

Column-Oriented Databases

Themis Palpanas
University of Paris

Data Intensive and Knowledge Oriented Systems



1



- thanks for slides to
 - Christoph Freytag
 - Stratos Idreos
 - Daniel Abadi

2.2

2

Overview

- ☐ Physical layout of data on pages
- ☐ Tuple Cracking
- ☐ Example Systems

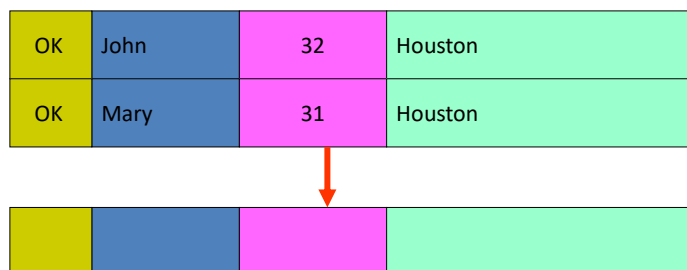


8.3

3

Physical data organization

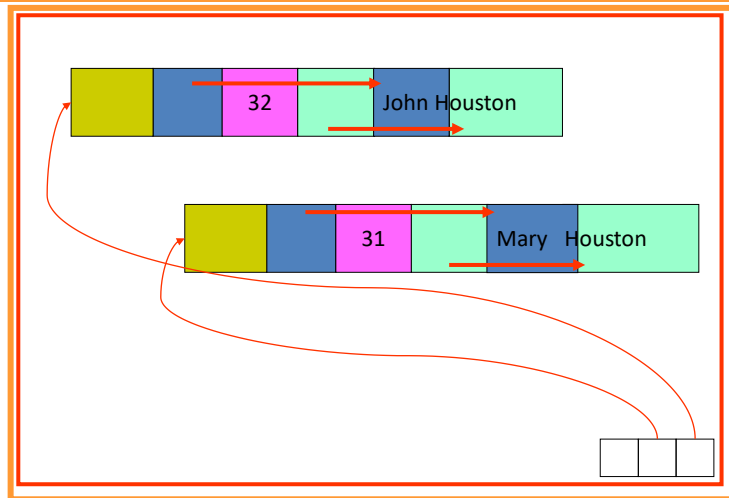
- Early 80s: Tuple oriented storage in most DBMS
 - Easy to access at the cost of wasted space



8.4

4

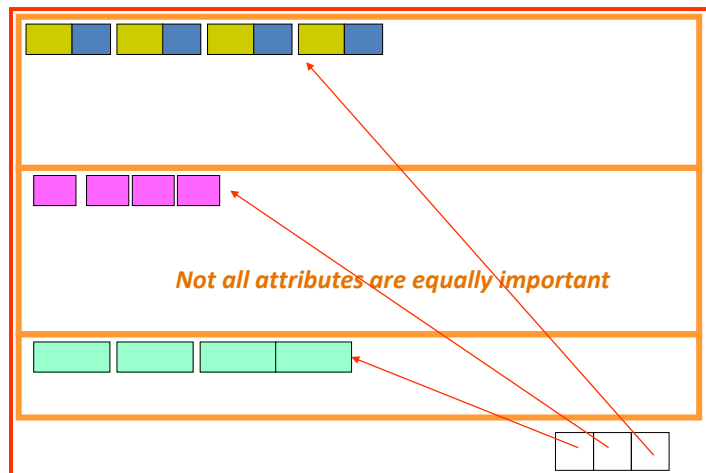
Physical data organization (2)



8.5

5

Physical data organization (3)



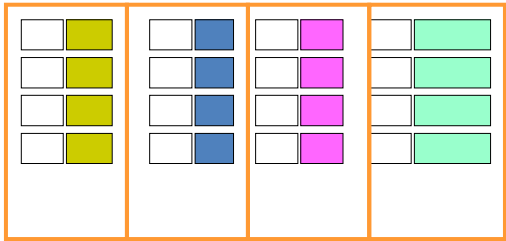
8.6

6



Alternative - Column Store

- A column orientation is as simple and acts like an array
- Attributes of a tuple are correlated by offset or ID



8. 7

7



Row vs. Column-Stores

- Example Customer Relation

Row-Store

Last Name	First Name	E-mail	Phone #	Street Address

Column-Store

Last Name	First Name	E-mail	Phone #	Street Address

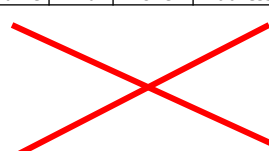
8. 8

8

Example 1 – Simple Selection

- SELECT LastName FROM Customer

Last Name	First Name	E-mail	Phone #	Street Address



- Only relevant data are read
 - Found densely stored on page(s)

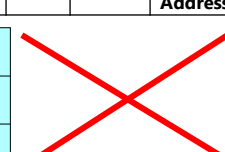
8. 9

9

Example 2 – More Complex Selection

- SELECT LastName, FirstName FROM Customer

RID	Last Name	First Name	E-mail	Phone#	Street Address
1					
2					
3					



- Read from two different file
 - How to bring both columns together?
 - Join needed (not specified in the query!)
 - Internal ID (RID) needed

8. 10

10

Column Store – Processing needs

- Different optimization for column oriented database
 - Compression
 - Late Materialization
 - Late Joins with Dictionary and between columns
 - Block Iteration:
 - Process all entries per block
 - Block oriented Join
 - “Invisible” Join

8. 13

13

Column oriented storage

- Advantage
 - Each column can be more easily optimized for storage using compression schemes
 - Each column can be replicated for read-only access
 - Brings a much tighter packaging and improves transport through the memory hierarchy (still to prove!)
 - Column orientation benefits data warehousing

8. 14

14

Row vs. Column-Stores

Row Store

- Easy to add a new record
- Might read in unnecessary data
- Logical group equals physical grouping

Column Store

- Only need to read in relevant data
- Tuple writes might require multiple seeks
- More Joins needed
- High compression possible

Column Store Systems

- Used in
 - Monet DB (CWI, Amsterdam):
 - X-Store/Vertica (MIT, Stonebraker)
 - Vertica now owned by HP – see www.vertica.com
 - SAP Sysbase IQ: for a long time

Tuple Cracking (from Monet DB)

17

17

Tuple Cracking (introduced by Monet)

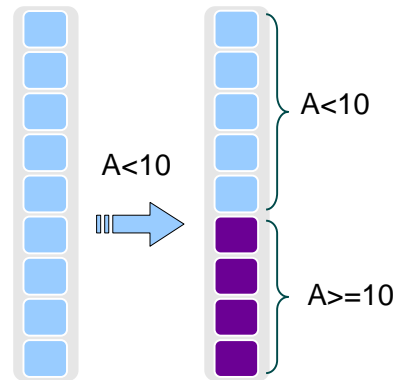
- Basic idea
 - Try to use previous access pattern to optimize for the future
 - Constantly reorganize
 - Challenge: How to do it with minimal overhead

8. 18

18

Cracking algorithms

- Physical reorganization happens *per column* based on selection predicates
- Split a piece of a column in *two* new pieces

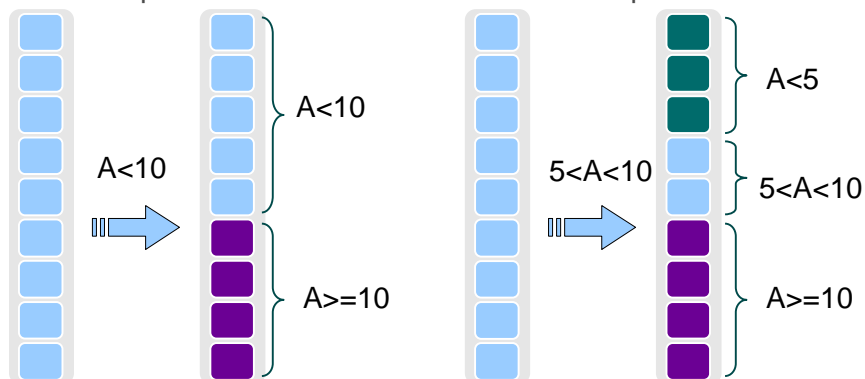


8. 19

19

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



8. 20

20

Cracking algorithms

select $A > 5$ and $A < 10$

17
3
8
6
2
12
13
4
15

8. 21

21

Cracking example

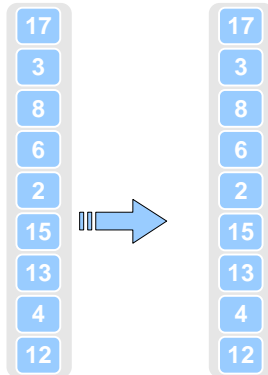
select $A > 5$ and $A < 10$

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

8. 22

22

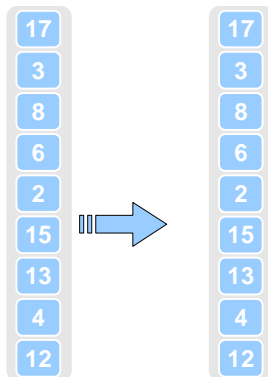
Cracking example

select $A > 5$ and $A < 10$ we need all values < 5 at the beginning of the table

8. 23

23

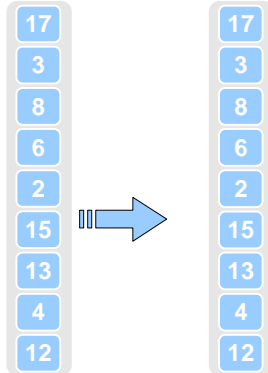
Cracking example

select $A > 5$ and $A < 10$ we need all values < 5 at the beginning of the tablewe need all values > 10 at the end of the table

8. 24

24

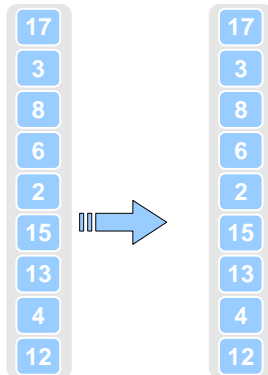
Cracking example

select $A > 5$ and $A < 10$ we need all values < 5 at the beginning of the tablewe need all values > 5 AND < 10 in the middle of the tablewe need all values > 10 at the end of the table

