

# Hash-based Grouping and Aggregation

R

co-grouping vs join processing

→ sort merge join

→ generalized co-grouped join

# Hash-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

```
SELECT R.x , aggregate ( R.y )  
FROM R  
GROUP BY R.x;
```

↑  
scalar

# Hash-based Grouping and Aggregation

R

grouping attribute:  $R.x$

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**HashBasedGrouping**(  $R$ ,  $\text{aggregate}()$  ):

HashMap hm = new HashMap();

List group = NULL;

ForEach  $r$  in  $R$ :

42

If NOT hm.contains(  $r.x$  ):

    group = new List();

Else:

//initialize hash map

//handle to group of  $R.x$

//for every tuple in  $R$

//check if already seen this key

//create new group

//i.e. key already in hash map (key seen before)

# Hash-based Grouping and Aggregation

R

grouping attribute:  $R.x$

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**HashBasedGrouping**(  $R$ ,  $\text{aggregate}()$  ):

```
HashMap hm = new HashMap();
```

```
List group = NULL;
```

```
ForEach  $r$  in  $R$ :
```

```
    If NOT hm.contains(  $r.x$  ):
```

```
        group = new List();
```

```
    Else:
```

```
        group = hm.get(  $r.x$  );
```

$\hookrightarrow \text{key} \mapsto \text{List} \langle \text{Values} \rangle$   
 $\uparrow$   
 $[R]$

```
//initialize hash map
```

```
//handle to group of  $R.x$ 
```

```
//for every tuple in  $R$ 
```

```
//check if already seen this key
```

```
//create new group
```

```
//i.e. key already in hash map (key seen before)
```

```
//get existing group from hash map
```

# Hash-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**HashBasedGrouping( **R**, aggregate() ):**

HashMap hm = new HashMap();

List group = NULL;

ForEach **r** in **R**:

    If NOT hm.contains( **r.x** ):

        group = new List();

    Else:

        group = hm.get( **r.x** );

    group.append( **r** );

    // hm.put( **r.x**, group );

//initialize hash map

//handle to group of R.x

//for every tuple in R

//check if already seen this key

//create new group

//i.e. key already in hash map (key seen before)

//get existing group from hash map

//append element to group

//insert/replace group in hash map



# Hash-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**HashBasedGrouping( **R**, aggregate() ):**

HashMap hm = new HashMap();

List group = NULL;

ForEach **r** in **R**:

If NOT hm.contains( **r.x** ):

    group = new List();

Else:

    group = hm.get( **r.x** );

group.append( **r** );

hm.put( **r.x**, group );

ForEach key in hm:

//initialize hash map

//handle to group of R.x

//for every tuple in R

//check if already seen this key

//create new group

//i.e. key already in hash map (key seen before)

//get existing group from hash map

//append element to group

//insert/replace group in hash map

//loop over all existing keys in hash map

# Hash-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**HashBasedGrouping( R, aggregate() ):**

HashMap hm = new HashMap();

List group = NULL;

ForEach **r** in R:

  If NOT hm.contains( **r.x** ):

    group = new List();

  Else:

    group = hm.get( **r.x** );

  group.append( **r** );

  hm.put( **r.x**, group );

ForEach key in hm:

  group = hm.get( key );

//initialize hash map

//handle to group of R.x

//for every tuple in R

//check if already seen this key

//create new group

//i.e. key already in hash map (key seen before)

//get existing group from hash map

//append element to group

//insert/replace group in hash map

//loop over all existing keys in hash map

//retrieve result group for that key

# Hash-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**HashBasedGrouping( R, aggregate() ):**

HashMap hm = new HashMap();

List group = NULL;

ForEach **r** in R:

  If NOT hm.contains( **r.x** ):

    group = new List();

  Else:

    group = hm.get( **r.x** );

  group.append( **r** );

  hm.put( **r.x**, group );

ForEach key in hm:

  group = hm.get( key );

  aggregationResult = aggregate( group );

  output( key, aggregationResult );

//initialize hash map

//handle to group of R.x

//for every tuple in R

//check if already seen this key

//create new group

//i.e. key already in hash map (key seen before)

//get existing group from hash map

//append element to group

//insert/replace group in hash map

//loop over all existing keys in hash map

//retrieve result group for that key

//call aggregation function on that group

//output result pairs



# Sort-based Grouping and Aggregation

R

grouping attribute:  $R.x$

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**SortBasedGrouping**(  $R$ ,  $\text{aggregate}()$  ):

sort(  $R$  on  $R.x$  );

Pointer  $PR = R[0]$ ;

Value  $\text{currentGroupValue} = R[0].x$ ;

List  $\text{group} = \text{new List}()$ ;

Do:

    If  $PR.x \neq \text{currentGroupValue}$ :

$R.x$   
2  
2  
2  
4  
4  
7  
7  
7  
5

//sort R on grouping attribute  $R.x$

//initialize pointer to first tuple of R

//value of  $R.x$  for **this** group

//create new group

//start loop

//if current tuple belongs to same group

# Sort-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**SortBasedGrouping( R, aggregate() ):**

sort( R on R.x );

Pointer PR = R[0];

Value currentGroupValue = R[0].x;

List group = new List();

Do:

If PR.x != currentGroupValue:

group closing | aggregationResult = aggregate( group );  
output( currentGroupValue, aggregationResult );  
group = new List();  
currentGroupValue = PR.x;

group.append( PR );

PR++;

While PR != R.end;

//sort R on grouping attribute R.x

//initialize pointer to first tuple of R

//value of R.x for **this** group

//create new group

//start loop

//if current tuple belongs to same group

//aggregate existing result group

//output result pairs

//create new group

//re-init value of R.x for this group

//append this tuple to result group anyway

//move pointer forward to next tuple

//end of table?

# Sort-based Grouping and Aggregation

R

grouping attribute: **R.x**

aggregation function:  $\text{aggregate}(\text{MultiSet}\langle R' \rangle): \text{MultiSet}\langle R' \rangle \mapsto \langle \text{ScalarType} \rangle, [R'] \subseteq [R]$

**SortBasedGrouping( R, aggregate() ):**

```
sort( R on R.x );
Pointer PR = R[0];
Value currentValue = R[0].x;
List group = new List();
Do:
    If PR.x != currentValue:
        aggregationResult = aggregate( group );
        output( currentValue, aggregationResult );
        group = new List();
        currentValue = PR.x;
    group.append( PR );
    PR++;
While PR != R.end;
aggregationResult = aggregate( group );
output( currentValue, aggregationResult );
```

```
//sort R on grouping attribute R.x
//initialize pointer to first tuple of R
//value of R.x for this group
//create new group
//start loop
//if current tuple belongs to same group
//aggregate existing result group
//output result pairs
//create new group
//re-init value of R.x for this group
//append this tuple to result group anyway
//move pointer forward to next tuple
//end of table?
//aggregate existing result group
//output result pairs
```