# Introduction to Database Systems
## CSE 414

### Lecture 16: Query Evaluation

---

## Announcements

- HW5 + WQ5 due tomorrow

- Midterm this Friday in class!
  - Review session this Wednesday evening
  - See course website

- HW6 will be released later this week
  - Due on Friday 5/11
  - No WQ6 (yet)!

YOU GOT THIS
GOOD LUCK!

---

## Class Overview

- Unit 1: Intro
- Unit 2: Relational Data Models and Query Languages
- Unit 3: Non-relational data
- Unit 4: RDMBS internals and parallel query processing
- Unit 5: DBMS usability, conceptual design
- Unit 6: Transactions
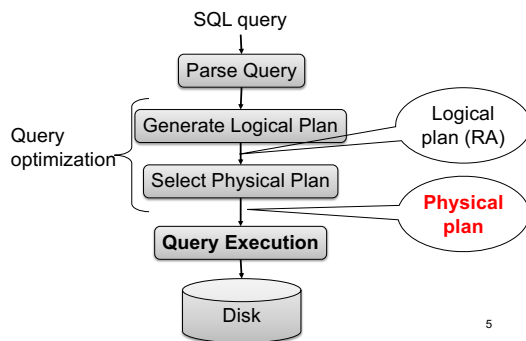- Unit 7: Advanced topics

---

## From Logical RA Plans
## to Physical Plans

---

## Query Evaluation Steps Review

SQL query

Parse Query

Generate Logical Plan → Logical plan (RA)

Query optimization

Select Physical Plan → **Physical plan**

**Query Execution**

Disk

---

## Logical vs Physical Plans

- Logical plans:
  - Created by the parser from the input SQL text
  - Expressed as a relational algebra tree
  - Each SQL query has many possible logical plans

- Physical plans:
  - Goal is to choose an efficient implementation for each operator in the RA tree
  - Each logical plan has many possible physical plans

## Slide 7

Supplier(`sid`, sname, scity, sstate)
Supply(`sid, pno`, quantity)

# Review: Relational Algebra

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and  y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'
```

3. $\pi_{sname}$

2. $\sigma_{scity= 'Seattle' \text{ and } sstate= 'WA' \text{ and } pno=2}$

1. $\bowtie_{sid = sid}$

Supplier        Supply

Relational algebra expression is also called the "logical query plan"

CSE 414 - Spring 2018          7

## Slide 8

Supplier(`sid`, sname, scity, sstate)
Supply(`sid, pno`, quantity)

# Physical Query Plan 1

(On the fly)        $\pi_{sname}$

(On the fly)

$\sigma_{scity= 'Seattle' \text{ and } sstate= 'WA' \text{ and } pno=2}$

(Nested loop)

$\bowtie_{sid = sid}$

Supplier          Supply
(File scan)        (File scan)

A physical query plan is a logical query plan annotated with physical implementation details

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'
```

8

## Slide 9

Supplier(`sid`, sname, scity, sstate)
Supply(`sid, pno`, quantity)

# Physical Query Plan 2

(On the fly)        $\pi_{sname}$

(On the fly)

$\sigma_{scity= 'Seattle' \text{ and } sstate= 'WA' \text{ and } pno=2}$

(Hash join)

$\bowtie_{sid = sid}$

Supplier          Supply
(File scan)        (File scan)

Same logical query plan
Different physical plan

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'
```

9

## Slide 10

Supplier(`sid`, sname, scity, sstate)
Supply(`sid, pno`, quantity)

# Physical Query Plan 3

(On the fly)        $\pi_{sname}$  (d)

(Sort-merge join)   $\bowtie_{sid = sid}$  (c)

(Scan & write to T1)

(a) $\sigma_{scity= 'Seattle' \text{ and } sstate= 'WA'}$      (b) $\sigma_{pno=2}$  (Scan & write to T2)

Supplier          Supply
(File scan)        (File scan)

Different but equivalent logical query plan; different physical plan

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
    and y.pno = 2
    and x.scity = 'Seattle'
    and x.sstate = 'WA'
```

