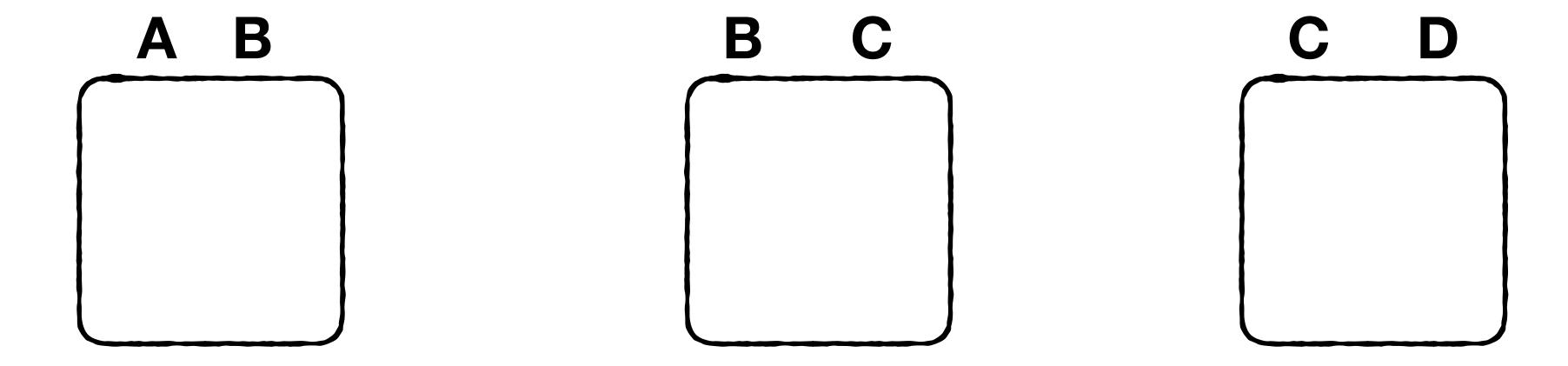
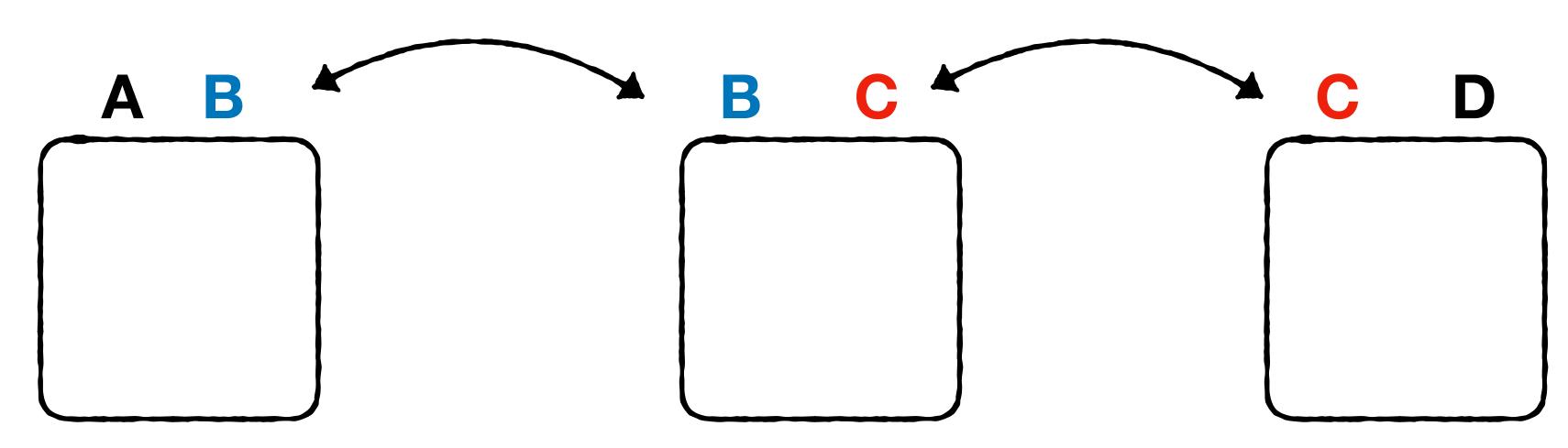
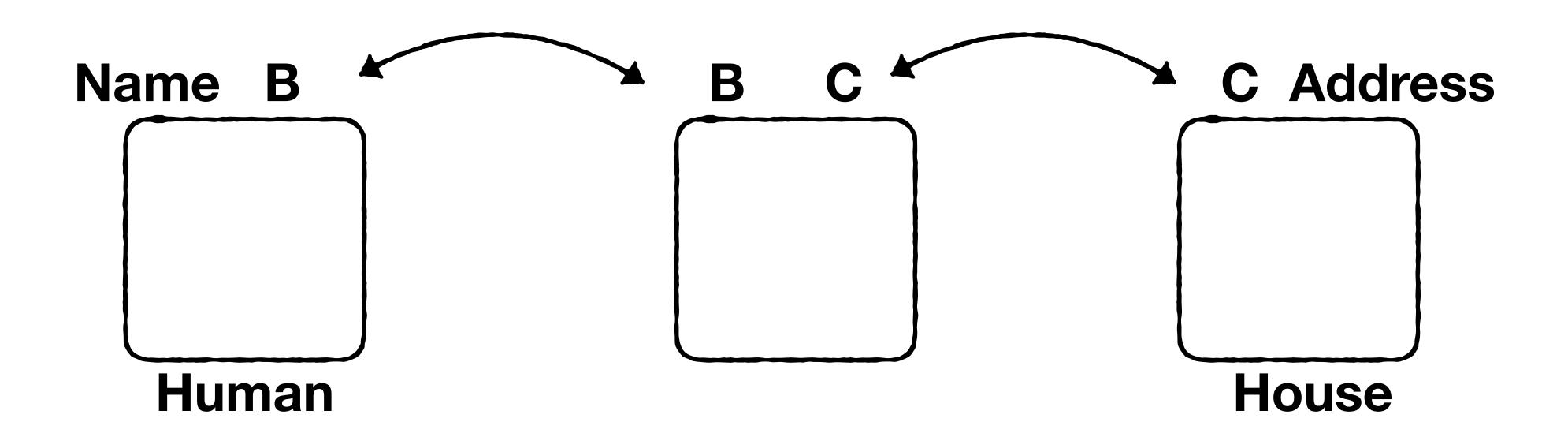
Query Evalaution

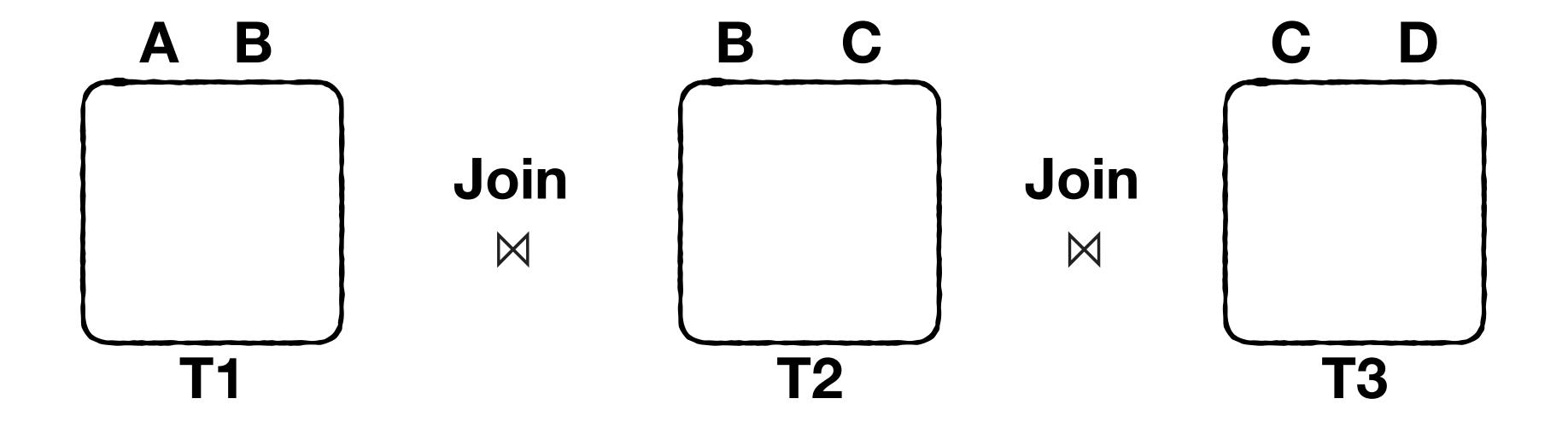
CSC443H1 Database System Technology



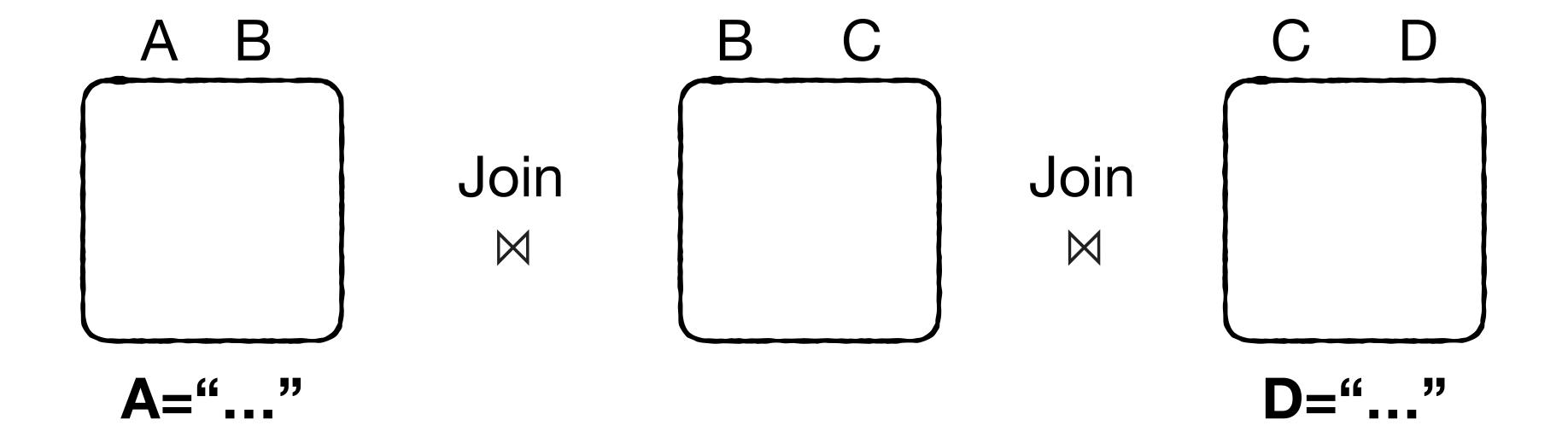




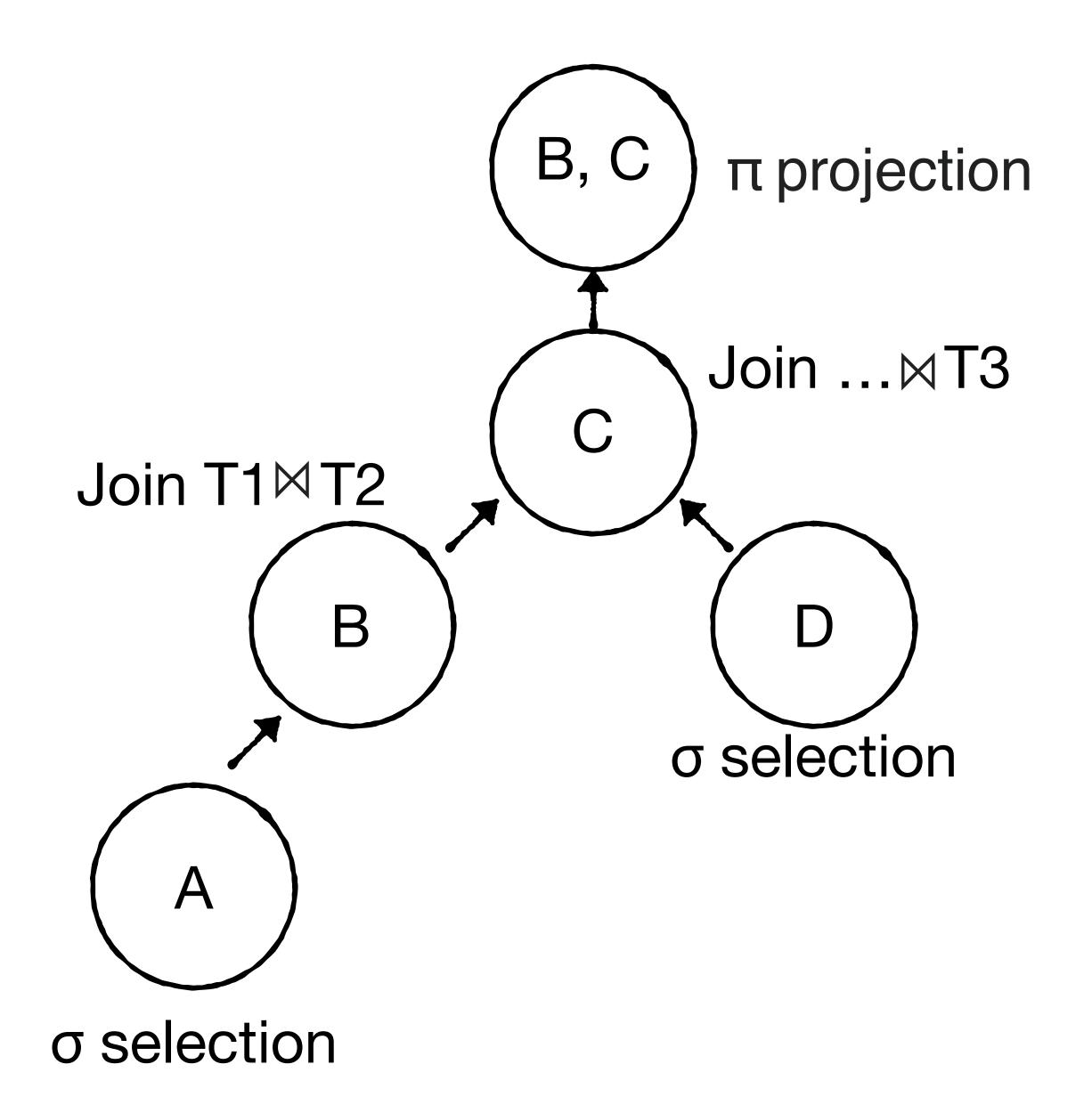


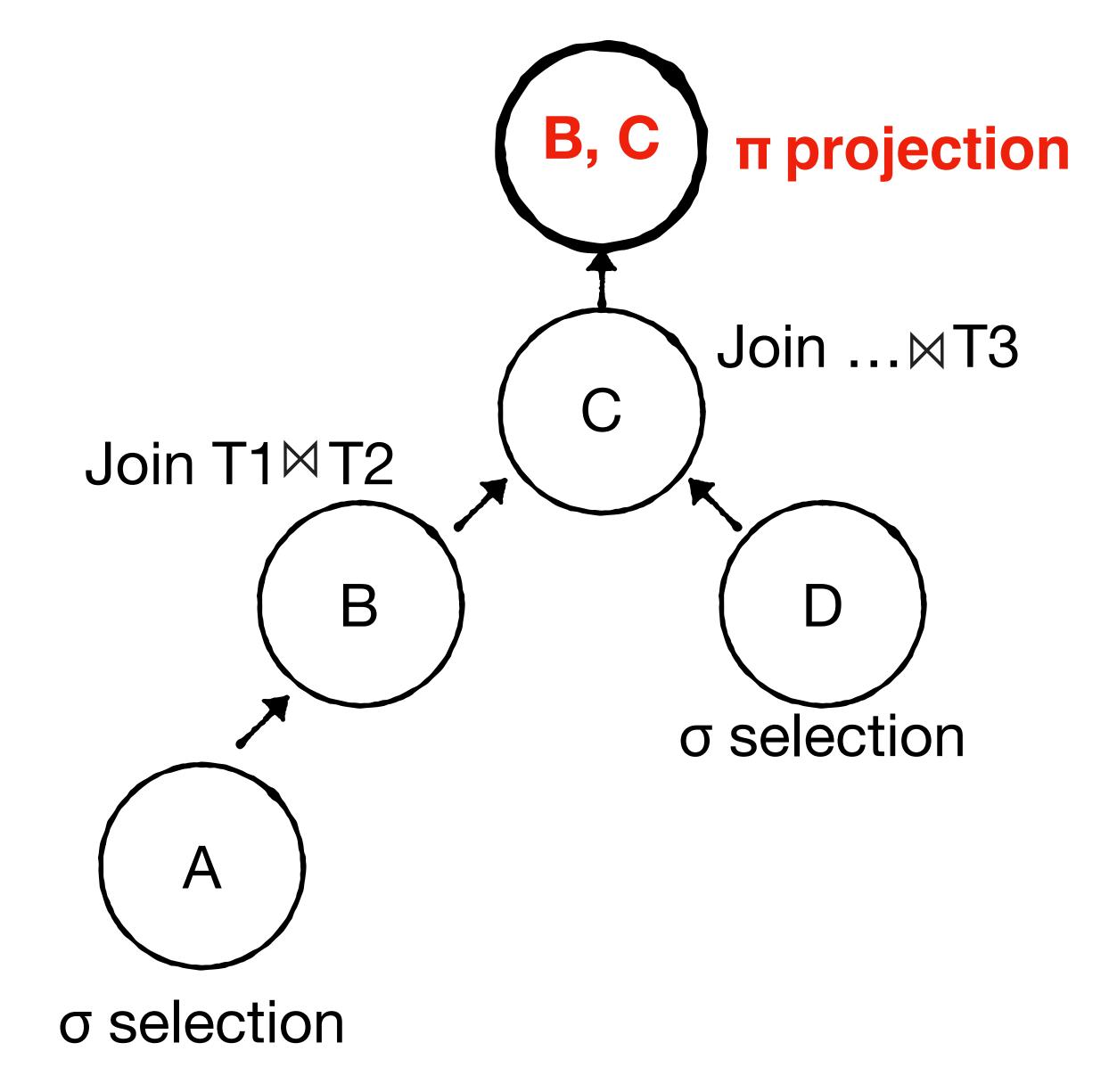


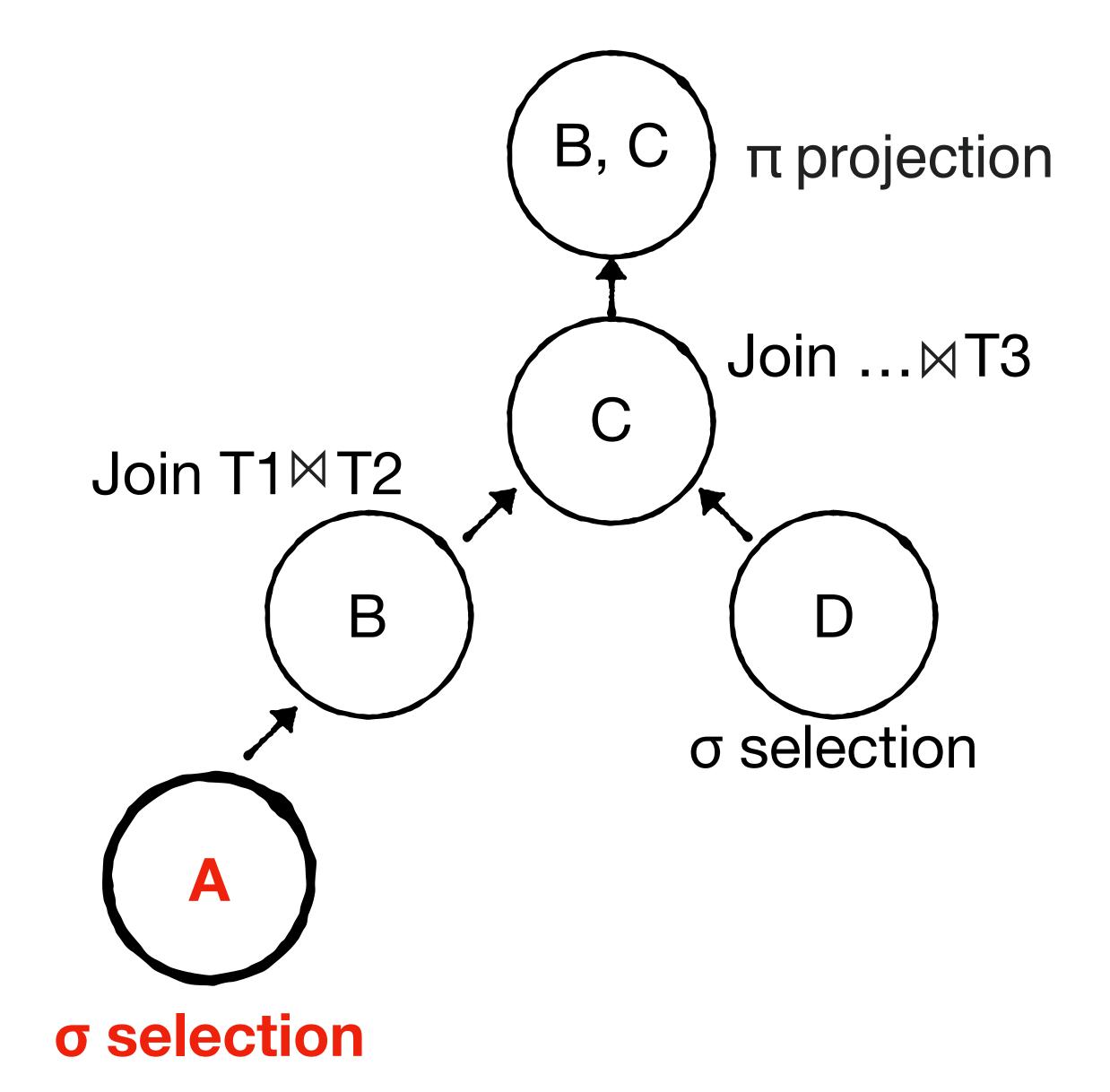
Select A, D from T1, T2, T3 where T1.B = T2.B and T2.C = T3.C

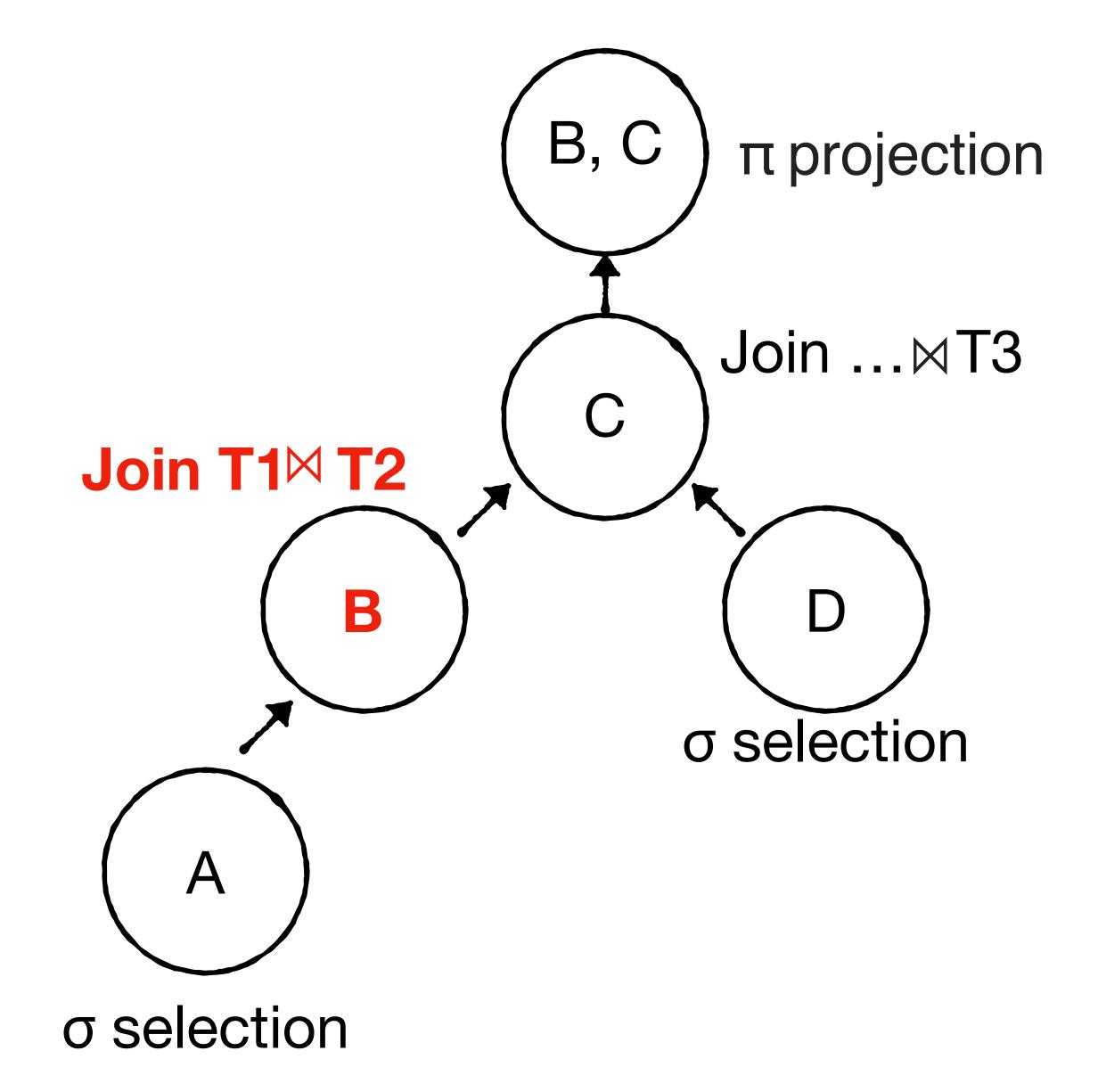


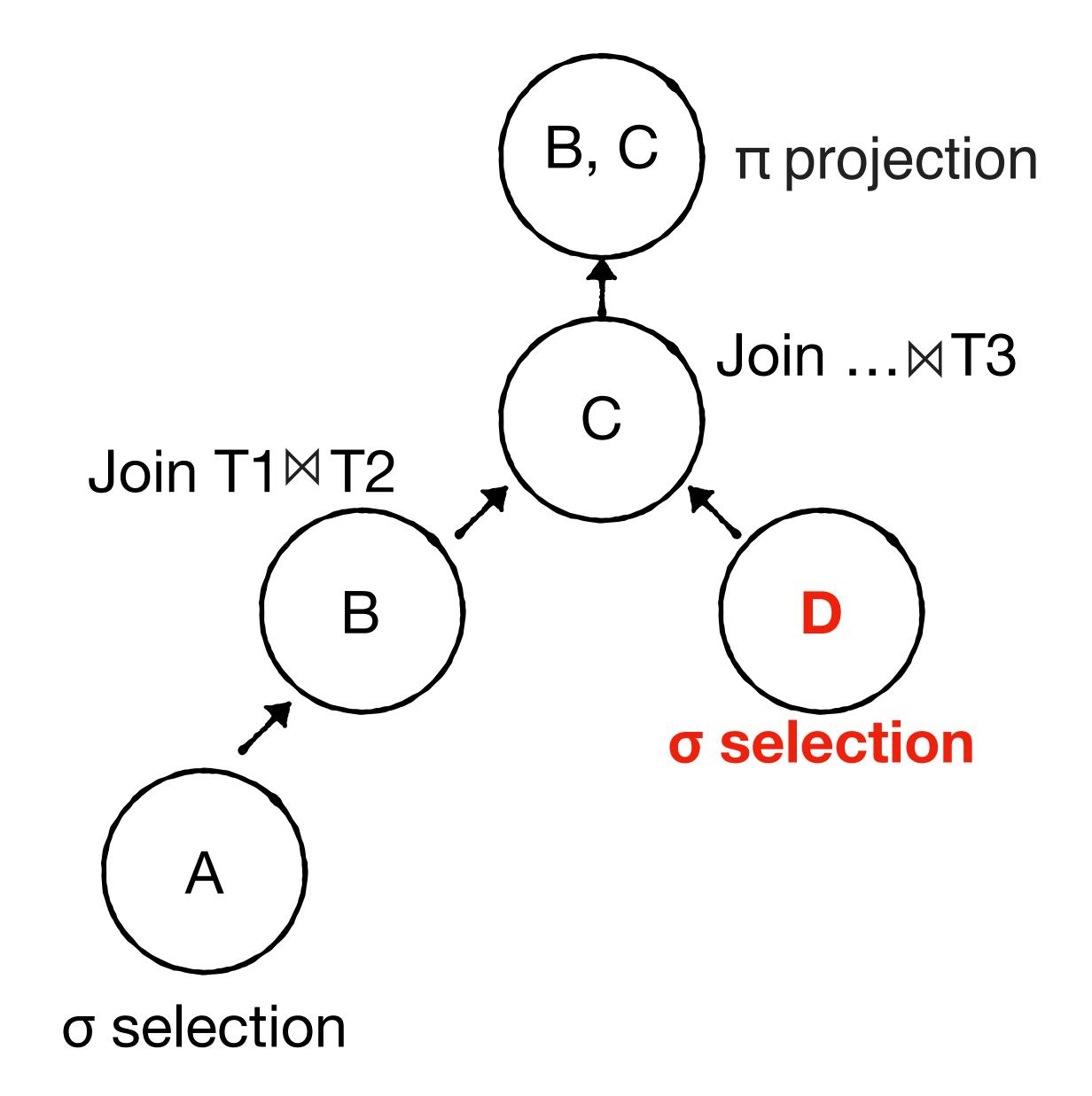
Converted into a query plan of logical operators

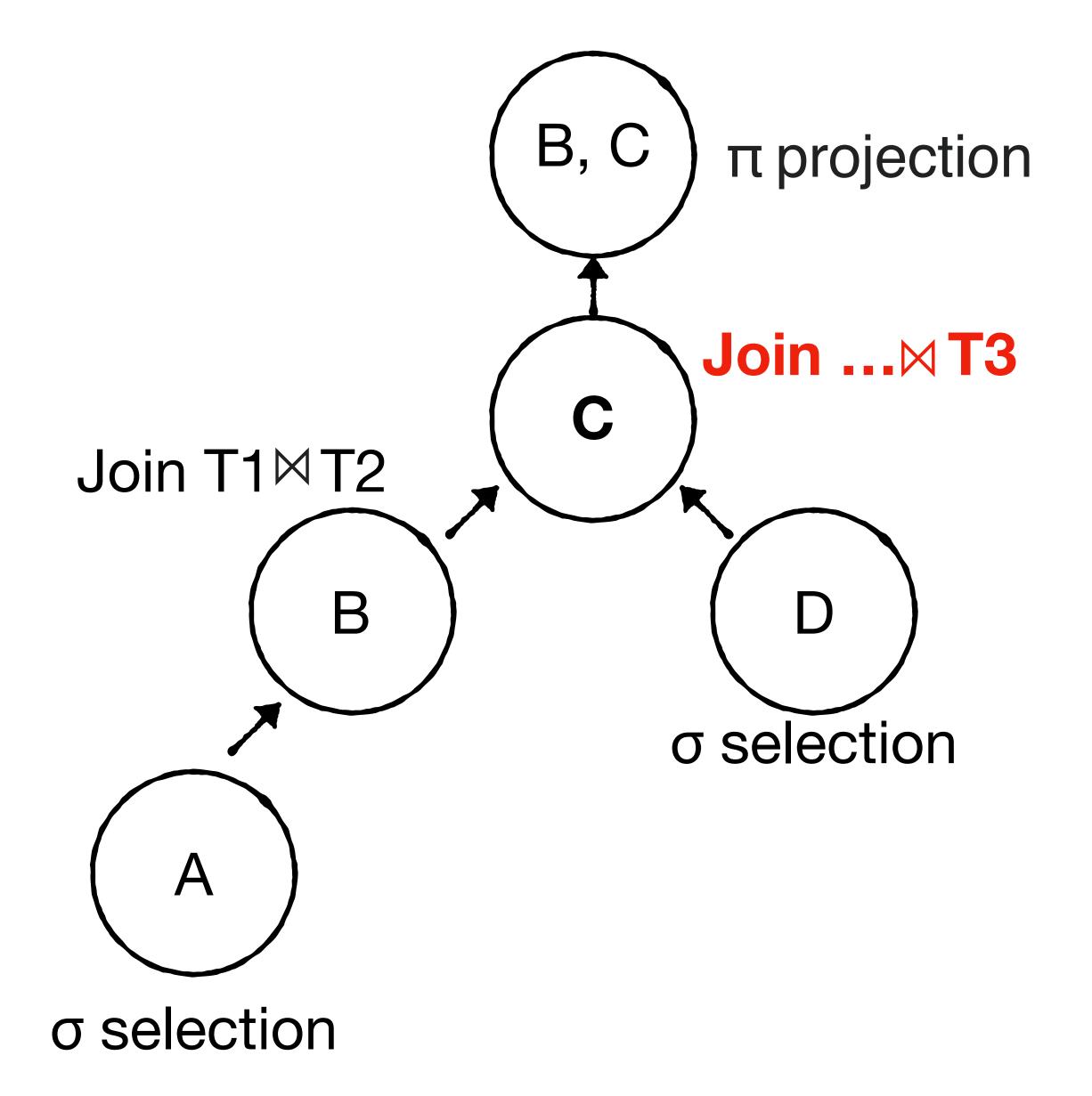




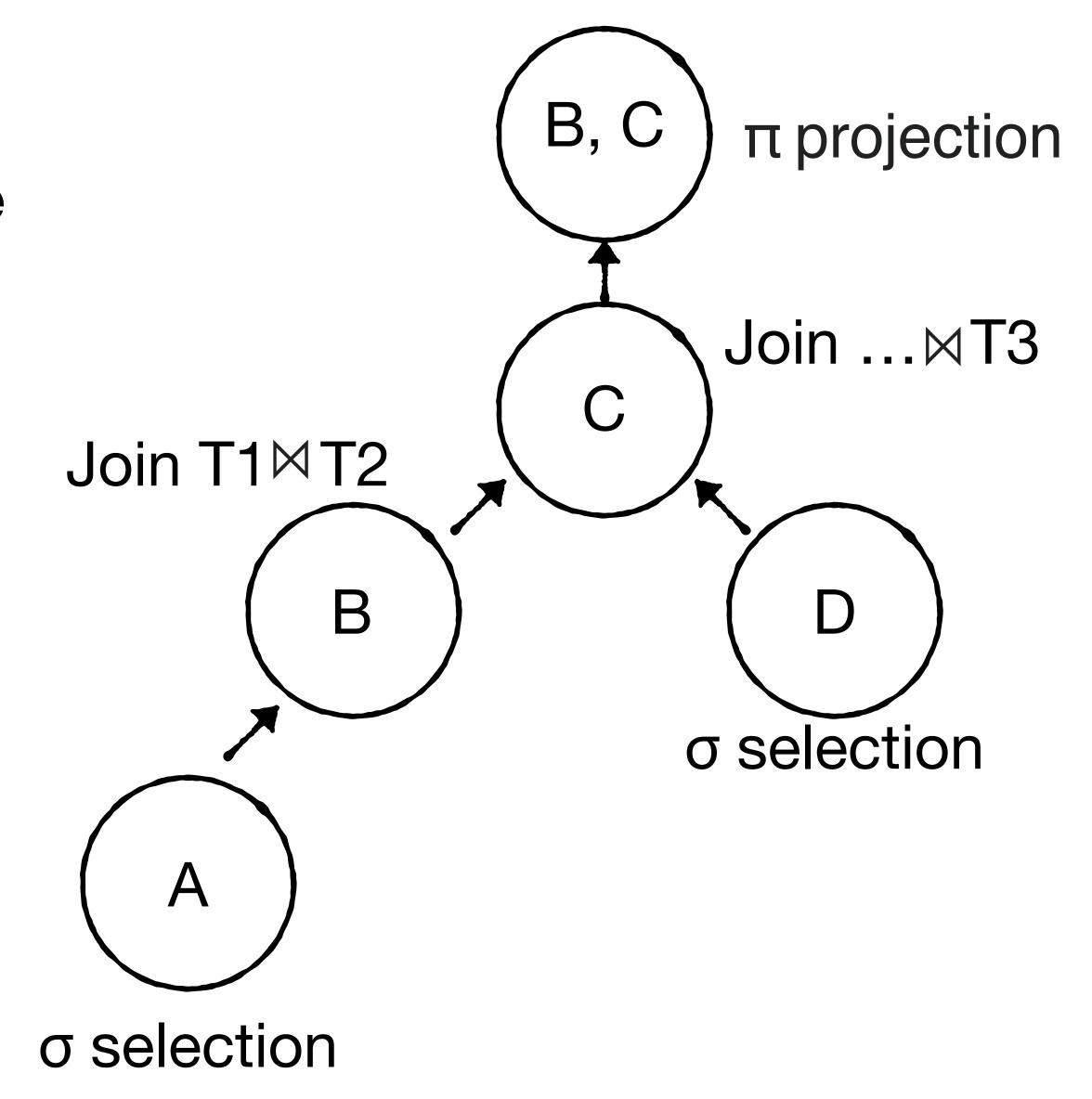








All implement an iterator interface (the glue)



Selection



Order by



Projection



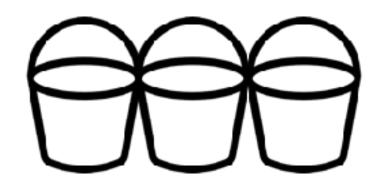
Distinct



Join



Group by



Each can be implemented using different algorithms

Selection



Order by



Projection



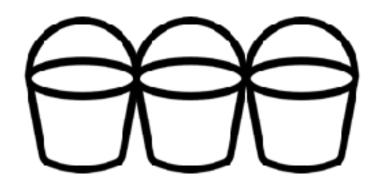
Distinct



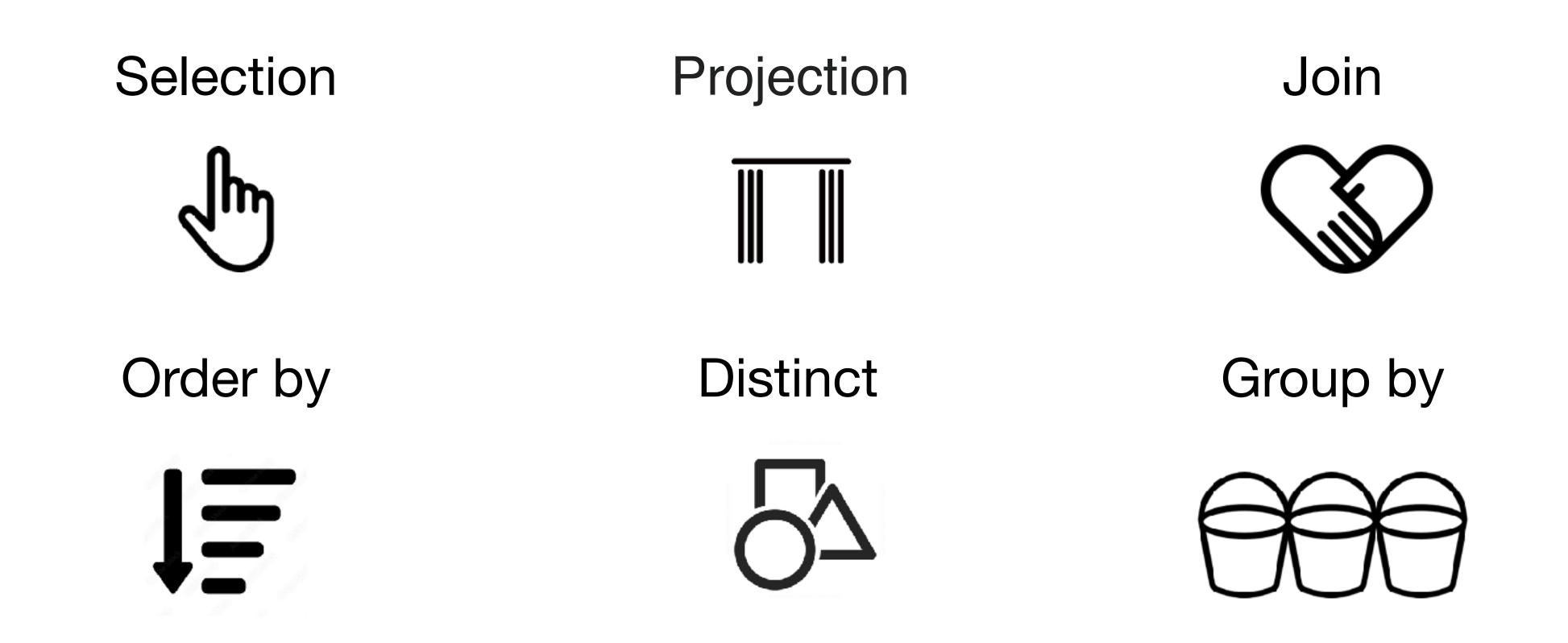
Join



Group by

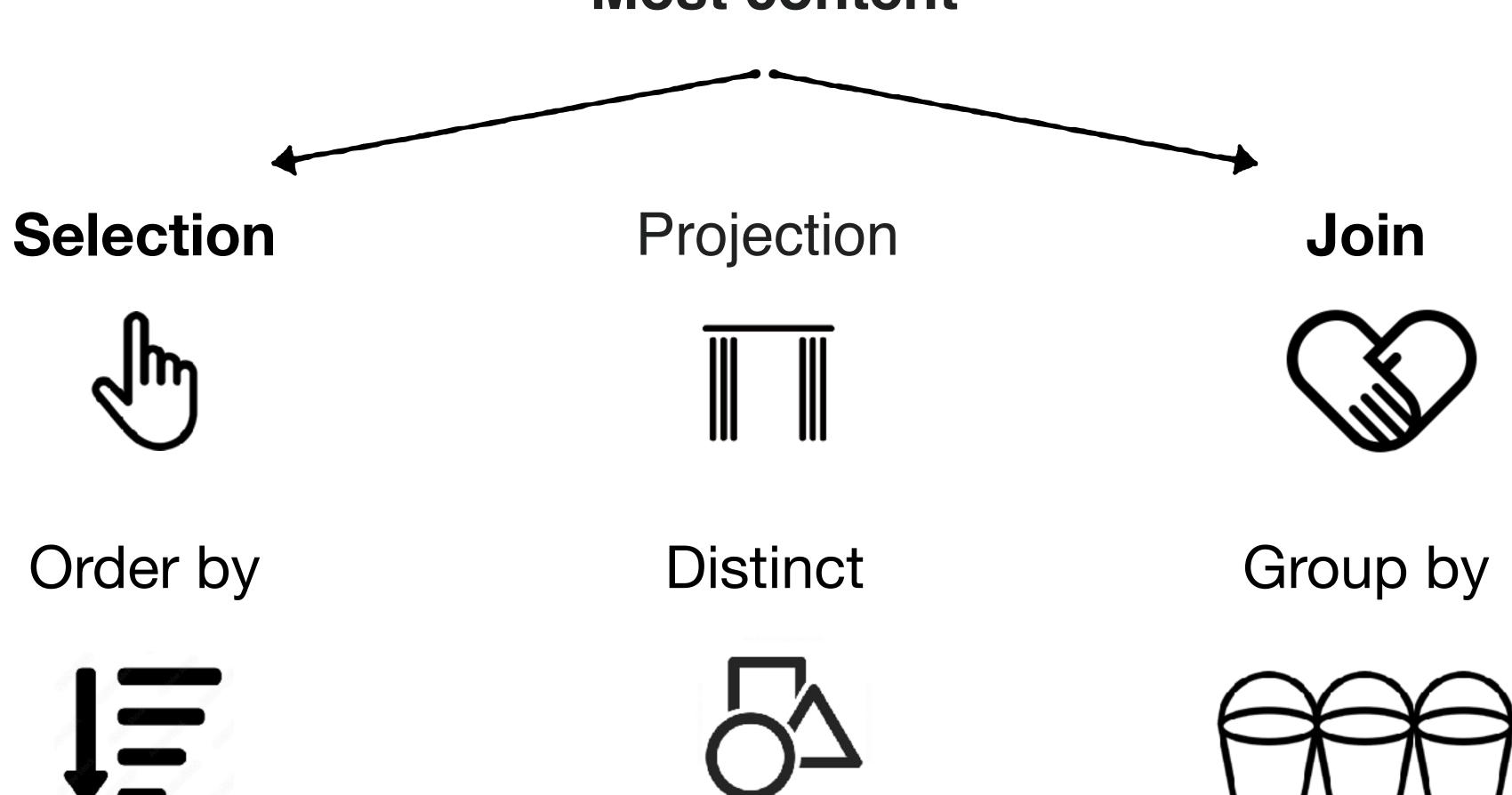


Each can be implemented using different algorithms



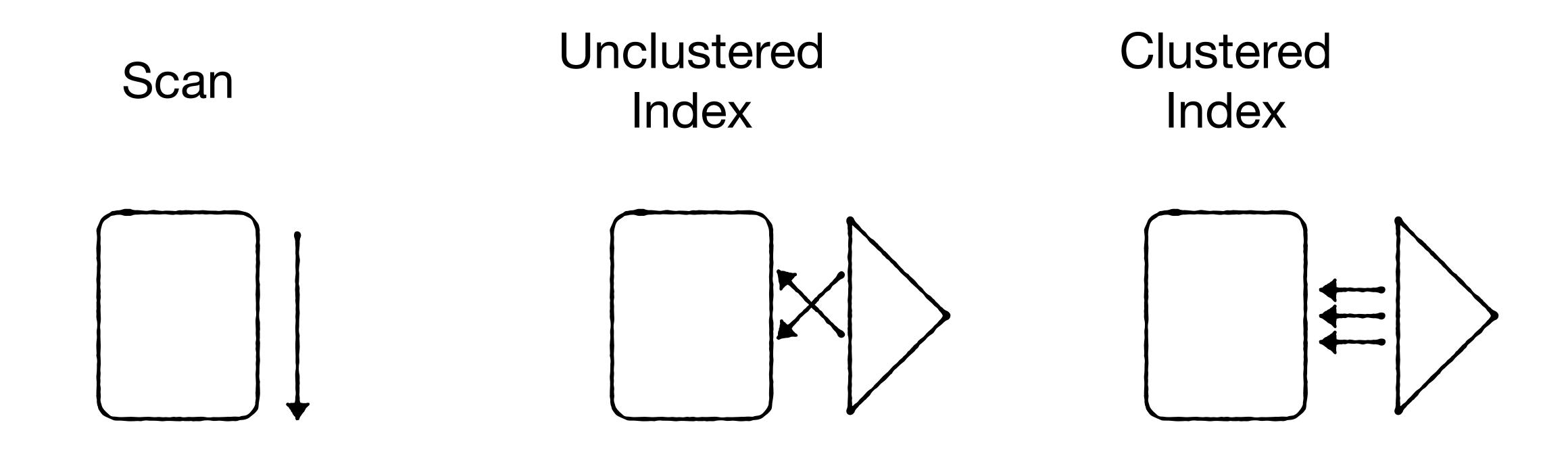
different trade-offs for different contexts

