# Storage & Indexing in Modern Databases

## ECS 165A – Winter 2022
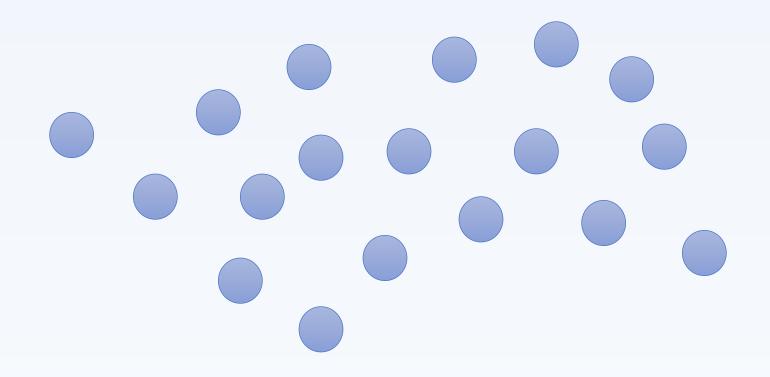
**Mohammad Sadoghi**
Exploratory Systems Lab
Department of Computer Science

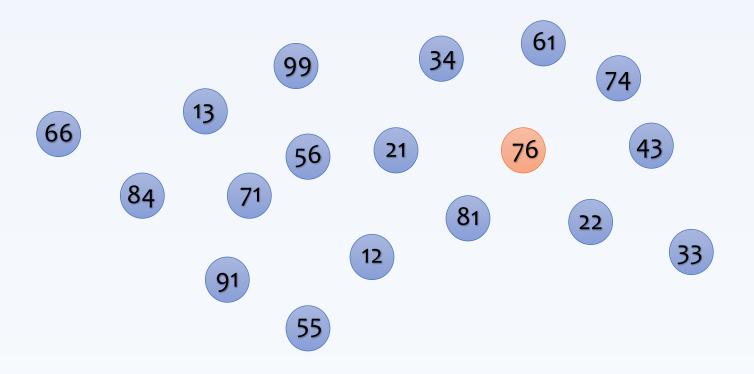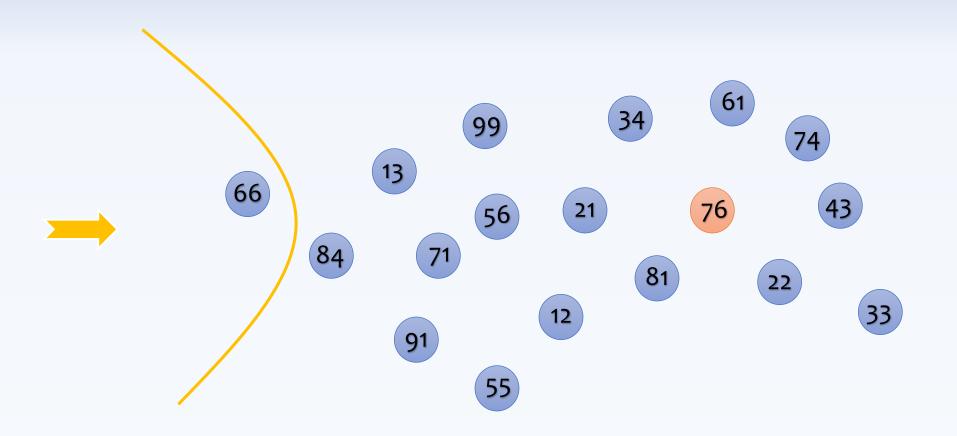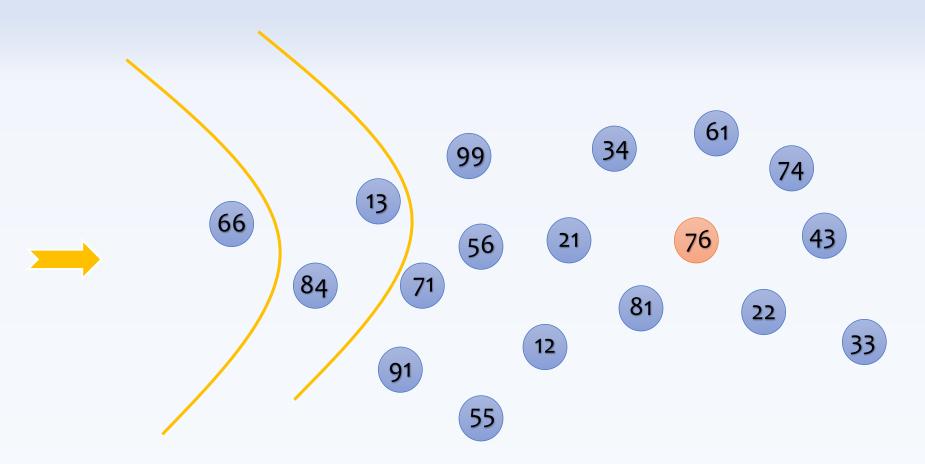**UCDAVIS**
UNIVERSITY OF CALIFORNIA

