



# MonetDB, Cracking and recycling

Martin Kersten

*CWI*

*Amsterdam*



M.Kersten 2008



*Try to maximize performance*

Present

Paste

Materialized  
Views

Potency

B-tree,  
Hash  
Indices

Cracking

MONETDB



M.Kersten 2008



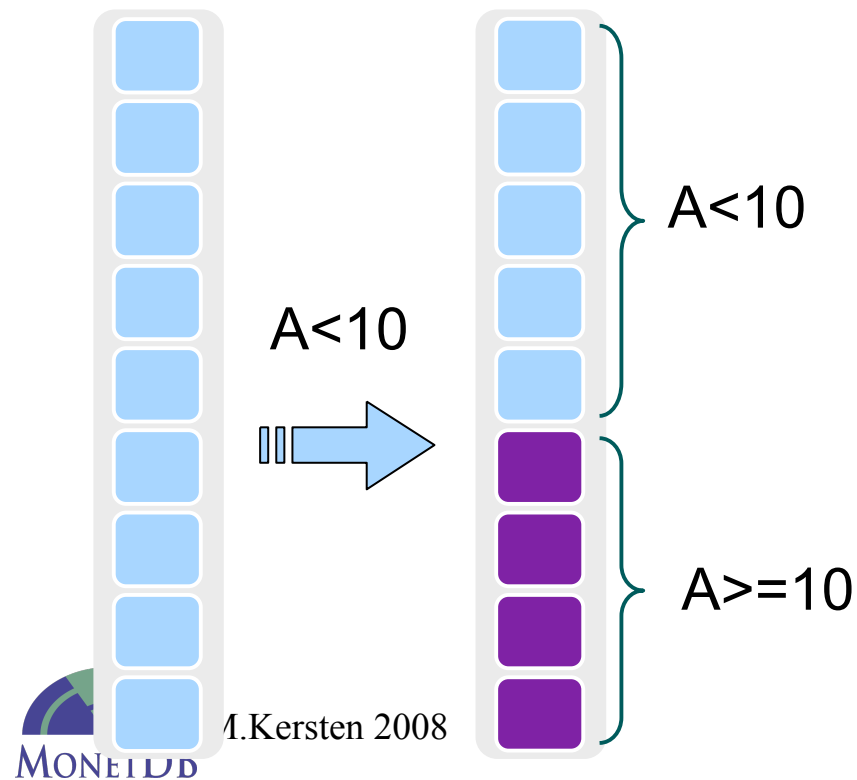
*Find a trusted fortune teller*

- Indices in database systems focus on:
  - All tuples are equally important for fast retrieval
  - There are ample resources to maintain indices
- MonetDB cracks the database into pieces based on actual query load

# Cracking algorithms

Physical reorganization happens per column based on selection predicates.

Split a piece of a column in **two** new pieces

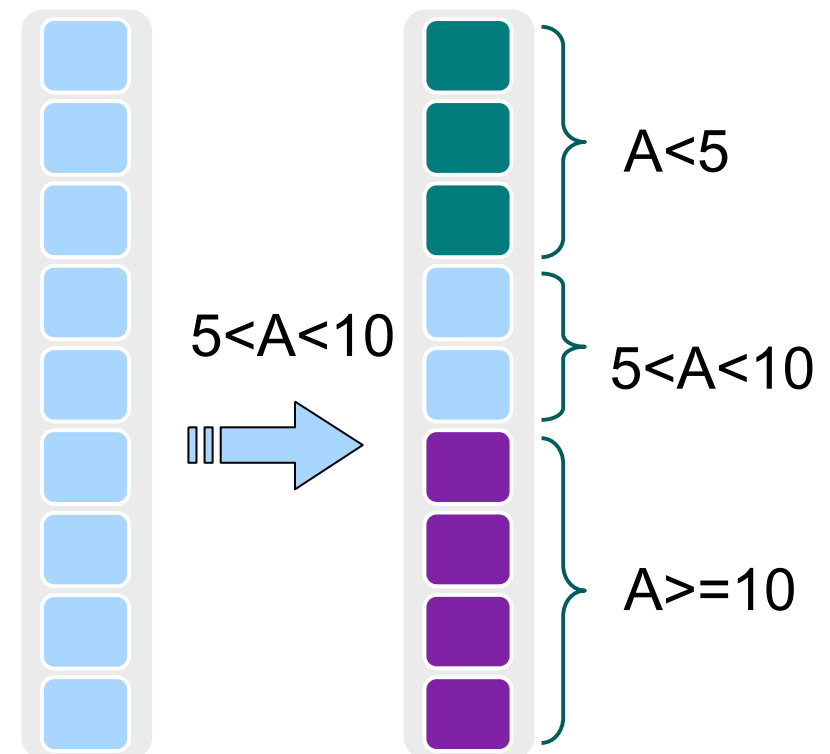
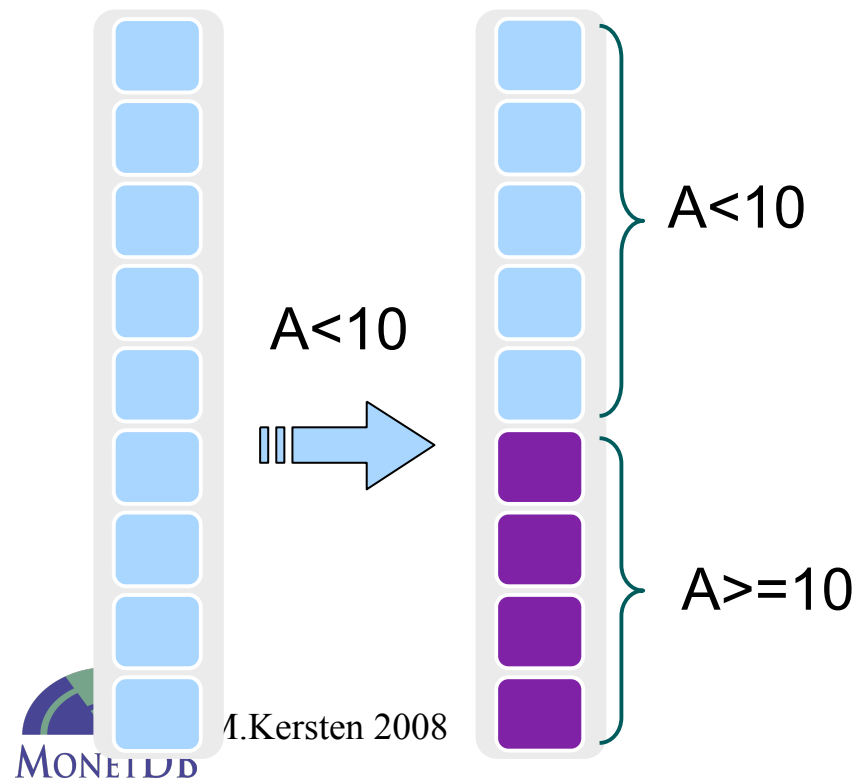


