# Query Processing
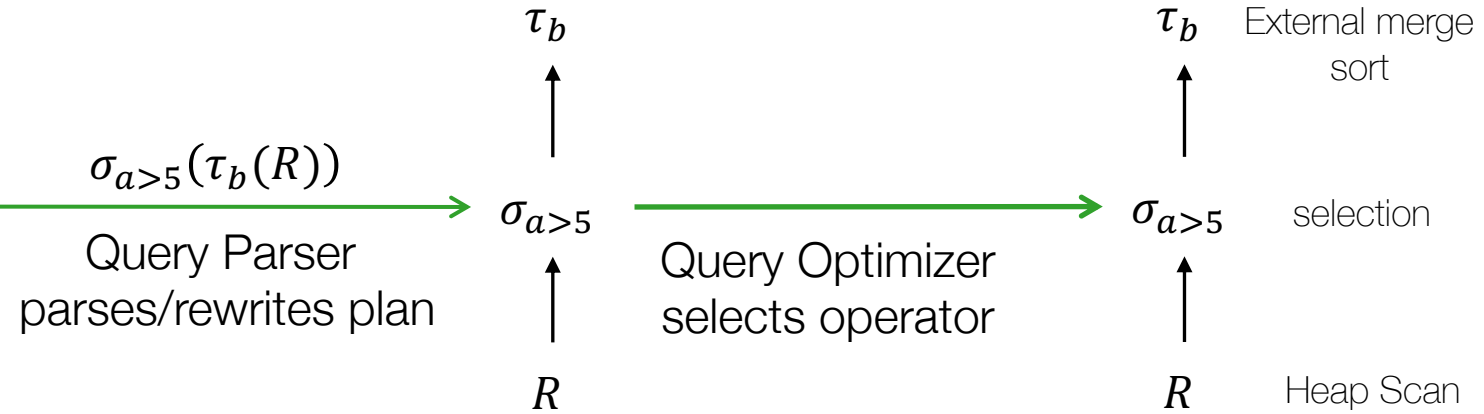
SQL Query

Logical Query Plan
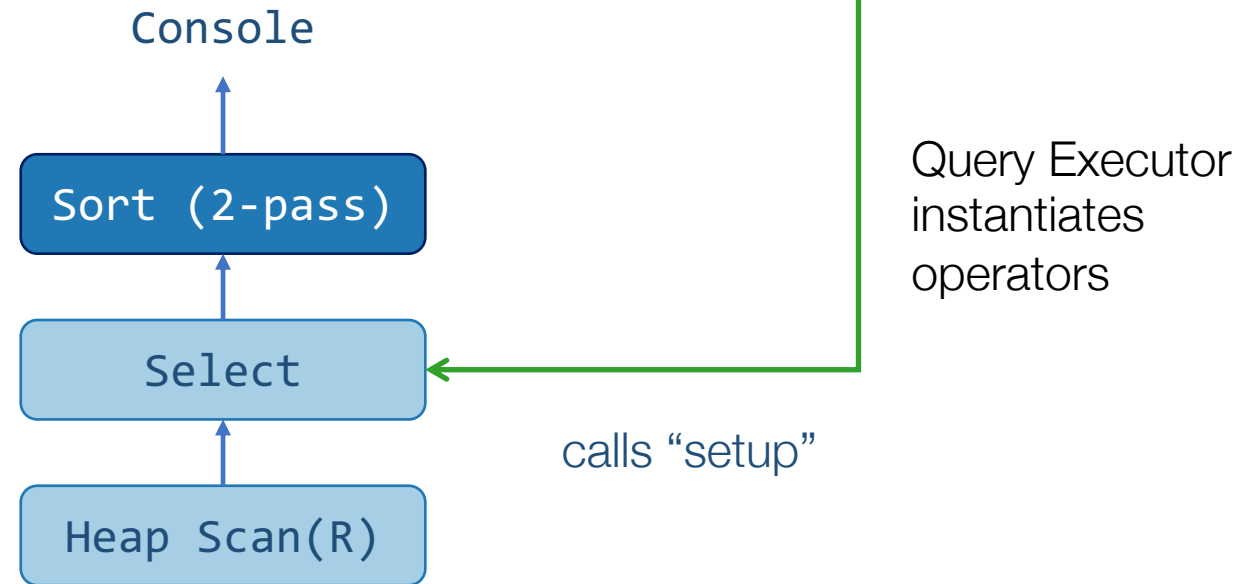
Optimized Physical Query Plan

```
select * from R
where a > 5
order by b;
```

$\sigma_{a>5}(\tau_b(R))$

Query Parser
parses/rewrites plan

$\tau_b$

$\sigma_{a>5}$

$R$

Query Optimizer
selects operator

$\tau_b$  External merge
sort

$\sigma_{a>5}$  selection

$R$  Heap Scan

Each operator is a
subclass of an *iterator*

```
abstract class iterator
  void setup(List<Iterator> children);
  void init(args);
  tuple next();
  void close();
}
```

Console

Sort (2-pass)

Select

Heap Scan(R)

Query Executor
instantiates
operators

calls "setup"

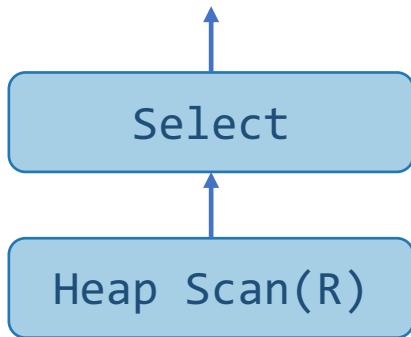# Query Processing

Each operator is a
subclass of an *iterator*

```
abstract class iterator
  void setup(List<Iterator> children);
  void init(args);
  tuple next();
  void close();
}
```

Console

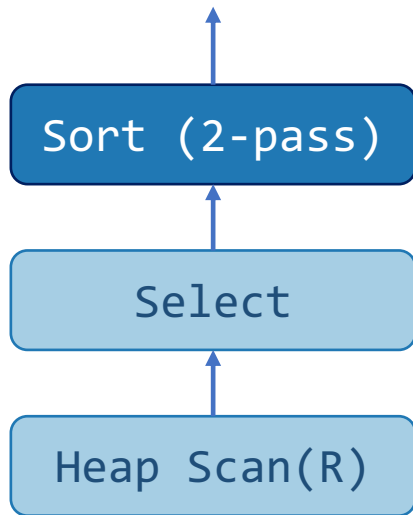Sort (2-pass)

