# Indexes Part 2 and Query Execution

Instructor: Matei Zaharia

cs245.stanford.edu

# From Last Time: Indexes

Conventional indexes

B-trees

Hash indexes

Multi-key indexing

# Conventional Indexes

**Pros:**

- Simple
- Index is sequential file (good for scans or binary search)

**Cons:**

- Inserts expensive, and/or
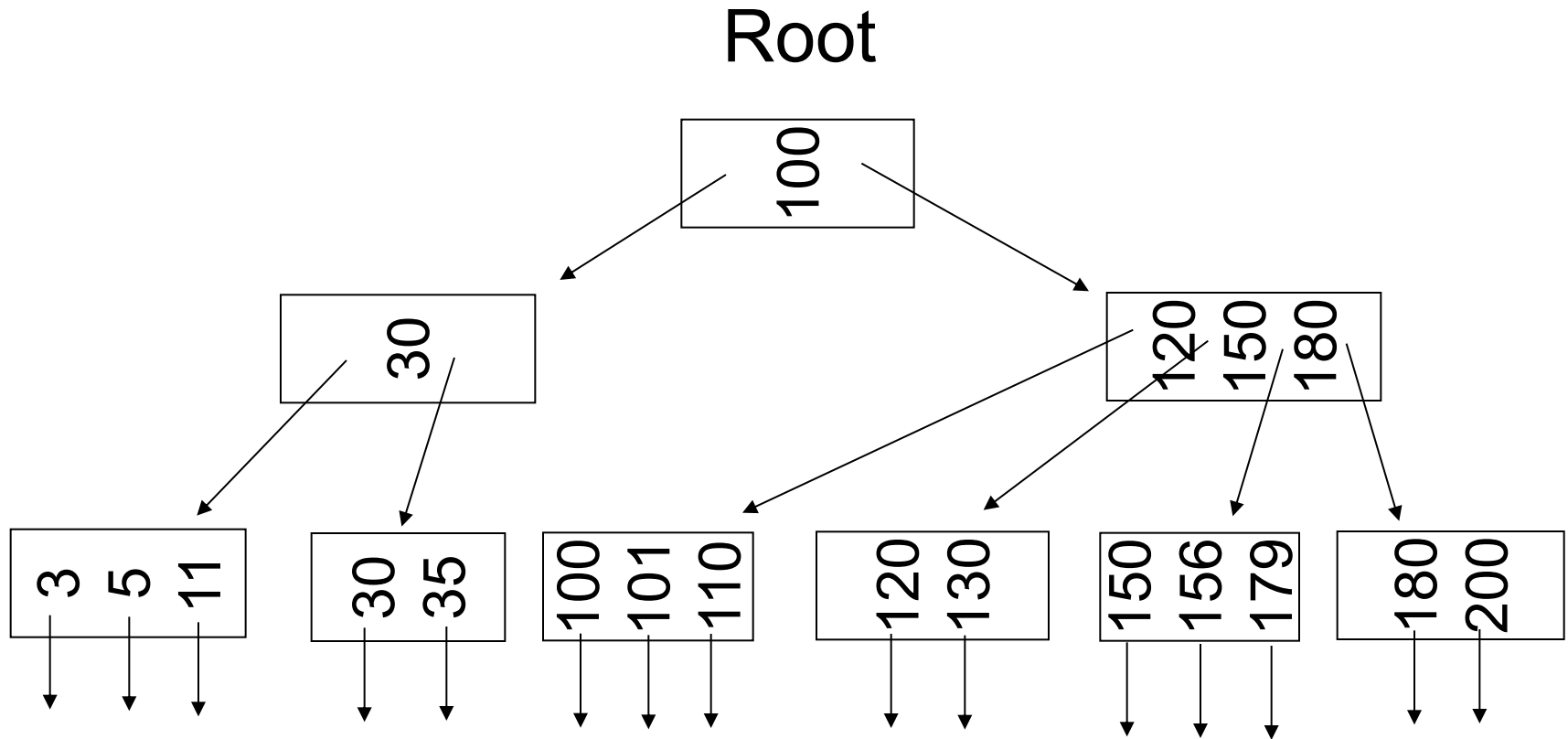- Lose sequentiality & balance

# B-Trees

Another type of index
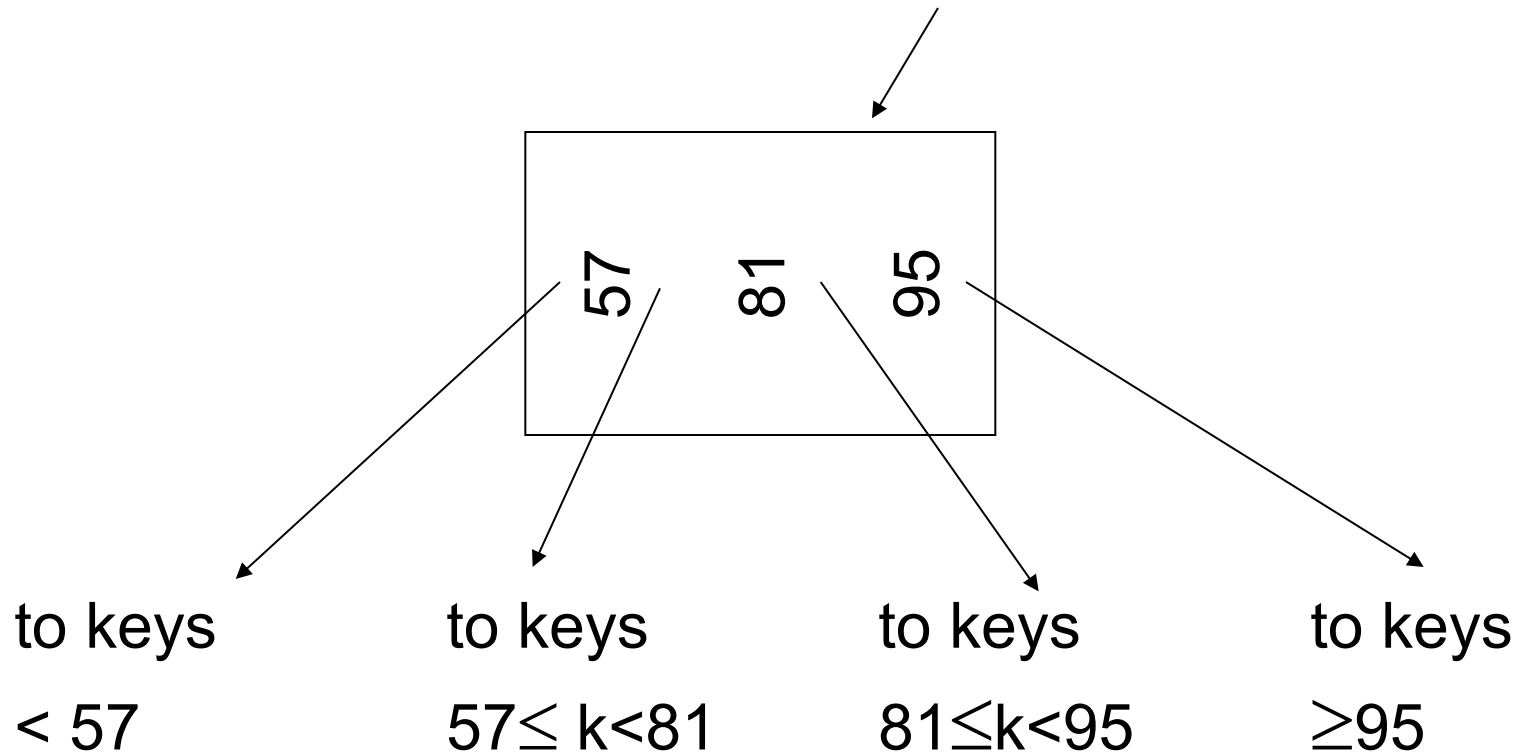   » Give up on sequentiality of index
   » Try to get "balance"