# Cracking algorithms

Physical reorganization happens per column

Split a piece of a column  
in **two** new pieces

Split a piece of a column  
in **three** new pieces





# Cracking example

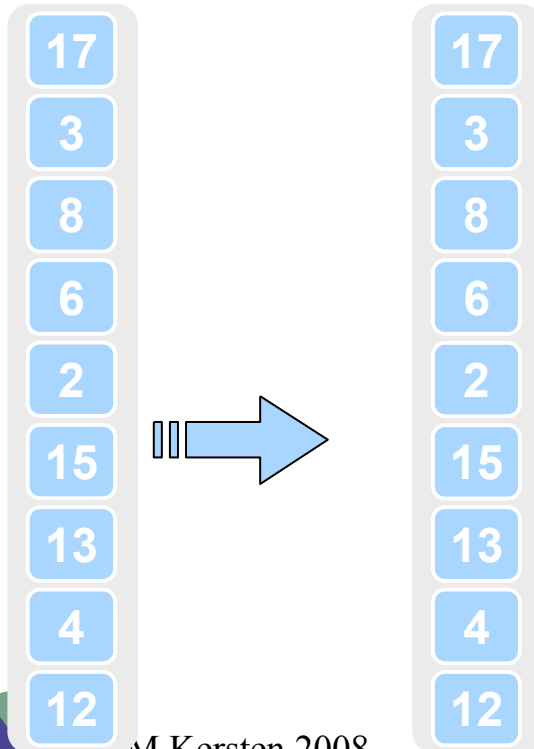
select A>5 and A<10

- 17
- 3
- 8
- 6
- 2
- 12
- 13
- 4
- 15



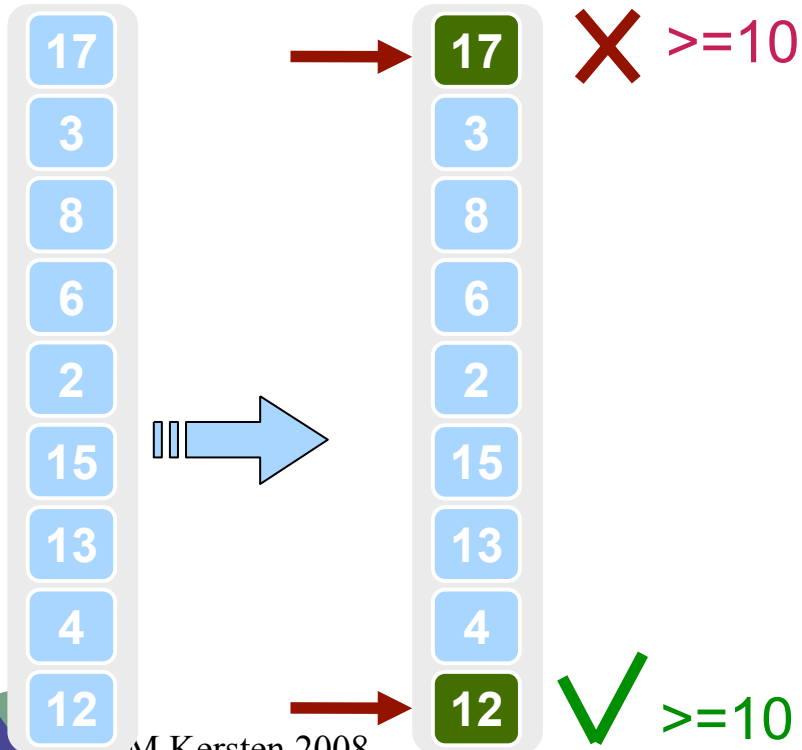
# Cracking example

select A>5 and A<10



# Cracking example

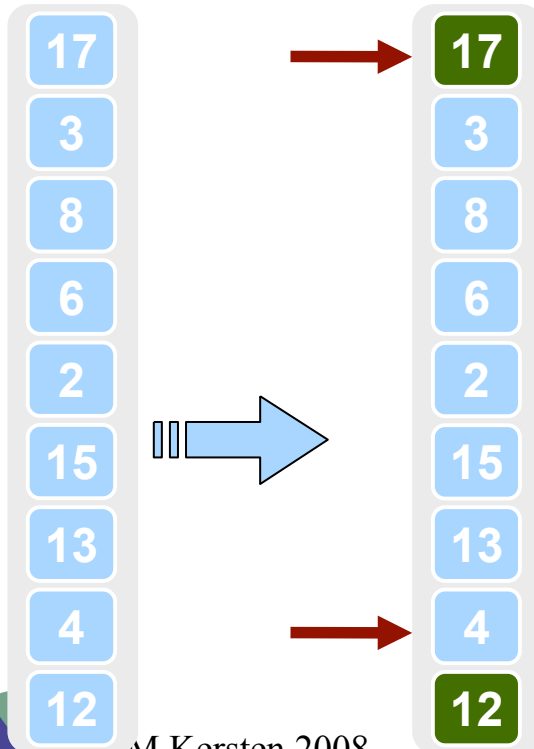
select A>5 and A<10





# Cracking example

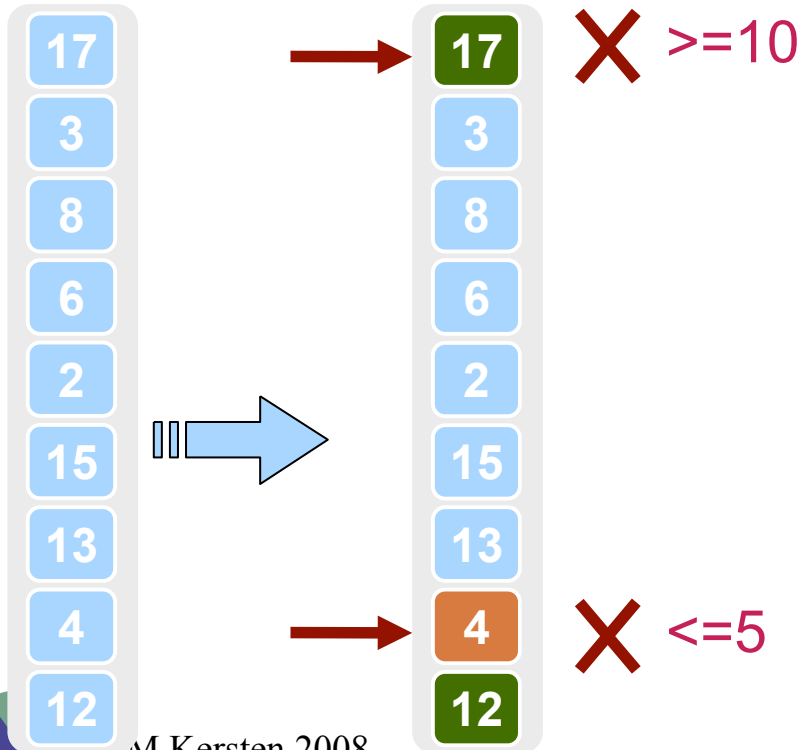
select A>5 and A<10



✗  $\geq 10$

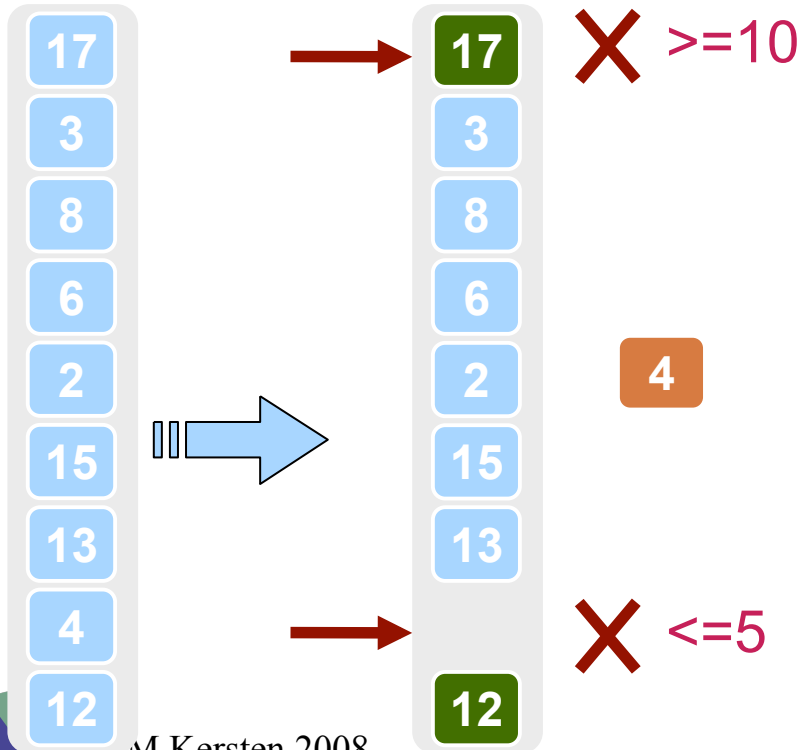
# Cracking example

select A>5 and A<10



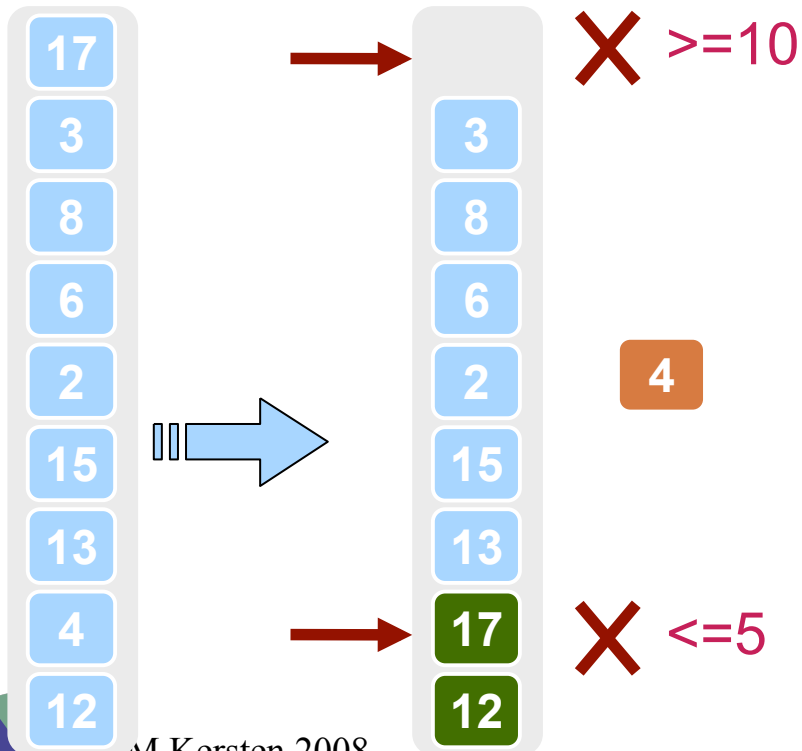
# Cracking example

select A>5 and A<10



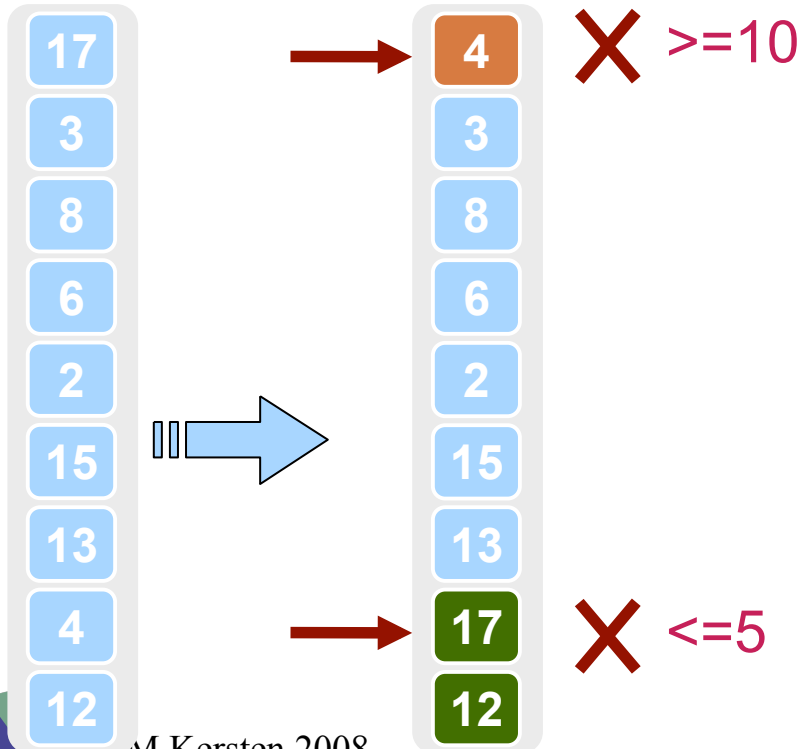
# Cracking example

select A>5 and A<10



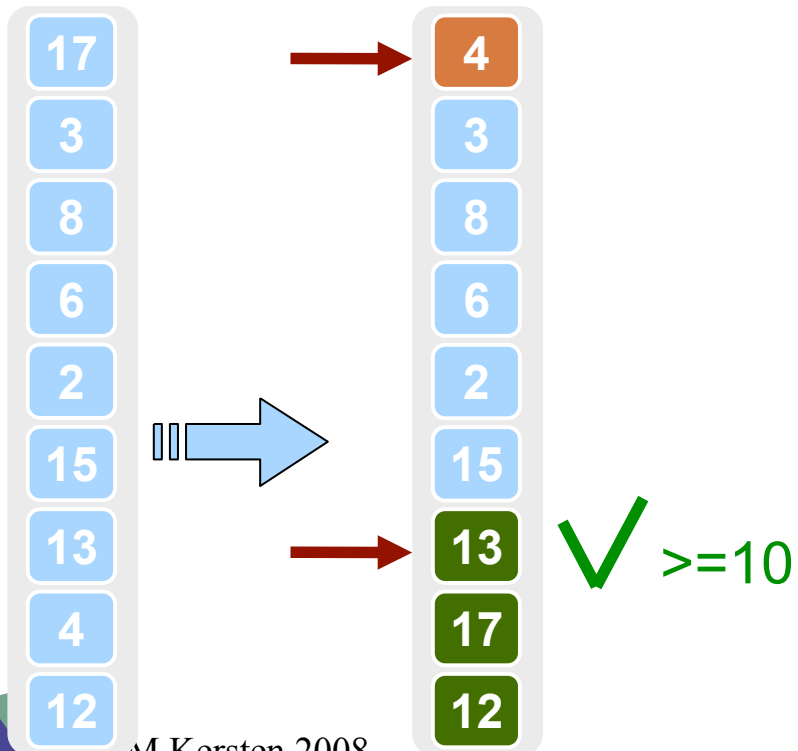
# Cracking example

select A>5 and A<10



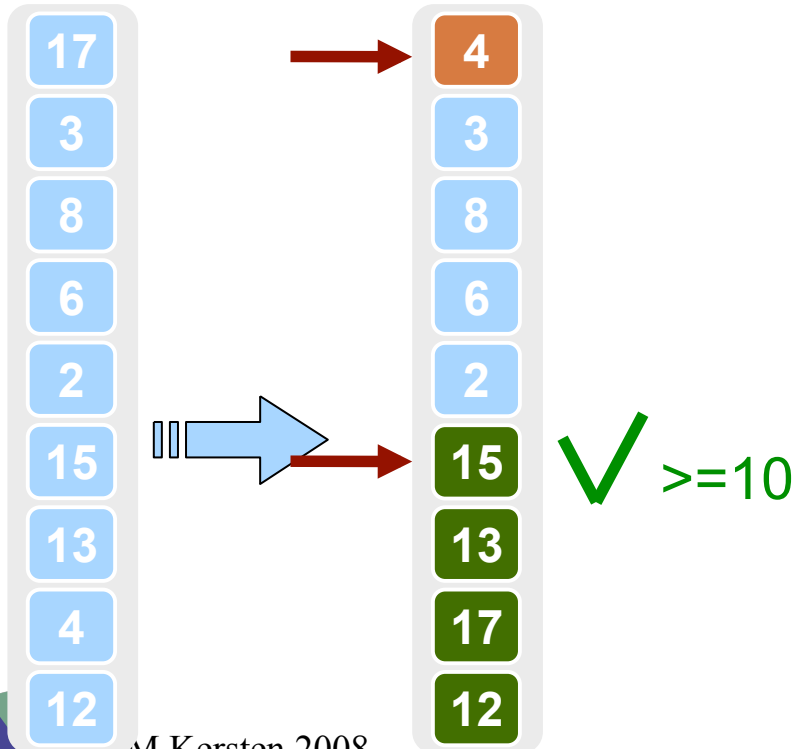
# Cracking example

select A>5 and A<10



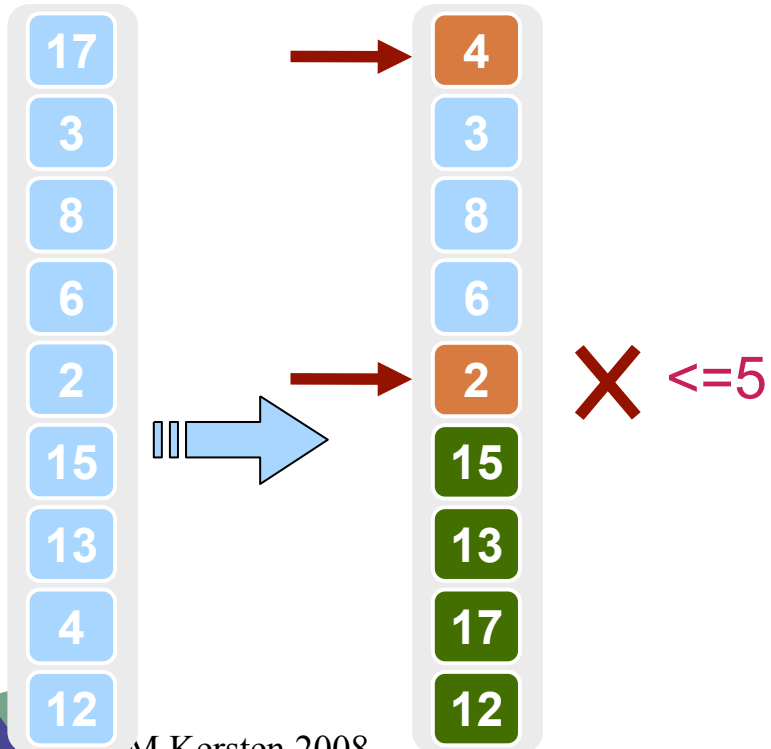
# Cracking example

select A>5 and A<10



# Cracking example

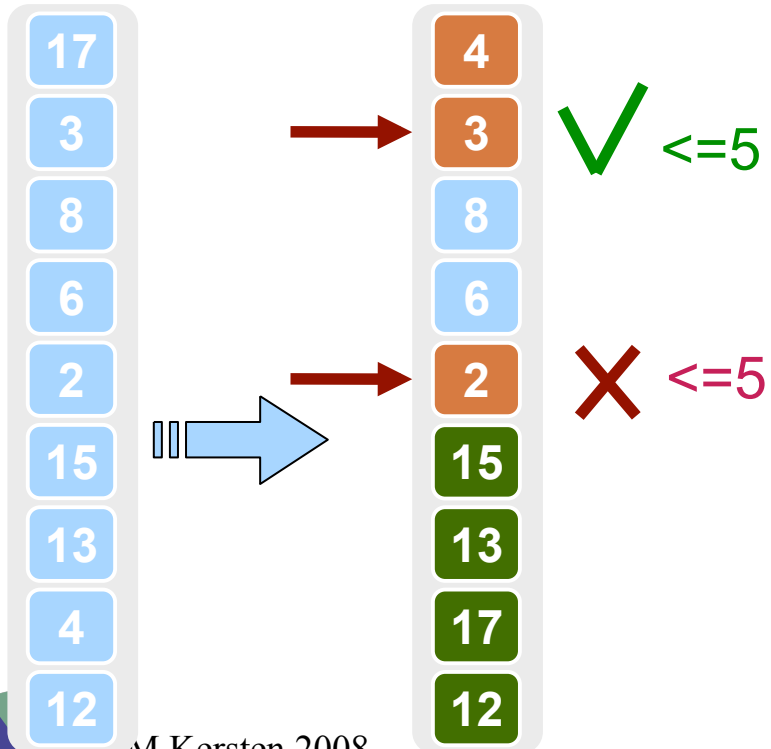
select A>5 and A<10





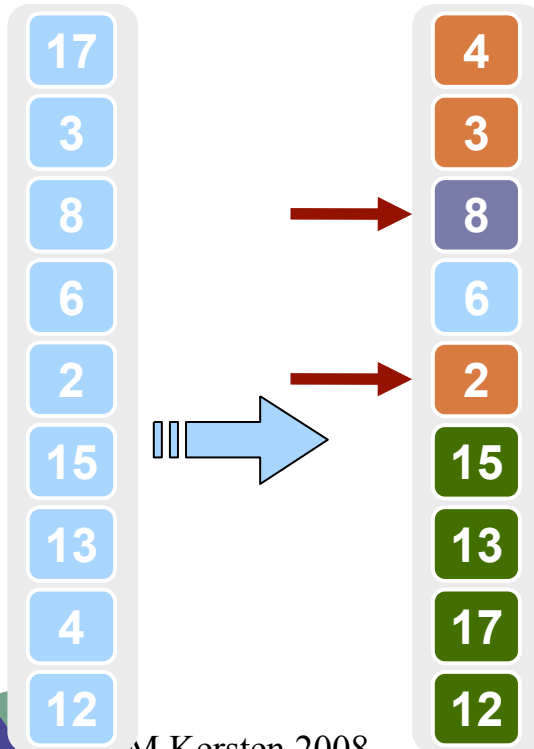
# Cracking example

select A>5 and A<10



# Cracking example

select A>5 and A<10

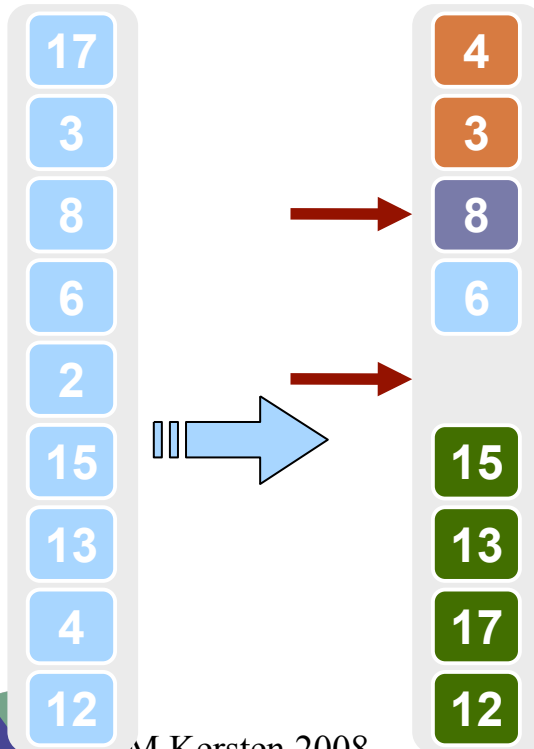


✗ >5 and <10

✗ <=5

# Cracking example

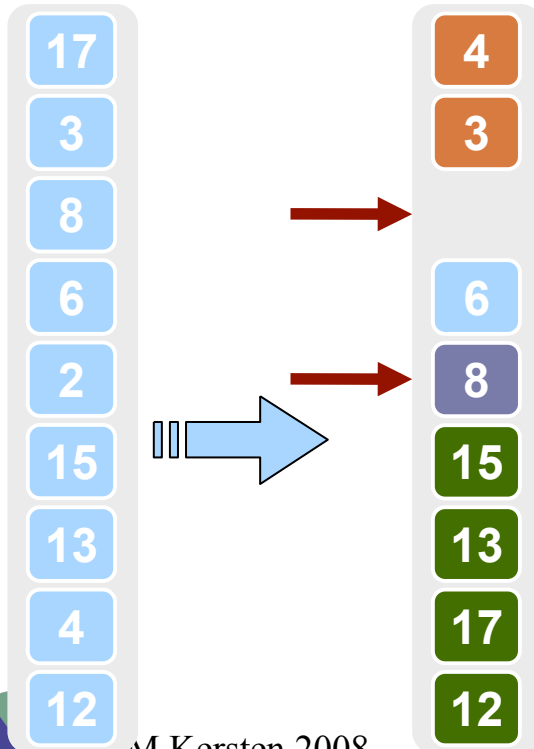
select A>5 and A<10



~~X~~ >5 and <10  
2  
~~X~~ <=5

# Cracking example

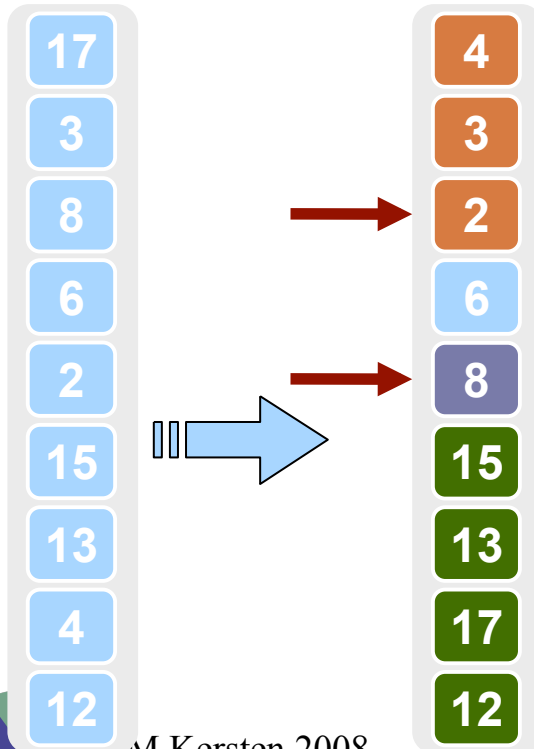
select A>5 and A<10



~~X~~ >5 and <10  
2  
~~X~~ <=5

# Cracking example

select A>5 and A<10

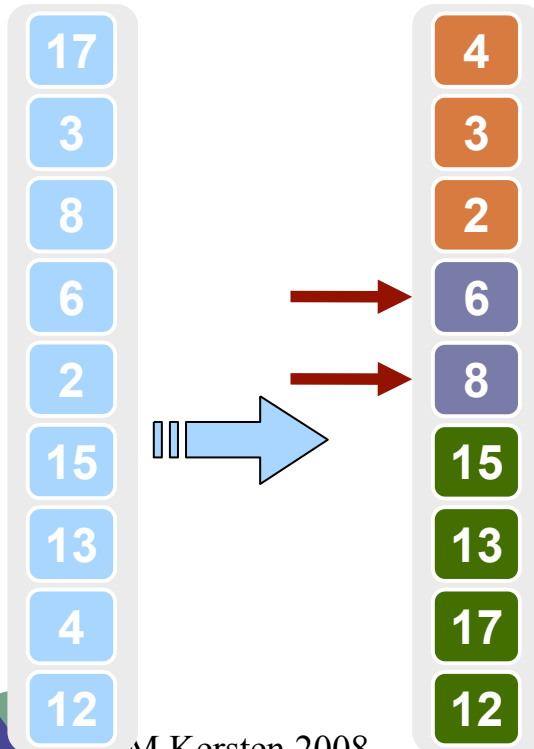


✗ >5 and <10

✗ <=5

# Cracking example

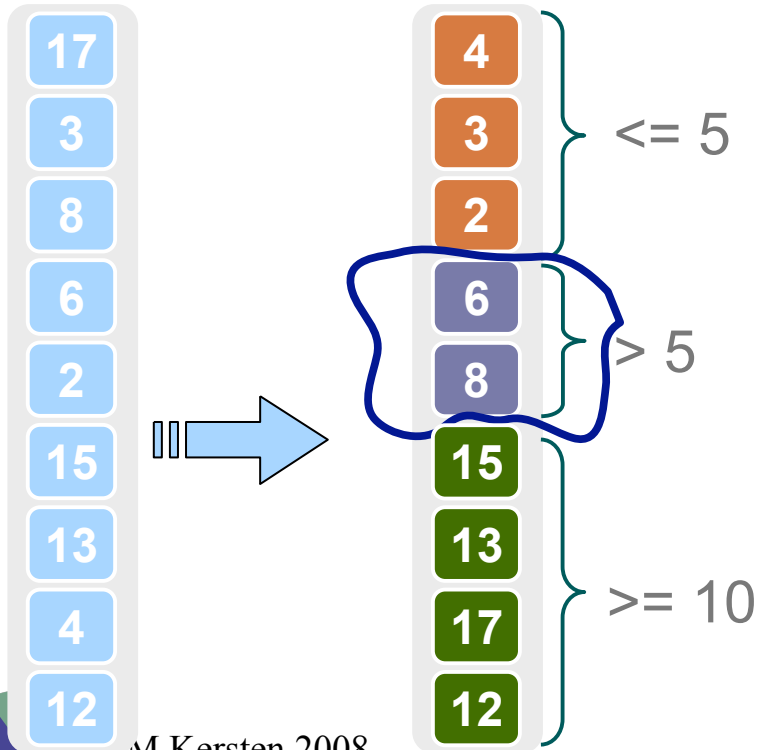
select A>5 and A<10



