
Vermont	1

Elwood	1
Kew	1
Balwyn	1
Malvern	1
Clayton	1

Local Aggregation Phase



Processor 1



Processor 2



Processor 3



Processor 4

Adam	Clayton
Ben	Clayton
Chris	Caulfield
Dennis	Malvern
Eric	Vermont

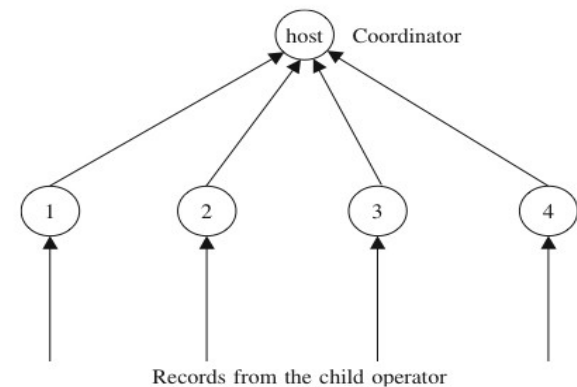
Fred	Hawthorn
George	Richmond
Harold	Elwood
Irene	Malvern
Jessica	Kew

Kelly	Balwyn
Lesley	Hawthorn
Megan	Kew
Naomi	Richmond
Oscar	Vermont

Peter	Elwood
Quin	Kew
Roger	Balwyn
Sarah	Malvern
Tracy	Clayton

4.4. Parallel GroupBy (cont'd)

- Traditional Method: Merge All



Global Aggregation Phase

Clayton	2
Hawthorn	1
Balwyn	1
Elwood	1
...	...
...	...
Vermont	1
Clayton	1

Clayton	2
Caulfield	1
Malvern	1
Vermont	1

Processor 1

Hawthorn	1
Richmond	1
Elwood	1
Malvern	1
Kew	1

Processor 2

Balwyn	1
Hawthorn	1
Kew	1
Richmond	1
Vermont	1

Processor 3

Elwood	1
Kew	1
Balwyn	1
Malvern	1
Clayton	1

Processor 4

4.4. Parallel GroupBy (cont'd)

- Traditional Method: Merge All

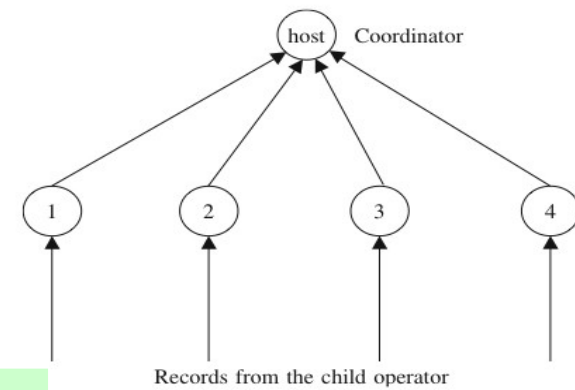
Clayton	3
Hawthorn	2
Balwyn	2
...	...
...	...
Vermont	2

Final results

Clayton	2
Hawthorn	1
Balwyn	1
Elwood	1
...	...
...	...
Vermont	1
Clayton	1



Global Aggregation Phase



• **Exercise 6 (FLUX Quiz)**

- The limitations of the Traditional Approach (Merge All) to process a Group By query are:
 - A. Global aggregation is carried out by one processor
 - B. Network bottleneck when sending the local aggregation results to the coordinator
 - C. No parallelism in the global aggregation phase
 - D. All of the above
 - E. Some of the above

4.4. Parallel GroupBy (cont'd)

▪ Two-Phase Method

- Step 1: local aggregate in each processor. Each processor groups local records according to the groupby attribute
- Step 2: global aggregation where all temp results from each processor are redistributed and then final aggregate is performed in each processor