Note: the exact data structure we'll look at is a **B+ tree**, but plain old "B-trees" are similar
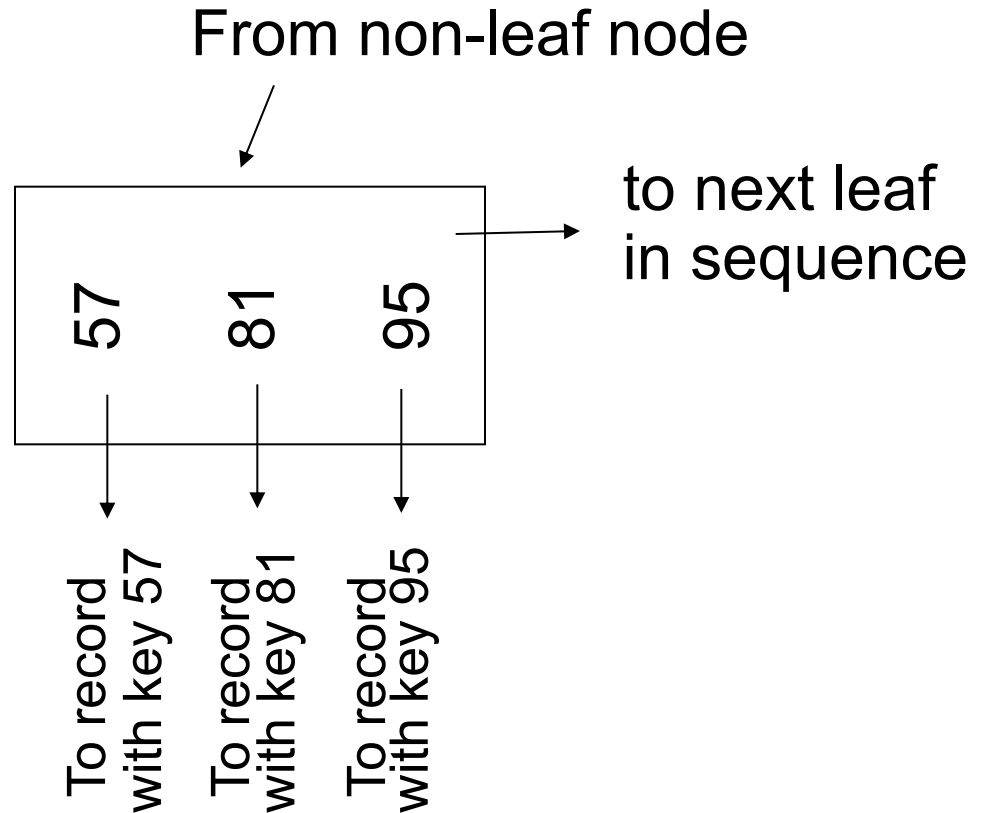
# B+ Tree Example

(n = 3)