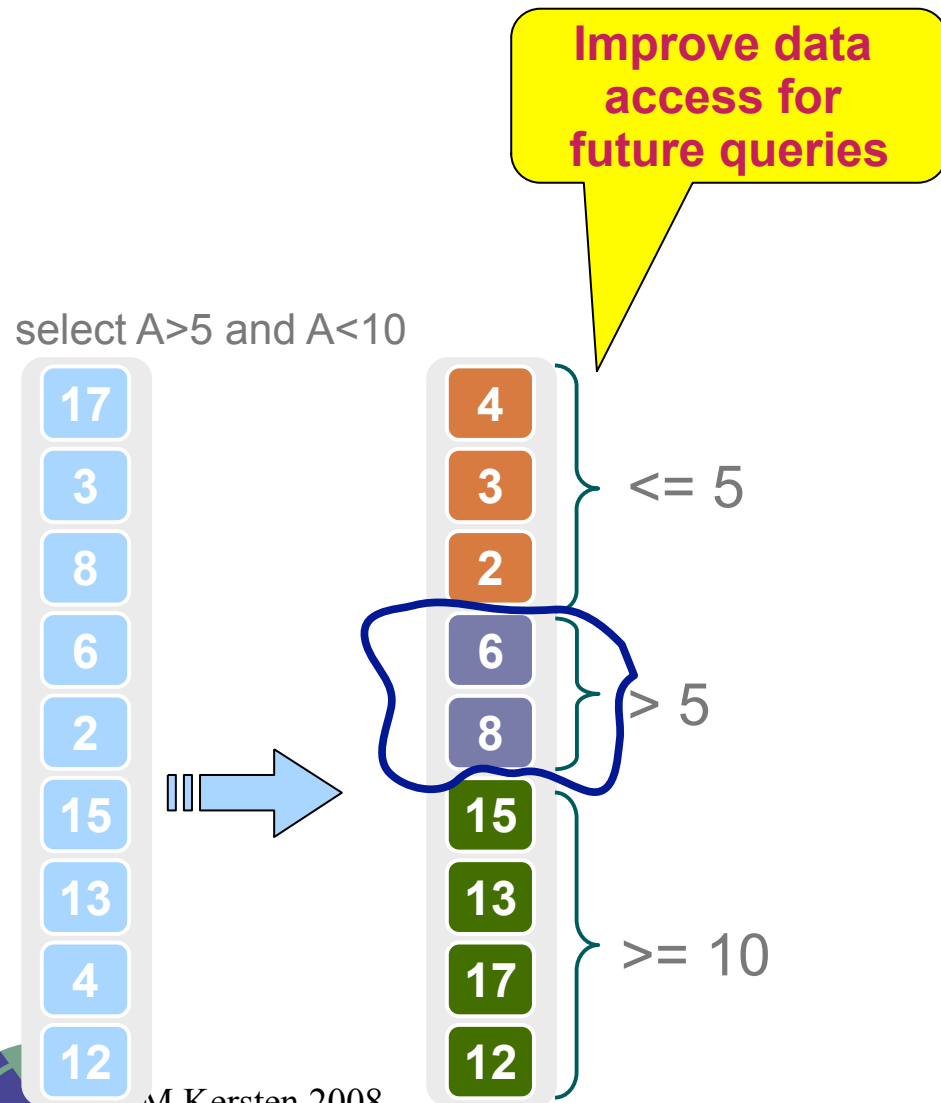
✓ >5 and <10

# Cracking example

select A>5 and A<10



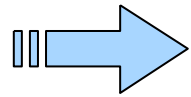
# Cracking example





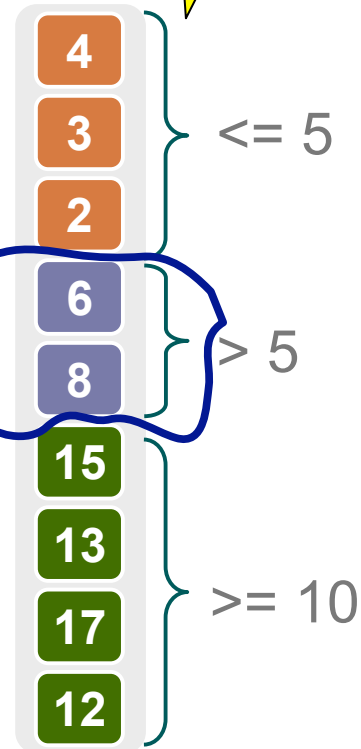
# Cracking example

select A>5 and A<10



Improve data access for future queries

select A>3 and A<14



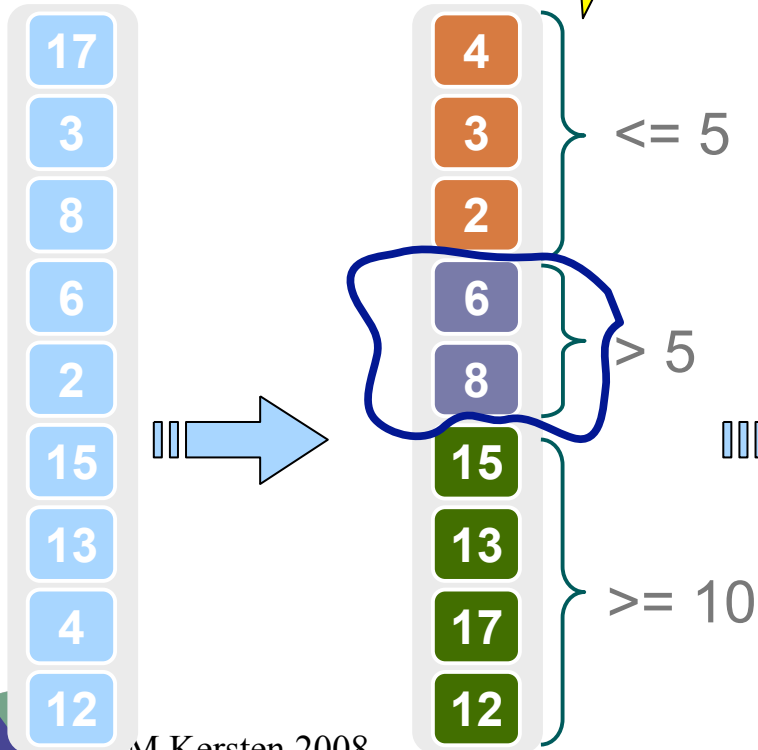
# Cracking example



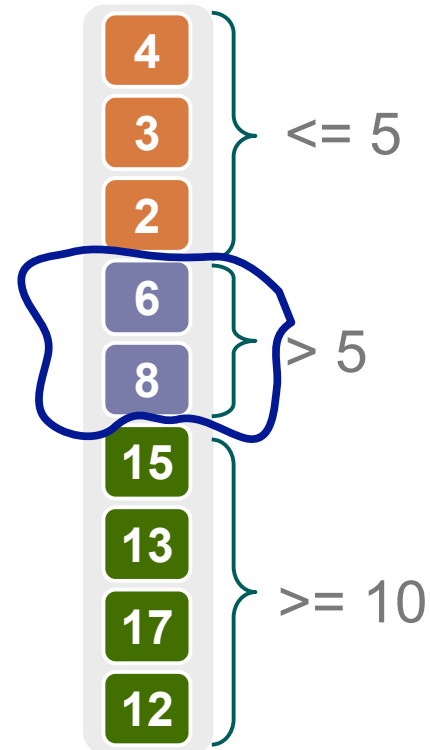
# Cracking example

Improve data access for future queries

select A>5 and A<10



select A>3 and A<14



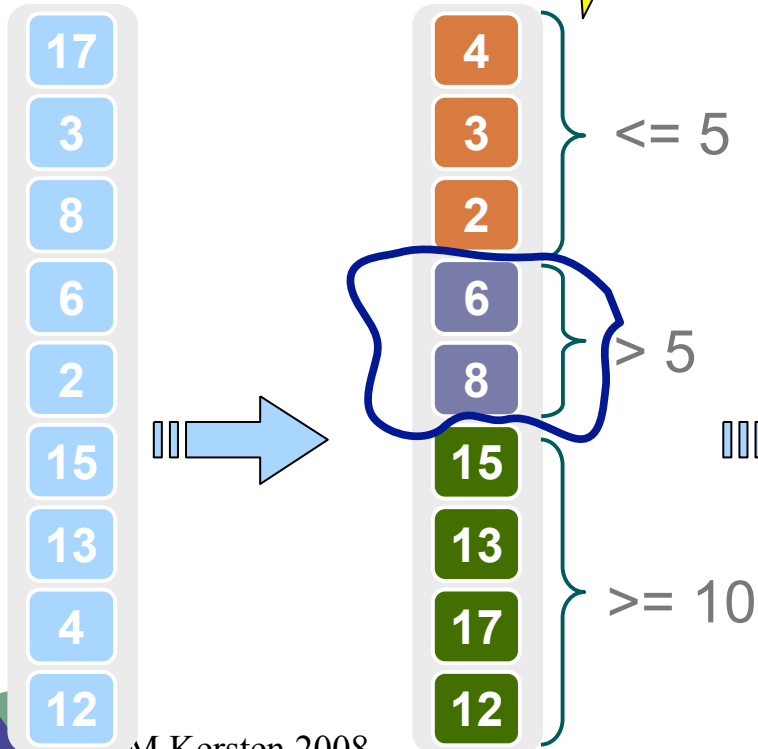
# Cracking example



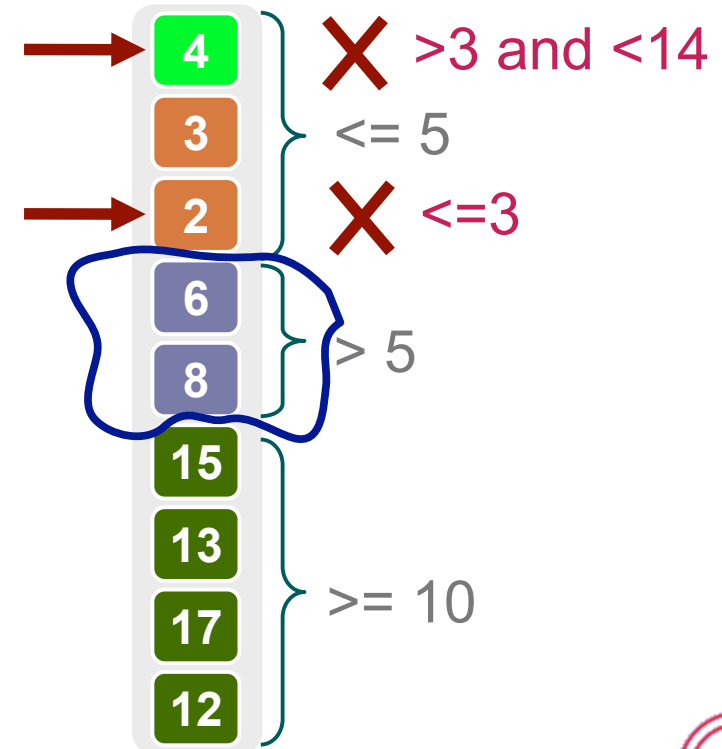
# racking example

Improve data access for future queries

select A>5 and A<10



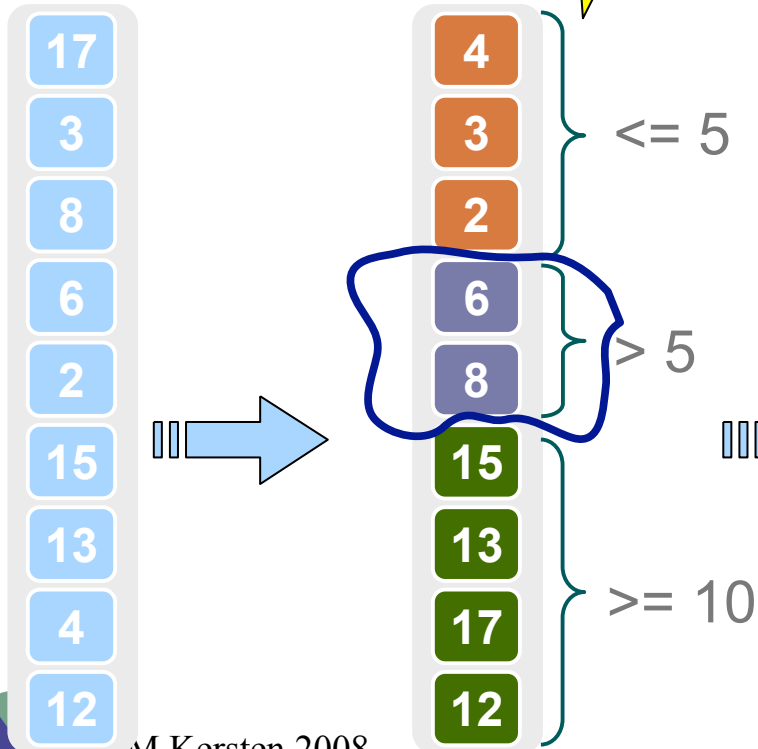
select A>3 and A<14



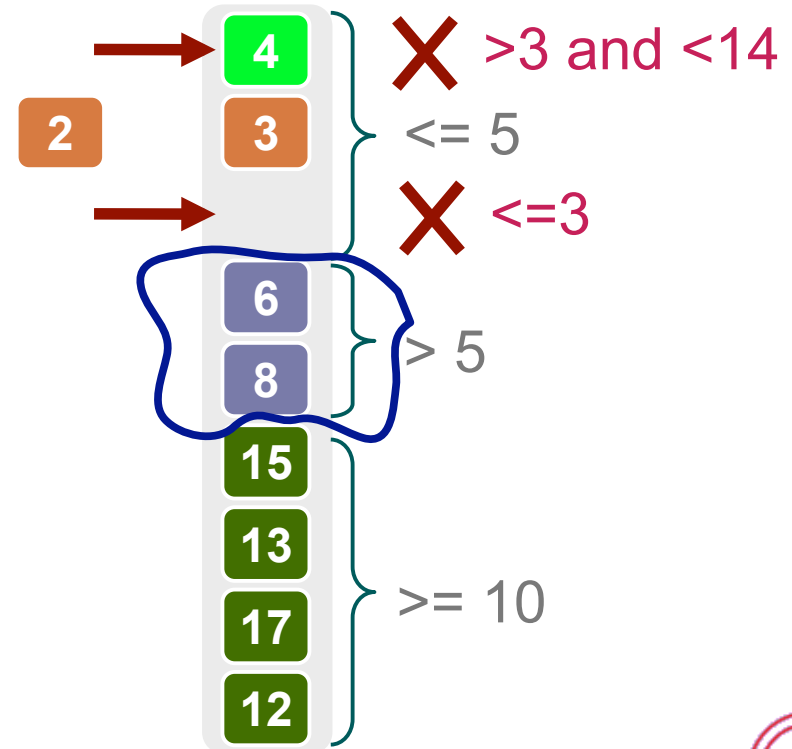
# Cracking example

Improve data access for future queries

select A>5 and A<10



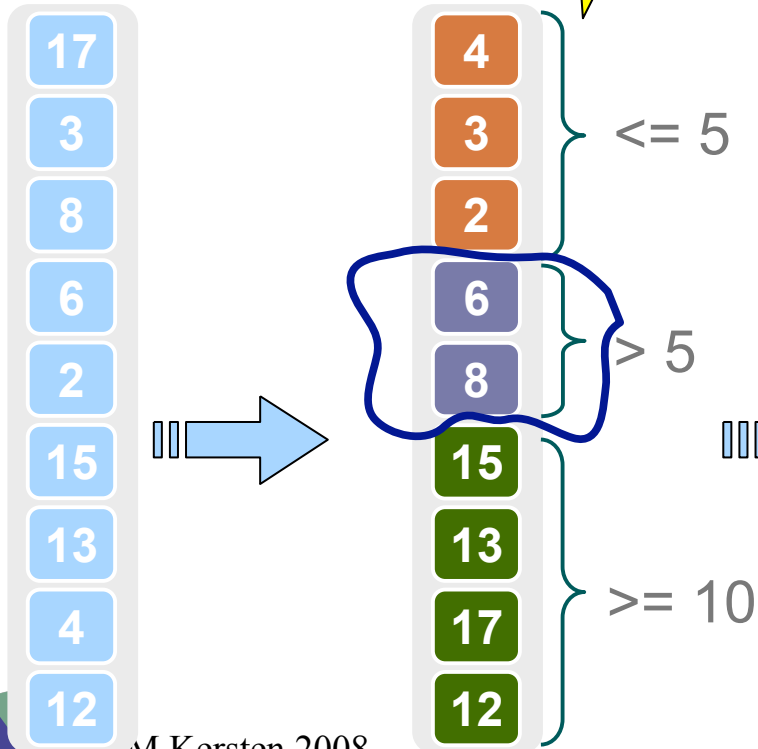
select A>3 and A<14



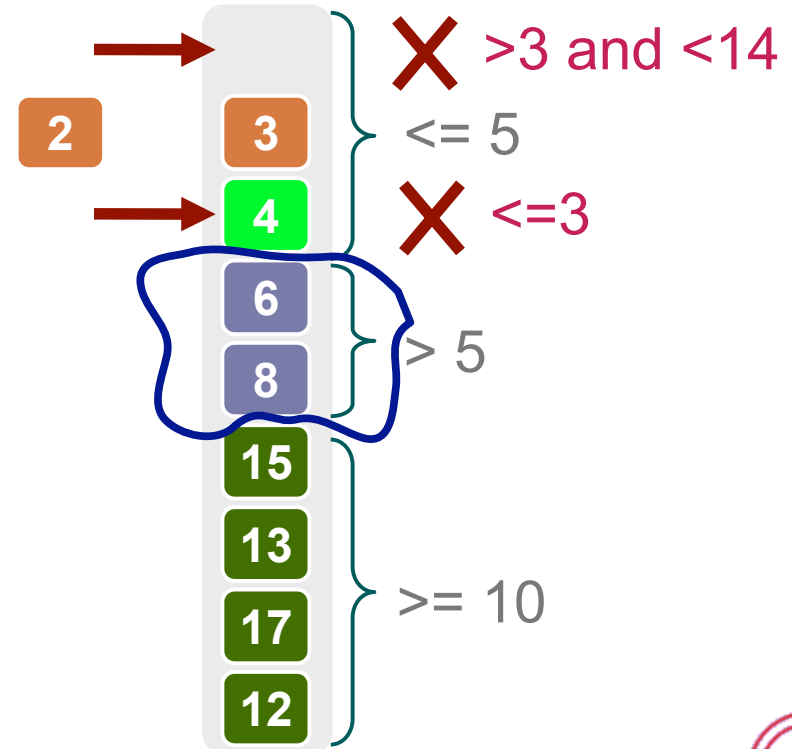
# Cracking example

Improve data access for future queries

select A>5 and A<10



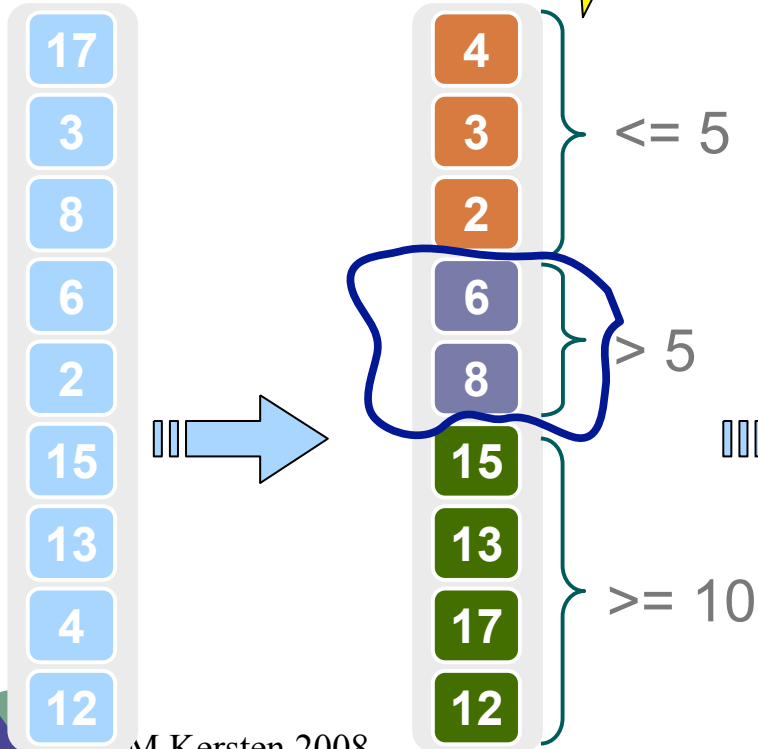
select A>3 and A<14



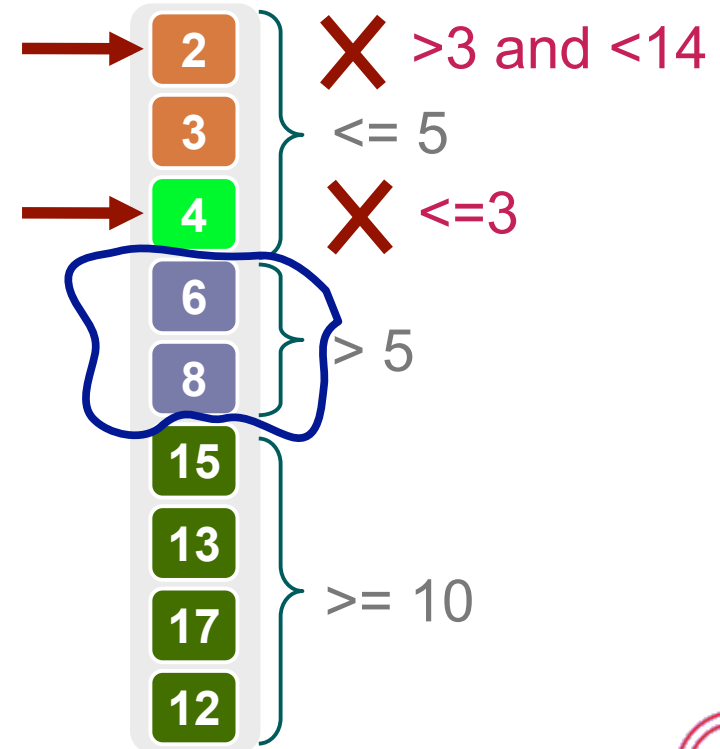
# Cracking example

Improve data access for future queries

select A>5 and A<10



select A>3 and A<14



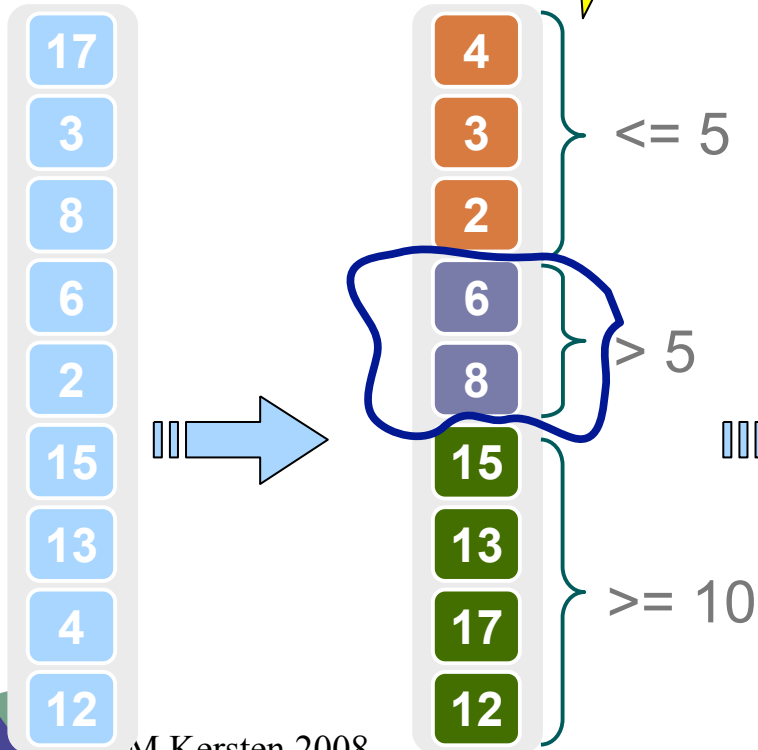




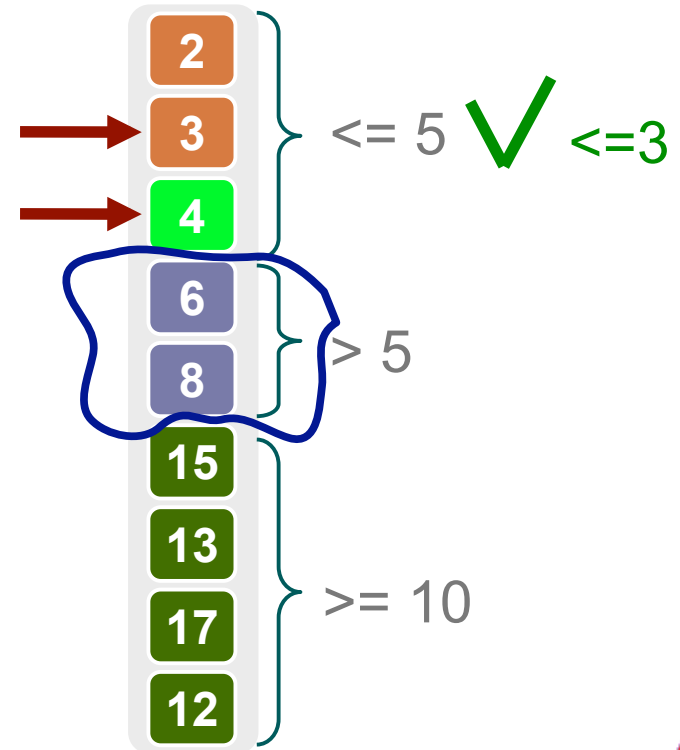
# Cracking example

Improve data access for future queries

select A>5 and A<10



select A>3 and A<14



# Cracking example



# Cracking example



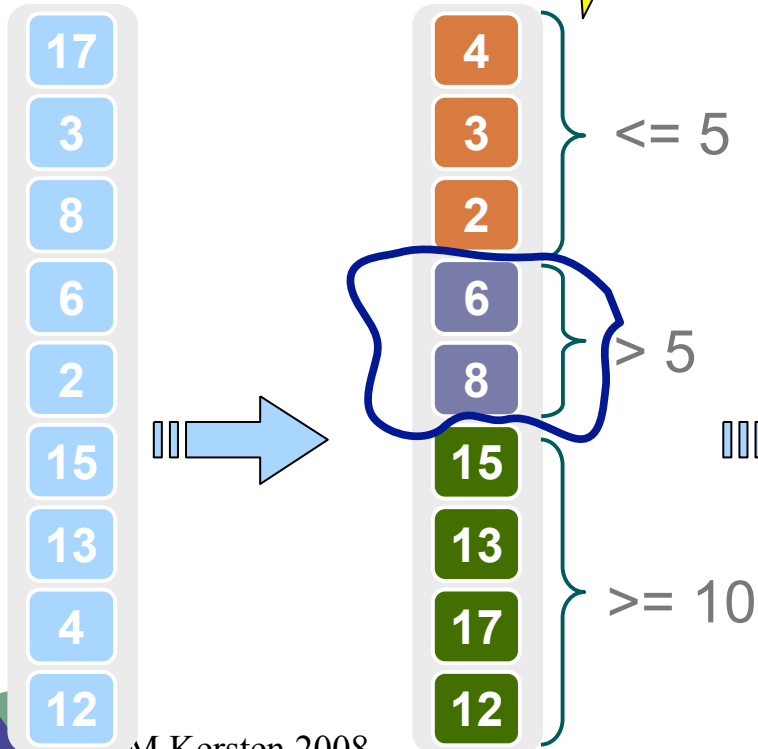
# Cracking example



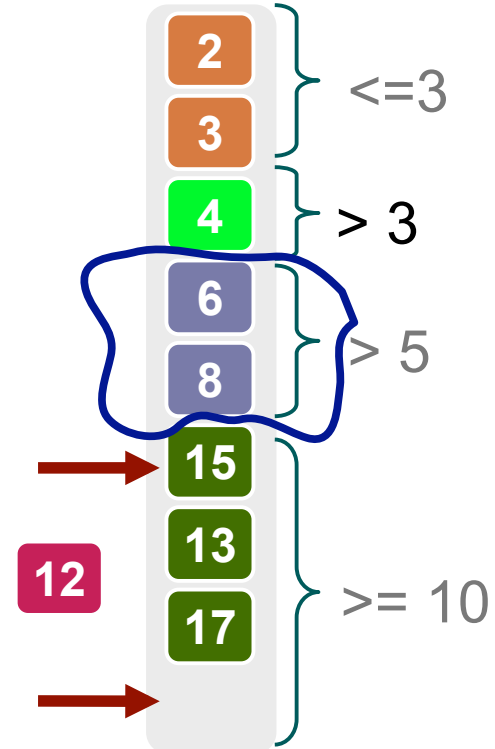
# Cracking example

Improve data access for future queries

select A>5 and A<10



select A>3 and A<14