8. 25

25

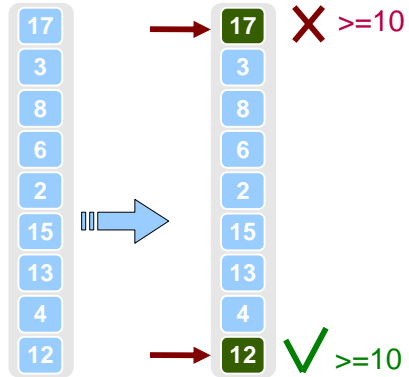
Cracking example

select $A > 5$ and $A < 10$ 

8. 26

26

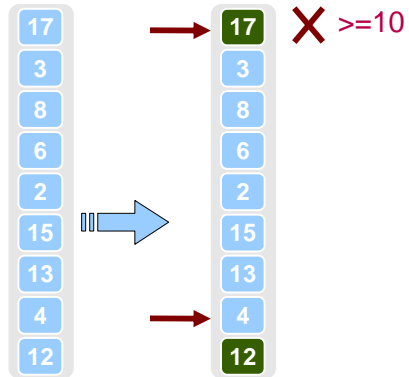
select A>5 and A<10



8. 27

27

select A>5 and A<10



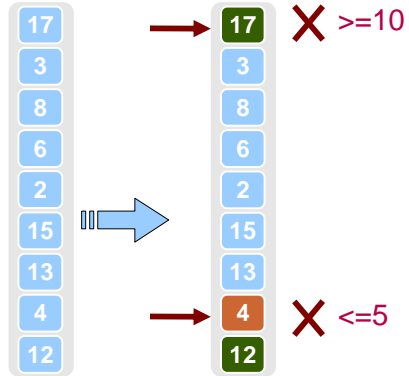
8. 28

28

Cracking example



select A>5 and A<10



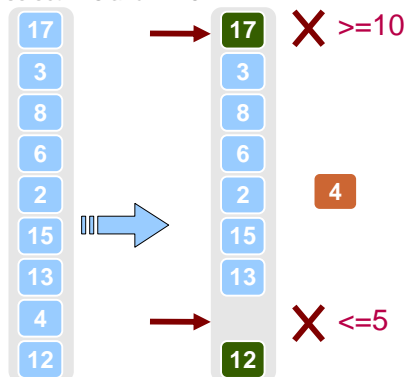
8. 29

29

Cracking example



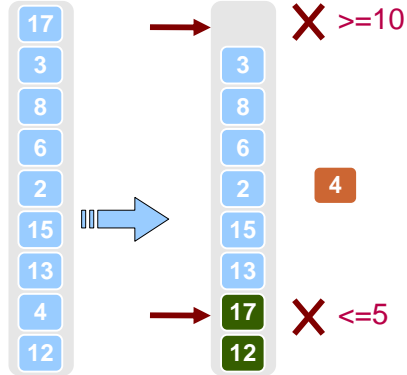
select A>5 and A<10



8. 30

30

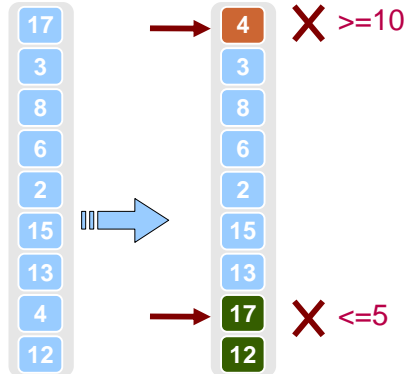
select A>5 and A<10



8. 31

31

select A>5 and A<10

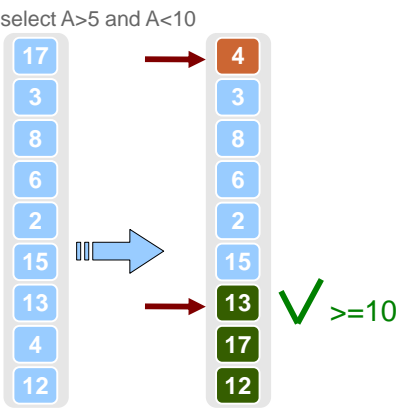


8. 32

32



Cracking example

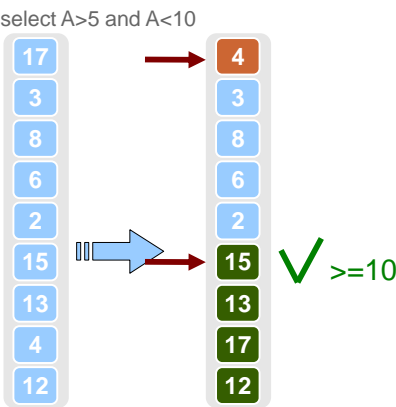


8. 33

33



Cracking example



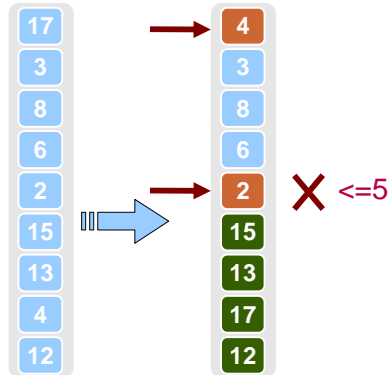
8. 34

34

Cracking example



select A>5 and A<10



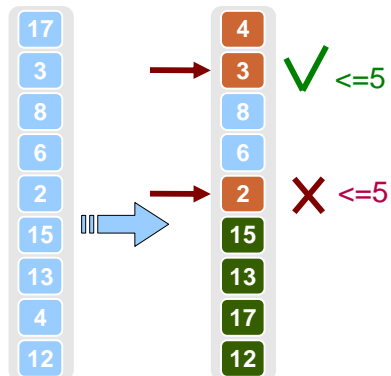
8. 35

35

Cracking example



select A>5 and A<10



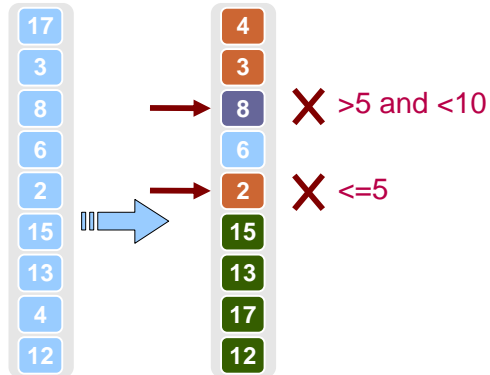
8. 36

36

Cracking example



select A>5 and A<10