Select

Heap Scan(R)

```
for each child in children:
  child.init()
```

child.init()

child.init()

child.next()          Blocking

child.next()          Streaming

child.next()          Streaming

# Query Processing – Iterator Model

```
Select

Heap Scan(R)
```

```
new(function predicate): //constructor given predicate
    p = predicate; //{return tuple.getValue("a")>5}
init():
    child.init(); //Initializes the heap scan op
    current = null;
next():
    while (current != EOF && !p(current)){
        current = child.next();
    }
    return current;
```

```
init(): //the op was setup to access the heap file
    current_page = file.getPage(0);
    current_slot = current_page.getSlot(0);
next():
    if(current_page == null) return EOF;
    current = current_slot.getTuple();
    current_slot.next();
    if(current_slot == null):
        current_page.next();
        if(current_page != null):
            current_slot = current_page.getSlot(0);
    return current;
```

```
Sort (2-pass)

Select

Heap Scan(R)
```

```
init():
    child.init();
    repeatedly call child.next() to generate sorted
            runs on disk until EOF

    open each sorted run / load into input buffer

next():
    output = min tuple across all input buffers (remove
        min tuple)

    if no tuples remain:
        return EOF;

    if min tuple was last one in its buffer:
        fetch next page from that run into buffer;

    return output;
```

# The Join Operator

# What is a Join?

$$R \bowtie_\sigma S = \sigma(R \times S)$$

# Cartesian Product ×

$R \times S$

Each row in R is paired with each row in S to produce $nm$ rows.

Loans ($R$ has $n$ rows)

| sid | eid | date | duration | ... |
|-----|-----|------|----------|-----|
| 72 | 981 | 3/8/20 | 2 weeks | |
| 76 | 786 | 3/18/21 | 2 days | |

Students ($S$ has $m$ rows)

| sid | name | major | ... |
|-----|------|-------|-----|
| 72 | Ibn Sina | Bio | |
| 73 | Plato | Phil | |
| 76 | Al Khawarizmi | CS | |

Loans × Students

| sid | eid | date | duration | sid | name | major | ... |
|-----|-----|------|----------|-----|------|-------|-----|
| 72 | 981 | 3/8/20 | 2 weeks | 72 | Ibn Sina | Bio | |
| 72 | 981 | 3/8/20 | 2 weeks | 73 | Plato | Phil | |
| 72 | 981 | 3/8/20 | 2 weeks | 76 | Al Khawarizmi | CS | |
| 76 | 786 | 3/18/21 | 2 days | 72 | Ibn Sina | Bio | |
| 76 | 786 | 3/18/21 | 2 days | 73 | Plato | Phil | |
| 76 | 786 | 3/18/21 | 2 days | 76 | Al Khawarizmi | CS | |

Loans ($R$ has $n$ rows)

| sid | eid | date | duration | ... |
|-----|-----|------|----------|-----|
| 72 | 981 | 3/8/20 | 2 weeks | |
| 76 | 786 | 3/18/21 | 2 days | |

Students ($S$ has $m$ rows)

| sid | name | major | ... |
|-----|------|-------|-----|
| 72 | Ibn Sina | Bio | |
| 73 | Plato | Phil | |
| 76 | Al Khawarizmi | CS | |

# Join ⋈

$$R \bowtie_{R.\text{sid}=S.\text{sid}} S$$

Each row in R is matched to a row in S that satisfies the join condition.

Loans ⋈ Students $\equiv \sigma_{\text{sid}=\text{sid}}(\text{Loans} \times \text{Students})$

| sid | eid | date | duration | sid | name | major | ... |
|-----|-----|------|----------|-----|------|-------|-----|
| 72 | 981 | 3/8/20 | 2 weeks | 72 | Ibn Sina | Bio | |
| 76 | 786 | 3/18/21 | 2 days | 76 | Al Khawarizmi | CS | |
| 72 | 981 | 3/8/20 | 2 weeks | 76 | Al Khawarizmi | CS | |
| 76 | 786 | 3/18/21 | 2 days | 72 | Ibn Sina | Bio | |
| 76 | 786 | 3/18/21 | 2 days | 73 | Plato | Phil | |
| 76 | 786 | 3/18/21 | 2 days | 76 | Al Khawarizmi | CS | |

# The Join Analysis Set Up

**Left / Outer Relation R**

a, …
a, …
b, …

2

3

4

N

\# of pages: N
\# of tuples: n

**Buffer Pool Size B**

1

…

…

B

**Right / Inner Relation S**

a, …
a, …
b, …

2

3

4

M

\# of pages: M
\# of tuples: m

Each table is broken down into pages.

Each table has a join key attribute and a value that could be a record id, or a tuple, etc. We only show the join keys.

When computing join costs, *we will ignore the output costs:*

(a) For now, we don't know how many tuples will join
(b) Across all implementations, the output cost is the same

| | R – Loans |
|---|---|
| | (*sid*, eid, date, duration, …) |
| | $N = 500$ |
| | $n = 40{,}000$ |
| | Tuples per page: $80$ |

| | S – Students |
|---|---|
| | (*sid*, name, major, …) |
| | $M = 1000$ |
| | $m = 100{,}000$ |
| | Tuples per page: $100$ |

| | B buffer size |
|---|---|
| | $B \leq 102$ |

| Algorithm | Naïve NLJ | Page NLJ | Block NLJ | Index NLJ | … |
|---|---|---|---|---|---|
| Cost (inner: S) | | | | | |
| Cost (inner: R) | | | | | |
| Effect of Buffer Size | | | | | |

# The Equipment Loans Application

# Nested Loops Join

# Naïve Nested Loops Join

**R**
# of pages: N
# of tuples: n

| a, … |
| b, … |
| p, … |

| w, … |
| b, … |
| c, … |

| a, … |
| p, … |
| x, … |

**S**
# of pages: M
# of tuples: m

| z, … |
| f, … |
| r, … |

| a, … |
| d, … |
| w, … |

| p, … |
| x, … |
| f, … |

| p, … |
| c, … |
| a, … |

### Buffer Pool Size B

| w, … |
| b, … |
| p, … |

| p, … |
| 2 |
| w, … |

| p,p,… |
| w,w… |
| p,p,… |

| a,a … |
| a,a … |
| p,p … |

| p,p … |
| w,w… |
| c,c … |

| a,a … |
| a,a… |
| p,p … |

| p,p … |
| x,x… |

```
for each tuple r in R:
  for each tuple s in S:
    if(key(r) == key(s)):
      emit(r, s)
```