# Cracking example



# Cracking example



# Cracking example

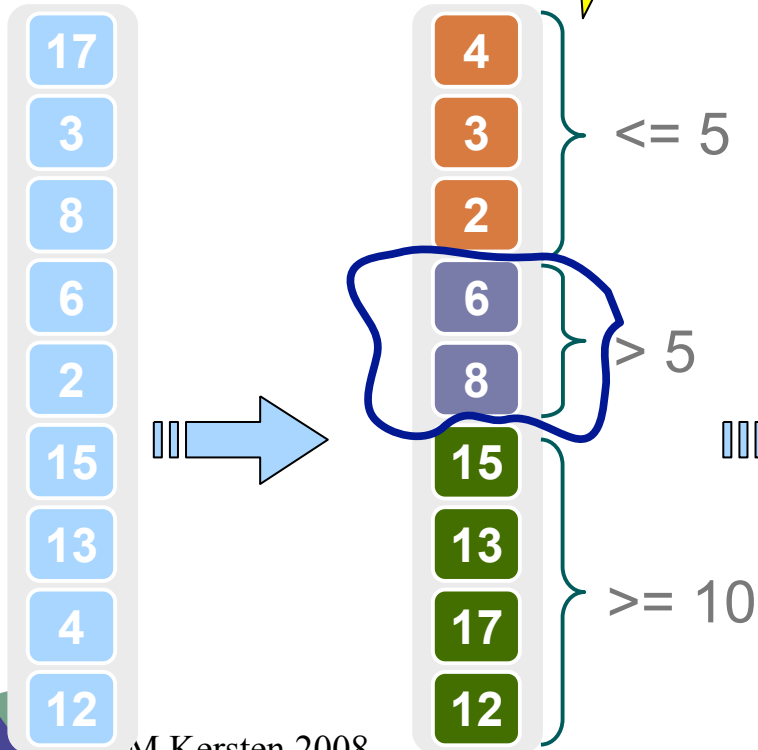




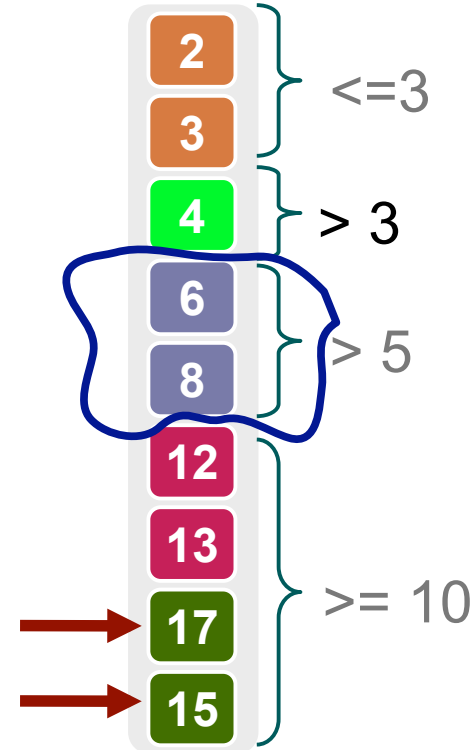
# Cracking example

Improve data access for future queries

select A>5 and A<10



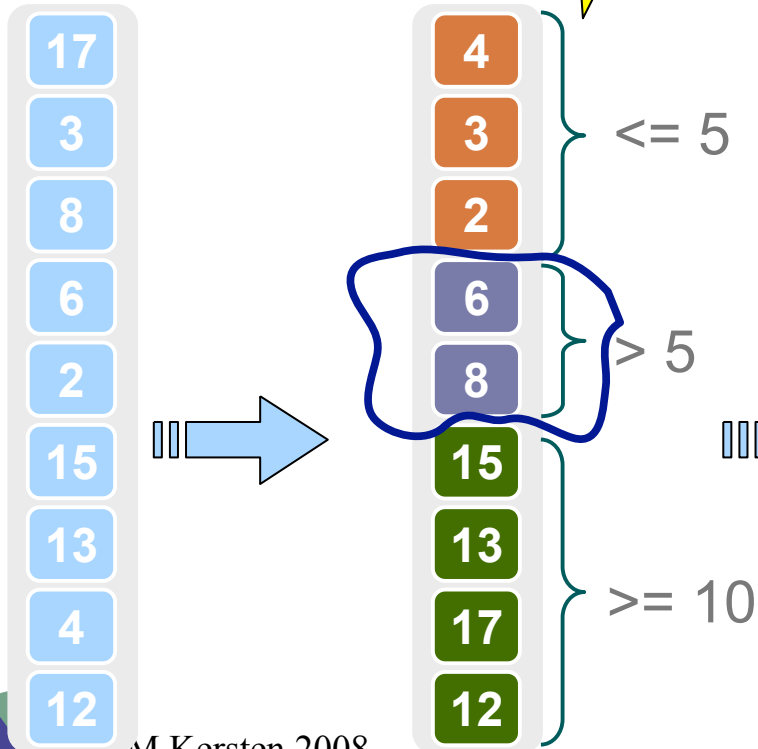
select A>3 and A<14



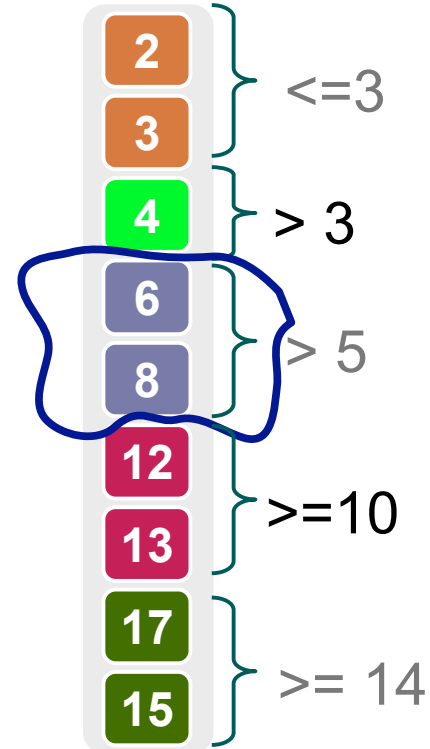
# Cracking example

Improve data access for future queries

select A>5 and A<10



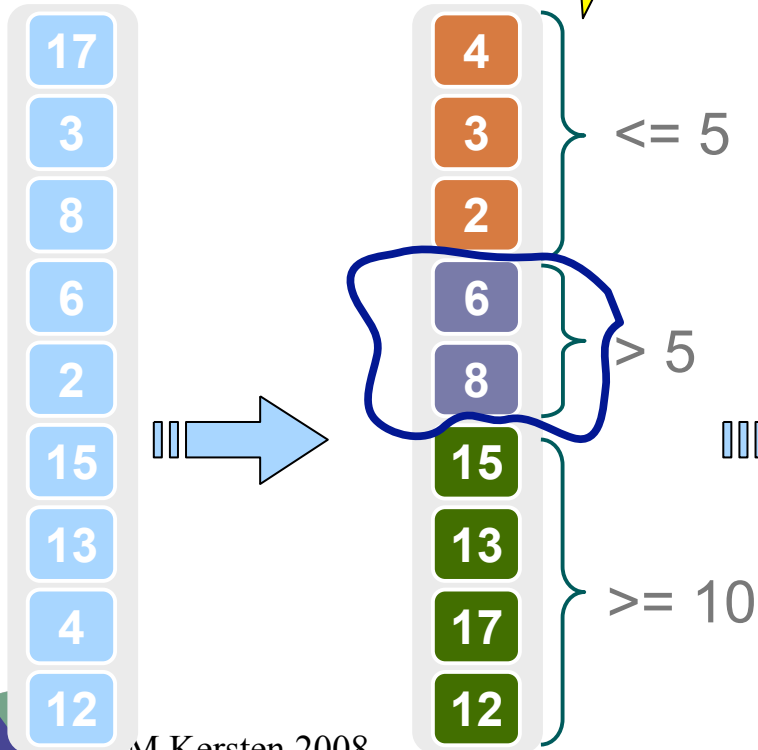
select A>3 and A<14



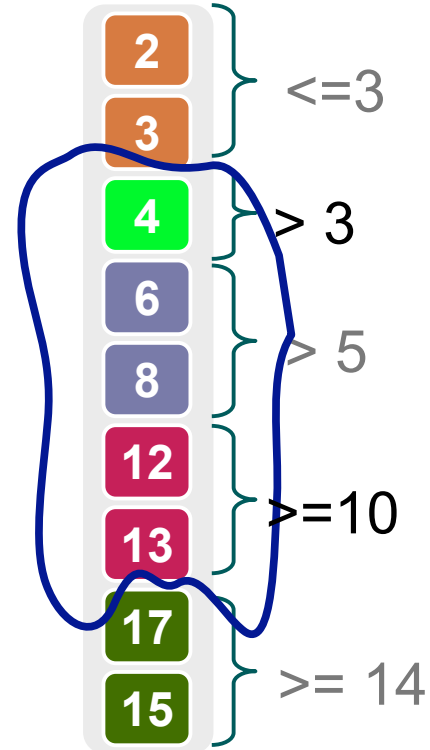
# Cracking example

Improve data access for future queries

select A>5 and A<10

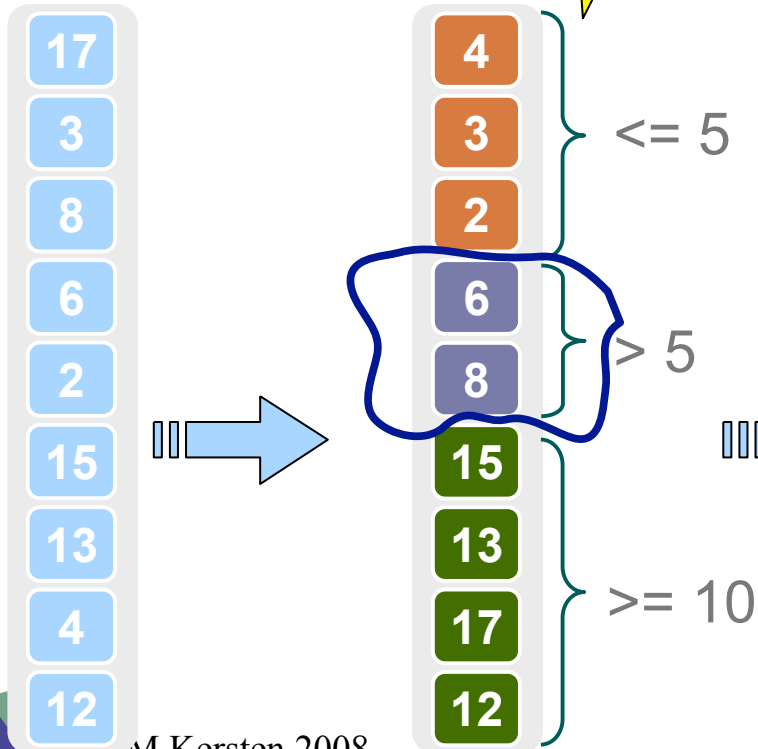


select A>3 and A<14



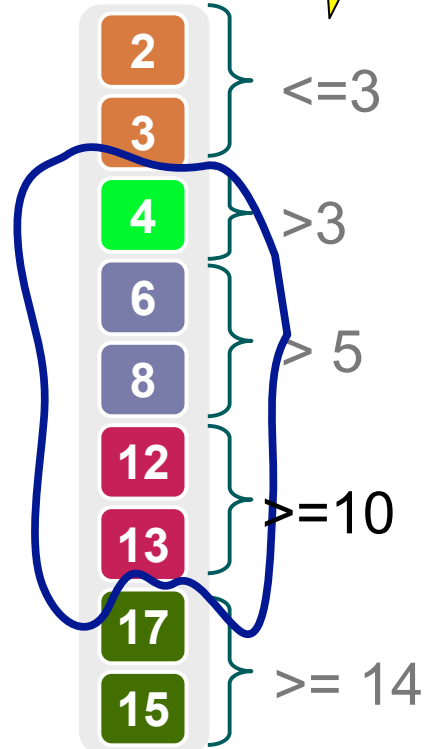
# Cracking example

select A>5 and A<10



Improve data access for future queries

select A>3 and A<14

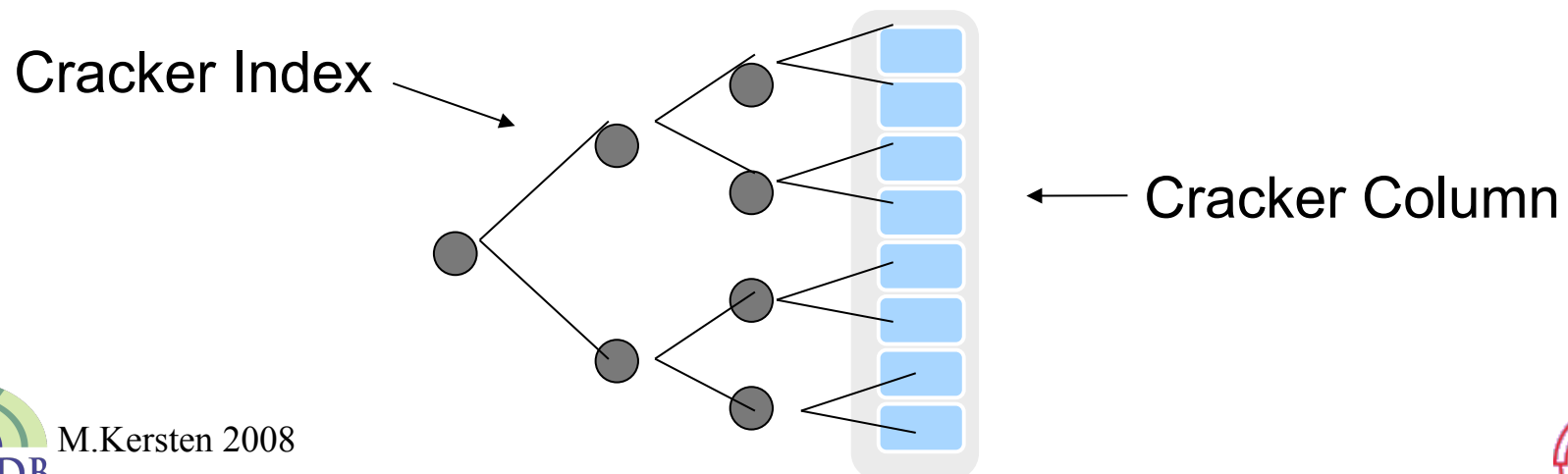


The more we crack the more we learn

The **first** time a range query is posed on an attribute  $A$ , a cracking DBMS makes a **copy** of column  $A$ , called the *cracker column* of  $A$

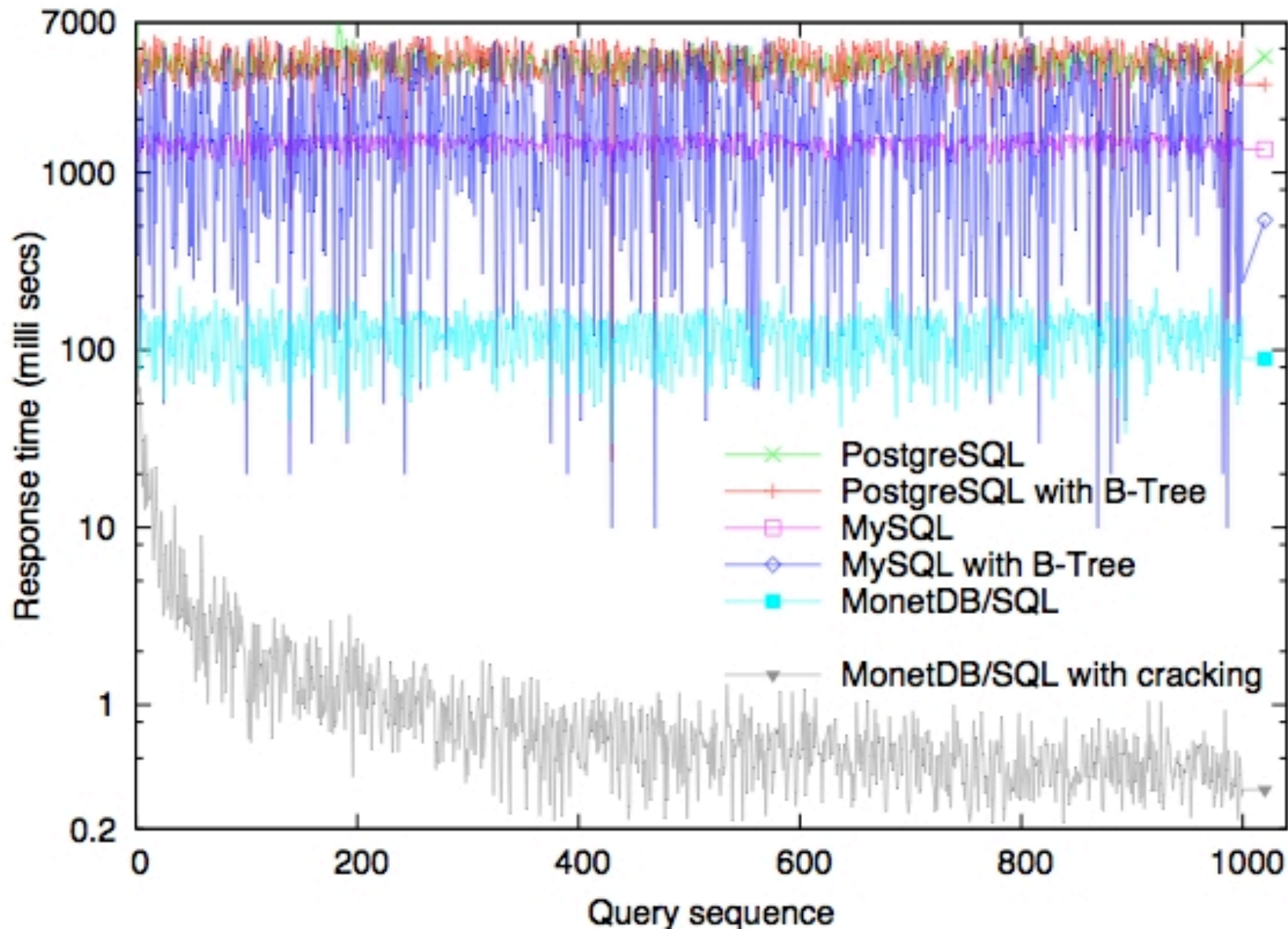
A cracker column is **continuously** physically reorganized based on queries that **need** to touch attribute such as the result is in a contiguous space

For each cracker column, there is a *cracker index*





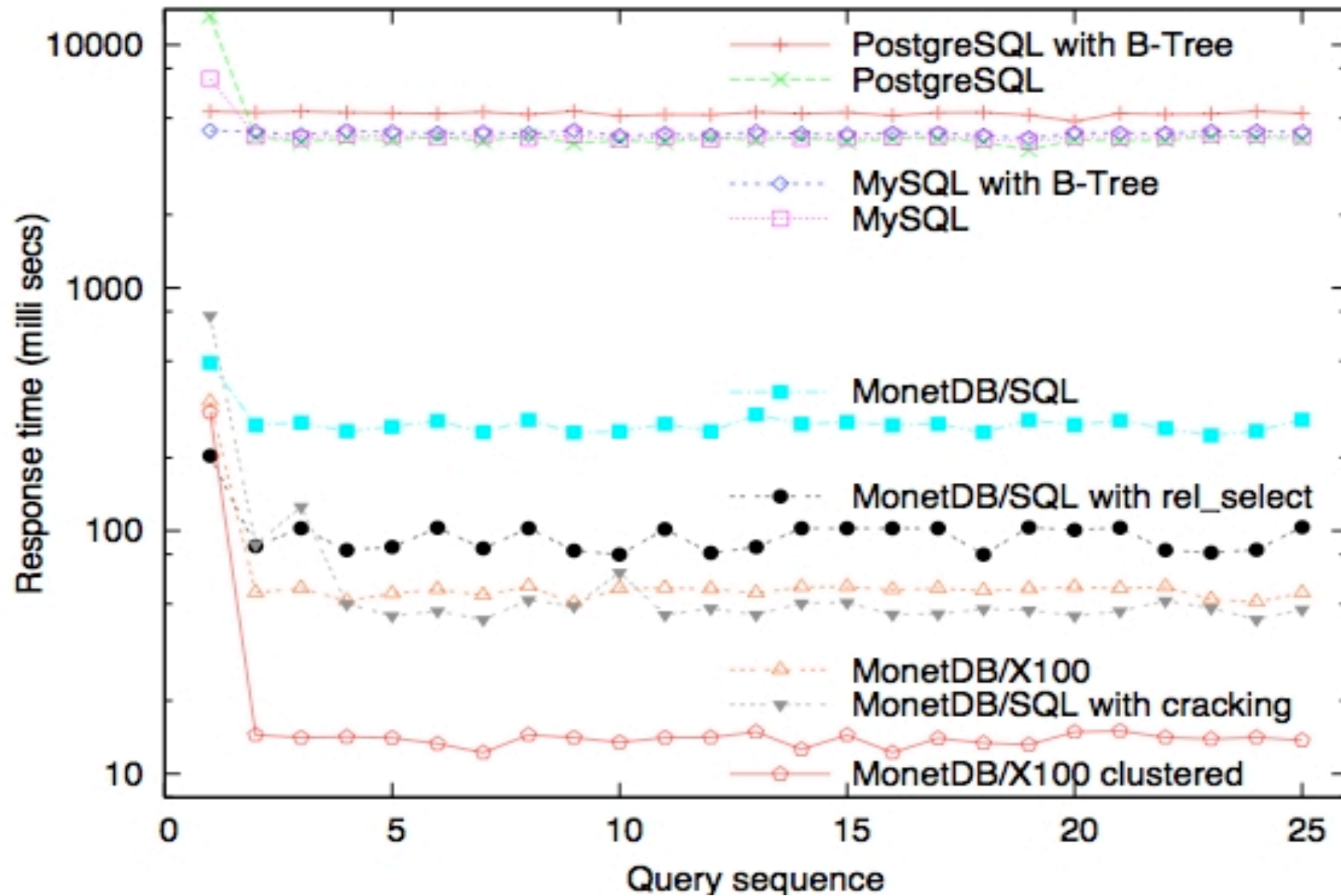
*Try to avoid useless investments*  
A simple range query





*Try to avoid useless investments*

## TPC-H query 6





*Try to avoid useless investments*

- Cracking is easy in a column store and is part of the critical execution path
- Cracking works under high volume updates





# Updates

- Base columns are updated as normally
- We need to update the cracker column and the cracker index
- Efficiently
- Maintain the self-organization properties
- Two issues:
  - When
  - How