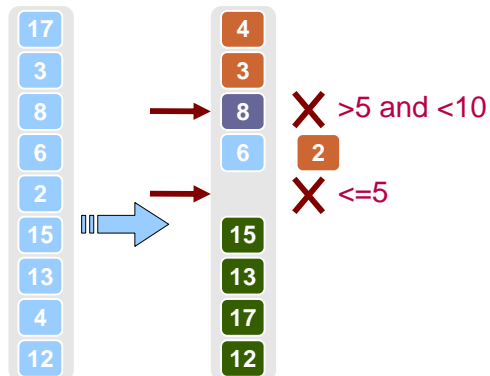
8. 37

37

Cracking example



select A>5 and A<10



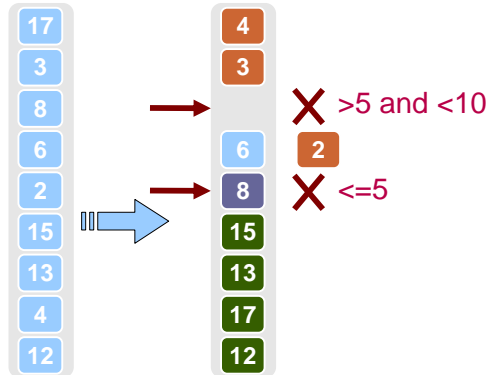
8. 38

38

Cracking example



select A>5 and A<10



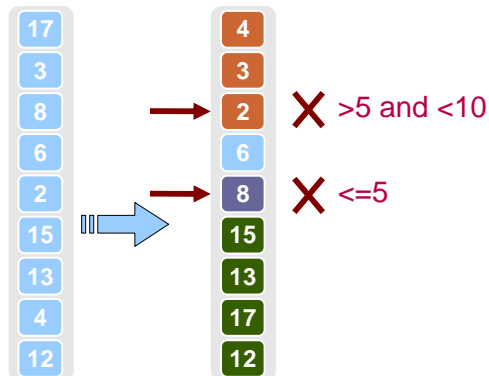
8. 39

39

Cracking example



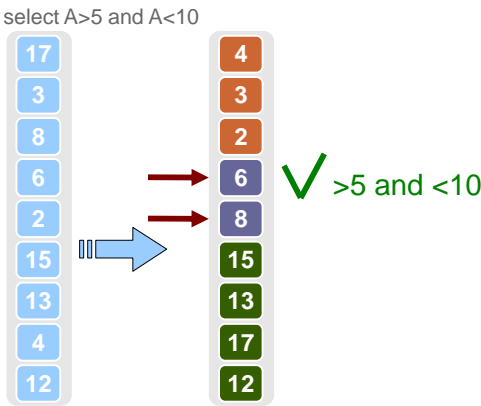
select A>5 and A<10



8. 40

40

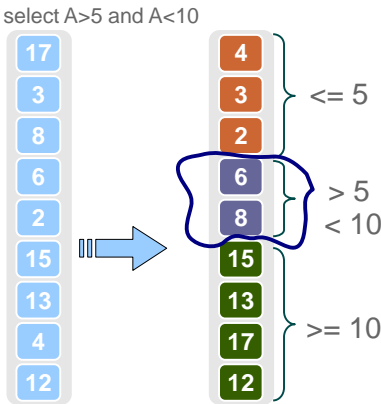
Cracking example



8. 41

41

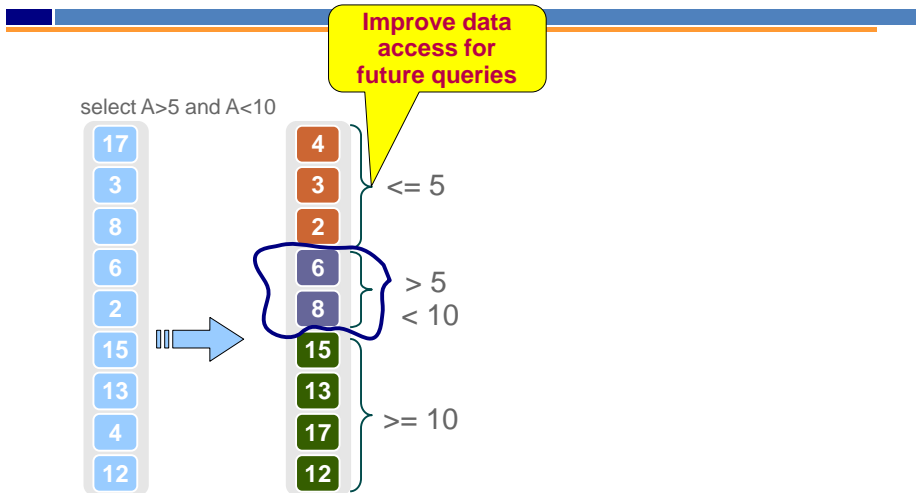
Cracking example



8. 42

42

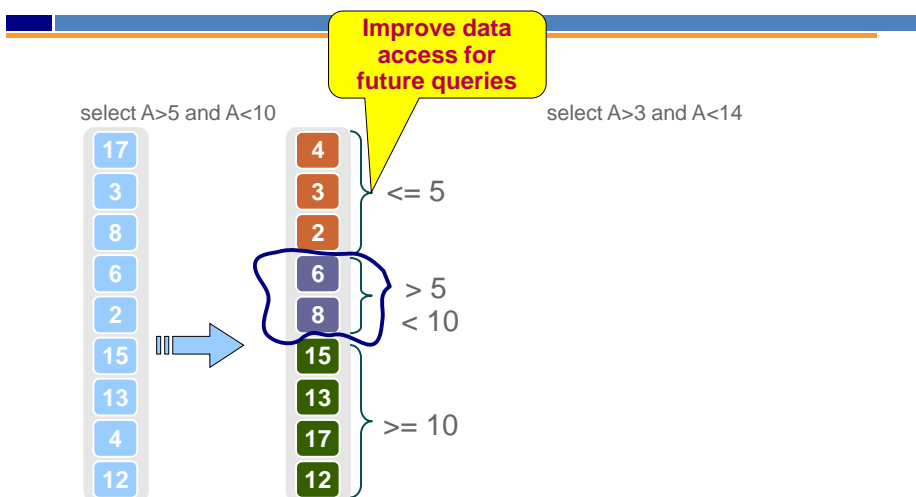
Cracking example



8. 43

43

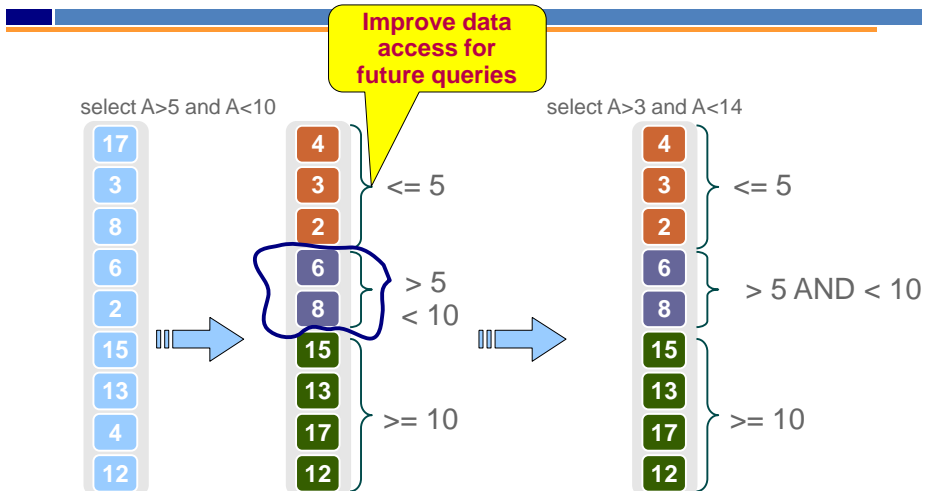
Cracking example



8. 44

44

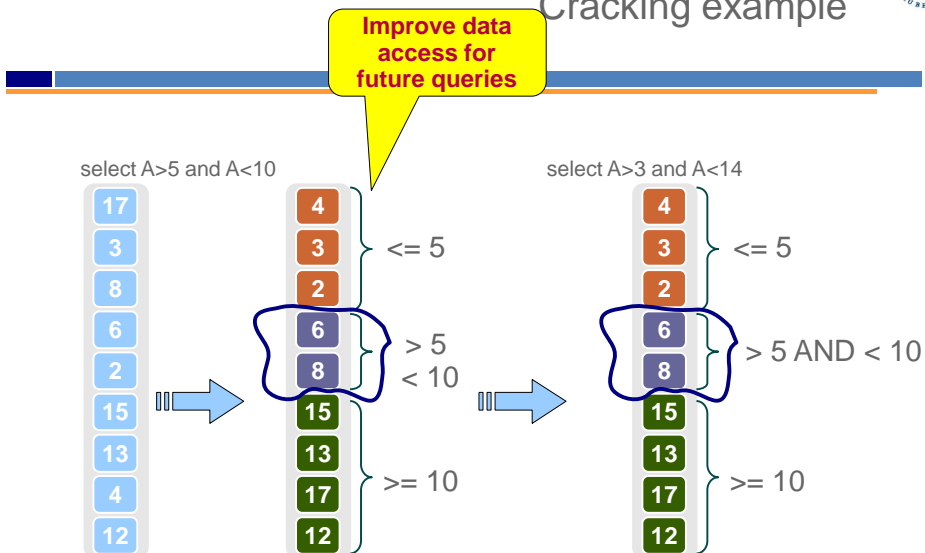
Cracking example



8. 45

45

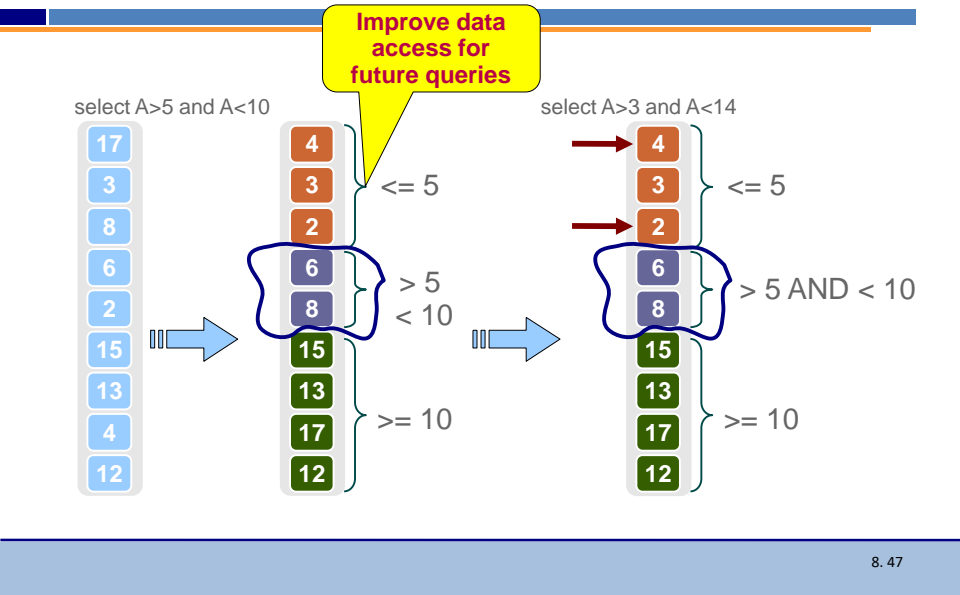
Cracking example



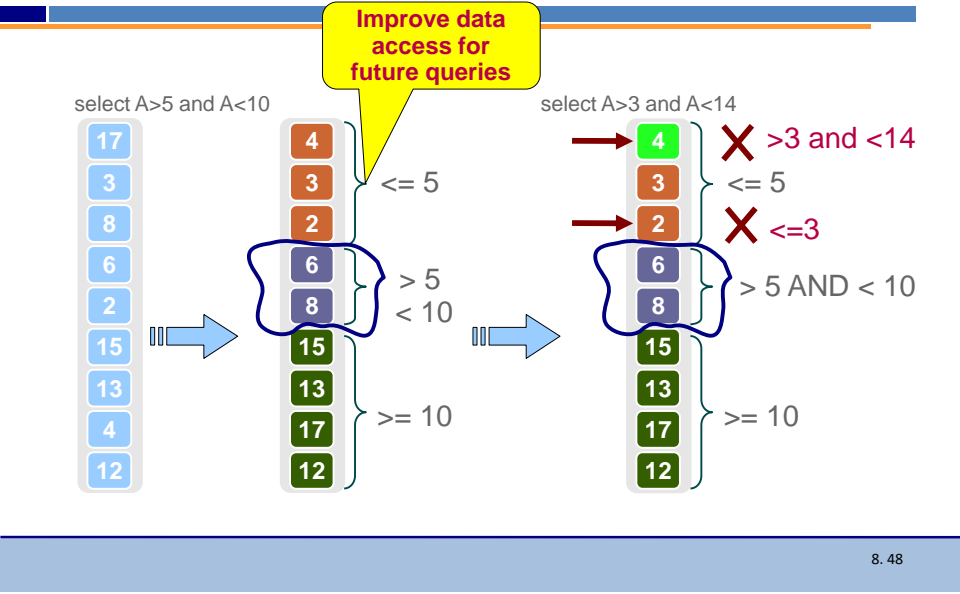
8. 46

46

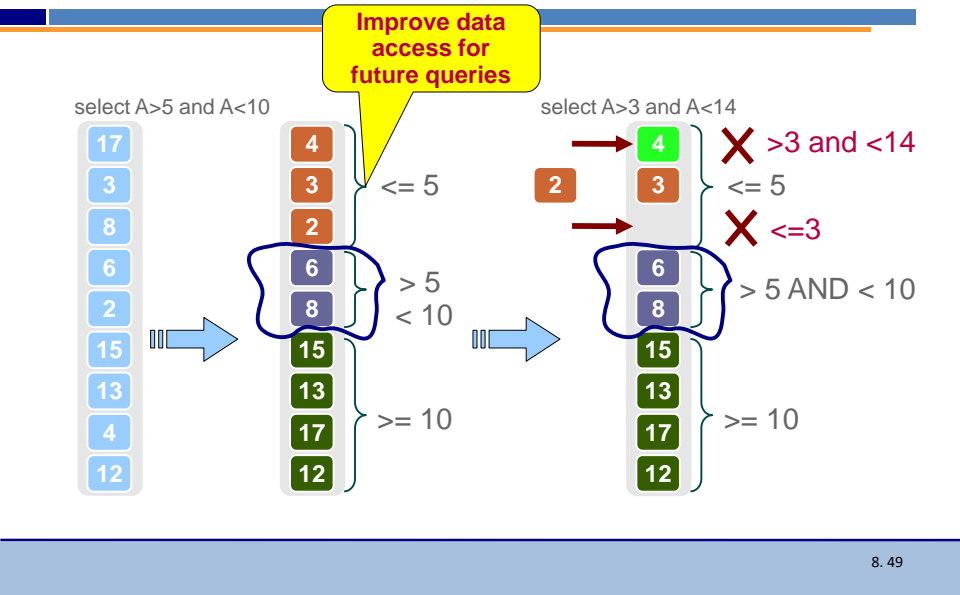
Cracking example



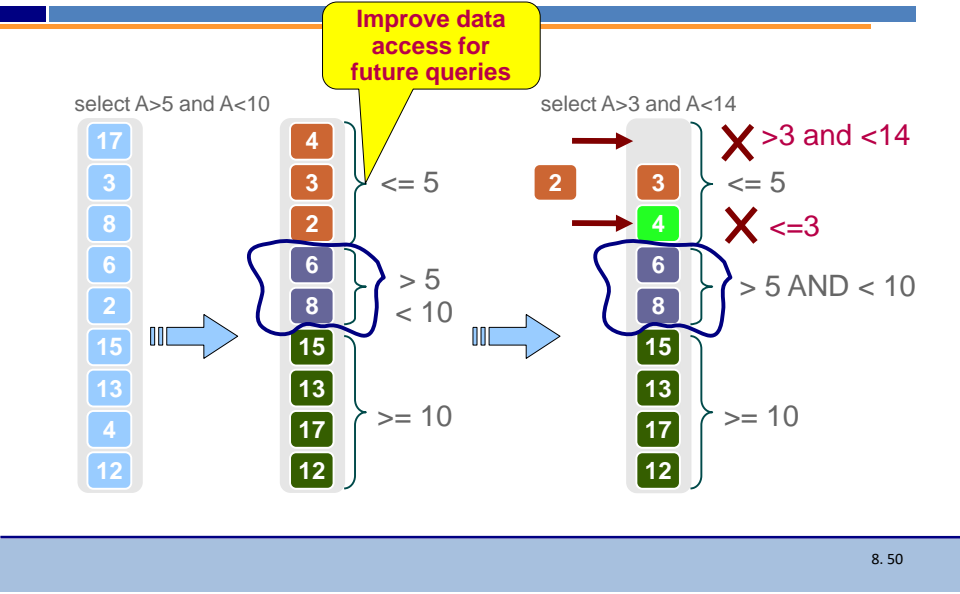
racking example



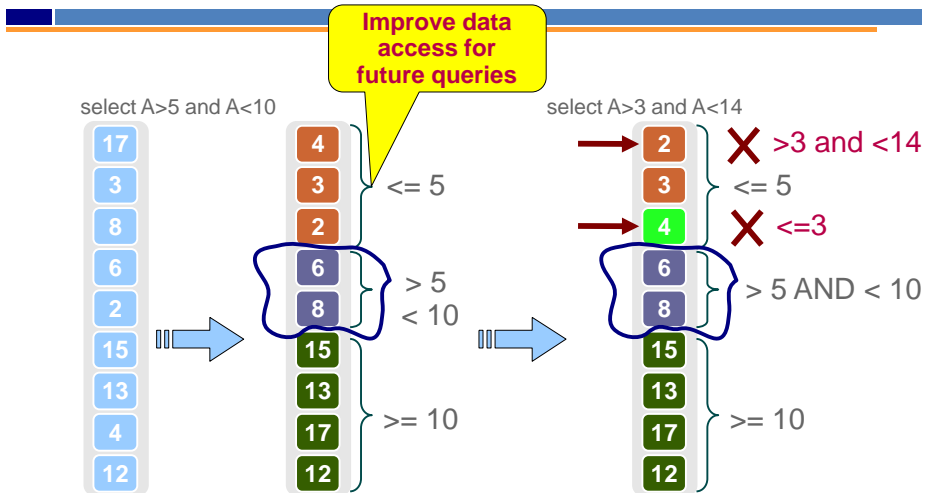
Cracking example