Root

# Sample Non-Leaf

57   81   95

to keys
< 57

to keys
57≤ k<81

to keys
81≤k<95

to keys
≥95

# Sample Leaf Node

From non-leaf node

to next leaf
in sequence

57  81  95

To record
with key 57

To record
with key 81

To record
with key 95

# Size of Nodes on Disk

$\left\{\begin{array}{l} \text{n + 1 pointers} \\ \\ \text{n keys} \end{array}\right.$
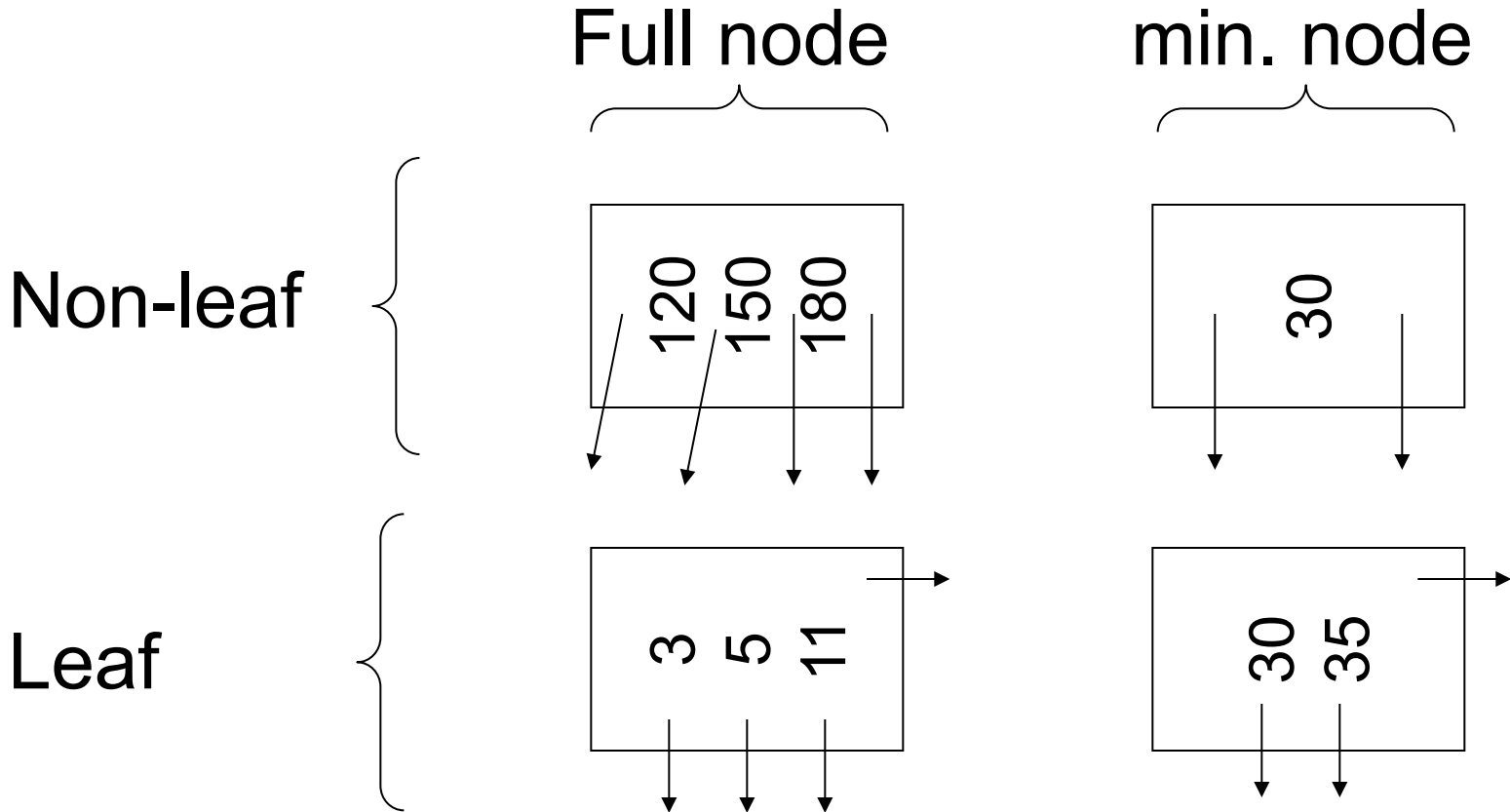
(Fixed size nodes)

# Don't Want Nodes to be Too Empty

Use at least

Non-leaf: $\lceil(n+1)/2\rceil$ pointers

Leaf: $\lfloor(n+1)/2\rfloor$ pointers to data

# Example: n = 3

Full node          min. node

Non-leaf

120 150 180

30

Leaf

3 5 11

30 35

# B+ Tree Rules          (tree of order n)

1. All leaves are at same lowest level (balanced tree)

2. Pointers in leaves point to records, except for "sequence pointer"

# B+ Tree Rules       (tree of order n)

(3) Number of pointers/keys for B+ tree:

|  | Max ptrs | Max keys | Min ptrs→data | Min keys |
|---|---|---|---|---|
| Non-leaf (non-root) | n+1 | n | $\lceil(n+1)/2\rceil$ | $\lceil(n+1)/2\rceil$-1 |
| Leaf (non-root) | n+1 | n | $\lfloor(n+1)/2\rfloor$ | $\lfloor(n+1)/2\rfloor$ |
| Root | n+1 | n | 2* | 1 |

* When there is only one record in the B+ tree, min pointers in the root is 1 (the other pointers are null)