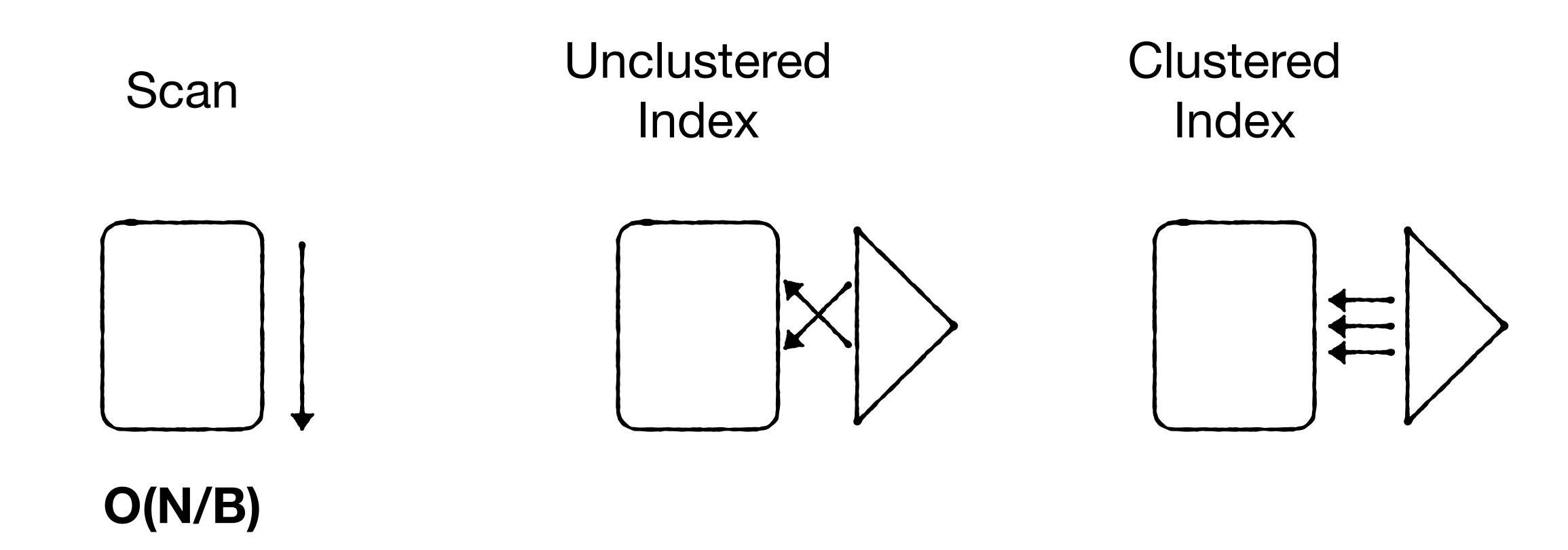
Most content

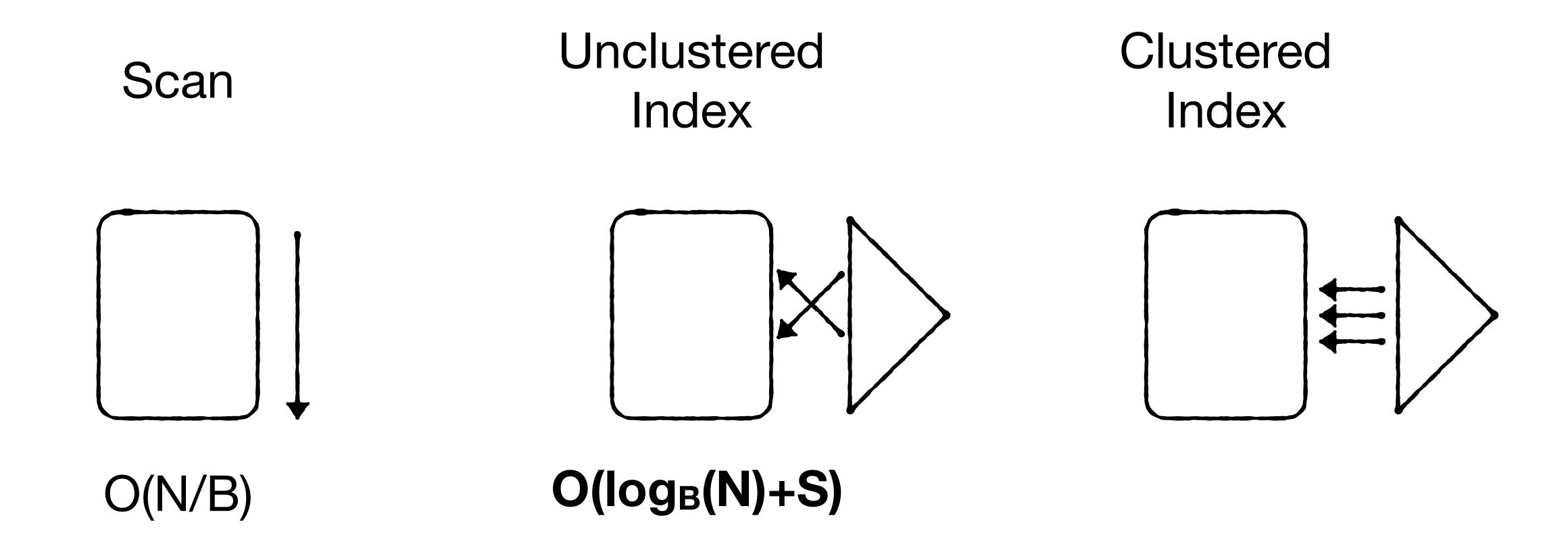


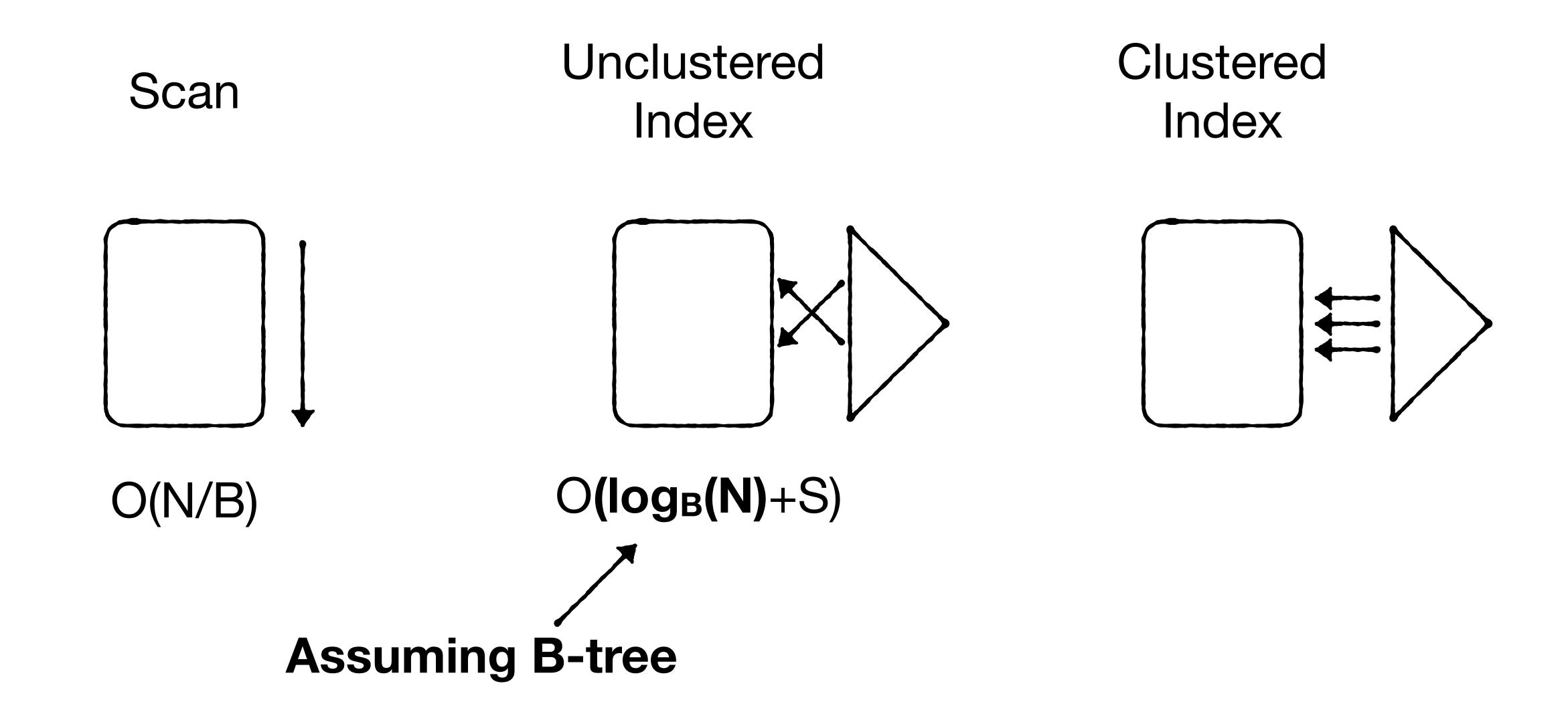
Selection

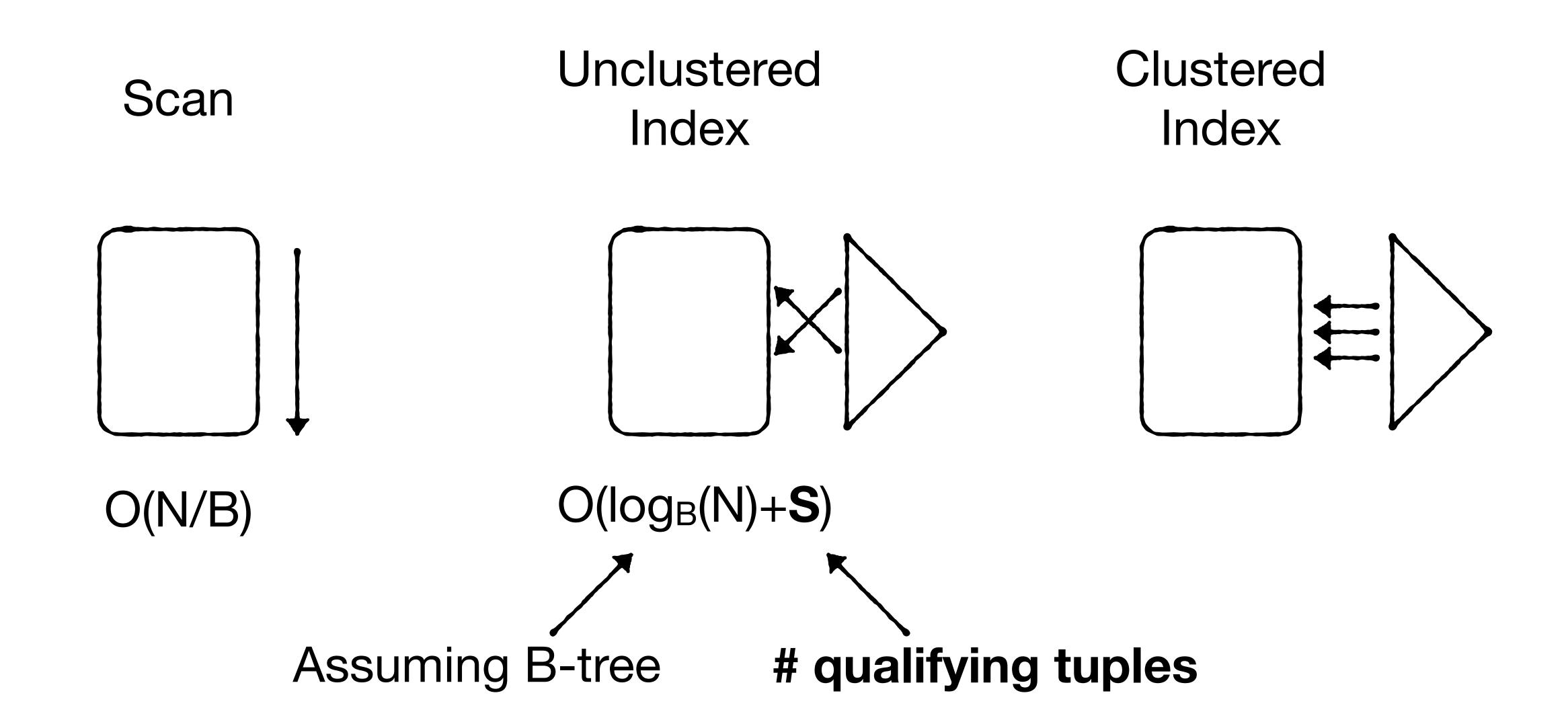
A B

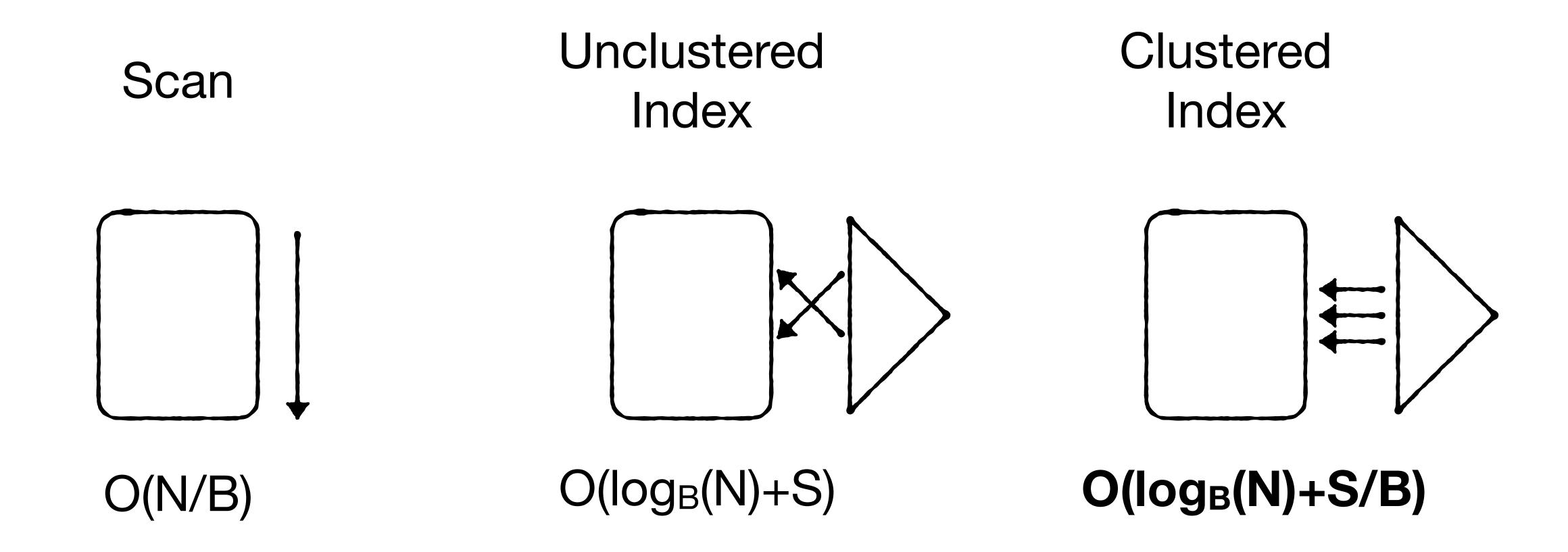




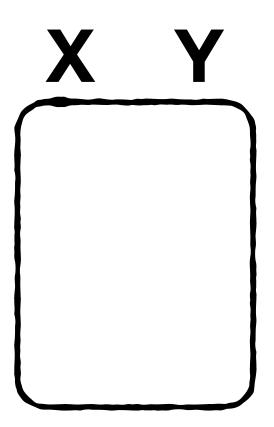






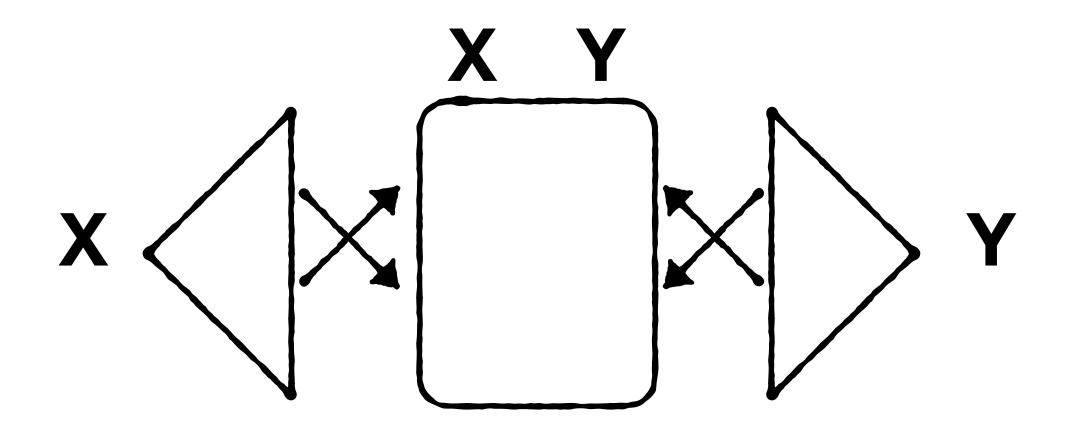


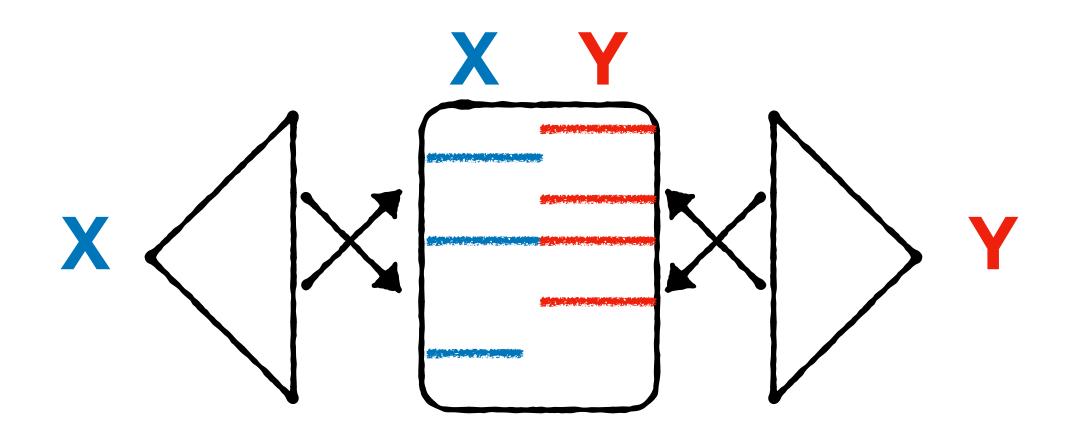
How about when we have multiple selection conditions?



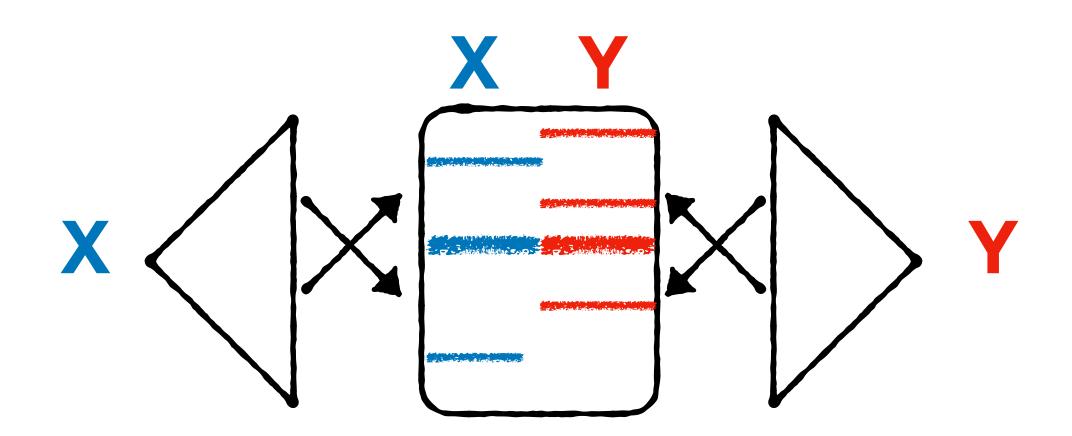
Select * from ... where X = "..." and Y = "..."

Assume two unclustered indexes



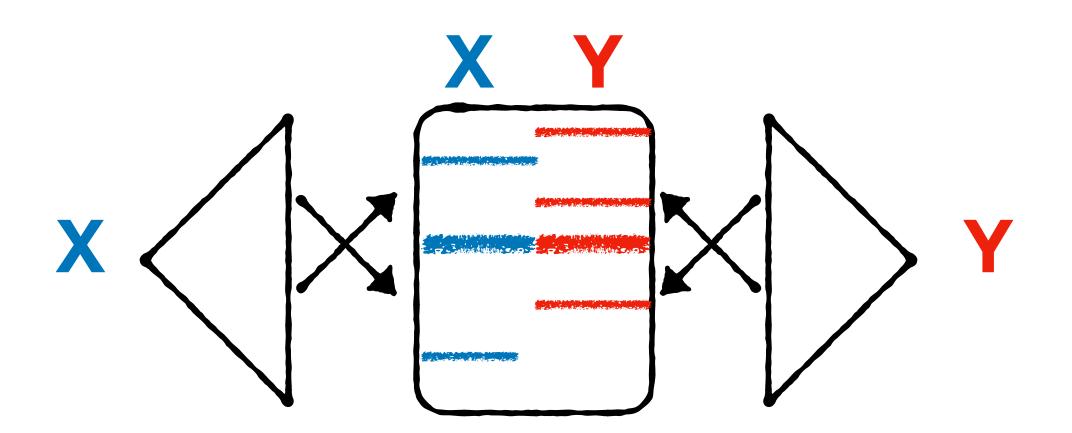


Let $|X_i|$ and $|Y_j|$ denote the number of matching rows to A and to B, resp. e.g., $|X_i| = 3$ and $|Y_j| = 4$



Let $|X_i|$ and $|Y_j|$ denote the number of matching rows to A and to B, resp. e.g., $|X_i| = 3$ and $|Y_j| = 4$

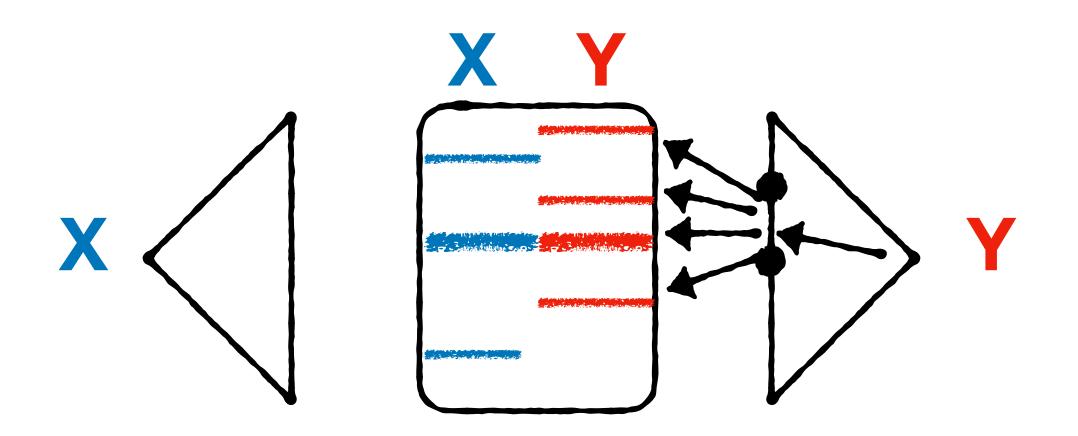
The query is looking for the intersection: $|X_i \cap Y_j| = 1$



Let $|X_i|$ and $|Y_j|$ denote the number of matching rows to A and to B, resp. e.g., $|X_i| = 3$ and $|Y_j| = 4$

The query is looking for the intersection: $|X_i \cap Y_j| = 1$

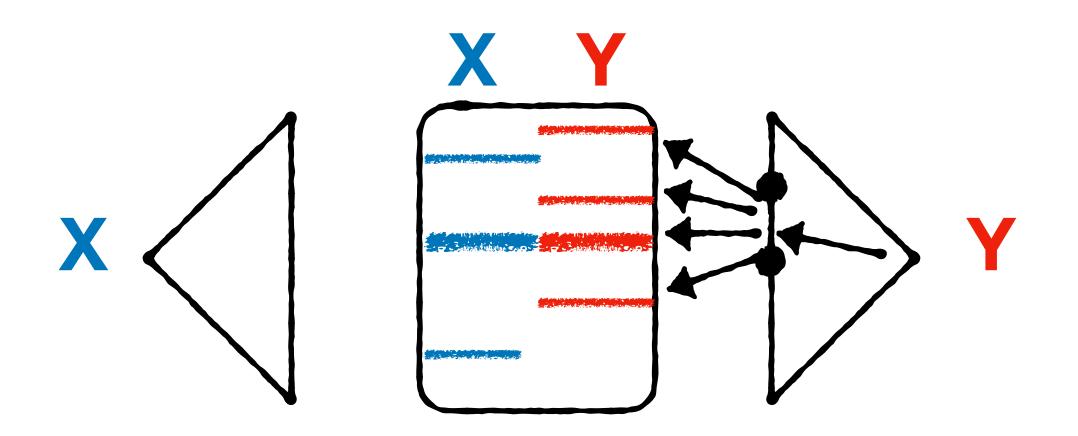
Some ideas?



Algorithm 1: Search index Y

For each row in table, check if X matches

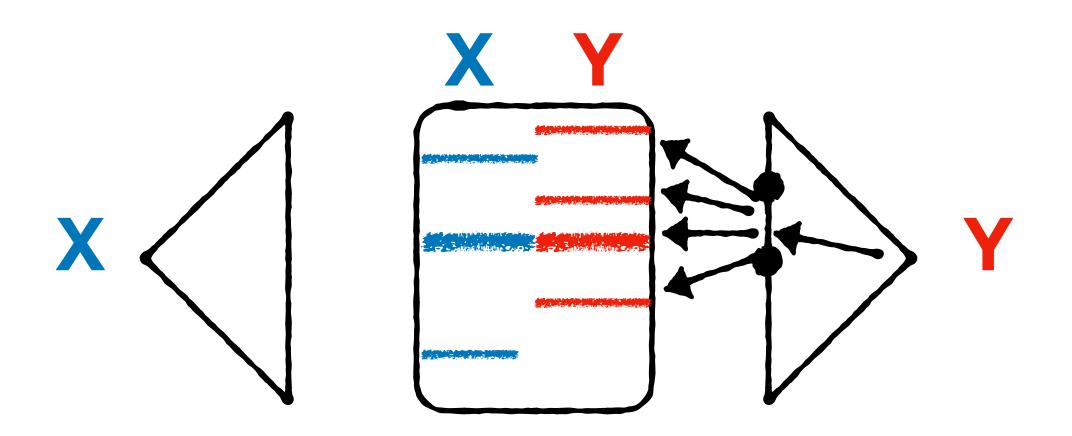
Cost:



Algorithm 1: Search index Y

For each row in table, check if X matches

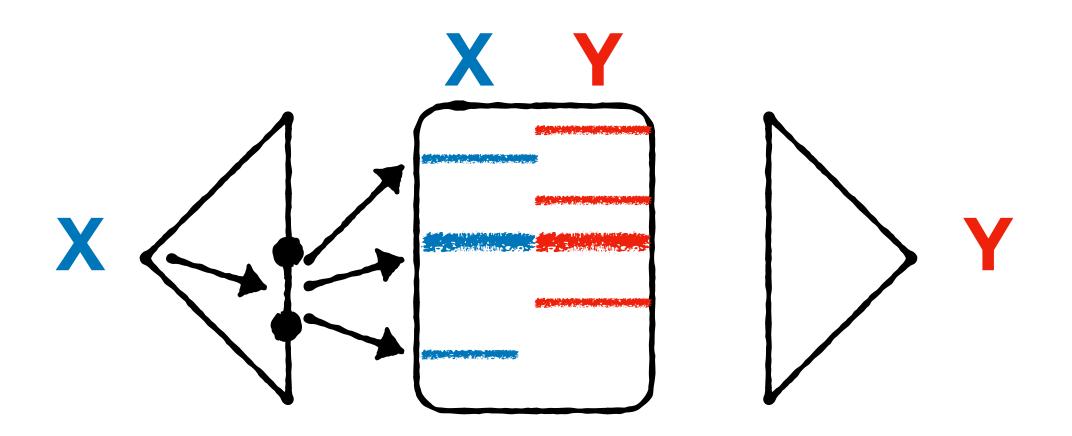
Cost: $log_B(N) + |Y_j|/B + |Y_j| I/O$



Algorithm 1: Search index Y

For each row in table, check if X matches

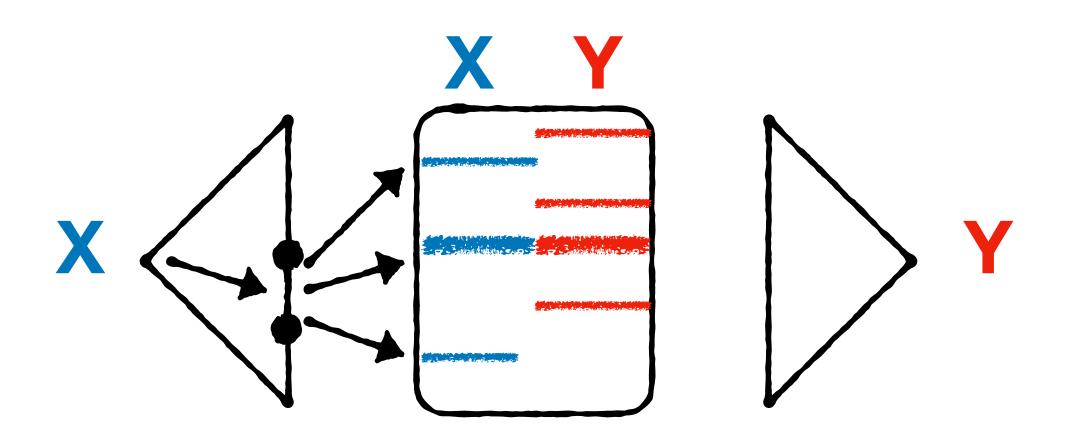
Cost: $log_B(N) + |Y_j| I/O$



Algorithm 2: Search index X

For each row in table, check if Y matches

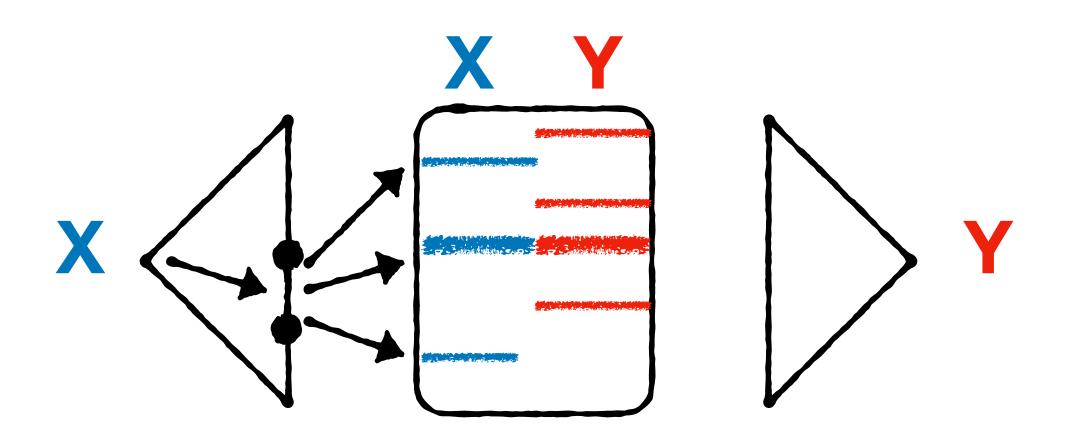
Cost:



Algorithm 2: Search index X

For each row in table, check if Y matches

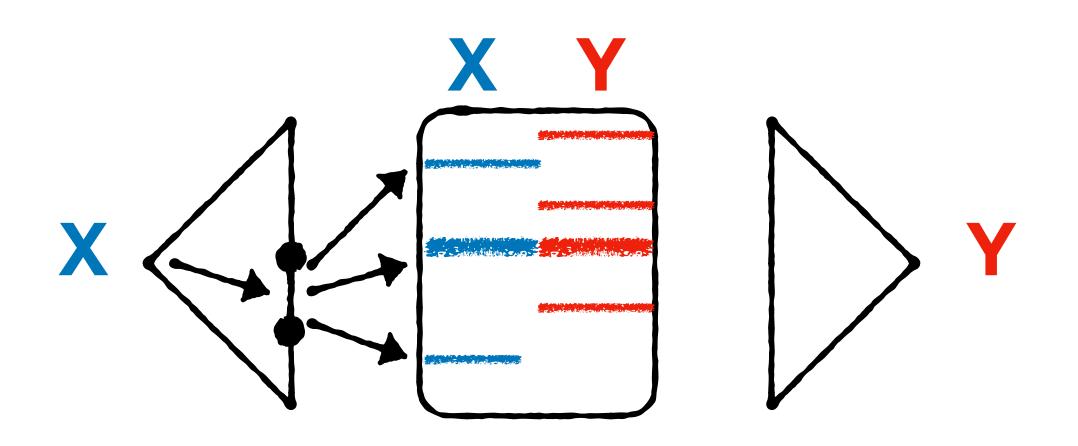
Cost: $log_B(N) + |X_i|/B + |X_i| I/O$



Algorithm 2: Search index X

For each row in table, check if Y matches

Cost: $log_B(N) + |X_i| I/O$

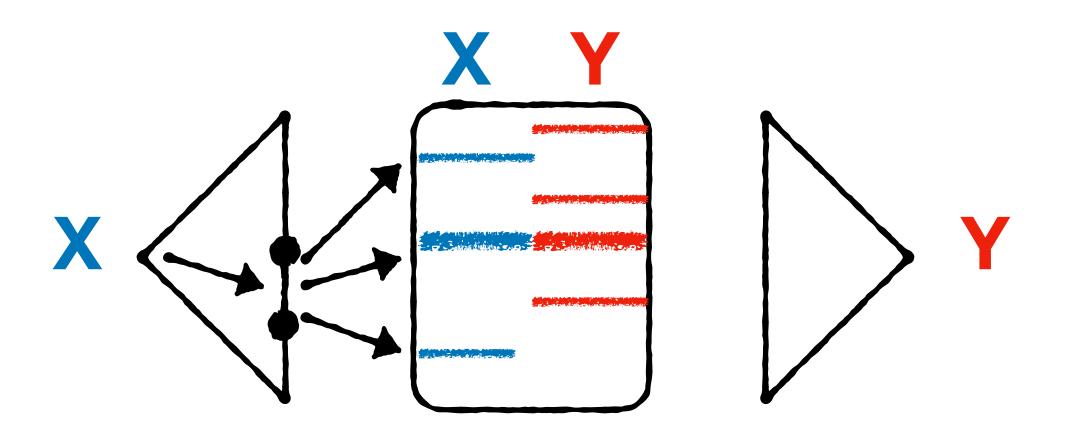


Algorithm 2: Search index X

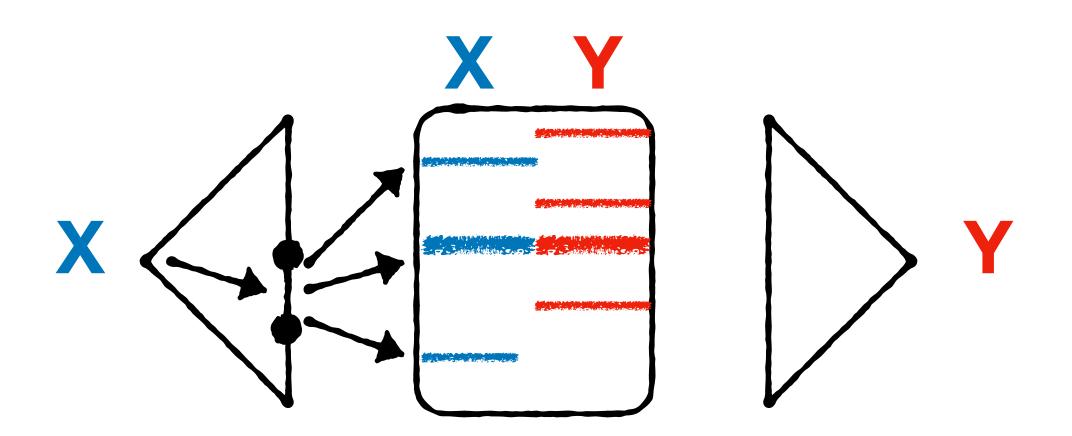
For each row in table, check if Y matches

Cost: $log_B(N) + |X_i| I/O$

Less expensive since $|X_i| < |Y_j|$ ($|X_i| = 3$ and $|Y_j| = 4$)

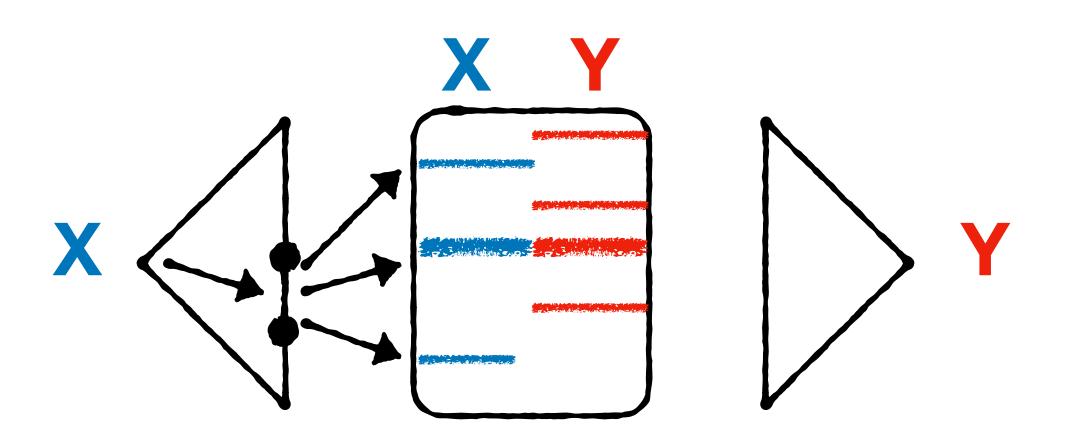


Principle: filter based on most selective predicate first

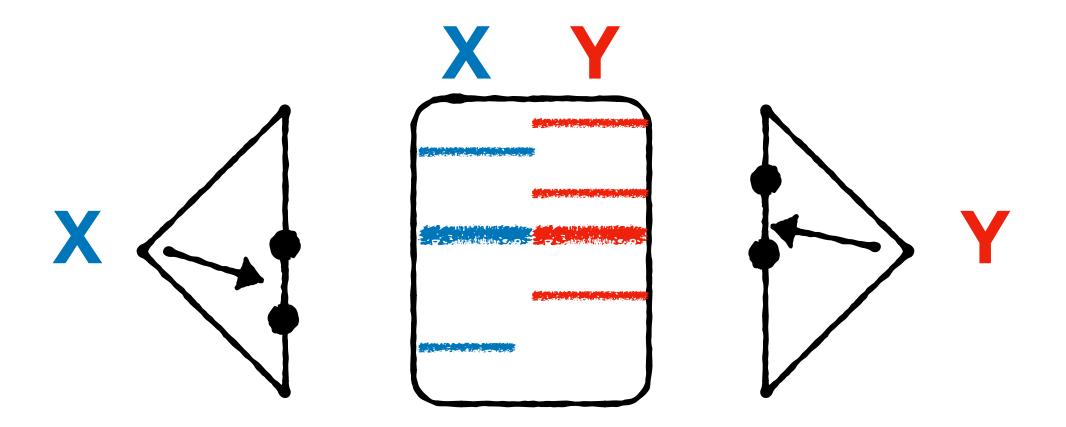


Principle: filter based on most selective predicate first

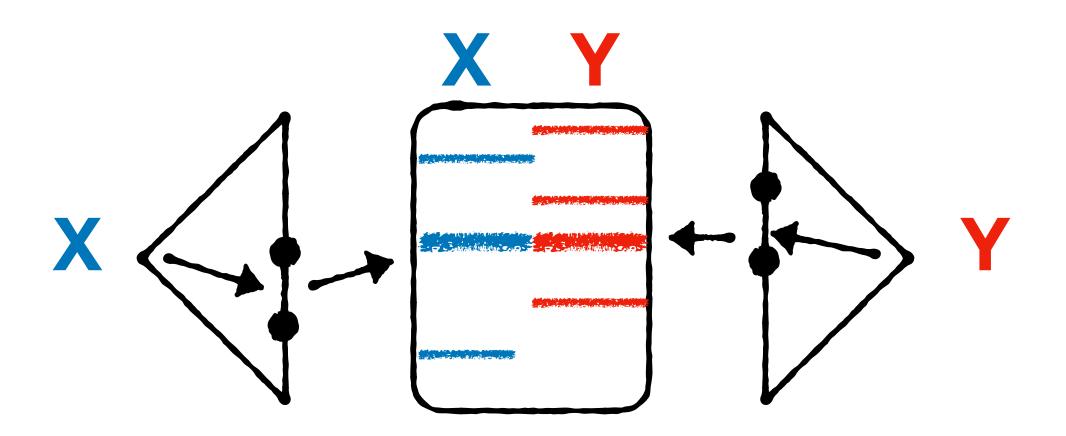
This requires frequency estimation, e.g., estimating |Xi| or |Yj|



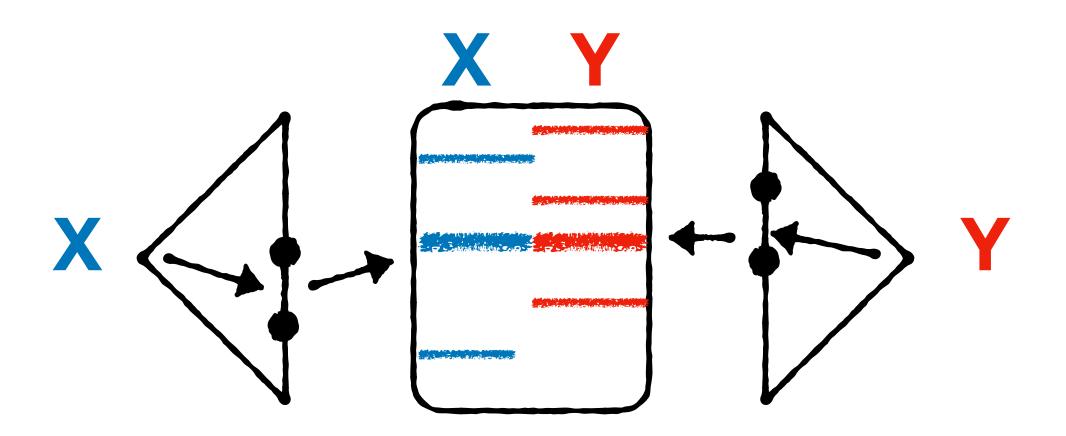
Principle: filter based on most selective predicate first This requires frequency estimation, e.g., estimating $|X_i|$ or $|Y_j|$ Is there yet another alternative?



Algorithm 3: Search index A and B for matching row IDs

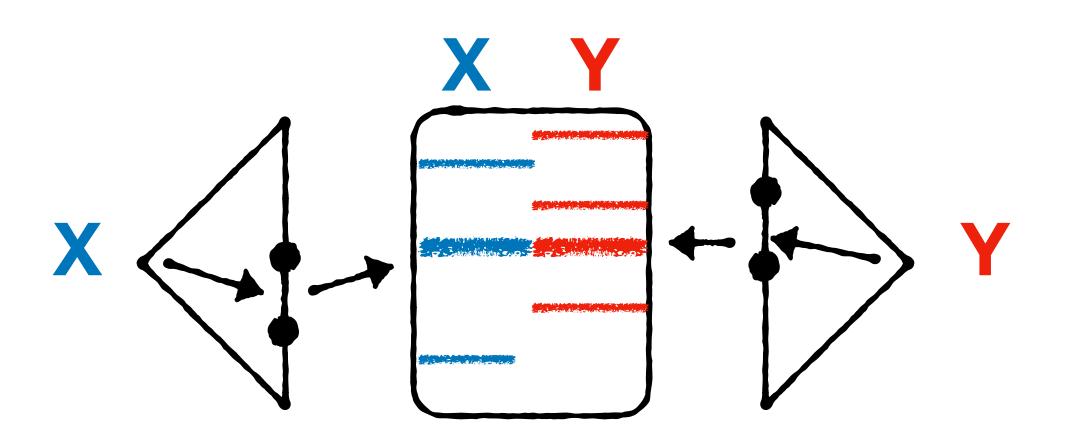


Algorithm 3: Search index A and B for matching row IDs



Algorithm 3: Search index A and B for matching row IDs

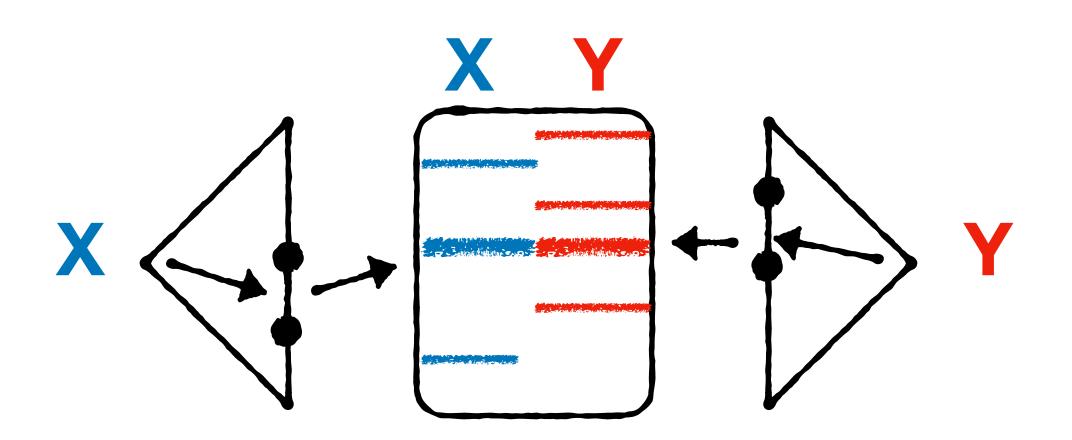
Access table for intersection



Algorithm 3: Search index A and B for matching row IDs

Access table for intersection

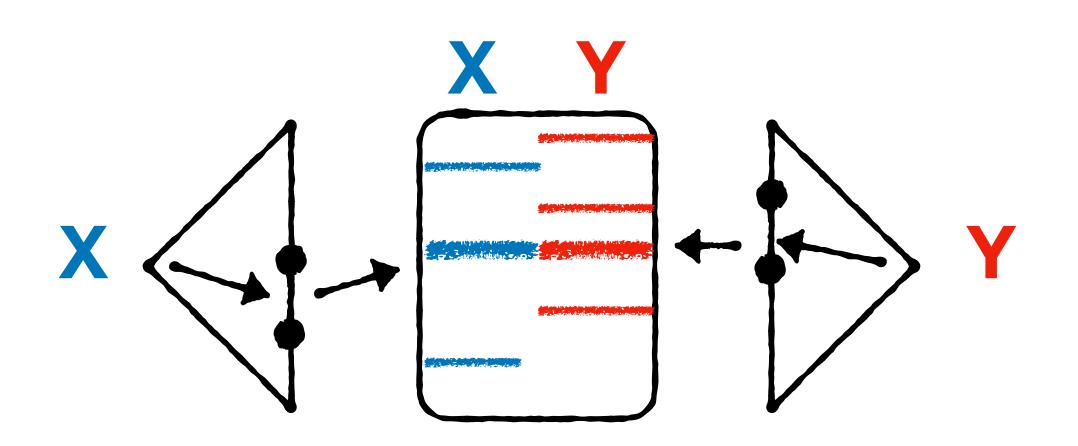
Cost:



Algorithm 3: Search index A and B for matching row IDs

Access table for intersection

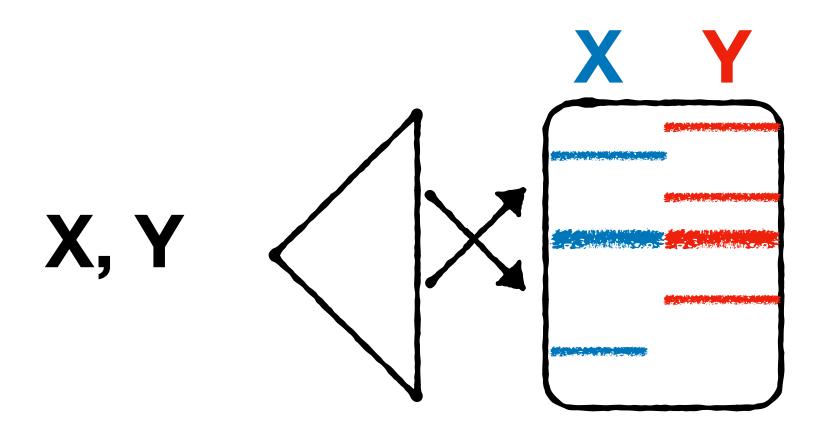
Cost: $2 \cdot \log_B(N) + |X_i|/B + |Y_j|/B + |X_i \cap Y_j| I/O$



Algorithm 3: Search index A and B for matching row IDs
Only search table for intersection

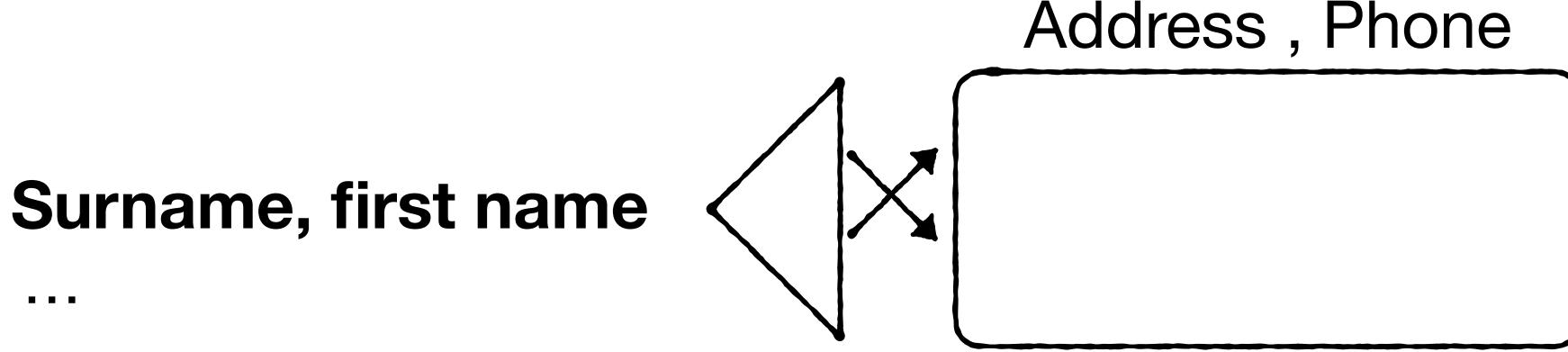
Cost: $2 \cdot \log_B(N) + |X_i|/B + |Y_j|/B + |X_i \cap Y_j| I/O$

Better if $|X_i \cap Y_j|$ is small relative to $|X_i|$, $|Y_j|$



What if we combine these indexes into a composite index?

Select * from ... where surname = "Lovelace" and firstname="Ada"



Long, Justin

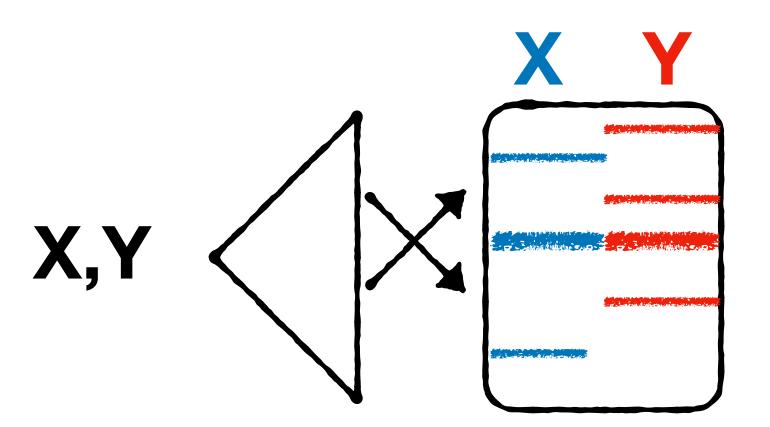
Lovelace, Abby

Lovelace, Ada

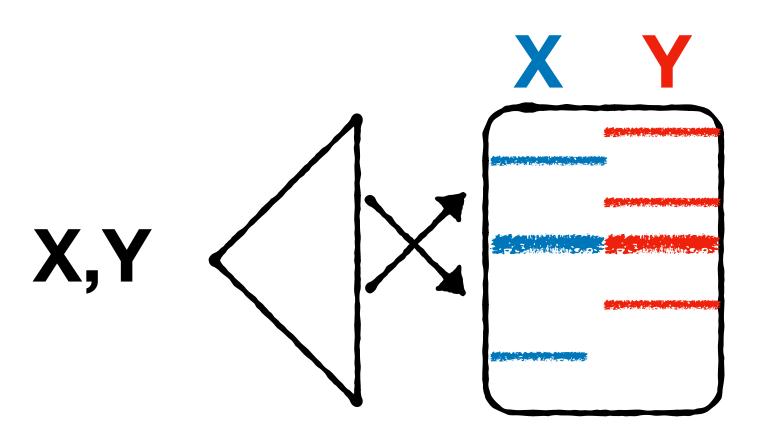
Lovelace, Adam

Lopez, Jenniffer

. . .

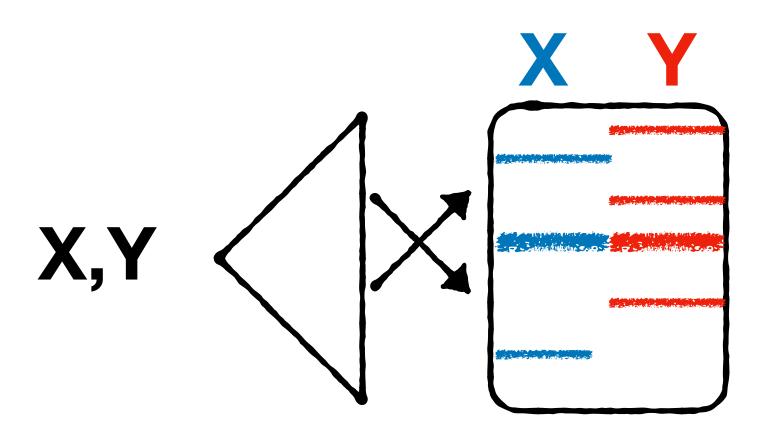


Algorithm 4: Search composite index A and B for matching row IDs



Algorithm 4: Search composite index A and B for matching row IDs

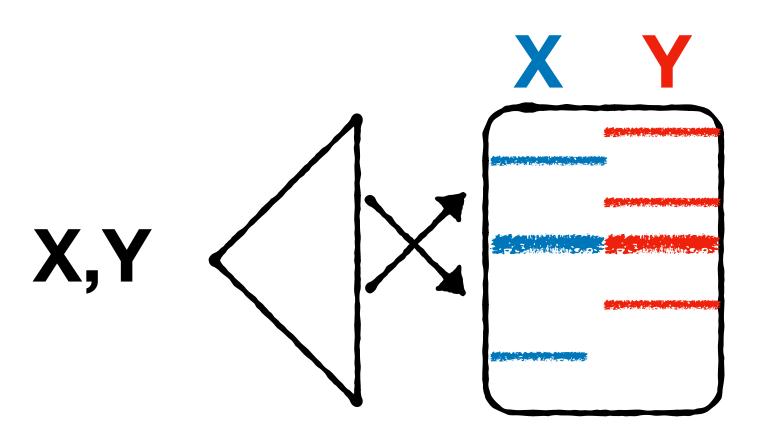
Cost: $log_B(N) + |X_i \cap Y_j| I/O$



Algorithm 4: Search composite index A and B for matching row IDs

Cost: $log_B(N) + |X_i \cap Y_j| I/O$

Cheapest option so far!