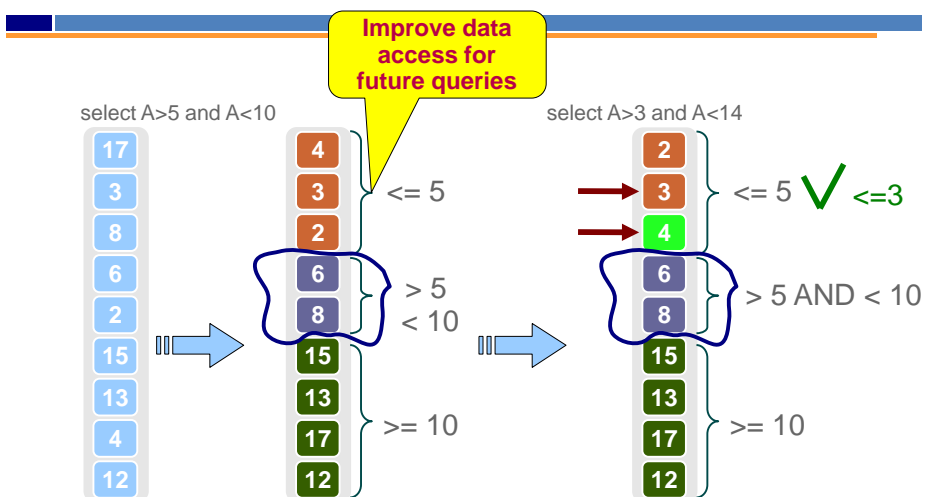
Cracking example



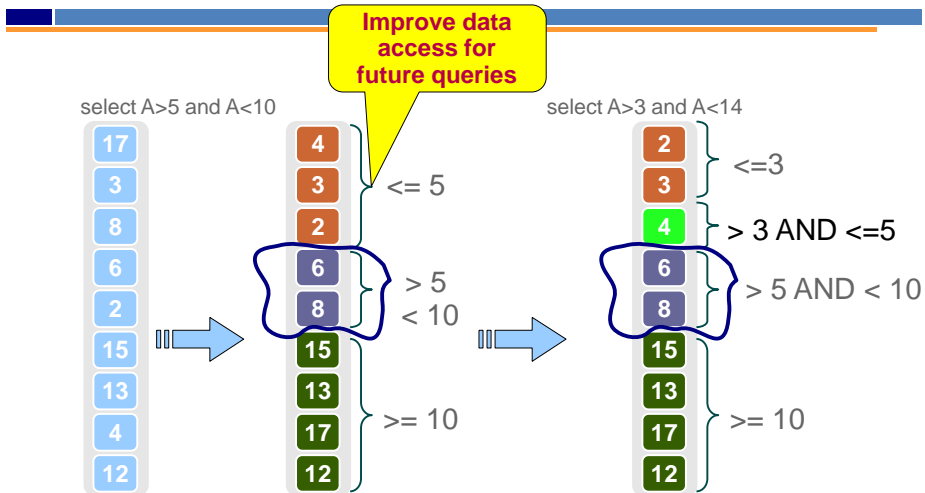
Cracking example



Cracking example



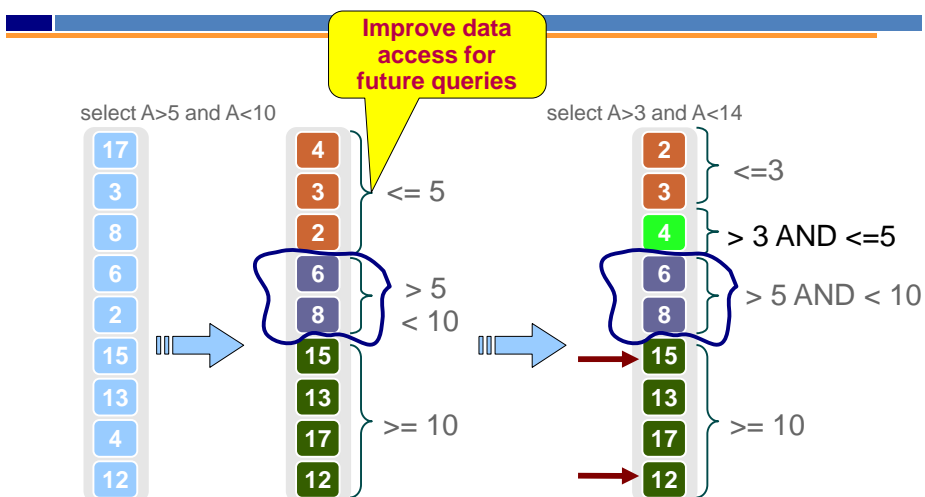
Cracking example



8. 53

53

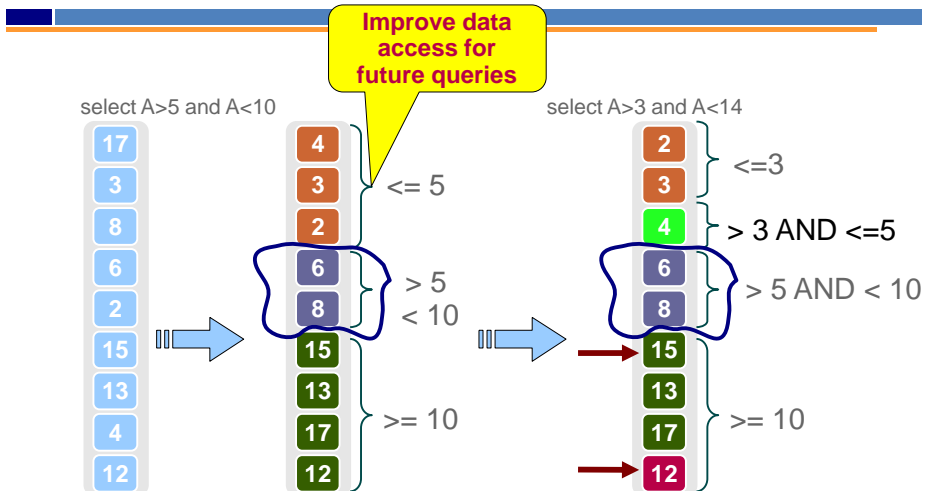
Cracking example



8. 54

54

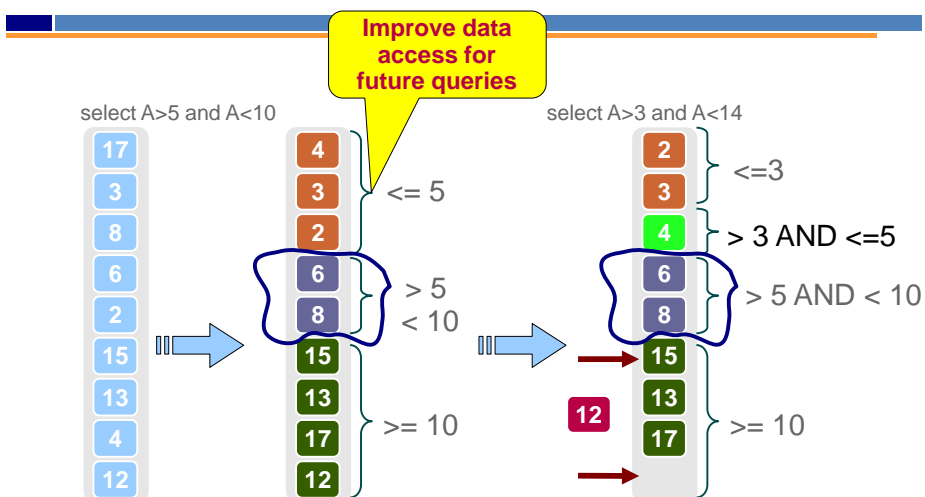
Cracking example



8. 55

55

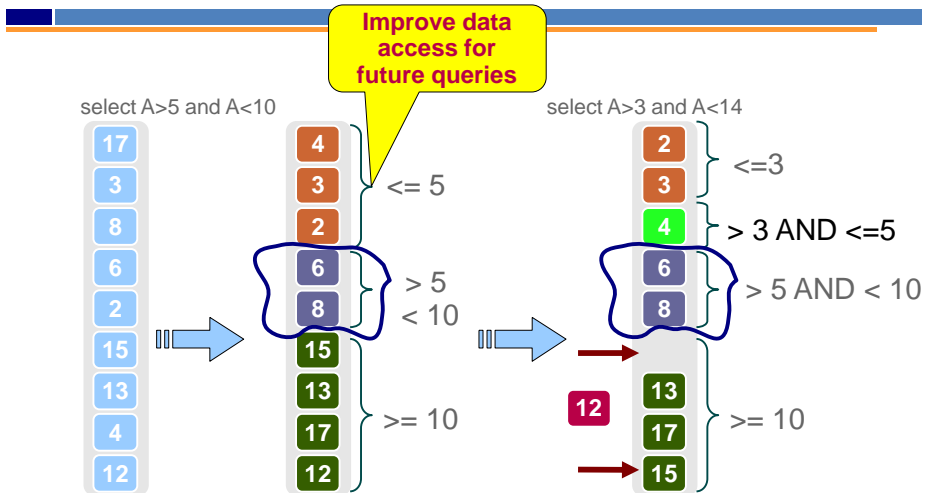
Cracking example



8. 56

56

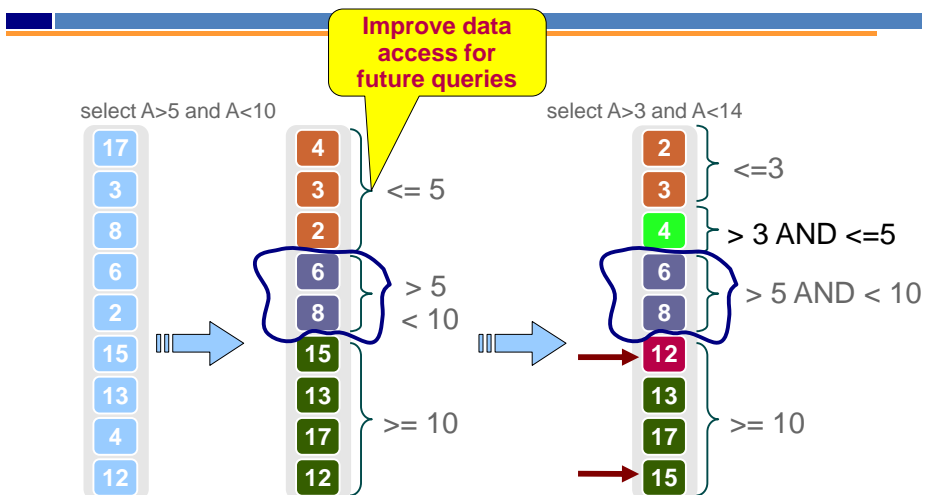
Cracking example



8. 57

57

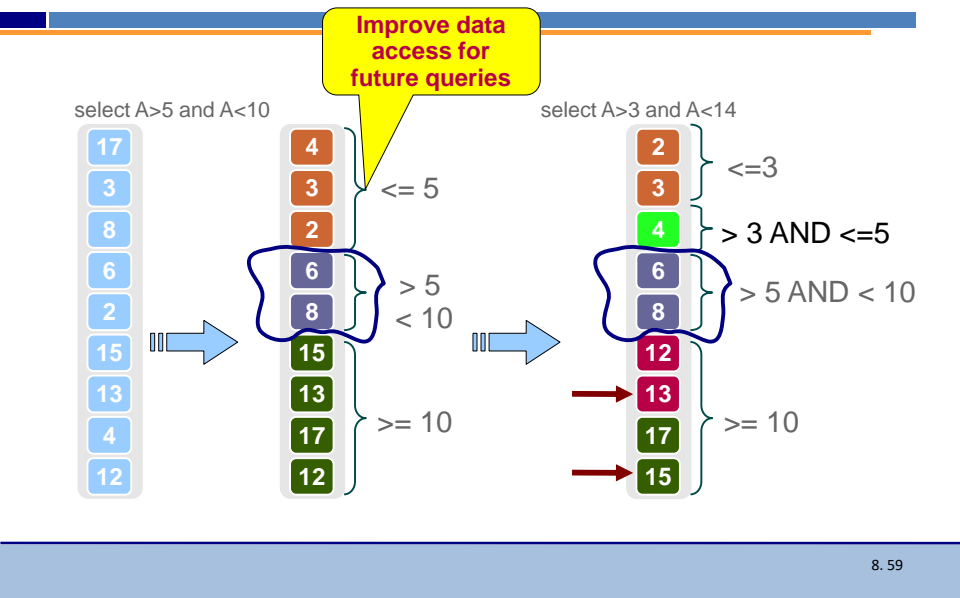
Cracking example



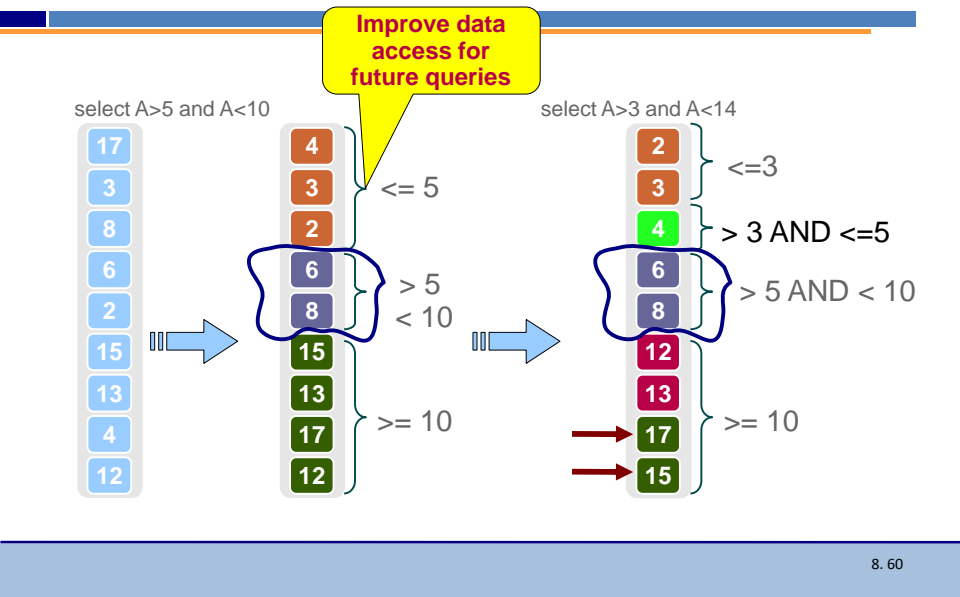
8. 58

58

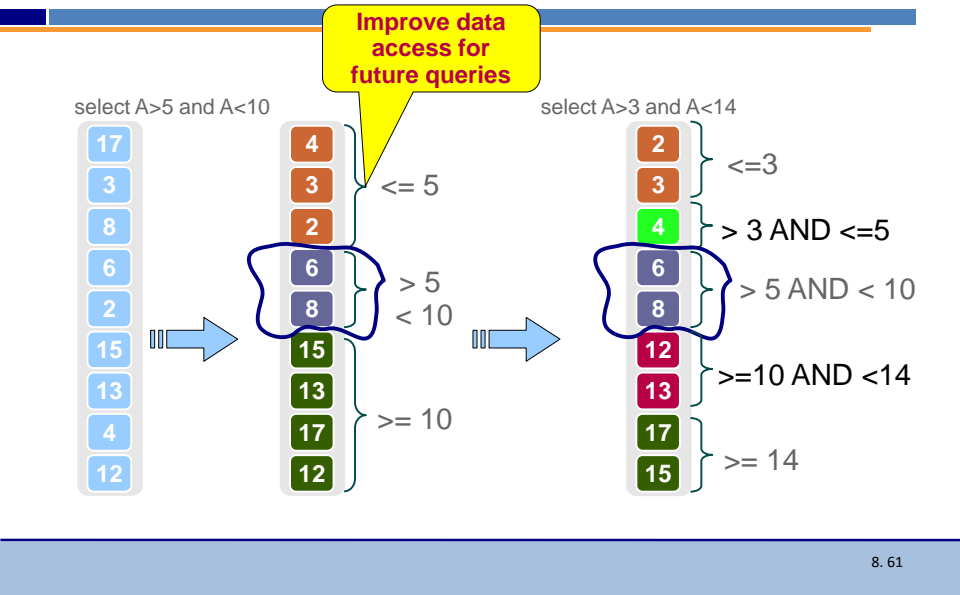
Cracking example



Cracking example

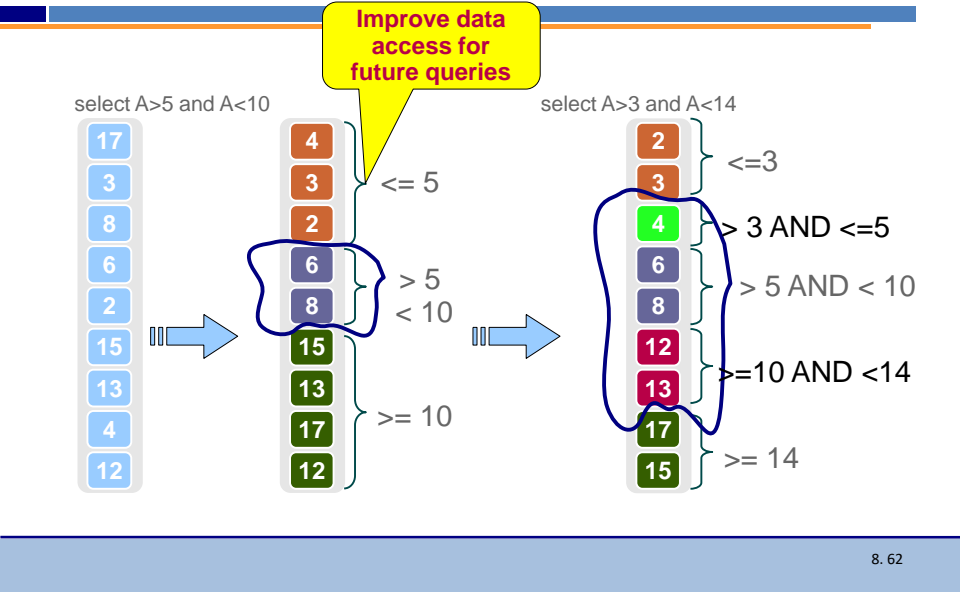


Cracking example



61

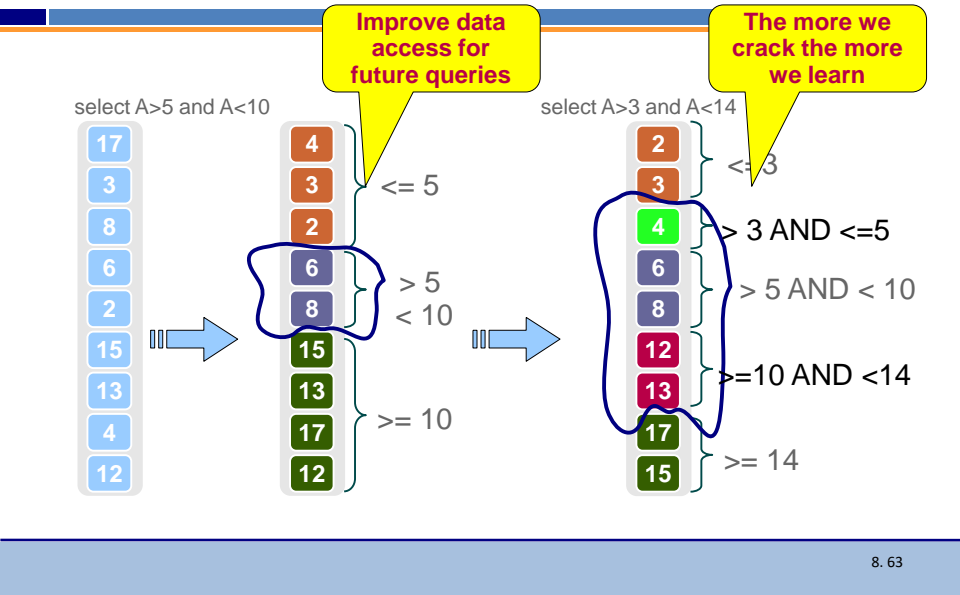
Cracking example



62



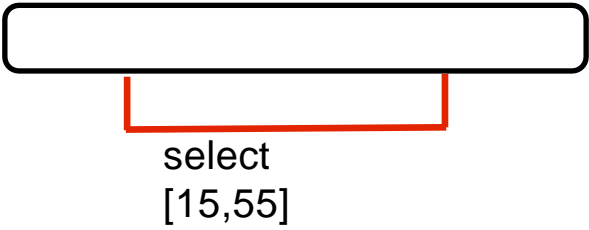
Cracking example



63

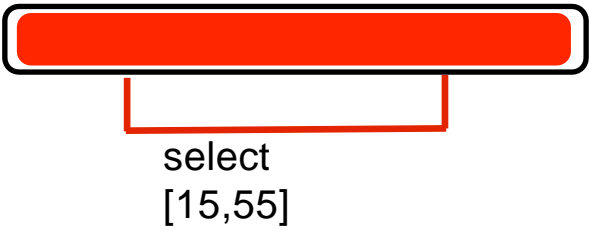


64



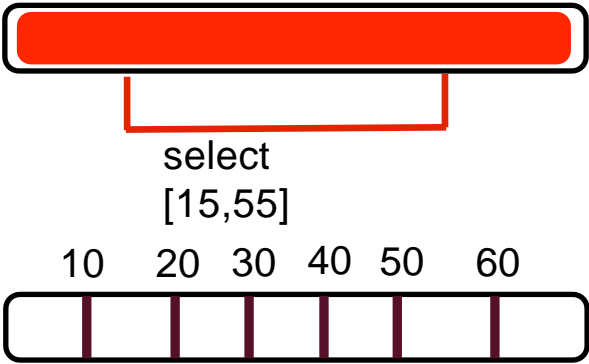
Database Cracking CIDR 2007