R – Loans
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40,000$
Tuples per page: $80$

S – Students
(*sid*, name, major, …)
$M = 1000$
$m = 100,000$
Tuples per page: $100$

B buffer size
$B \leq 102$

| Algorithm | Naïve NLJ | Page NLJ | Block NLJ | Index NLJ |
|---|---|---|---|---|
| Cost (inner: S) | $N + nM$ <br><br> 500 + 40,000×1000 <br> = 40,000,500 | | | |
| Cost (inner: R) | $M + mN$ <br><br> 1000 + 100,000×500 <br> = 50,001,000 | | | |
| Effect of Buffer Size | $B > 500 \rightarrow N + M$ <br> = 1500 | | | |

*Does it matter what the inner relation is?*
Yes! We want the larger relation inside.
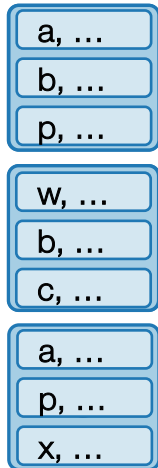
*Poor implementation*
We could match multiple tuples at a time for pages that are loaded
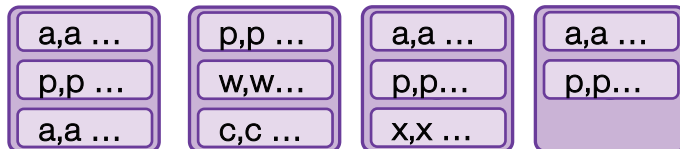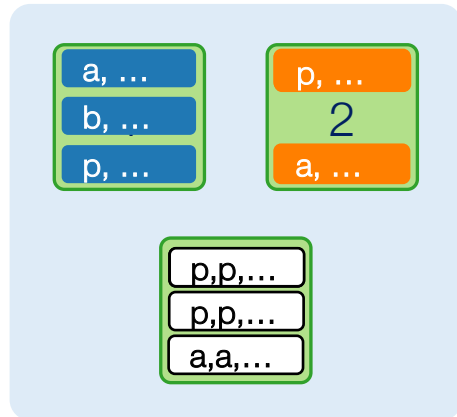
# Cost Analysis – Naïve NLJ

# Page Nested Loops Join
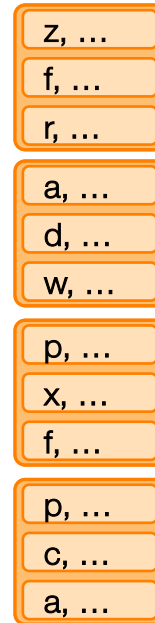
R
# of pages: N
# of tuples: n

| a, … |
| b, … |
| p, … |

| w, … |
| b, … |
| c, … |

| a, … |
| p, … |
| x, … |

Buffer Pool Size B

| a, … |
| b, … |
| p, … |

| p, … |
| 2 |
| a, … |

| p,p,… |
| p,p,… |
| a,a,… |

| a,a … |
| p,p … |
| a,a … |

| p,p … |
| w,w… |
| c,c … |

| a,a … |
| p,p… |
| x,x … |

| a,a … |
| p,p… |

S
# of pages: M
# of tuples: m

| z, … |
| f, … |
| r, … |

| a, … |
| d, … |
| w, … |

| p, … |
| x, … |
| f, … |

| p, … |
| c, … |
| a, … |

```
for each page p_r in R:
  for each page p_s in S:
    for each tuple r in p_r:
      for each tuple s in p_s:
        if(key(r) == key(s)):
          emit(r, s)
```

R − Loans
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40,000$
Tuples per page: 80

S − Students
(*sid*, name, major, …)
$M = 1000$
$m = 100,000$
Tuples per page: 100

B buffer size
$B \leq 102$

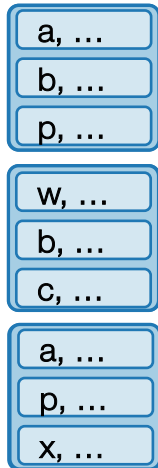| Algorithm | Naïve NLJ | Page NLJ | Block NLJ | Index NLJ |
|---|---|---|---|---|
| Cost (inner: S) | $N + nM$ | $N + NM$ | | |
| | 40,000,500 | 500,500 | | |
| Cost (inner: R) | $M + mN$ | $M + MN$ | | |
| | 50,001,000 | 501,000 | | |
| Effect of Buffer Size | $B > 500 \rightarrow N + M = 1500$ | | | |

*Better implementation but…*

We could match multiple tuples at a time for multiple pages of the outer relation that are loaded
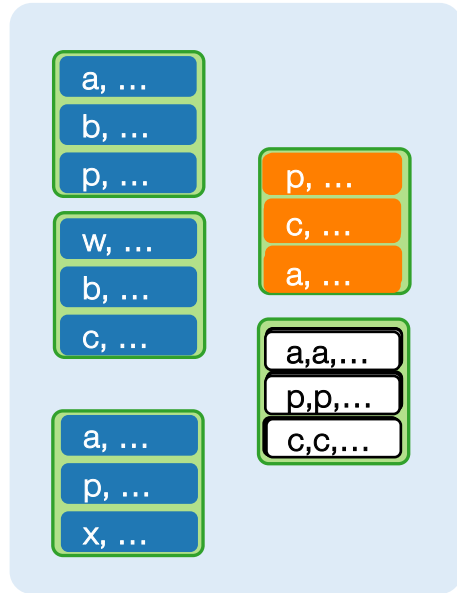
# Cost Analysis – Page NLJ
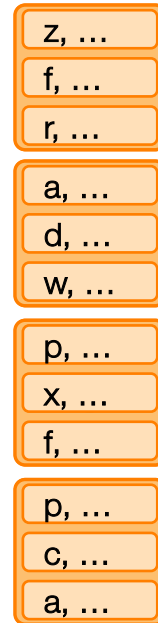
# Block Nested Loops Join

R
# of pages: N
# of tuples: n
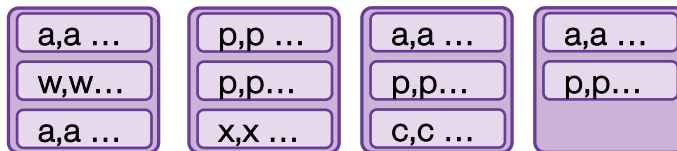
Buffer Pool Size B
B=5

