

 Intro to Database Systems (15-445/645)

13 Query Execution

Part 2

Carnegie
Mellon
University

SPRING
2023

Charlie
Garrod

ADMINISTRIVIA

Project 2 is ongoing

- First checkpoint was Friday, March 3rd (15% of P2 grade)
- Overall due Wednesday, March 22nd (85% of P2 grade)

LAST TIME: QUERY EXECUTION

Processing models

- Iterator model
- Materialization model
- Vectorized / batch model

Modification queries

- The Halloween problem

Expression evaluation

- JIT compilation

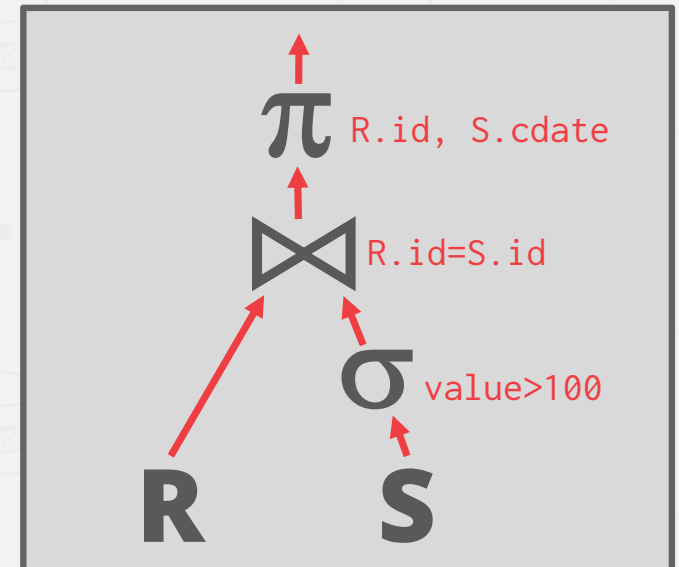
QUERY EXECUTION

We discussed in the last class how to compose operators together into a plan to execute an arbitrary query.

We assumed that the queries execute with a single worker (e.g., a thread).

We will now discuss how to execute queries using multiple workers.

```
SELECT R.id, S.cdate  
FROM R JOIN S  
ON R.id = S.id  
WHERE S.value > 100
```



WHY CARE ABOUT PARALLEL EXECUTION?

Increased performance for potentially the same hardware resources.

- Higher Throughput
- Lower Latency

Increased responsiveness of the system.

Potentially lower *total cost of ownership* (TCO)

- Fewer machines means less parts / physical footprint / energy consumption.

PARALLEL AND DISTRIBUTED DATABASES

Database is spread out across multiple resources to improve different aspects of the DBMS.

Appears as a single logical database instance to the application, regardless of physical organization.

→ SQL query for a single-resource DBMS should generate same result on a parallel or distributed DBMS.

PARALLEL VS. DISTRIBUTED

Parallel DBMSs

- Resources are physically close to each other.
- Resources communicate over high-speed interconnect.
- Communication is assumed to be cheap and reliable.

Distributed DBMSs

- Resources can be far from each other.
- Resources communicate using slow(er) interconnect.
- Communication cost and problems cannot be ignored.

TODAY'S AGENDA

Process Models
Execution Parallelism
I/O Parallelism

PROCESS MODEL

A DBMS's process model defines how the system is architected to support concurrent requests from a multi-user application.

A worker is the DBMS component that is responsible for executing tasks on behalf of the client and returning the results.

PROCESS MODEL

Approach #1: Process per DBMS Worker

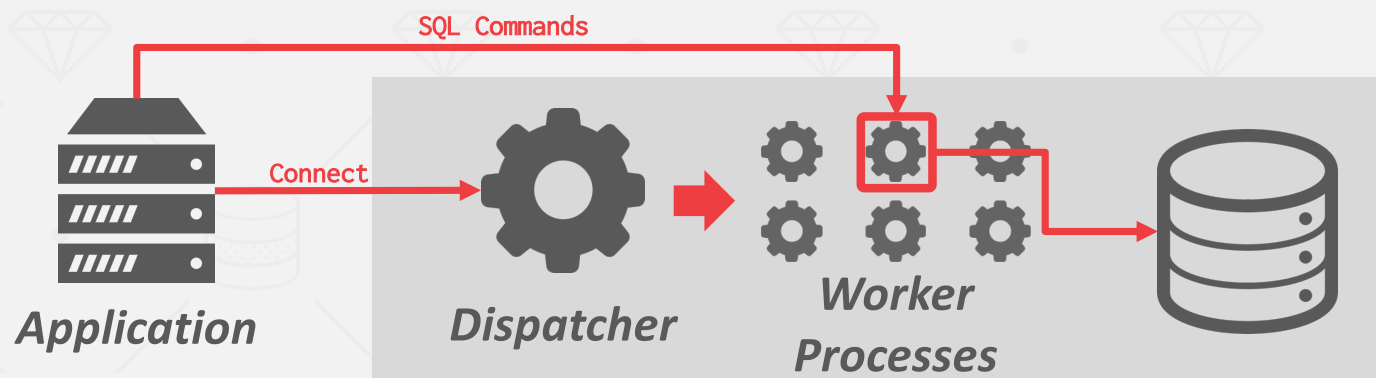
Approach #2: Thread per DBMS Worker

Approach #3: Embedded DBMS

PROCESS PER WORKER

Each worker is a separate OS process.