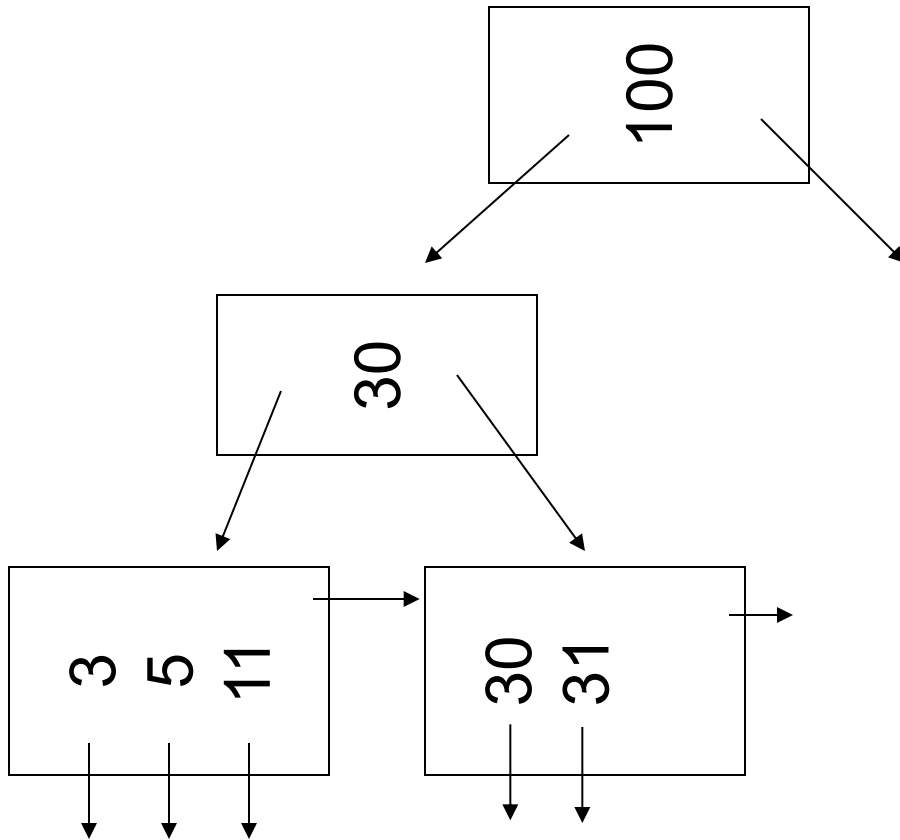# Insertion Into B+ Tree

(a) simple case: have space in leaf

(b) leaf overflow
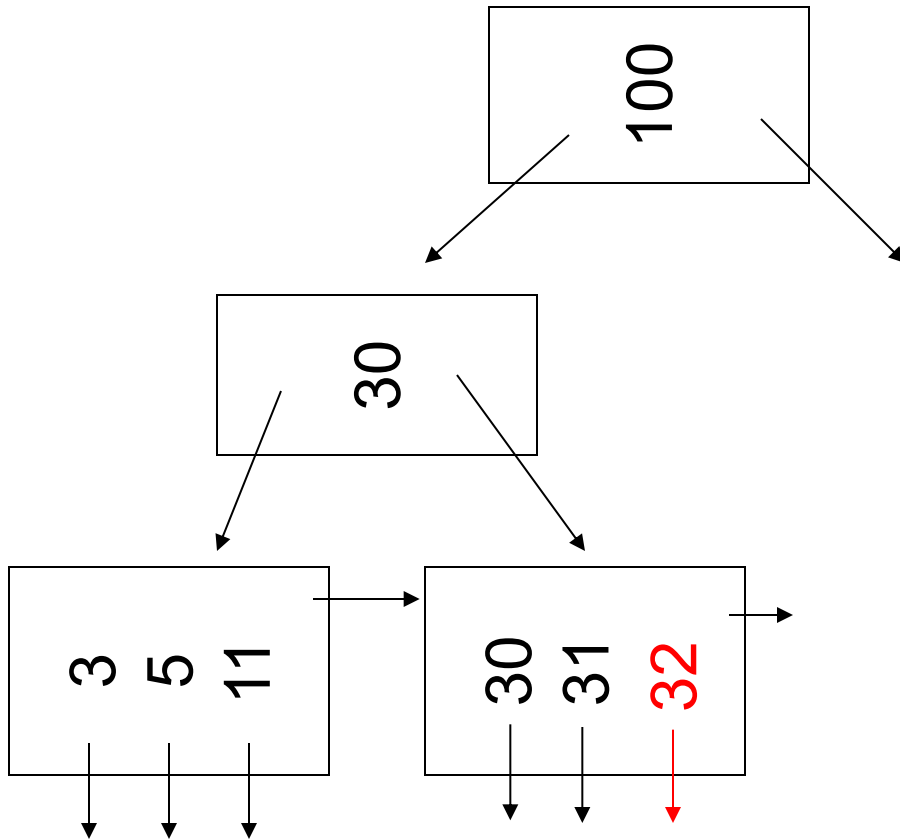
(c) non-leaf overflow

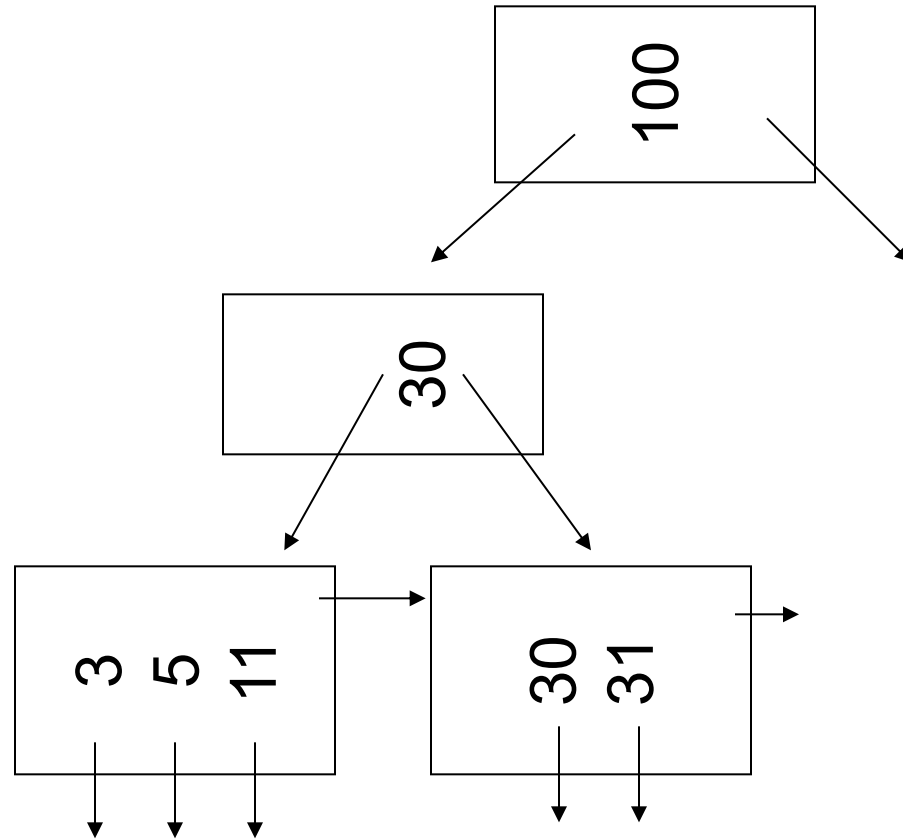(d) new root

# (a) Insert key = 32

n=3



The B-tree diagram shows a root node containing 100, pointing to an internal node containing 30. The internal node points to two leaf nodes: one containing 3, 5, 11 and another containing 30, 31.

# (a) Insert key = 32

n=3

# (b) Insert key = 7

n=3



100

30

3
5
11

30
31

# (b) Insert key = 7

n=3



100

30

3 5

3 5 7 11

30 31

# (a) Insert key = 7

n=3

# (c) Insert key = 160

n=3

100

120 150 180

150 156 179

180 200

# (c) Insert key = 160

n=3

100

120 150 180

150 156 179

160 179

180 200

# (c) Insert key = 160

n=3

100

120 150 180

180

150 156 179

160 179

180 200

# (c) Insert key = 160

n=3

# (d) New root, insert 45

n=3

```
                    ┌─────────────┐
                    │  10  20  30 │
                    └─────────────┘
        ┌──────────────┼──────┼──────────────────┐
        ▼              ▼      ▼                   ▼
   ┌─────────┐   ┌─────────┐  ┌─────────┐   ┌─────────┐
   │ 1  2  3 │→  │ 10  12  │→ │ 20  25  │→  │ 30 32 40│
   └─────────┘   └─────────┘  └─────────┘   └─────────┘
     ↓  ↓  ↓        ↓  ↓         ↓  ↓          ↓  ↓  ↓
```