S
# of pages: M
# of tuples: m

```
block := [B-2] pages
for each block B_r in R:
  for each page p_s in S:
    for each tuple r in B_r:
      for each tuple s in p_s:
        if(key(r) == key(s)):
          emit(r, s)
```

R pages:
a, …
b, …
p, …

w, …
b, …
c, …

a, …
p, …
x, …

Buffer pool:
a, …
b, …
p, …

w, …
b, …
c, …

a, …
p, …
x, …

p, …
c, …
a, …

a,a,…
p,p,…
c,c,…

S pages:
z, …
f, …
r, …

a, …
d, …
w, …

p, …
x, …
f, …

p, …
c, …
a, …

Output pages:
a,a …
w,w…
a,a …

p,p …
p,p…
x,x …

a,a …
p,p…
c,c …

a,a …
p,p…

**R – Loans**
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40{,}000$
Tuples per page: $80$

**S – Students**
(*sid*, name, major, …)
$M = 1000$
$m = 100{,}000$
Tuples per page: $100$

B buffer size
$B \leq 102$

| Algorithm | Naïve NLJ | Page NLJ | Block NLJ | Index NLJ |
|---|---|---|---|---|
| Cost (inner: S) | $N + nM$ | $N + NM$ | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | |
| (B = 102) | 40,000,500 | 500,500 | $500 + \left\lceil \dfrac{500}{100} \right\rceil 1000 = 5500$ | |
| Cost (inner: R) | $M + mN$ | $M + NM$ | $M + \left\lceil \dfrac{M}{B-2} \right\rceil N$ | |
| (B = 102) | 50,001,000 | 501,000 | $1000 + \left\lceil \dfrac{1000}{100} \right\rceil 500 = 6000$ | |
| Effect of Buffer Size | | $B > 500 \rightarrow N + M = 1500$ | | |

*Even better implementation but …*
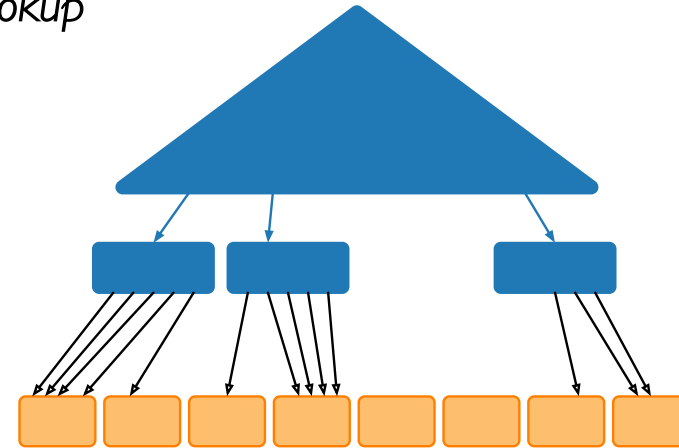What if we have an index on the inner relation?

# Cost Analysis – Block NLJ

# Index Nested Loops Join

```
for each tuple r in R:
  s = lookup(key(r), index(S))
  if(s):
    emit(r, s)
```
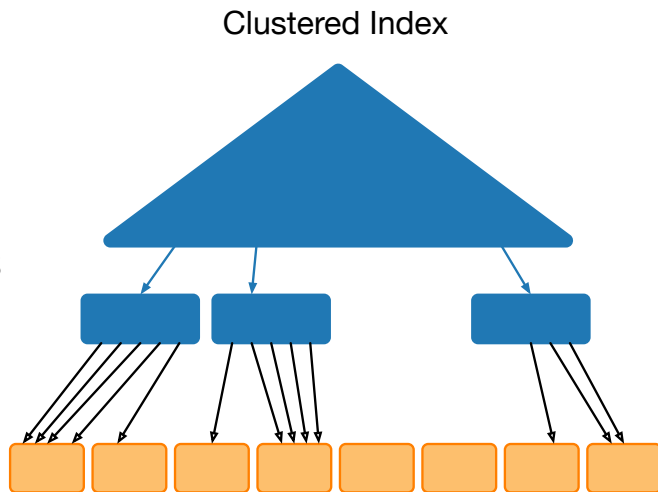
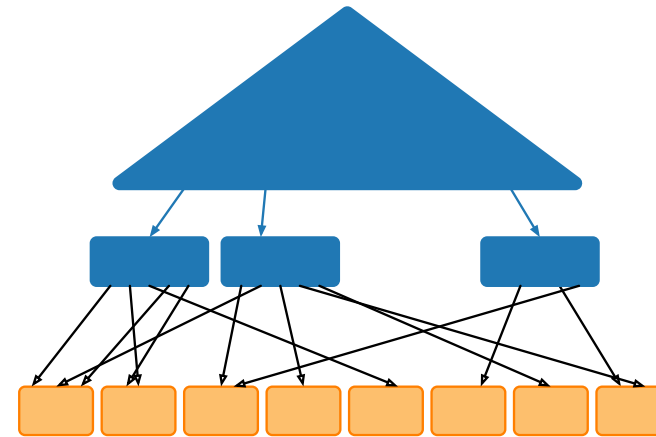*Cost of a Lookup*

Index
traversal

Height of a
tree (2-4)

+ Find page
in file by RID

+ 1

*Multiple matches*

Clustered Index

Unclustered Index

Index traversal +
# of matching pages

Index traversal +
# of matching tuples

R – Loans
(*sid*, eid, date, duration, ...)
$N = 500$
$n = 40{,}000$
Tuples per page: $80$

S – Students
(*sid*, name, major, ...)
$M = 1000$
$m = 100{,}000$
Tuples per page: $100$

B buffer size
$B \leq 102$

| Algorithm | Naïve NLJ | Page NLJ | Block NLJ | Index NLJ |
|---|---|---|---|---|
| Cost (inner: S) | $N + nM$ | $N + NM$ | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | $N + n(k)$ |
| (B = 102) | 40,000,500 | 500,500 | 5500 | $500 + 40{,}000(2+1)$ $= 120{,}500$ |
| Cost (inner: R) | $N + nM$ | $N + NM$ | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | No index on R by sid |
| (B = 102) | 50,001,000 | 501,000 | 6000 | |
| Effect of Buffer Size | | $B > 500 \rightarrow N + M$ | | |

Compute the full cross product: Quadratic!