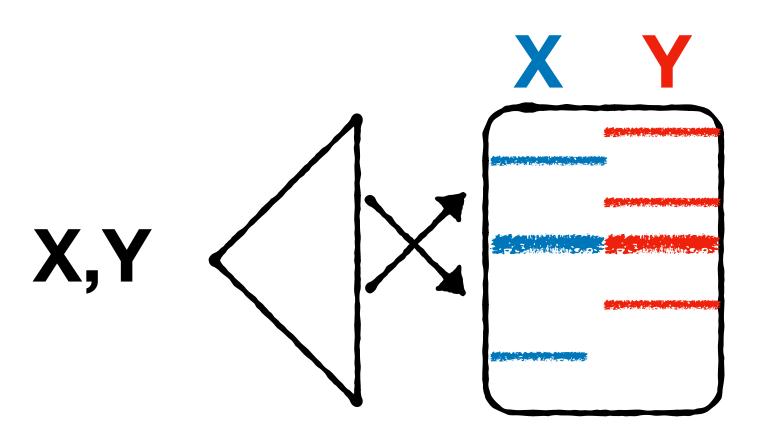
Algorithm 4: Search composite index A and B for matching row IDs

Cost: $\log_B(N) + |X_i \cap Y_j| I/O$

Cheapest option so far!

Downsides?

Select * from ... where Y="..."



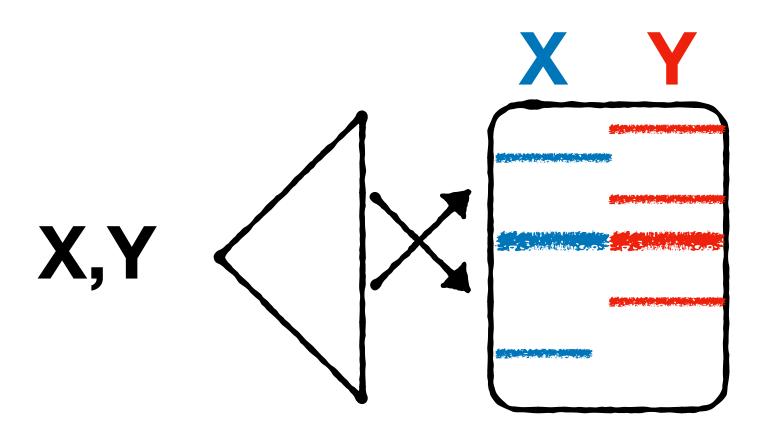
Algorithm 4: Search composite index A and B for matching row IDs

Cost: $log_B(N) + |X_i \cap Y_j| I/O$

Cheapest option so far!

Downsides? Cannot handle queries just based on Y

Select * from ... where Y="..."



Composite indexes can achieve better performance for some queries but are less generic. Choose them carefully:)

Selection



Order by



Projection



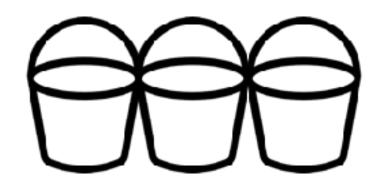
Distinct



Join

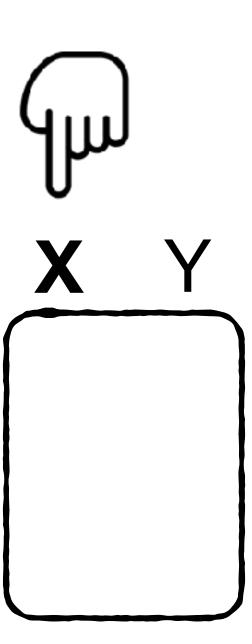


Group by



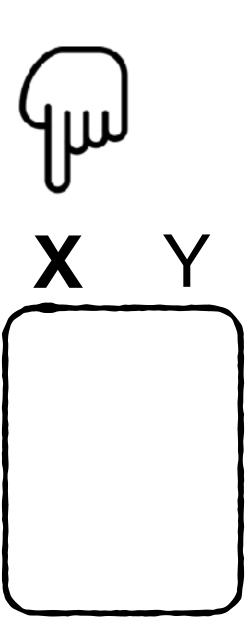
Projection

Select X from table



Projection

Select X from table

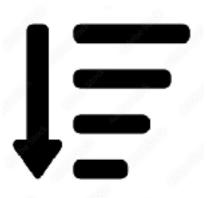


Trivial in row-stores. More interesting in column-stores (next week).

Selection



Order By



Projection



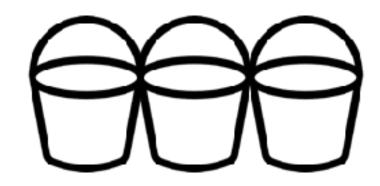
Distinct



Join

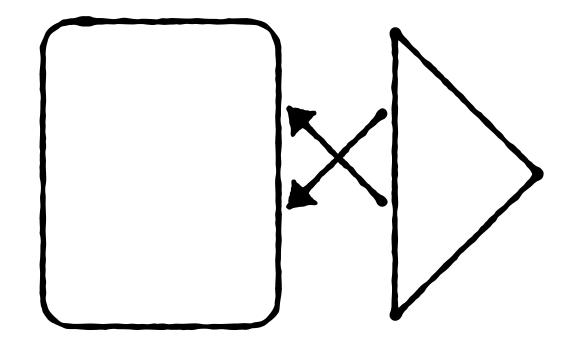


Group by

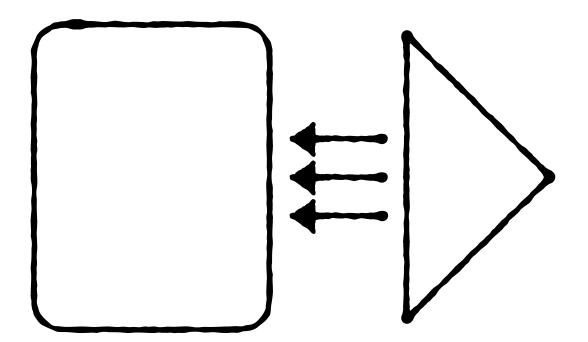


If we have an applicable index

Unclustered index scan

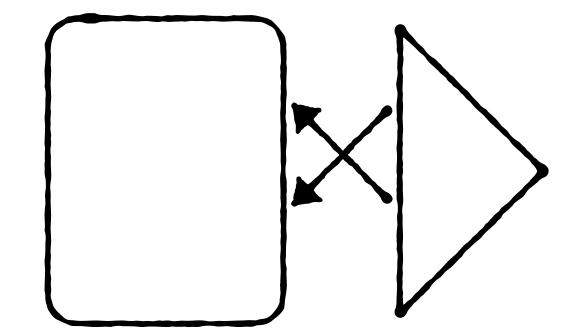


Clustered index scan



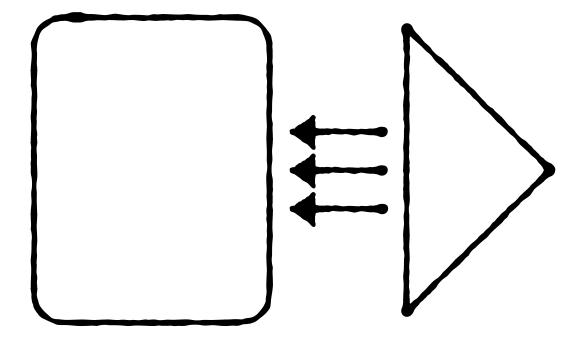
If we have an applicable index

Unclustered index scan

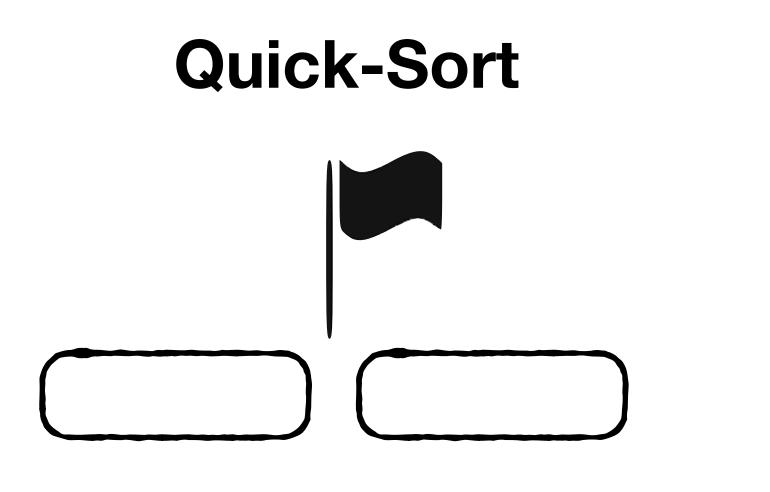


(Good if result set is small)

Clustered index scan

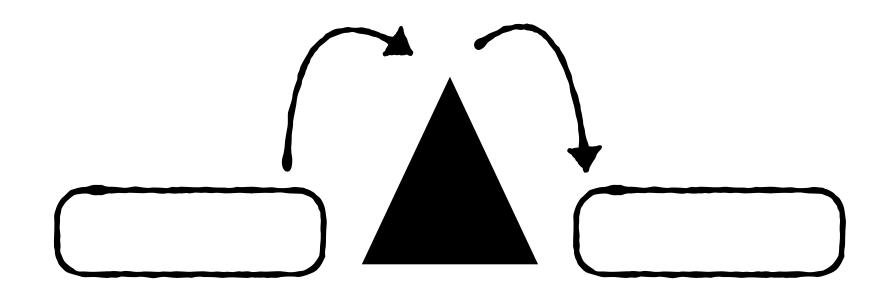


(Good in all cases)

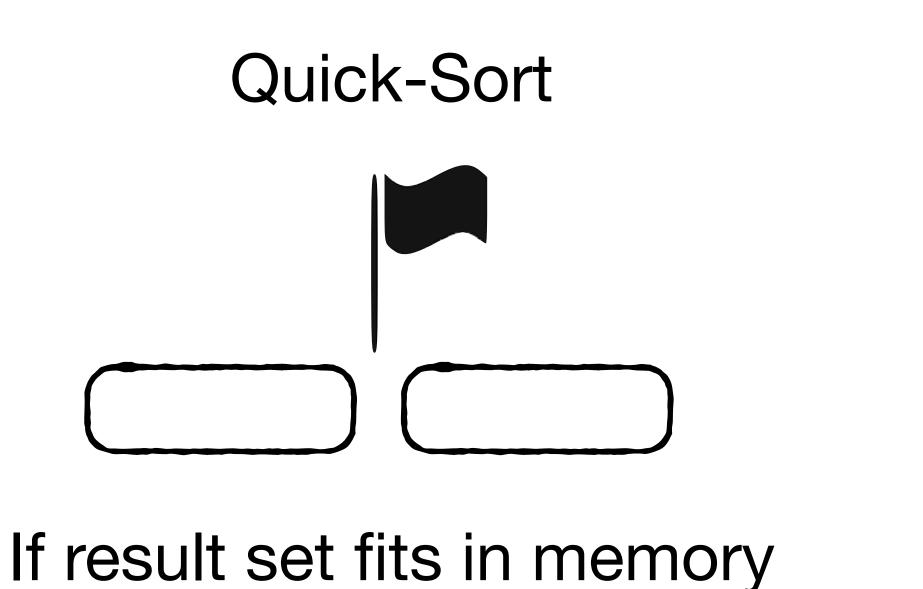


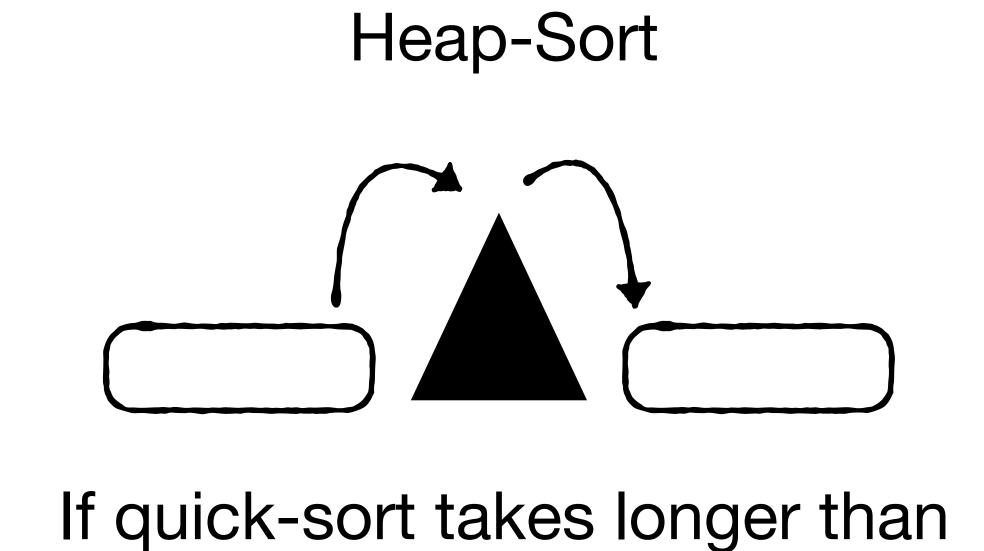






If quick-sort takes longer than expected

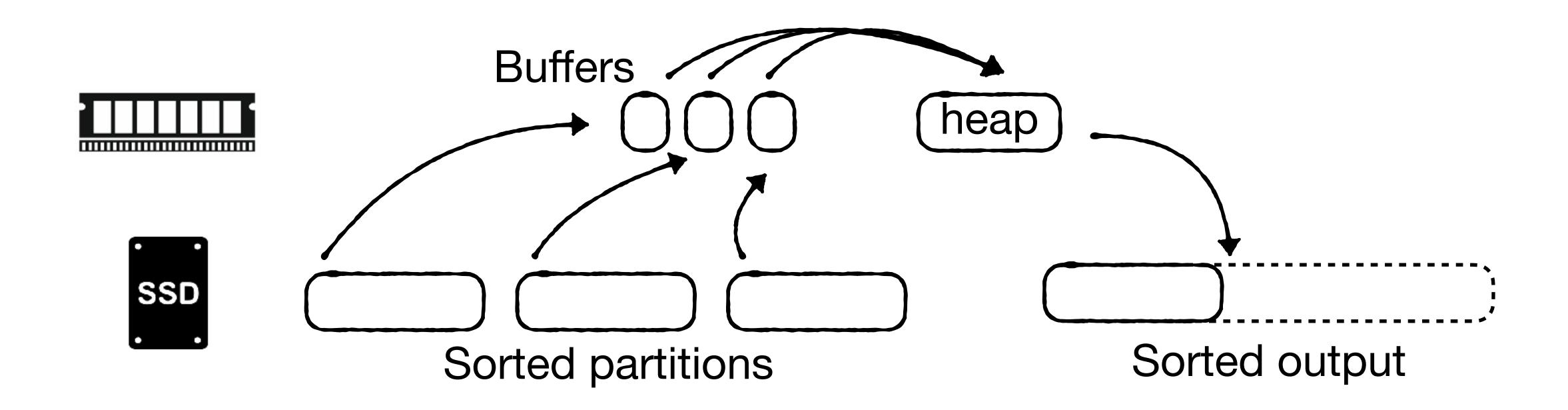




expected

Introsort:)

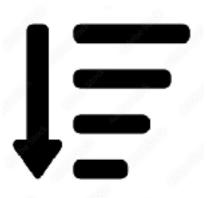
Or external sort if data does not fit in memory



Selection



Order By



Projection



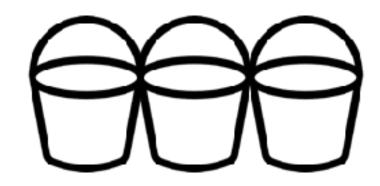
Distinct



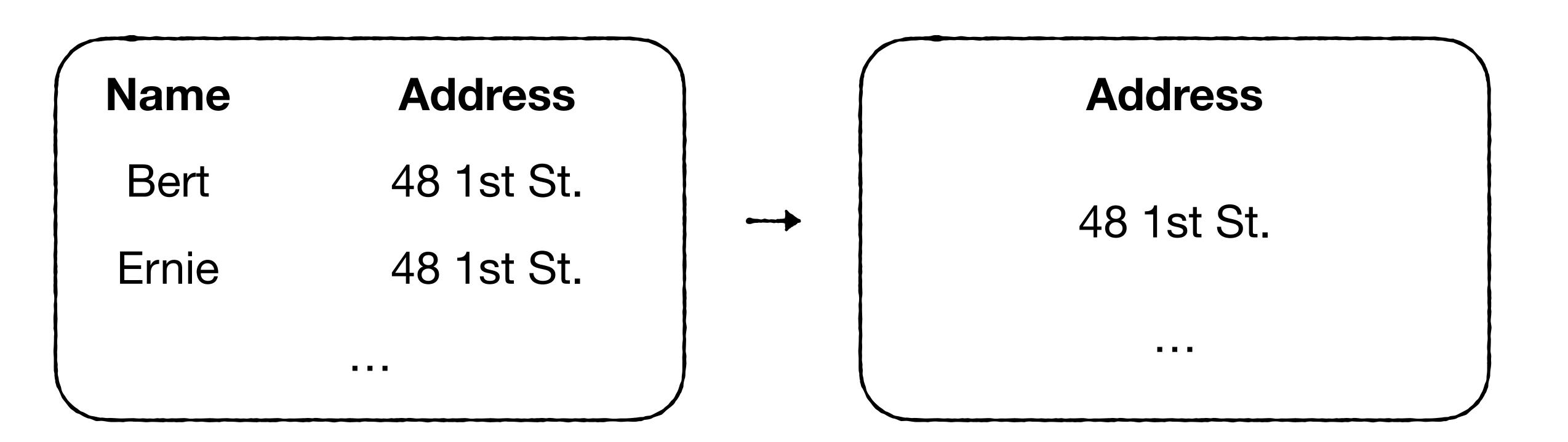
Join



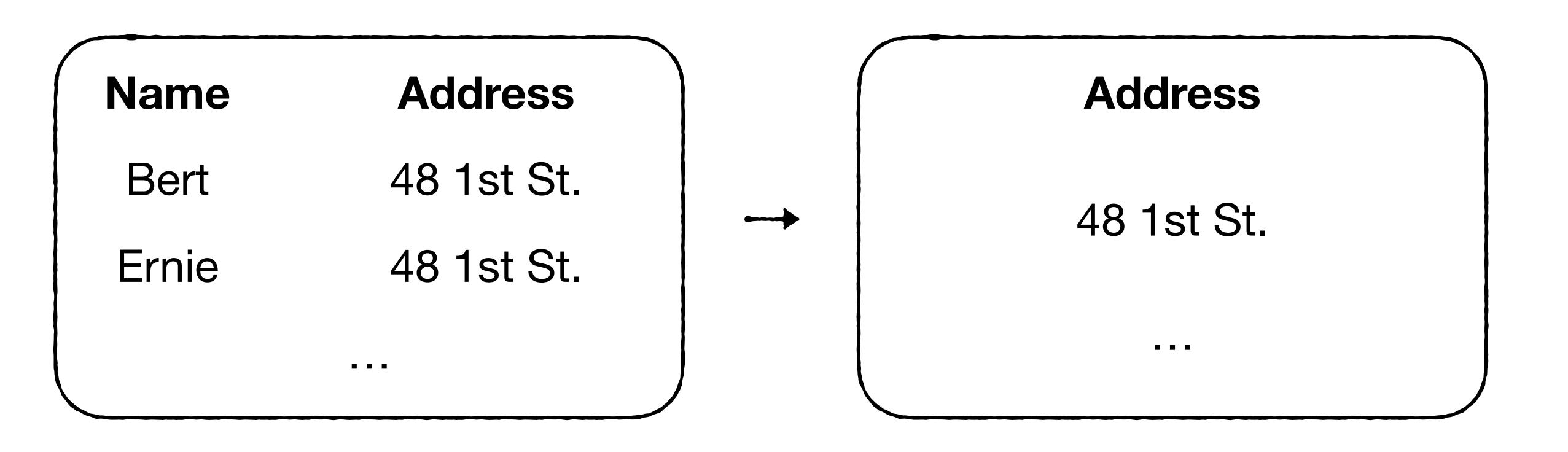
Group by



Select Distinct address FROM people;



Select Distinct address FROM people;



How to implement?

Sort & eliminate adjacent identical items



Sort & eliminate adjacent identical items

Quick-Sort
Heap-Sort
External Sort
Index scan
Etc.



Sort & eliminate adjacent identical items

Quick-Sort
Heap-Sort
External Sort
Index scan
Etc.

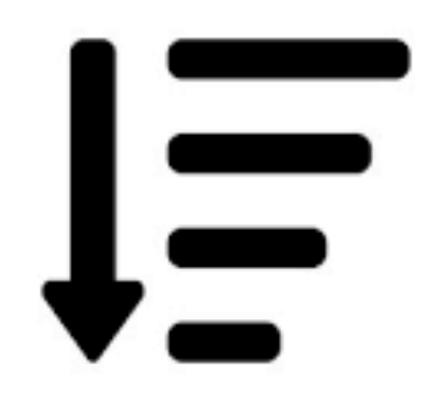


Hash to identify identical items



Sort & eliminate adjacent identical items

Quick-Sort
Heap-Sort
External Sort
Index scan
Etc.

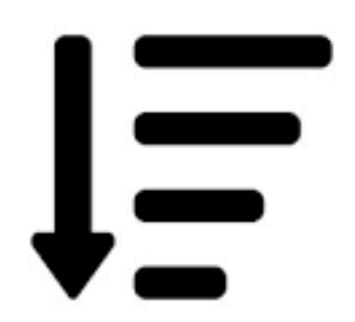


Hash to identify identical items



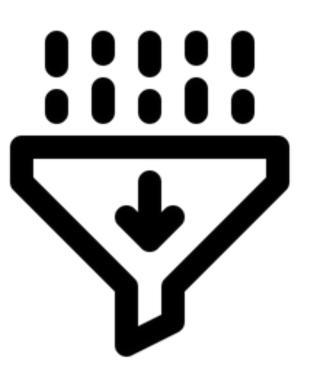
Linear probing
Chaining
Extendible hashing
Cuckoo hashing
Etc.

Sort & eliminate adjacent identical items



CPU: O(N log₂ N)

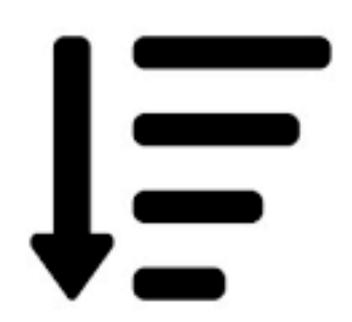
Hash to identify identical items



O(N)

Select Distinct c1, c2, ... FROM ... order by c1, c2 ...;

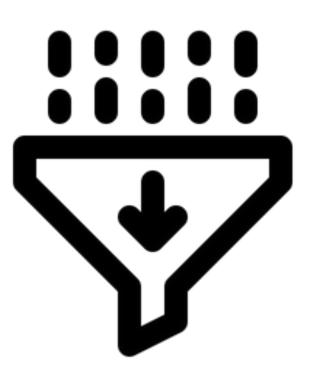
Sort & eliminate adjacent identical items



CPU: O(N log₂ N)

Better if we later need to sort anyways

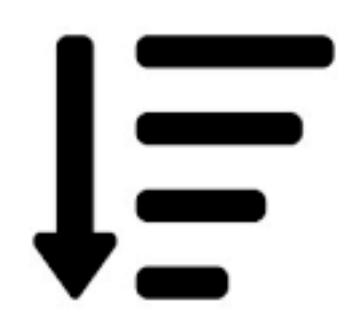
Hash to identify identical items



O(N)

Select Distinct c1, c2, ... FROM ...;

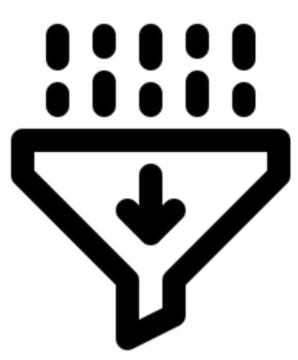
Sort & eliminate adjacent identical items



CPU: O(N log₂ N)

Better if we later need to sort anyways

Hash to identify identical items



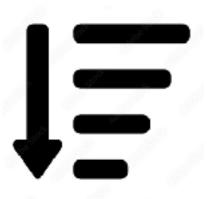
(N)C

Good if user is fine with unordered output

Selection



Order by



Projection



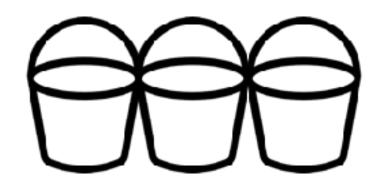
Distinct



Join

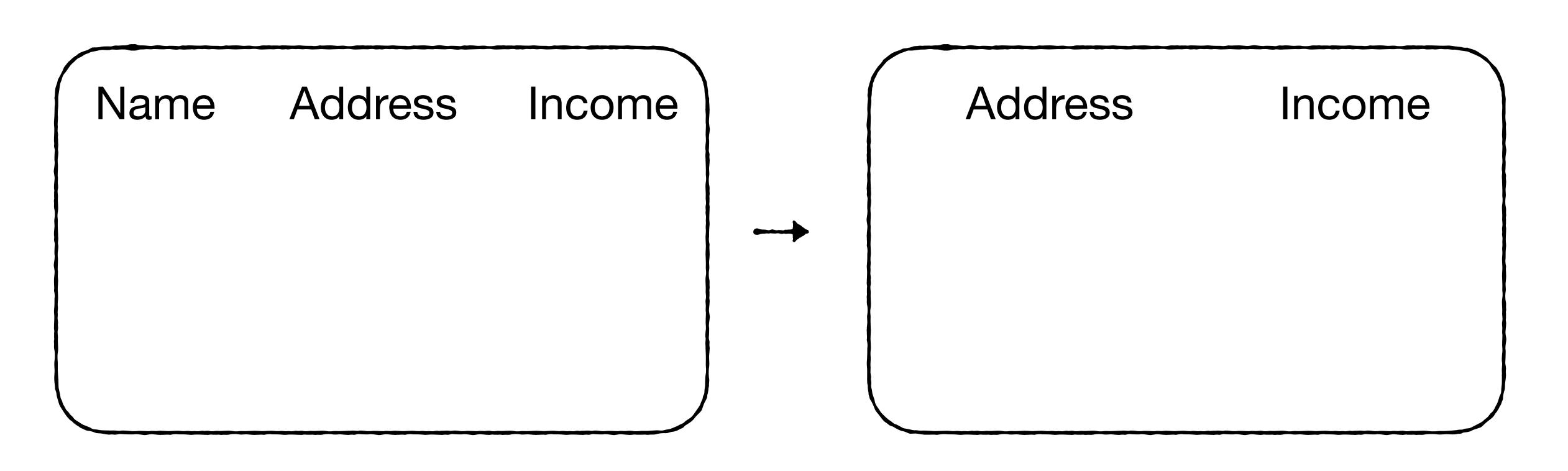


Group By



Select c1, c2, ... FROM ... Group By

Select address, sum(income) FROM people **Group By** address (Income per household)



Select address, sum(income) FROM people **Group By** address (Income per household)

Name Address Income

Bert 48 1st St. 100K

Ernie 48 1st St. 100K

Address Income

48 1st St. 200K

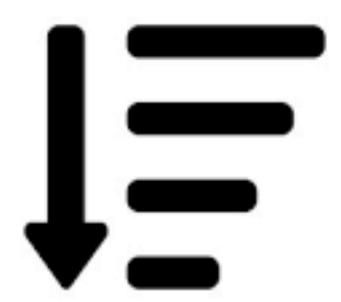
Select address, sum(income) FROM people **Group By** address (Income per household)

Name	Address	Income	Address	Income
Bert	48 1st St.	100K	48 1st St.	200K
Ernie	48 1st St.	100K		

How to execute?

Select address, sum(income) FROM people Group By address

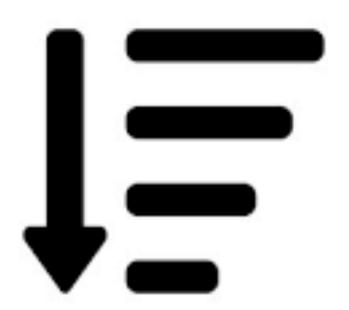
Index/Sort by Category



e.g., B-tree (address -> sum(income))

Select address, sum(income) FROM people Group By address

Index/Sort by Category



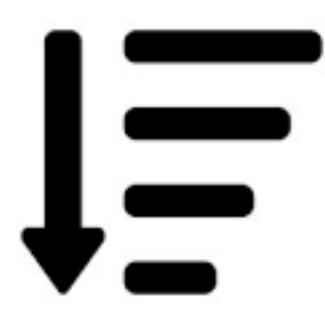
Hash by Category



e.g., also (address -> sum(income))

Select c1, c2, ... FROM ... Group By

Index/Sort by Category



CPU: O(N log₂ N)

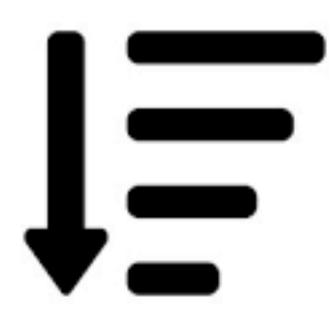
Hash by Category



O(N)

Select c1, c2, ... FROM ... Group By

Index/Sort by Category



CPU: $O(N log_2 N)$

Better if we later need to sort anyways

Hash by Category



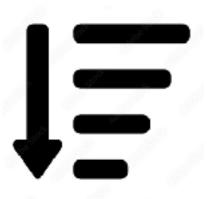
O(N)

Good if user is fine with unordered output

Selection



Order by



Projection



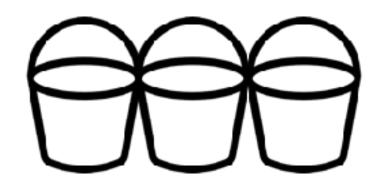
Distinct



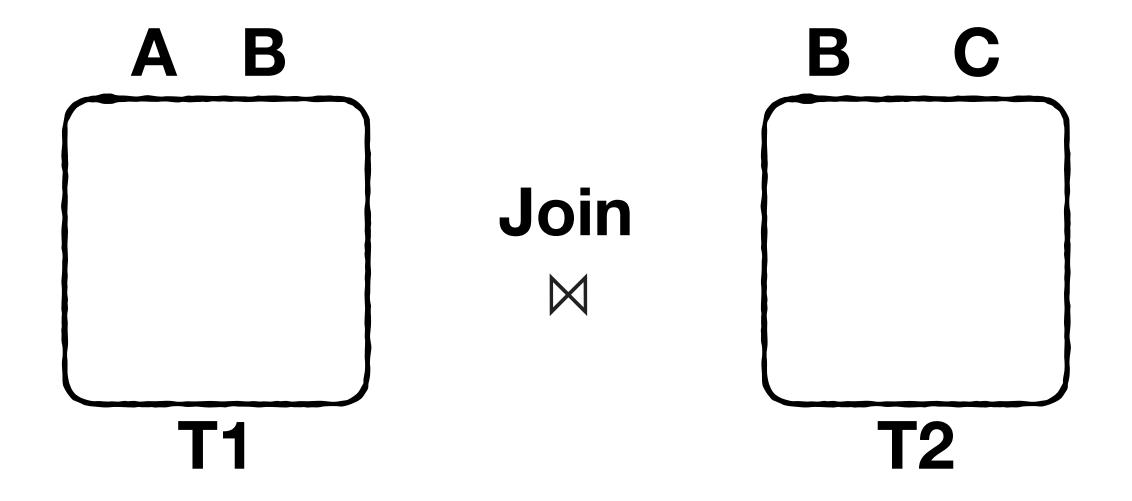
Join



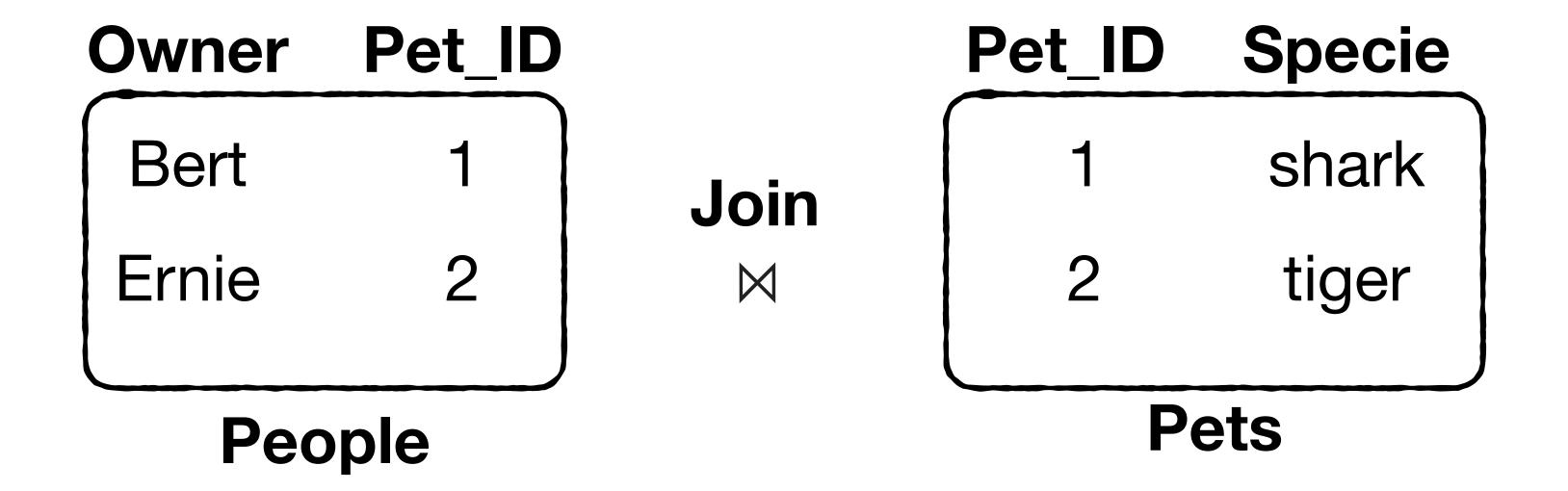
Group By



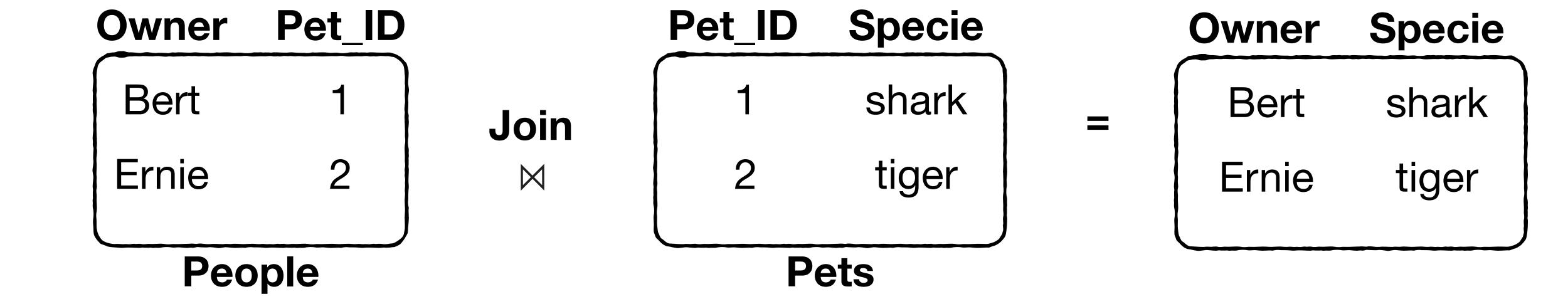
Join: Select ... FROM T1, T2 where **T1.B** = **T2**=**B**



Select Owner, Specie FROM People, Pets where T1.Pet_ID = T2=Pet_ID



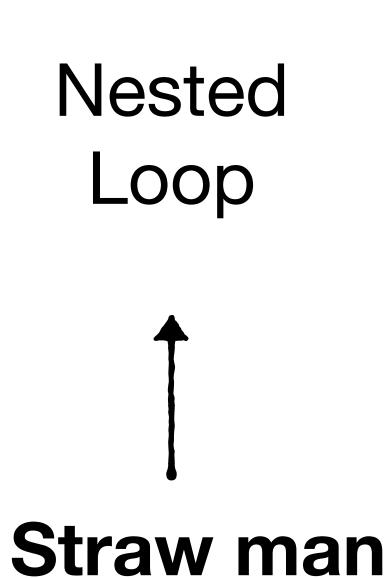
Select Owner, Specie FROM People, Pets where T1.Pet_ID = T2=Pet_ID