# When to propagate updates in cracking

- Follow the workload to maintain self-organization
- Updates become part of query processing
- When an update arrives, it is **not** applied
- For each cracker column there is
  - a pending insertions column
  - and a pending deletions column
- Pending updates are **applied only** when a query needs the specific values





# Updates aware select

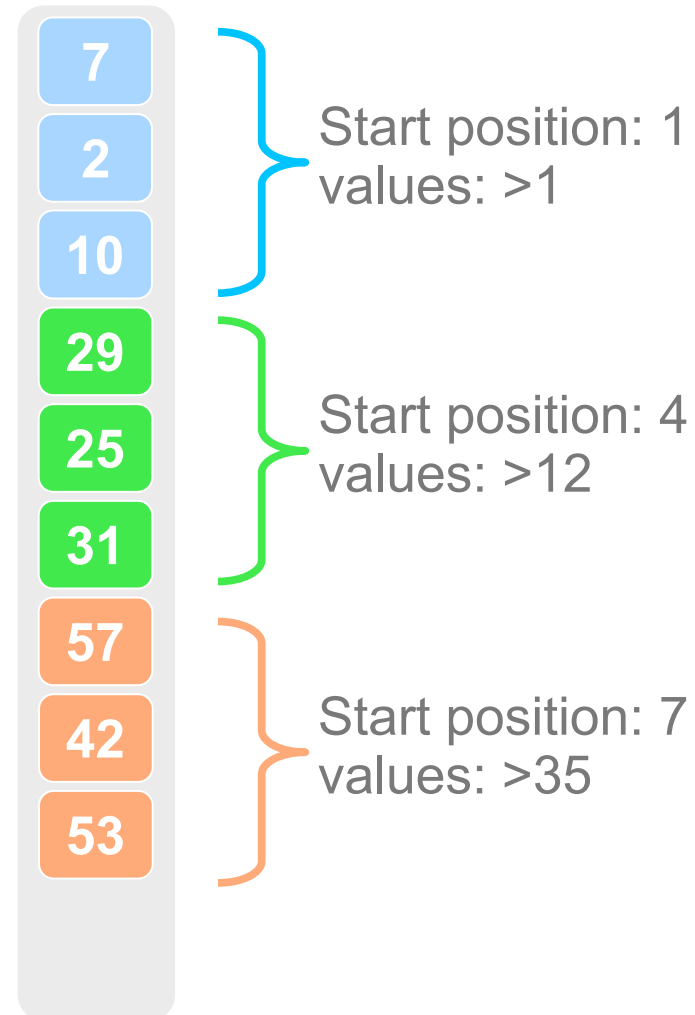
- We extended the cracker select operator to apply the needed updates **before** cracking
  
- The select operator:
  1. **Search the pending insertions column**
  2. **Search the pending deletions column**
  3. **If Steps 1 or 2 find tuples run an update algorithm**
  4. Search the cracker index
  5. Physically reorganize the cracker column
  6. Update the cracker index
  7. Return a slice of the cracker column



Insert a new tuple with value 9

The new tuple belongs to  
the blue piece

9



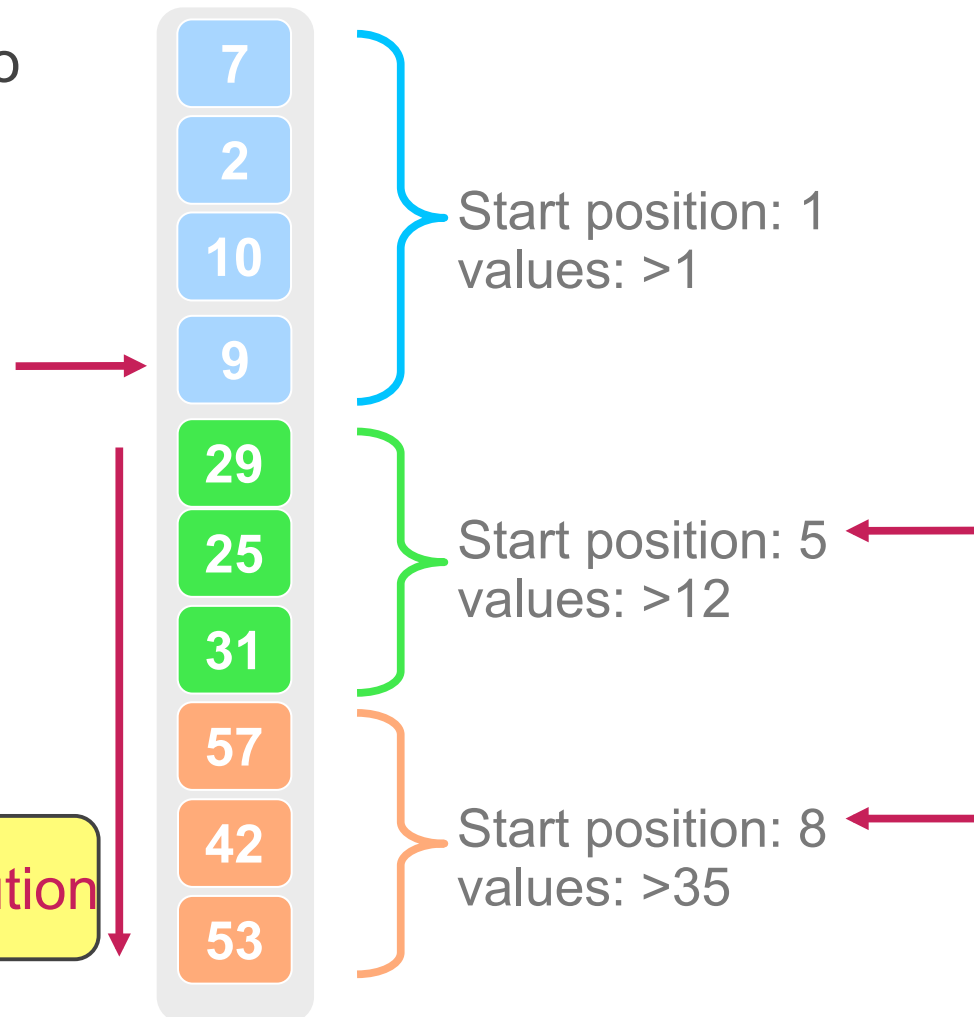
Insert a new tuple with value 9

The new tuple belongs to the blue piece

Pieces in the cracker column are **ordered**

Tuples inside a piece are **not ordered**

Shifting is not a viable solution



# Merging by Hopping

Insert a new tuple with value 9

9

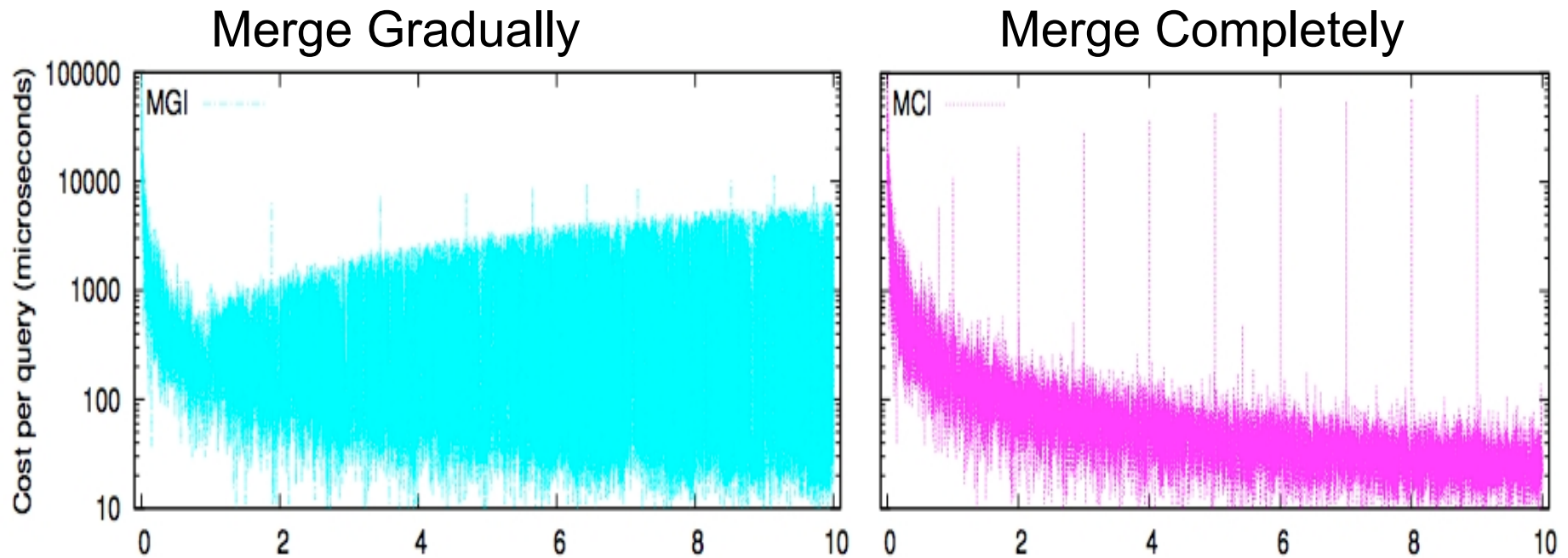
We need to make **enough room** to fit the new tuples

29



# Merge Gradually

- A query merges **only** the qualifying values, i.e., only the values that it needs for a correct and complete result



We avoid the large peaks  
but...

Average cost increases significantly



# The Ripple

Touch only the pieces that are **relevant** for the **current** query



M.Kersten 2008





# The Ripple

Touch only the pieces that are **relevant** for the **current** query





# The Ripple

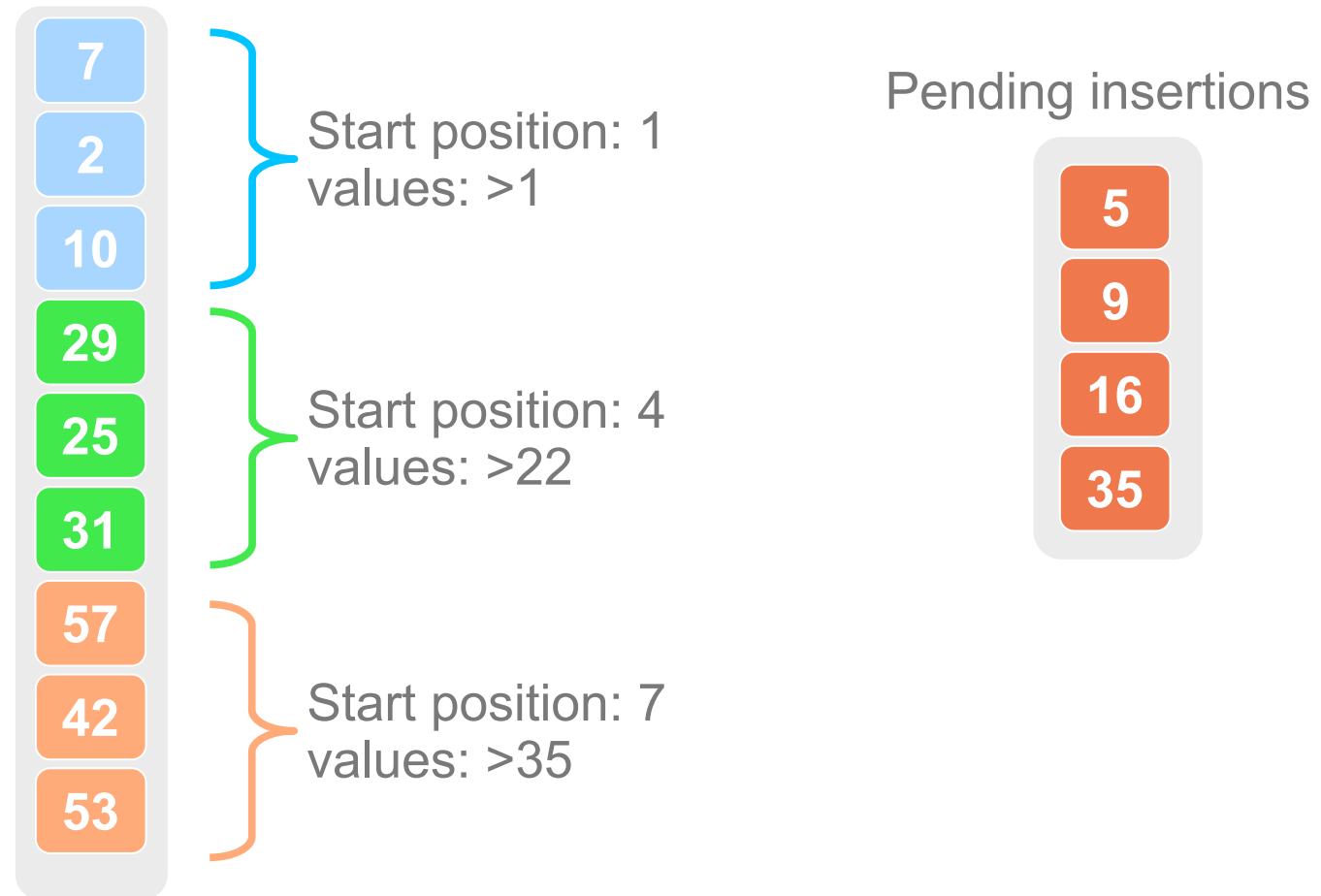
Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$



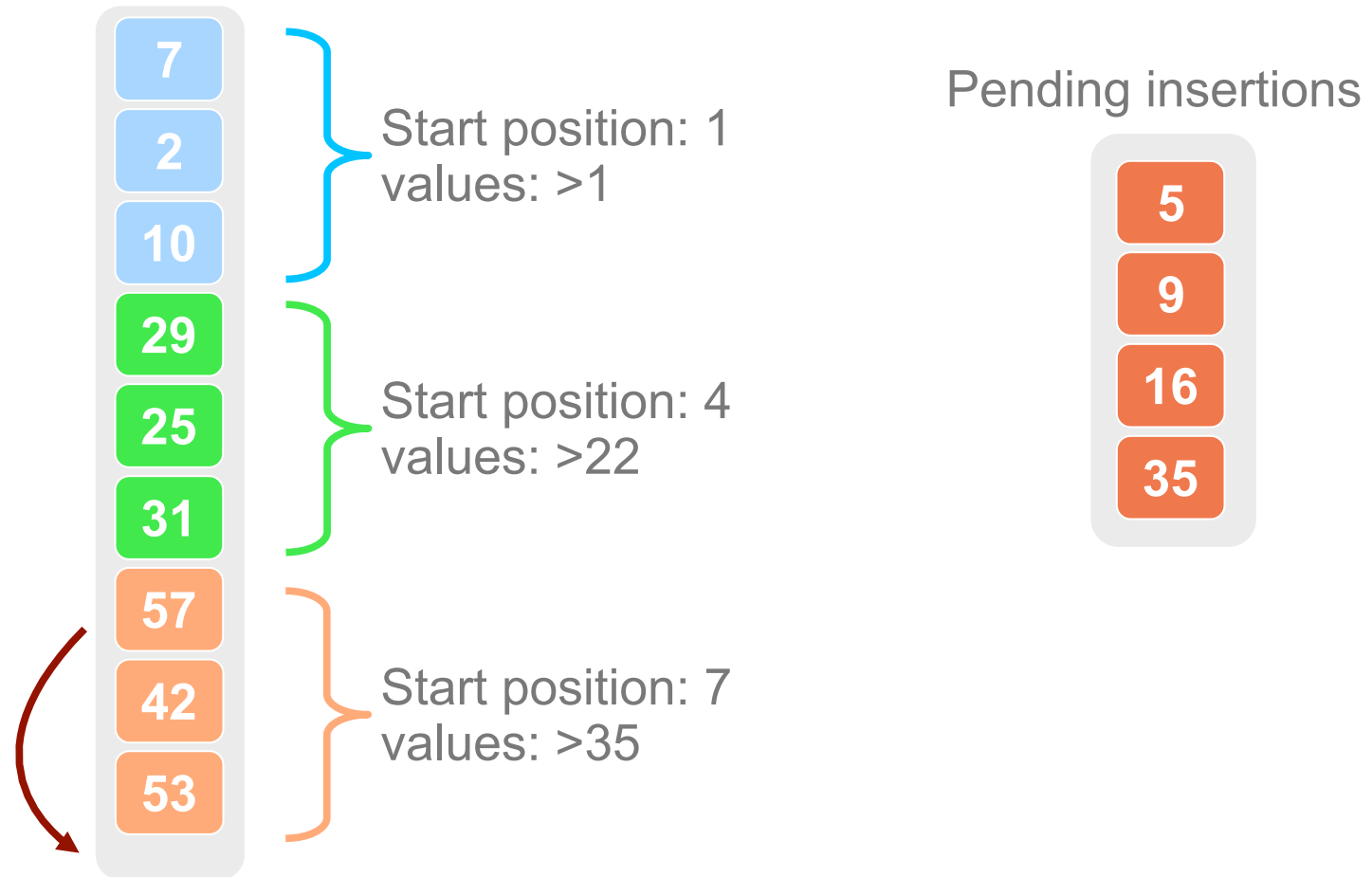


# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$



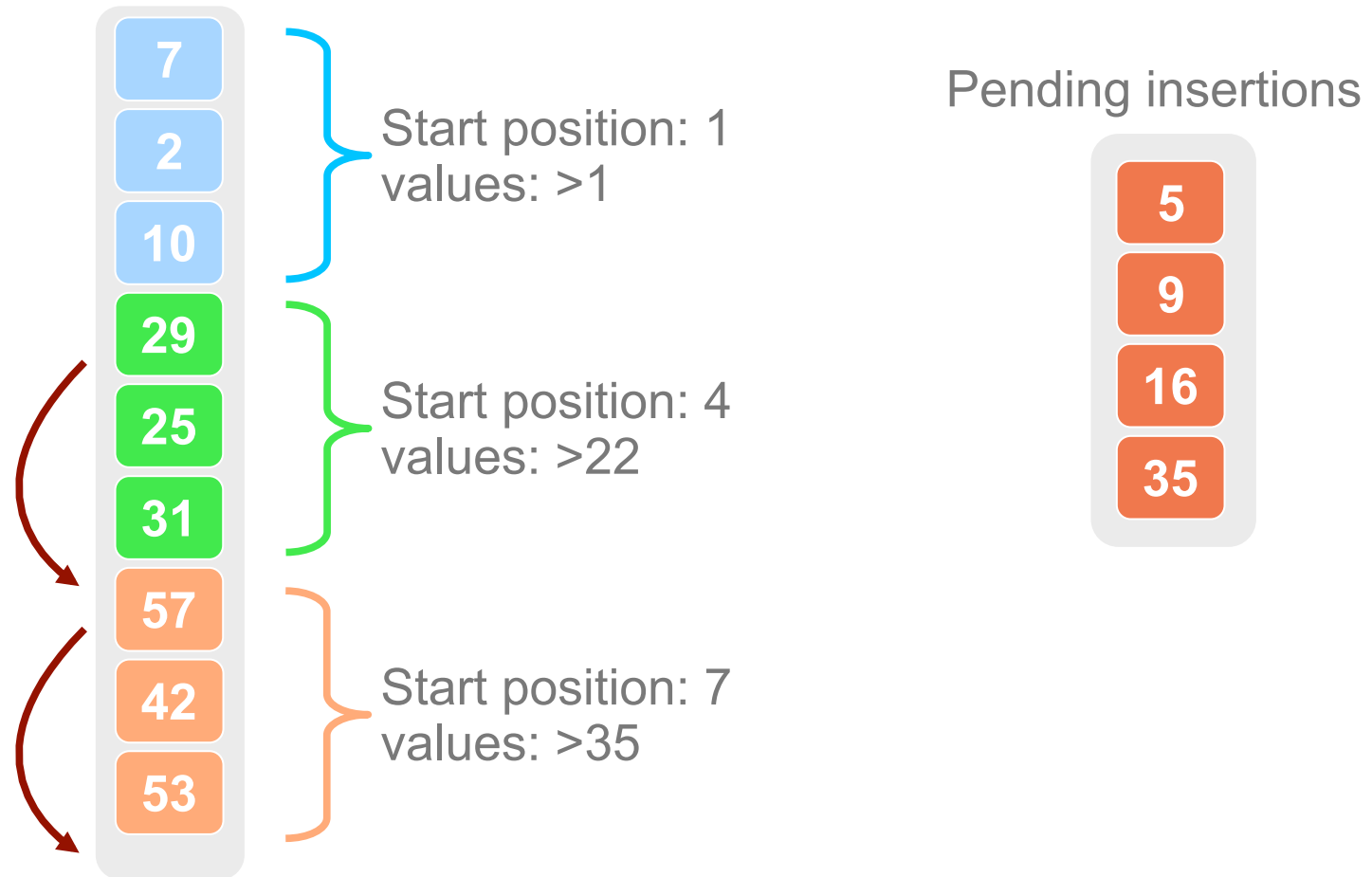
Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$