Linear

*INLJ uses structure to overcome the need to check all other tuples*
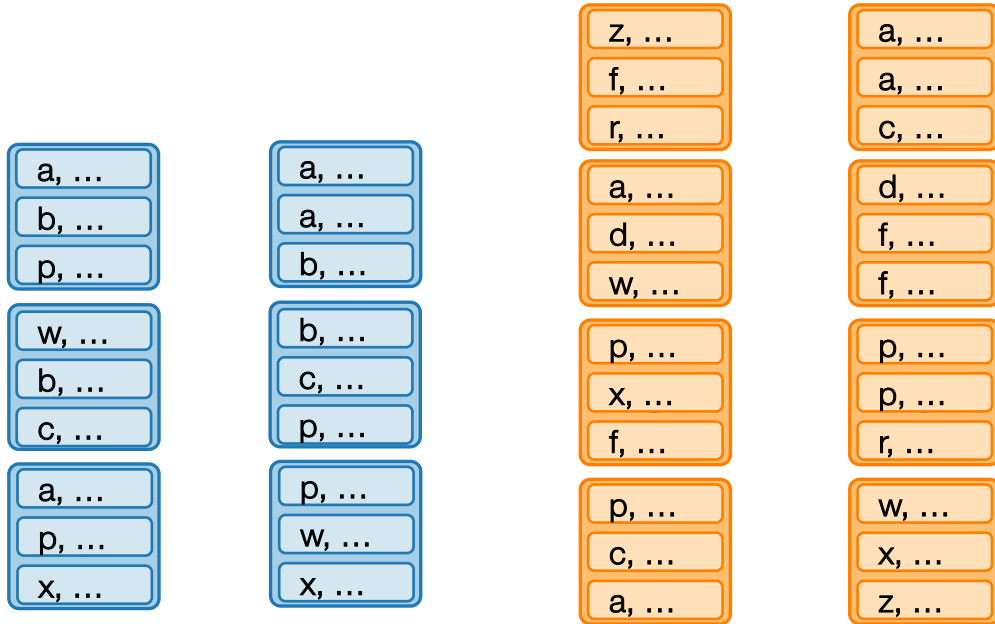*When does an index help?*
- $n \ll MN$
- *Small number of lookups*
- *Small buffer*

# Cost Analysis – Index NLJ

# Sort-Merge Join

# Sort-Merge Join

R

| a, … |
| b, … |
| p, … |

| w, … |
| b, … |
| c, … |

| a, … |
| p, … |
| x, … |

| a, … |
| a, … |
| b, … |

| b, … |
| c, … |
| p, … |

| p, … |
| w, … |
| x, … |

**R**
# of pages: N
# of tuples: n

S

| z, … |
| f, … |
| r, … |

| a, … |
| d, … |
| w, … |

| p, … |
| x, … |
| f, … |

| p, … |
| c, … |
| a, … |

| a, … |
| a, … |
| c, … |

| d, … |
| f, … |
| f, … |

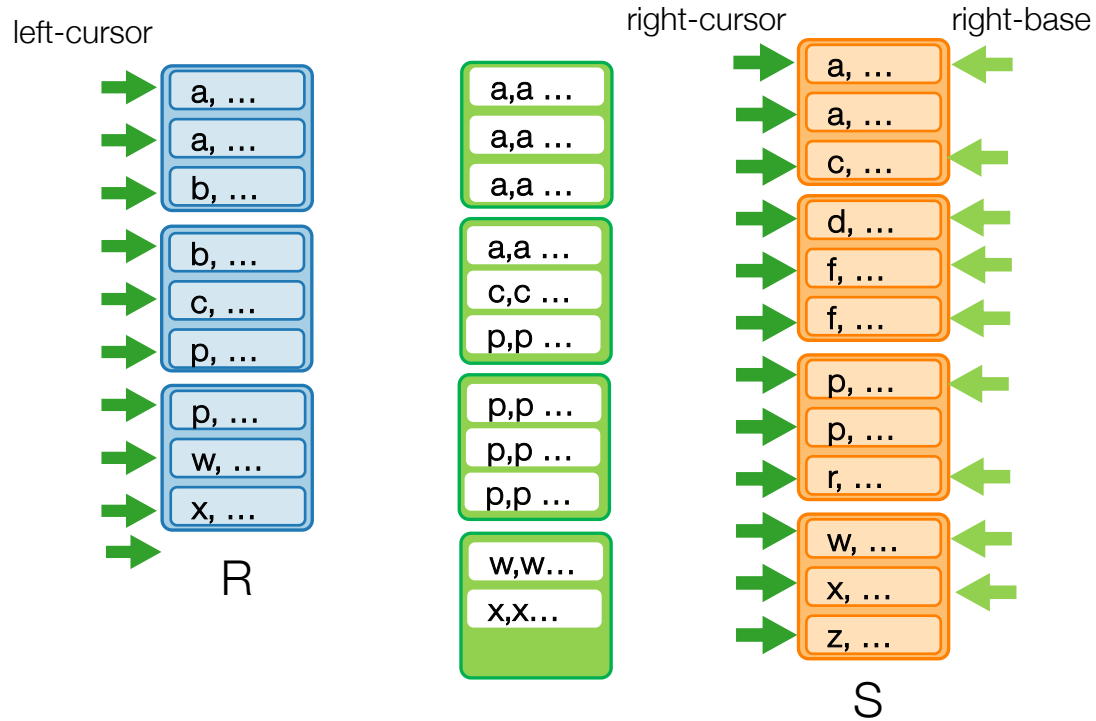| p, … |
| p, … |
| r, … |

| w, … |
| x, … |
| z, … |

**S**
# of pages: M
# of tuples: m

1. Sort R, S on join key using *external merge sort*

2. Scan sorted files and "merge"

# Sort-Merge Join – Merge Pass



```
do{
    key(left-cursor) == key(right-base):
        right-cursor = right-base
        do{
            emit(r, s)
            advance(right-cursor)
        } while(key(left-cursor) == key(right-cursor))
        advance left-cursor

    key(left-cursor) > key(right-base):
        if(right-cursor > right-base):
            right-base = right-cursor
        else
            advance(right-base, right-cursor)

    key(left-cursor) < key(right-base):
        advance(left-cursor)
}
while(left-cursor != EOF)
```

R – Loans
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40,000$
Tuples per page: 80

S – Students
(*sid*, name, major, …)
$M = 1000$
$m = 100,000$
Tuples per page: 100

B buffer size
$B \leq 102$

Assuming files are sorted

*Best Case Cost:*    $N + M$
- Single scan advancing cursors left and right!
- Equality Join with no duplicates!