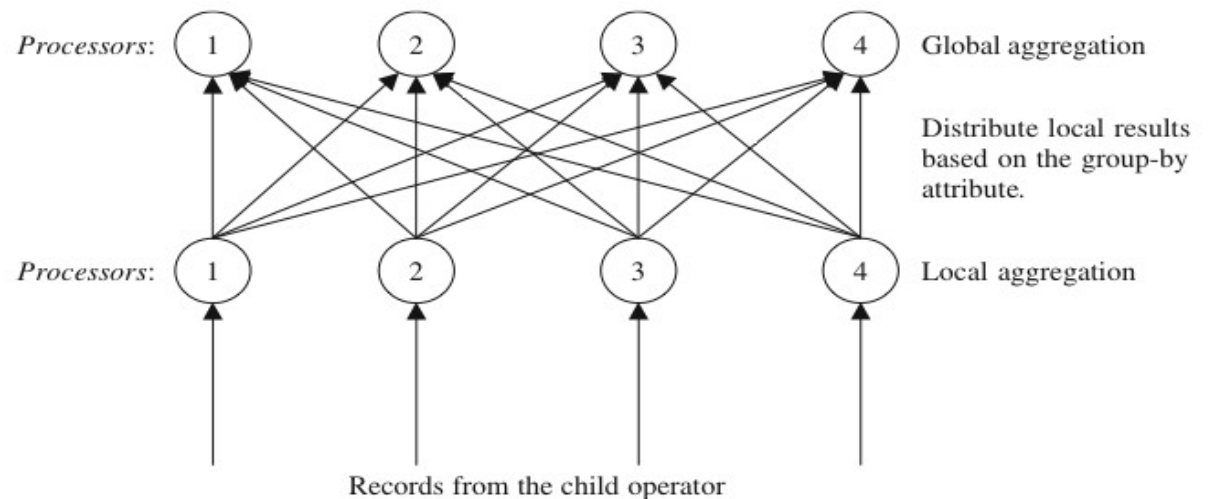
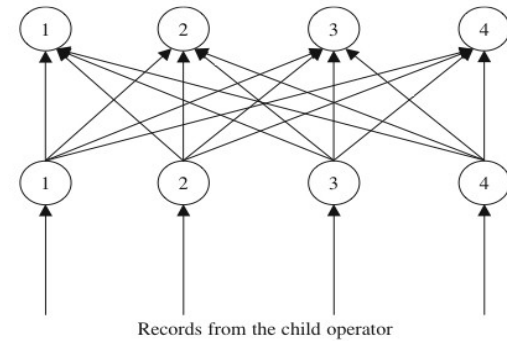