ExpoLab
Creativity Unfolded

ResilientDB

# How to quickly search for the desired information?

**Searching for 44?**
**(what-if the value does not exist)**
**(could we have an early termination?)**

# Could we impose an order to improve the search?

12 13 21 22 33 34 43 55 56 61 66 71 74 76 81 84 91 99

12  13  21  22  33  34  43  55  56  61  66  71  74  76  81  84  91  99

# Could we impose a structure to further improve the search?

# Could we spread the data cleverly to improve the search?

# hashtable

bucket

Hashing ( ⬤ ) = ?

**(returns a value between 1 to n, where n is the number of buckets)**

Hashing ( 81 ) = 6          81

81

Hashing ( 43 ) = 10    43

Hashing ( 76 ) = 8

81

76

43

Hashing ( 91 ) = 10

81

76

43 91

collisions
(when multiple values
hash to the same bucket)

collisions
(when multiple values
hash to the same bucket)

**(now we can have a constant lookup cost)**

| 56 | 12 |
|----|----|
| 99 |    |
| 71 |    |

| 33 | 61 | 74 |
|----|----|----|

| 55 |
| 81 |
| 84 |

Hashing ( 76 ) = 8

| 76 |
| 13 |

| 43 | 91 |
|----|----|

| 34 |

**Could we instead search for 76, 77, 78, …, 90, 91?**

56 12
99
71
33 61 74
55
81
84
76
13
43 91
34

Hashing ( 76 ) = 8

Hashing ( 77 ) = 1

Hashing ( 78 ) = 3

⋮

Hashing ( 81 ) = 6

⋮

Hashing ( 84 ) = 7

⋮

Hashing ( 90 ) = 8

Hashing ( 91 ) = 10

| 56 | 12 |    |
| 99 |    |    |
| 71 |    |    |
| 33 | 61 | 74 |
| 55 |    |    |
| 81 |    |    |
| 84 |    |    |
| 76 |    |    |
| 13 |    |    |
| 43 | 91 |    |
| 34 |    |    |

**Searching for 76-91**
**Could we instead search for**
**76, 77, 78, ..., 90, 91?**

Hashing ( 76 ) = 8

Hashing ( 77 ) = 1

Hashing ( 78 ) = 3

Hashing ( 81 ) = 6

Hashing ( 84 ) = 7

Hashing ( 90 ) = 8

Hashing ( 91 ) = 10

**Searching for 76-91**

**How about 76.01, 76.02, 76.03, …?**
**(simply not practical)**

**Could we imagine a new design to support searching for a range of values efficiently?**

Let's promote a subset of values as seeds

Let's promote a subset of values as seeds

Suppose every value points to its next larger value

R-Hash- *In-memory latch-free index structure* B Bhattacharjee, M Canim, M. Sadoghi, *US Patent 9,858,303*

sorted seeds

R-Hash- *In-memory latch-free index structure* B Bhattacharjee, M Canim, M. Sadoghi, *US Patent 9,858,303*

sorted seeds

34  71  91

Find the largest seed smaller than 76:  71

Hashing ( 71 ) = **3**

sorted seeds

34   71   91

Find the largest seed smaller than 76:   71

then simply follow the pointers to
find all values between 76-91

56 | 12

99

71

33 | 61 | 74

55

81

84

76

13

Hashing ( 79 ) = **10**

43 | 91 | 79

34

R-Hash- *In-memory latch-free index structure* B Bhattacharjee, M Canim, M. Sadoghi, *US Patent 9,858,303*

Hashing ( 71 ) = **3**

sorted seeds

34  71  91

Find the largest seed smaller than 79:  71

56  12

99

71

33  61  74

55

81

84

76

13

Hashing ( 79 ) = **10**

43  91  79

34

R-Hash- *In-memory latch-free index structure* B Bhattacharjee, M Canim, M. Sadoghi, *US Patent 9,858,303*

Hashing ( 71 ) = **3**

sorted seeds

34   71   91

Find the largest seed smaller than 79:   71

adjust the pointers accordingly

Hashing ( 79 ) = **10**

56   12

99

71

33   61   74

55

81

84

76

13

43   91   79

34

# Database Storage Layouts
## (how likely that we need an index for range queries?)

database pages
(containing a set of records)

[Name: Alice, Age:21, Major: CS]

a database record, e.g.,
[Name: Alice, Age:21, Major: CS]