65



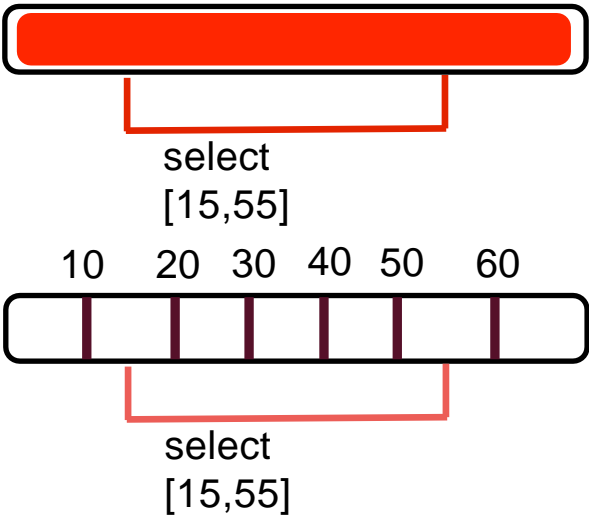
Database Cracking CIDR 2007

66



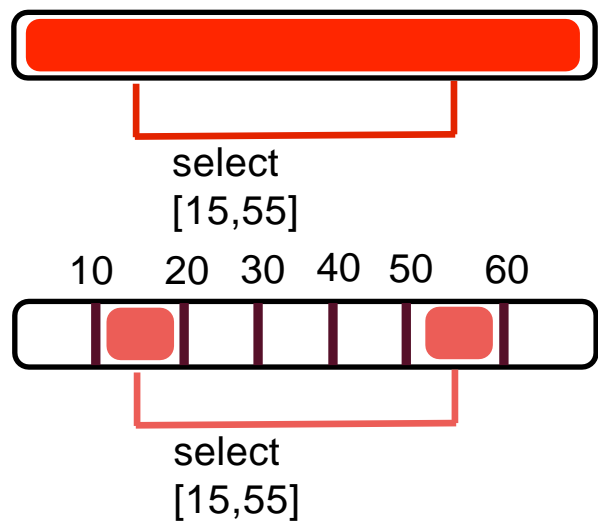
Database Cracking CIDR 2007

67



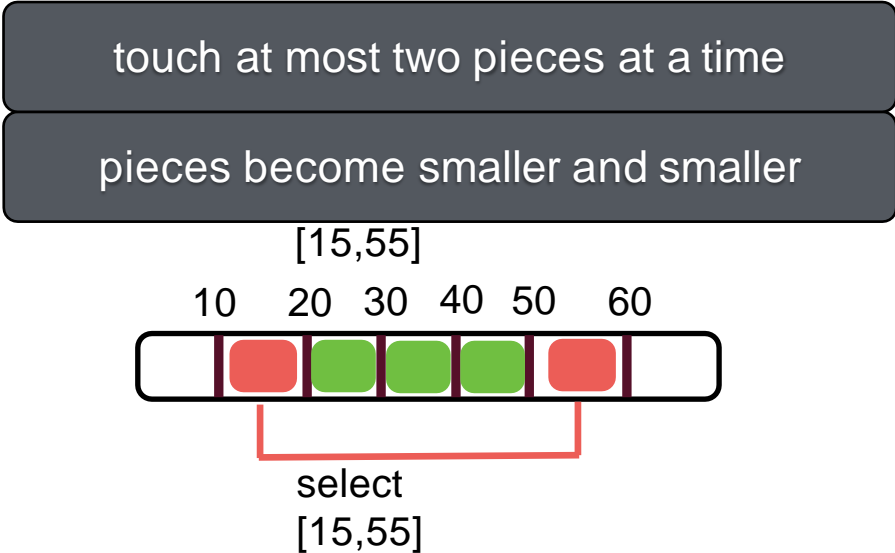
Database Cracking CIDR 2007

68



Database Cracking CIDR 2007

69



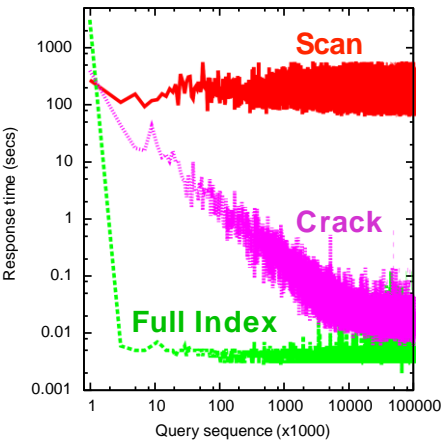
Database Cracking CIDR 2007

70

continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column



Database Cracking CIDR 2007

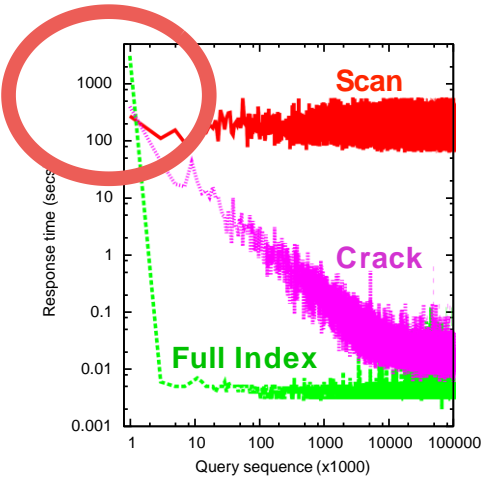
71

continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

almost no
initialization overhead



Database Cracking CIDR 2007

72

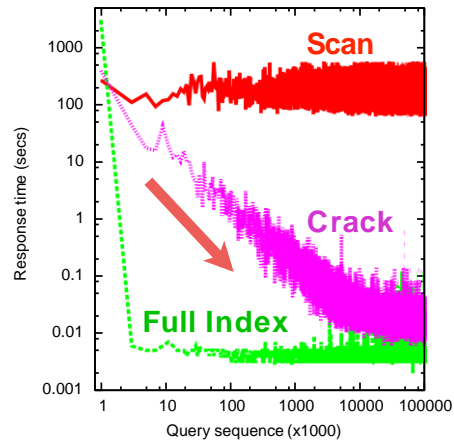
continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

almost no
initialization overhead

continuous improvement



Database Cracking CIDR 2007

73

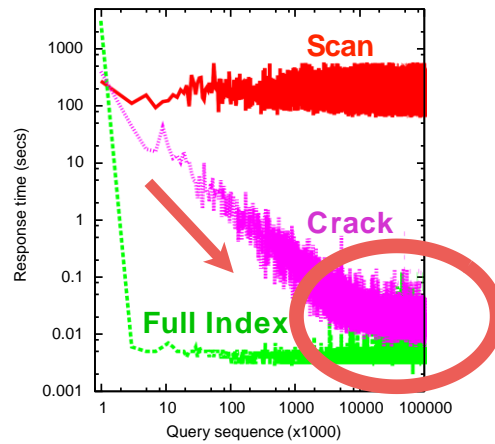
continuous adaptation

set-up

100K random selections
random selectivity
random value ranges
in a 10 million integer column

almost no