Figure 4.12 Two-phase method

4.4. Parallel GroupBy (cont'd)

- Two-Phase Method



Initial Data Placement

Processor 1

A dam	Clayton
B en	Clayton
C hris	Caulfield
D ennis	Malvern
E ric	Vermont

Processor 2

F red	Hawthorn
G eorge	Richmond
H arold	Elwood
I rene	Malvern
J essica	Kew

Processor 3

K elly	Balwyn
L esley	Hawthorn
M egan	Kew
N aomi	Richmond
O scar	Vermont

Processor 4

P eter	Elwood
Q uin	Kew
R oger	Balwyn
S arah	Malvern
T racy	Clayton

4.4. Parallel GroupBy (cont'd)

- Two-Phase Method



Clayton	2
Caulfield	1
Malvern	1
Vermont	1

Hawthorn	1
Richmond	1
Elwood	1
Malvern	1
Kew	1

Balwyn	1
Hawthorn	1
Kew	1
Richmond	1
Vermont	1

Elwood	1
Kew	1
Balwyn	1
Malvern	1
Clayton	1

Local Aggregation Phase



Processor 1



Processor 2



Processor 3



Processor 4

Adam	Clayton
Ben	Clayton
Chris	Caulfield
Dennis	Malvern
Eric	Vermont

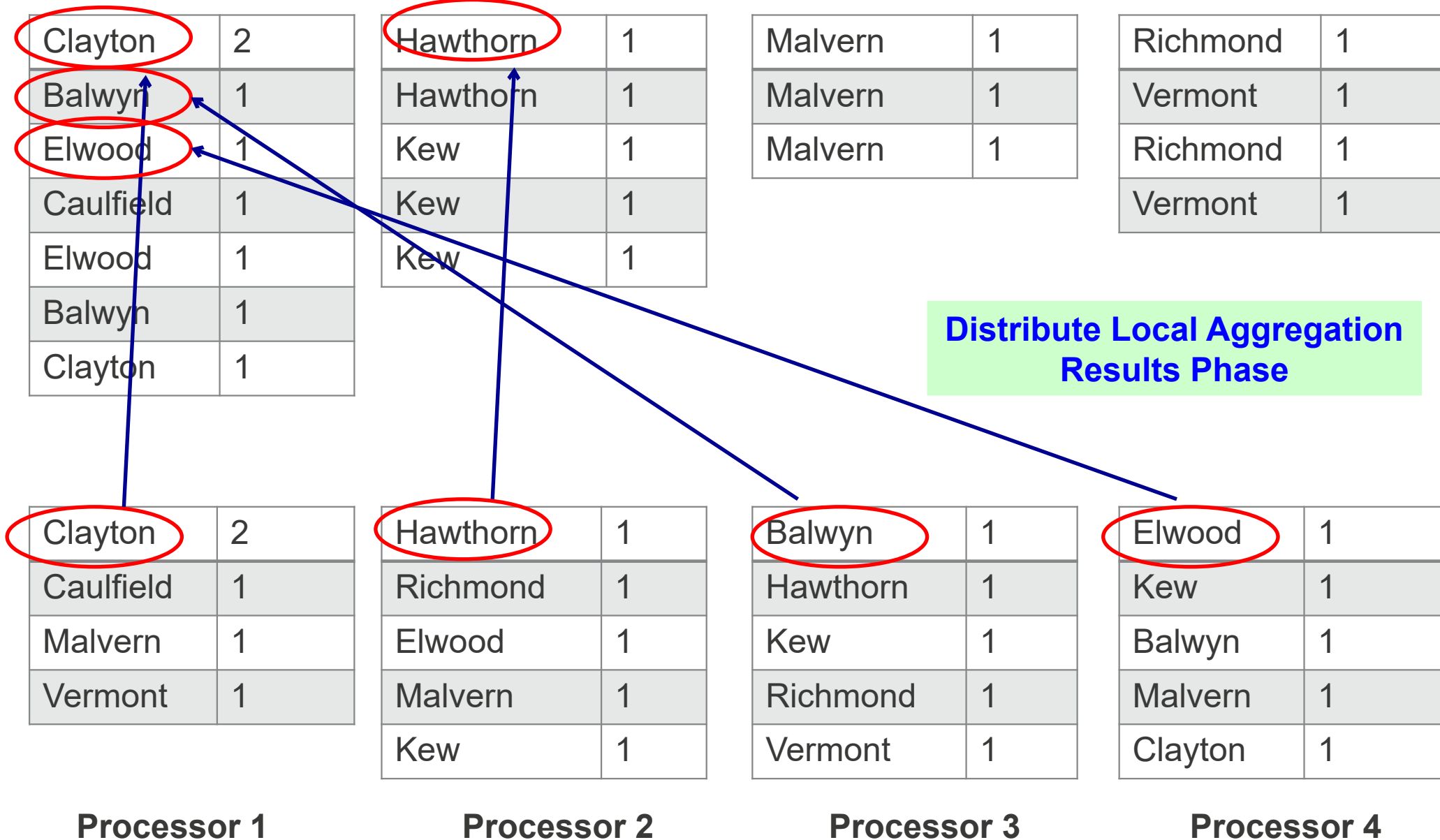
Fred	Hawthorn
George	Richmond
Harold	Elwood
Irene	Malvern
Jessica	Kew