CSE 414 - Spring 2018          10

## Slide 11

# Query Optimization Problem

- For each SQL query… many logical plans

- For each logical plan… many physical plans

- Choosing the best one among them is the goal of *query optimization*

- More on this later in the quarter

CSE 414 - Spring 2018          11

## Slide 12

# Distributed query processing

CSE 414 - Spring 2018          12

## Why compute in parallel?

- Multi-cores:
  - Most processors have multiple cores
  - This trend will likely increase in the future

- Big data: too large to fit in main memory
  - Distributed query processing on 100x-1000x servers
  - Widely available now using cloud services
  - Recall HW3 and motivation for NoSQL!

## Performance Metrics for Parallel DBMSs

Nodes = processors, computers
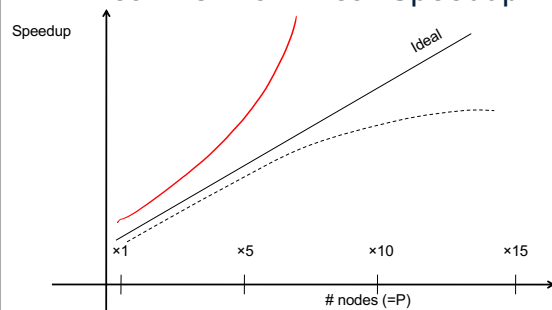
- Speedup:
  - More nodes, same data ➔ higher speed
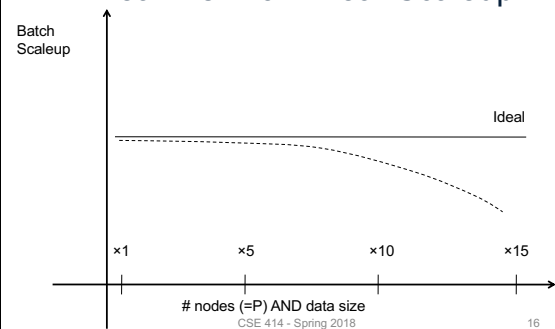
- Scaleup:
  - More nodes, more data ➔ same speed

## Linear v.s. Non-linear Speedup



Speedup

Ideal

×1     ×5     ×10     ×15

# nodes (=P)

## Linear v.s. Non-linear Scaleup



Batch Scaleup

Ideal

×1     ×5     ×10     ×15

# nodes (=P) AND data size

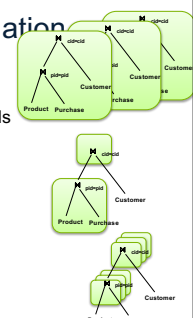## Why Sub-linear Speedup and Scaleup?

- **Startup cost**
  - Cost of starting an operation on many nodes

- **Interference**
  - Contention for resources between nodes

- **Skew**
  - Slowest node becomes the bottleneck

## Approaches to Parallel Query Evaluation



- Inter-query parallelism
  - One query per node
  - Good for transactional (OLTP) workloads

- Inter-operator parallelism
  - Operator per node
  - Good for analytical (OLAP) workloads

- Intra-operator parallelism
  - Operator on multiple nodes
  - Good for both?

We study only intra-operator parallelism: most scalable

# Parallel Data Processing in the 20th Century

---

## Let's parallelize RDBMS

- Data is horizontally partitioned on many servers

- Operators may require data reshuffling

- First let's discuss how to distribute data across multiple nodes / servers

---

## Horizontal Data Partitioning

Data:

Servers:

| K | A | B |
|---|---|---|
| … | … | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |
| | | |

| 1 | 2 | . . . | P |
|---|---|-------|---|

---

## Horizontal Data Partitioning

Data:

Servers:

| K | A | B |
|---|---|---|
| … | … | |

| 1 | 2 | . . . | P |
|---|---|-------|---|

Which tuples go to what server?