initialization overhead

continuous improvement



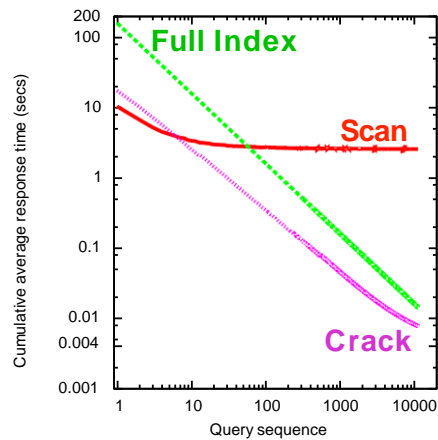
Database Cracking CIDR 2007

74

continuous adaptation

set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column



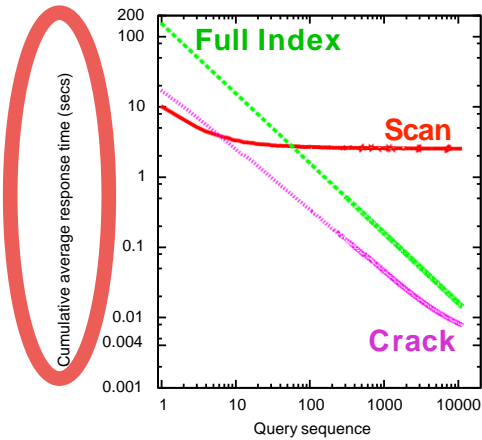
Database Cracking CIDR 2007

75

continuous adaptation

set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column



Database Cracking CIDR 2007

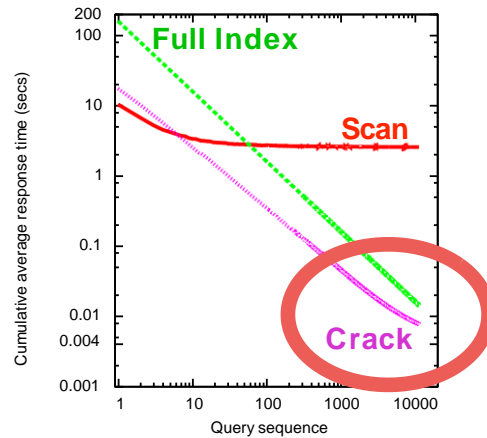
76

continuous adaptation

set-up

10K random selections
selectivity 10%
random value ranges
in a 30 million integer column

10K queries later,
Full Index still has not
amortized the initialization costs



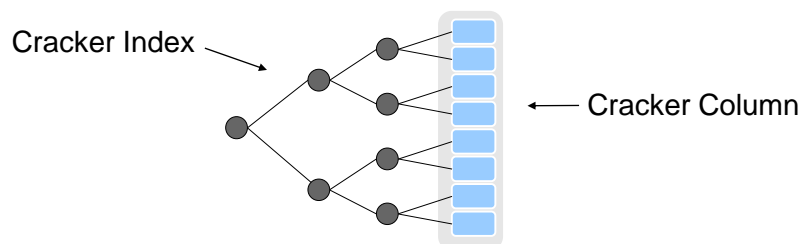
Database Cracking CIDR 2007

77



Cracking – further Design

- 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



8. 78

78

Cracking - Updates

- Update to
 - cracker column
 - cracker index
- Properties:
 - Maintain the self-organization properties
 - Two issues to solve:
 - When
 - How