Kelly	Balwyn
Lesley	Hawthorn
Megan	Kew
Naomi	Richmond
Oscar	Vermont

Peter	Elwood
Quin	Kew
Roger	Balwyn
Sarah	Malvern
Tracy	Clayton

4.4. Parallel GroupBy (cont'd)

Two-Phase Method

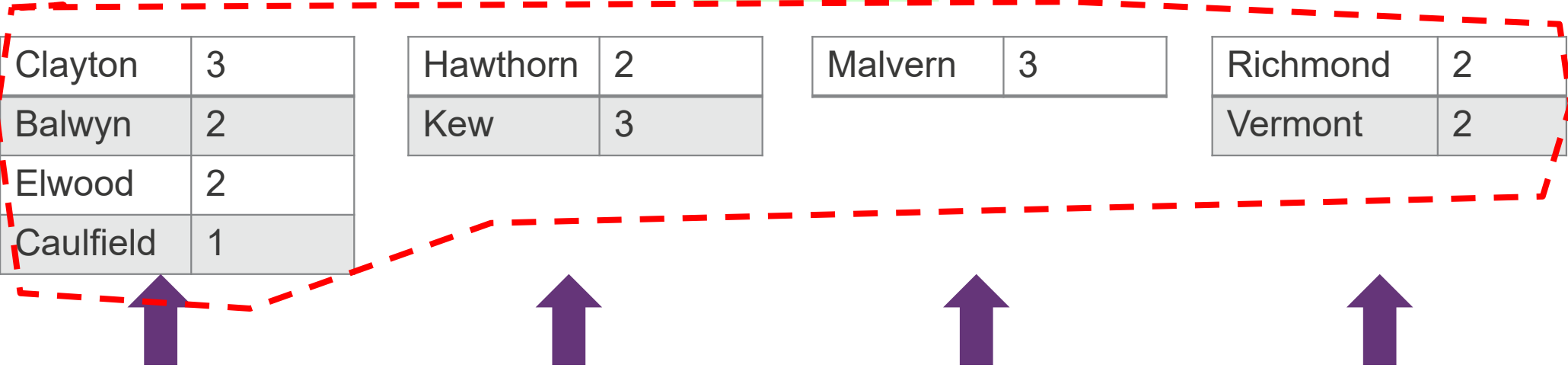


4.4. Parallel GroupBy (cont'd)

- Two-Phase Method



Final results



Global Aggregation Phase

Clayton	2
Balwyn	1
Elwood	1
Caulfield	1
Elwood	1
Balwyn	1
Clayton	1

Processor 1

Hawthorn	1
Hawthorn	1
Kew	1
Kew	1
Kew	1

Processor 2

Malvern	1
Malvern	1
Malvern	1

Processor 3

Richmond	1
Vermont	1
Richmond	1
Vermont	1

Processor 4

4.4. Parallel GroupBy (cont'd)

- **Redistribution Method**

- Step 1 (Partitioning phase): redistribute raw records to all processors
- Step 2 (Aggregation phase): each processor performs a local aggregation

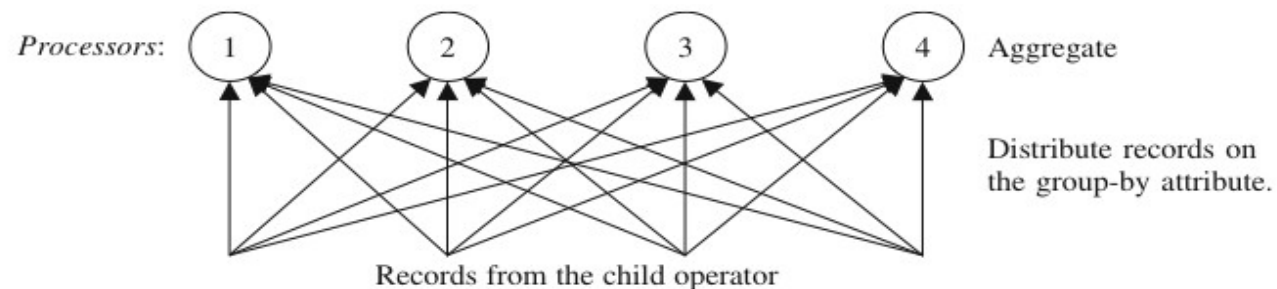
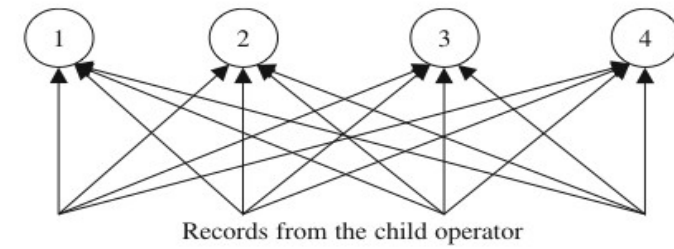