---

## Recall: Horizontal Data Partitioning

- **Block Partition:**
  - Partition tuples arbitrarily s.t. size($R_1$)≈ … ≈ size($R_P$)

- **Hash partitioned on attribute A:**
  - Tuple t goes to chunk i, where i = h(t.A) mod P + 1
  - Recall: calling hash fn's is free in this class

- **Range partitioned on attribute A:**
  - Partition the range of A into  $-\infty = v_0 < v_1 < … < v_P = \infty$
  - Tuple t goes to chunk i, if $v_{i-1} < t.A < v_i$

---

## Uniform Data v.s. Skewed Data

- Let R(K,A,B,C); which of the following partition methods may result in skewed partitions?

- **Block partition**      Uniform

- **Hash-partition**
  - On the key K      Uniform      Assuming good hash function
  - On the attribute A      May be skewed      E.g. when all records have the same value of the attribute A, then all records end up in the same partition

**Keep this in mind in the next few slides**

---

## Parallel Execution of RA Operators: Grouping

Data: $R(\underline{K},A,B,C)$
Query: $\gamma_{A,sum(C)}(R)$

How to compute group by if:

- R is hash-partitioned on A ?
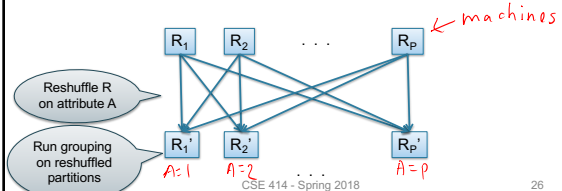- R is block-partitioned ?
- R is hash-partitioned on K ?

CSE 414 - Spring 2018    25

---

## Parallel Execution of RA Operators: Grouping

Data: $R(\underline{K},A,B,C)$
Query: $\gamma_{A,sum(C)}(R)$

- R is block-partitioned or hash-partitioned on K

*machines*

$R_1$   $R_2$   . . .   $R_P$

Reshuffle R on attribute A

Run grouping on reshuffled partitions

$R_1'$   $R_2'$   . . .   $R_P'$

$A=1$   $A=2$   $A=P$

CSE 414 - Spring 2018    26

---

## Speedup and Scaleup

- Consider:
  - Query: $\gamma_{A,sum(C)}(R)$
  - Runtime: only consider I/O costs
- If we double the number of nodes P, what is the new running time?
  - Half (each server holds ½ as many chunks)
- If we double both P and the size of R, what is the new running time?
  - Same (each server holds the same # of chunks)

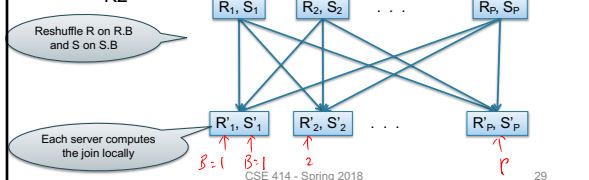**But only if the data is without skew!** 27

---

## Skewed Data

- $R(\underline{K},A,B,C)$
- Informally: we say that the data is skewed if one server holds much more data that the average
- E.g., we hash-partition on A, and some value of A occurs very many times ("Justin Bieber")
- Then the server holding that value will be skewed

CSE 414 - Spring 2018    28

---

## Parallel Execution of RA Operators: Partitioned Hash-Join

- Data: $R(\underline{K1}, A, B), S(\underline{K2}, B, C)$
- Query: $R(\underline{K1}, A, B) \bowtie S(\underline{K2}, B, C)$
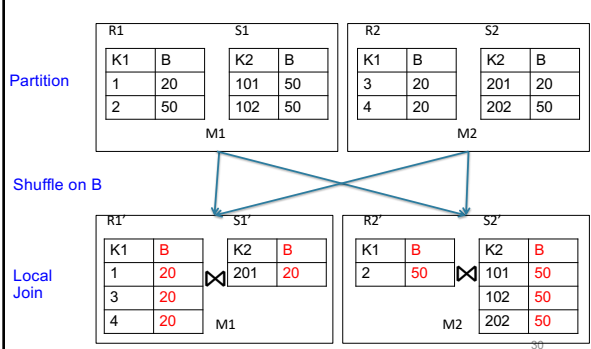  - Initially, both R and S are partitioned on K1 and K2