**Row-based Layout**

database pages
(containing a set of records)

[Name: Alice, Age:21, Major: CS]

a database record, e.g.,
[Name: Alice, Age:21, Major: CS]

[Name]

[Alice]

[Age]

[21]

[Major]

[CS]

**Row-based Layout**

**Column-based Layout**

database pages
(containing a set of records)

[Name]   [Age]   [Major]

[Alice]   [21]   [CS]

a database record, e.g.,
[Name: Alice, Age:21, Major: CS]

[Name: Alice, Age:21, Major: CS]

[Bob]   [21]   [CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]

[Name: Alex, Age:24, Major: EE]

[Name: Sally, Age:25, Major: EE]

[Joe]   [23]   [EE]
[Alex]   [24]   [CS]
[Sally]   [25]   [EE]

**Row-based Layout**

**Column-based Layout**

# Searching for all students between the age of 21 to 24 (may return many students)

**Row-based Layout**

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]

[Name: Alex, Age:24, Major: EE]

[Name: Sally, Age:25, Major: EE]

**Column-based Layout**

| [Name] | [Age] | [Major] |
|--------|-------|---------|
| [Alice] | [21] | [CS] |
| [Bob] | [21] | [CS] |

| [Name] | [Age] | [Major] |
|--------|-------|---------|
| [Joe] | [23] | [EE] |
| [Alex] | [24] | [CS] |
| [Sally] | [25] | [EE] |

# Searching for all students between the age of 21 to 24 (may return many students)

**Index on Age**

[21, 23, 24]

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]

[Name: Alex, Age:24, Major: EE]

[Name: Sally, Age:25, Major: EE]

**Row-based Layout**

[Name]  [Age]  [Major]

[Alice]  [21]  [CS]

[Bob]  [21]  [CS]

[Name]  [Age]  [Major]

[Joe]  [23]  [EE]

[Alex]  [24]  [CS]

[Sally]  [25]  [EE]

**Column-based Layout**

58

Searching for all students between the age of 21 to 24 (may return many students)

Index on Age

[21, 23, 24]

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]

[Name: Alex, Age:24, Major: EE]

[Name: Sally, Age:25, Major: EE]

Row-based Layout

[Name]  [Age]  [Major]

[Alice]  [21]  [CS]

[Bob]  [21]  [CS]

[Joe]  [23]  [EE]
[Alex]  [24]  [CS]
[Sally]  [25]  [EE]

Column-based Layout

59

# Searching for all students between the age of 21 to 24 (may return many students)

**Index on Age**

[21, 23, 24]

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]
[Name: Alex, Age:24, Major: EE]
[Name: Sally, Age:25, Major: EE]

**Row-based Layout**

[Name]    [Age]    [Major]

[Alice]    [21]    [CS]

[Bob]    [21]    [CS]

[Joe]    [23]    [EE]
[Alex]    [24]    [CS]
[Sally]    [25]    [EE]

Alternatively read only the Age column to find the relevant values

**Column-based Layout**

# Searching for all students between the age of 21 to 24
## (may return many students)

[Name]        [Age]        [Major]

[Alice]       [21]         [CS]

[Bob]         [21]         [CS]

**Index on Age**

[21, 23, 24]

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]

[Name: Alex, Age:24, Major: EE]

[Name: Sally, Age:25, Major: EE]

**Is an index really useful here?**

[Joe]         [23]         [EE]
[Alex]        [24]         [CS]
[Sally]       [25]         [EE]

**Row-based Layout**

**Column-based Layout**

Searching for all students over the age of 24
(may return only a few students)

Index on Age

[24+]

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]
[Name: Alex, Age:24, Major: EE]
[Name: Sally, Age:25, Major: EE]

Row-based Layout

[Name]  [Age]  [Major]
[Alice]  [21]  [CS]
[Bob]  [21]  [CS]

[Joe]  [23]  [EE]
[Alex]  [24]  [CS]
[Sally]  [25]  [EE]

Column-based Layout

Searching for all students over the age of 24
(may return only a few students)

Index on Age

[24+]

Could we instead employ
hashing with the seeding idea?

Row-based Layout

Column-based Layout

[Name: Alice, Age:21, Major: CS]

[Name: Bob, Age:21, Major: CS]

[Name: Joe, Age:23, Major: EE]
[Name: Alex, Age:24, Major: EE]
[Name: Sally, Age:25, Major: EE]

[Name]   [Age]   [Major]

[Alice]   [21]    [CS]

[Bob]    [21]    [CS]

[Joe]    [23]    [EE]
[Alex]   [24]    [CS]
[Sally]  [25]    [EE]

# Thank You
# Questions?