Nested Loop Block Nested Loop

Index-Join

Sort-merge Join Grace hash Join

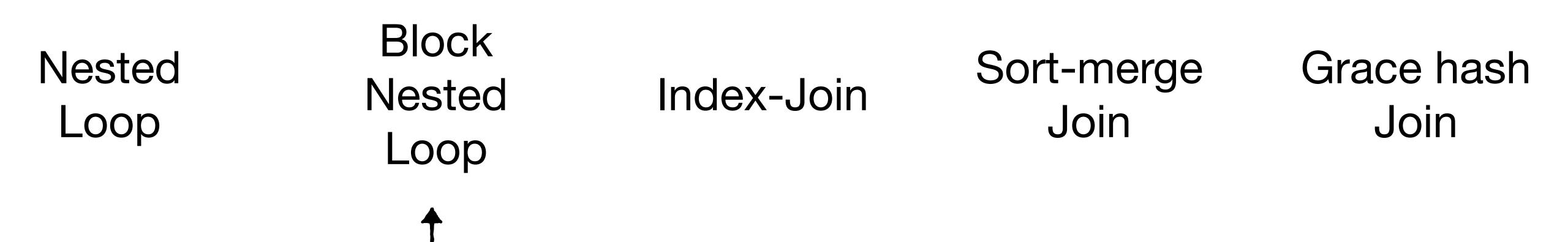


Block Nested Loop

Index-Join

Sort-merge Join

Grace hash Join



Good when one relation almost or entirely fits in memory



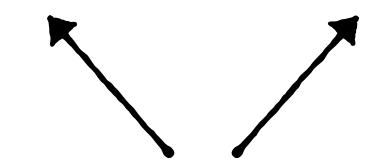
There are index/es on the join keys

Nested Loop Block Nested Loop

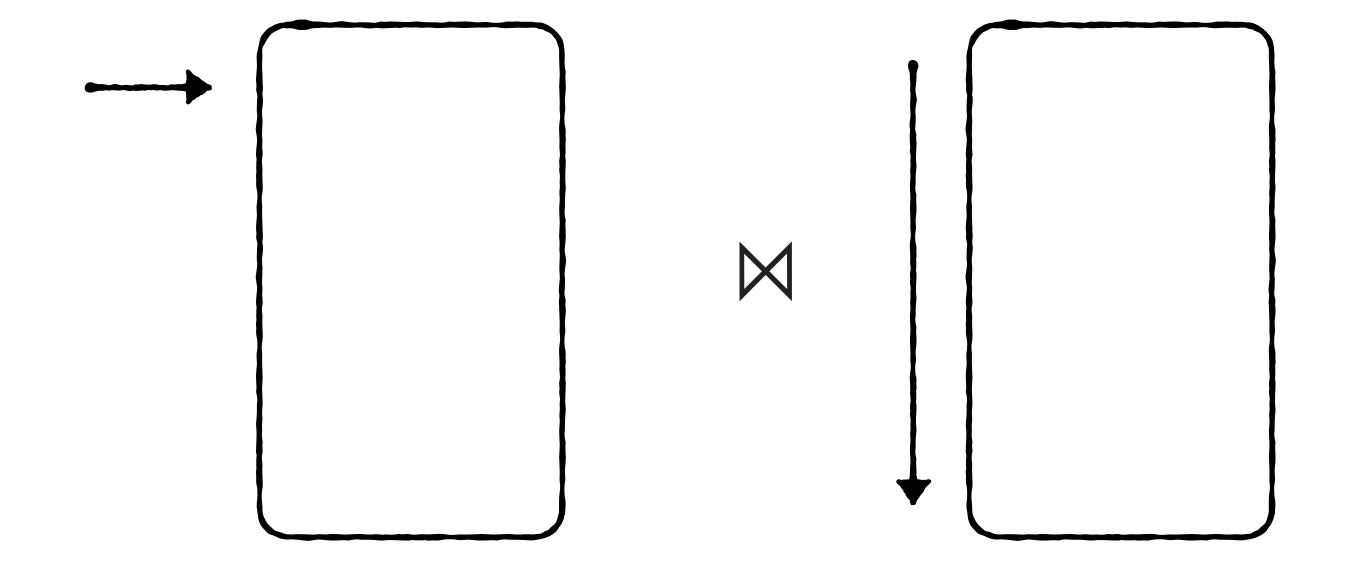
Index-Join

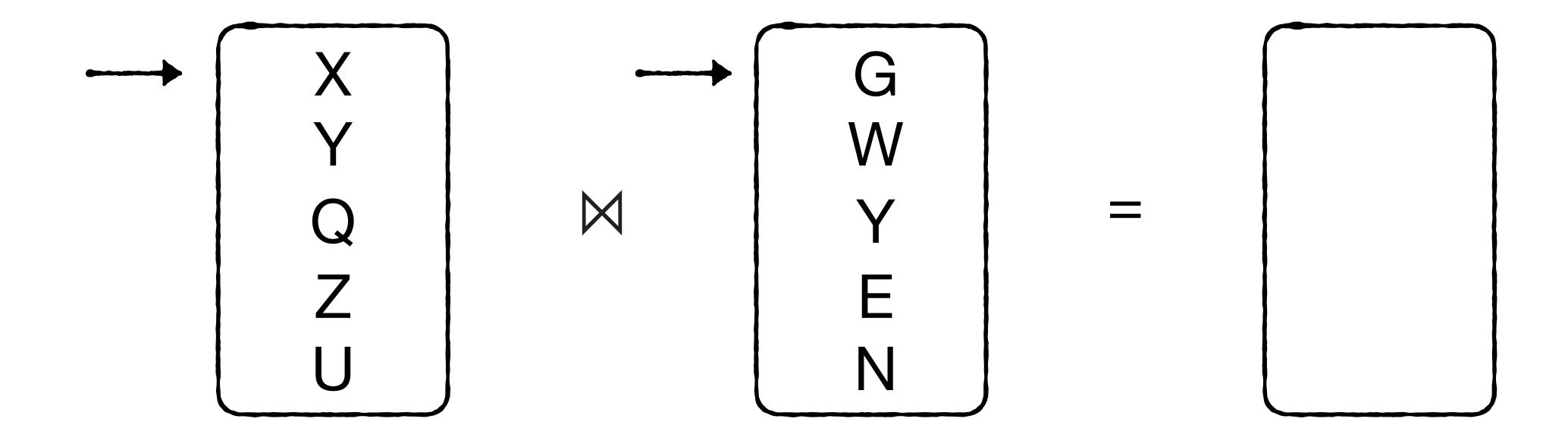
Sort-merge Join

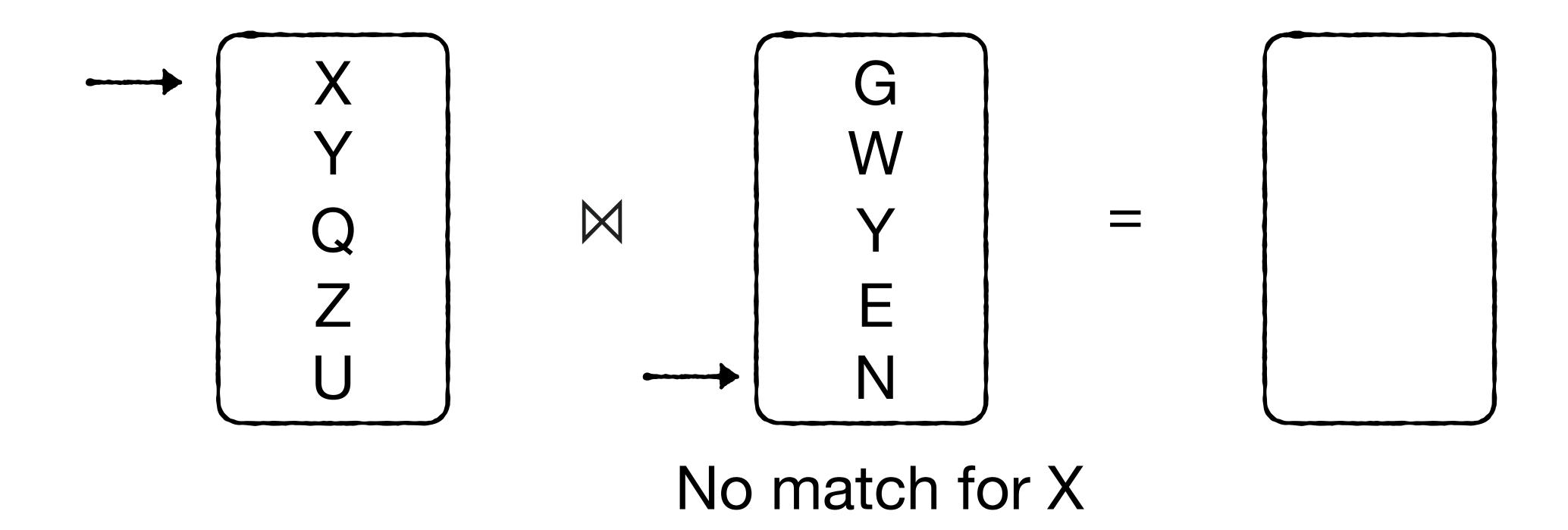
Grace hash Join

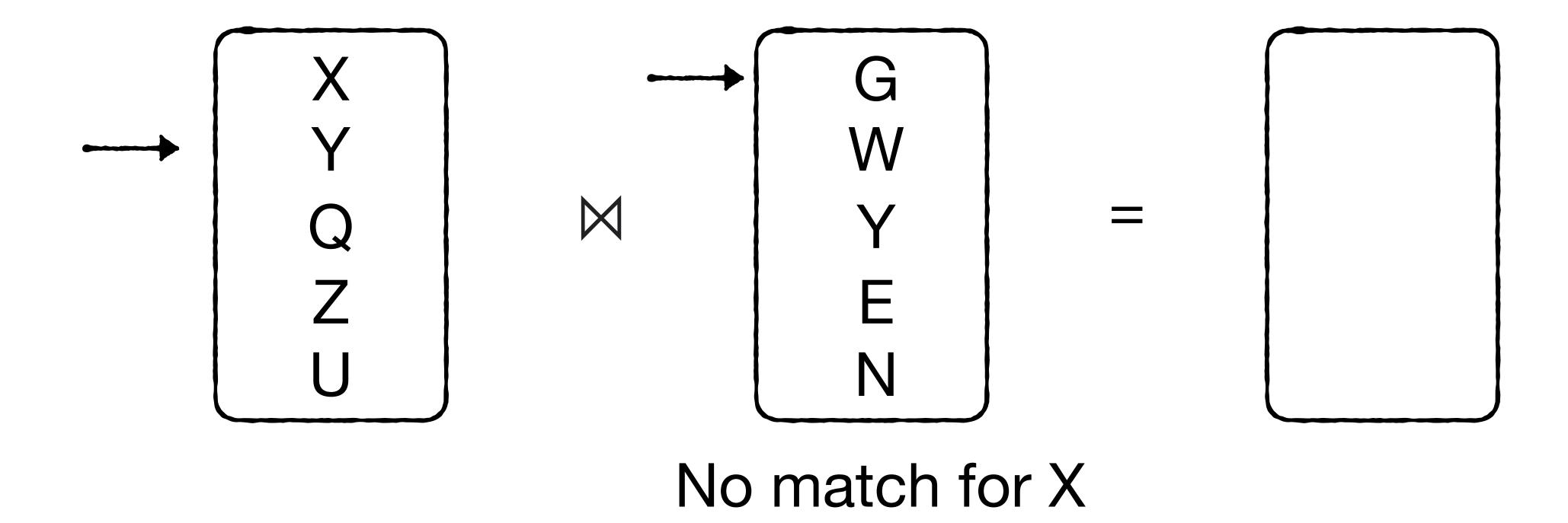


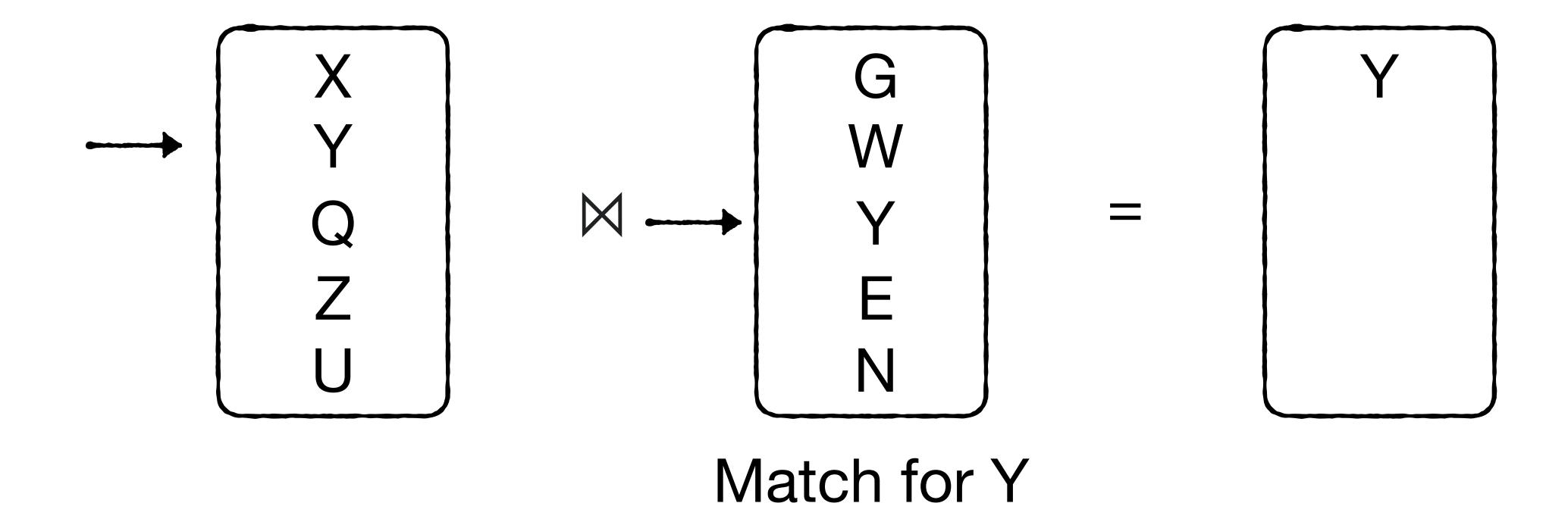
Both relations must bigger than memory

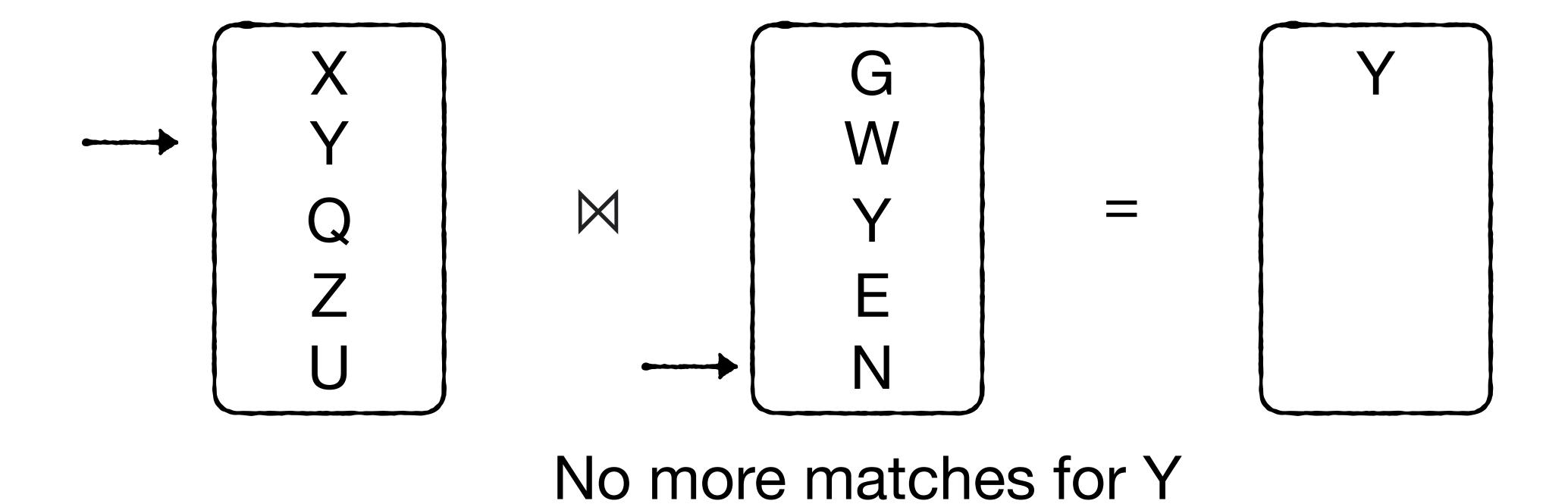


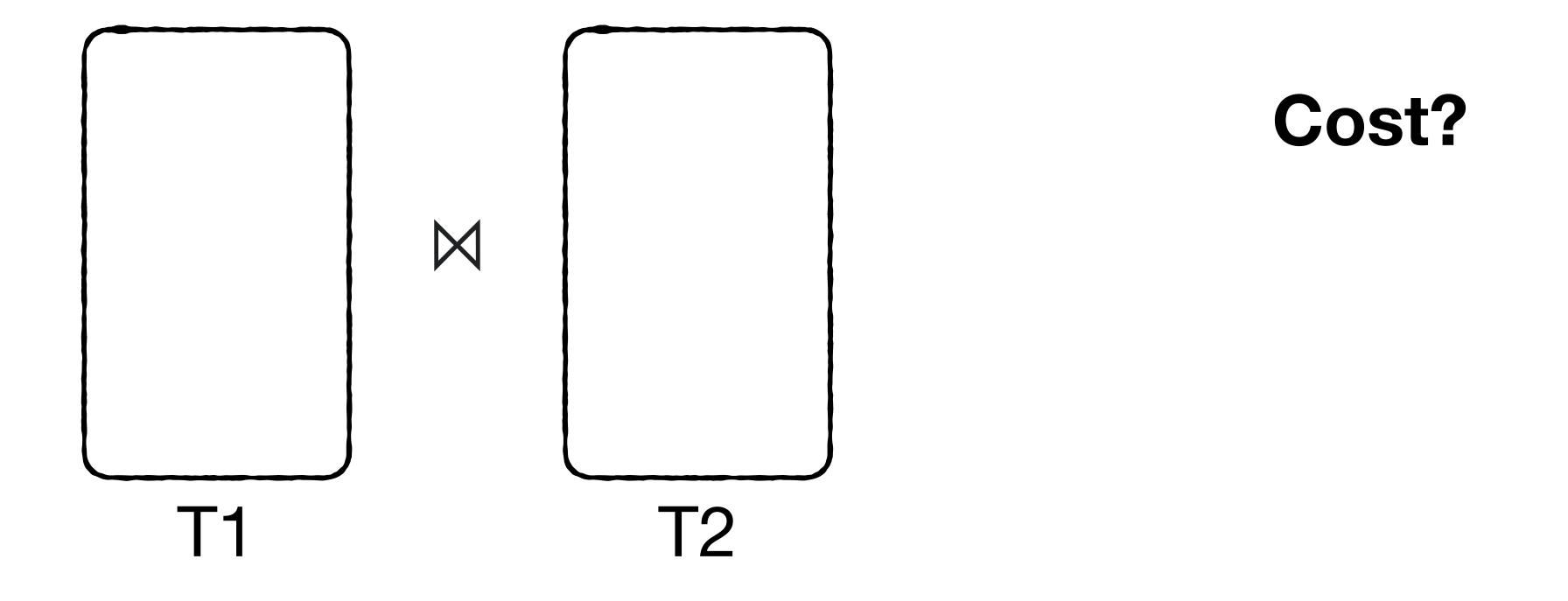


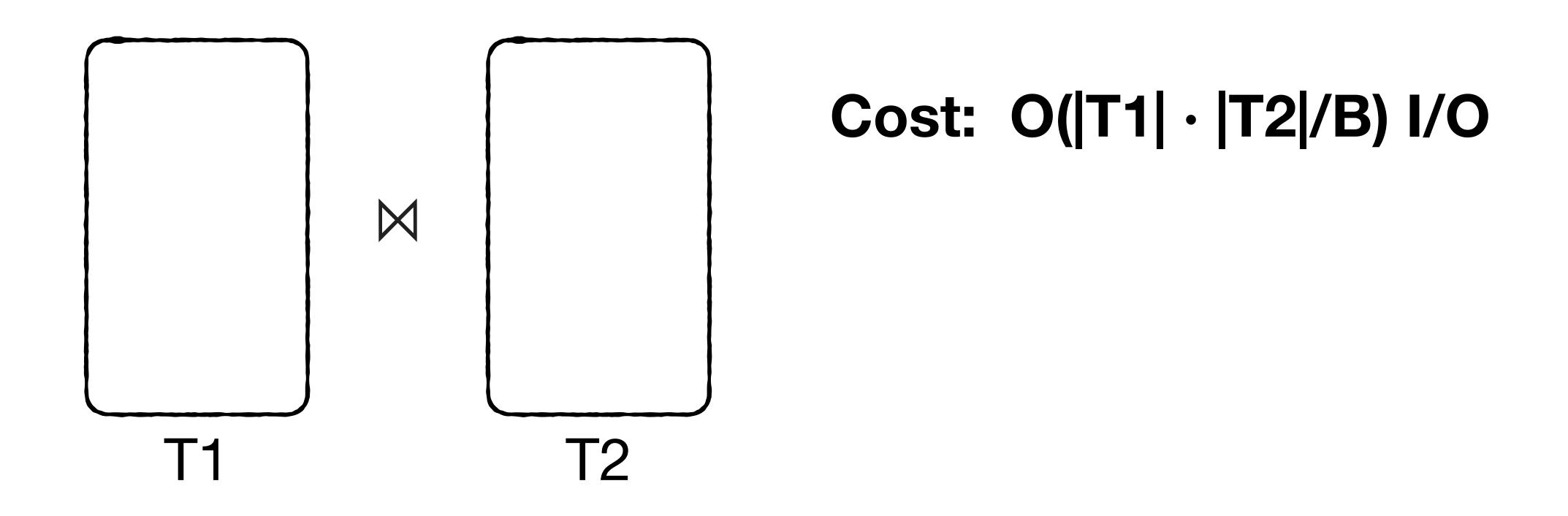


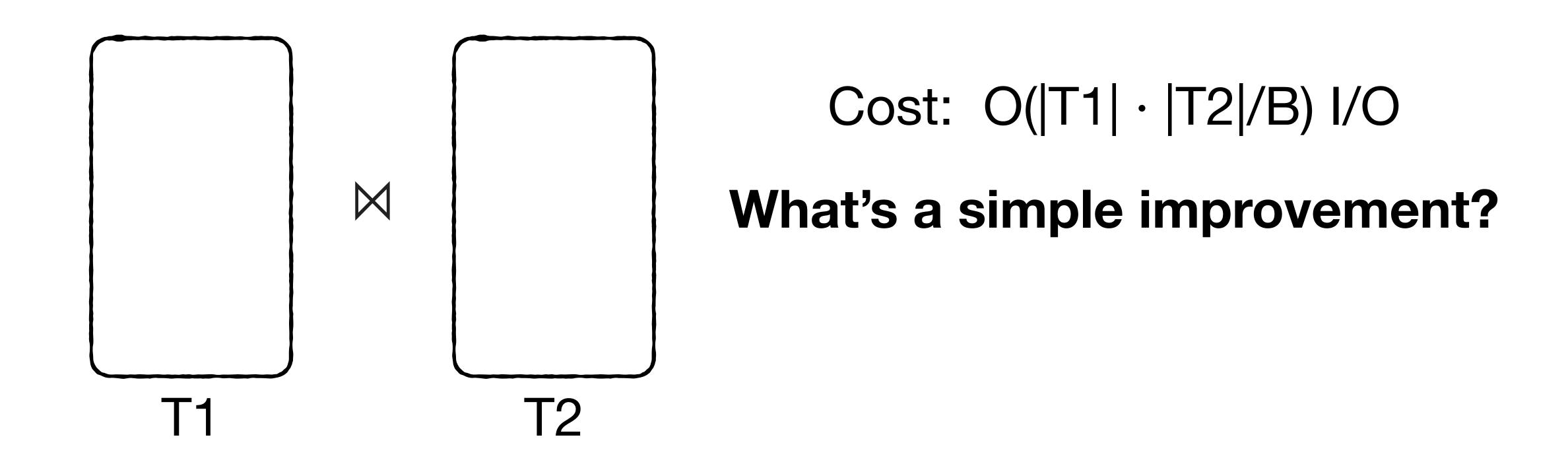


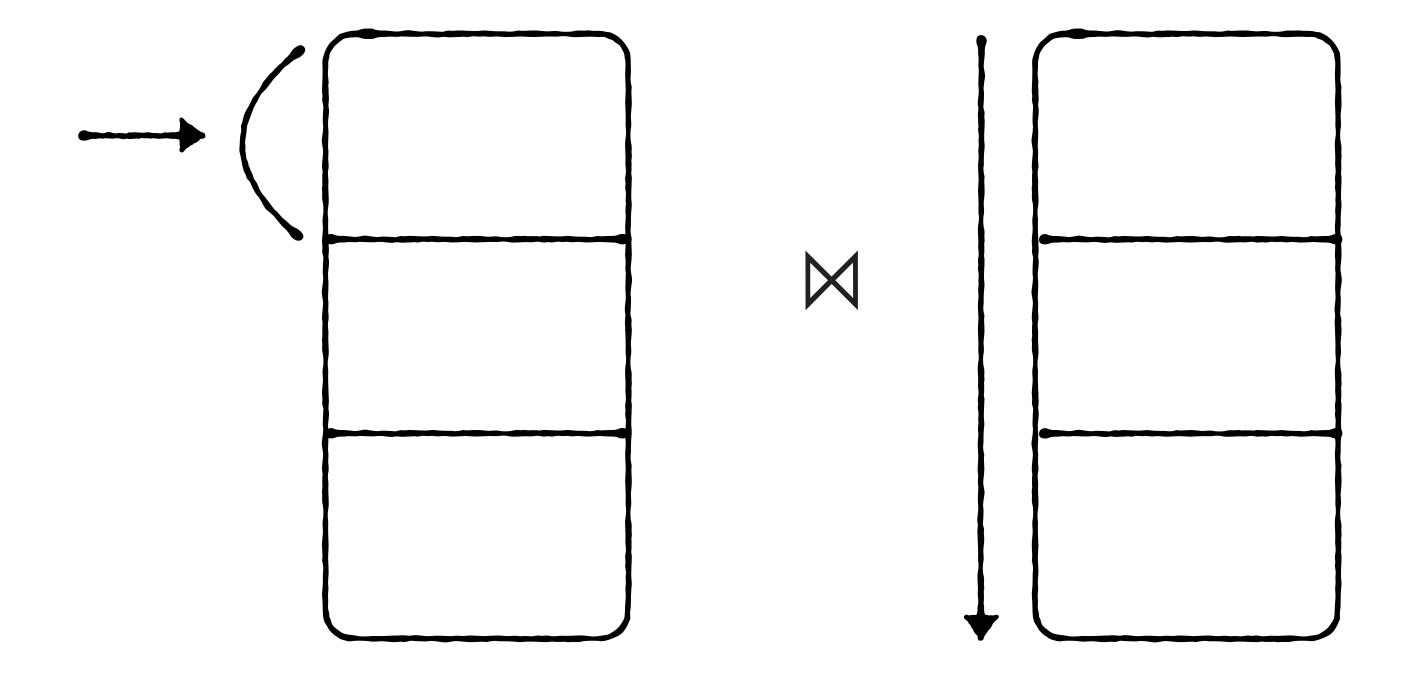


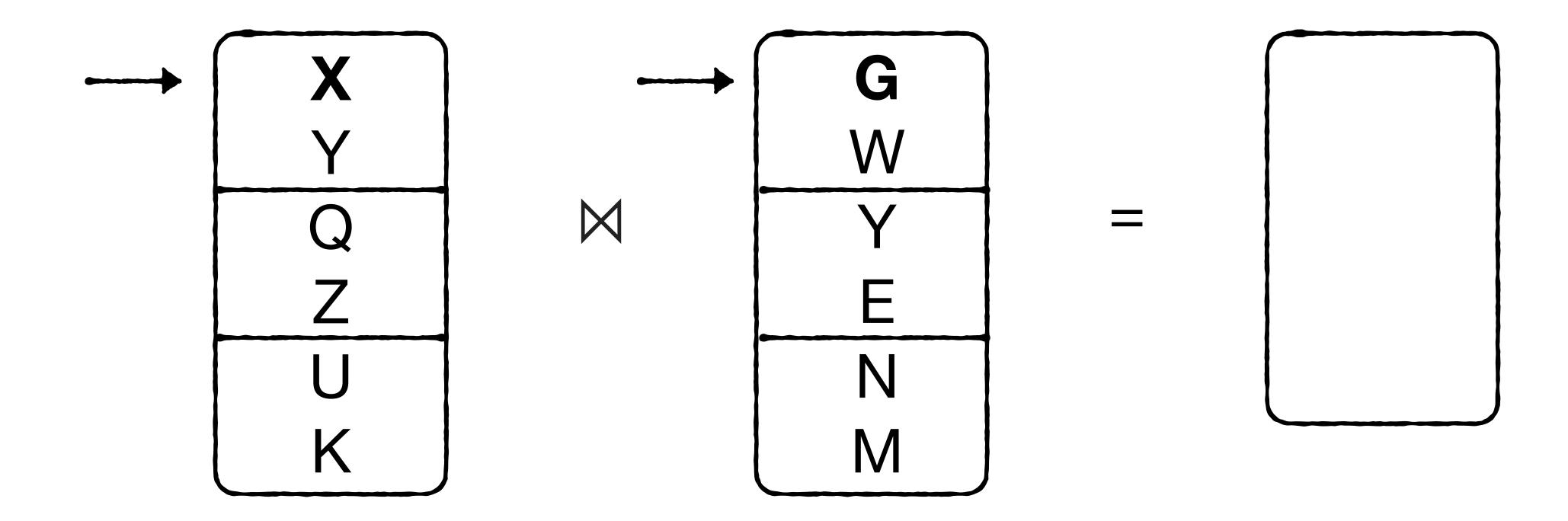


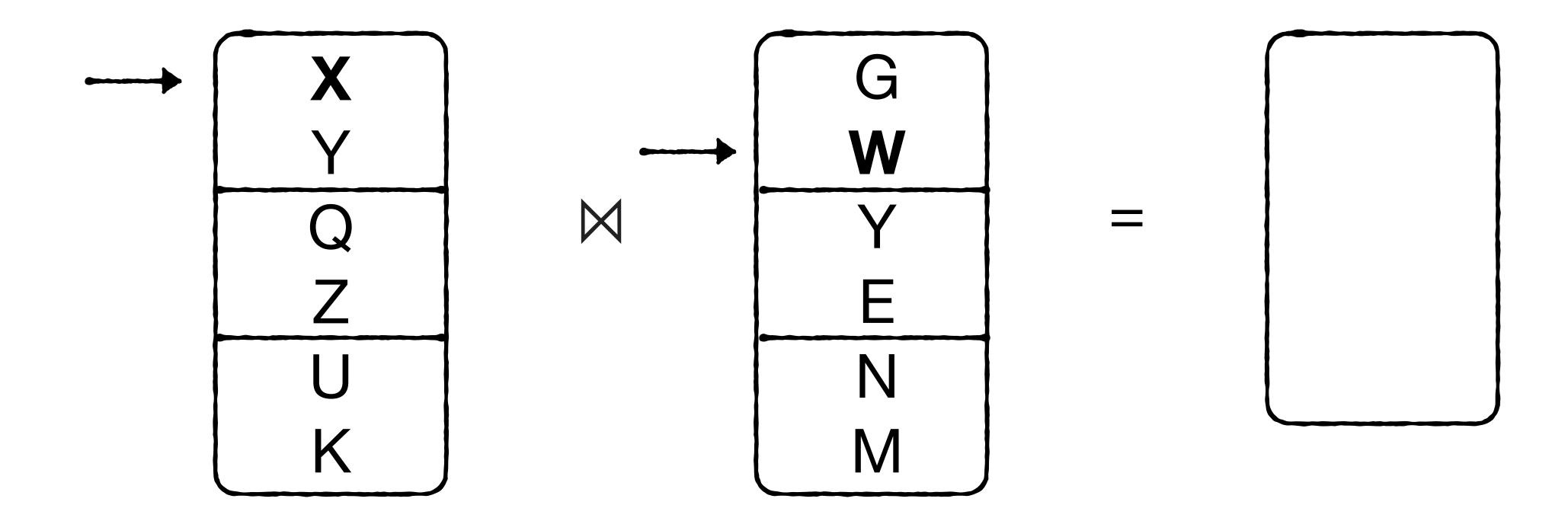


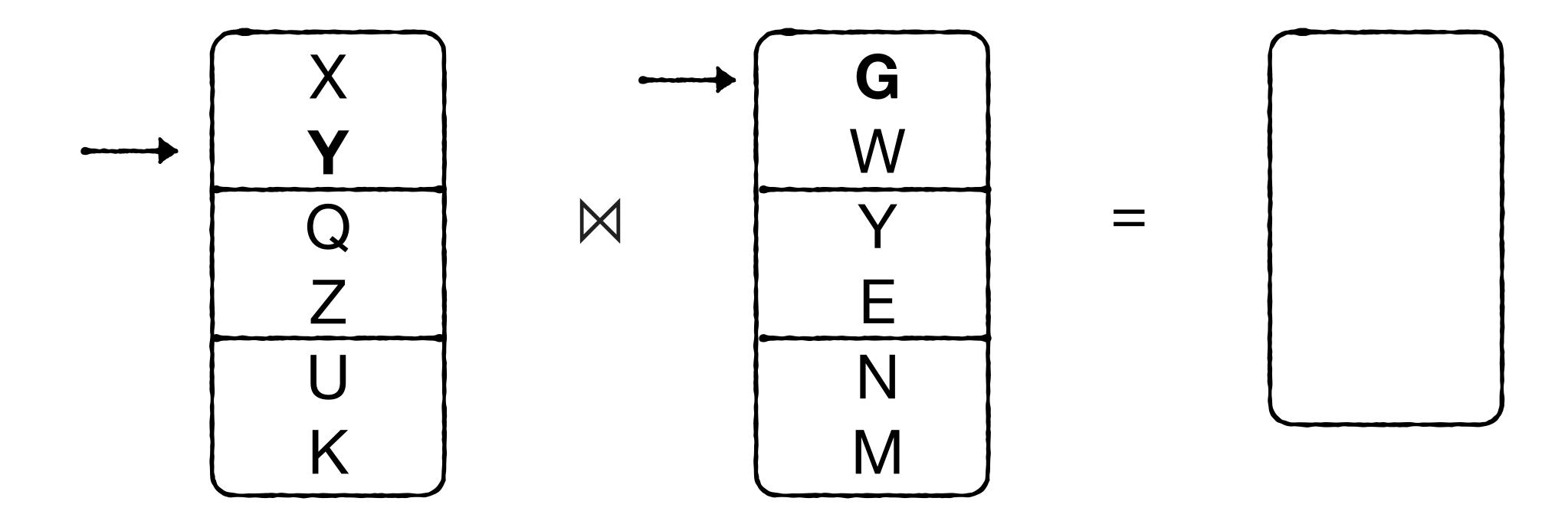


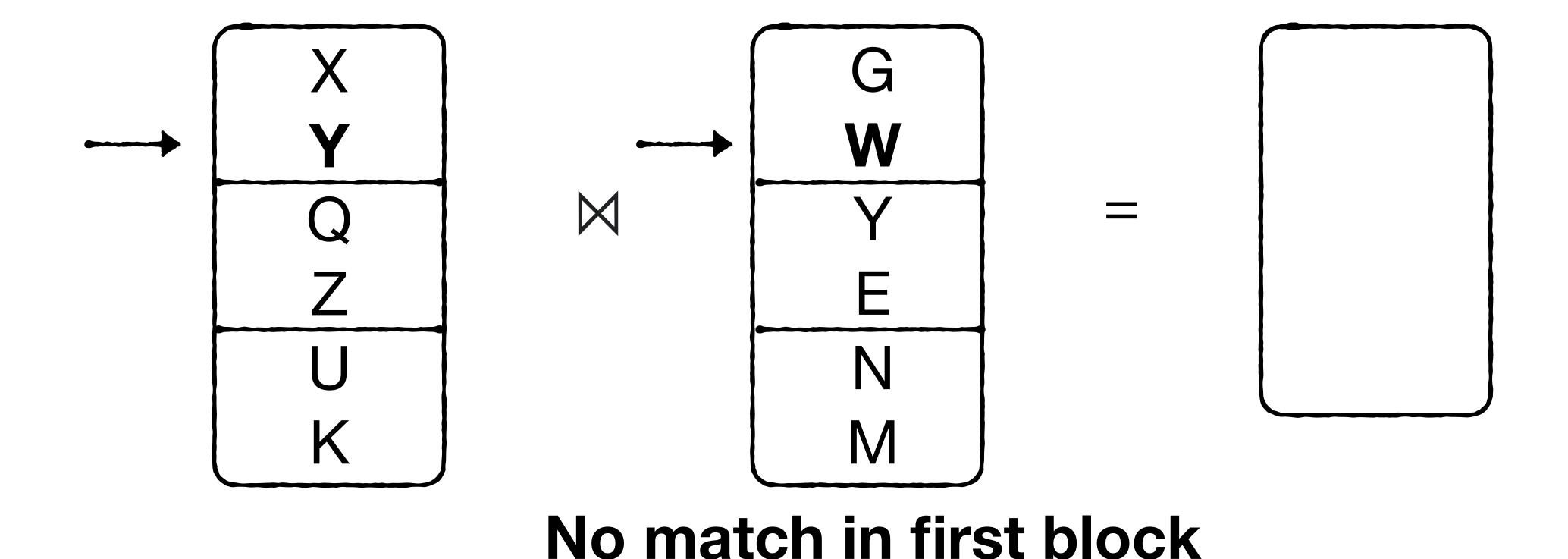


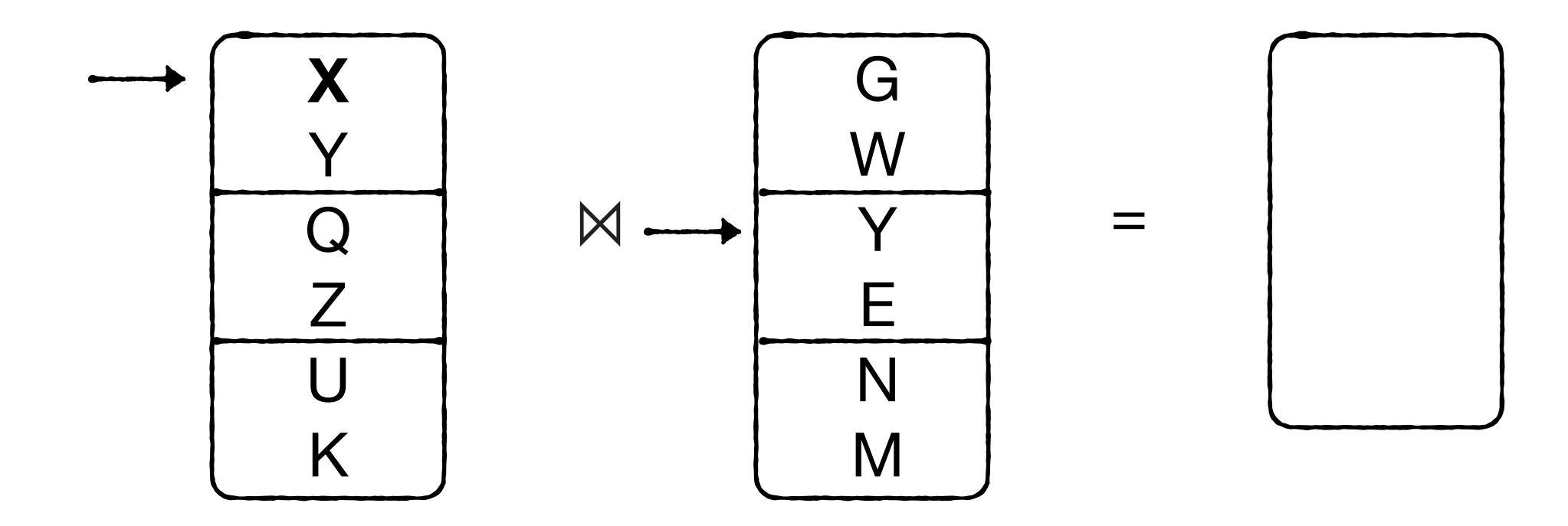


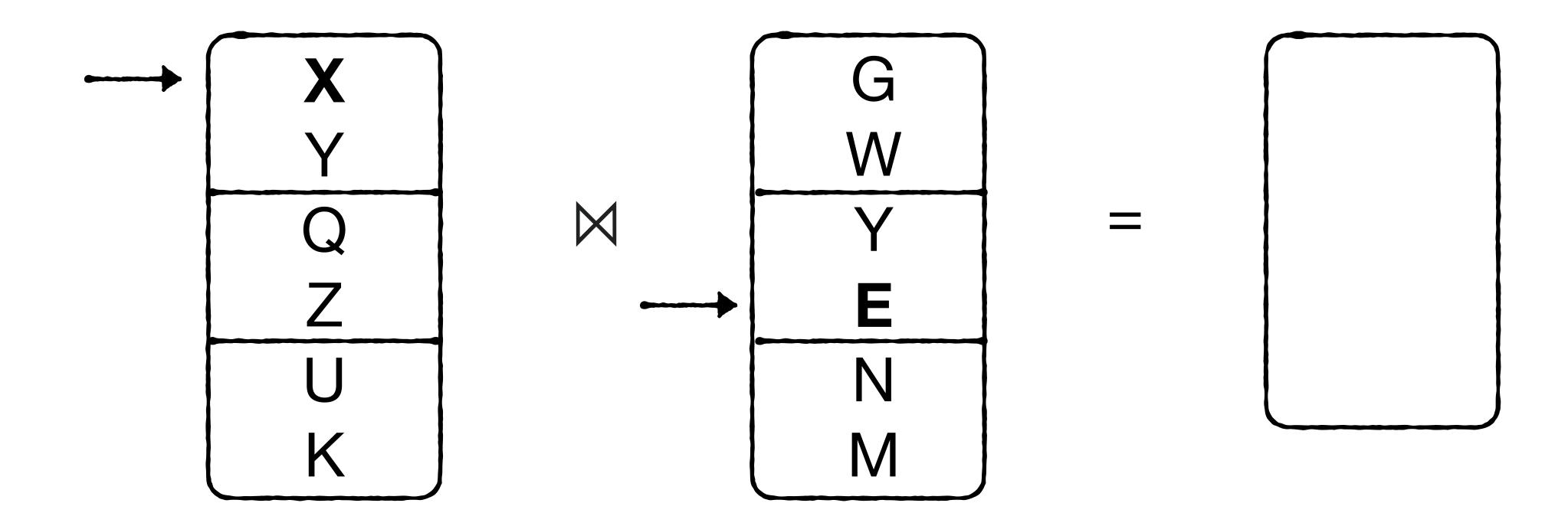


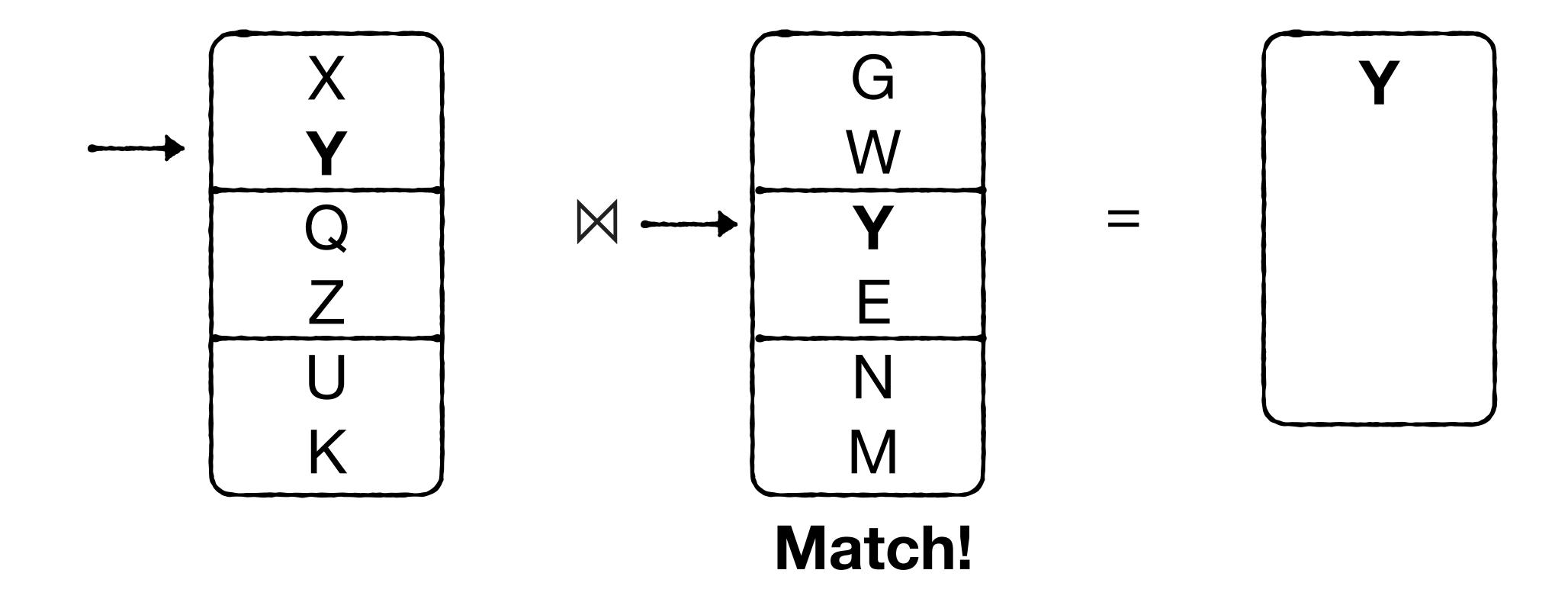


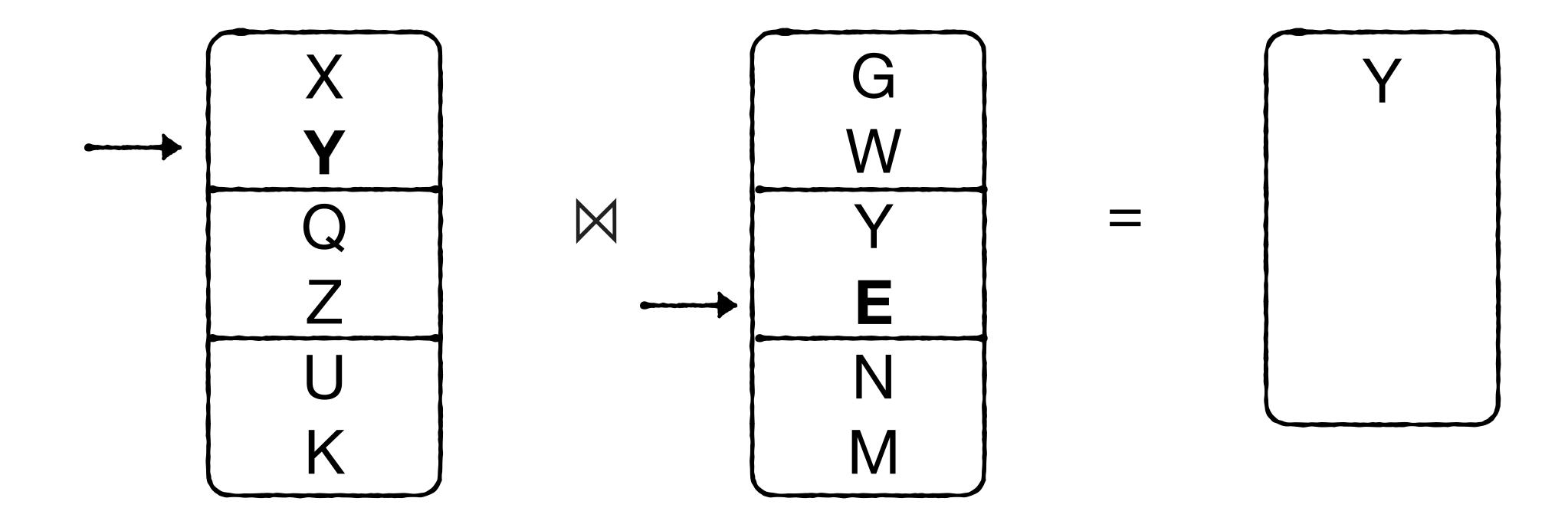


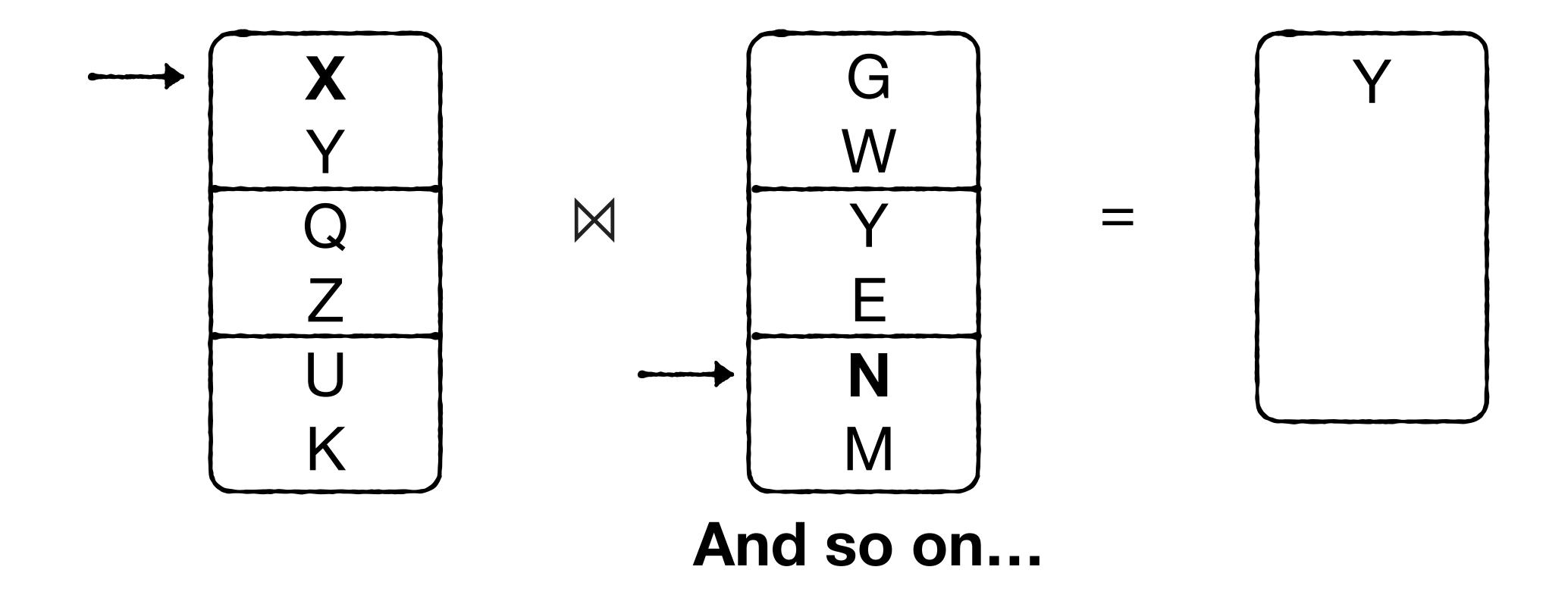




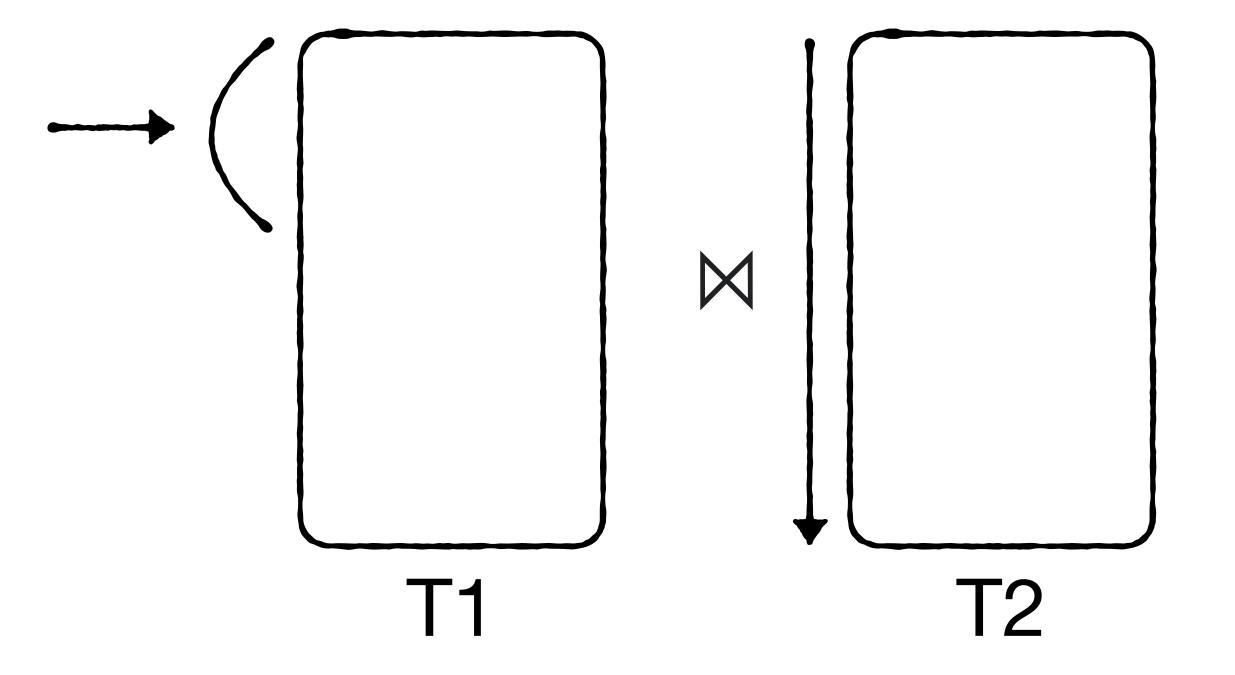




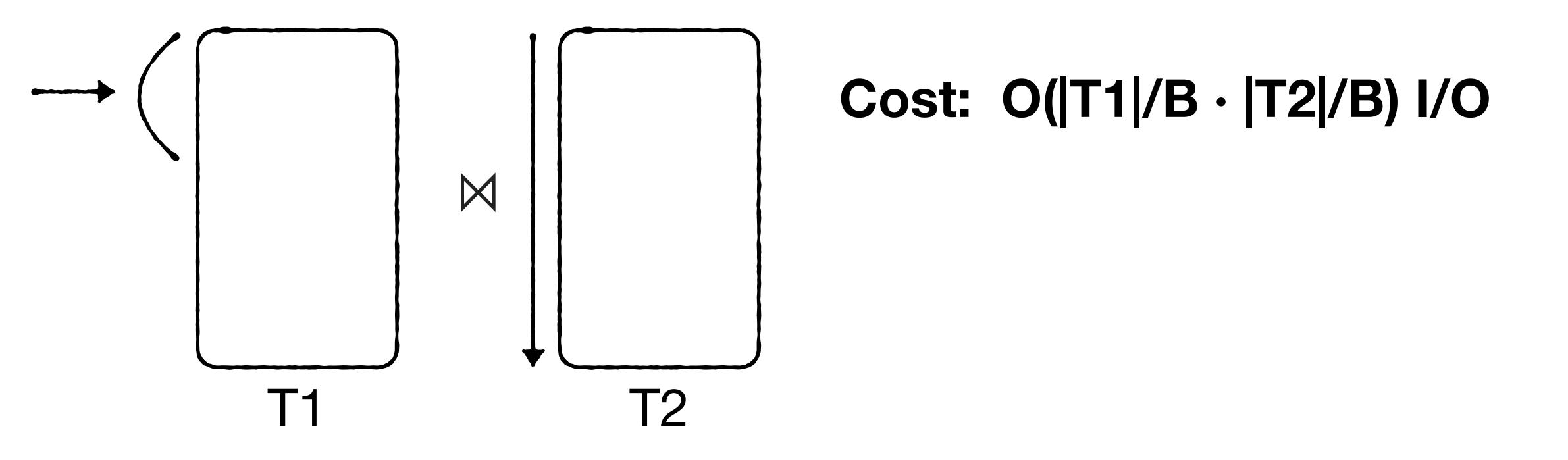


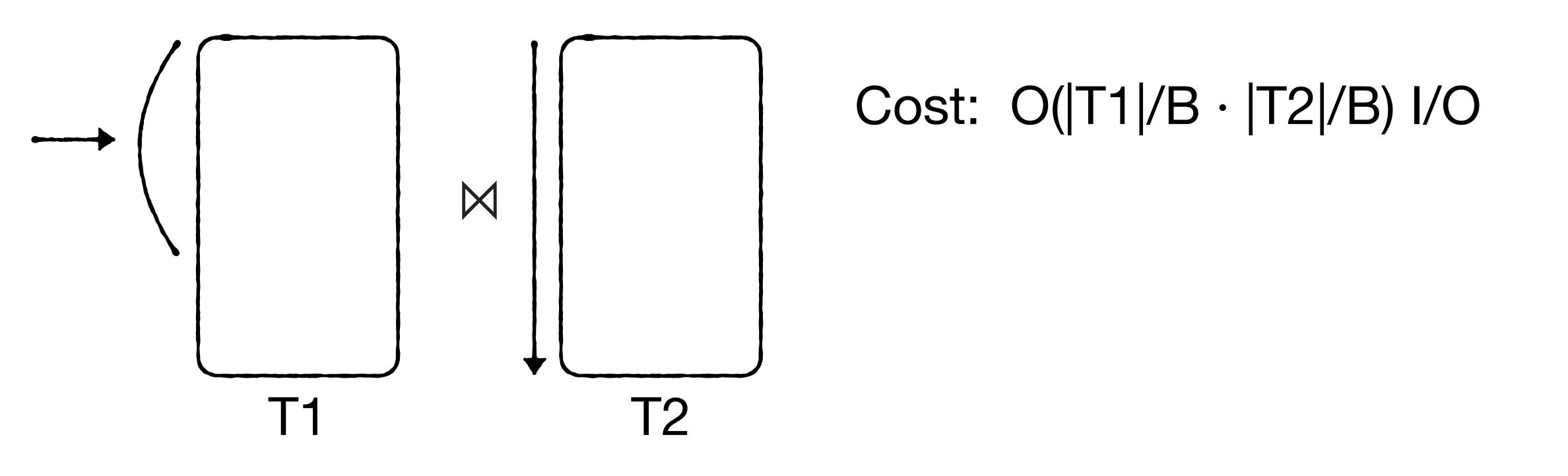


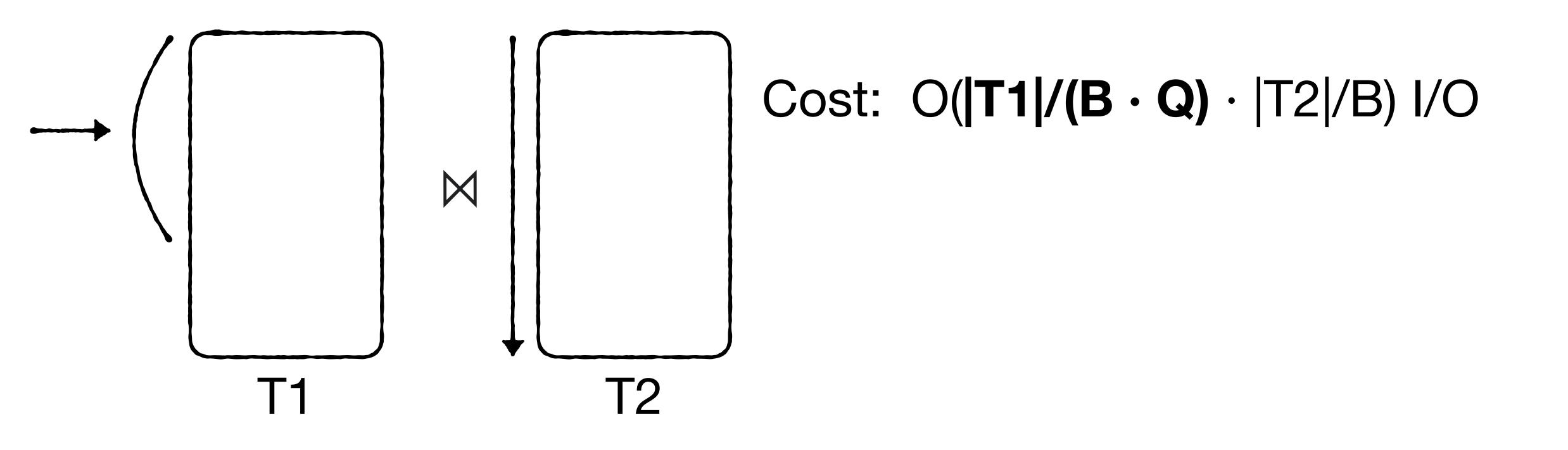
For each block in one relation, scan whole other relation

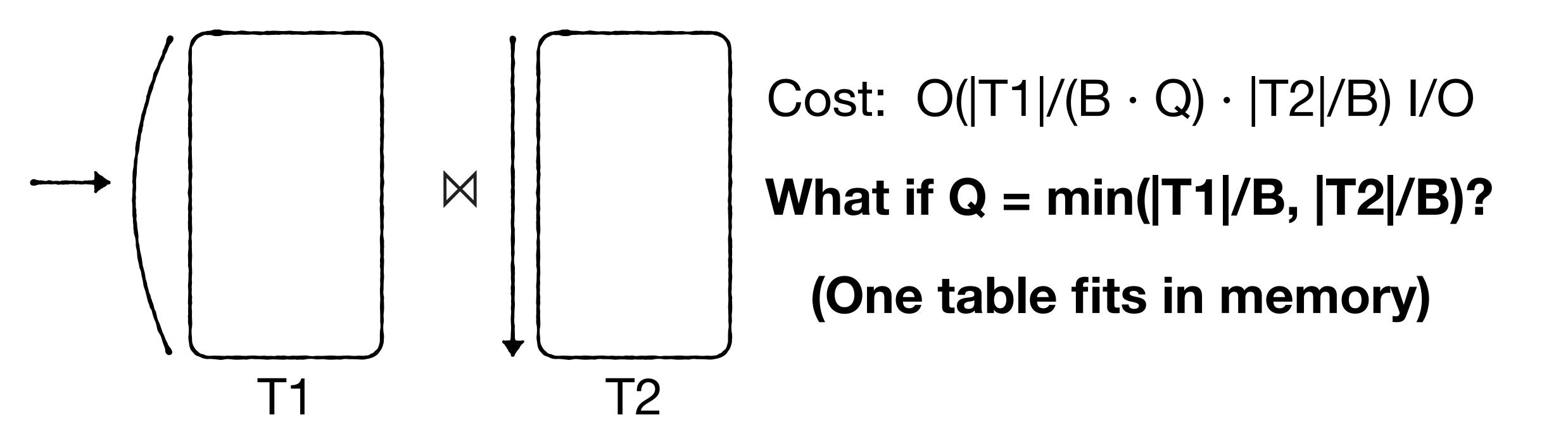


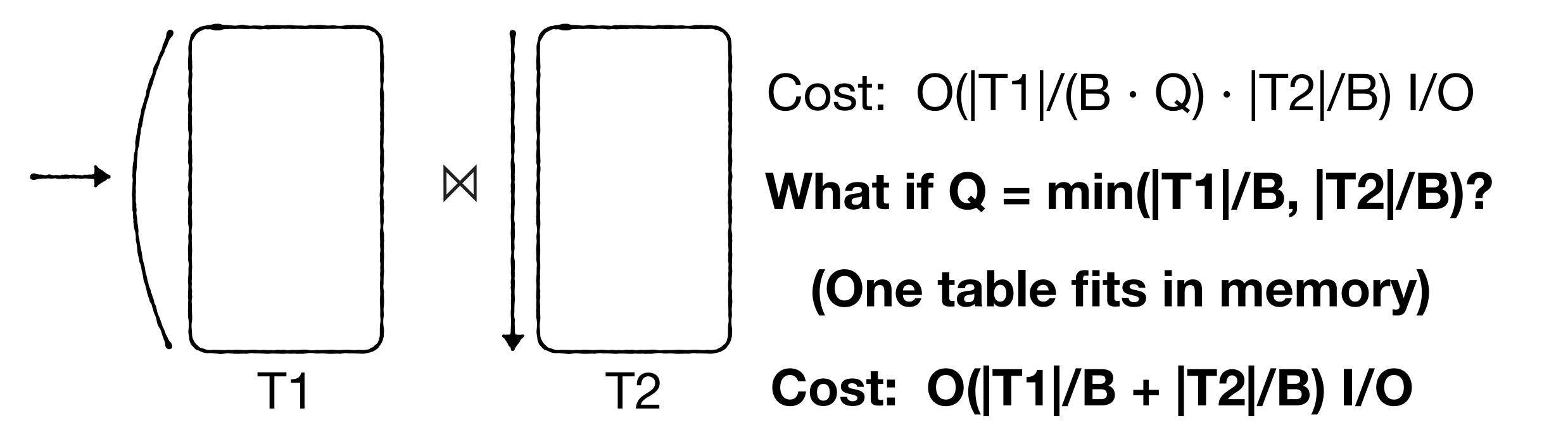
Cost?