8. 79

79

When to propagate updates

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

8. 80

80

Update aware Select Operator

- Approach:
 - Select operator applies needed updates before cracking
- Steps of the Select operator:
 - Search the pending insertions column
 - Search the pending deletions column
 - If Steps 1 or 2 find tuples run update algorithm
 - Search the cracker index
 - Physically reorganize the cracker column
 - Update the cracker index
 - Return a slice of the cracker column

8. 81

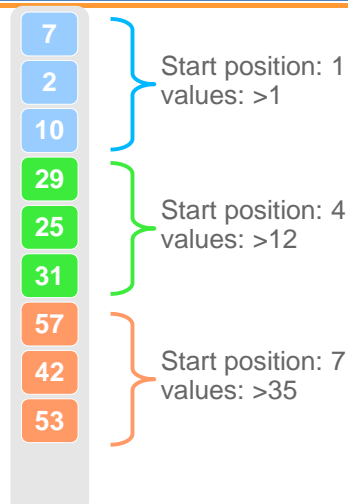
81

Problem: Merging in Update

Insert a new tuple with value 9

9

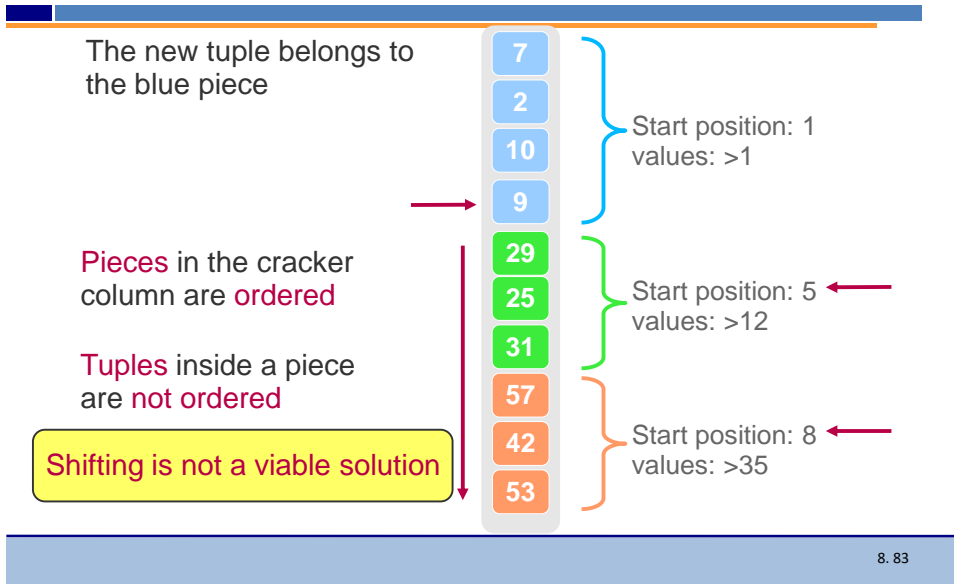
The new tuple belongs to
the blue piece



8. 82

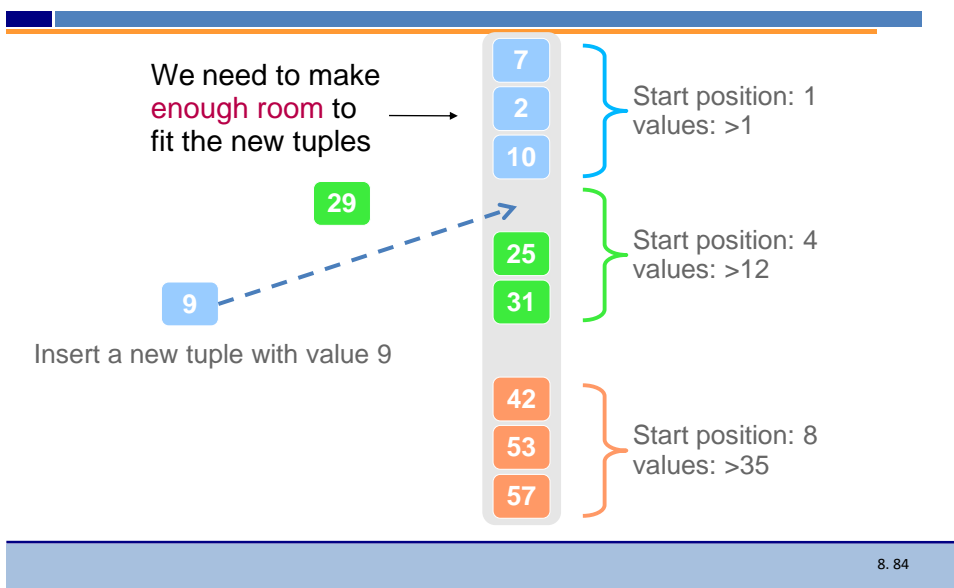
82

First Approach: Merging



83

Second Approach: Merging by Hopping

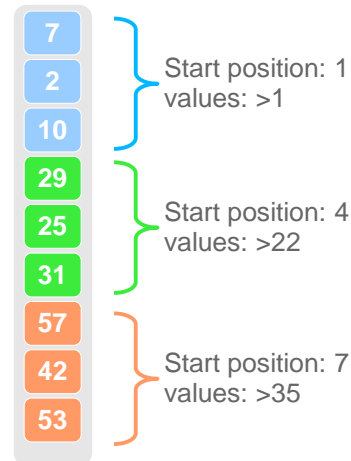


84

Third Approach: Minimalism Principle

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

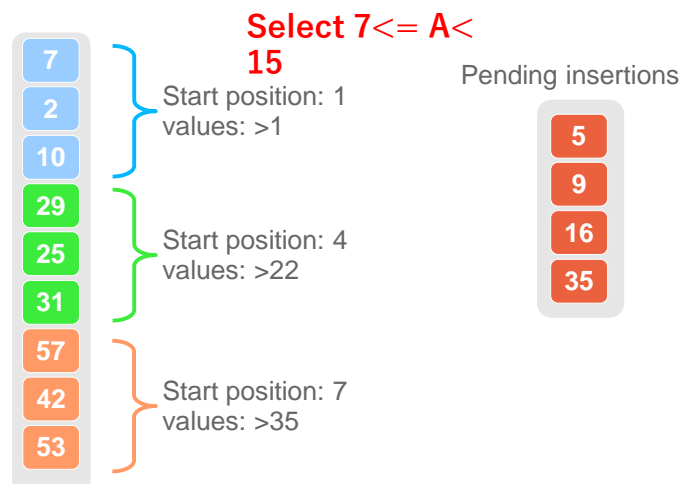
Select $7 \leq A < 15$