# (d) New root, insert 45

n=3

10 20 30

1 2 3    10 12    20 25    30 32 40    40 45

# (d) New root, insert 45

n=3

# (d) New root, insert 45

n=3

new root

30

10  20  30

40

1  2  3        10  12        20  25        30  32  40        40  45
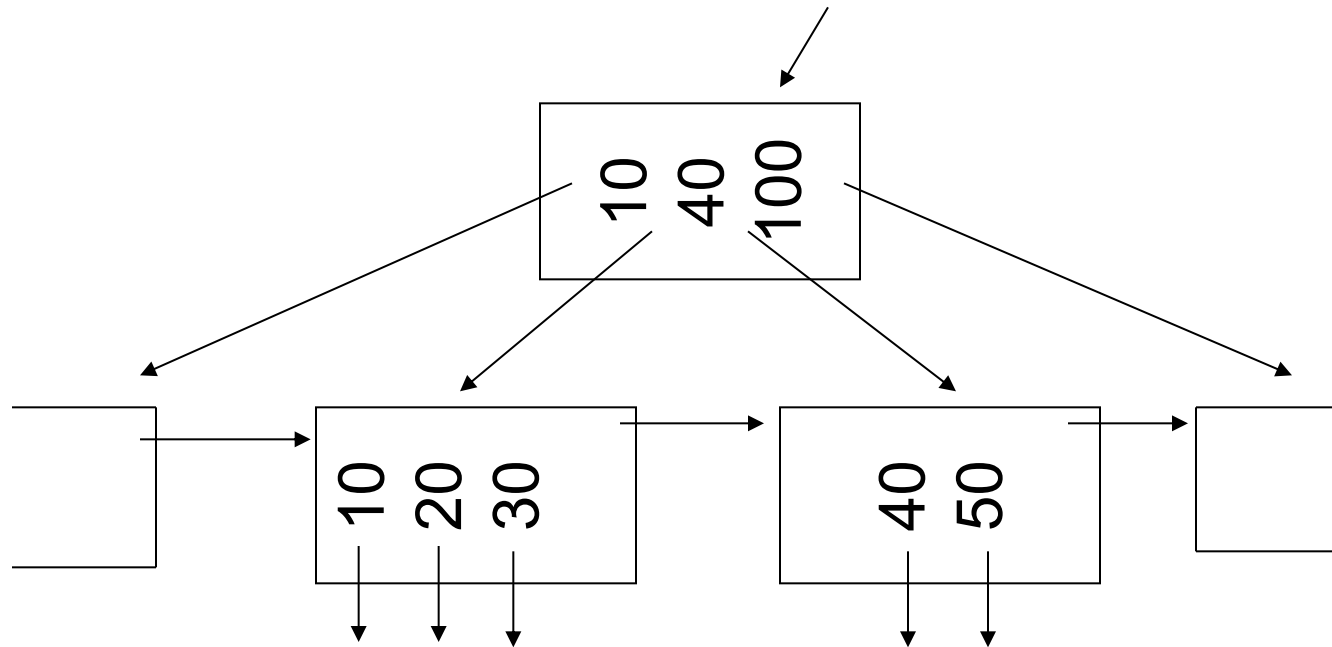
# Deletion from B+ Tree

(a) Simple case: no example

(b) Coalesce with neighbor (sibling)

(c) Re-distribute keys

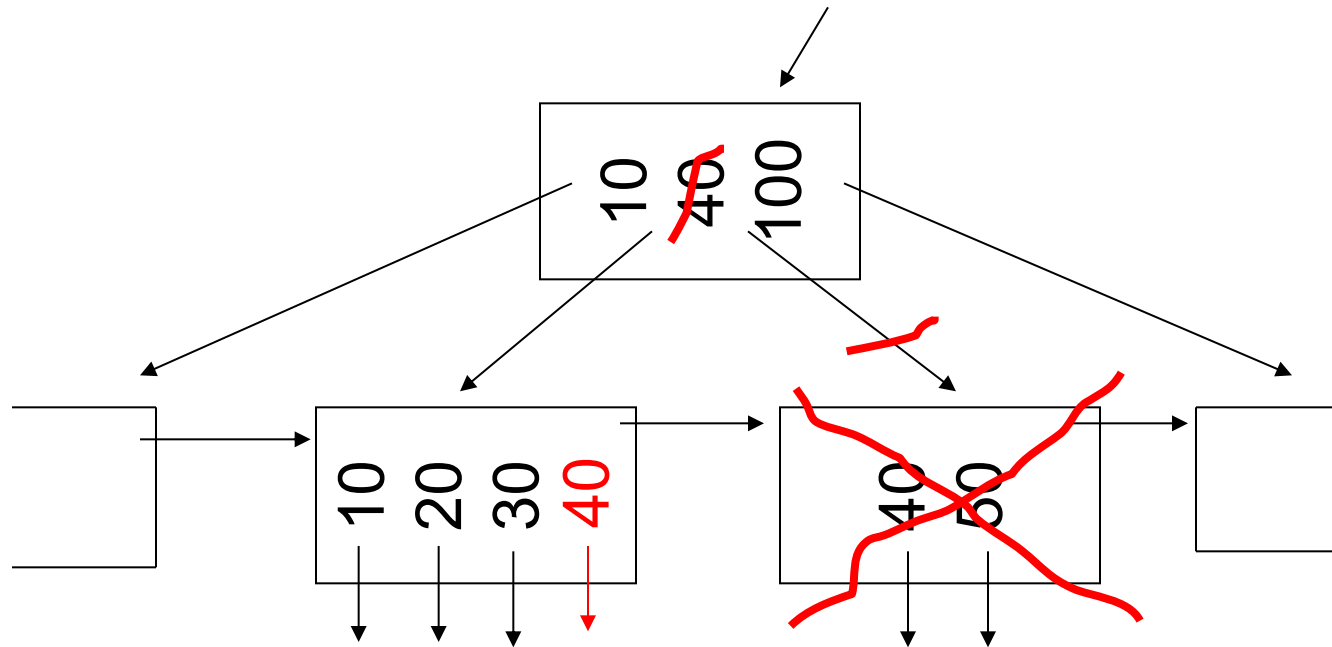(d) Cases (b) or (c) at non-leaf

# (b) Coalesce with sibling
» Delete 50

n=4

# (b) Coalesce with sibling
» Delete 50

n=4

# (c) Redistribute keys
» Delete 50

n=4

# (c) Redistribute keys
  » Delete 50

n=4

# (d) Non-leaf coalesce
– Delete 37

$n=4$

# (d) Non-leaf coalesce
– Delete 37

n=4

# (d) Non-leaf coalesce
## – Delete 37

n=4

# (d) Non-leaf coalesce
  – Delete 37

n=4



new root

# B+ Tree Deletion in Practice

Often, coalescing is not implemented
  » Too hard and not worth it! (Most datasets generally grow in size over time.)

# Interesting Problem:

For B+ tree, how large should n be?

n is number of keys / node

# Sample Assumptions:

(1) Time to read node from disk is
$(S + Tn)$ msec.

# Sample Assumptions:

(1) Time to read node from disk is
     $(S + Tn)$ msec.

(2) Once block in memory, use binary
     search to locate key:
     $(a + b \log_2 n)$ msec.