Figure 4.13 Redistribution method

4.4. Parallel GroupBy (cont'd)

- Redistribution Method



Initial Data Placement

Processor 1

A dam	Clayton
B en	Clayton
C hris	Caulfield
D ennis	Malvern
E ric	Vermont

Processor 2

F red	Hawthorn
G eorge	Richmond
H arold	Elwood
I rene	Malvern
J essica	Kew

Processor 3

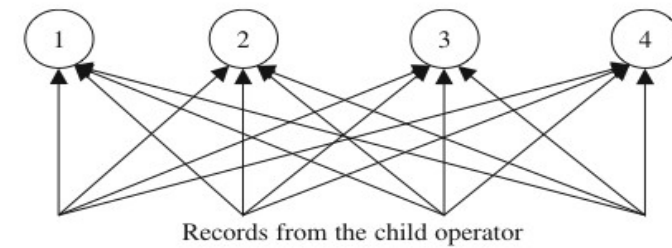
K elly	Balwyn
L esley	Hawthorn
M egan	Kew
N aomi	Richmond
O scar	Vermont

Processor 4

P eter	Elwood
Q uin	Kew
R oger	Balwyn
S arah	Malvern
T racy	Clayton

4.4. Parallel GroupBy (cont'd)

- Redistribution Method



Partitioning Phase

Adam	Clayton
Kelly	Balwyn
Peter	Elwood
Ben	Clayton
Chris	Caulfield
Harold	Elwood
Roger	Balwyn
Tracy	Clayton

Fred	Hawthorn
Lesley	Hawthorn
Quin	Kew
Megan	Kew
Jessica	Kew

Dennis	Malvern
Irene	Malvern
Sarah	Malvern

George	Richmond
Naomi	Richmond
Eric	Vermont
Oscar	Vermont

Processor 1

Processor 2

Processor 3

Processor 4

Adam	Clayton
Ben	Clayton
Chris	Caulfield
Dennis	Malvern
Eric	Vermont

Fred	Hawthorn
George	Richmond
Harold	Elwood
Irene	Malvern
Jessica	Kew

Kelly	Balwyn
Lesley	Hawthorn
Megan	Kew
Naomi	Richmond
Oscar	Vermont

Peter	Elwood
Quin	Kew
Roger	Balwyn
Sarah	Malvern
Tracy	Clayton

4.4. Parallel GroupBy (cont'd)

- Redistribution Method



Final results

Clayton	3
Balwyn	2
Elwood	2
Caulfield	1

Hawthorn	2
Kew	3

Malvern	3
---------	---

Richmond	2
Vermont	2

Adam	Clayton
Kelly	Balwyn
Peter	Elwood
Ben	Clayton
Chris	Caulfield
Harold	Elwood
Roger	Balwyn
Tracy	Clayton

Fred	Hawthorn
Lesley	Hawthorn
Quin	Kew
Megan	Kew
Jessica	Kew

Dennis	Malvern
Irene	Malvern
Sarah	Malvern

George	Richmond
Naomi	Richmond
Eric	Vermont
Oscar	Vermont

Aggregation Phase

Processor 1

Processor 2

Processor 3

Processor 4

4.4. Parallel GroupBy (cont'd)

- Redistribution Method

Clayton	3
Balwyn	2
Elwood	2
Caulfield	1

Hawthorn	2
Kew	3

Malvern	3
---------	---

Richmond	2
Vermont	2

Adam	Clayton
Kelly	Balwyn
Peter	Elwood
Ben	Clayton
Chris	Caulfield
Harold	Elwood
Roger	Balwyn
Tracy	Clayton

Fred	Hawthorn
Lesley	Hawthorn
Quin	Kew
Megan	Kew
Jessica	Kew

Dennis	Malvern
Irene	Malvern
Sarah	Malvern

George	Richmond
Naomi	Richmond
Eric	Vermont
Oscar	Vermont

What is the problem here?