# The Ripple

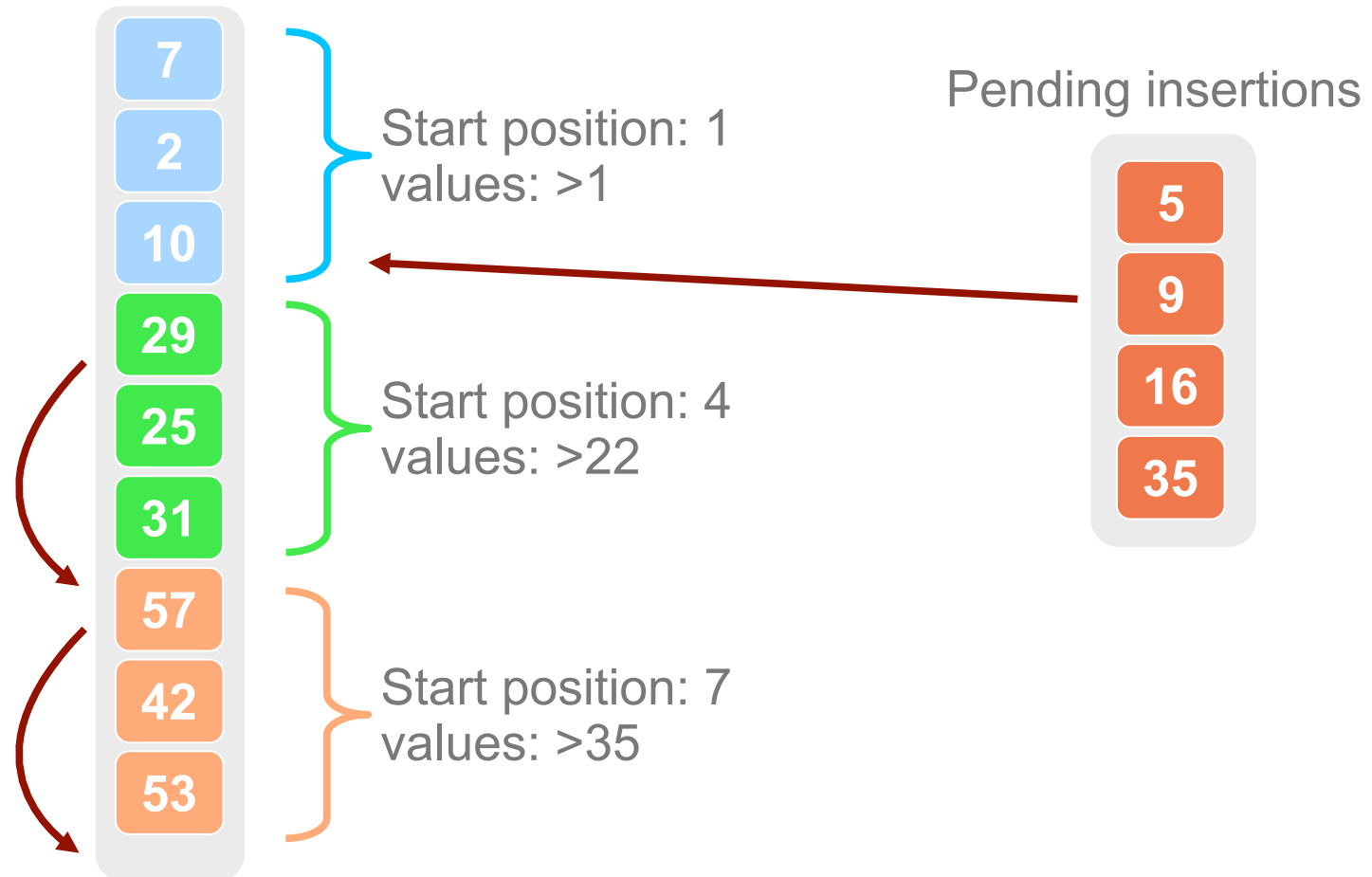
Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$





# The Ripple

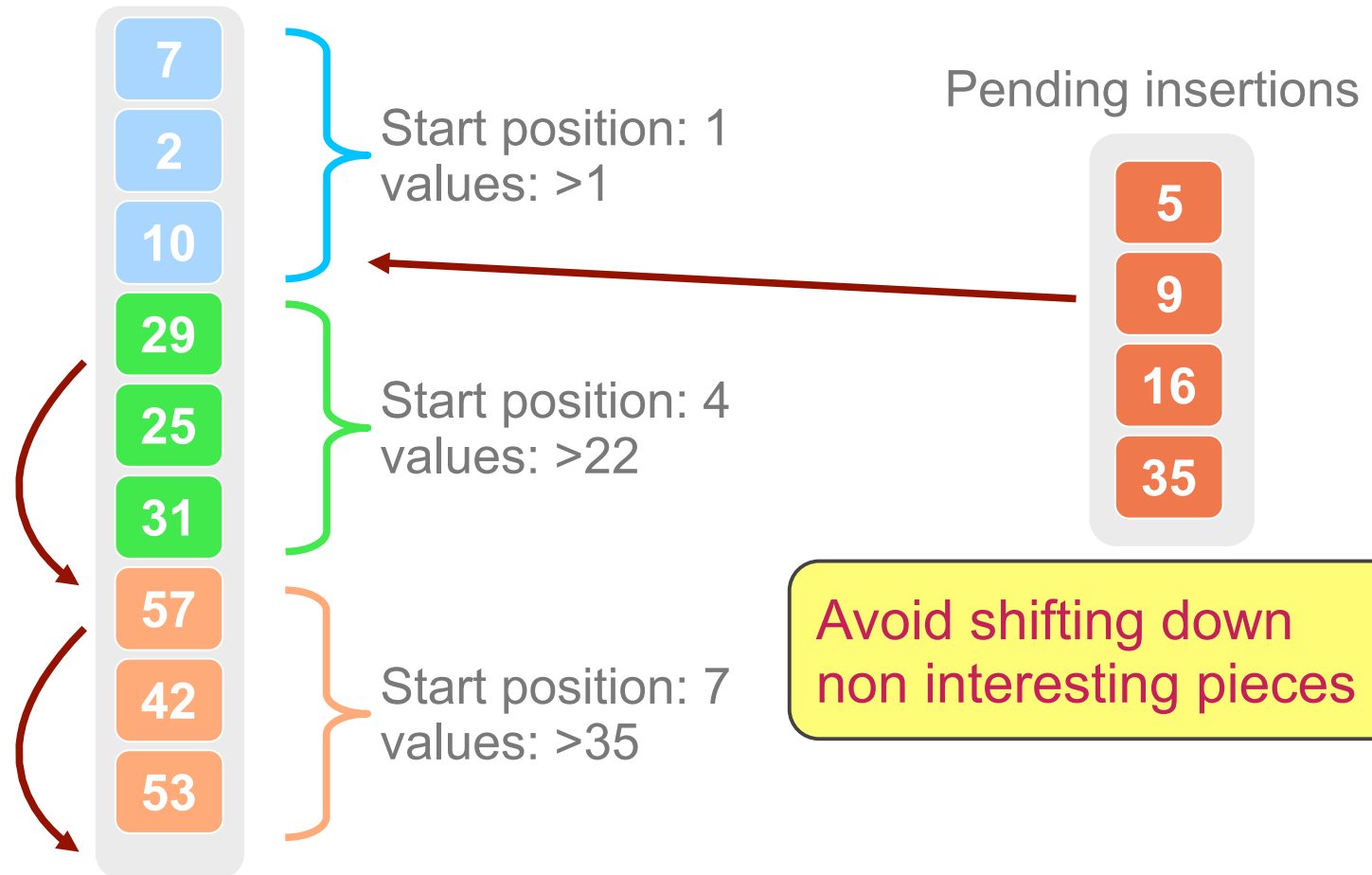
Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$





# The Ripple

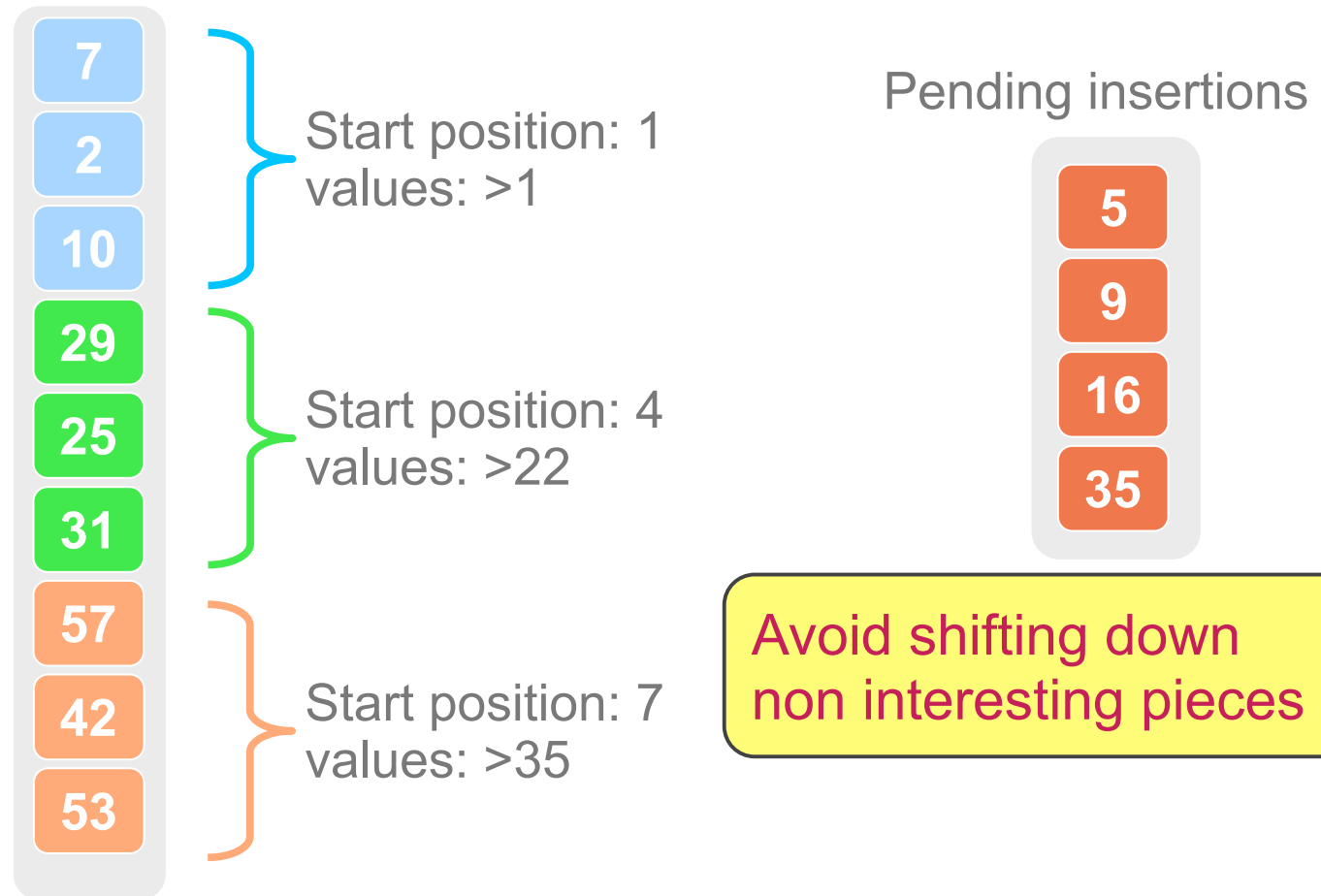
Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$





# The Ripple

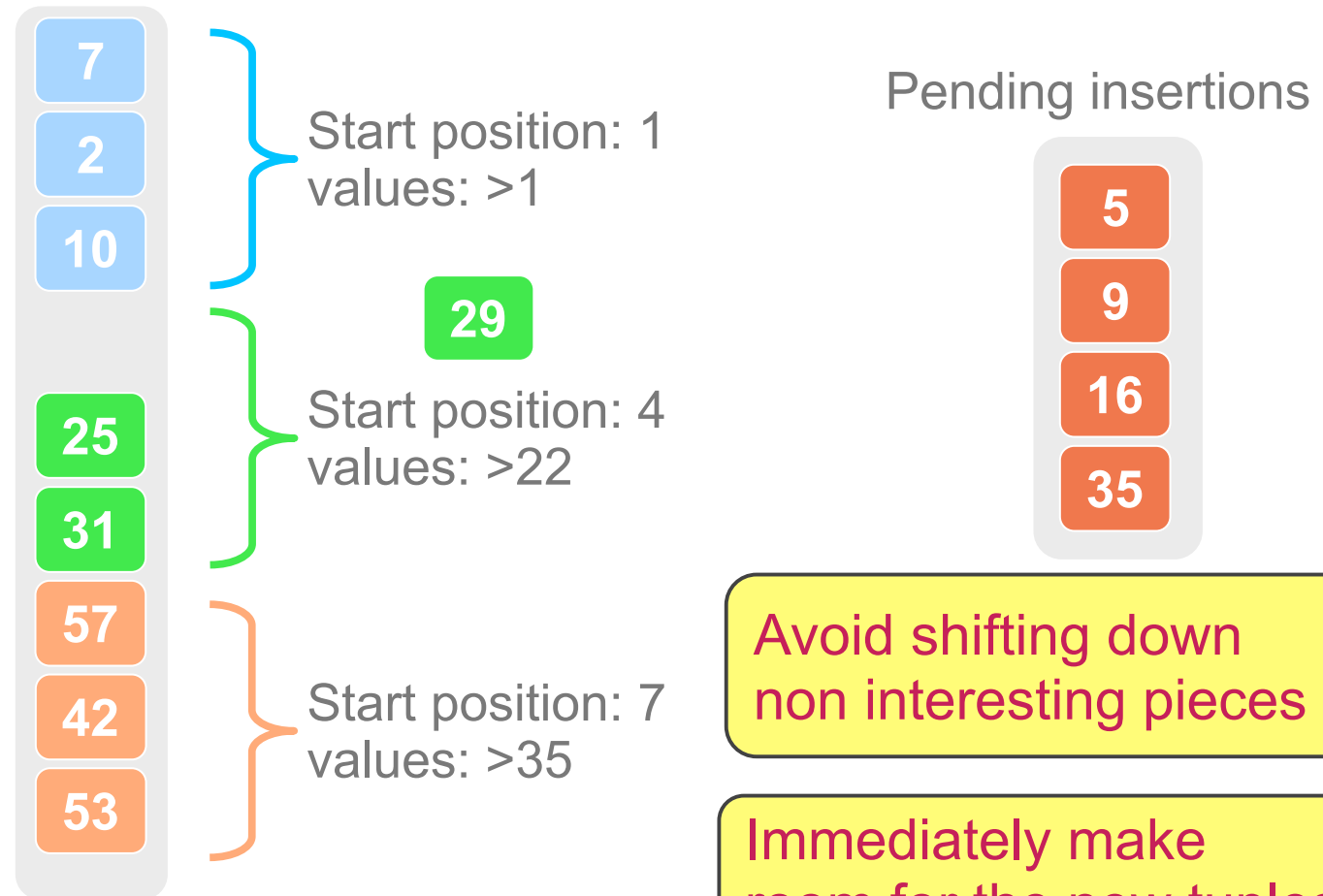
Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$





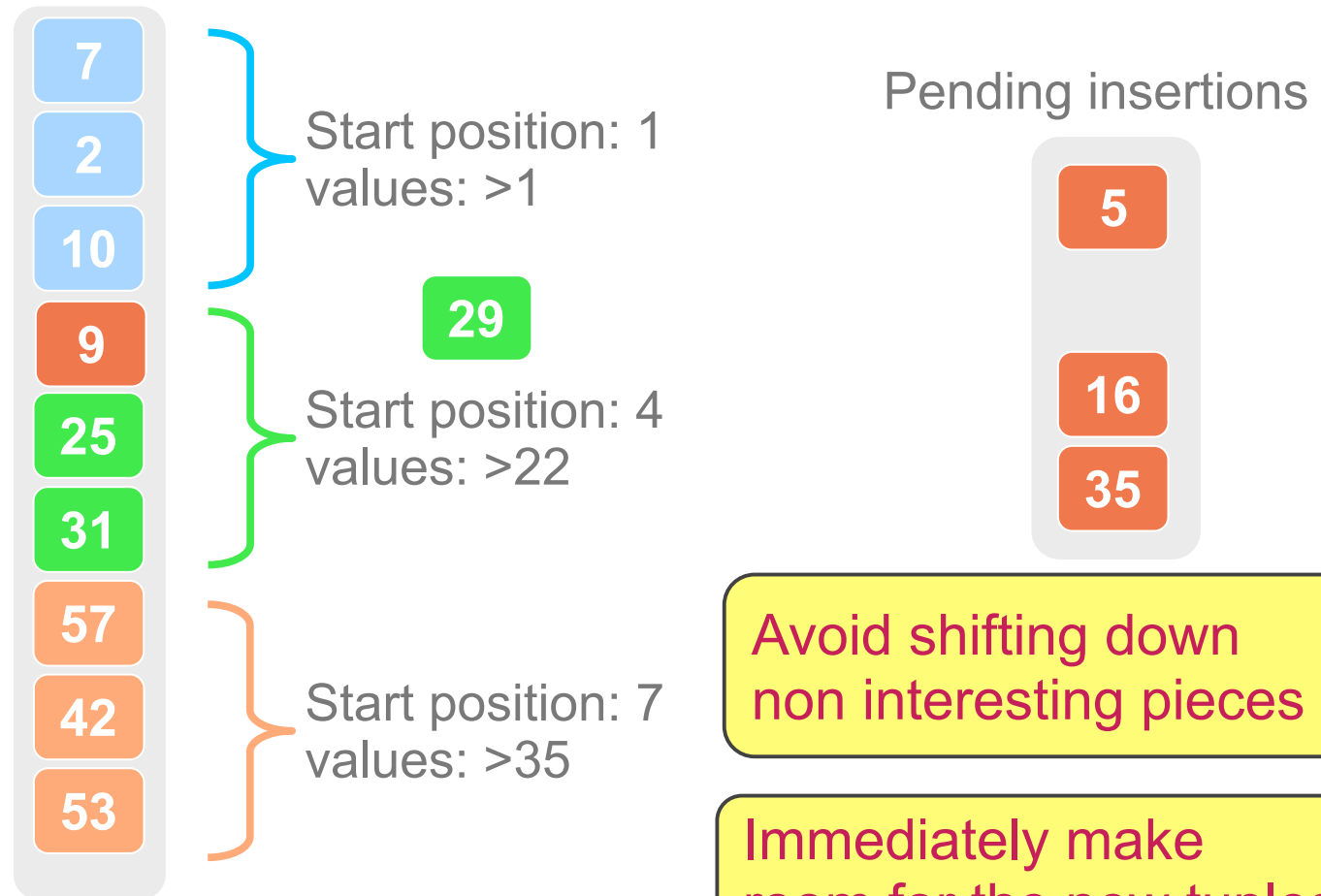
# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
 Select  $7 \leq A < 15$



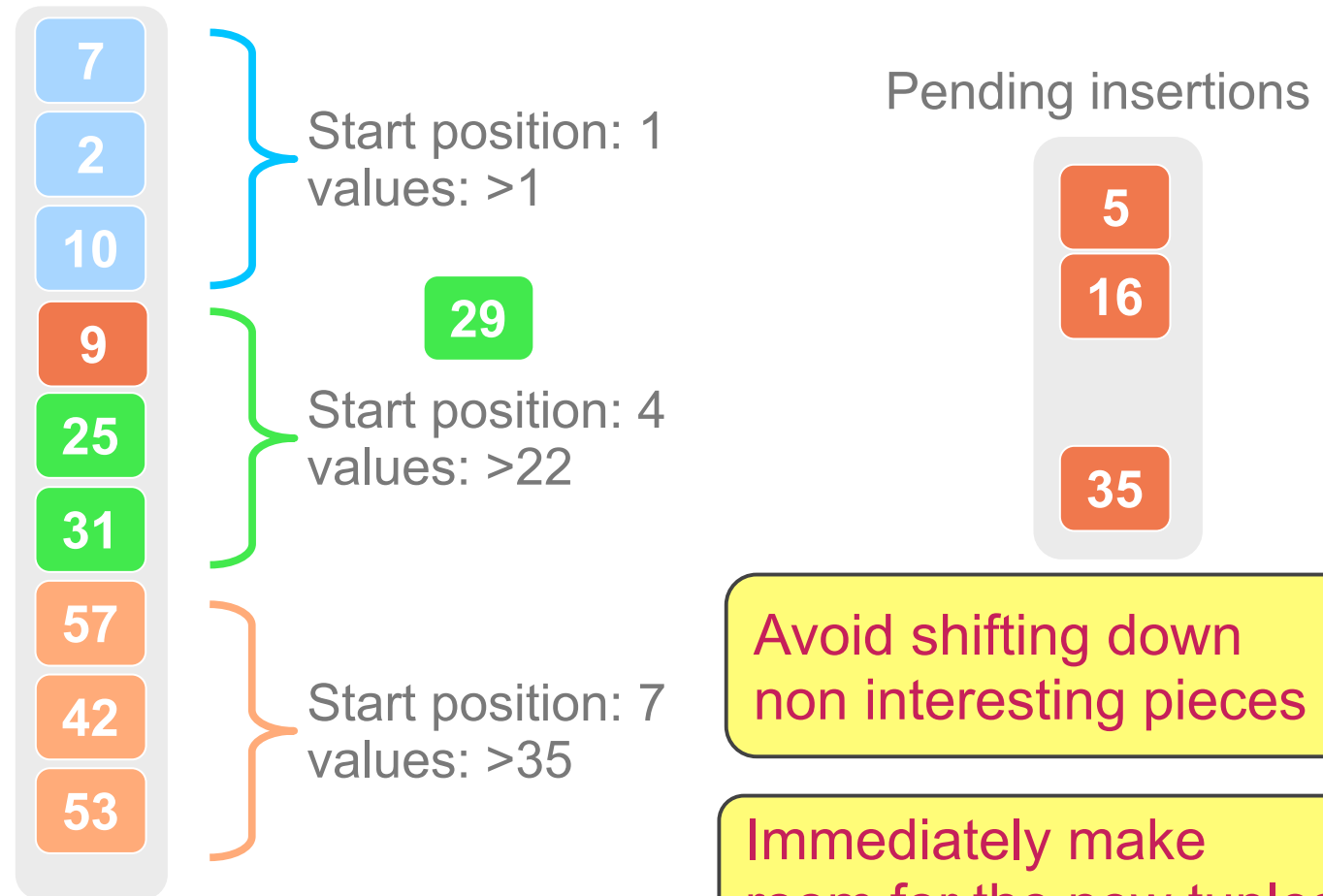
# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
 Select  $7 \leq A < 15$