    For some constants a, b;  Assume a << S

# Sample Assumptions:

(1) Time to read node from disk is
  $(S + Tn)$ msec.

(2) Once block in memory, use binary
  search to locate key:
  $(a + b \log_2 n)$ msec.

  For some constants a, b;   Assume a << S

(3) Assume B+tree is full, i.e., # nodes to
  examine is $\log_n N$ where N = # records

# Can Get:
## f(n) = time to find a record

$$f(n)$$

$$n_{opt}$$    n

# Find $n_{opt}$ by setting f'(n) = 0

Answer is $n_{opt}$ = "a few hundred" in practice

# Exercise

| | |
|---|---|
| S = | 14000 µs |
| T = | 0.2 µs |
| b = | 0.002 µs |
| a = | 0 µs |
| N = | 10,000,000 |

$$f(n) = \log_n N * (S + T\,n + a + b \log_2 n)$$

# N = 10 Million Records

| | |
|---|---|
| S= | 14000 |
| T= | 0.2 |
| b= | 0.002 |
| a= | 0 |
| N= | 10,000,000 |



**find time**

n

times in microseconds

# N = 100 Million Records

| | |
|---|---:|
| S= | 14000 |
| T= | 0.2 |
| b= | 0.002 |
| a= | 0 |
| N= | 100,000,000 |



n

times in microseconds

# N = 10 Million Records

| | |
|---|---:|
| S= | varies |
| T= | 0.2 |
| b= | 0.002 |
| a= | 0 |
| N= | 10,000,000 |



## Effect of S

times in microseconds

# Some Types of Indexes

Conventional indexes

B-trees

Hash indexes

Multi-key indexing

# Hash Indexes

key $\longrightarrow$ h(key)

record / ptr - - → overflow bucket

Buckets
(block sized)

Chaining is used to handle bucket overflow

# Hash vs Tree Indexes

+ O(1) instead of O(log N) disk accesses

– Can't efficiently do range queries

# Challenge: Resizing

Hash tables try to keep occupancy in a fixed range (50-80%) and slow down beyond that
  » Too much chaining

How to resize the table when this happens?
  » **In memory:** just move everything, amortized cost is pretty low
  » **On disk:** moving everything is expensive!

# Extendible Hashing

Tree-like design for hash tables that allows cheap resizing while requiring 2 IOs / access

# Extendible Hashing: 2 Ideas

(a) Use i of b bits output by hash function

$$\longleftarrow \quad b \quad \longrightarrow$$

h(K) → 00110101

i

i will grow over time; the first i bits of each key's hash are used to map it to a bucket

# **Extendible Hashing: 2 Ideas**

(b) Use a directory with pointers to buckets

h(K)[0..i]          to bucket

# Example: 4-bit h(K), 2 keys/bucket

global depth

local depth

i = 1

| 1 |
|---|
| 0001 |
| 0010 |

| 1 |
|---|
| 1001 |
| 1100 |

Insert 0010

# Example: 4-bit h(K), 2 keys/bucket

local depth

global depth

i = 1

| 1 |
|---|
| 0001 |
| |

| 1 |
|---|
| 1001 |
| 1100 |

Insert 1010

# Example: 4-bit h(K), 2 keys/bucket

i = 1

1
| 0001 |
| |

1
| 1001 |
| 1100 |

1010 1100

Insert 1010

1
| 1100 |
| |

# Example: 4-bit h(K), 2 keys/bucket

1

0001

i = 1

1 2

1001

1010 1100

1 2

1100

Insert 1010

i = 2

00

01

10

11

New directory

# Example

i = 2

00

01

10

11

1

0001

2

1001

1010

2

1100

Insert:

0111

0000

# Example

i = 2

00

01

10

11

Insert:

0111

0000

0000

0001

1

~~0001~~ 0111

~~0111~~

2

1001

1010

2

1100

# Example

i = 2

00

01

10

11

2
0000
0001

1 2
0001   0111
0111

2
1001
1010

2
1100

Insert:

0111

0000

# Example

i = 2

00
01
10
11

0000    2
0001

0111    2

1001    2
1010

1100    2

Note: still need chaining if values of h(K) repeat and fill a bucket

# Some Types of Indexes

Conventional indexes

B-trees

Hash indexes

Multi-key indexing

# Motivation

Example: find records where

DEPT = "Toy" AND SALARY > 50k

# Strategy I:

Use one index, say Dept.

Get all Dept = "Toy" records
      and check their salary

# Strategy II:

Use 2 indexes; intersect lists of pointers

Toy → ☐☐☐☐☐       ☐☐☐☐☐☐☐☐ ← Sal
                                    > 50k

# Strategy III:

Multi-key index

One design:

# k-d Trees

Split dimensions in any order to hold k-dimensional data

# k-d Trees



10    20

# k-d Trees

# k-d Trees

# k-d Trees



Efficient range queries in both dimensions

# Storage System Examples

**MySQL:** transactional DBMS
  » Row-oriented storage with 16 KB pages
  » Variable length records with headers, overflow
  » Index types:
    • B-tree
    • Hash (in memory only)
    • R-tree (spatial data)
    • Inverted lists for full text search
  » Can compress pages with Lempel-Ziv

# Storage System Examples

**Apache Parquet + Hive:** analytical data lake

» Column-oriented storage as set of ~1 GB files (each file has a slice of all columns)

» Various compression and encoding schemes at the level of pages in a file

- Special scheme for nested fields (Dremel)

» Header with statistics at the start of each file

- Min/max of columns, nulls, Bloom filter

» Files partitioned into directories by one key

# Query Execution

Overview

Relational operators

Execution methods

# Query Execution Overview

Recall that one of our key principles in data intensive systems was **declarative APIs**

» Specify what you want to compute, not how

We saw how these can translate into many storage strategies

How to execute queries in a declarative API?