Processor 1

Processor 2

Processor 3

Processor 4

4.4. Parallel GroupBy (cont'd)

- Redistribution Method (Task Stealing)

Create 5 buckets, instead of 4

Adam	Clayton
Kelly	Balwyn
Ben	Clayton
Chris	Caulfield
Roger	Balwyn
Tracy	Clayton

Peter	Elwood
Harold	Elwood

Processor 1

Fred	Hawthorn
Lesley	Hawthorn
Quin	Kew
Megan	Kew
Jessica	Kew

Processor 2

Dennis	Malvern
Irene	Malvern
Sarah	Malvern

Processor 3

George	Richmond
Naomi	Richmond
Eric	Vermont
Oscar	Vermont

Processor 4

Adam	Clayton
Ben	Clayton
Chris	Caulfield
Dennis	Malvern
Eric	Vermont

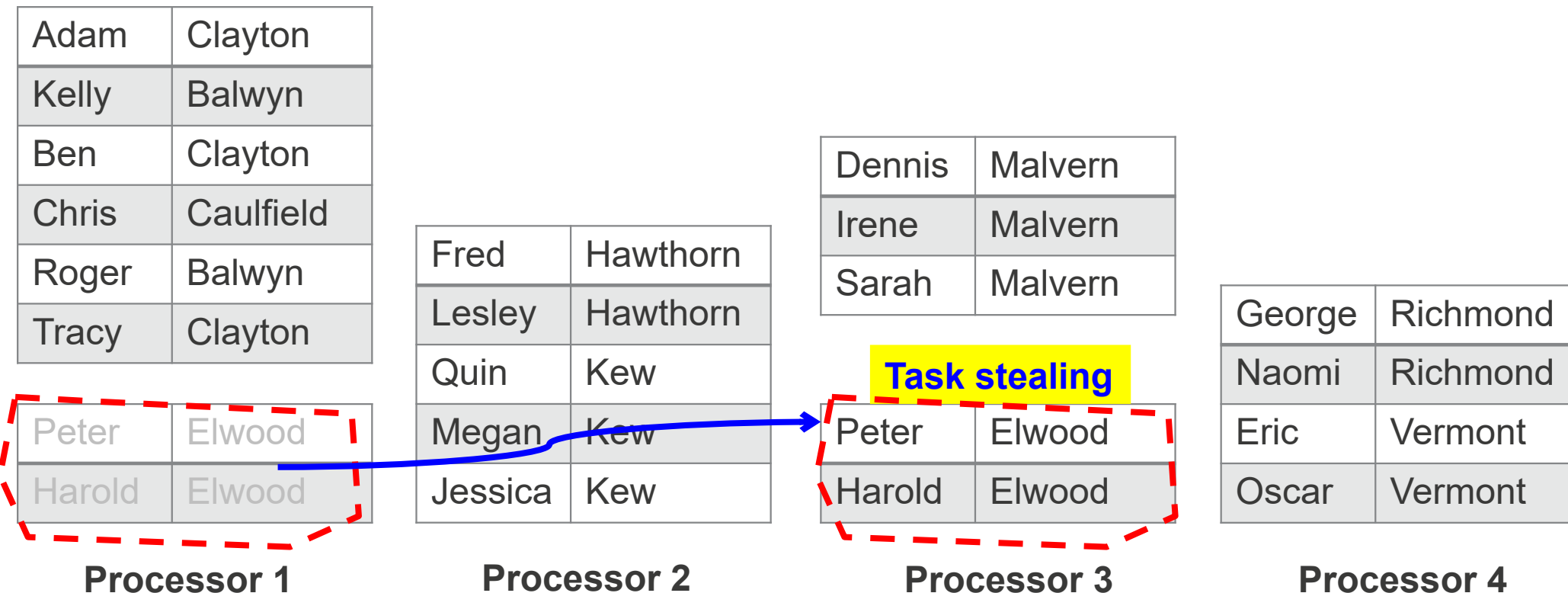
Fred	Hawthorn
George	Richmond
Harold	Elwood
Irene	Malvern
Jessica	Kew