*machine node*

$R_1, S_1$   $R_2, S_2$   . . .   $R_P, S_P$

Reshuffle R on R.B and S on S.B

Each server computes the join locally

$R'_1, S'_1$   $R'_2, S'_2$   . . .   $R'_P, S'_P$

$B=1$   $B=1$   $2$   $P$

CSE 414 - Spring 2018    29

---

## Parallel Join Illustration

Data: $R(\underline{K1},A, B), S(\underline{K2}, B, C)$
Query: $R(\underline{K1},A,B) \bowtie S(\underline{K2},B,C)$

**Partition**

| R1 | | | S1 | |
|----|----|----|----|----|
| K1 | B | | K2 | B |
| 1 | 20 | | 101 | 50 |
| 2 | 50 | | 102 | 50 |

M1

| R2 | | | S2 | |
|----|----|----|----|----|
| K1 | B | | K2 | B |
| 3 | 20 | | 201 | 20 |
| 4 | 20 | | 202 | 50 |

M2

**Shuffle on B**

**Local Join**

| R1' | | | S1' | |
|----|----|----|----|----|
| K1 | B | | K2 | B |
| 1 | 20 | ⋈ | 201 | 20 |
| 3 | 20 | | | |
| 4 | 20 | | | |

M1

| R2' | | | S2' | |
|----|----|----|----|----|
| K1 | B | | K2 | B |
| 2 | 50 | ⋈ | 101 | 50 |
| | | | 102 | 50 |
| | | | 202 | 50 |

M2

30

## Slide 31

Data: R(A, B), S(C, D)
Query: R(A,B) ⋈$_{B=C}$ S(C,D)

### Broadcast Join

Broadcast S

Reshuffle R on R.B

| R$_1$ | R$_2$ | . . . | R$_P$ | | S |

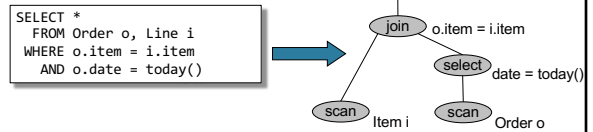| R'$_1$, S | R'$_2$, S | . . . | R'$_P$, S |

Why would you want to do this?

## Slide 32

Order(oid, item, date), Line(item, …)

### Putting it Together:
### Example Parallel Query Plan

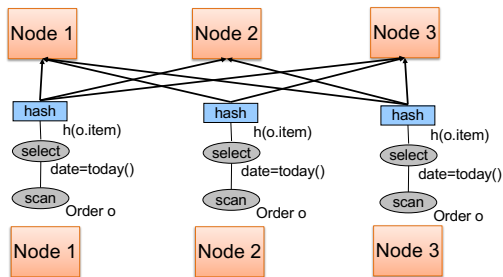*Find all orders from today, along with the items ordered*

```
SELECT *
  FROM Order o, Line i
 WHERE o.item = i.item
   AND o.date = today()
```

join   o.item = i.item

select   date = today()

scan   Item i          scan   Order o

## Slide 33

Order(oid, item, date), Line(item, …)

### Example Parallel Query Plan

join   o.item = i.item
select   date = today()
scan   Order o

| Node 1 | Node 2 | Node 3 |

hash   h(o.item)
select   date=today()
scan   Order o

hash   h(o.item)
select   date=today()
scan   Order o

hash   h(o.item)
select   date=today()
scan   Order o

| Node 1 | Node 2 | Node 3 |

## Slide 34

Order(oid, item, date), Line(item, …)

### Example Parallel Query Plan

join   o.item = i.item
date = today()
scan   Item i          Order o

| Node 1 | Node 2 | Node 3 |

hash   h(i.item)
scan   Item i

hash   h(i.item)
scan   Item i

hash   h(i.item)
scan   Item i

| Node 1 | Node 2 | Node 3 |

## Slide 35

Order(oid, item, date), Line(item, …)

### Example Parallel Query Plan

join   o.item = i.item        join   o.item = i.item        join   o.item = i.item

| Node 1 | Node 2 | Node 3 |

contains all orders and all
lines where hash(item) = 3

contains all orders and all
lines where hash(item) = 2

contains all orders and all
lines where hash(item) = 1