8. 85

85

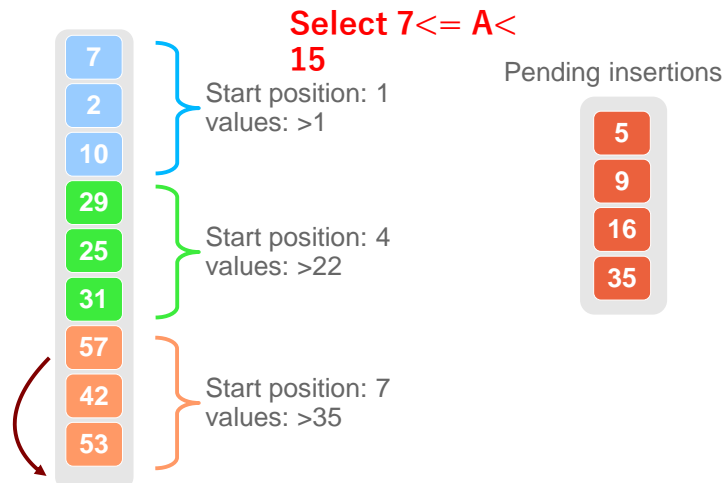
Minimalism Principle: Ripple Approach



8. 86

86

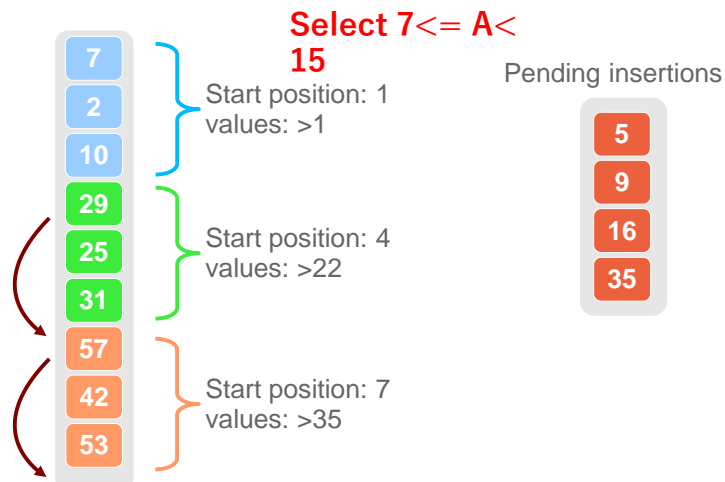
Minimalism Principle: Ripple Approach



8. 87

87

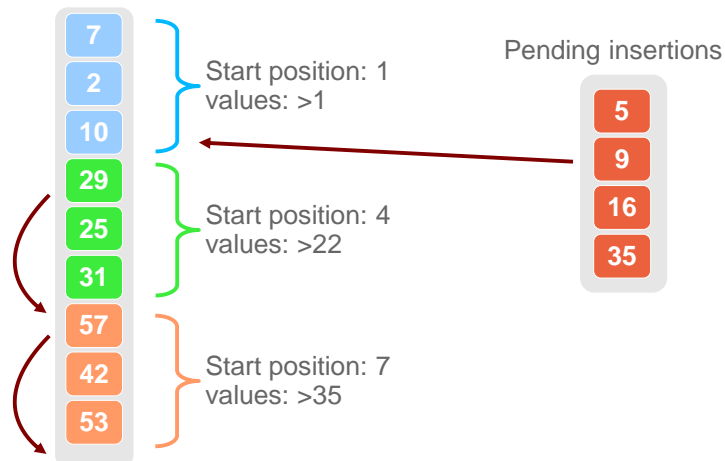
Minimalism Principle: Ripple Approach



8. 88

88

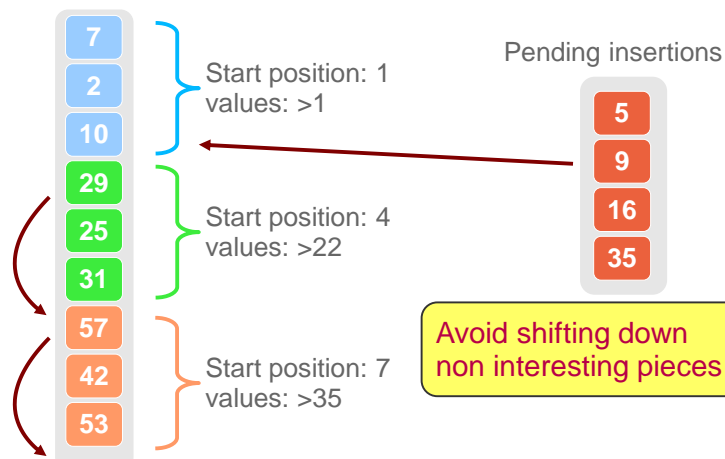
Minimalism Principle: Ripple Approach



8. 89

89

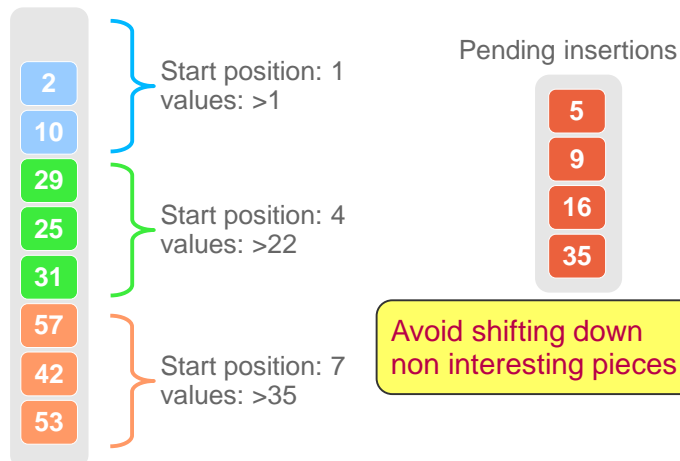
Minimalism Principle: Ripple Approach



8. 90

90

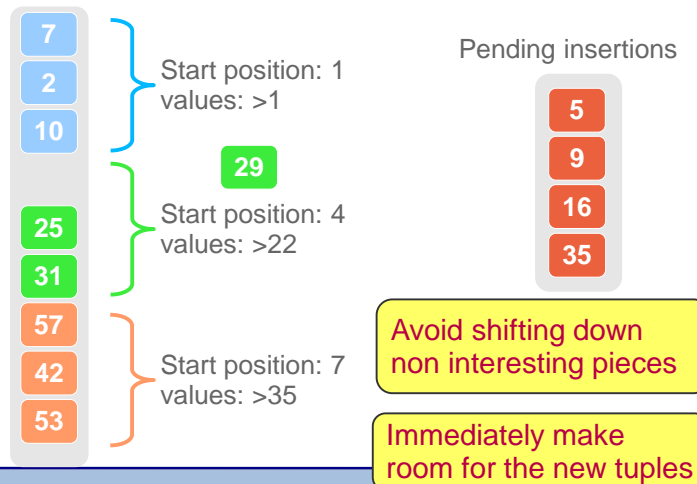
Minimalism Principle: Ripple Approach



8. 91

91

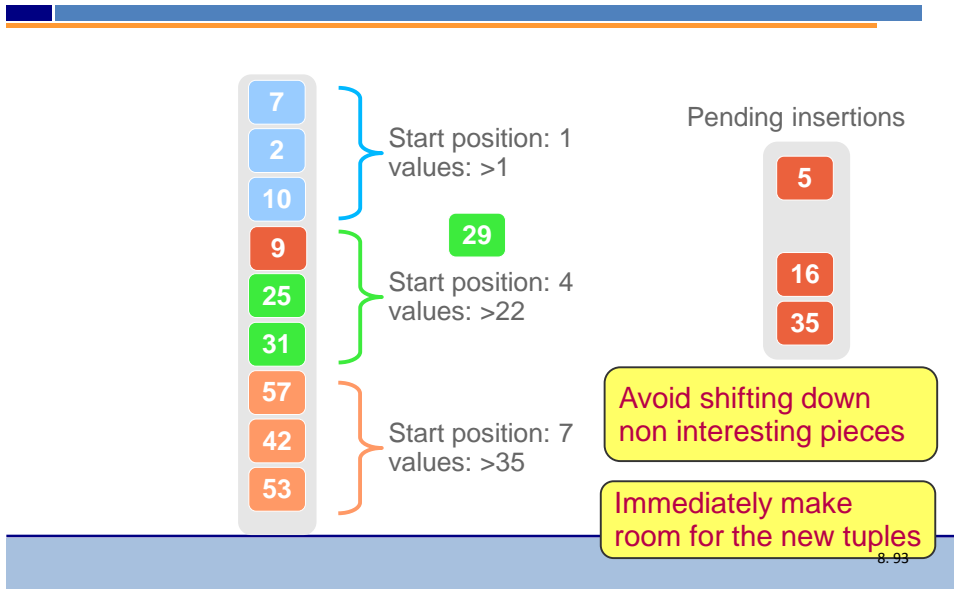
Minimalism Principle: Ripple Approach



8. 92

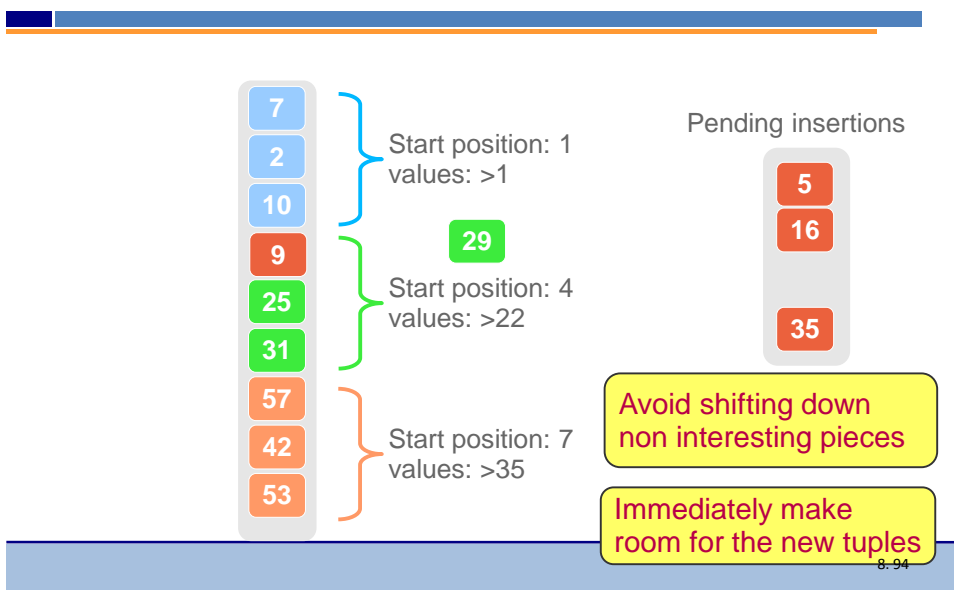
92

The Ripple



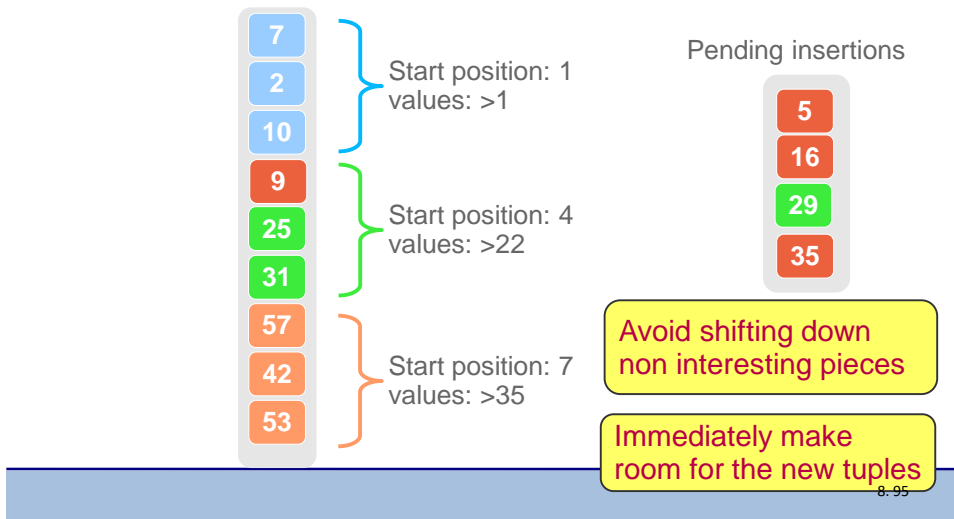
93

The Ripple

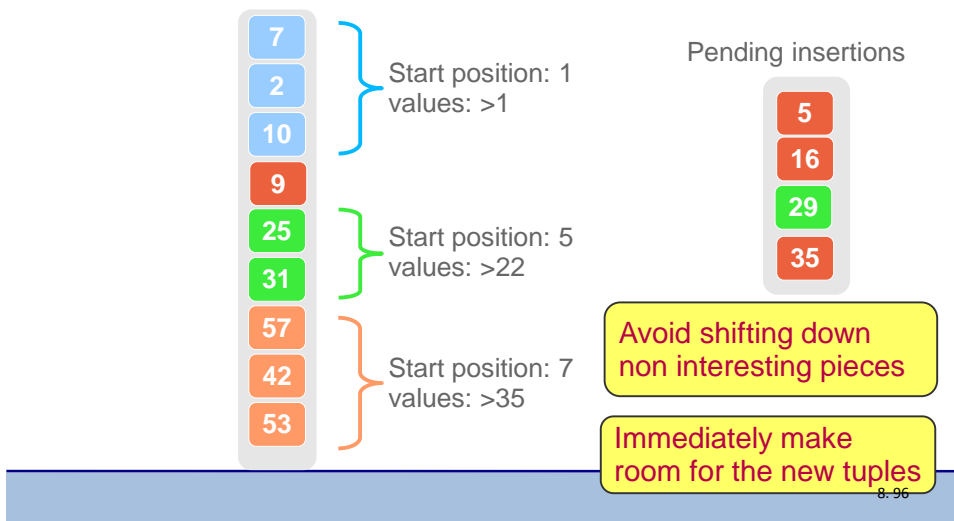


94

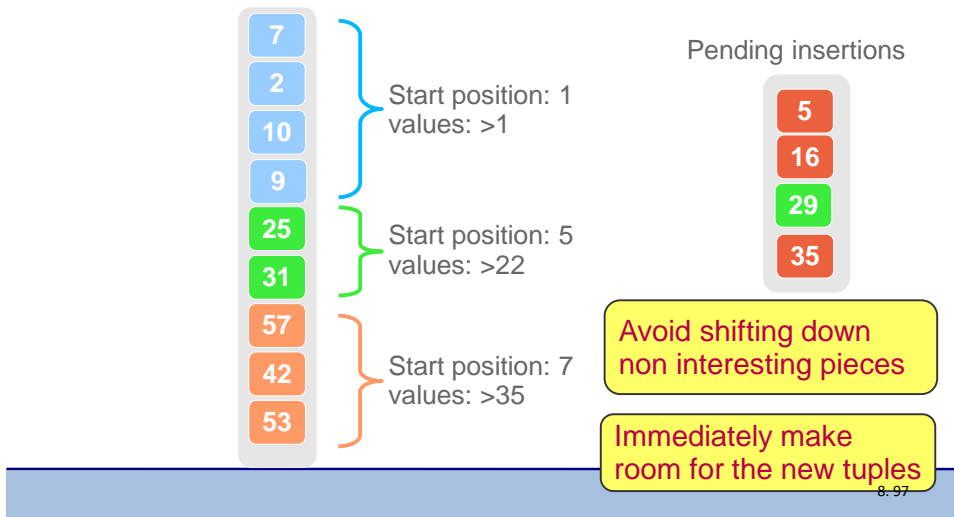
Minimalism Principle: Ripple Approach



Minimalism Principle: Ripple Approach



Minimalism Principle: Ripple Approach



97

Questions???



8. 98

98