# Query Execution Overview

```
┌─────────────────────────┐
│   Query representation   │
│       (e.g. SQL)         │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│    Logical query plan    │
│ (e.g. relational algebra)│
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│  Optimized logical plan  │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│      Physical plan       │
│  (code/operators to run) │
└─────────────────────────┘
```

Many execution methods: per-record exec, vectorization, compilation

# Plan Optimization Methods

**Rule-based:** systematically replace some expressions with other expressions
  » Replace X OR TRUE with TRUE
  » Replace M*A + M*B with M*(A+B) for matrices

**Cost-based:** propose several execution plans and pick best based on a **cost model**

**Adaptive:** update execution plan at runtime

# Execution Methods

**Interpretation:** walk through query plan operators for each record

**Vectorization:** walk through in batches

**Compilation:** generate code (like System R)

# Typical RDBMS Execution

SQL query

( parse )

parse tree

( convert )

logical query plan

( apply rules )

"improved" l.q.p

( estimate result sizes )

l.q.p. +sizes

( consider physical plans )

$\{P_1, P_2, \dots\}$

result

( execute )

statistics

$P_i$

( pick best )

$\{(P_1,C_1), (P_2,C_2), \dots\}$

( estimate costs )

CS 245

81

# Query Execution

Overview

Relational operators

Execution methods

# The Relational Algebra

Collection of operators over tables (relations)
» Each table has named attributes (fields)

Codd's original RA: tables are **sets of tuples** (unordered and tuples cannot repeat)

SQL's RA: tables are **bags (multisets) of tuples**; unordered but each tuple may repeat

# Relational Algebra Operators

Basic set operators:

**Intersection:** R ∩ S

**Union:** R ∪ S          for tables with same schema

**Difference:** R – S

**Cartesian Product:** R × S    { (r, s) | r ∈ R, s ∈ S }

# Relational Algebra Operators

Basic set operators:

**Intersection:** $R \cap S$

**Union:** $R \cup S$ $\longleftarrow$ consider both distinct (set union) and non-distinct (bag union)

**Difference:** $R - S$

**Cartesian Product:** $R \times S$

# Relational Algebra Operators

Special query processing operators:

**Selection:** $\sigma_{condition}(R)$    { r ∈ R | condition(r) is true }

**Projection:** $\Pi_{expressions}(R)$   { expressions(r) | r ∈ R }

**Natural Join:** R ⋈ S   { (r, s) ∈ R × S) | r.key = s.key }
                           where key is the common fields

# Relational Algebra Operators

Special query processing operators:

**Aggregation:** $_{keys}G_{agg(attr)}(R)$     SELECT agg(attr)
                                                FROM R
                                                GROUP BY keys

Examples:   $_{department}G_{Max(salary)}(Employees)$

            $G_{Max(salary)}(Employees)$

# Algebraic Properties

Many properties about which combinations of operators are equivalent

» That's why it's called an algebra!

# Properties: Unions, Products and Joins

R ∪ S = S ∪ R

R ∪ (S ∪ T) = (R ∪ S) ∪ T

Tuple order in a relation doesn't matter (unordered)

R × S = S × R

(R × S) × T = R × (S × T)

Attribute order in a relation doesn't matter either

R ⋈ S = S ⋈ R

(R ⋈ S) ⋈ T = R ⋈ (S ⋈ T)

# Properties: Selects

$\sigma_{p \land q}(R) =$

$\sigma_{p \lor q}(R) =$

# Properties: Selects

$$\sigma_{p \wedge q}(R) = \sigma_p(\sigma_q(R))$$

$$\sigma_{p \vee q}(R) = \sigma_p(R) \cup \sigma_q(R)$$

careful with repeated elements

# Bags vs. Sets

R = {a,a,b,b,b,c}

S = {b,b,c,c,d}

R ∪ S = ?

# Bags vs. Sets

R = {a,a,b,b,b,c}

S = {b,b,c,c,d}

R ∪ S = ?

- **Option 1:** SUM of counts
  R ∪ S = {a,a,b,b,b,b,b,c,c,c,d}
- **Option 2:** MAX of counts
  R ∪ S = {a,a,b,b,b,c,c,d}

# Executive Decision

Use "SUM" option for bag unions

Some rules that work for set unions cannot be used for bags

# Properties: Project

Let: X = set of attributes
     Y = set of attributes

$\Pi_{X \cup Y}$ (R) =

# Properties: Project

Let:  X = set of attributes

  Y = set of attributes

$\Pi_{X \cup Y} (R) = \Pi_X(\Pi_Y(R))$

# Properties: Project

Let:  X = set of attributes

   Y = set of attributes

$$\Pi_{X \cup Y}(R) = \Pi_X(\Pi_Y(R))$$

# Properties: σ + ⋈

Let p = predicate with only R attribs

    q = predicate with only S attribs

    m = predicate with only R, S attribs

$\sigma_p(R \bowtie S) =$

$\sigma_q(R \bowtie S) =$

# Properties: σ + ⋈

Let p = predicate with only R attribs

    q = predicate with only S attribs

    m = predicate with only R, S attribs

$\sigma_p(R \bowtie S) = \sigma_p(R) \bowtie S$

$\sigma_q(R \bowtie S) = R \bowtie \sigma_q(S)$

# **Properties: σ + ⋈**

Some rules can be derived:

$\sigma_{p \wedge q}(R \bowtie S) =$

$\sigma_{p \wedge q \wedge m}(R \bowtie S) =$

$\sigma_{p \vee q}(R \bowtie S) =$

# Properties: σ + ⋈

Some rules can be derived:

$$\sigma_{p \wedge q}(R \bowtie S) = \sigma_p(R) \bowtie \sigma_q(S)$$

$$\sigma_{p \wedge q \wedge m}(R \bowtie S) = \sigma_m(\sigma_p(R) \bowtie \sigma_q(S))$$

$$\sigma_{p \vee q}(R \bowtie S) = (\sigma_p(R) \bowtie S) \cup (R \bowtie \sigma_q(S))$$

# Prove One, Others for Practice

$$\sigma_{p \wedge q}(R \bowtie S) = \sigma_p (\sigma_q(R \bowtie S))$$

$$= \sigma_p (R \bowtie \sigma_q(S))$$

$$= \sigma_p (R) \bowtie \sigma_q(S)$$

# Properties: $\Pi$ + σ

Let x = subset of R attributes

z = attributes in predicate p
(subset of R attributes)

$\Pi_x(\sigma_p (R)) =$

# Properties: $\Pi$ + σ

Let x = subset of R attributes

z = attributes in predicate p
(subset of R attributes)

$$\Pi_x(\sigma_p (R)) = \sigma_p(\Pi_x(R))$$

# Properties: $\Pi$ + σ

Let x = subset of R attributes

z = attributes in predicate p
(subset of R attributes)

$$\Pi_x(\sigma_p (R)) = \Pi_x(\sigma_p(\Pi_{x\cup z}(R)))$$

# Properties: $\Pi$ + $\bowtie$

Let    x = subset of R attributes

        y = subset of S attributes

        z = intersection of R,S attributes

$\Pi_{x \cup y}(R \bowtie S) = \Pi_{x \cup y} ((\Pi_{x \cup z} (R)) \bowtie (\Pi_{y \cup z} (S)))$