Linear!

*Worst Case Cost:*    $nM$

- Effectively a nested loop join
- Single duplicate key on both sides === a cross product! (very unlikely)

Quadratic!

# Cost Analysis – Merge Join

R – Loans
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40{,}000$
Tuples per page: $80$

S – Students
(*sid*, name, major, …)
$M = 1000$
$m = 100{,}000$
Tuples per page: $100$
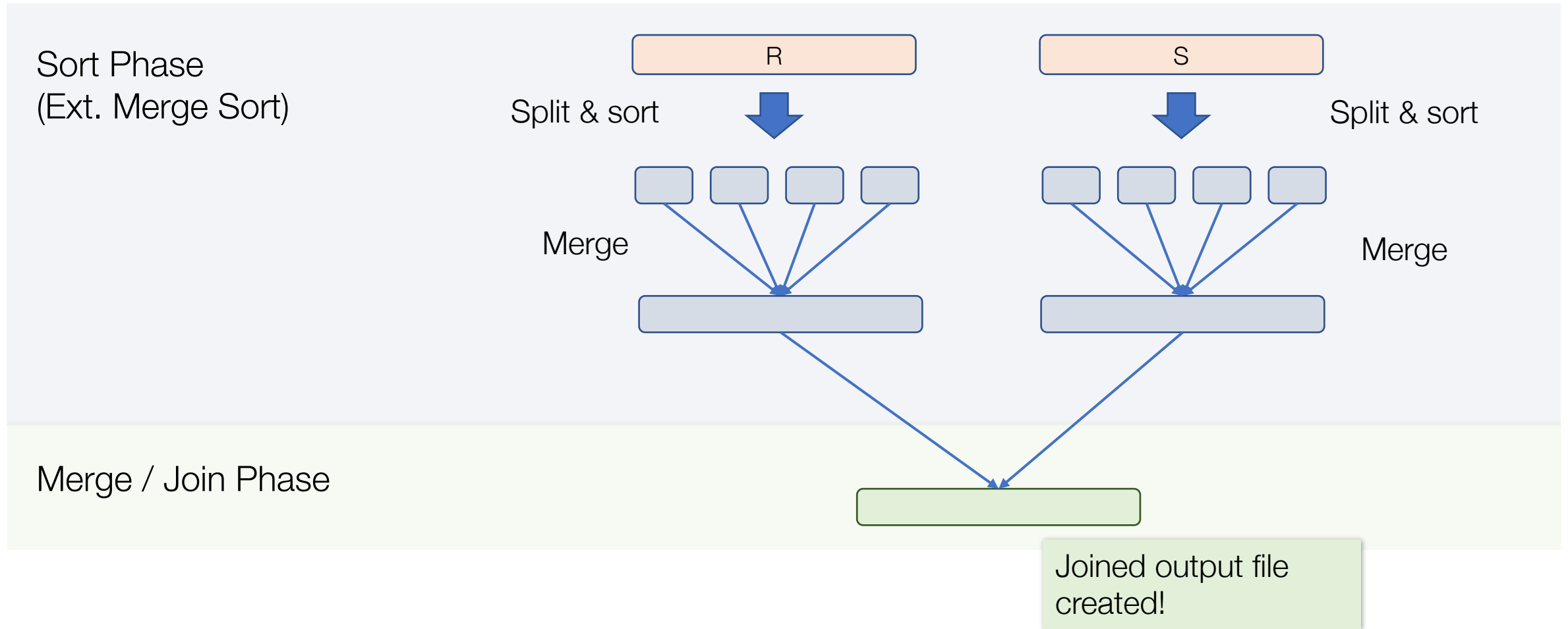
B buffer size
$B \leq 102$

| Algorithm | Block NLJ | Sort Merge Join |
|---|---|---|
| Cost | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | Sort + Merge $(4N + 4M) + (N + M)$ |
| $B = 102$ $\geq \max(\sqrt{N}, \sqrt{M})$ | 5500 | 7500 |
| $B = 35$ $\geq \max(\sqrt{N}, \sqrt{M})$ | 16500 | 7500 |

But we can do better if we can integrate merge join within the sort pass!
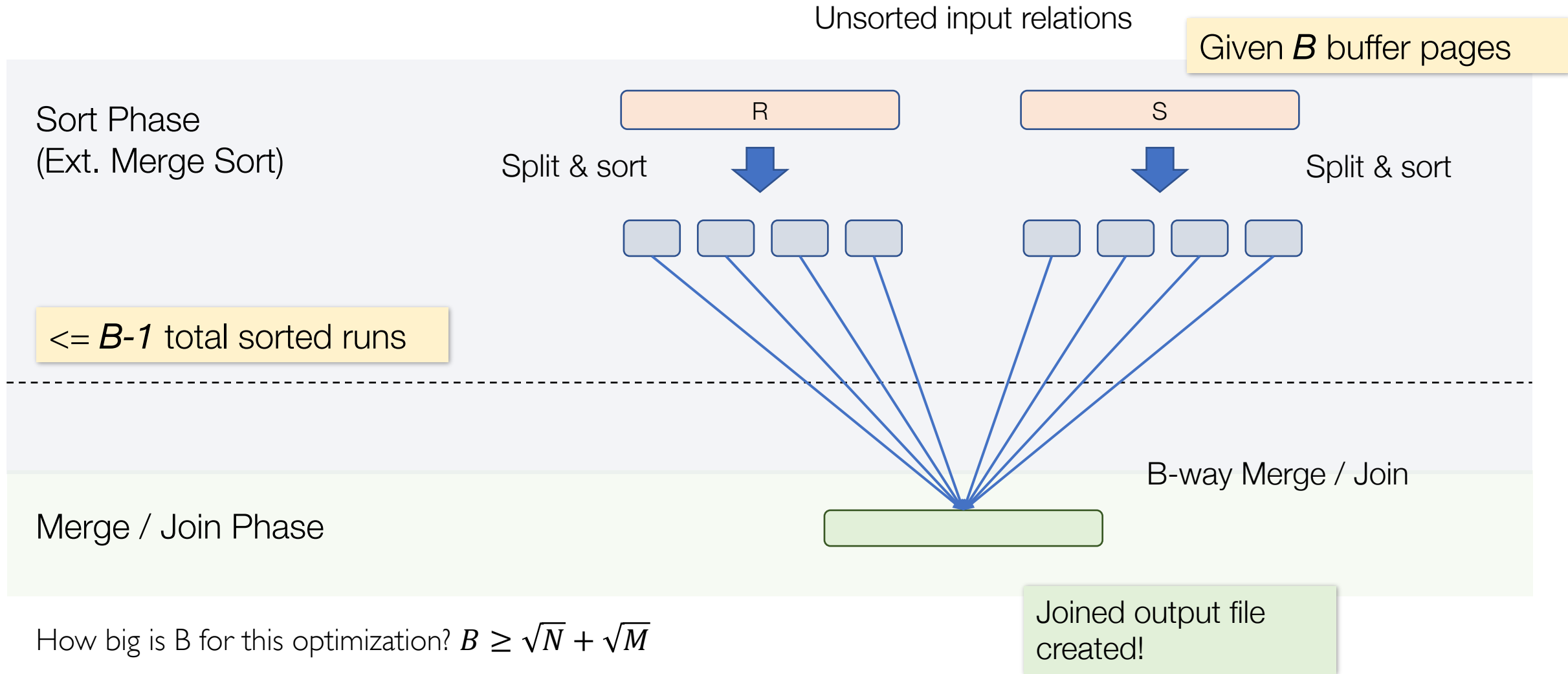
# Cost Analysis – SMJ