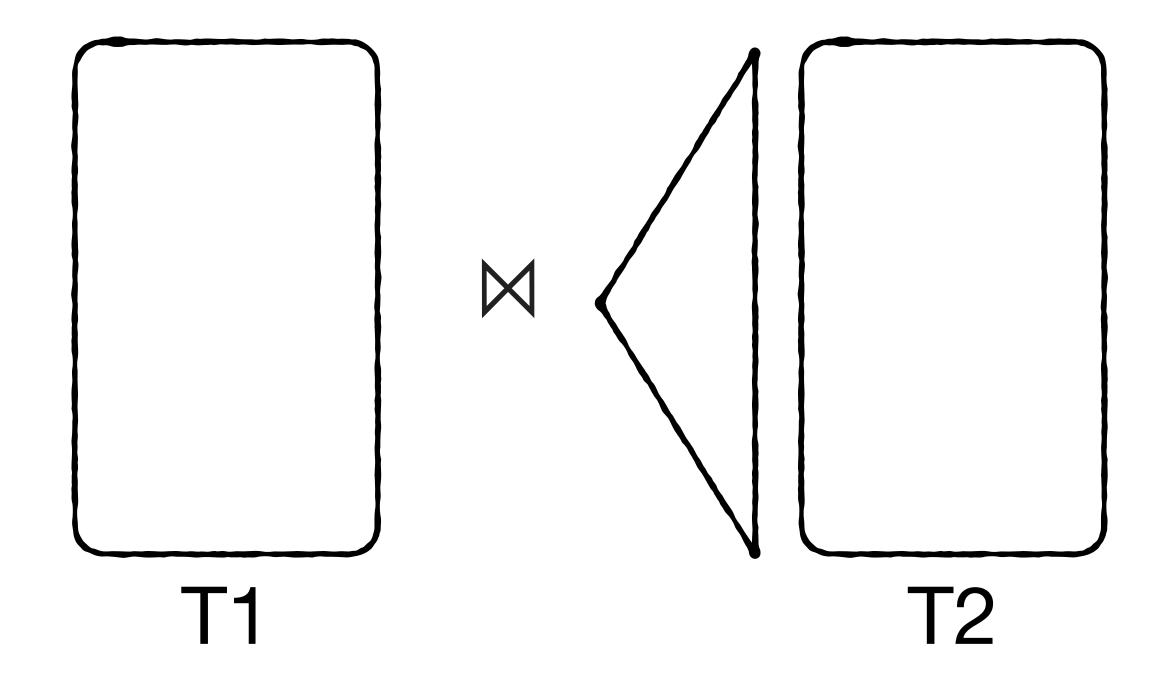
Join Algorithms

Nested Loop Block Nested Loop

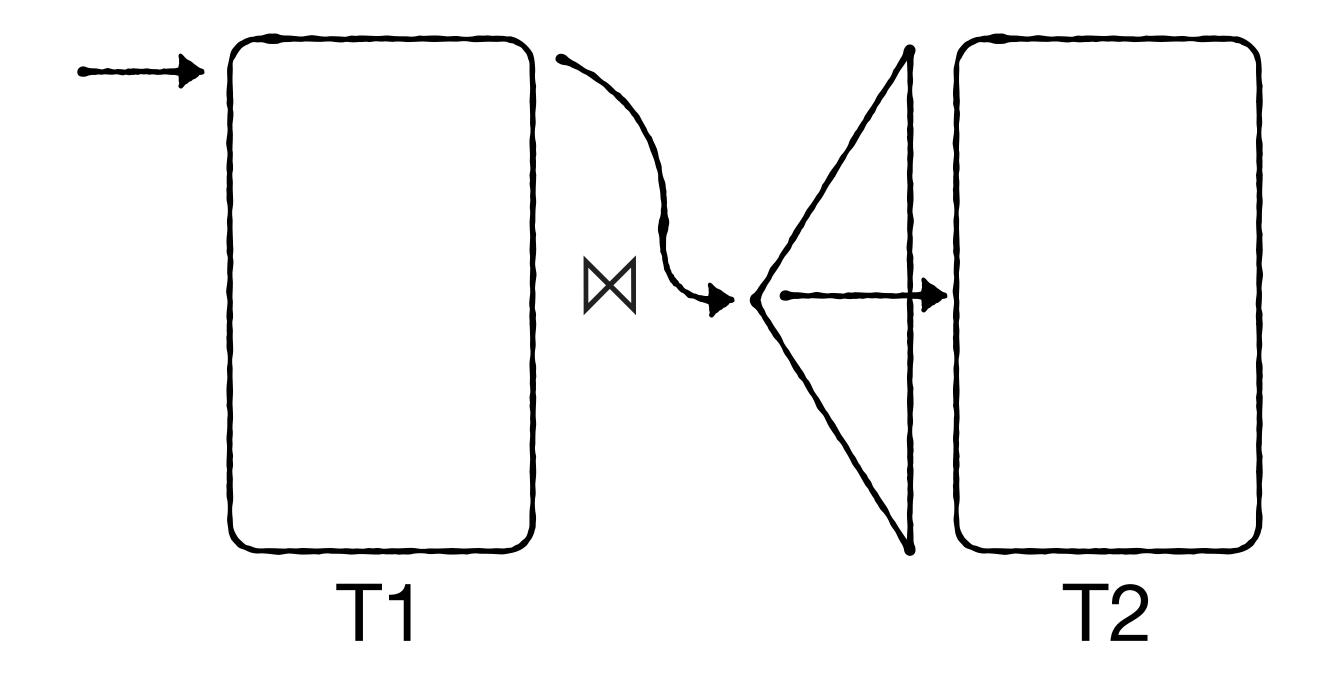
Index-Join

Sort-merge Join Grace hash Join

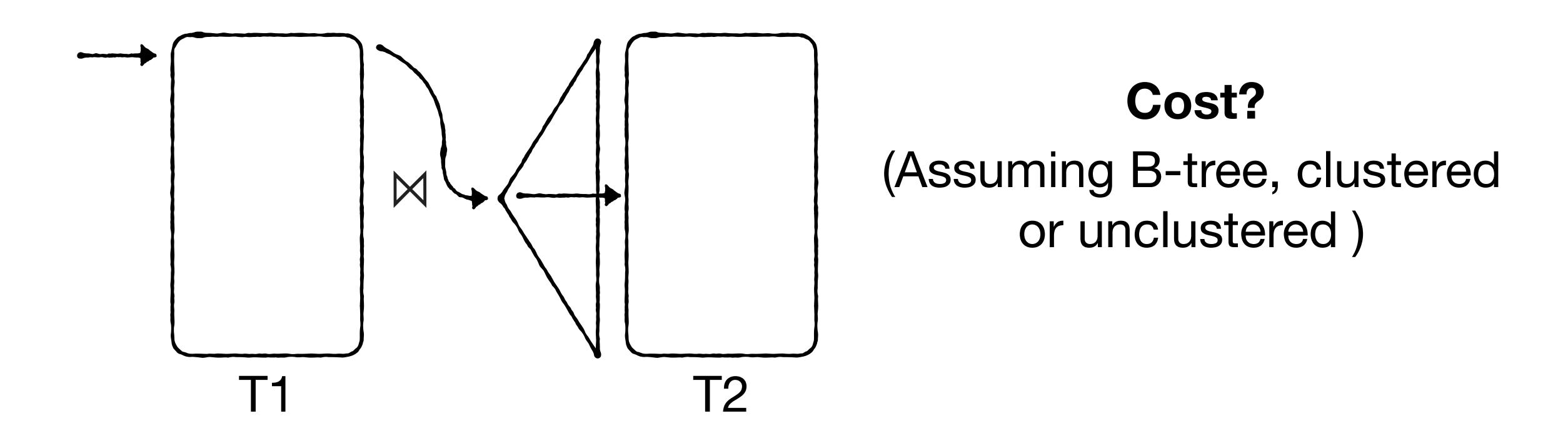
Suppose we have index on one of the relations



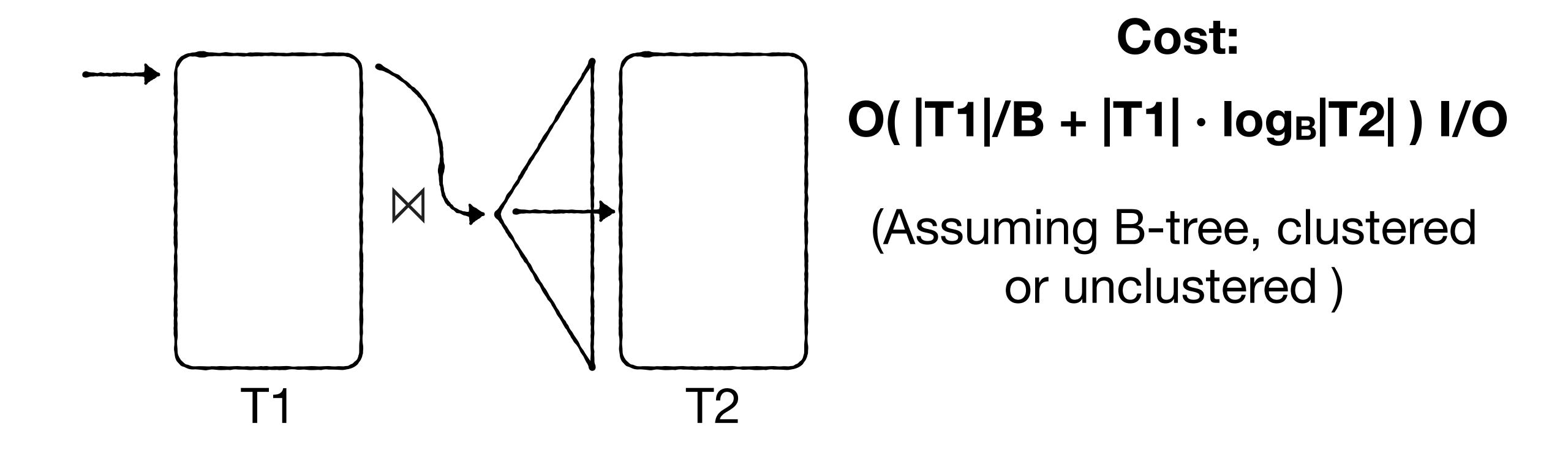
For each entry in one relation, search index for other relation



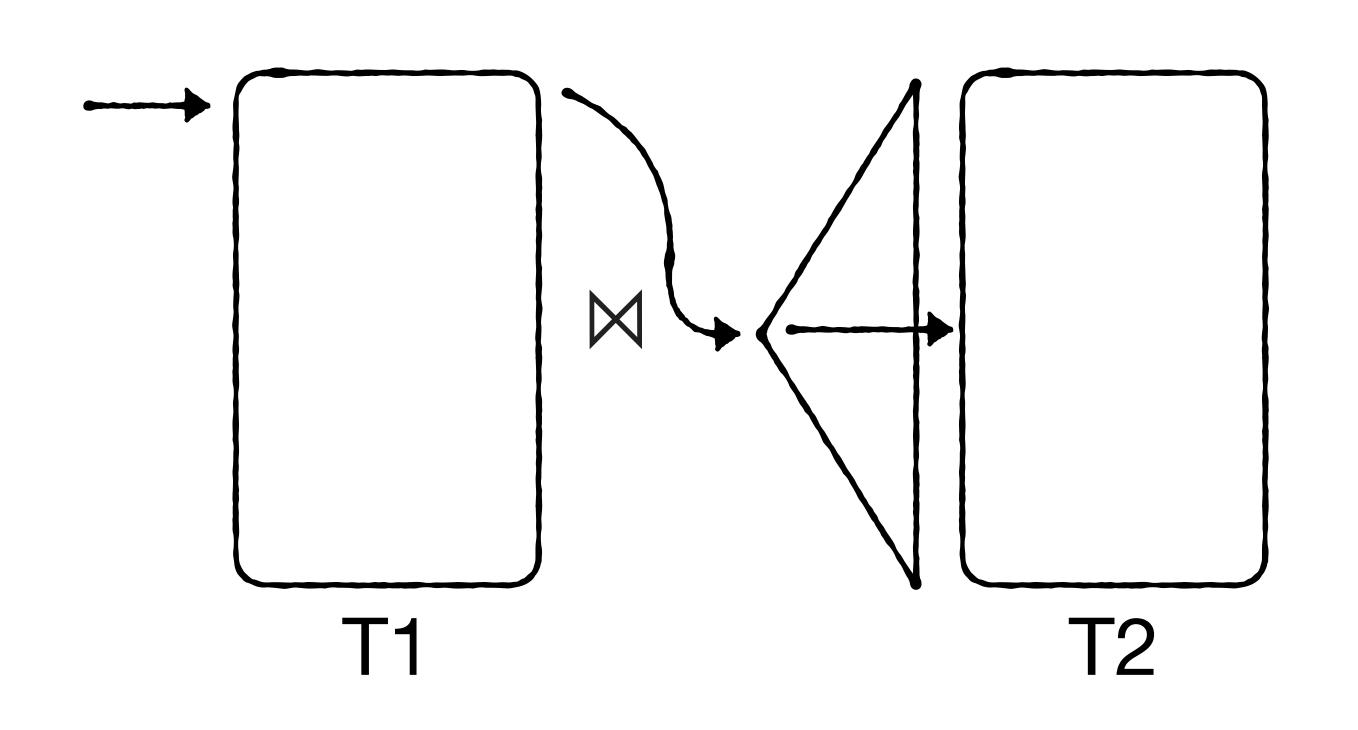
For each entry in one relation, search index for other relation



For each entry in one relation, search index for other relation



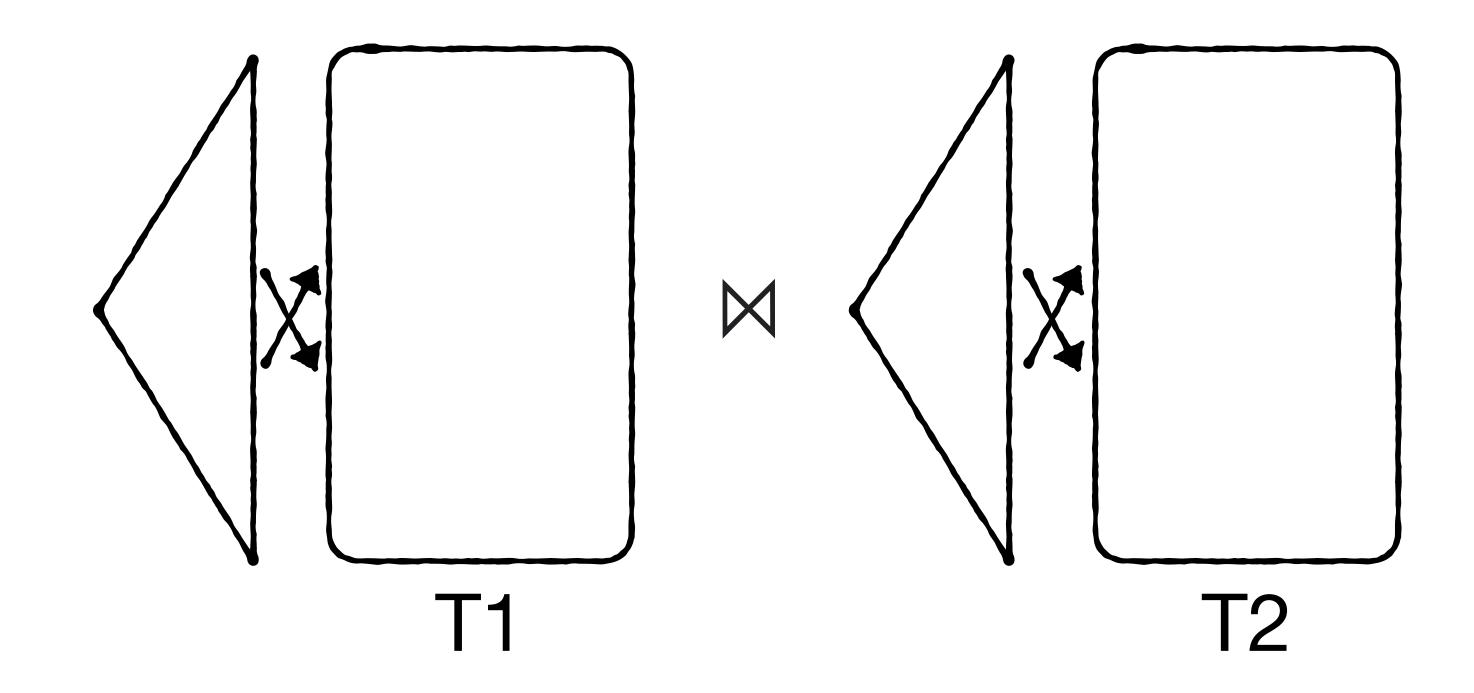
For each entry in one relation, search index for other relation

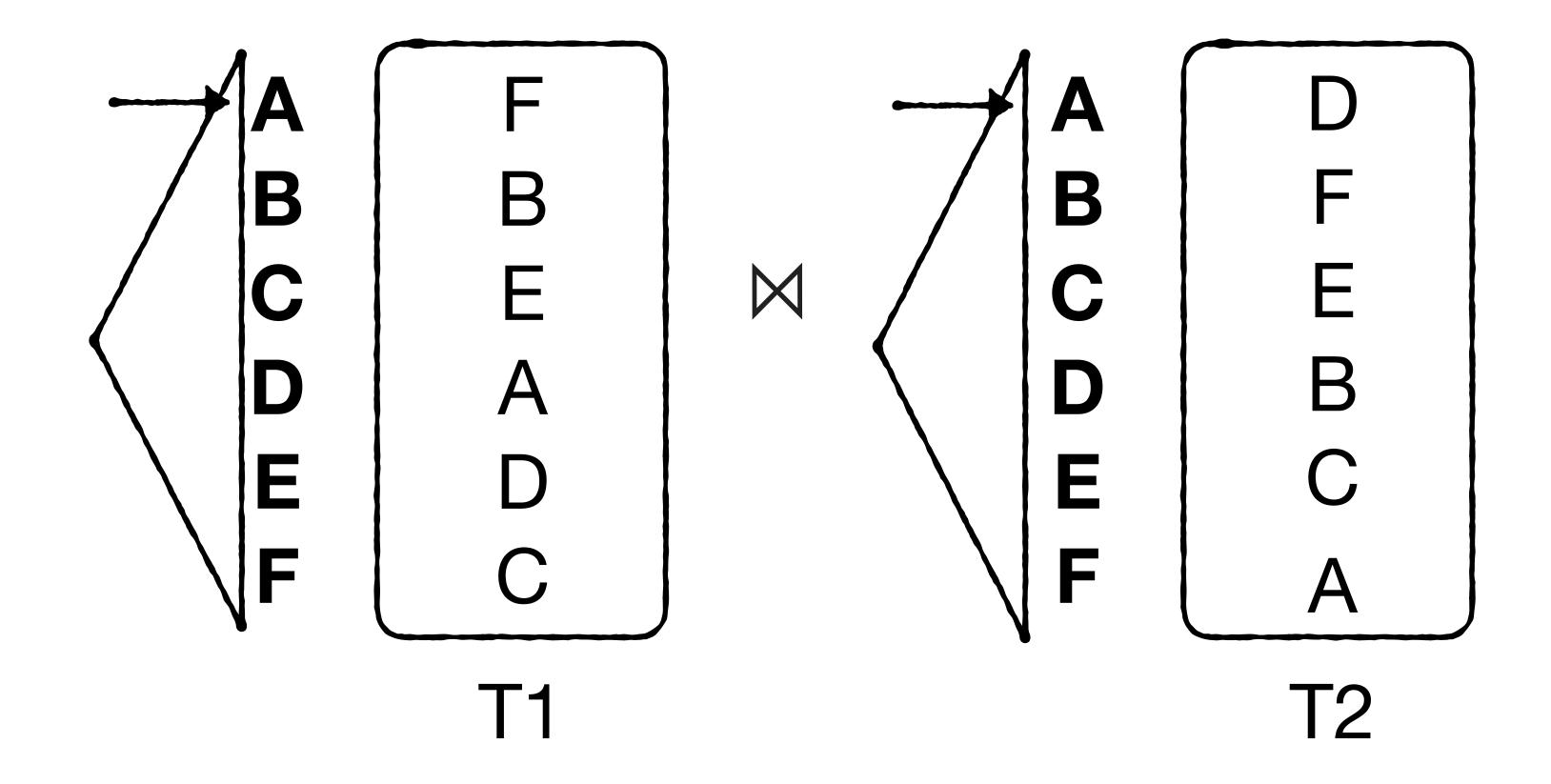


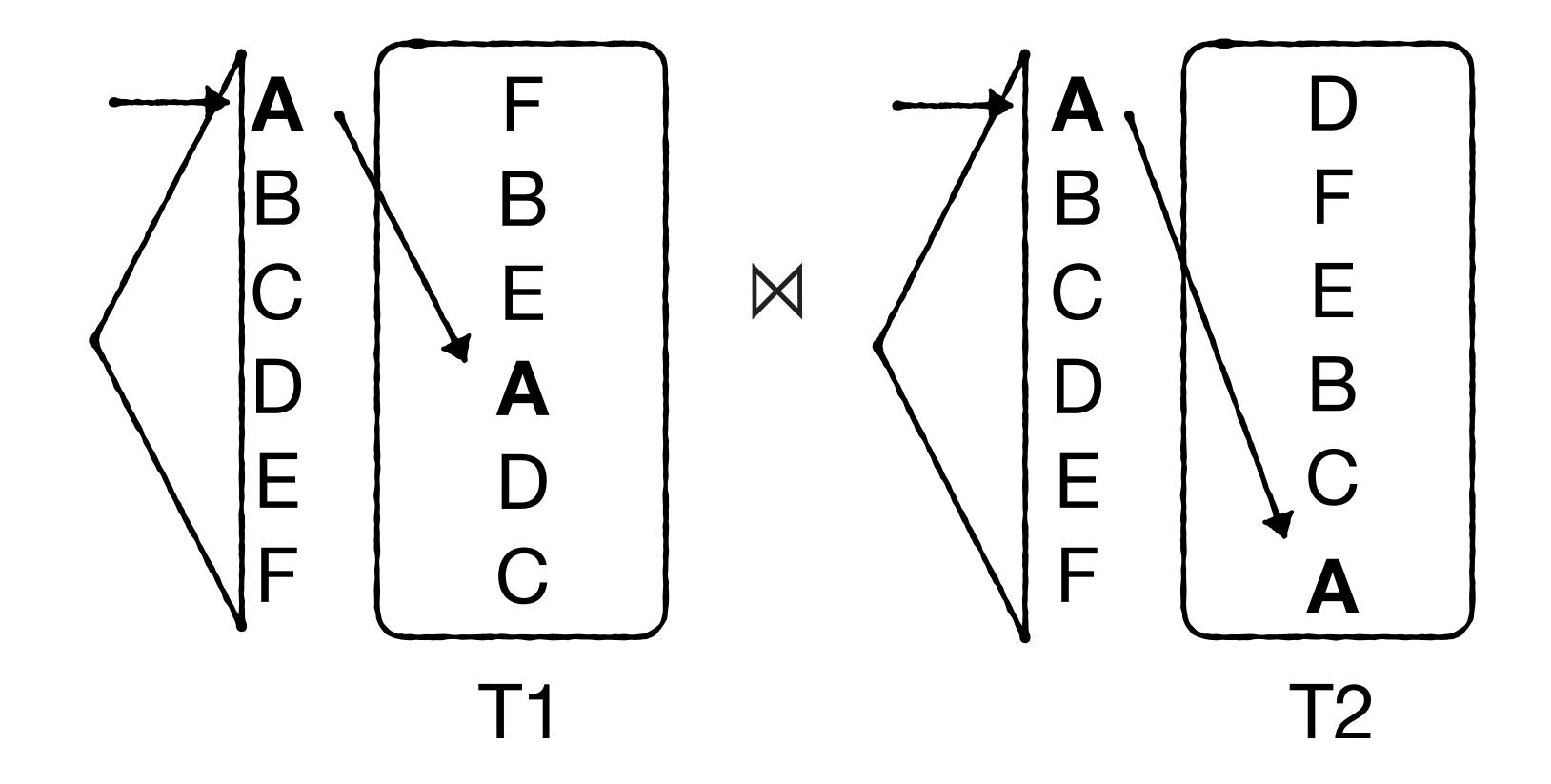
Cost:

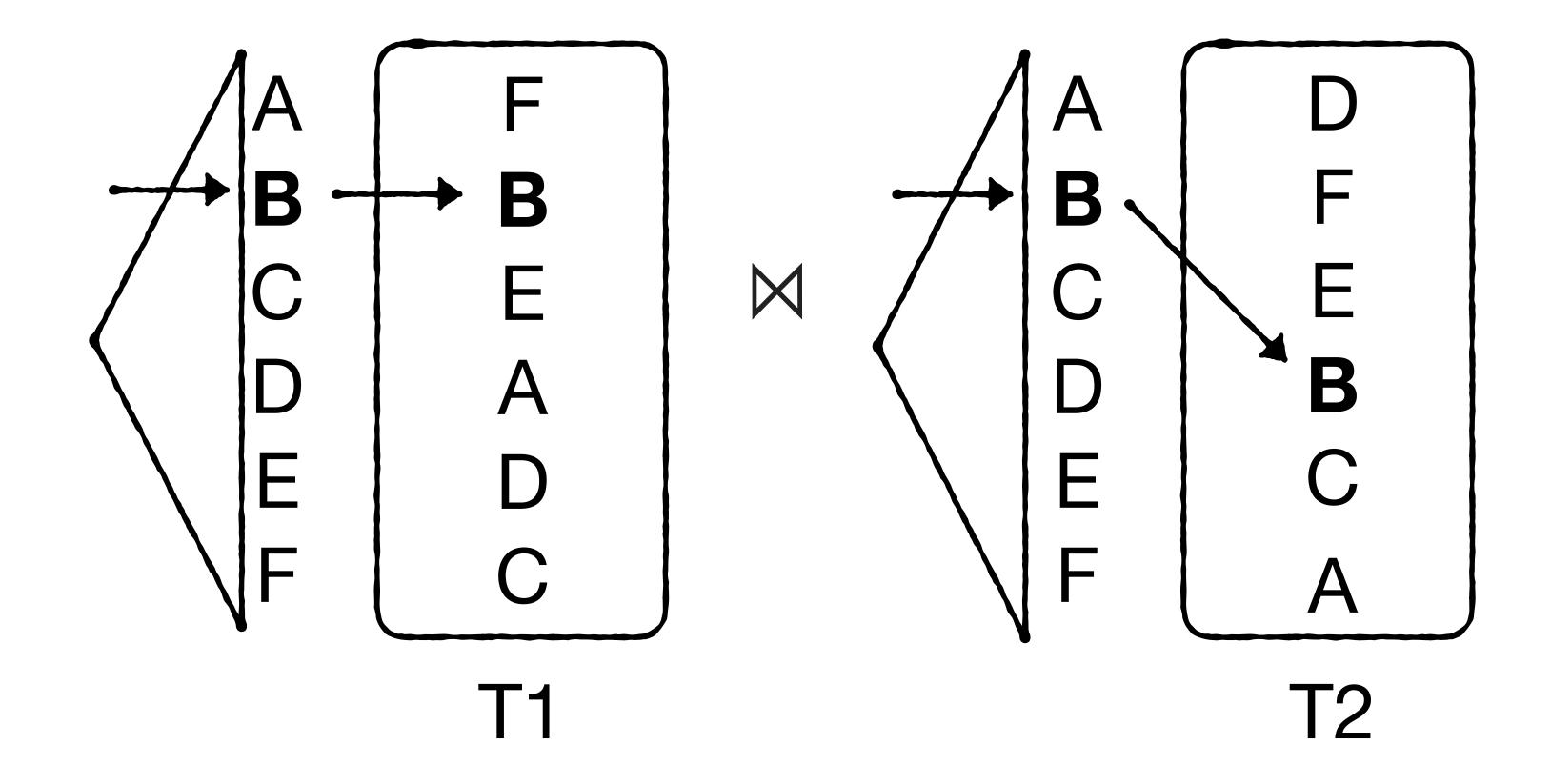
O(|T1| · log_B|T2|) I/O

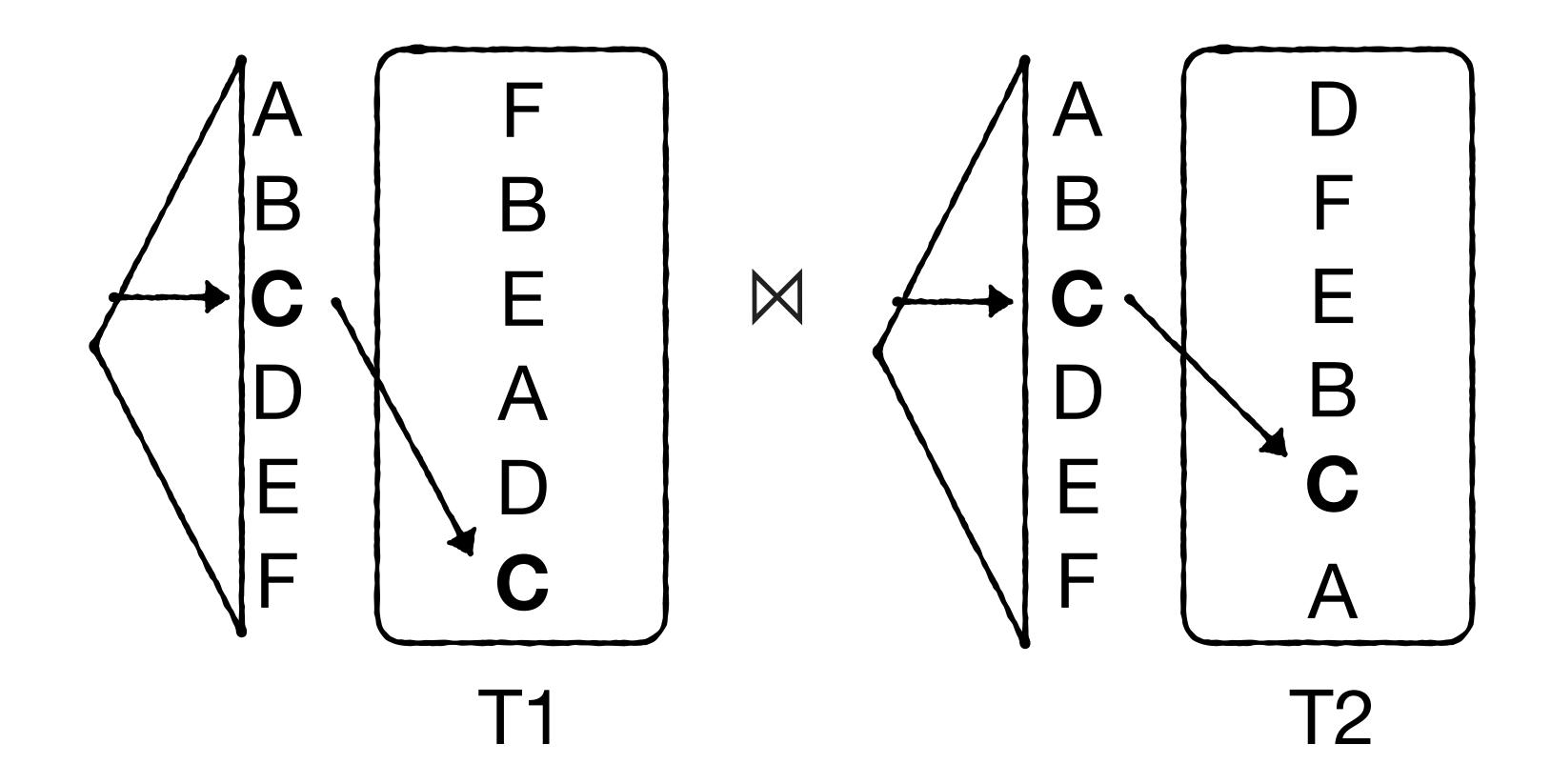
(Assuming B-tree, clustered or unclustered)

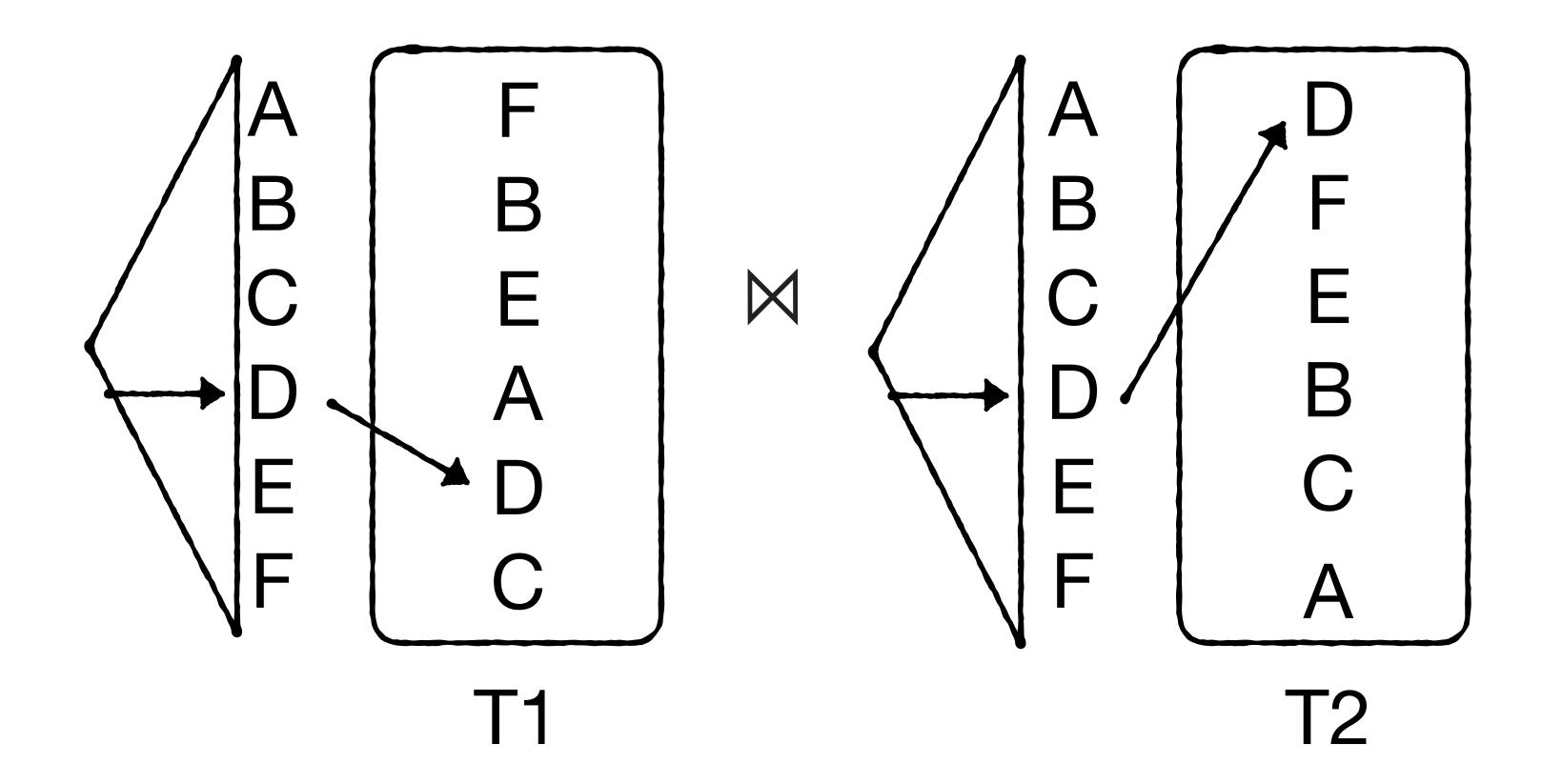




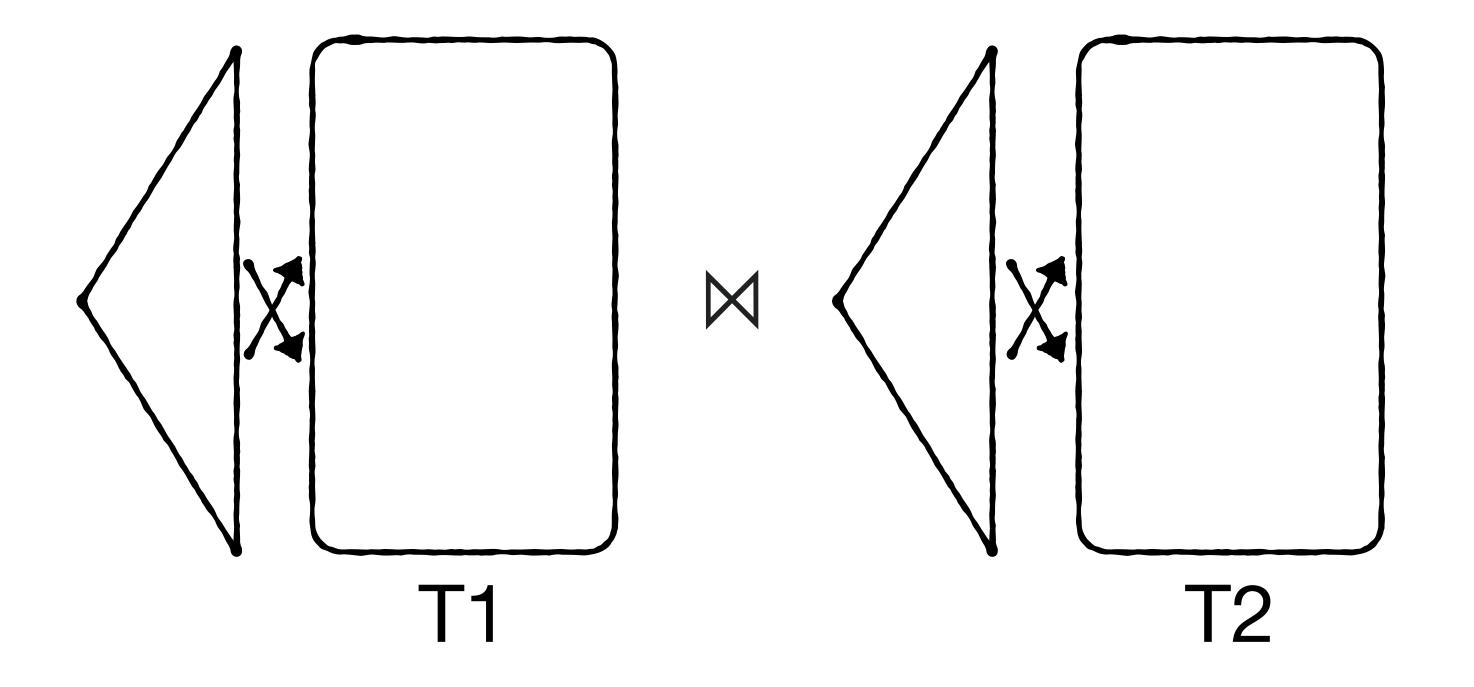




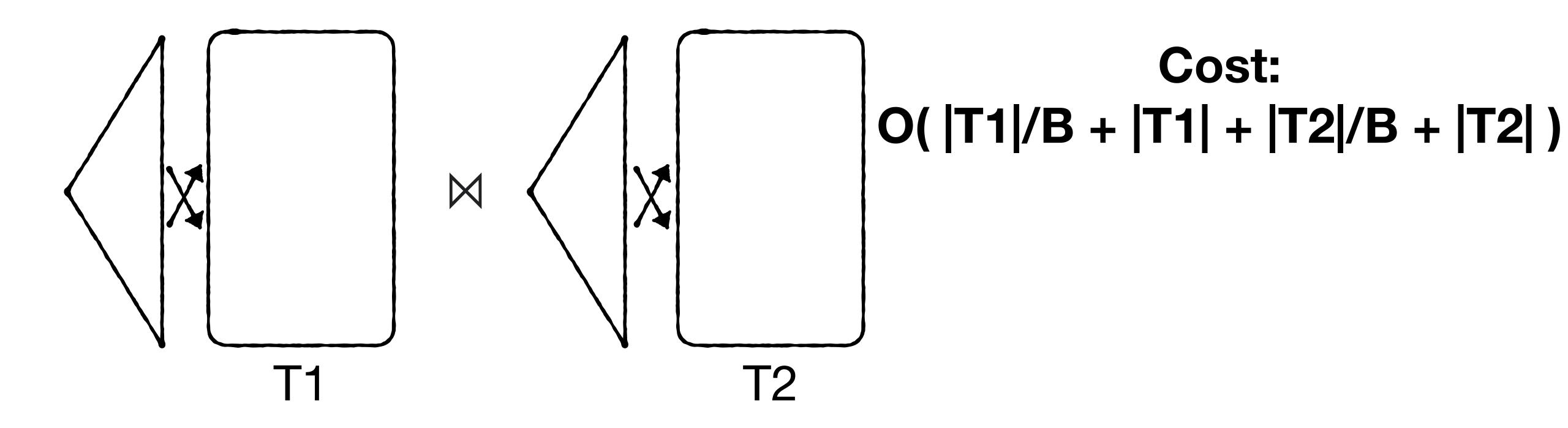




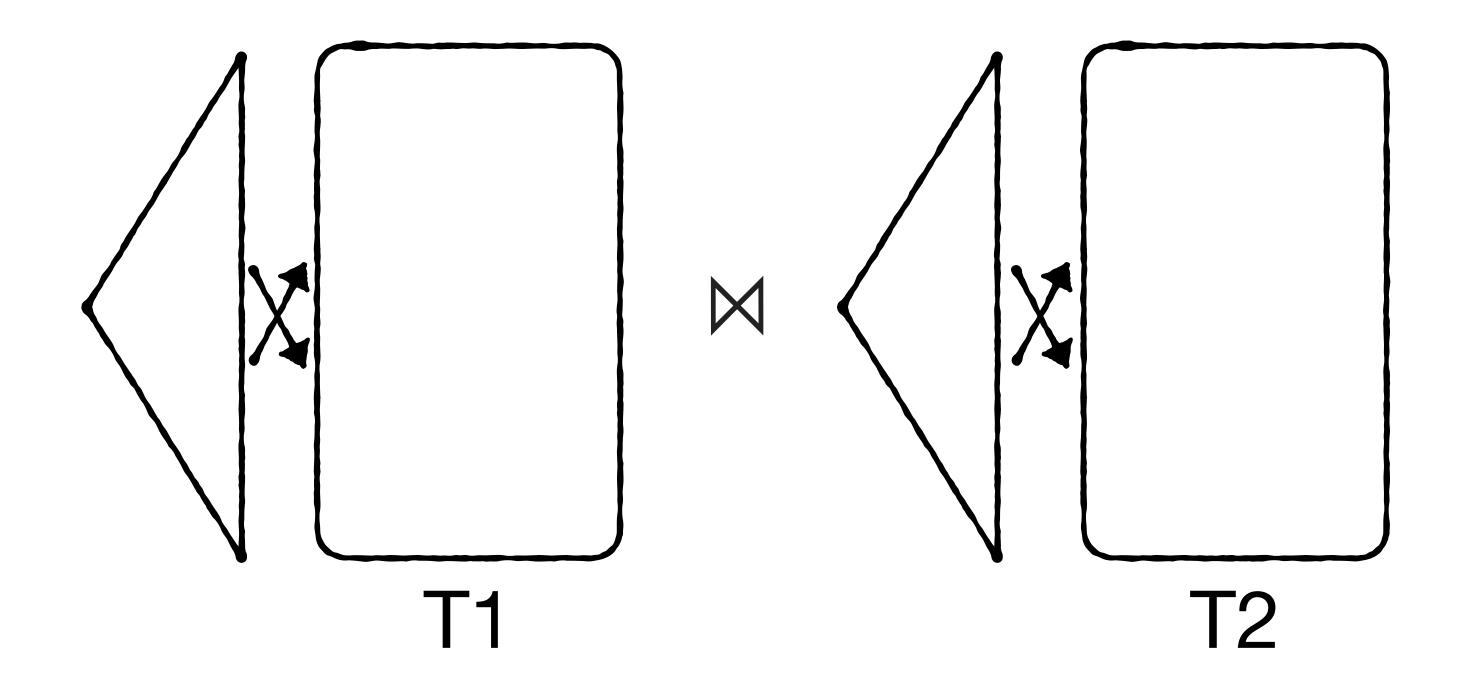
What if both relations have an unclustered index on the join key?



Cost?

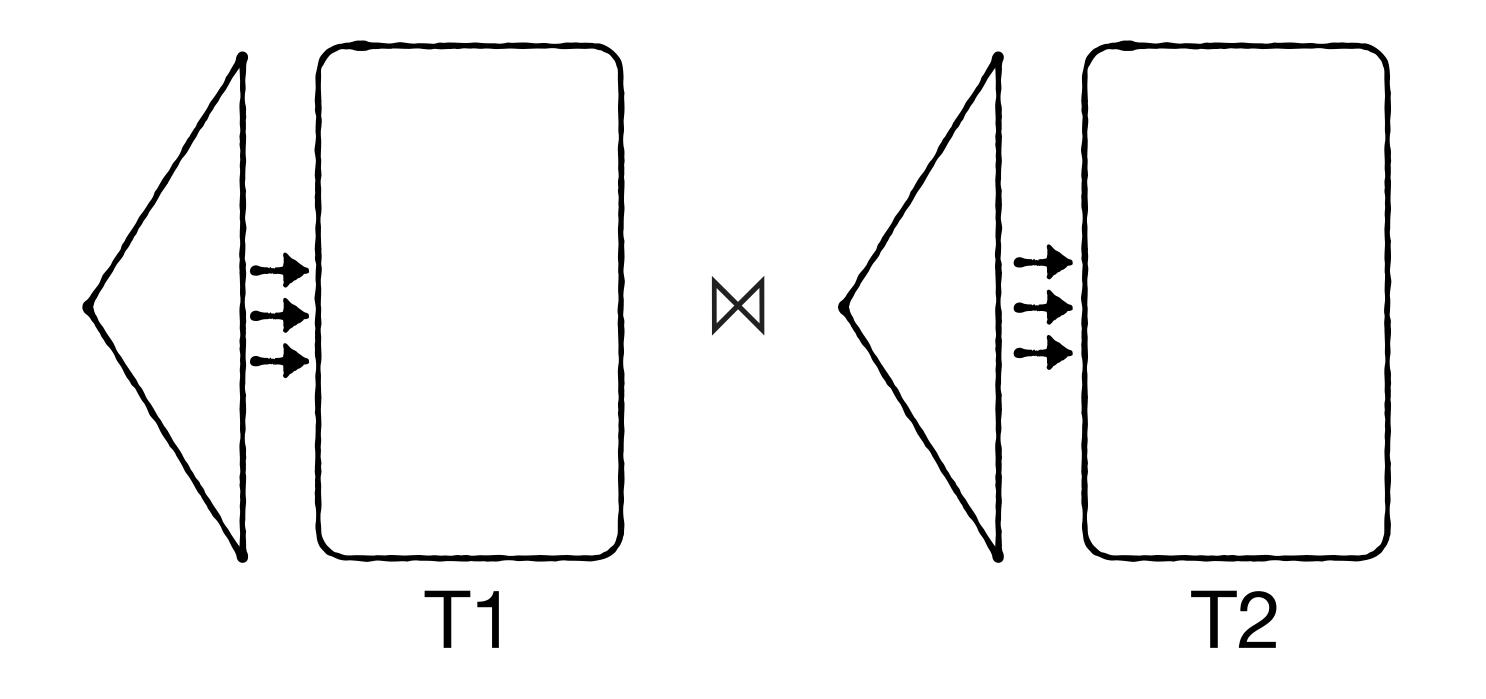


What if both relations have an unclustered index on the join key?



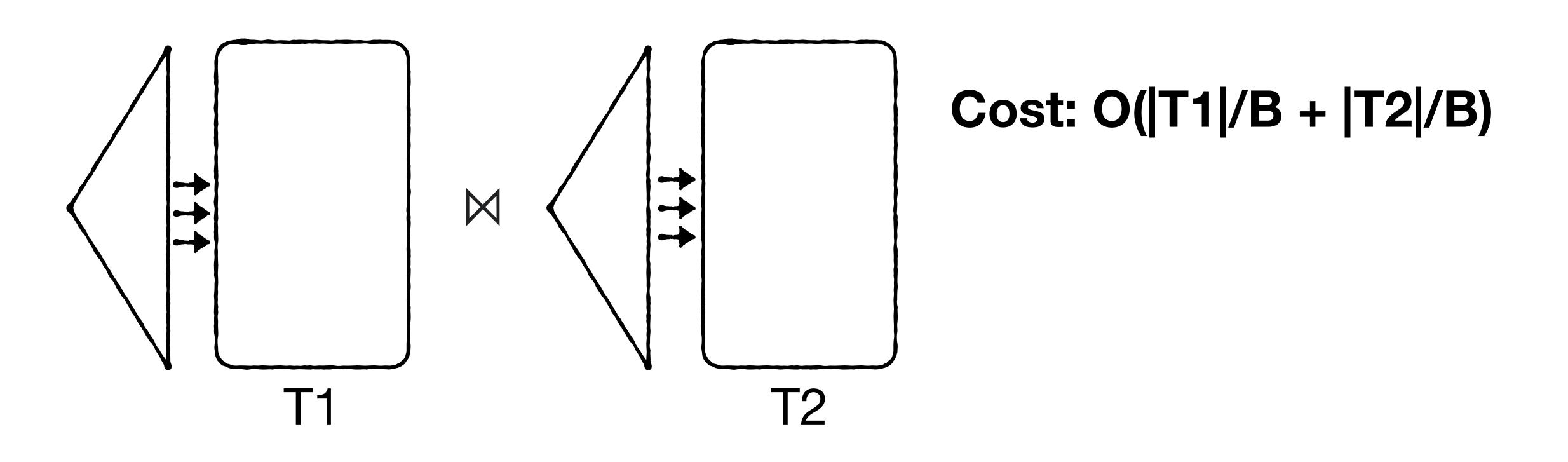
Cost:
O(|T1|+|T2|)

What if both indexes are clustered?



Cost?

What if both indexes are clustered?