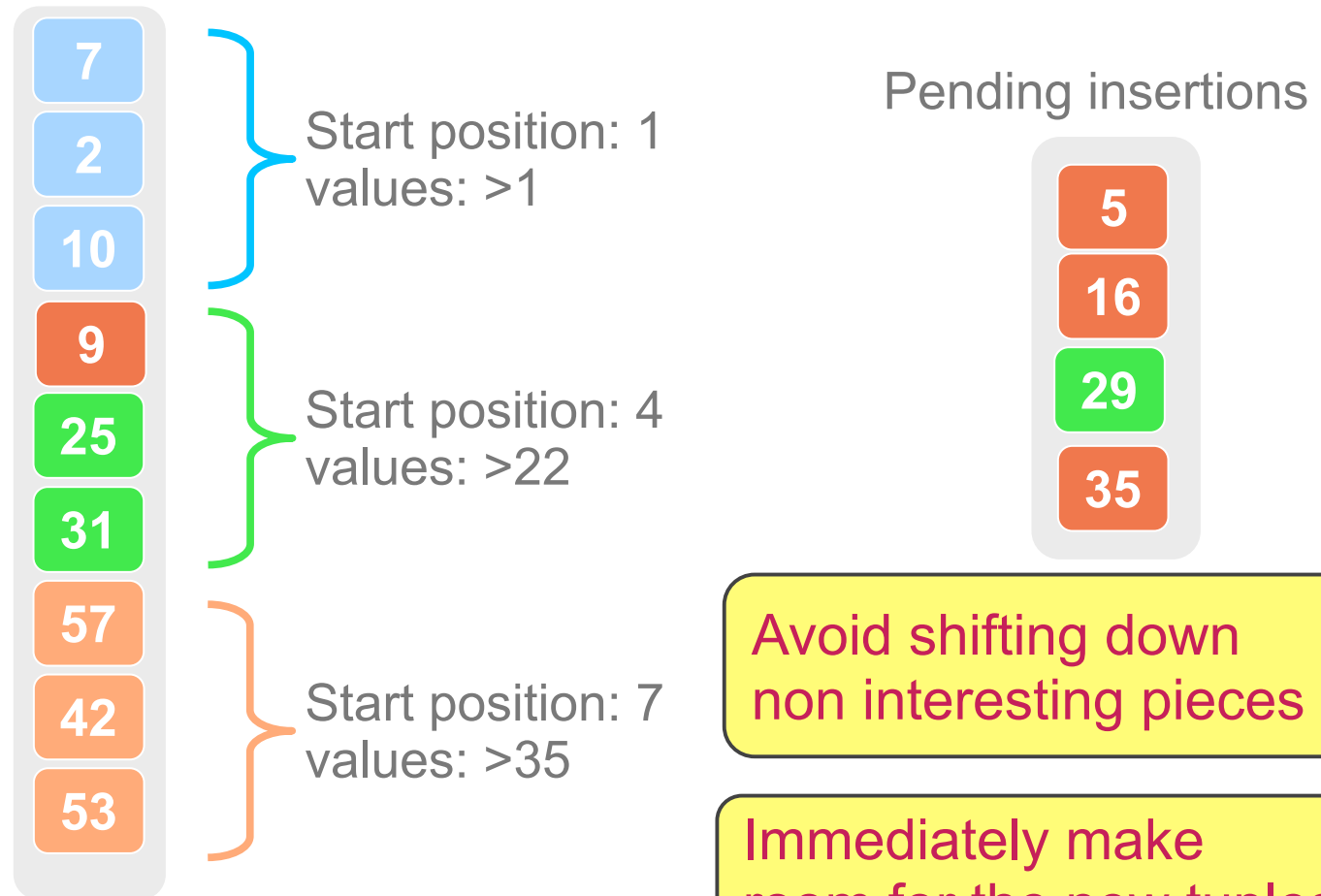
# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$



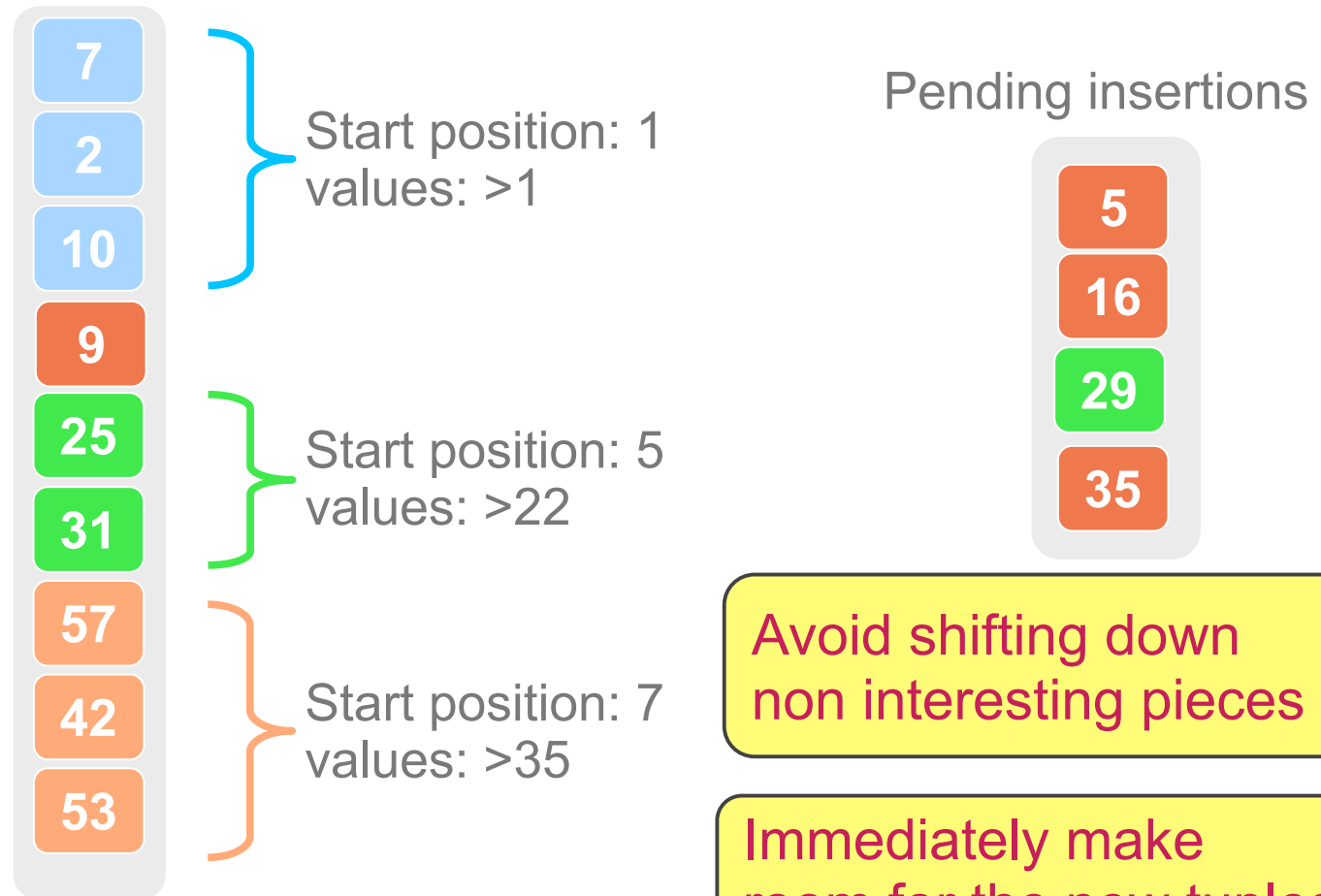
# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$



# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$



Avoid shifting down  
non interesting pieces

Immediately make  
room for the new tuples

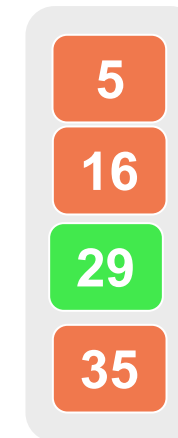


# The Ripple

Touch only the pieces that are **relevant** for the **current** query  
Select  $7 \leq A < 15$



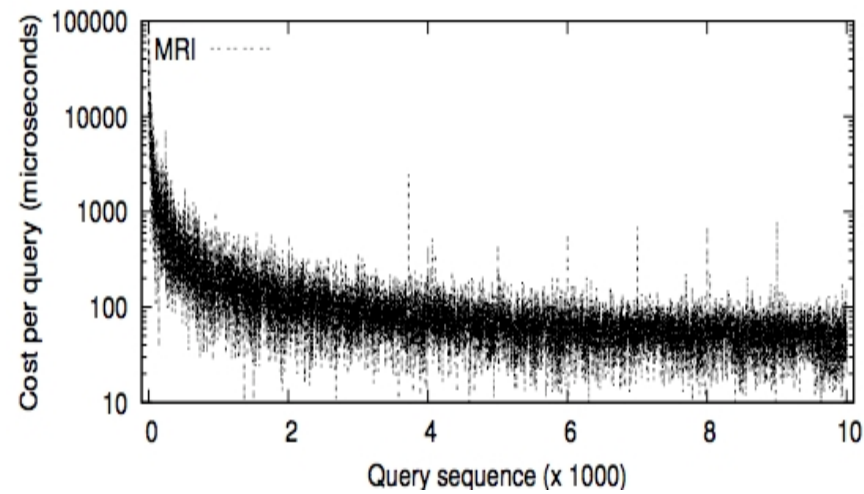
Pending insertions



Avoid shifting down  
non interesting pieces

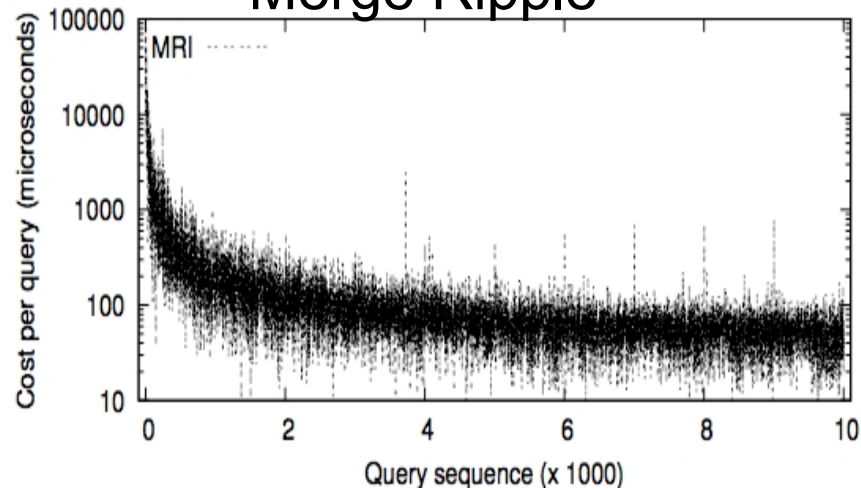
Immediately make  
room for the new tuples

# The Ripple



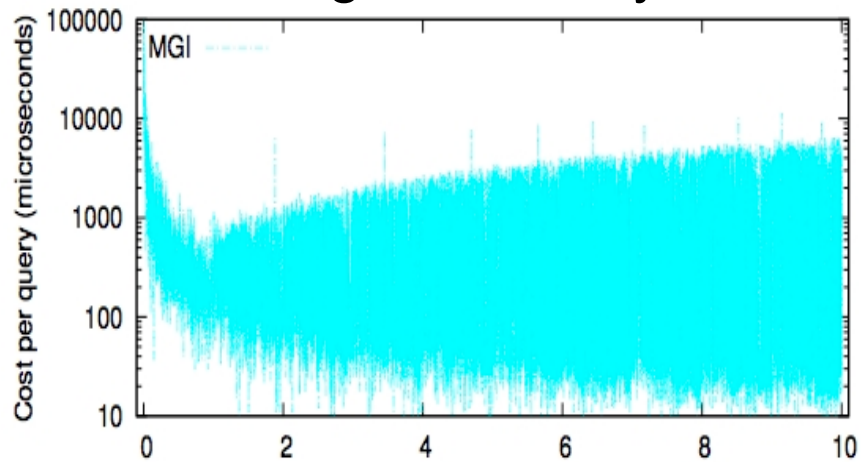
Maintain high performance through the whole query sequence in a self-organizing way

## Merge Ripple

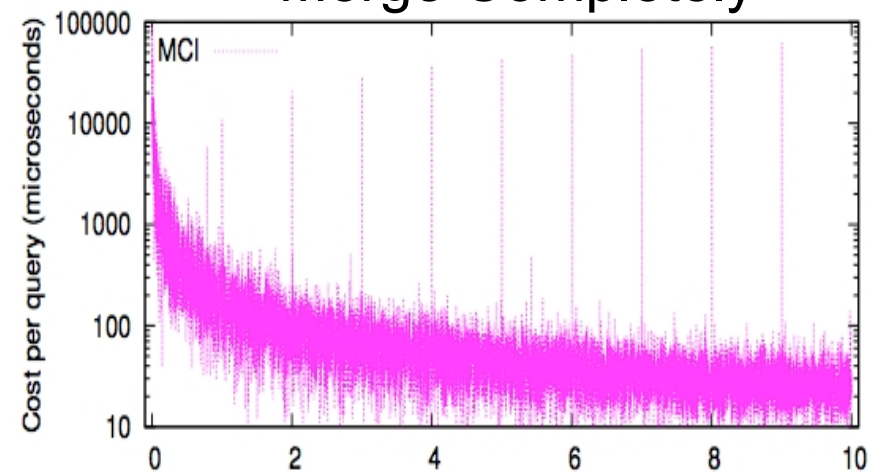


Maintain high performance through the whole query sequence in a self-organizing way

## Merge Gradually



## Merge Completely







# Recycling intermediates



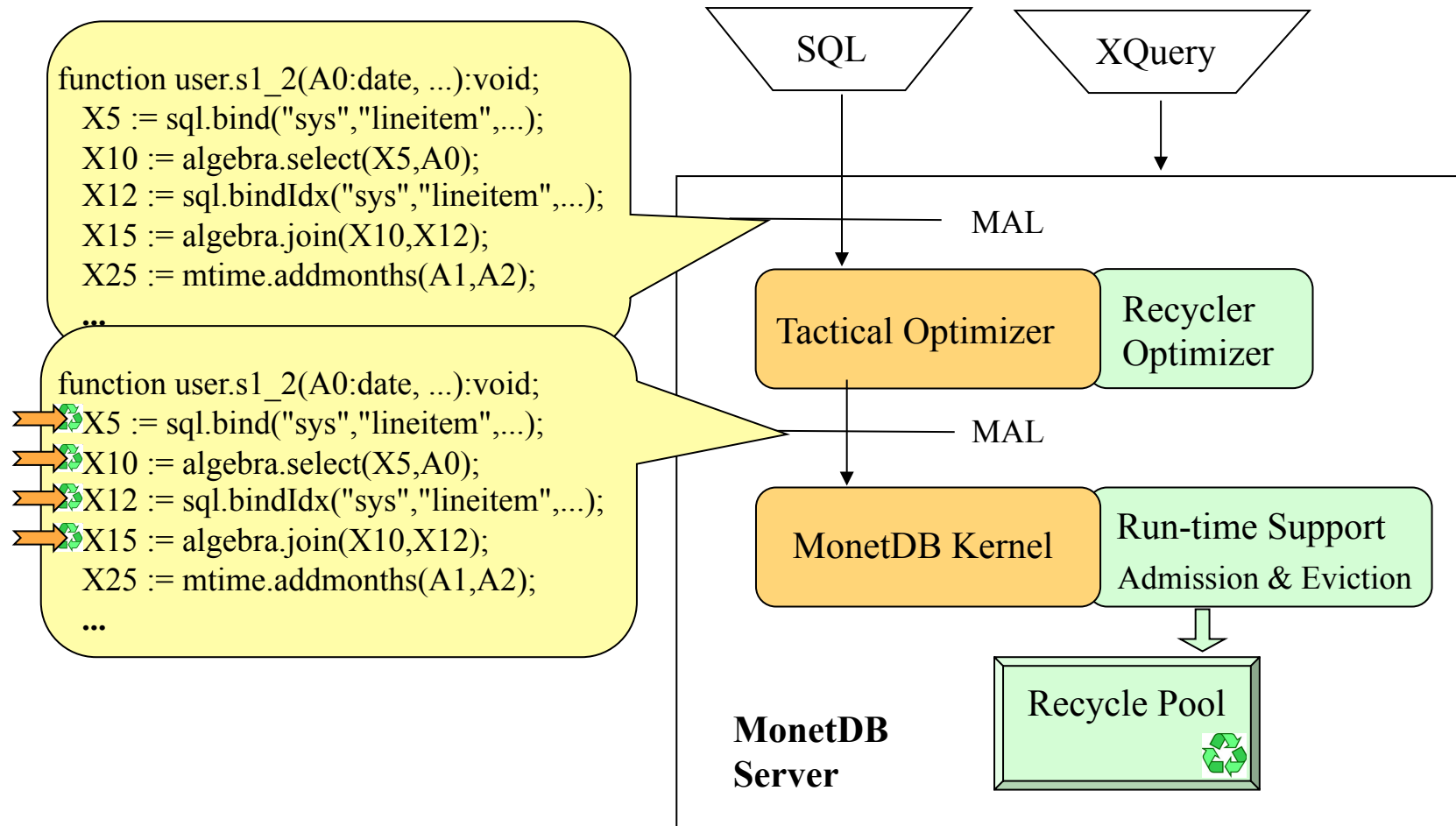
## MonetDB Background

- Operator-at-a-time execution paradigm
- Canonical implementation of a **column-store**
  - Reduced dimensionality
  - Finer granularity
  - Simplified overlap analysis
- Recycler extension of MonetDB engine





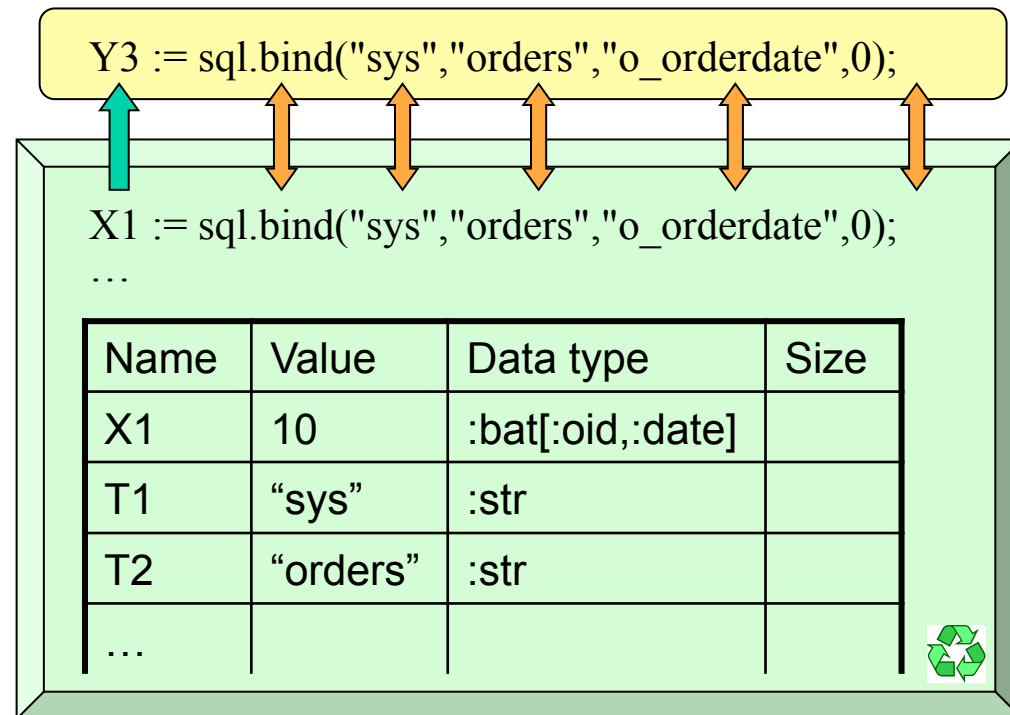
# MonetDB Architecture



Run time comparison of

- instruction types
- argument values

Exact  
matching



# Instruction Subsumption

Y3 := algebra.select(X1,20,45);

X3 := algebra.select(X1, 0,80)

...