# Sort-Merge Join (No Refinement)

Unsorted input relations

## Sort Phase (Ext. Merge Sort)

R

S

Split & sort

Split & sort

Merge

Merge

## Merge / Join Phase

Joined output file created!

# Sort-Merge Join (With Refinement)

Unsorted input relations

Sort Phase
(Ext. Merge Sort)

| R | | S |

Split & sort     Split & sort

<= **B-1** total sorted runs

B-way Merge / Join

Merge / Join Phase

Joined output file created!

How big is B for this optimization? $B \geq \sqrt{N} + \sqrt{M}$

R – Loans
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40,000$
Tuples per page: 80

S – Students
(*sid*, name, major, …)
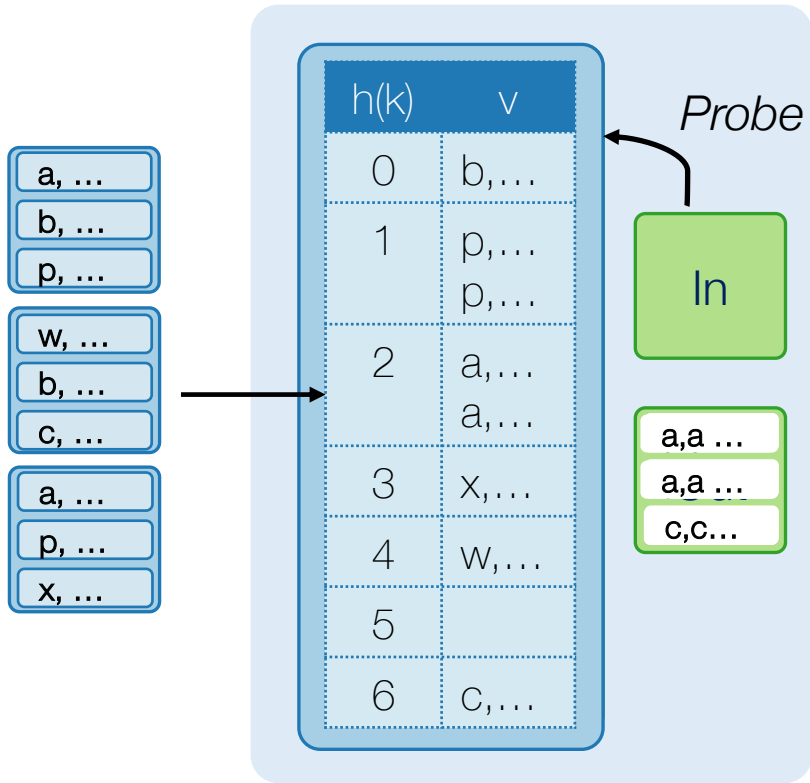$M = 1000$
$m = 100,000$
Tuples per page: 100

B buffer size
$B \leq 102$

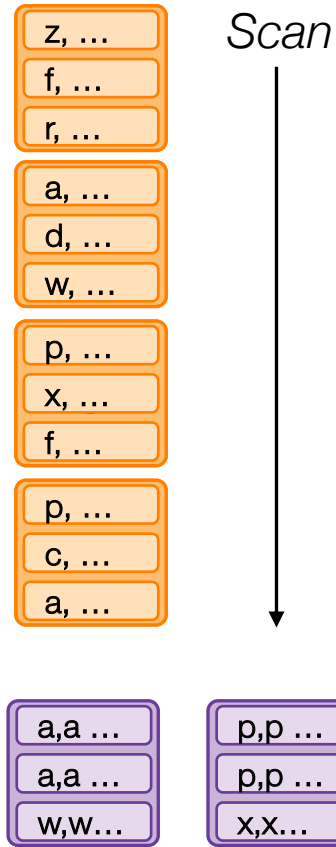| Algorithm | Block NLJ | Sort Merge Join | SMJ - Refined |
|---|---|---|---|
| Cost | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | Sort + Merge $(4N + 4M) + (N + M)$ | First Sort Pass + Merge Pass (exclude output) $3N + 3M$ |
| $B = 102$ $\geq \sqrt{N} + \sqrt{M}$ | 5500 | 7500 | 4500 |
| $B = 55$ $\geq \sqrt{N} + \sqrt{M}$ | 10500 | 7500 | 4500 |

# Cost Analysis – SMJ

# Grace Hash Join

# Naïve Hash Join



| h(k) | v |
|------|---|
| 0 | b, ... |
| 1 | p, ...<br>p, ... |
| 2 | a, ...<br>a, ... |
| 3 | x, ... |
| 4 | w, ... |
| 5 | |
| 6 | c, ... |

*Probe*

In

a,a ...
a,a ...
c,c...

*Scan*

z, ...
f, ...
r, ...
a, ...
d, ...
w, ...
p, ...
x, ...
f, ...
p, ...
c, ...
a, ...

Build in-memory hash table of size (B-2) using hashing function *h(k)*

a,a ...
a,a ...
w,w...

p,p ...
p,p ...
x,x...

p,p ...
p,p ...
c,c...

a,a ...
a,a ...