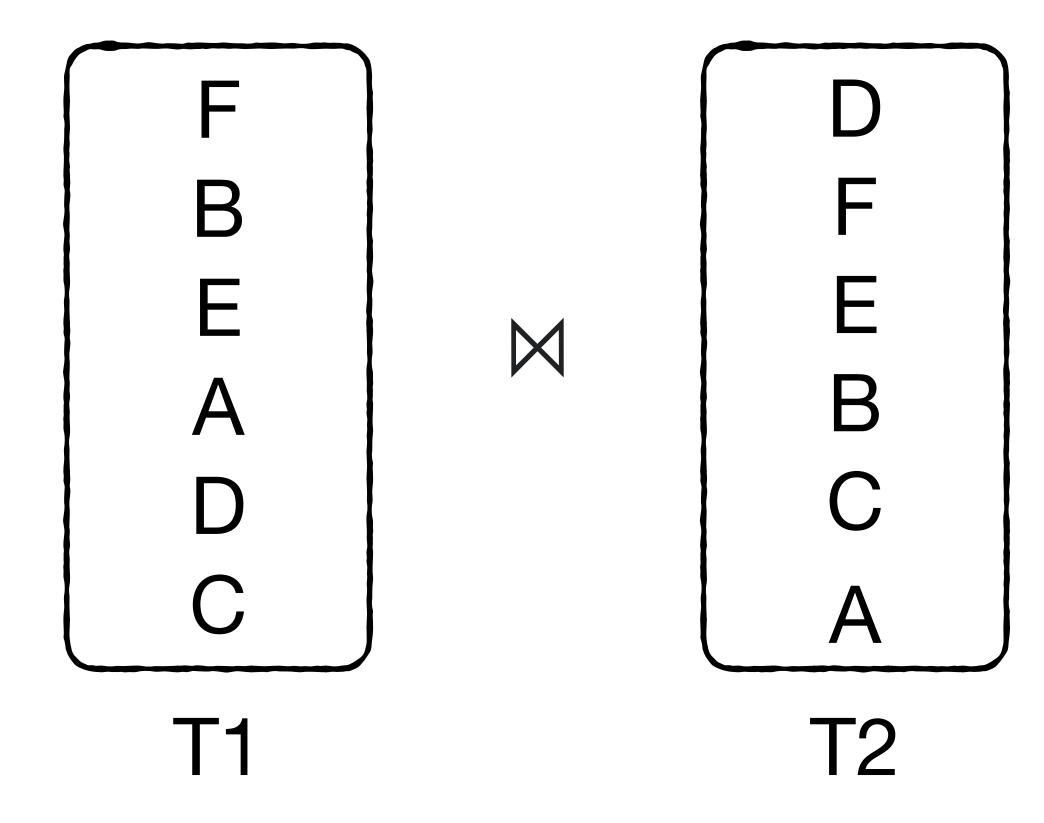
Nested Loop Block Nested Loop

Index-Join

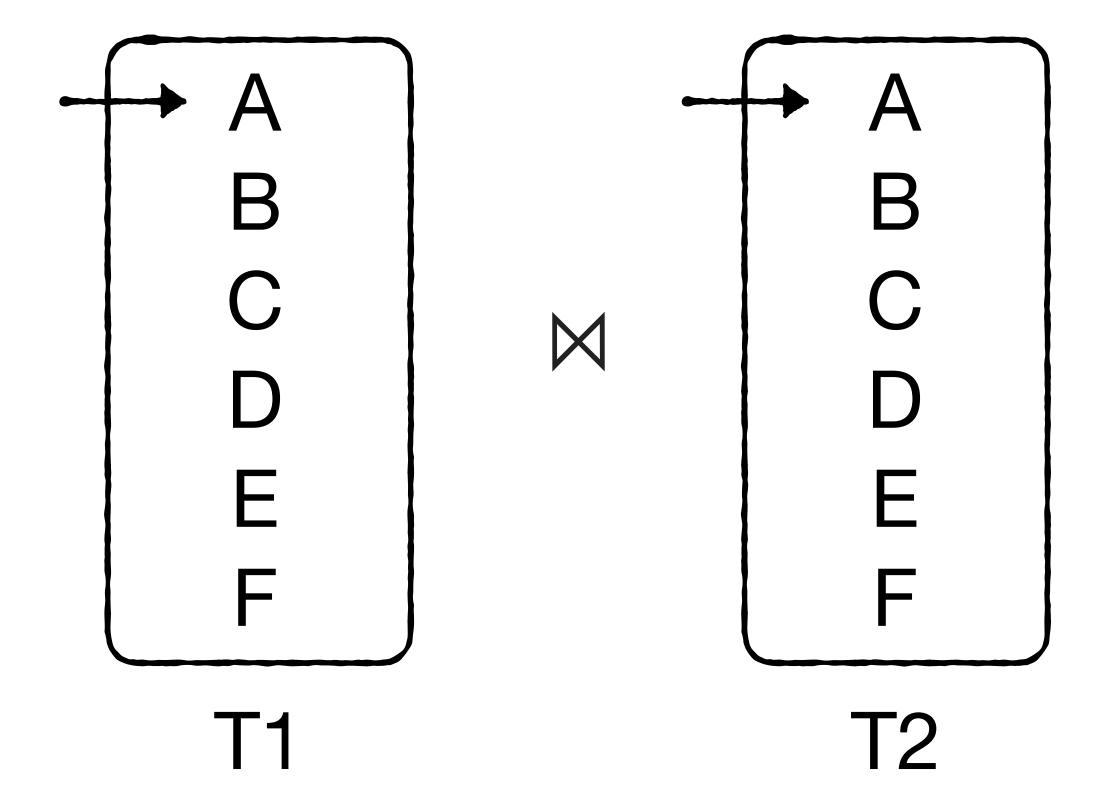
Sort-Merge Join

Grace hash Join

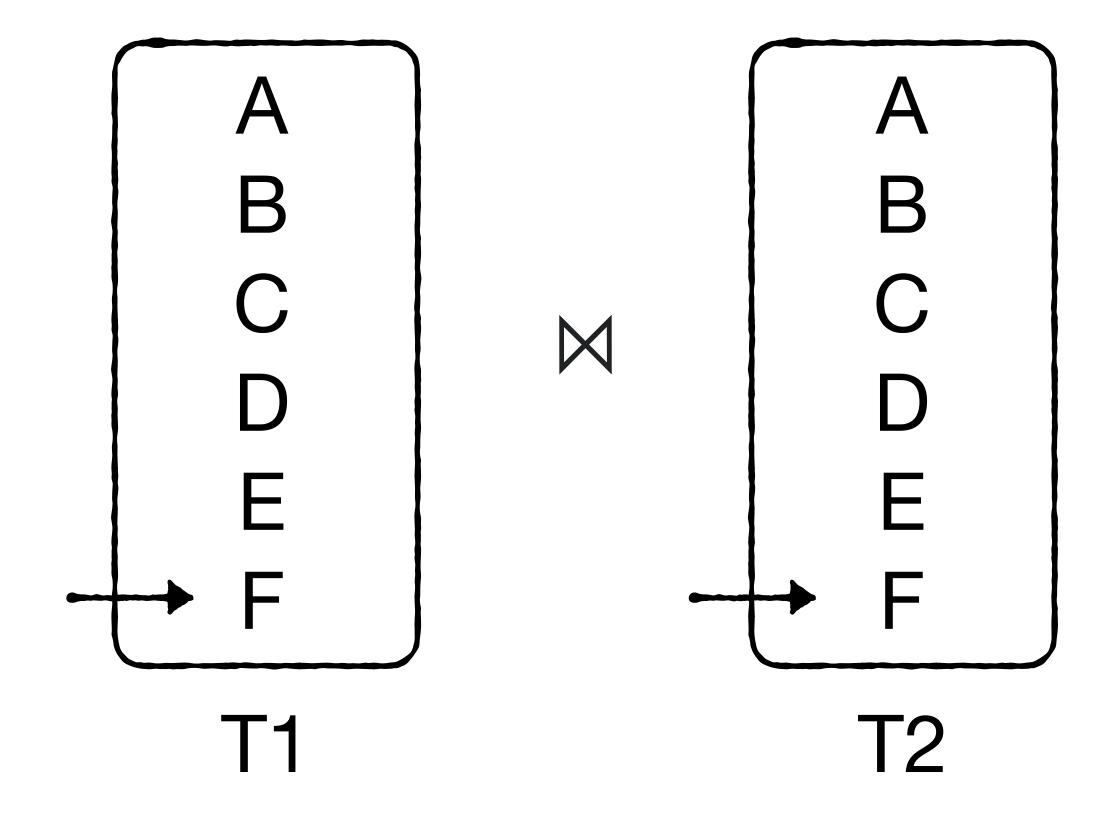
(A) Sort both relations based on join key



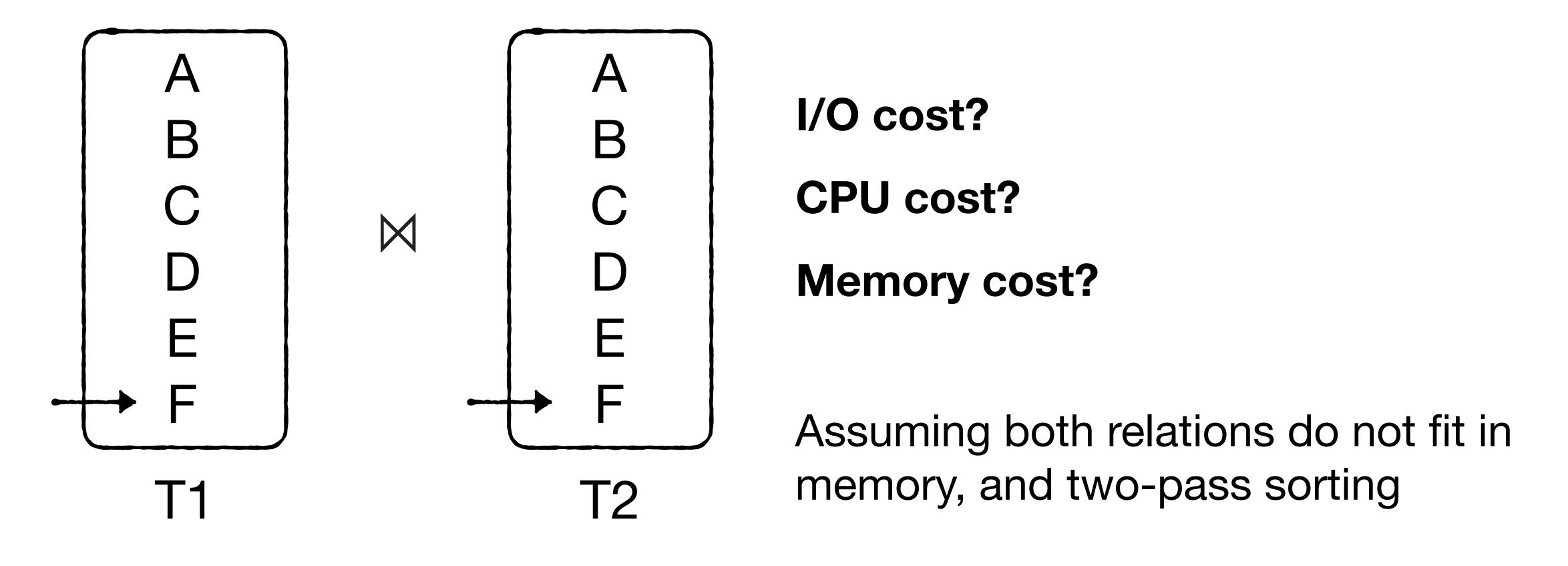
- (A) Sort both relations based on join key
- (B) Scan both relations linearly



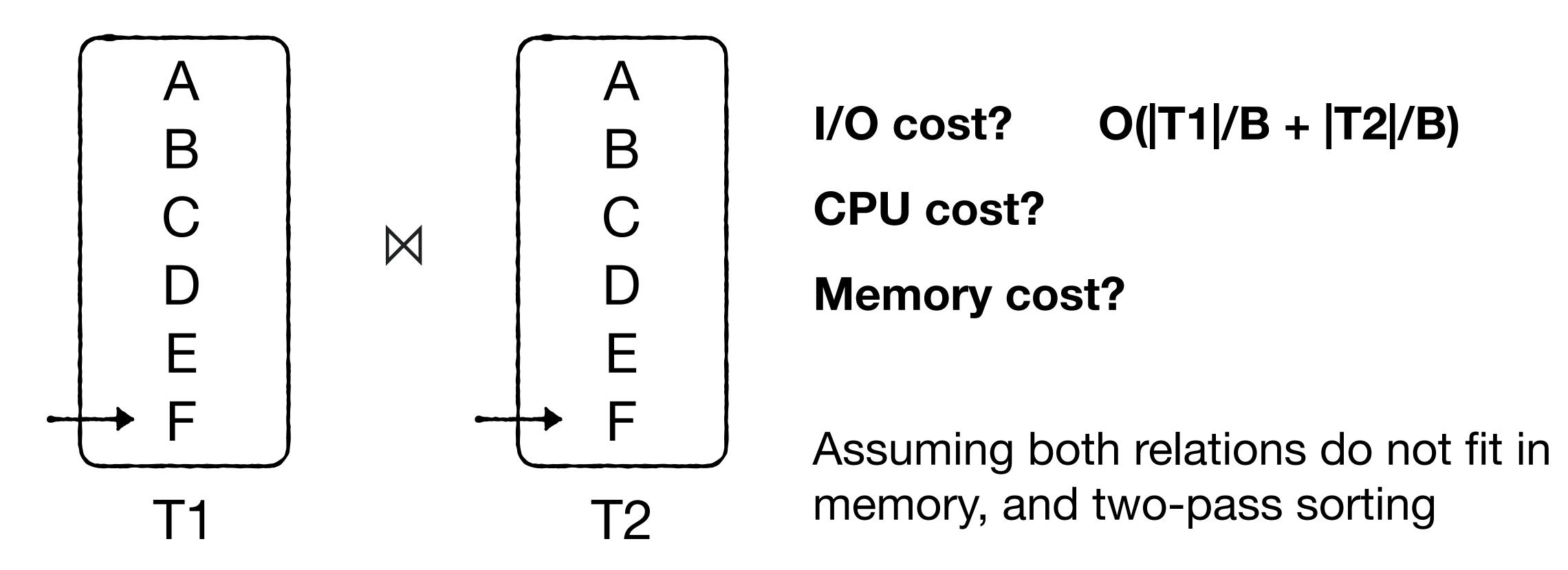
- (A) Sort both relations based on join key
- (B) Scan both relations linearly



- (A) Sort both relations based on join key
- (B) Scan both relations linearly

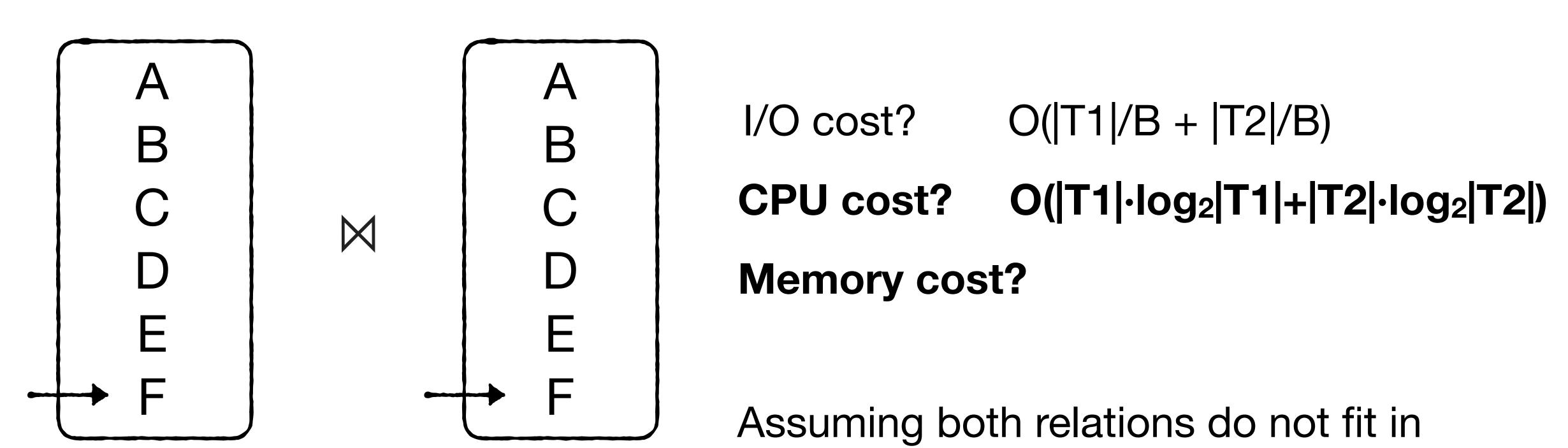


- (A) Sort both relations based on join key
- (B) Scan both relations linearly



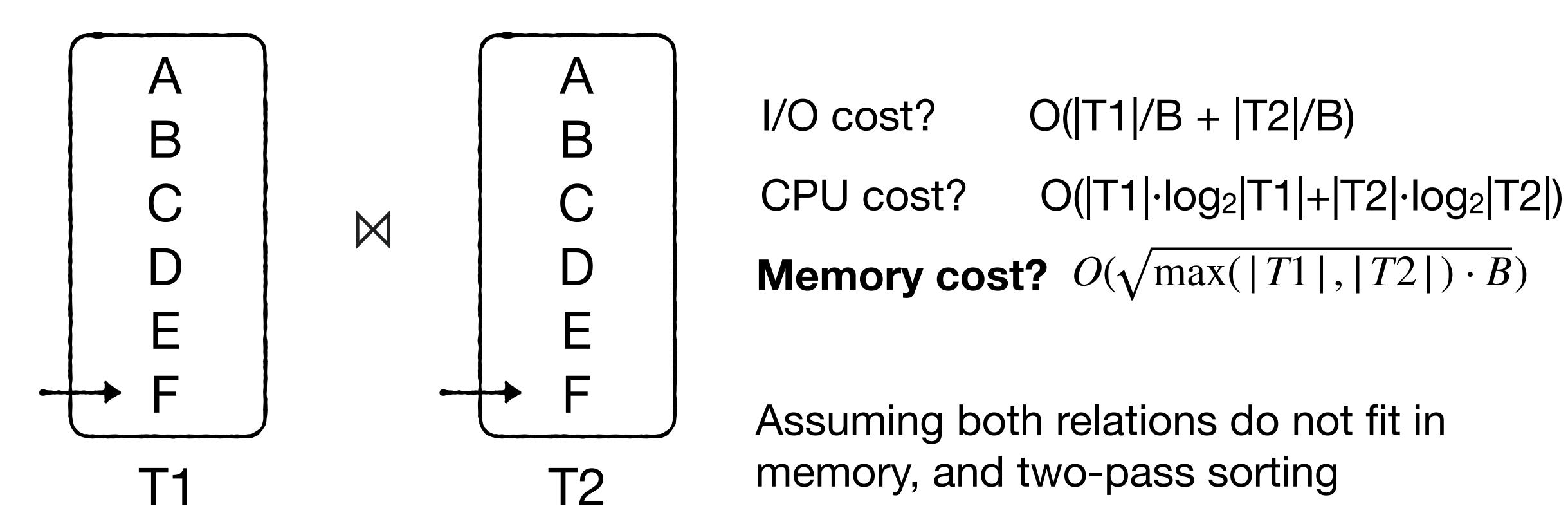
- (A) Sort both relations based on join key
- (B) Scan both relations linearly

T2



memory, and two-pass sorting

- (A) Sort both relations based on join key
- (B) Scan both relations linearly



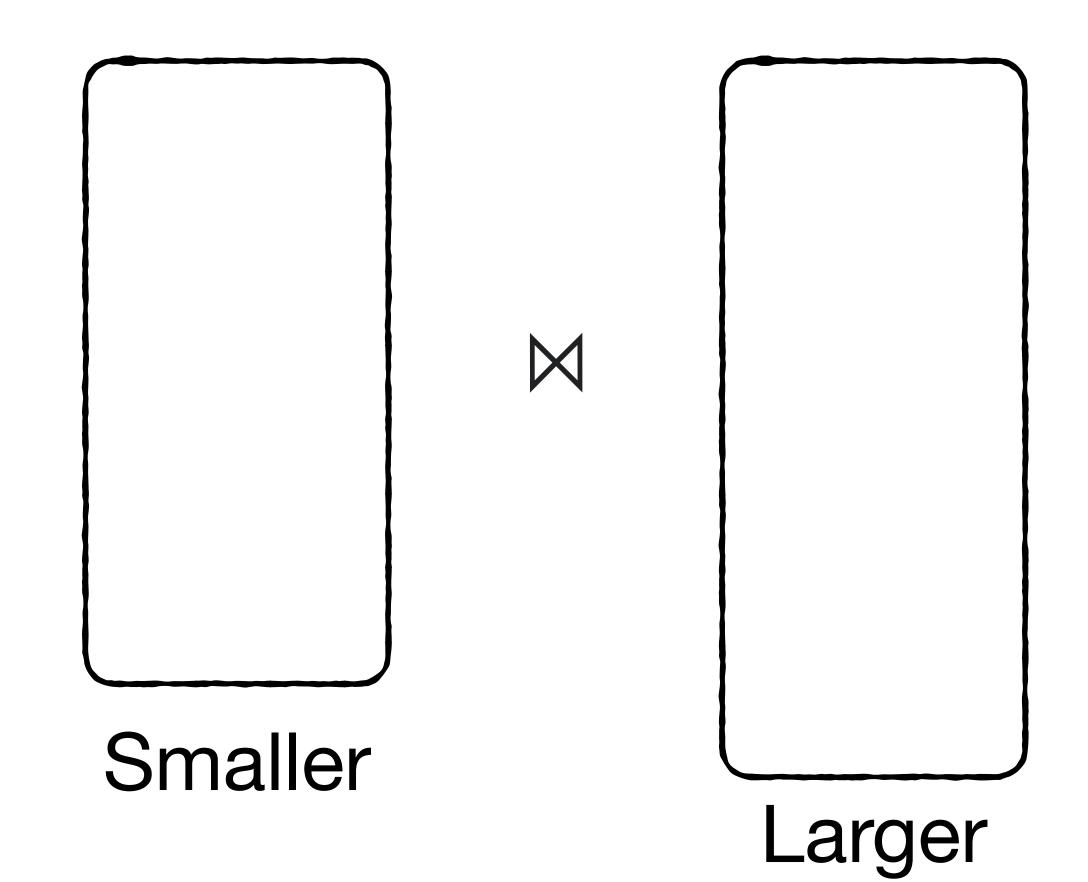
Nested Loop Block Nested Loop

Index-Join

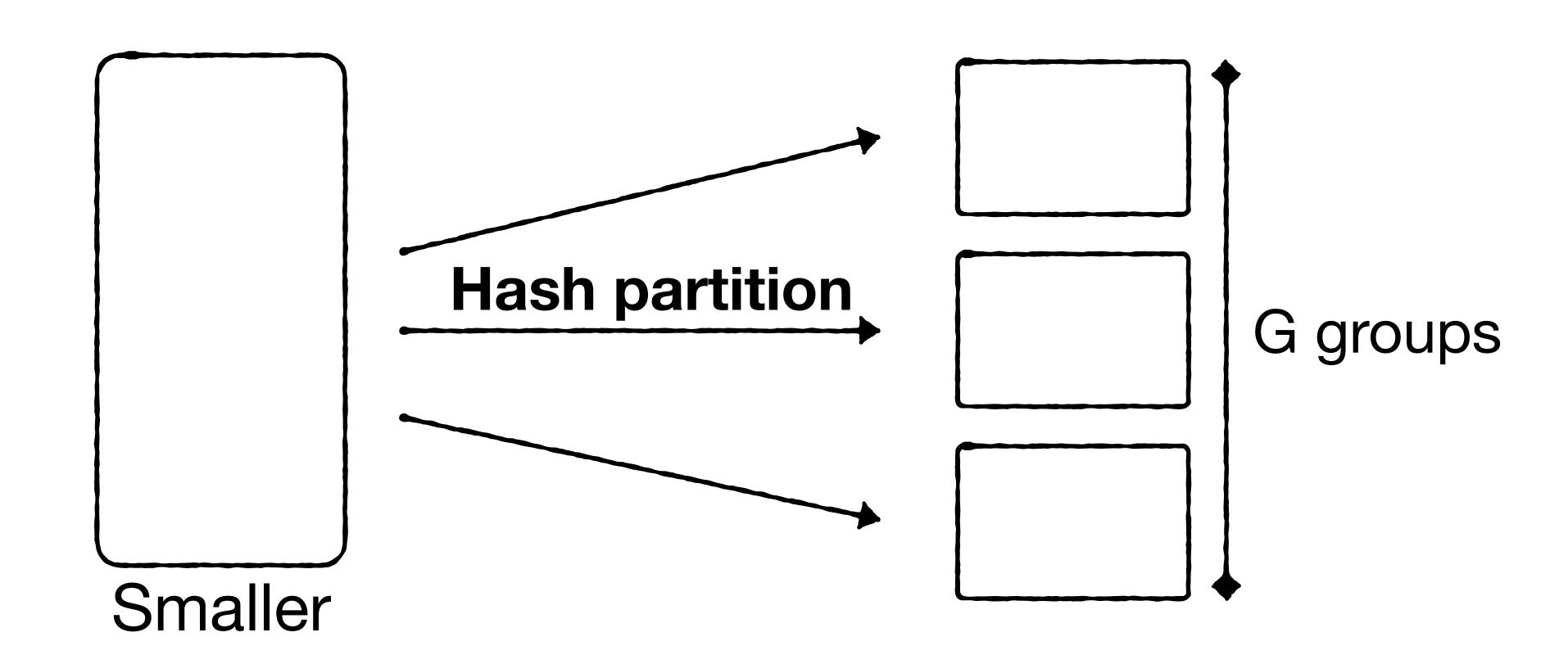
Sort-Merge Join Grace Hash Join

Grace Hash Join

Best for very large relations that do not fit in memory

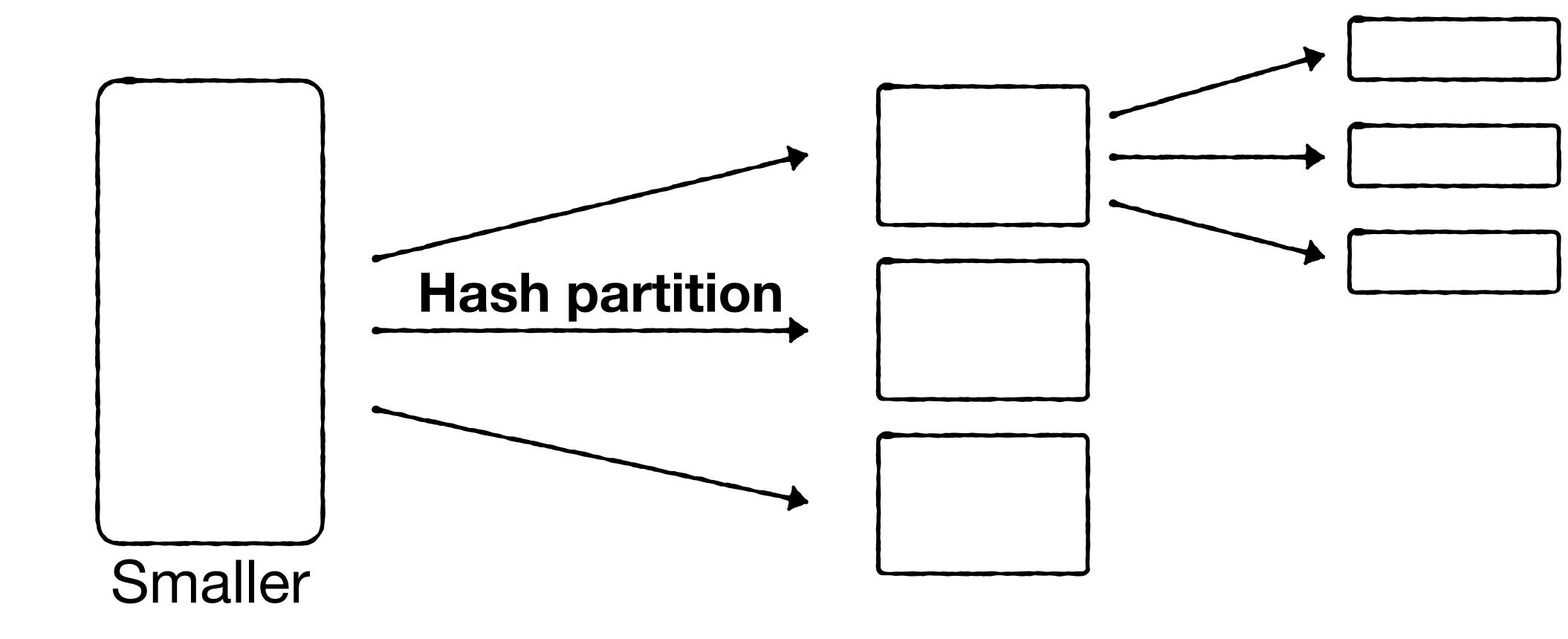


(A) Hash partition smaller table



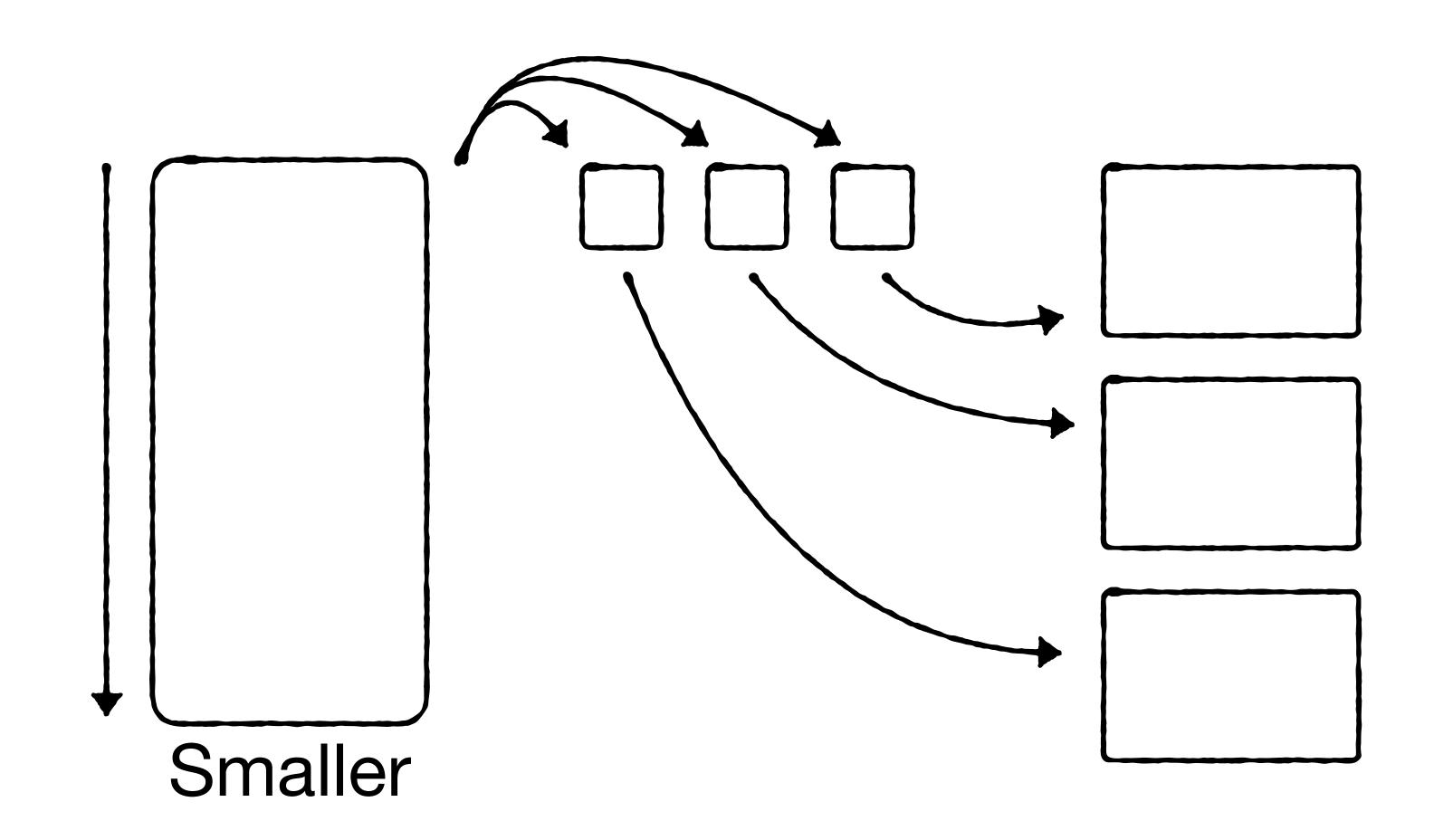
(A) Hash partition smaller table

Repartition until each partition fits in memory

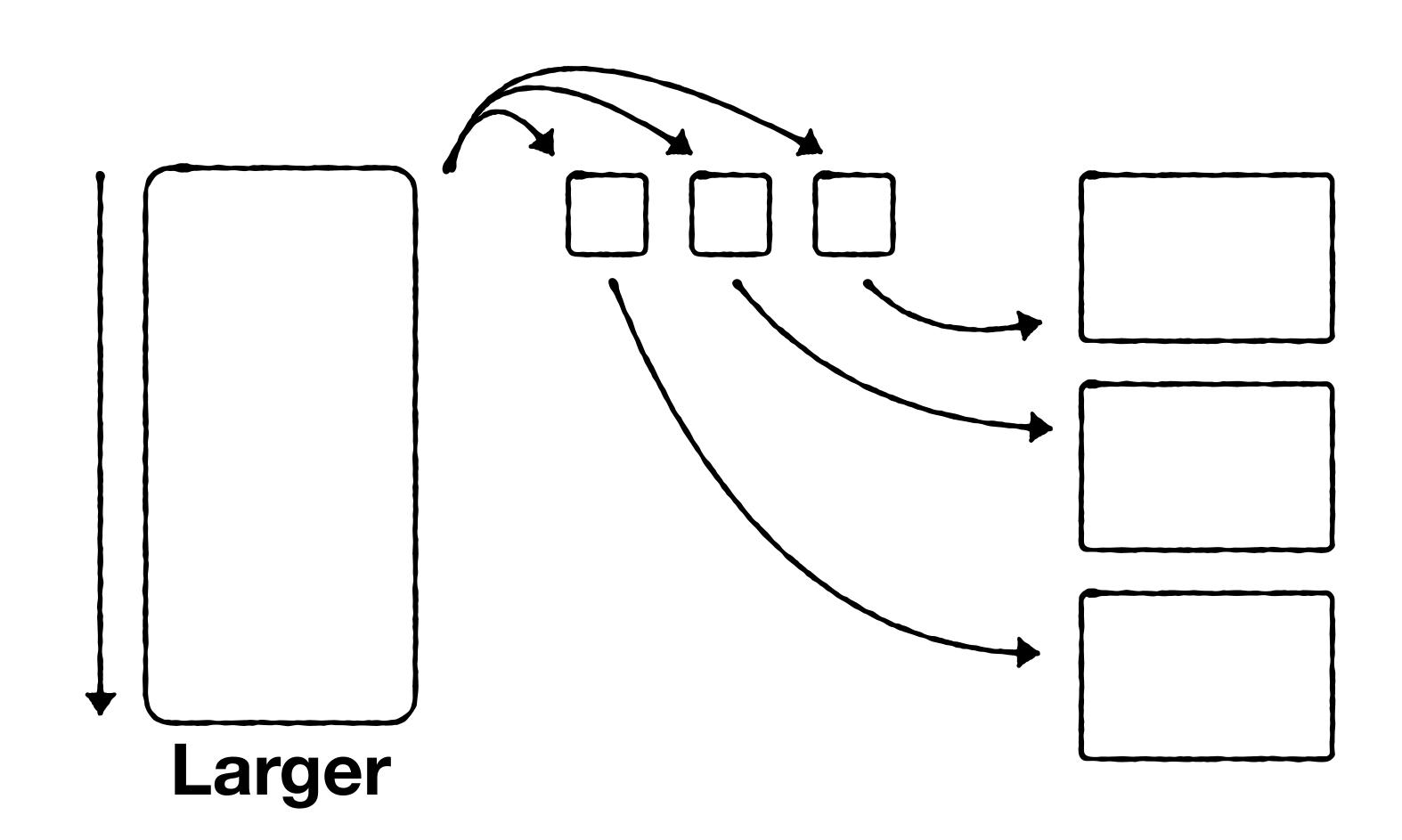


(A) Hash partition smaller table

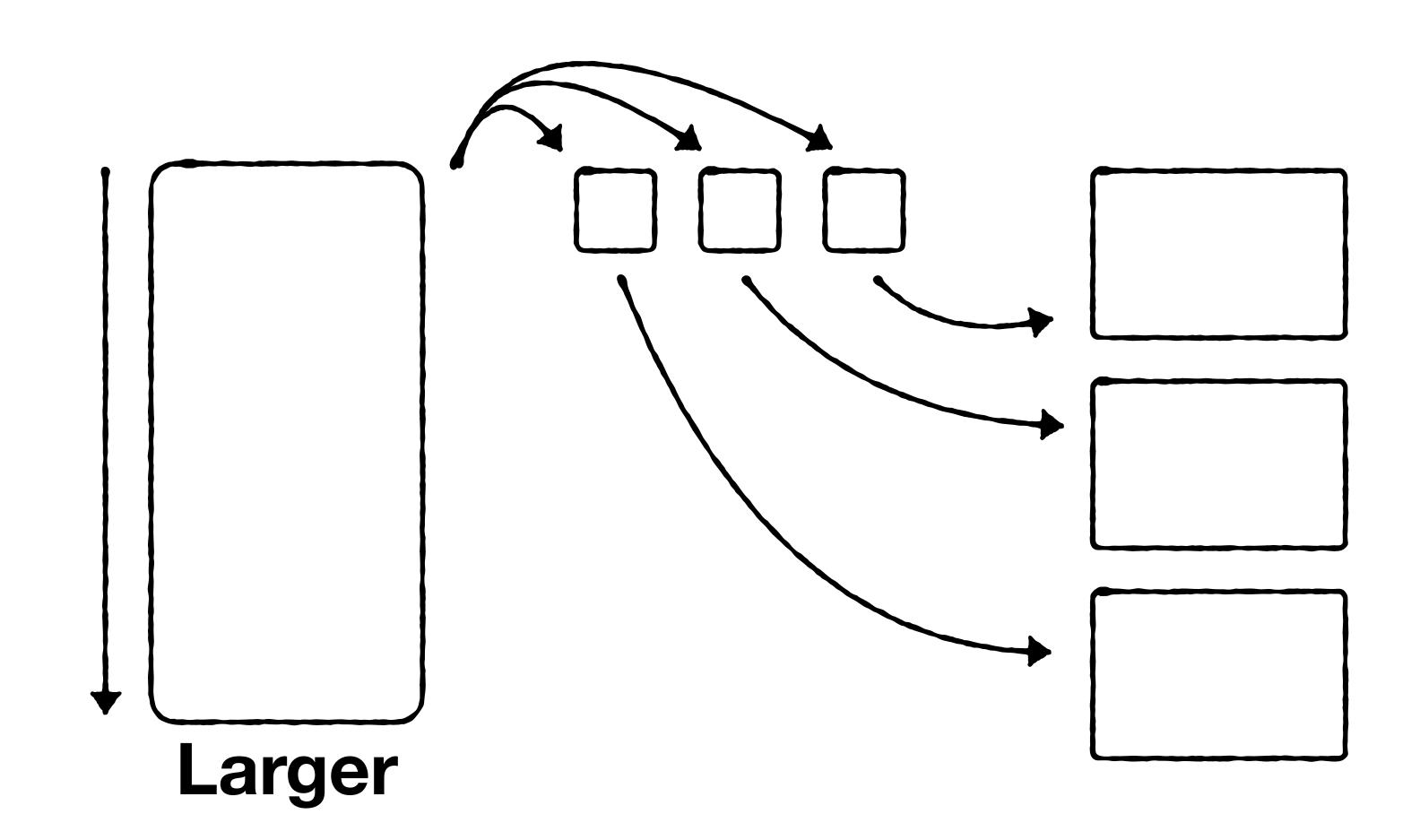
Requires one pass and multiple buffers



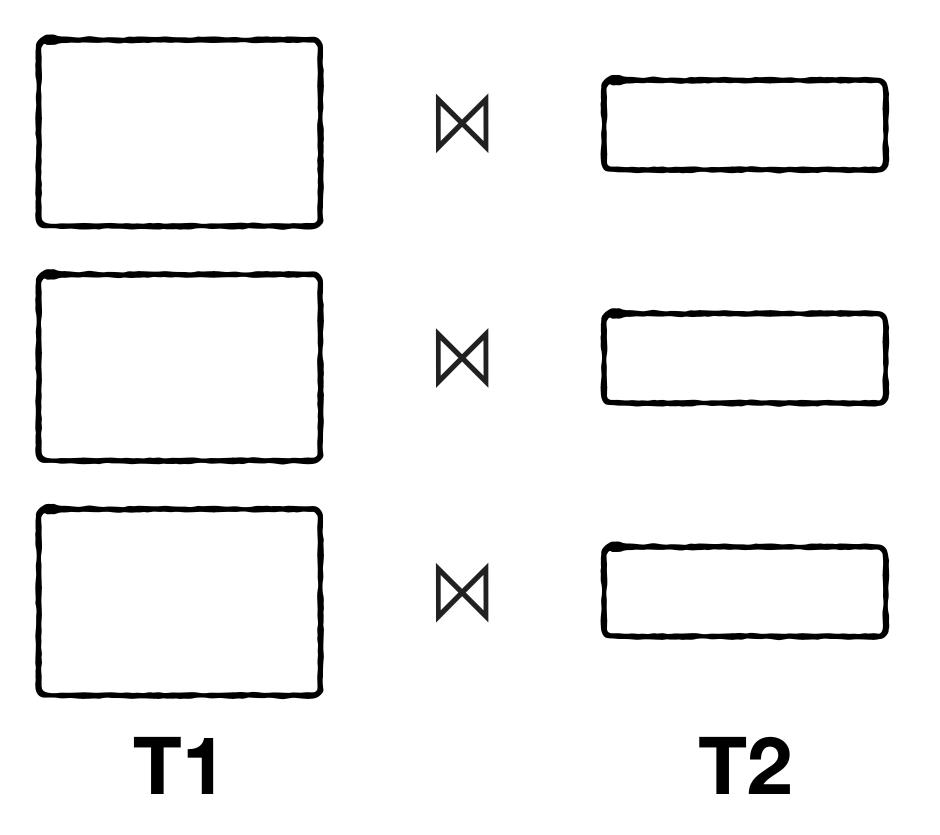
(B) Hash partition larger table using same hash function



(B) Hash partition larger table using same hash function For larger table, each partition can be larger than memory.



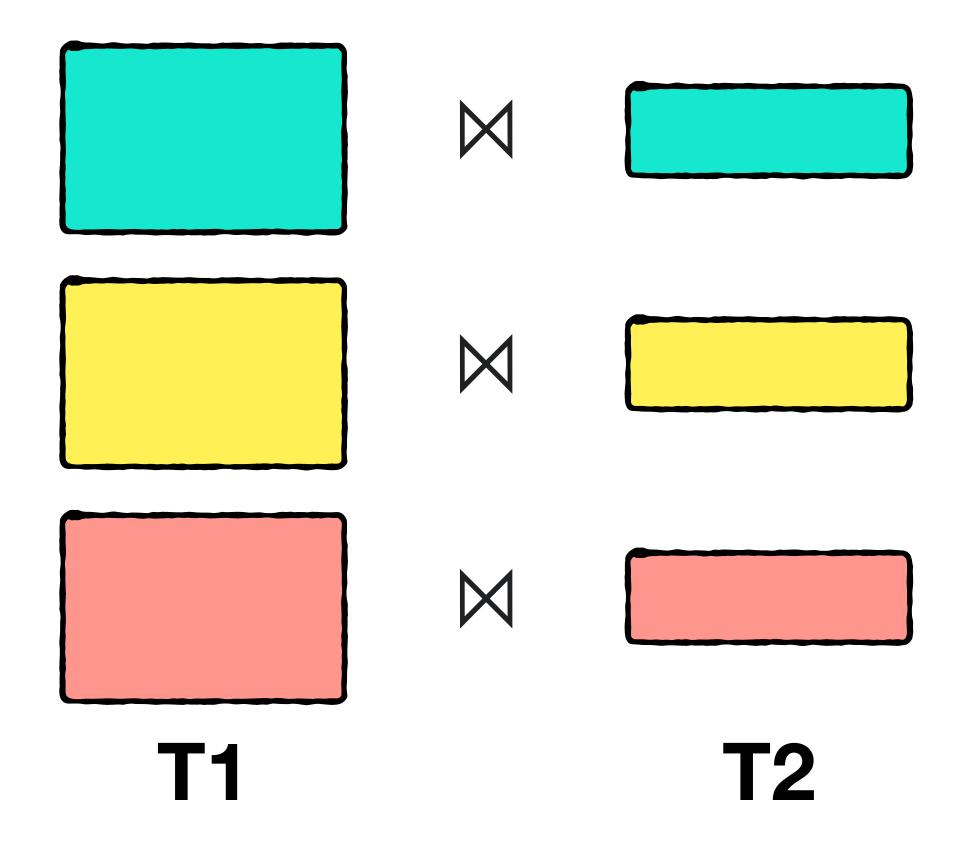
(C) Join each pair of matching partitions independently



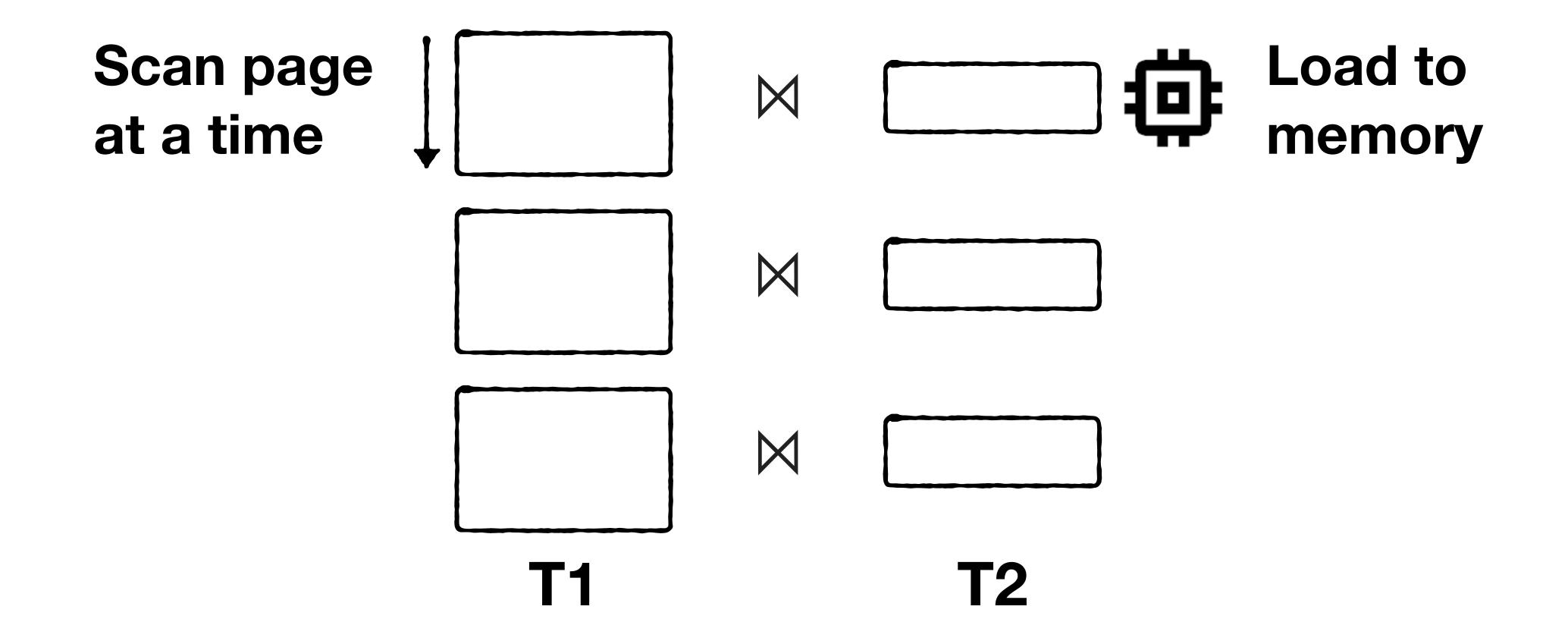
(C) Join each pair of matching partitions independently

Only a pair of matching partitions can have joining

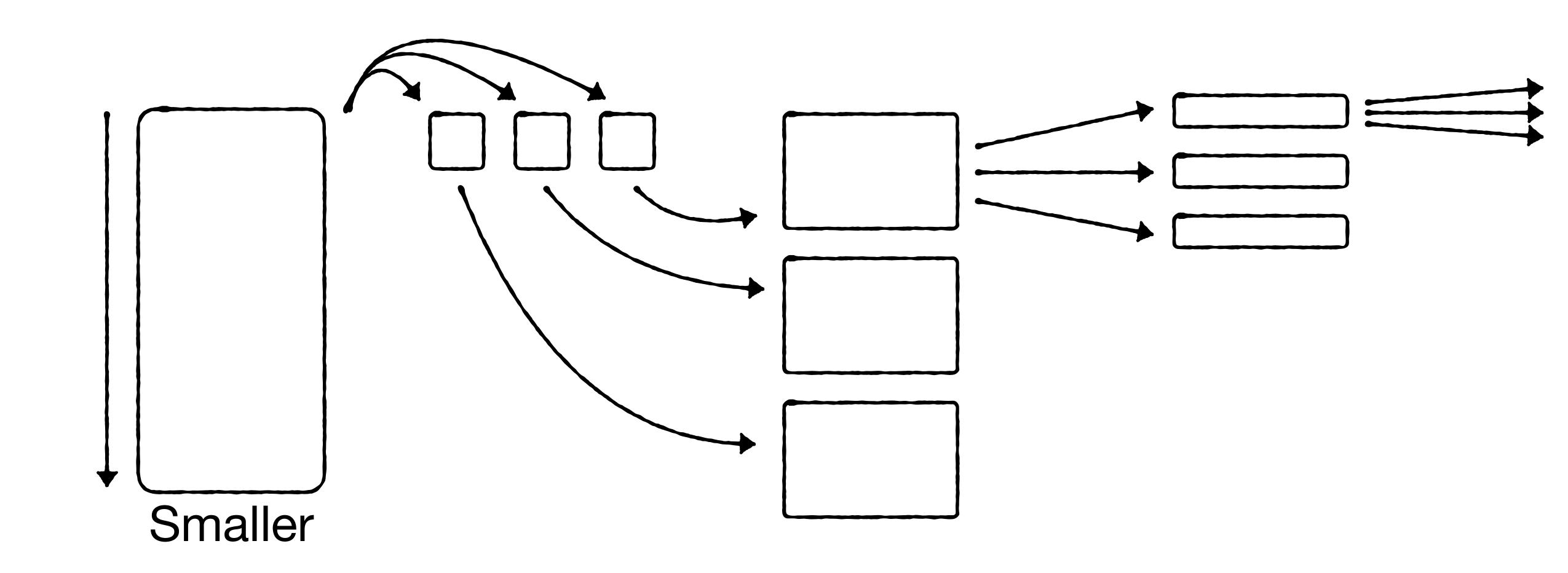
Only a pair of matching partitions can have joining keys since we hash partitioned



(C) Join each pair of matching partitions independently

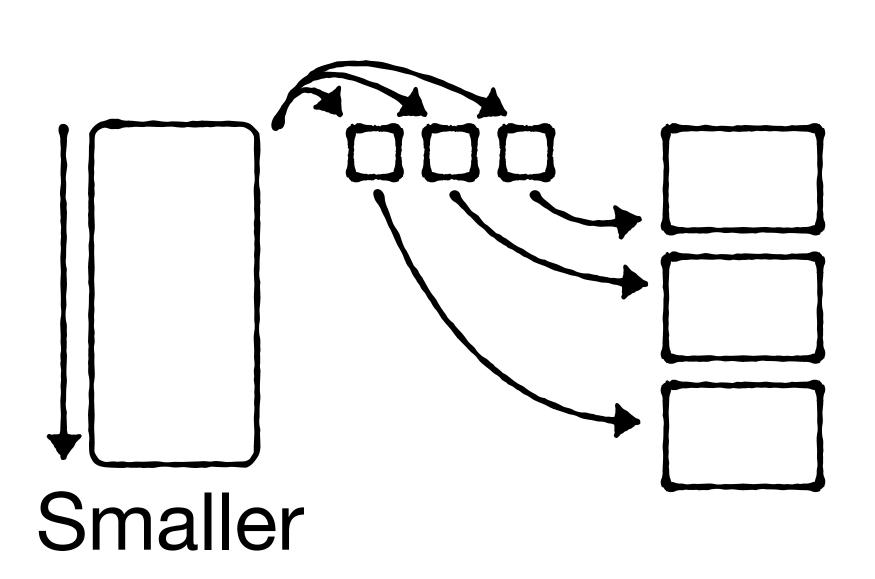


When partitioning smaller table, number of available buffers dictates how many iterations we need until all partitions it in memory



When partitioning smaller table, number of available buffers dictates how many iterations we need until all partitions it in memory

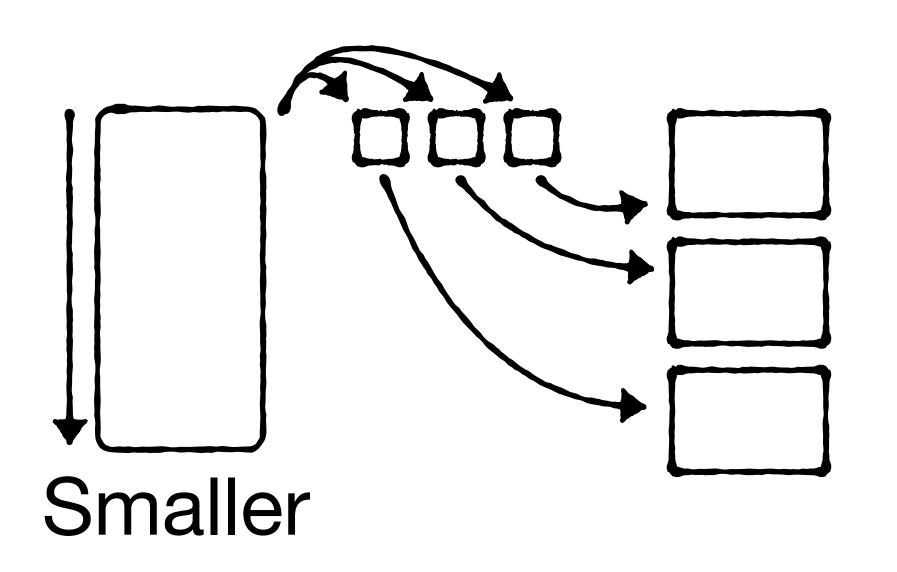
Assuming M entries fit in memory, # iterations needed is:

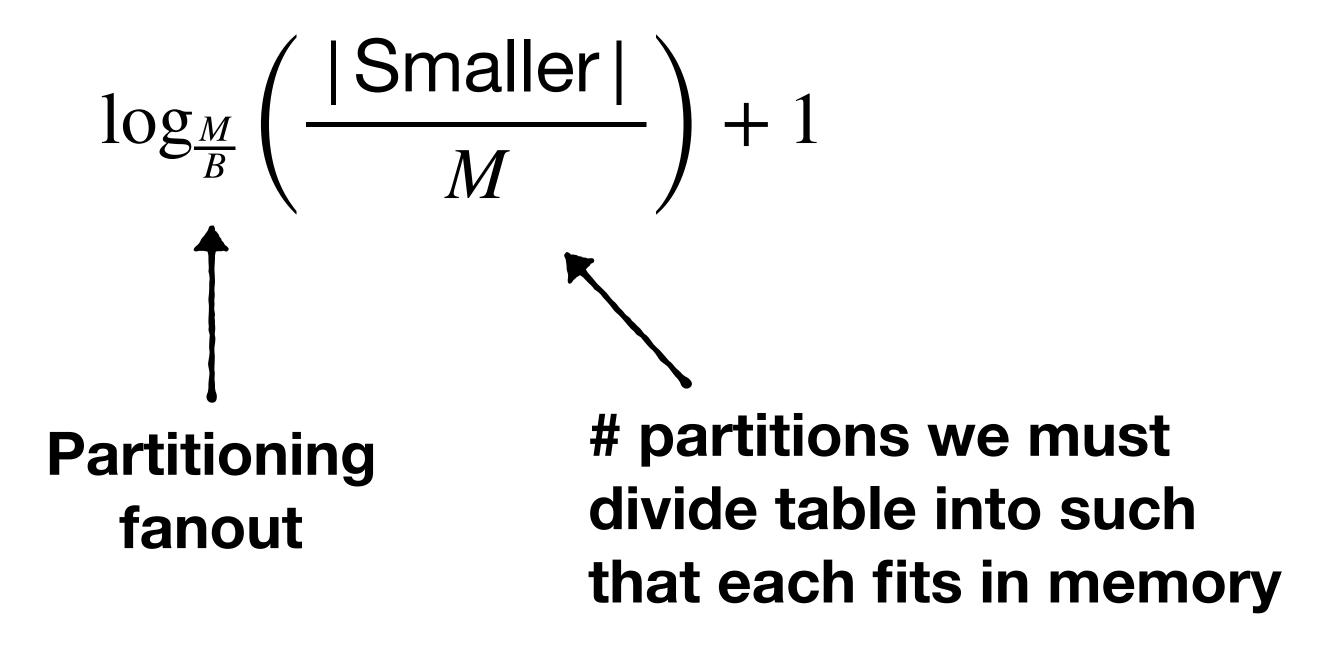


$$\log_{\frac{M}{B}} \left(\frac{|\mathsf{Smaller}|}{M} \right) + 1$$

When partitioning smaller table, number of available buffers dictates how many iterations we need until all partitions it in memory

Assuming M entries fit in memory, # iterations needed is:



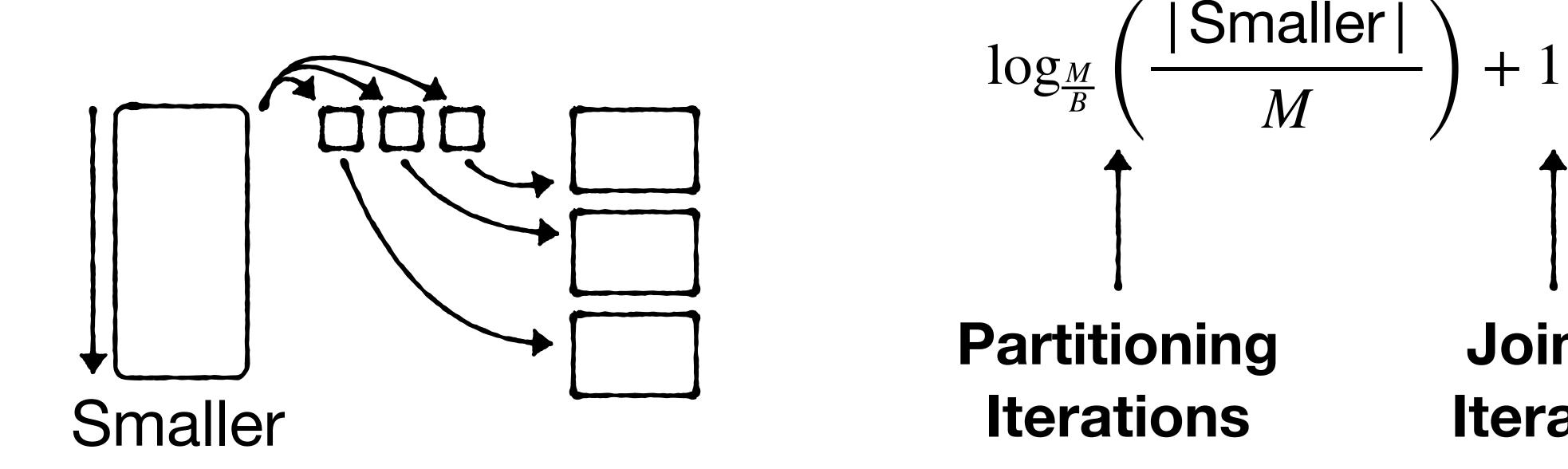


Joining

Iteration

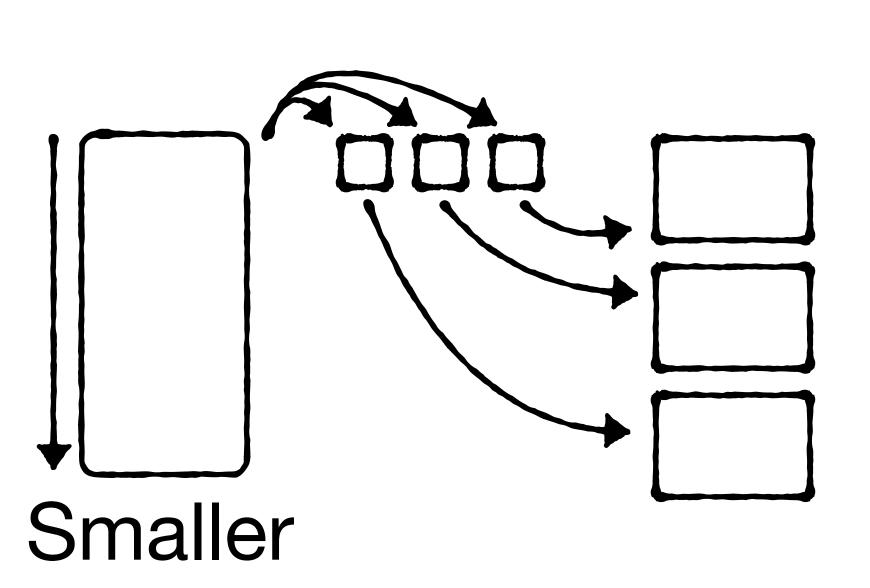
When partitioning smaller table, number of available buffers dictates how many iterations we need until all partitions it in memory

Assuming M entries fit in memory, # iterations needed is:



When partitioning smaller table, number of available buffers dictates how many iterations we need until all partitions it in memory

Assuming M entries fit in memory, # iterations needed is:

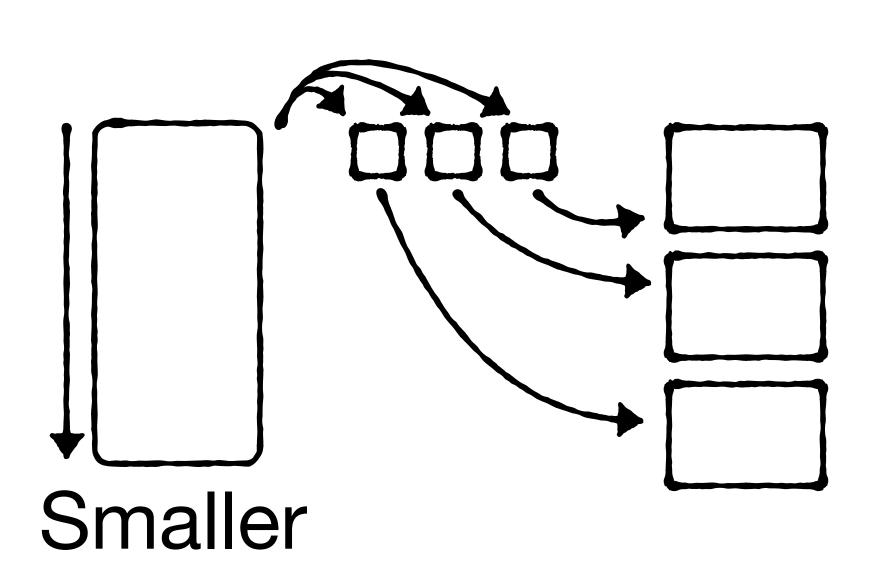


$$\log_{\frac{M}{B}} \left(\frac{|\mathsf{Smaller}|}{M} \right) + 1 = \log_{\frac{M}{B}} \left(\frac{|\mathsf{Smaller}|}{B} \right)$$

Same as external merge sort:)

When partitioning smaller table, number of available buffers dictates how many iterations we need until all partitions it in memory

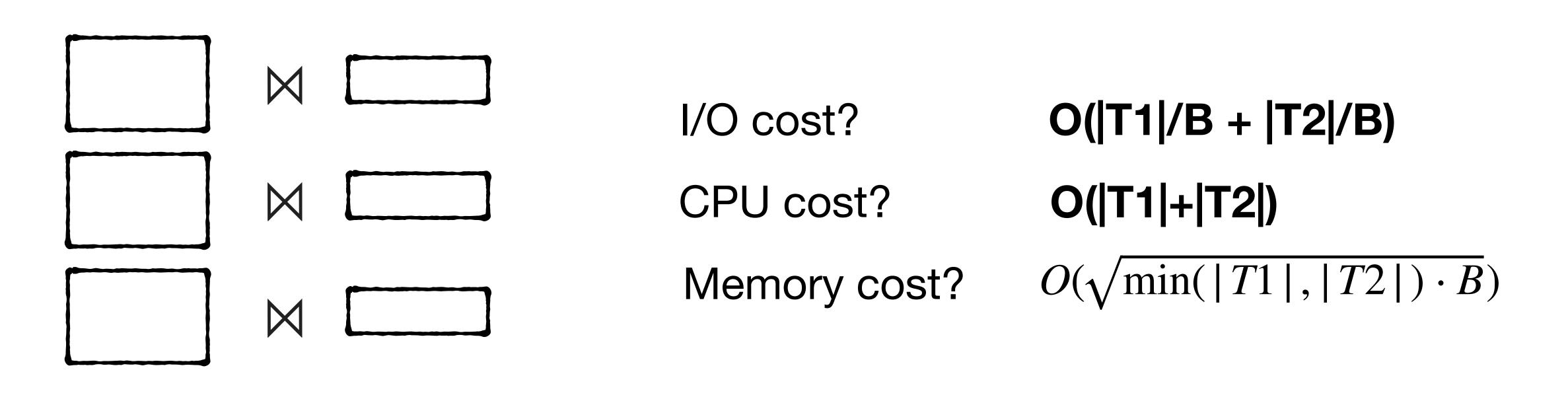
Assuming M entries fit in memory, # iterations needed is:



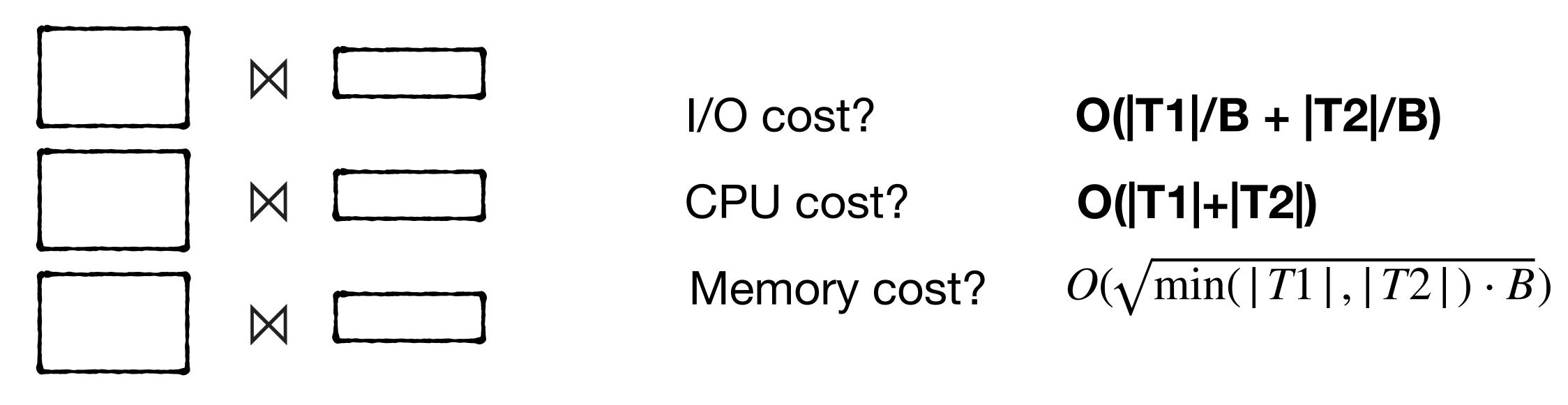
$$\log_{\frac{M}{B}} \left(\frac{|\mathsf{Smaller}|}{B} \right) = 2 \longrightarrow M = \sqrt{|\mathit{Smaller}| \cdot B}$$

Equate to 2 and solve for M to obtain memory needed to join in two passes

Overall analysis including both tables, assuming two-pass partitioning

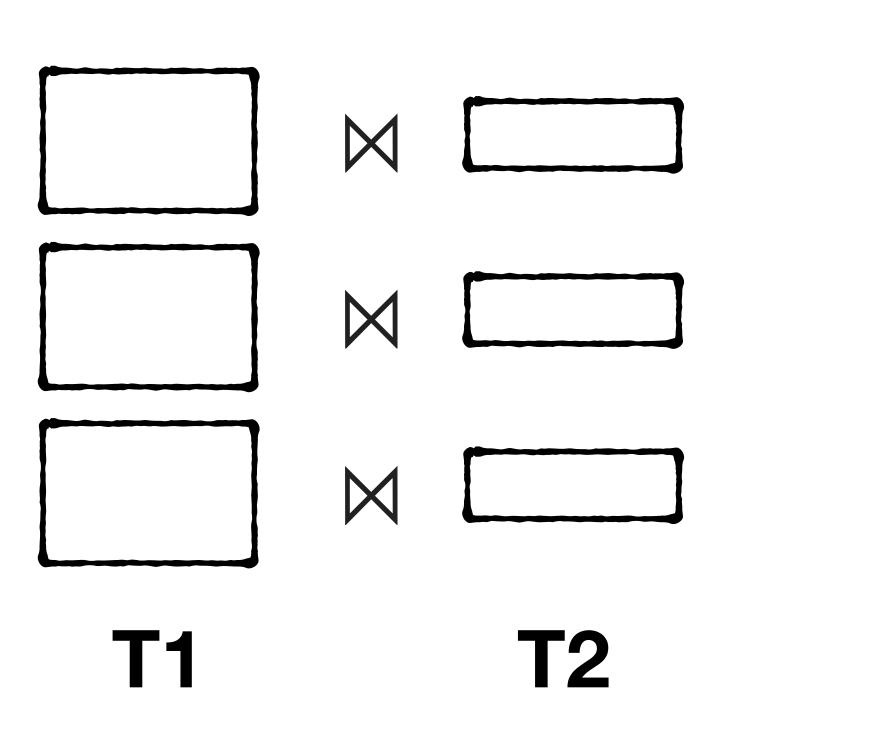


Overall analysis including both tables, assuming two-pass partitioning



T1 T2 Lower memory footprint and CPU cost than merge-sort join :)

Overall analysis including both tables, assuming two-pass partitioning



I/O cost? O(|T1|/B + |T2|/B)CPU cost? O(|T1|+|T2|)Memory cost? $O(\sqrt{\min(|T1|,|T2|)\cdot B})$

Lower memory footprint and CPU cost than merge-sort join:)

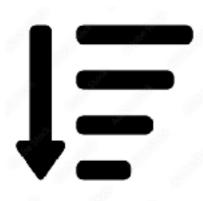
Sort join is still better if we have an order by clause

Query Operators

Selection



Order by



Projection



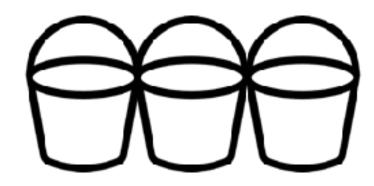
Distinct



Join



Group by

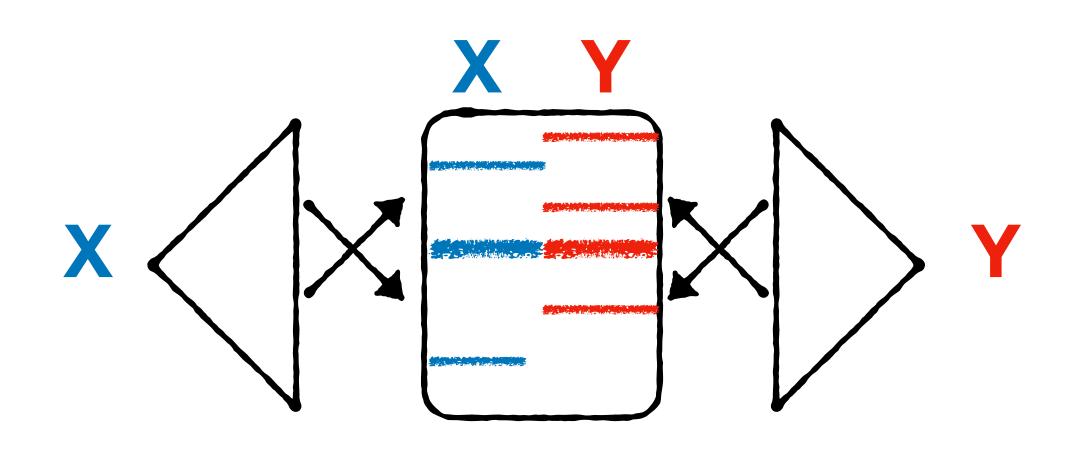


Query Operators Cardinality Estimation Query Optimization Query Operators

Cardinality Estimation

Query Optimization



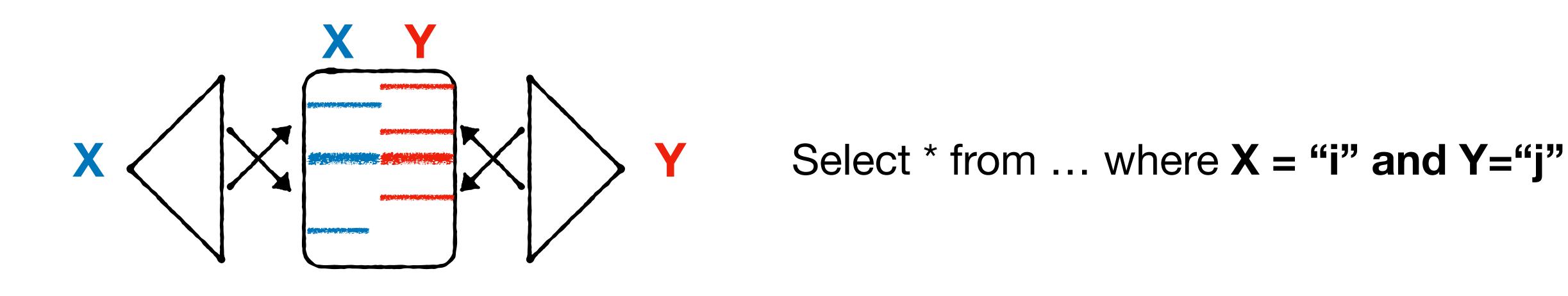


Select * from ... where **X** = "i" and **Y**="j"

Algo 1: Search index Y $log_B(N) + |Y_j| I/O$

Algo 2: Search index X $log_B(N) + |X_i| I/O$

Algo 3: Search both 2 · log_B(N) + |X_i|/B + |Y_j|/B + |X_i∩Y_j| I/O

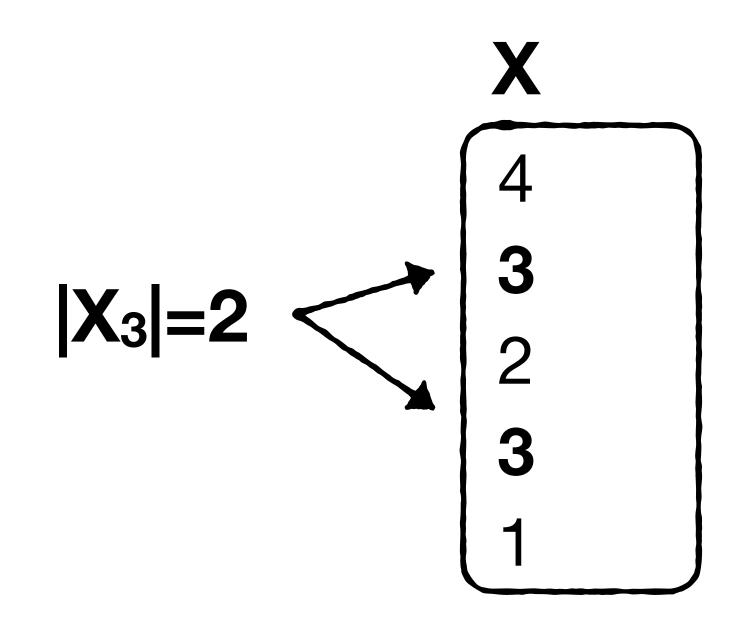


Algo 1: Search index Y $log_B(N) + |Y_j| I/O$

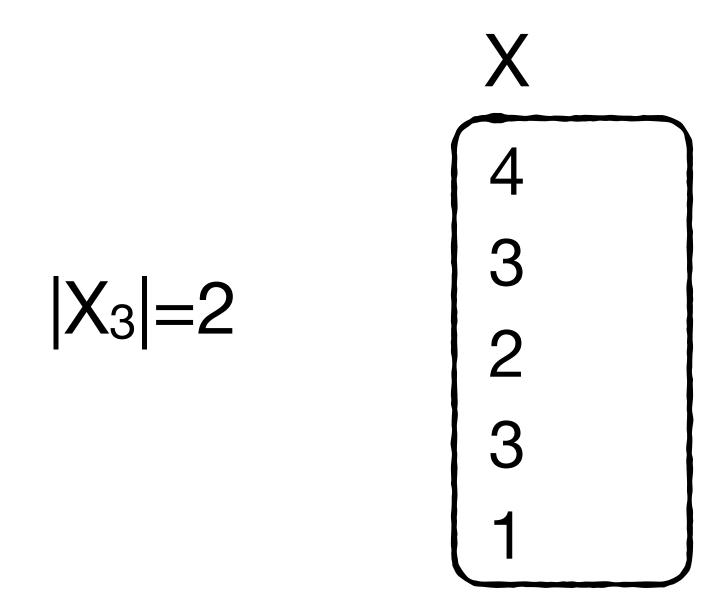
Algo 2: Search index X $log_B(N) + |X_i| I/O$

Algo 3: Search both $2 \cdot \log_B(N) + |X_i|/B + |Y_j|/B + |X_i \cap Y_j| I/O$

Determining best plan requires estimating |Xi|, |Yj| and |Xi∩Yj|



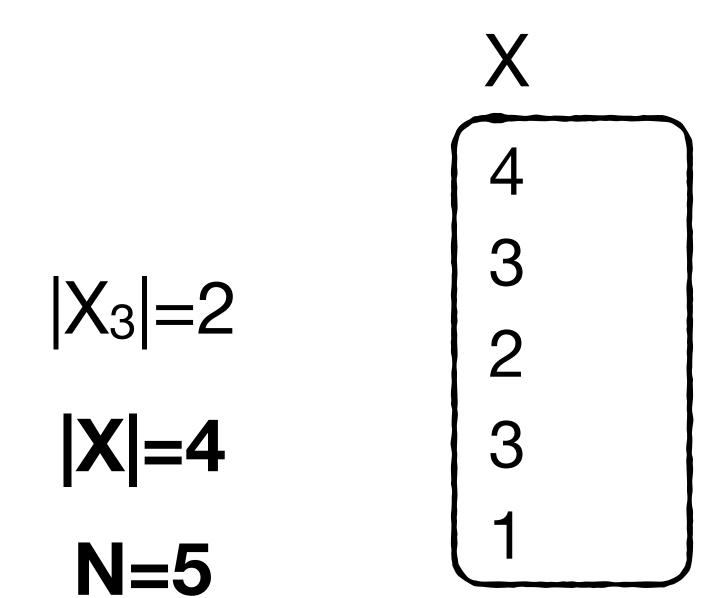
How many rows have a particular value X_i?



How many rows have a particular value Xi?

Approach 1: Estimate X_i as N / |X|

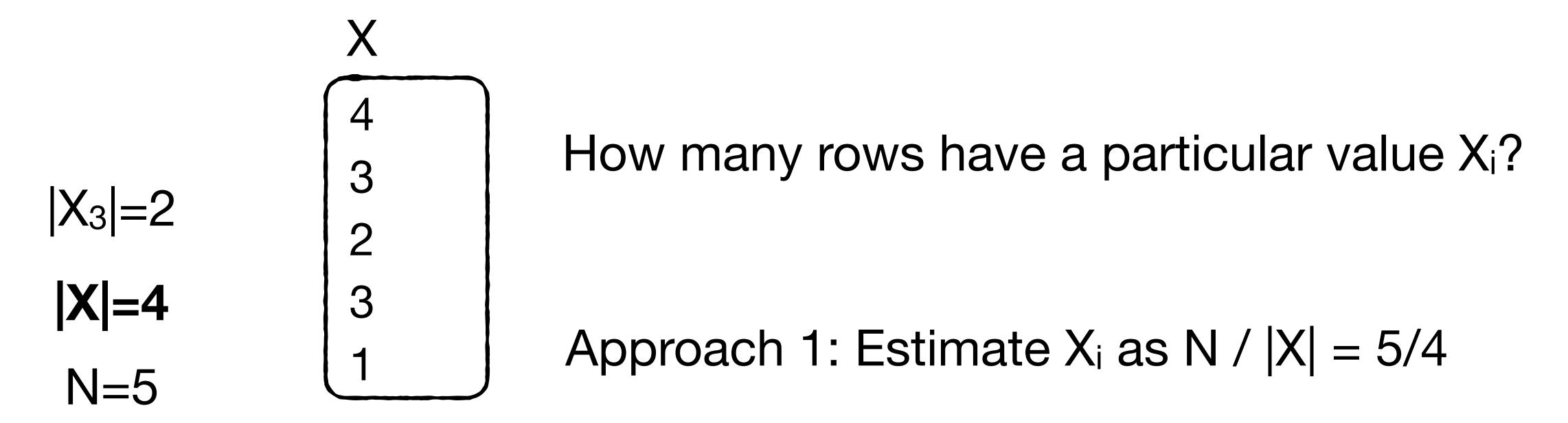
rows # unique X values



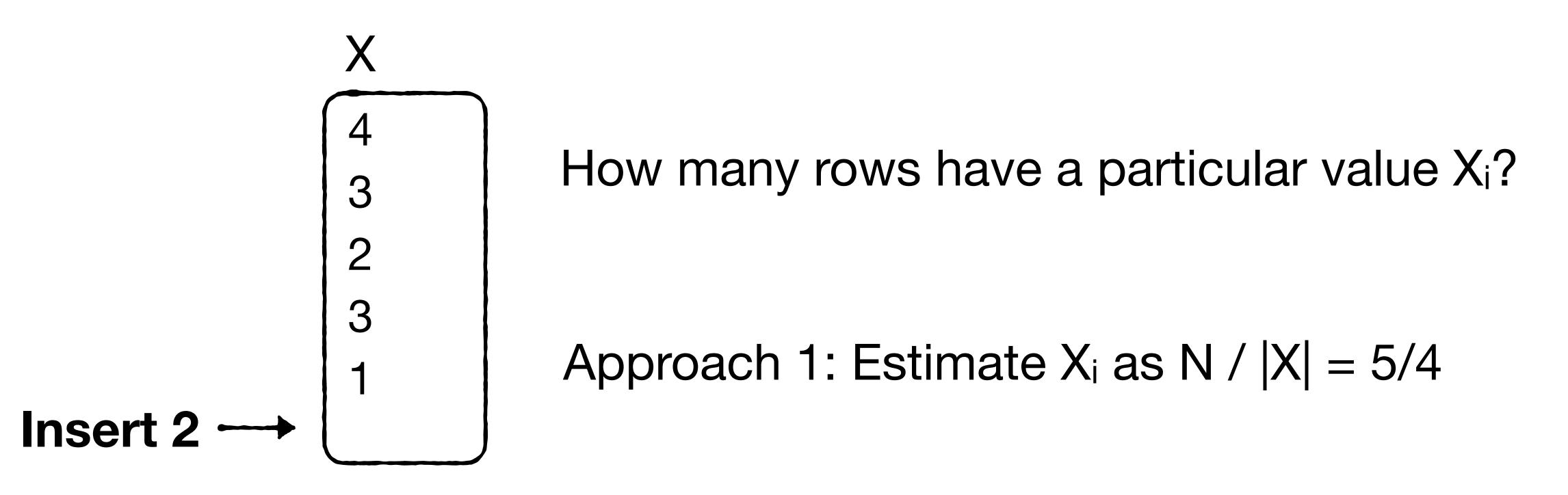
How many rows have a particular value Xi?

Approach 1: Estimate X_i as N / |X| = 5/4

rows # unique X values

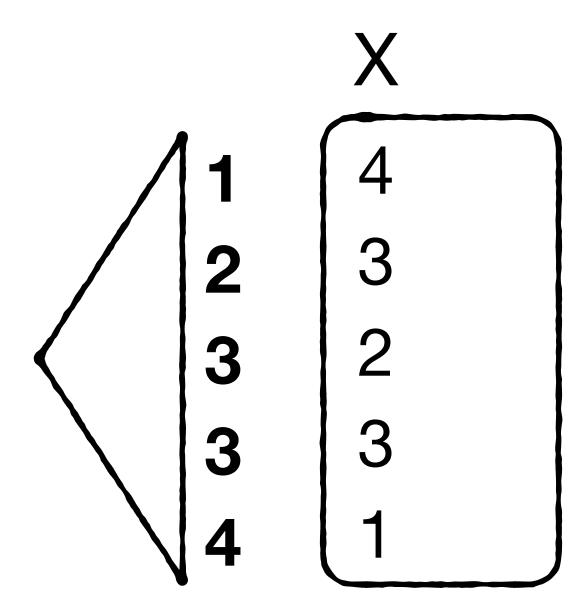


Problem: estimating |X| is non-trivial as well!



Problem: estimating |X| is non-trivial as well!

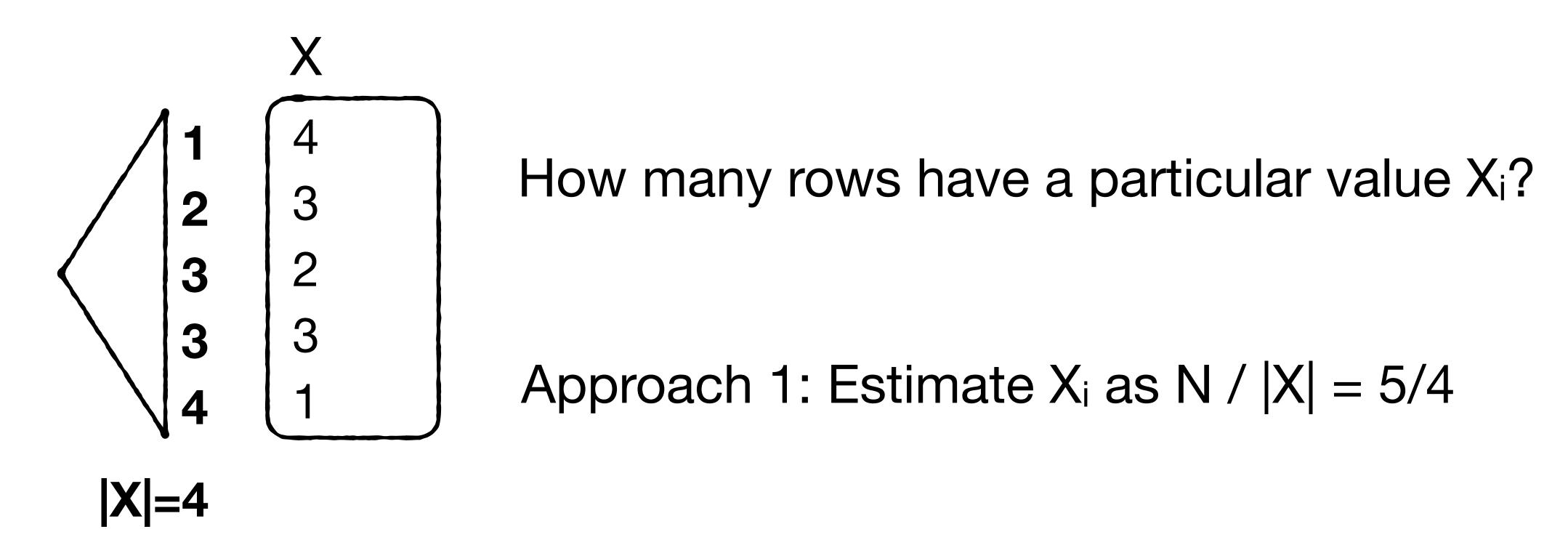
e.g., during insertion, we do not know if new value is unique or not



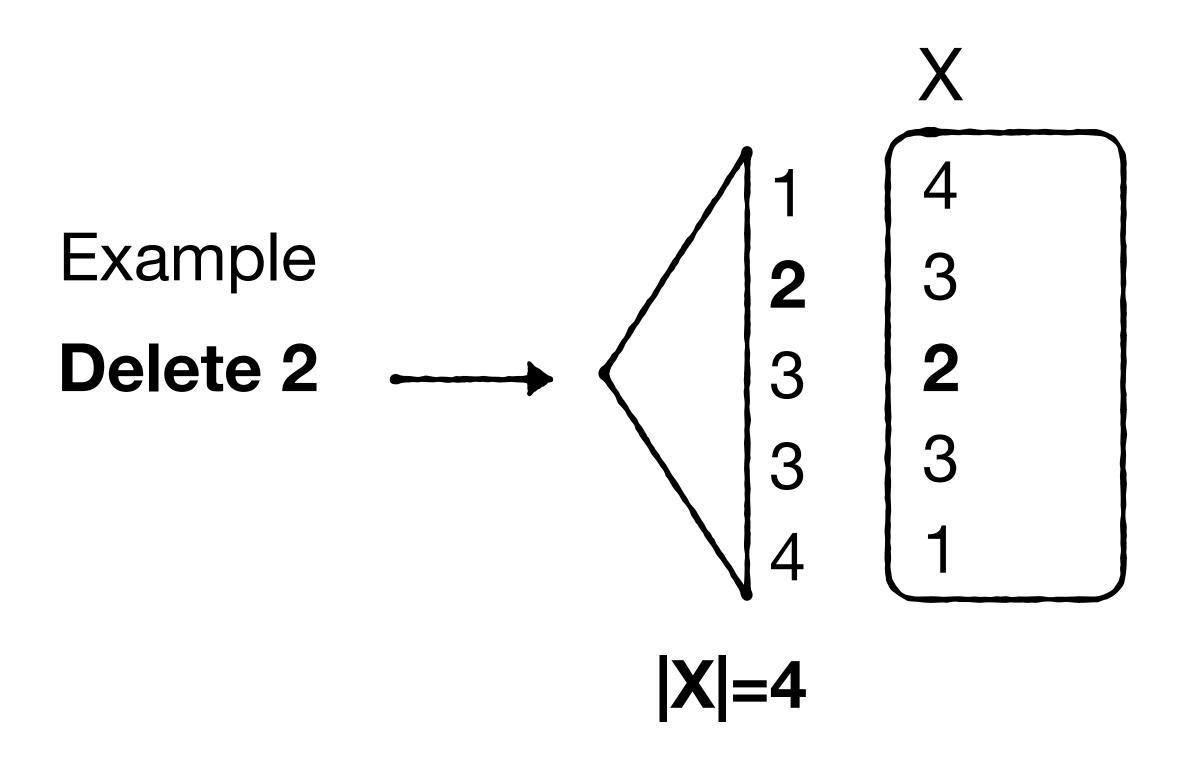
How many rows have a particular value Xi?

Approach 1: Estimate X_i as N / |X| = 5/4

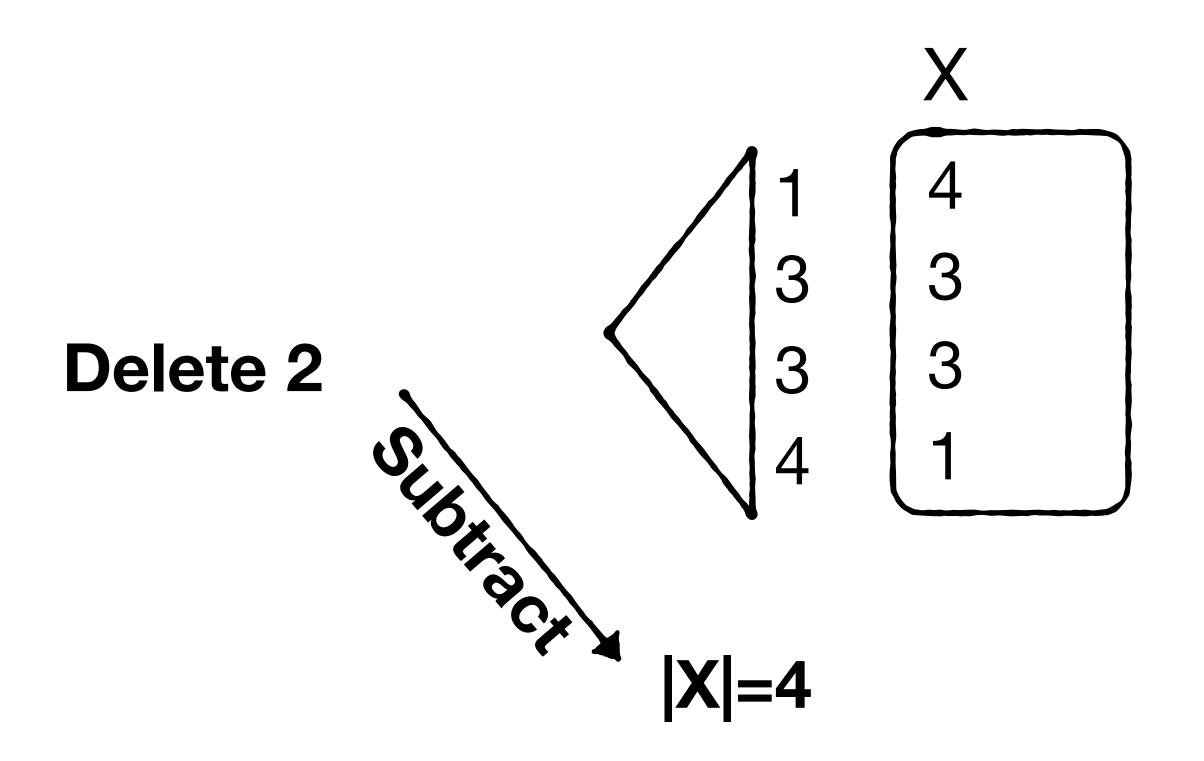
If we have an index on X, it becomes easy to tell if any new insert/delete/update adds or removes a unique value



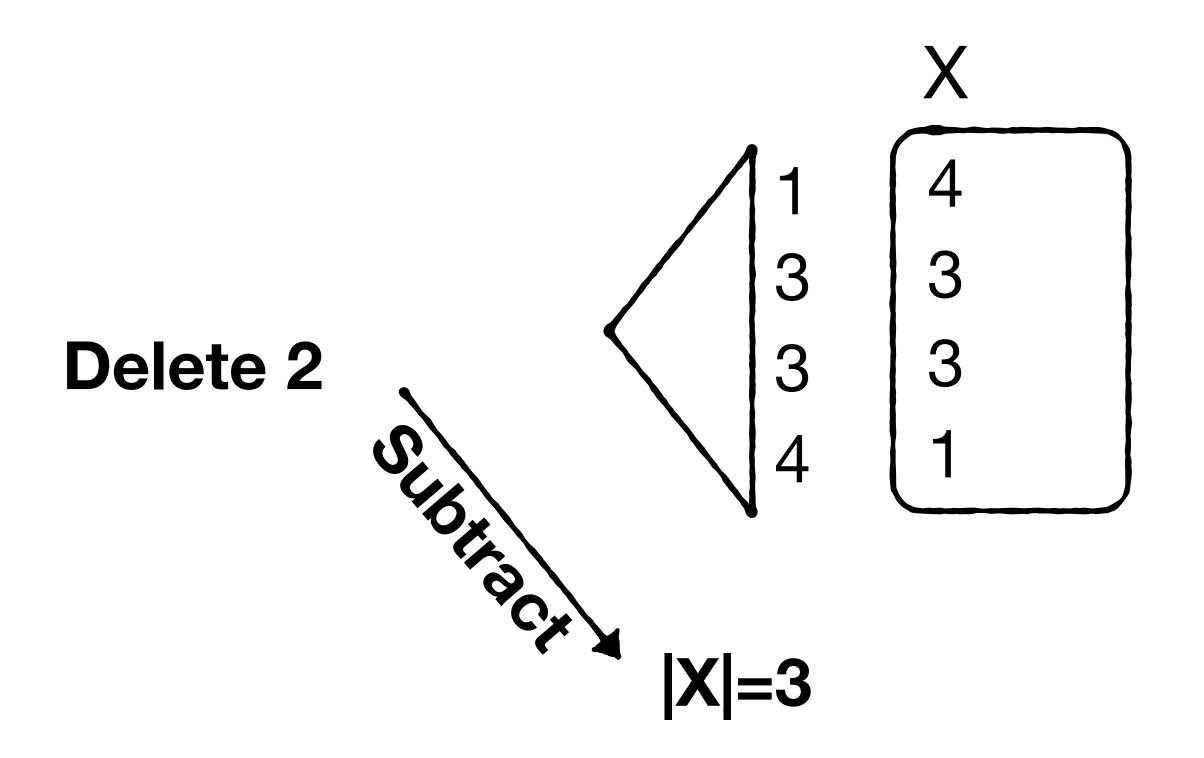
If we have an index on X, it becomes easy to tell if any new insert/delete/update adds or removes a unique value



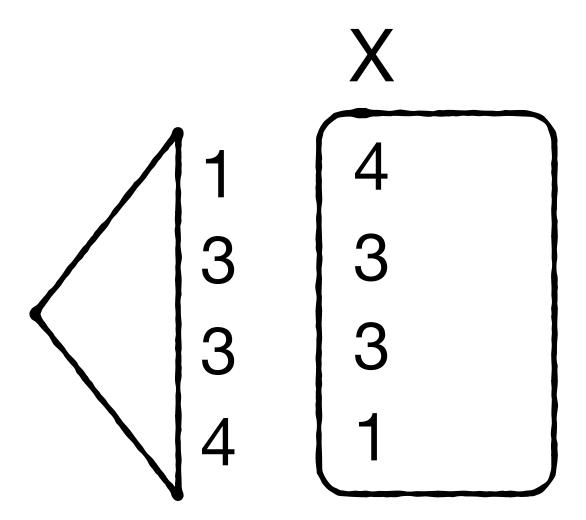
If we have an index on X, it becomes easy to tell if any new insert/delete/update adds or removes a unique value



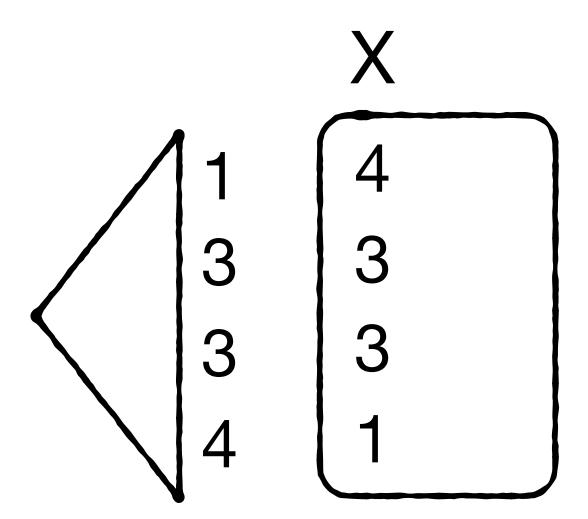
If we have an index on X, it becomes easy to tell if any new insert/delete/update adds or removes a unique value



If we have an index on X, it becomes easy to tell if any new insert/delete/update adds or removes a unique value

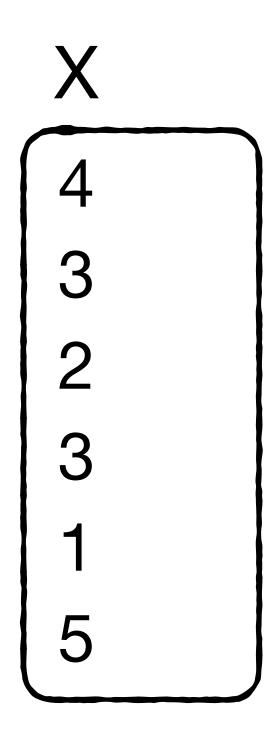


What if there is no index?



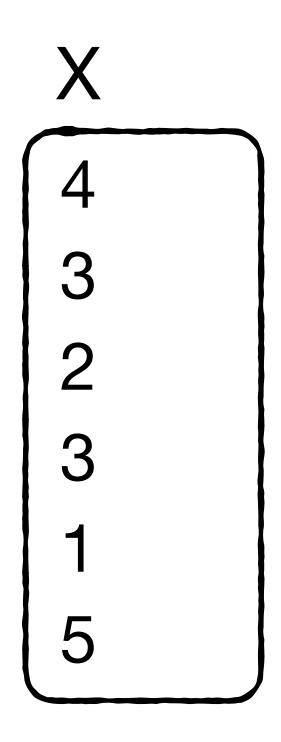
What if there is no index?

periodically scan and count #unique entries



estimate any |Xi| as N/|X|

Assumes counts of all values are uniform

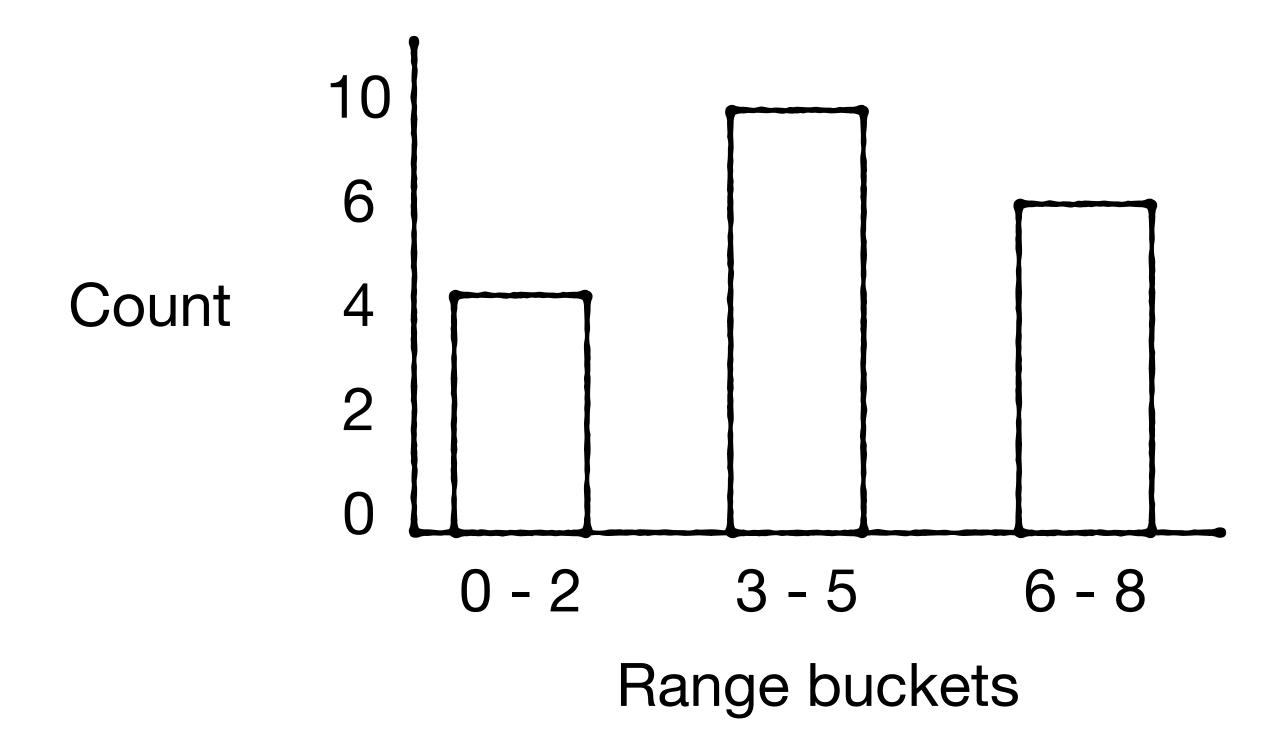


estimate any |Xi| as N/|X|

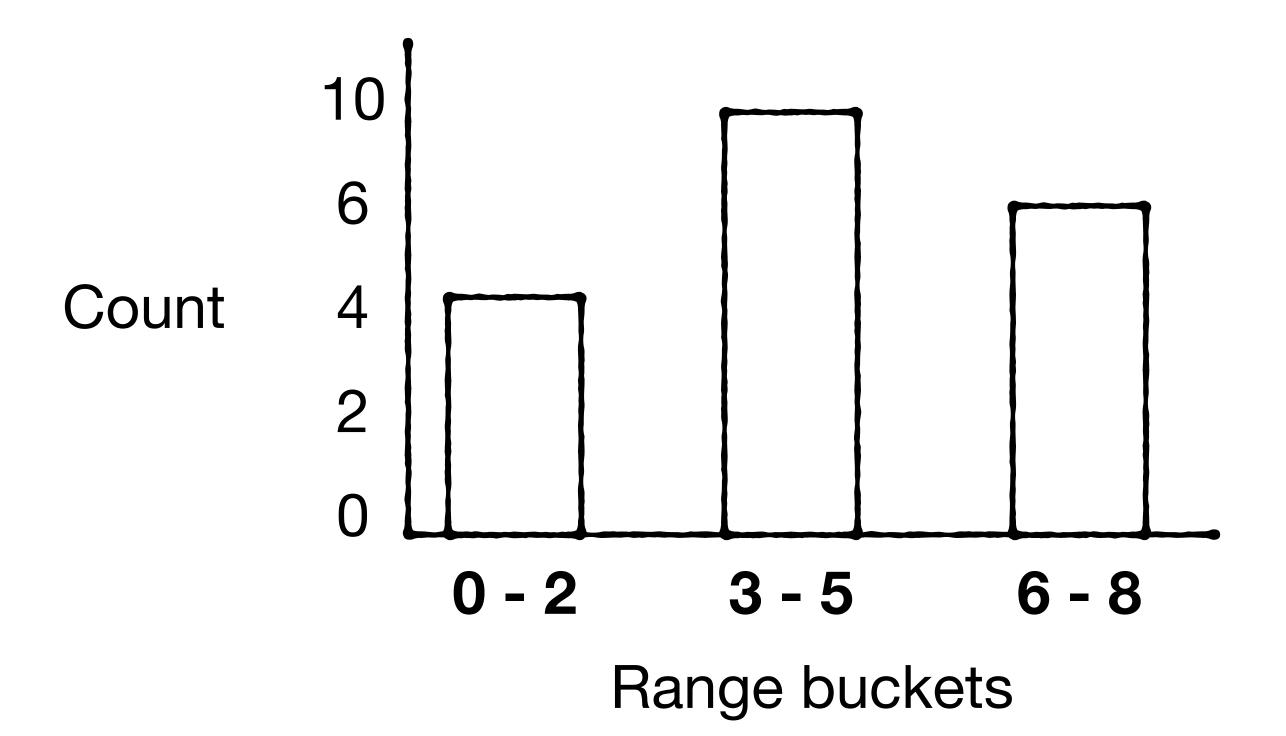
Assumes counts of all values are uniform

Can we do better?