# Typical RDBMS Execution

SQL query

parse

parse tree

convert

logical query plan

result

apply rules

execute

statistics

"improved" l.q.p

$P_i$

estimate result sizes

pick best

l.q.p. +sizes

$\{(P_1,C_1), (P_2,C_2), ...\}$

consider physical plans

estimate costs

$\{P_1, P_2, ...\}$

# Example SQL Query

```
SELECT title
FROM StarsIn
WHERE starName IN (
        SELECT name
        FROM MovieStar
        WHERE birthdate LIKE '%1960'
);
```

(Find the movies with stars born in 1960)

# Parse Tree

<Query>
|
<SFW>

SELECT  <SelList>  FROM  <FromList>  WHERE  <Condition>

<Attribute>  <RelName>  <Tuple>  IN  <Query>

title  StarsIn  <Attribute>  ( <Query> )

starName  <SFW>

SELECT  <SelList>  FROM  <FromList>  WHERE  <Condition>

<Attribute>  <RelName>  <Attribute>  LIKE  <Pattern>

name  MovieStar  birthDate  '%1960'

# Logical Query Plan

$$\Pi_{title}$$

$$\sigma_{starName=name}$$

$$\times$$

StarsIn            $$\Pi_{name}$$

$$\sigma_{birthdate\ LIKE\ '\%1960'}$$

MovieStar

# Improved Logical Query Plan

$$\Pi_{\text{title}}$$

$$\bowtie$$ starName=name

StarsIn          $\Pi_{\text{name}}$

$\sigma_{\text{birthdate LIKE '%1960'}}$

MovieStar

Question:
Push $\Pi_{\text{title}}$
to StarsIn?

# Estimate Result Sizes

⋈          Need expected size

StarsIn

Π
┆
σ
┆
MovieStar

# One Physical Plan



Hash join → Parameters: join order, memory size, project attributes, ...

Seq scan

Index scan → Parameters: select condition, ...

H

StarsIn

MovieStar

# Another Physical Plan

Hash join → Parameters: join order, memory size, project attributes, ...

H

Index scan

Seq scan → Parameters: select condition, ...

StarsIn

MovieStar

# Another Physical Plan



Which plan is likely to be better?

# Estimating Plan Costs

Logical plan

$P_1$      $P_2$      ...      $P_n$      Physical plan candidates

$C_1$      $C_2$      ...      $C_n$

Pick best!

Covered in next few lectures!

# Query Execution

Overview

Relational operators

Execution methods

# Now That We Have a Plan, How Do We Run it?

Several different options that trade between complexity, setup time & performance

# Example: Simple Query

```
SELECT quantity * price
  FROM orders
 WHERE productId = 75
```

$$\Pi_{\text{quanity*price}} \left( \sigma_{\text{productId=75}} \left( \text{orders} \right) \right)$$

# Method 1: Interpretation

```
interface Operator {
  Tuple next();
}

class TableScan: Operator {
  String tableName;
}

class Select: Operator {
  Operator parent;
  Expression condition;
}

class Project: Operator {
  Operator parent;
  Expression[] exprs;
}
```

```
interface Expression {
  Value compute(Tuple in);
}

class Attribute: Expression {
  String name;
}

class Times: Expression {
  Expression left, right;
}

class Equals: Expression {
  Expression left, right;
}
```

# Example Expression Classes

```
class Attribute: Expression {
  String name;

  Value compute(Tuple in) {
    return in.getField(name);
  }
}


class Times: Expression {
  Expression left, right;

  Value compute(Tuple in) {
    return left.compute(in) * right.compute(in);
  }
}
```

probably better to use a
numeric field ID instead

# Example Operator Classes

```
class TableScan: Operator {
  String tableName;

  Tuple next() {
    // read & return next record from file
  }
}

class Project: Operator {
  Operator parent;
  Expression[] exprs;

  Tuple next() {
    tuple = parent.next();
    fields = [expr.compute(tuple) for expr in exprs];
    return new Tuple(fields);
  }
}
```

# Running Our Query with Interpretation

```
ops = Project(
        expr = Times(Attr("quantity"), Attr("price")),
        parent = Select(
          expr = Equals(Attr("productId"), Literal(75)),
          parent = TableScan("orders")
        )
      );

while(true) {
  Tuple t = ops.next();
  if (t != null) {
    out.write(t);
  } else {
    break;
  }
}
```

recursively calls Operator.next()
and Expression.compute()

Pros & cons of this approach?

# Method 2: Vectorization

Interpreting query plans one record at a time is simple, but it's too slow

  » Lots of virtual function calls and branches for each record (recall Jeff Dean's numbers)

Keep recursive interpretation, but make Operators and Expressions run on **batches**

# Implementing Vectorization

```
class TupleBatch {
  // Efficient storage, e.g.
  // schema + column arrays
}

interface Operator {
  TupleBatch next();
}

class Select: Operator {
  Operator parent;
  Expression condition;
}

...
```

```
class ValueBatch {
  // Efficient storage
}

interface Expression {
  ValueBatch compute(
    TupleBatch in);
}

class Times: Expression {
  Expression left, right;
}

...
```

# Typical Implementation

Values stored in columnar arrays (e.g. int[]) with a separate bit array to mark nulls

Tuple batches fit in L1 or L2 cache

Operators use SIMD instructions to update both values and null fields without branching

# Pros & Cons of Vectorization

\+ Faster than record-at-a-time if the query processes many records

\+ Relatively simple to implement

– Lots of nulls in batches if query is selective

– Data travels between CPU & cache a lot

# Method 3: Compilation

Turn the query into executable code

# Compilation Example

$$\Pi_{\text{quanity*price}} (\sigma_{\text{productId=75}} (\text{orders}))$$

generated class with the right field types for orders table

```
class MyQuery {
  void run() {
    Iterator<OrdersTuple> in = openTable("orders");
    for(OrdersTuple t: in) {
      if (t.productId == 75) {
        out.write(Tuple(t.quantity * t.price));
      }
    }
  }
}
```

Can also theoretically generate vectorized code

# Pros & Cons of Compilation

\+ Potential to get fastest possible execution

\+ Leverage existing work in compilers

– Complex to implement

– Compilation takes time

– Generated code may not match hand-written

# What's Used Today?

Depends on context & other bottlenecks

**Transactional databases (e.g. MySQL):** mostly record-at-a-time interpretation

**Analytical systems (Vertica, Spark SQL):** vectorization, sometimes compilation

**ML libs (TensorFlow):** mostly vectorization (the records *are* vectors!), some compilation