Simple Algorithm

Cost: $N + M$

Memory requirement
- $\min(N, M) < B - 2$

*What if the hash table of the smaller relation doesn't fit?*

Divide / Partition Phase

- Use a partitioning hash function to divide each table into $B - 1$ uniform partitions

Conquer Phase

For each partition $P_i$

- Build an in-memory hash table using the smaller partition(R): $P_i(R)$
- Hash table has to fit in $B - 2$ buffers
- Scan each page from S's partition $P_i(S)$ and probe the in-memory hash-table

# Grace Hash Join

Cost

Partitioning Phase + Matching Phase

$$(2N + 2M) + (N + M) = 3(N + M)$$

Memory Requirement

R is the smaller relation

Partitioning Phase divides R into (B-1) partitions of size $\frac{N}{B-1}$

Matching Phase requires each partition to be:

- $\frac{N}{B-1} < B - 2$
- $N < (B - 2)(B - 1)$
- $B \geq \sqrt{N}$

The probing relation S can be quite big, there are no restrictions on the size of its partitions!

# Cost and Memory Analysis of Grace Hash Join

## Cost Analysis – Grace Hash Join

**R – Loans**
(*sid*, eid, date, duration, …)
$N = 500$
$n = 40{,}000$
Tuples per page: $80$

**S – Students**
(*sid*, name, major, …)
$M = 1000$
$m = 100{,}000$
Tuples per page: $100$

**B buffer size**
$B \leq 102$

| Algorithm | Block NLJ | SMJ - Refined | Grace Hash Join |
|---|---|---|---|
| Cost | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | $3N + 3M$ | $3N + 3M$ |
| $B = 102$ $\geq \sqrt{N} + \sqrt{M}$ | 5500 | 4500 | 4500 |
| $B = 55$ $\geq \sqrt{N} + \sqrt{M}$ | 10500 | 4500 | 4500 |

## R – Loans

(*sid*, eid, date, duration, …)

$N = 500$

$n = 40,000$

Tuples per page: 80

## S – Students

(*sid*, name, major, …)
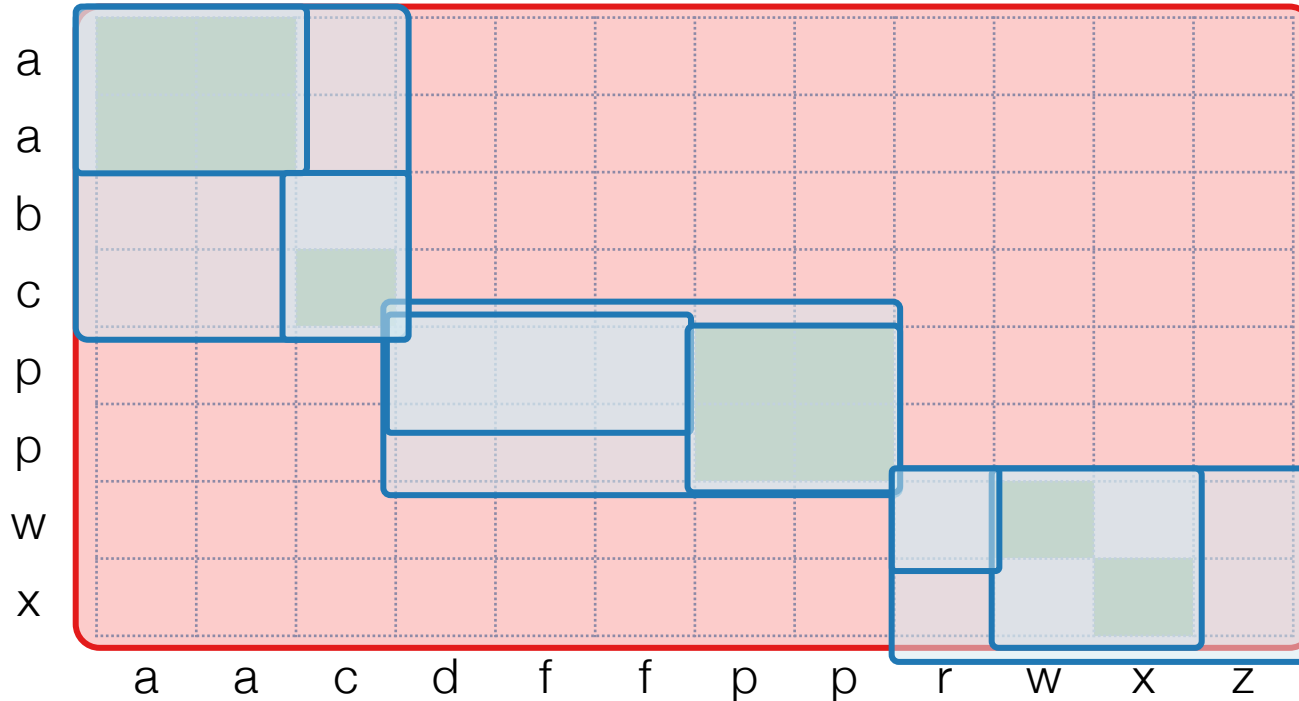
$M = 1000$

$m = 100,000$

Tuples per page: 100

## B buffer size

$B \leq 102$

| Algorithm | Block NLJ | SMJ - Refined | Grace Hash Join |
|---|---|---|---|
| Cost | $N + \left\lceil \dfrac{N}{B-2} \right\rceil M$ | $3N + 3M$ | $3N + 3M$ |
| $B = 102$ $\geq \sqrt{N} + \sqrt{M}$ | 5500 | 4500 | 4500 |
| $B = 55$ $\geq \sqrt{N} + \sqrt{M}$ | 10500 | 4500 | 4500 |
| $B = 25 \geq \sqrt{N}$ | 22500 | Needs more passes | 4500 |

# Cost Analysis – Grace Hash Join

# Summary

The Grace Hash Join partitioning breaks down the grid into smaller grids for further matching

The SMJ uses order to avoid searching the whole grid and establishing search boundaries

BNLJ doesn't take advantage of structure – we explore the whole grid for matches!



# Visual Comparison – BNLJ, SMJ vs. GHJ

*Nested Loops Join*    Works for arbitrary join conditions

*Index Nested Loops Join*    If you have an index, equi-join and a small number of lookups! $n < NM$

*Sort-Merge/Hash Join*
- Linear IO complexity
- No index required
- Hash is better if one of the relations is much smaller
- Sort-Merge is better if order is required or if the relations are already sorted (perhaps from a previous join).

A typical DBMS implements all of these and uses a query optimizer to select the best join for a given query plan!

# Key Takeaways