Kelly	Balwyn
Lesley	Hawthorn
Megan	Kew
Naomi	Richmond
Oscar	Vermont

Peter	Elwood
Quin	Kew
Roger	Balwyn
Sarah	Malvern
Tracy	Clayton

4.4. Parallel GroupBy (cont'd)

- Redistribution Method (Task Stealing)



4.4. Parallel GroupBy (cont'd)

- Redistribution Method (Task Stealing)

Clayton	3
Balwyn	2
Caulfield	1

Hawthorn	2
Kew	3

Malvern	3
Elwood	2

Richmond	2
Vermont	2



Adam	Clayton
Kelly	Balwyn
Ben	Clayton
Chris	Caulfield
Roger	Balwyn
Tracy	Clayton

Fred	Hawthorn
Lesley	Hawthorn
Quin	Kew
Megan	Kew
Jessica	Kew

Dennis	Malvern
Irene	Malvern
Sarah	Malvern

Peter	Elwood
Harold	Elwood

George	Richmond
Naomi	Richmond
Eric	Vermont
Oscar	Vermont

Processor 1

Processor 2

Processor 3

Processor 4

• **Exercise 7 (FLUX Quiz)**

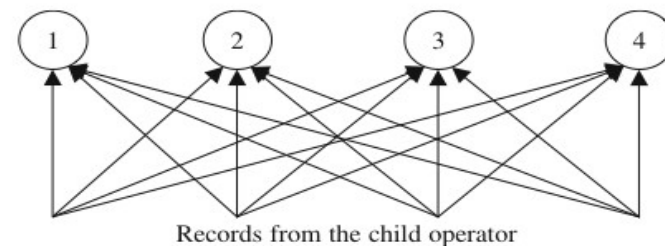
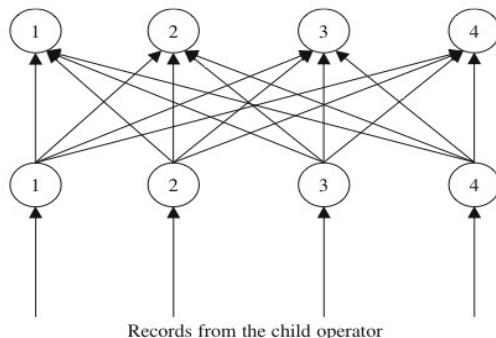
- The Redistribution Method has a load balancing option, through the Task Stealing method. The Two-Phase Method does not have a load balancing problem.
- A. TRUE
- B. FALSE

4.7. Summary

- Parallel groupby algorithms
 - Traditional methods (merge-all and hierarchical methods)
 - Two-phase method - – Local aggregation **before** data redistribution
 - Redistribution method - – Local aggregation **after** data redistribution
- Two-phase** and **Redistribution** methods perform better than the traditional and hierarchical merging methods
- Two-phase method** works well when the number of groups is small, whereas the **Redistribution method** works well when the number of groups is large

Ambuj and Naughton. "Adaptive parallel aggregation algorithms." (1995):

Why??



• Homework Exercises

- 1. Show how Load Balancing through Task Stealing be achieved in the Two Phase Method (using the same sample data as above) – EASY
- 2. Why is the Two-Phase Method good when the number of groups is small, whereas the Redistribution Method good when the number of groups is large? – MORE CHALLENGING
- 3. In what scenario may super linear speed up be achieved? – MORE CHALLENGING
- (Hints: See slides #8-#13 → the hash table cannot fit into main-memory)