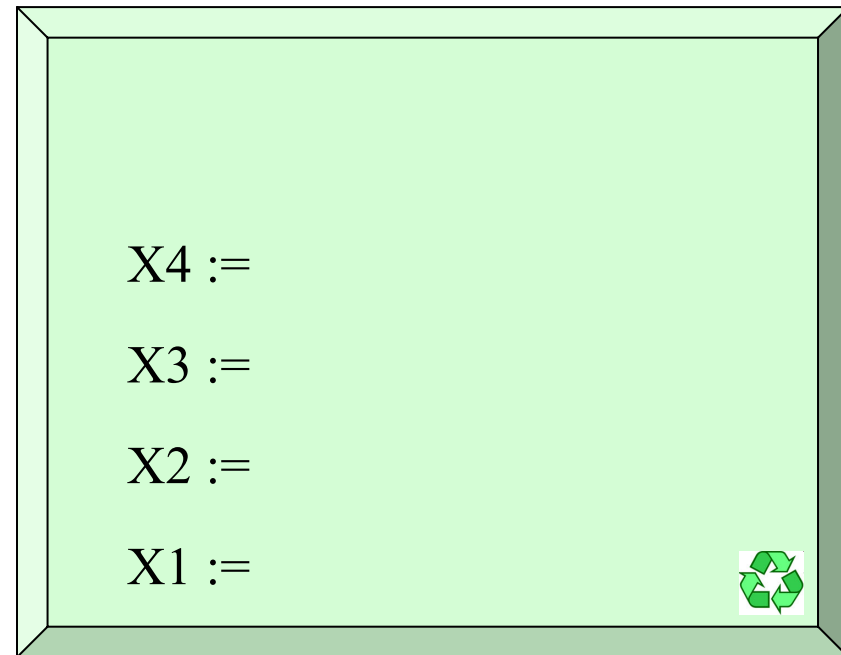
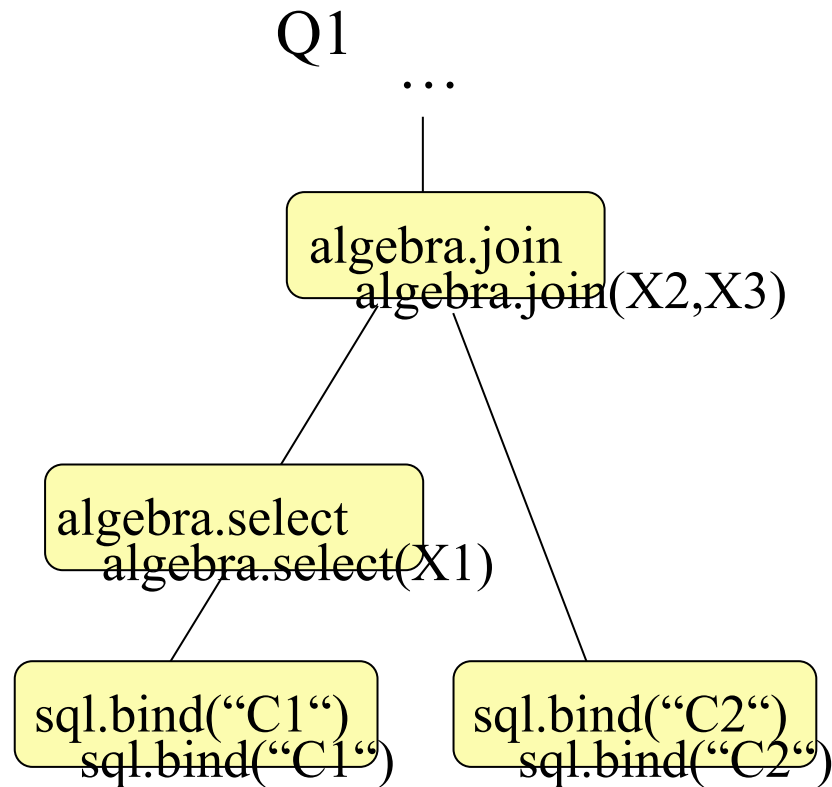
X5 = algebra.select(X1,20,60);

Name	Value	Data type	Size
X1	10	:bat[:oid,:int]	2000
X3	130	:bat[:oid,:int]	700
X5	150	:bat[:oid,:int]	350
...			



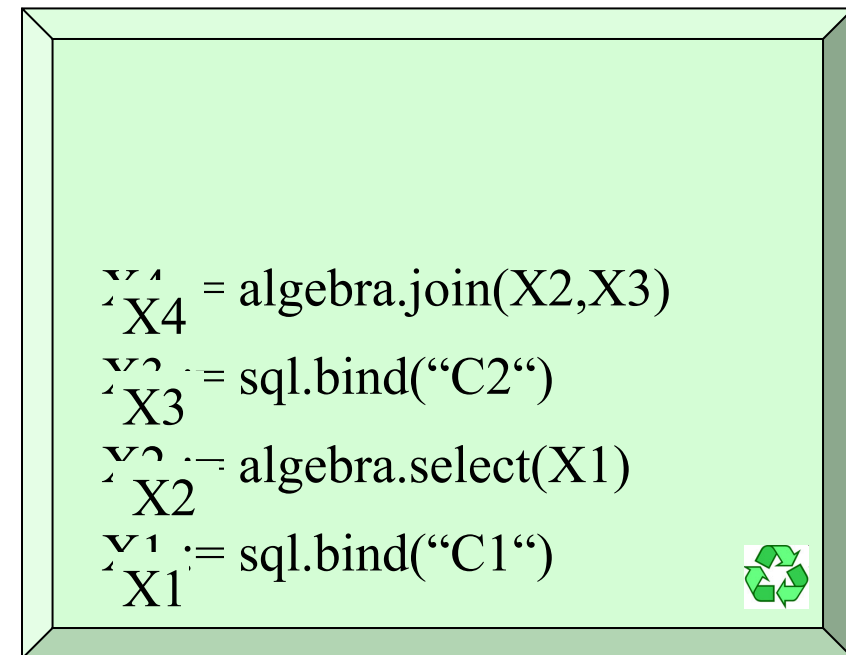
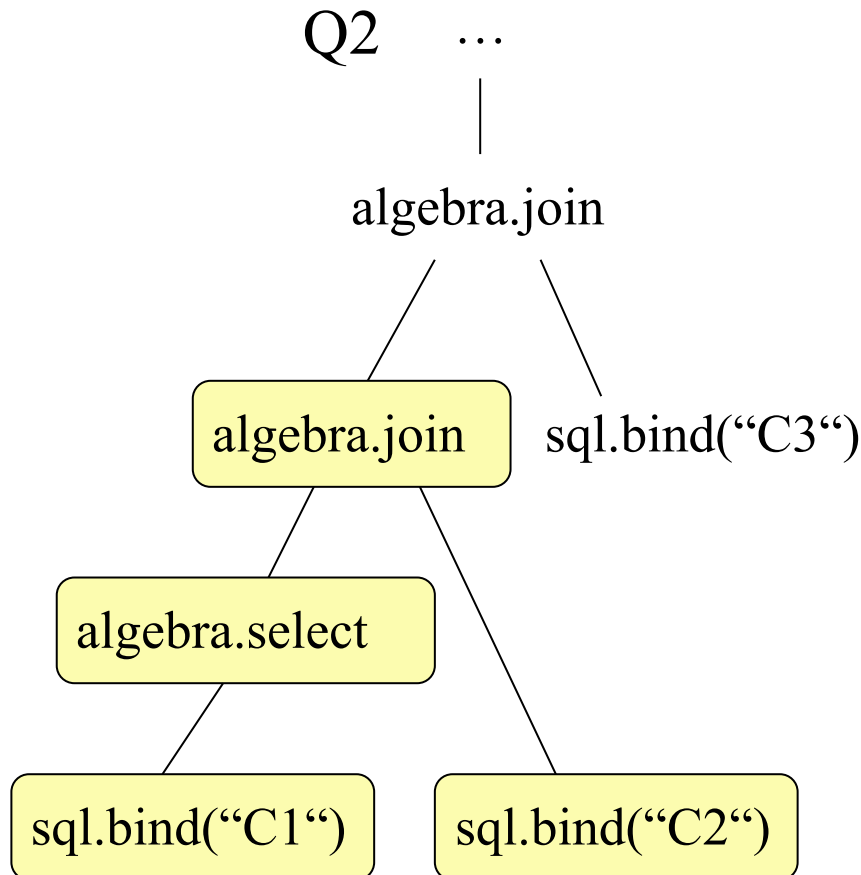


# Recycle Pool: a Cache with Lineage

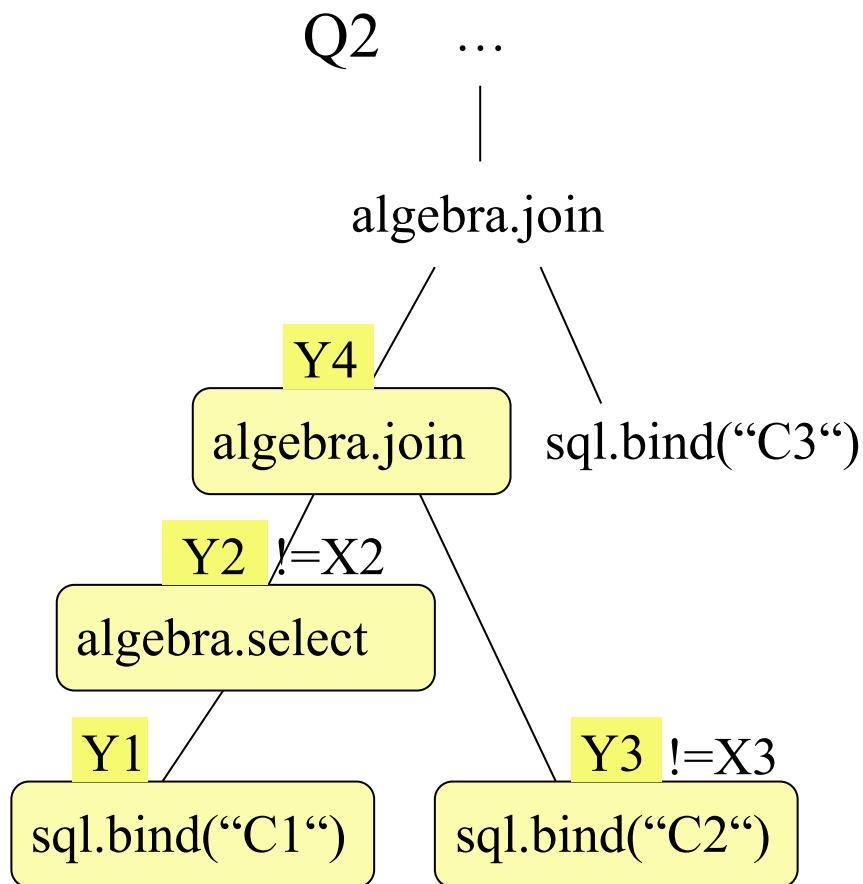




# Recycle Pool: a Cache with Lineage

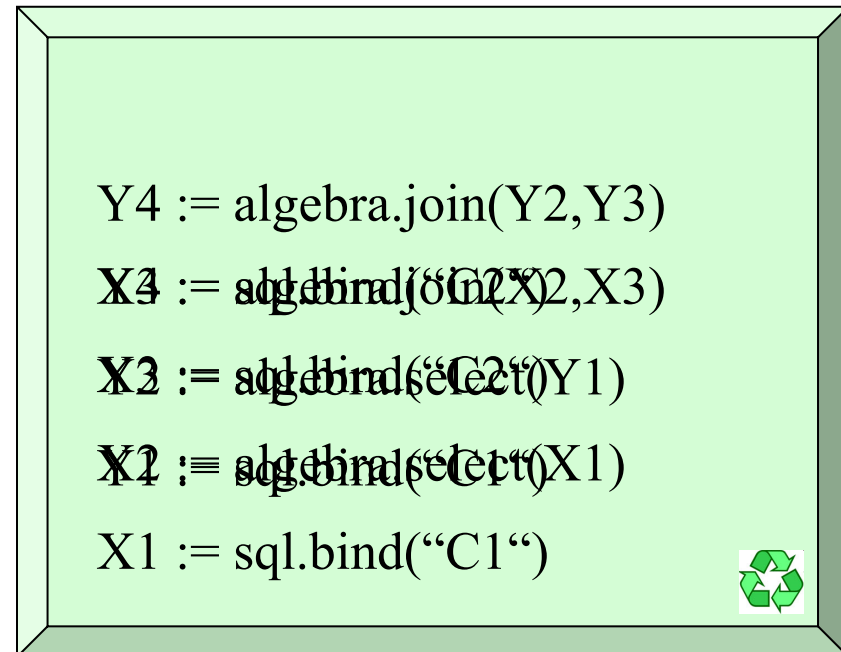


# Mismatching



```

Y4 := algebra.join(Y2, Y3)
X4 := algebra.join(Y2, X3)
X3 := algebra.select(Y1)
X2 := algebra.select(X1)
X1 := sql.bind("C1")
  
```







## Admission Policies

Decide about storing the results

- **KEEPALL**
  - all instructions advised by the optimizer
- **CREDIT**
  - instructions supplied with credits
  - storage 'paid' with 1 credit
  - reuse returns credits
  - lack of reuse limits admission and resource claims



## Cache Policies

- Decide about eviction of intermediates
- Filter 'top' instructions without dependents
- Pick instructions with smallest utility
  - LRU : time of computation or last reuse
  - BENEFIT : estimated contribution to performance: CPU and I/O costs, recycling
- Triggered by resource limitations (memory or entries)



# TPC-H Evaluation

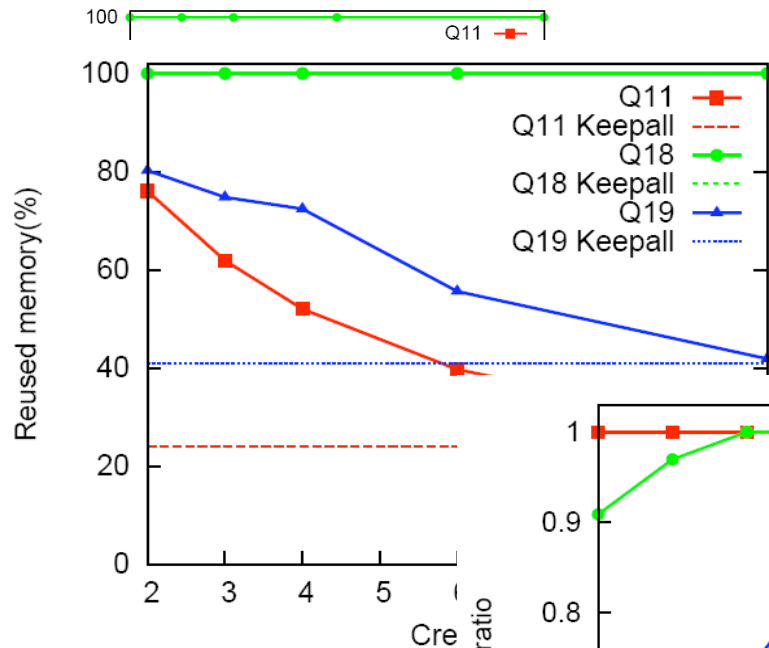
- SF1
- Baseline performance
- Impact of design choices
  - Admission policies
  - Cache policies



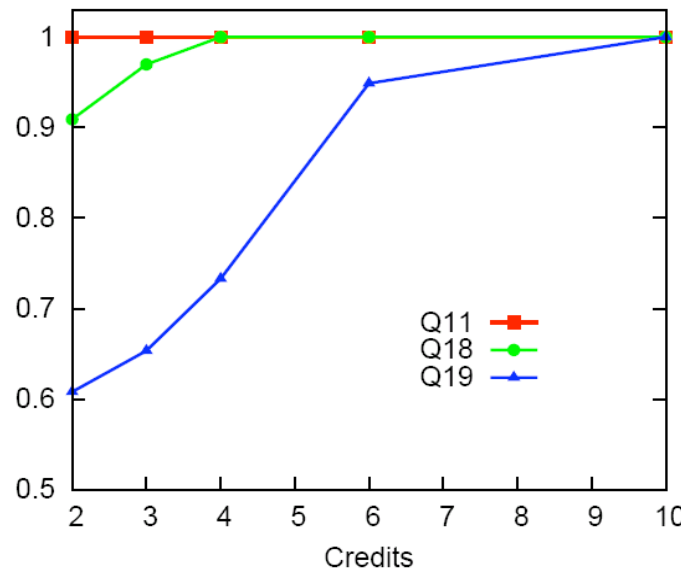
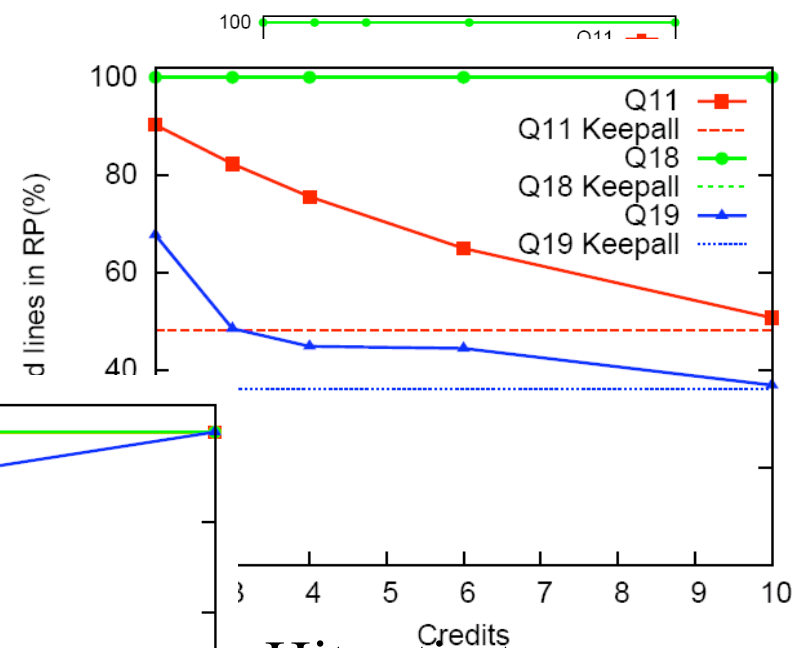


# CREDIT Admission Impact

## Reused memory



## Reused entries



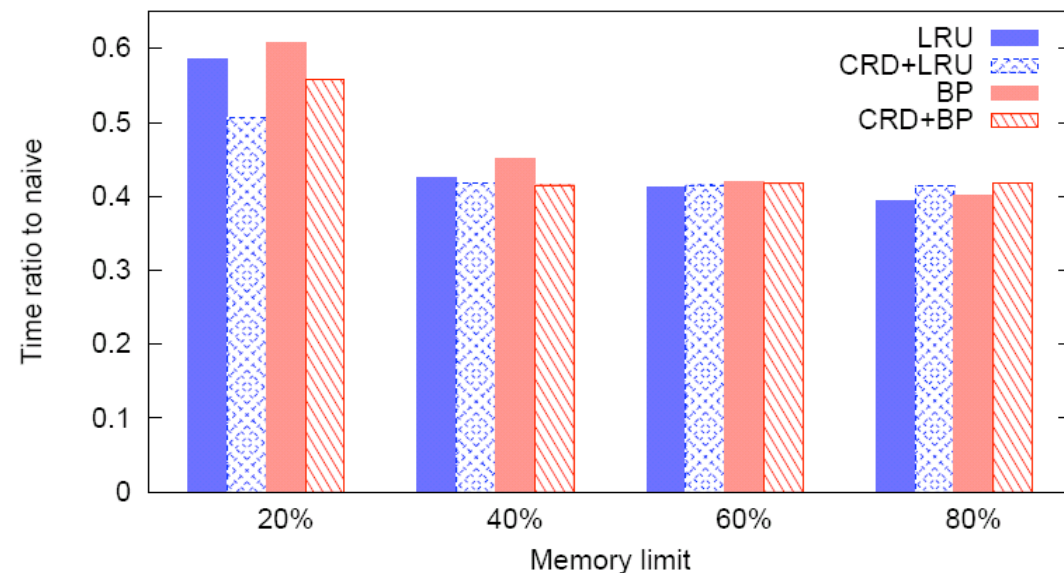
Hit ratio to KeepAll



# Cache Policies Evaluation

200 TPC-H queries

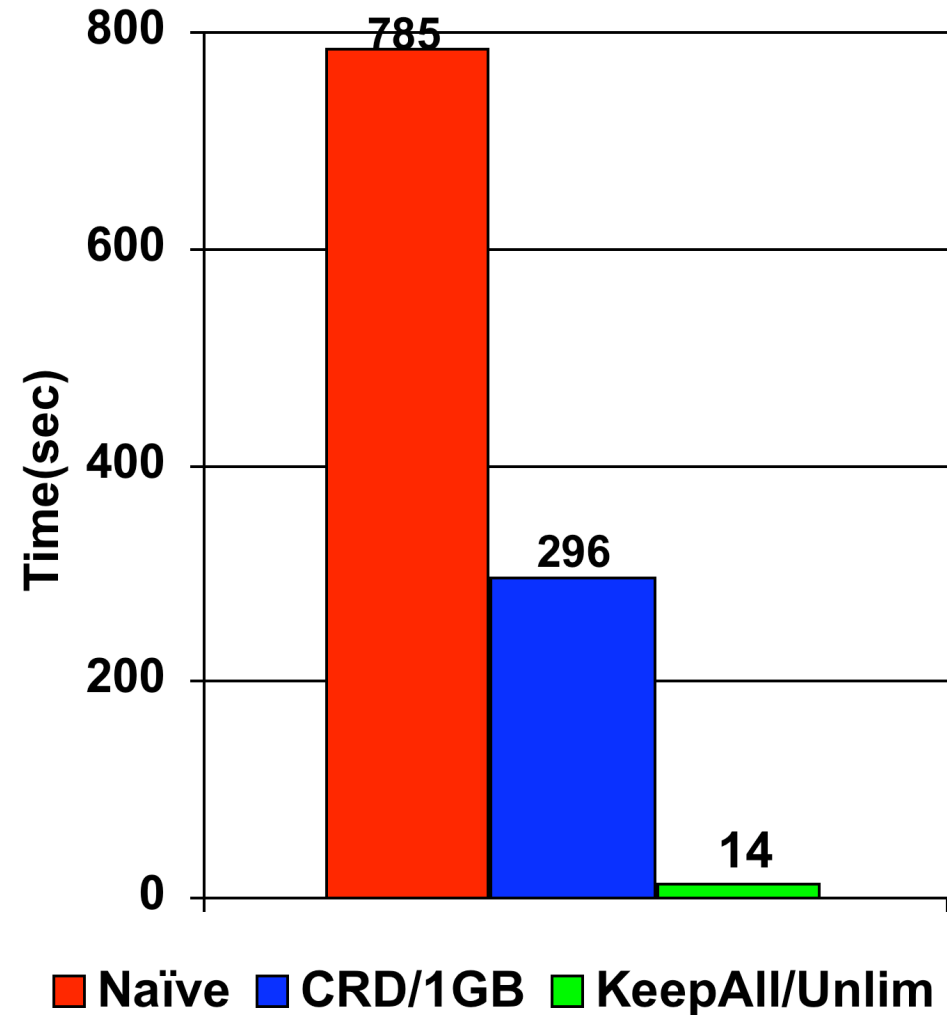
RP	Total	Reuse
Memory	4GB	42.7%
Entries	5219	28%





# SkyServer Evaluation

- 100 GB subset of DR4
- 100-query batch from January 2008 log
- 1.5GB intermediates, 99% reuse
- Join intermediates major contributor to savings





## Summary

- Database architecture augmented with recycling intermediates
- Significant performance benefits in SkyServer and TPC-H
- Self-organizing technique
- Extension to MonetDB transforming materialization overhead into benefit





## Future Work

- Refining and developing admission and cache policies
- Opportunities by query class recognition
- Automatic switch to suitable policies
- Application to pipelined architectures





# Recycling



Green

Is