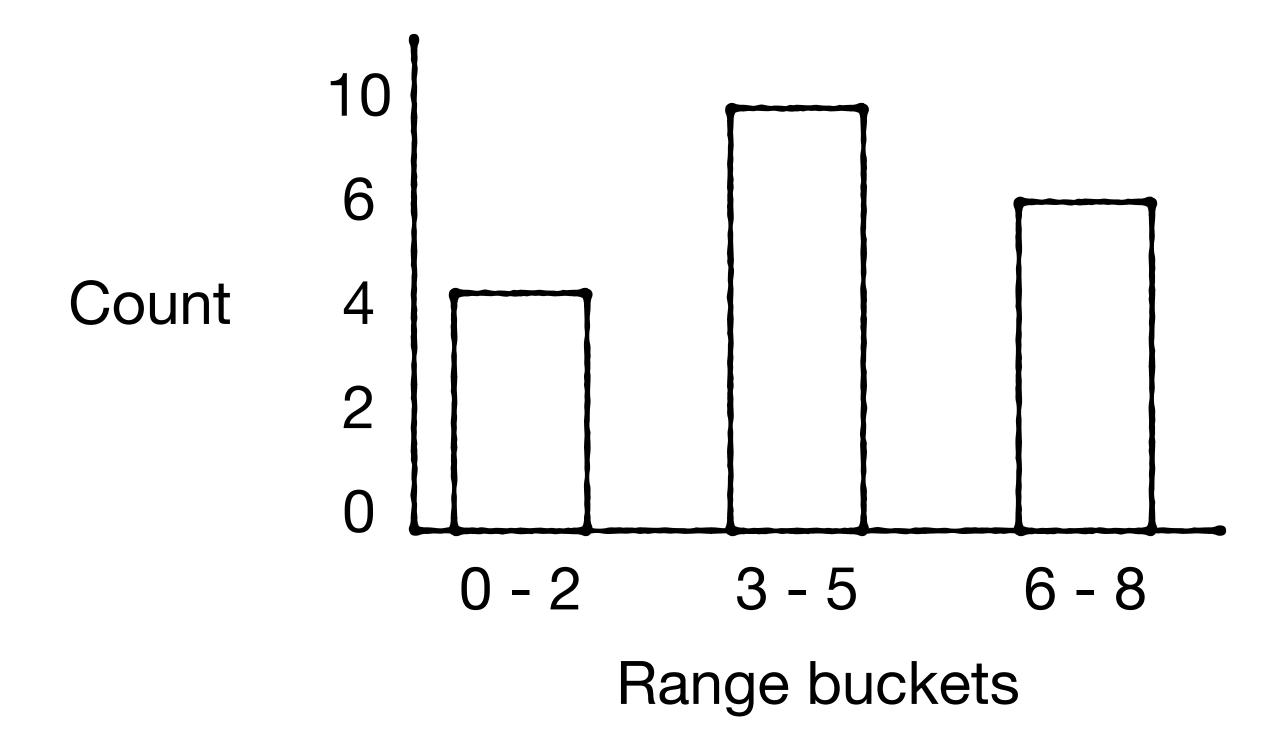
Histograms



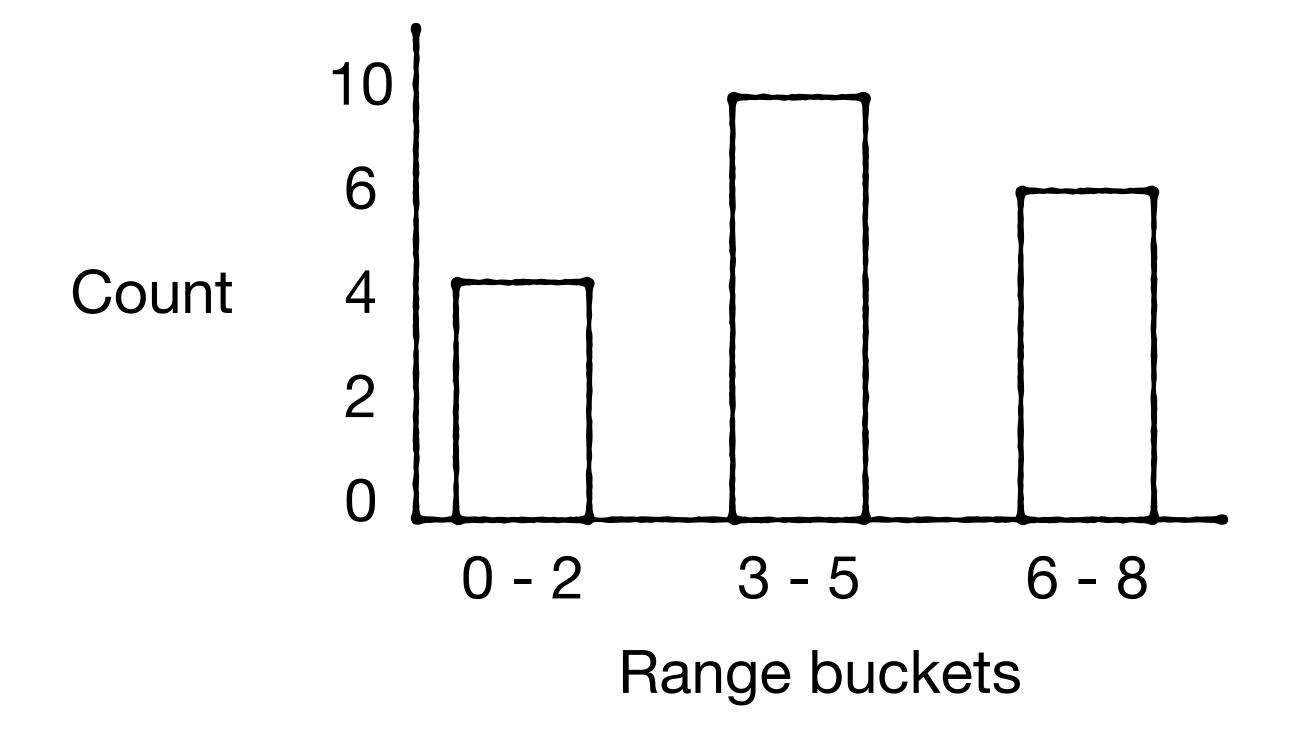
Histograms



Histograms



Estimate |Xi| as bucket count / bucket range



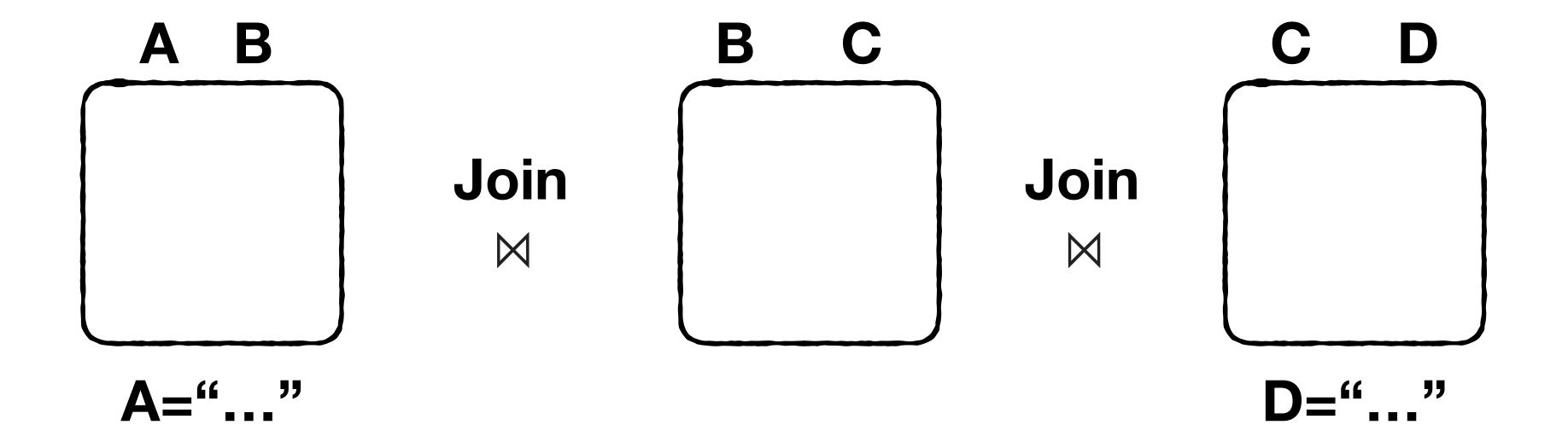
Estimate |Xi| as bucket count / bucket range

Estimate $|X_4| = as 10/3=3.3$

Query Operators

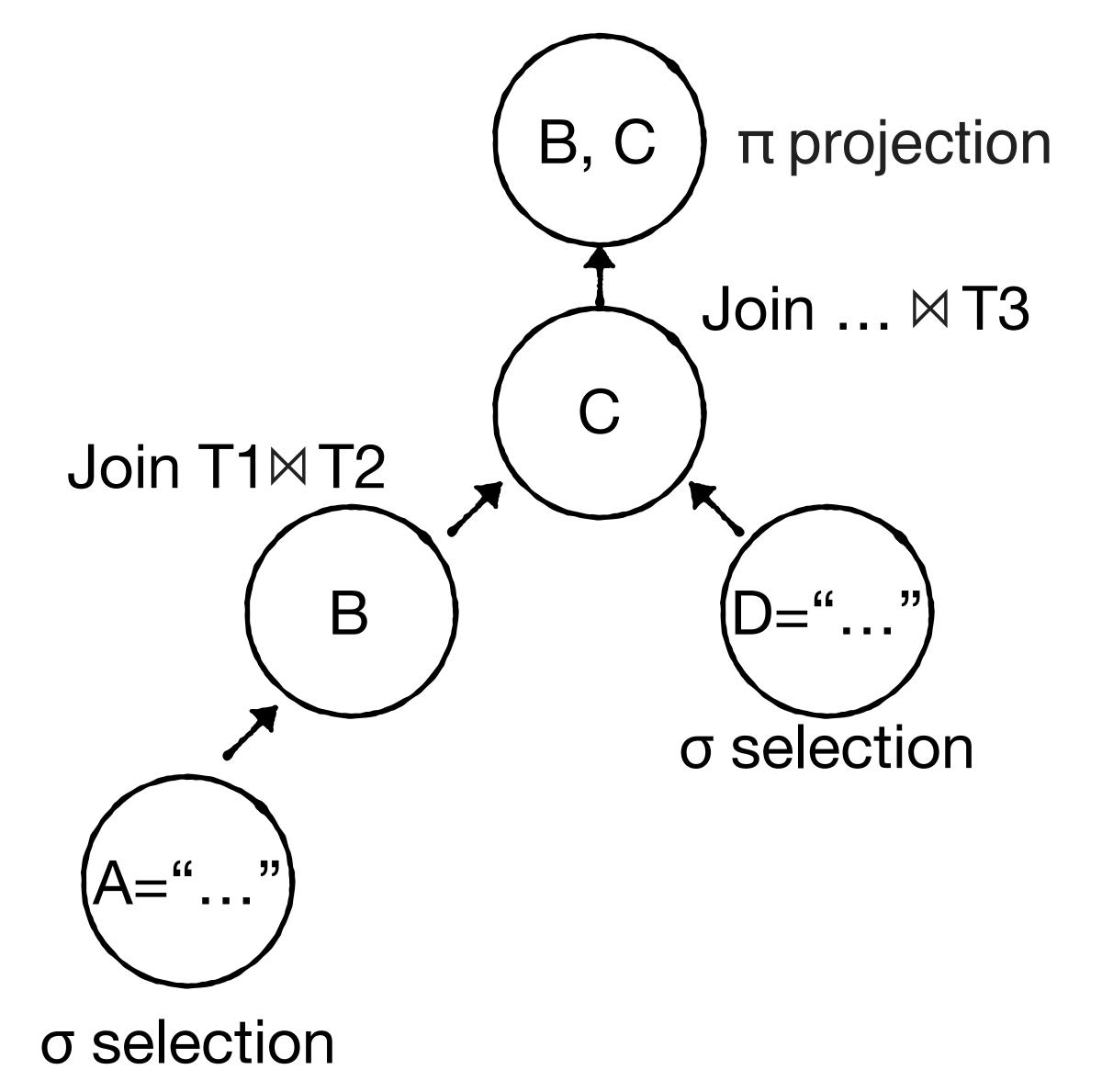
Cardinality Estimation Query Optimization

Query Optimization



Select A, D from T1, T2, T3 where T1.B = T2.B and T2.C = T3.C and A="..." and D="..."

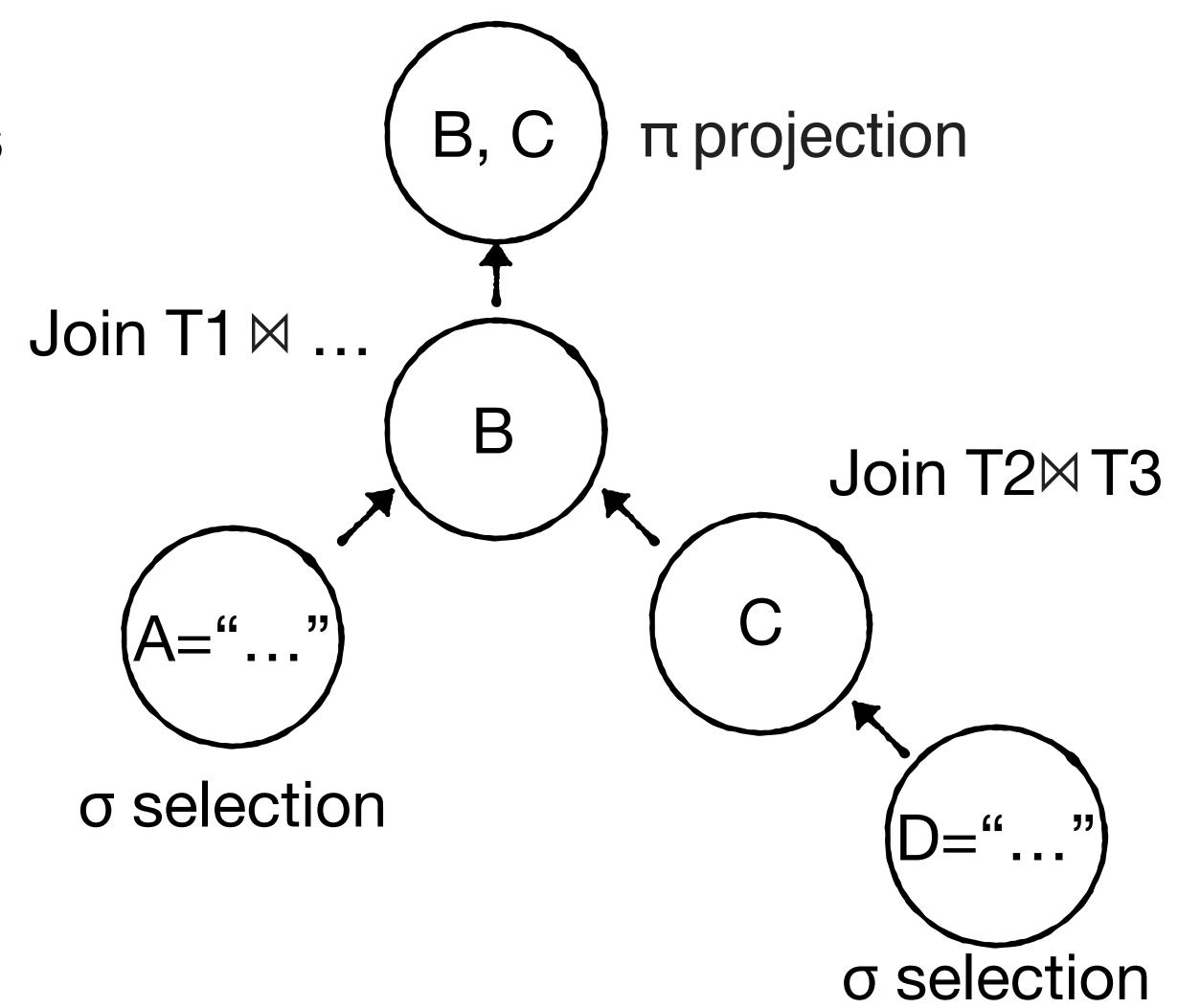
Select A, D from T1, T2, T3 where T1.B = T2.B and T2.C = T3.C and A="..." and D="..."



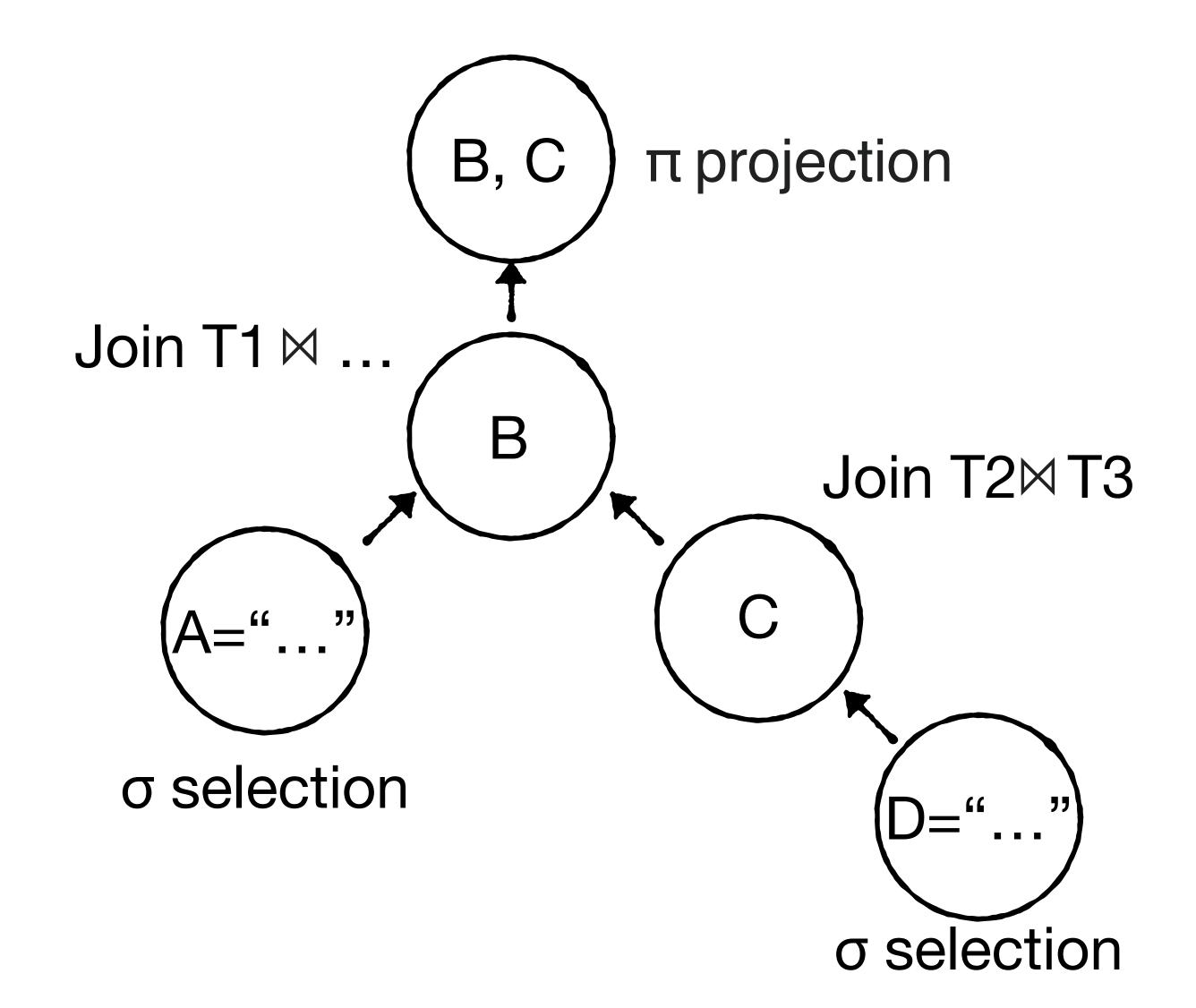
πprojection Can join tables in different orders Join ... ⋈ T3 C Join T1⋈T2 B σ selection

σ selection

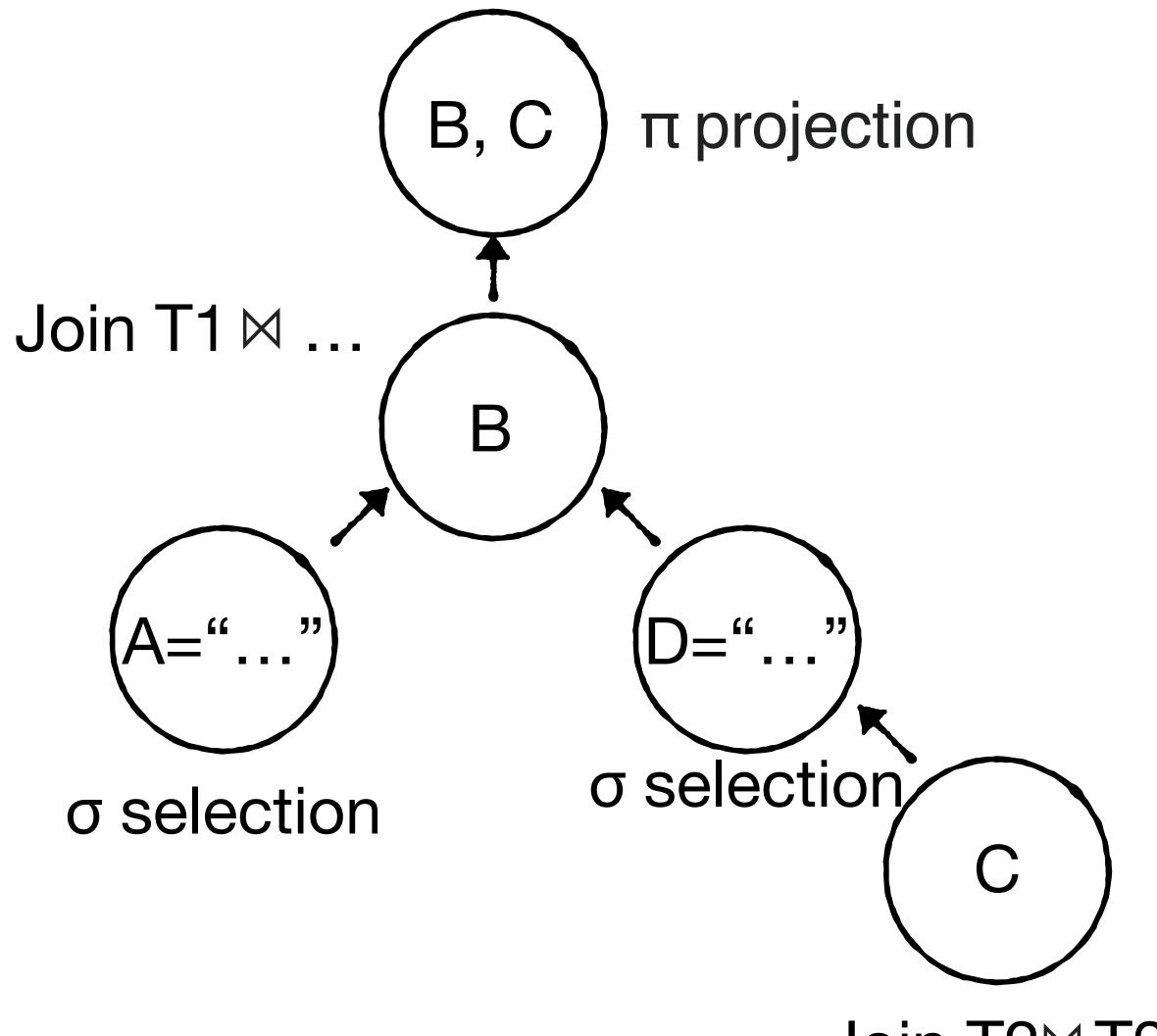
Can join tables in different orders



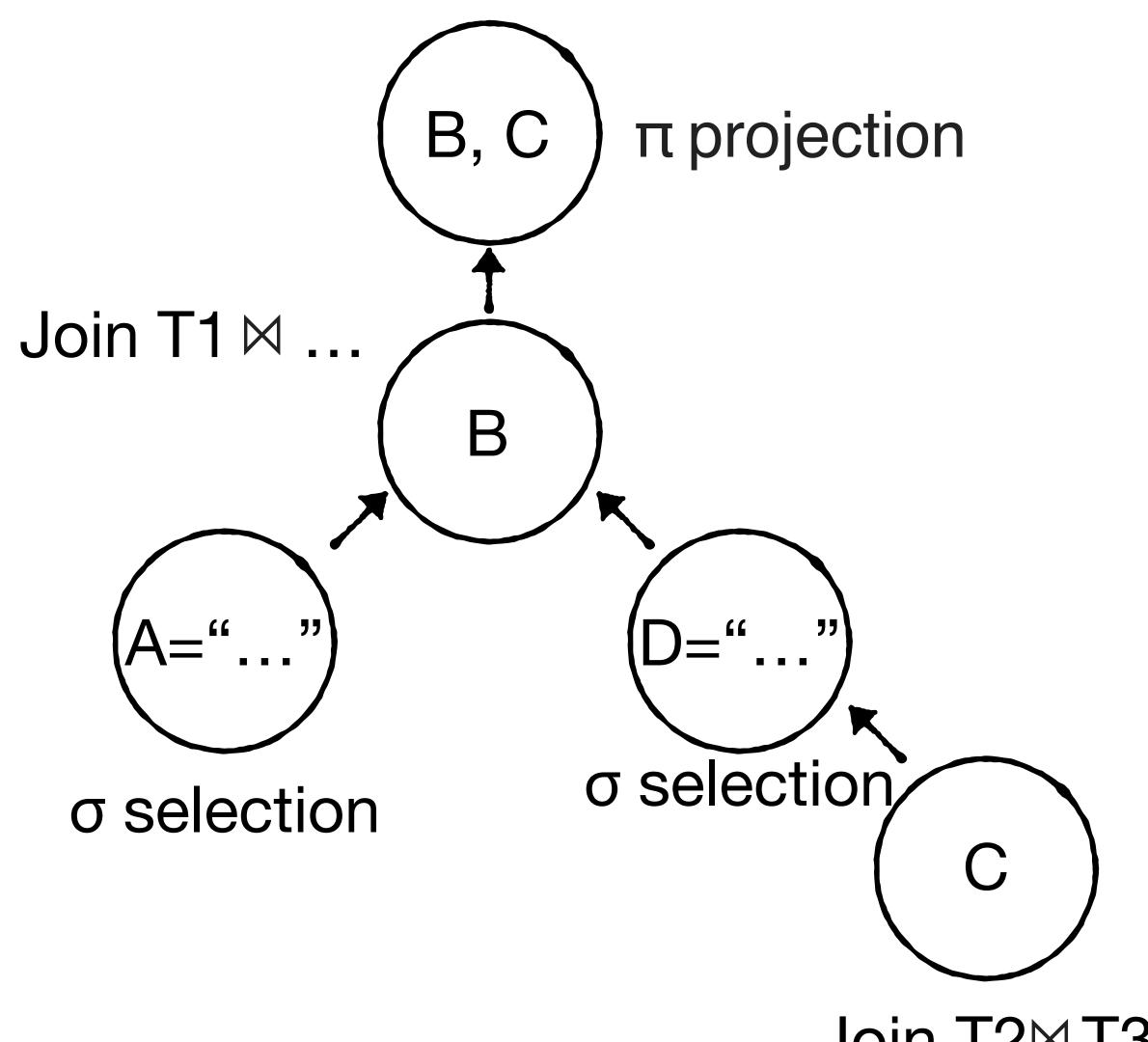
Can select in different orders

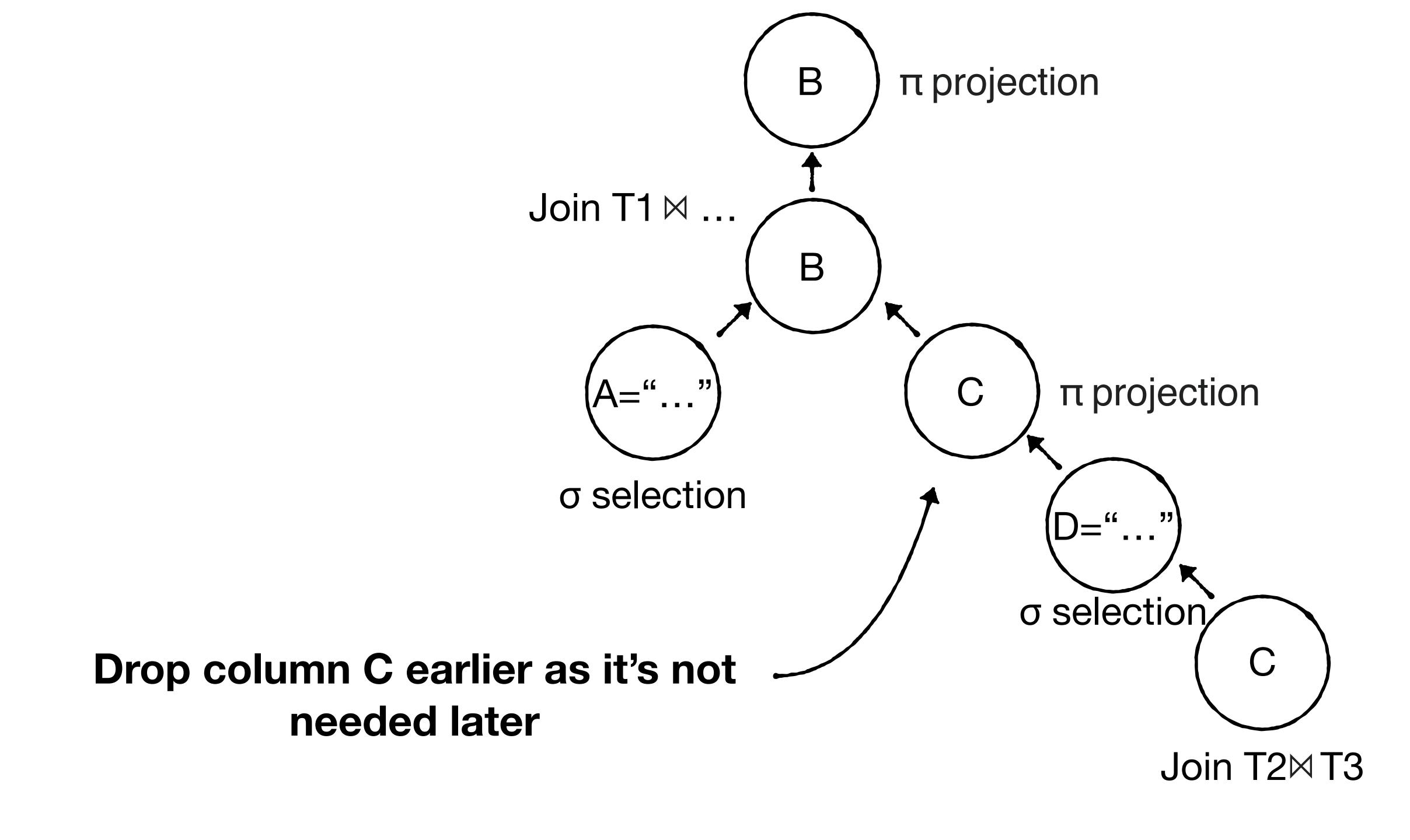


Can select in different orders

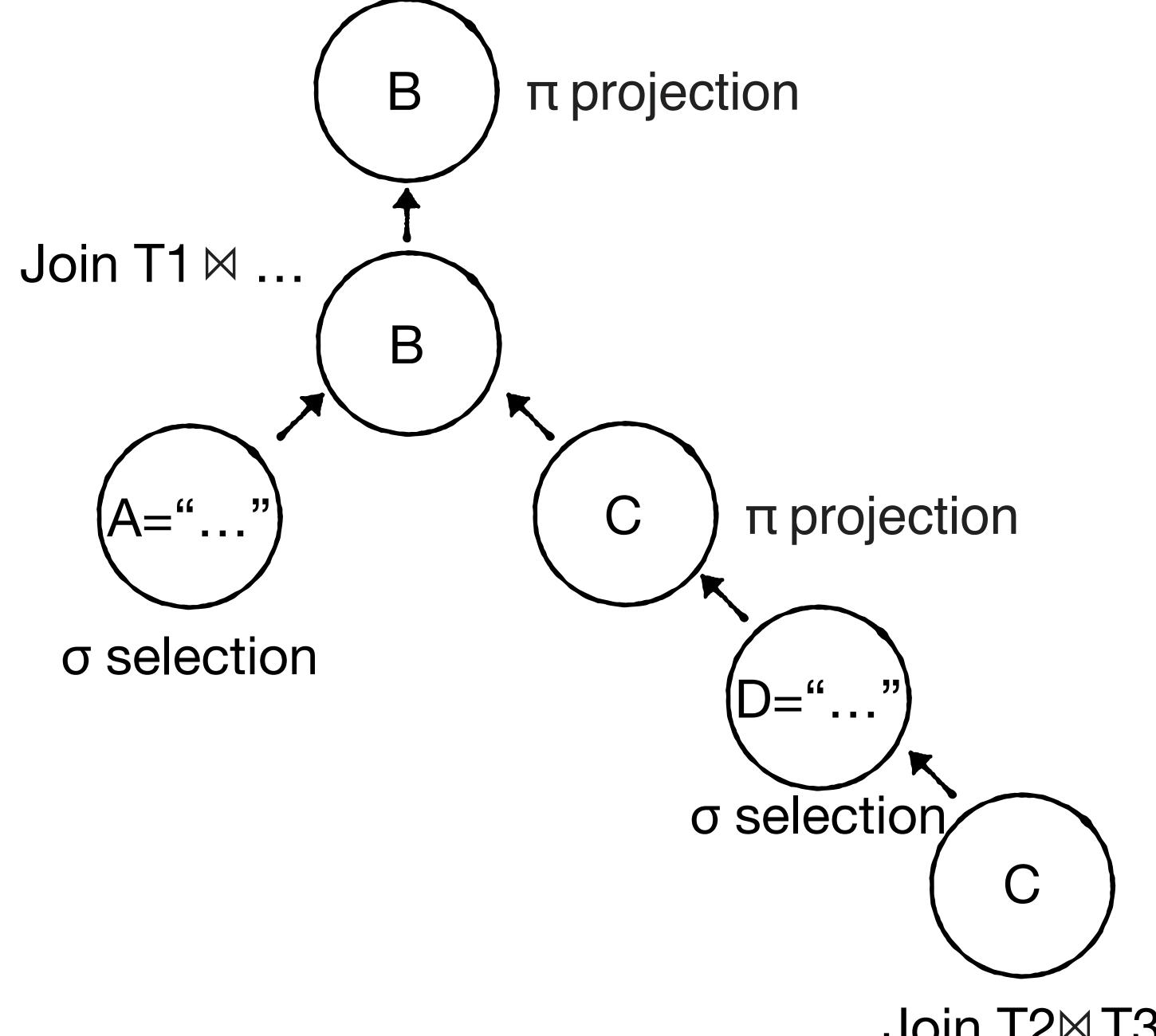


Can project columns in different orders as long as selections or joins do not rely on them



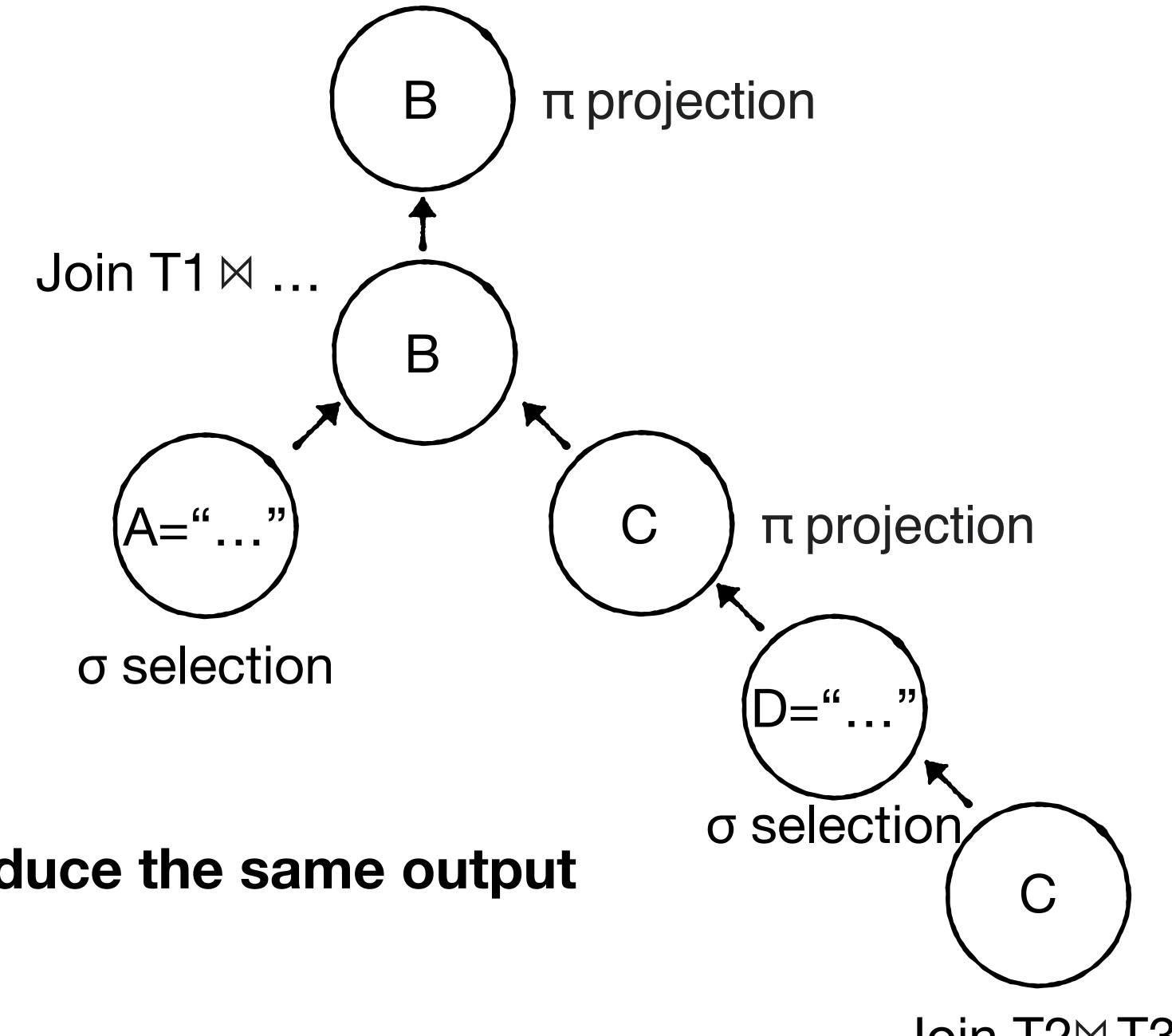


Can implement each operator in different ways

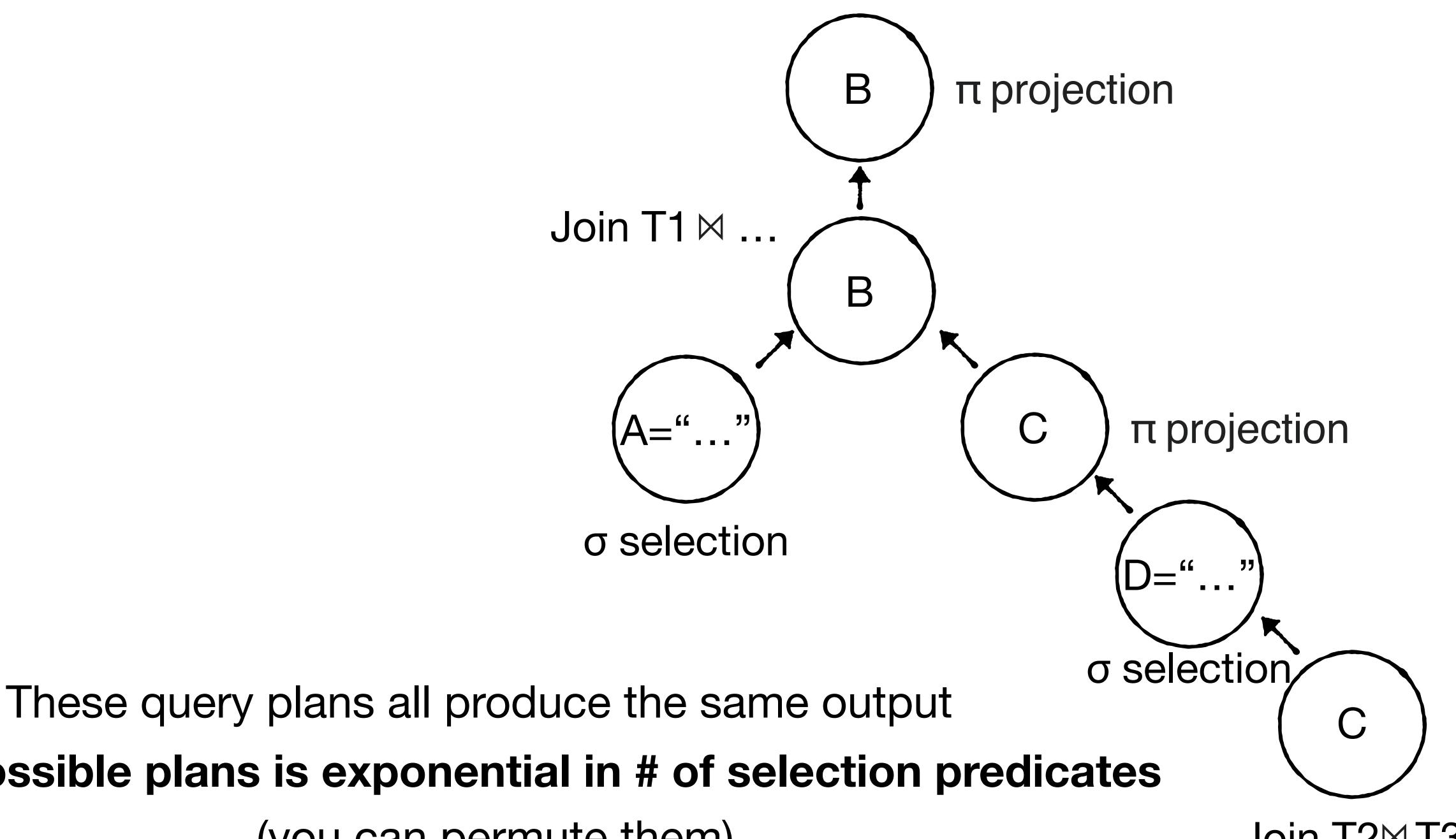


Can implement each πprojection B operator in different ways Join T1 ⋈ B Scan vs. index? π projection σ selection σ selection

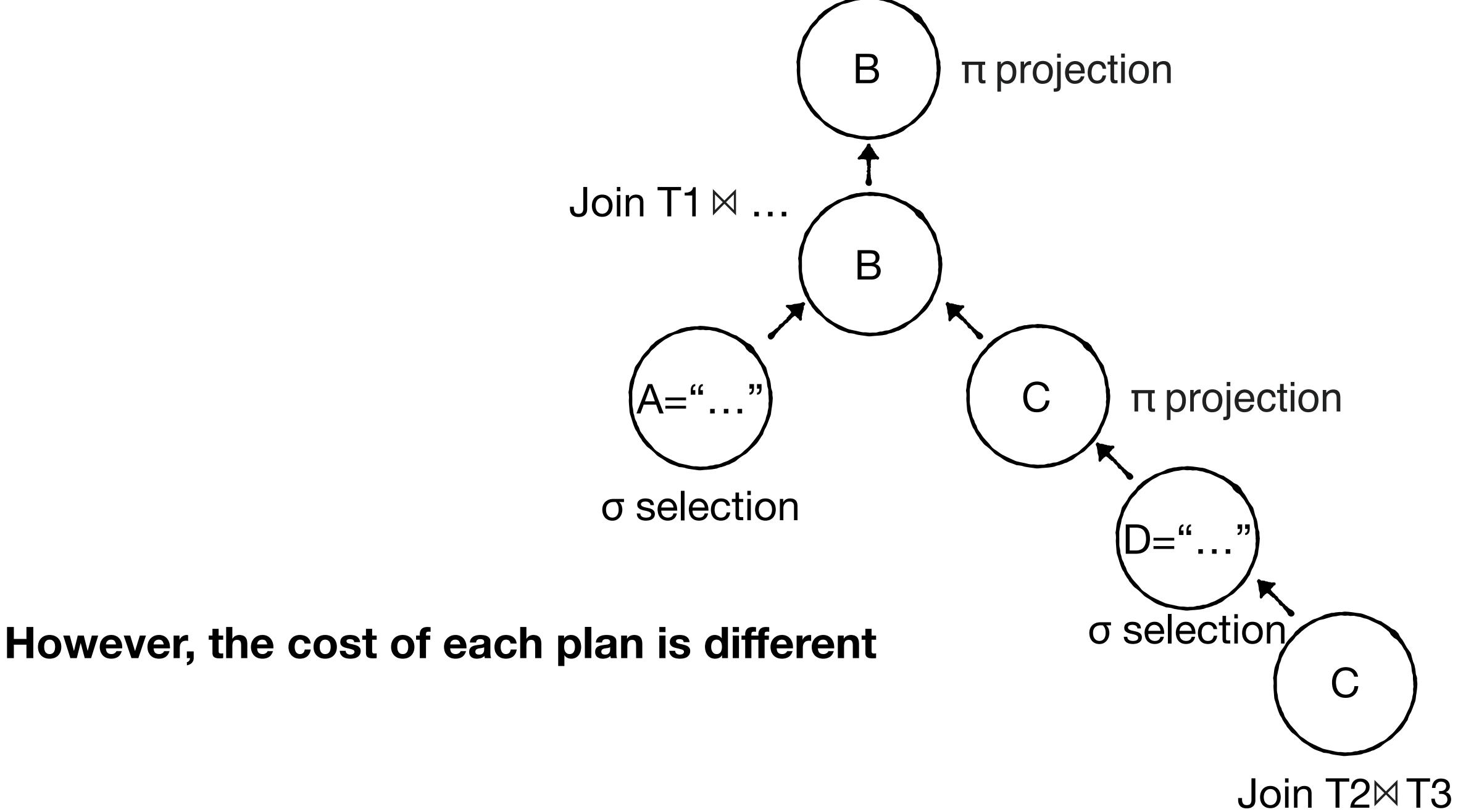
Can implement each πprojection B operator in different ways Join T1⋈ B π projection Block nested loop? Sort-merge? σ selection **Grace Hash?** Index? σ selection

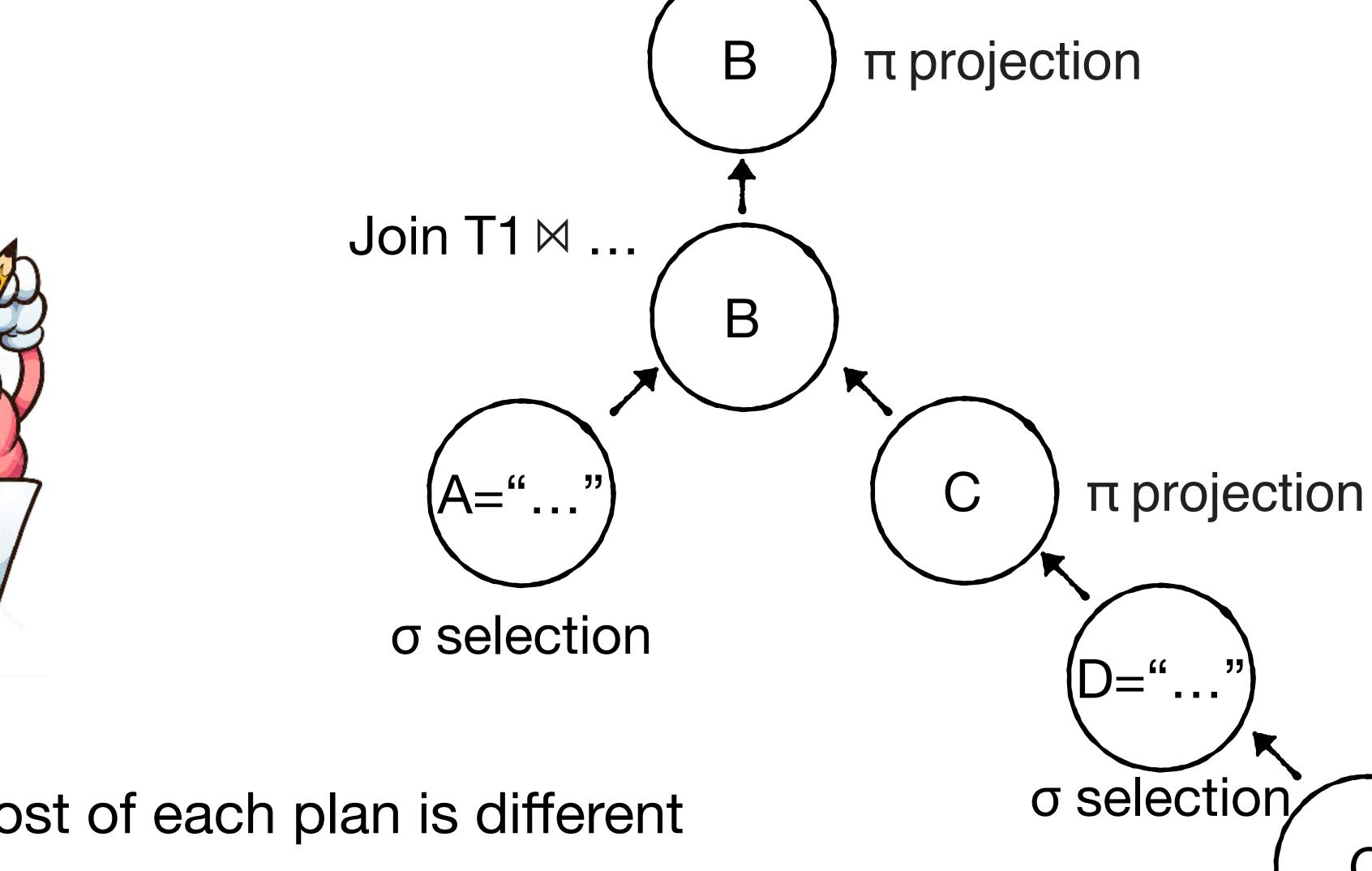


These query plans all produce the same output



possible plans is exponential in # of selection predicates (you can permute them)





However, the cost of each plan is different

The optimizer searches the space of possible plans to find a good one (not best but good enough)