- Relies on OS scheduler.
- Use shared-memory for global data structures.
- A process crash does not take down entire system.
- Examples: IBM DB2, Postgres, Oracle



THREAD PER WORKER

Single process with multiple worker threads.

- DBMS (mostly) manages its own scheduling.
- May or may not use a dispatcher thread.
- Thread crash (may) kill the entire system.
- Examples: **MSSQL**, MySQL, DB2, Oracle (2014)

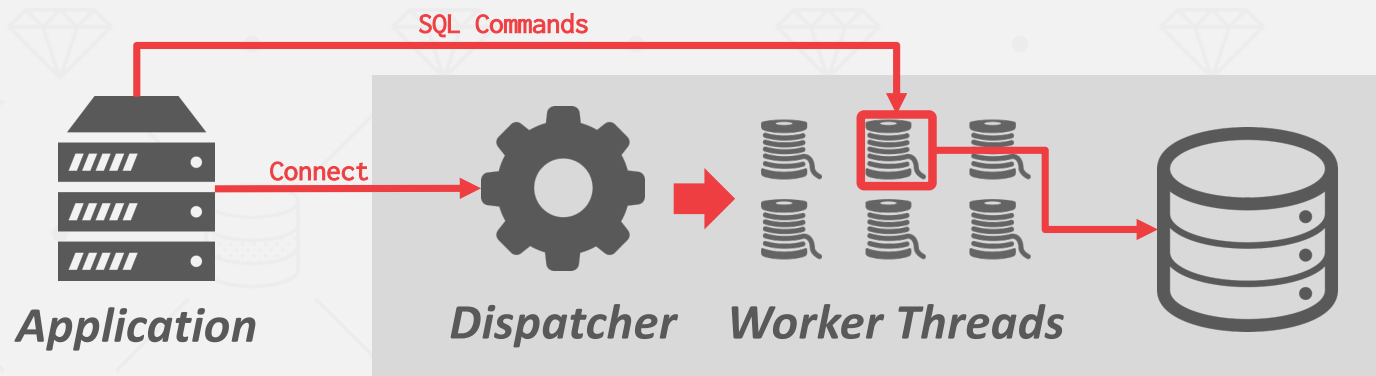
Almost every DBMS created in the last 20 years!

Microsoft®
SQL Server®

MySQL™

IBM® **DB2**®

ORACLE®



SCHEDULING

For each query plan, the DBMS decides where, when, and how to execute it.

- How many tasks should it use?
- How many CPU cores should it use?
- What CPU core should the tasks execute on?
- Where should a task store its output?

The DBMS *always* knows more than the OS.

SQL SERVER – SQLOS

SQLOS is a user-level OS layer that runs inside of the DBMS and manages provisioned hardware resources.

- Determines which tasks are scheduled onto which threads.
- Also manages I/O scheduling and higher-level concepts like logical database locks.

Non-preemptive thread scheduling through instrumented DBMS code.

SQL SERVER - SQLOS

SQLOS quantum is 4 ms but the scheduler cannot enforce that.

DBMS developers must add explicit yield calls in various locations in the source code.

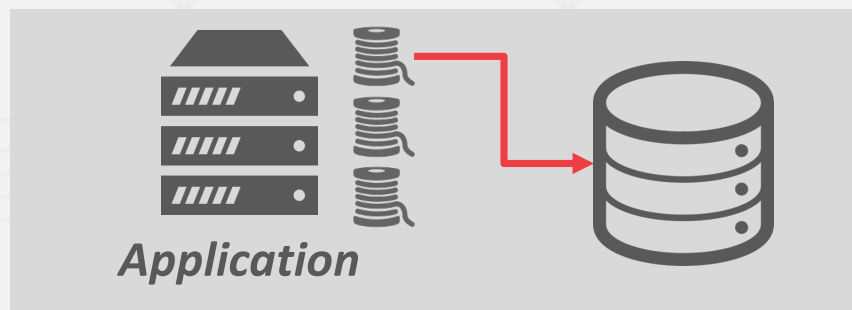
```
SELECT * FROM R WHERE R.val = ?
```

```
last = now()
for tuple in R:
    if now() - last > 4ms:
        ➡ yield
        last = now()
    if eval(predicate, tuple, params):
        emit(tuple)
```

EMBEDDED DBMS

DBMS runs inside of the same address space as the application. Application is (mostly) responsible for threads and scheduling.

The application may support outside connections.
→ Examples: BerkeleyDB, SQLite, RocksDB, LevelDB



PROCESS MODELS

Advantages of a multi-threaded architecture:

- Less overhead per context switch.
- Do not have to manage shared memory.

The thread per worker model does not mean that the DBMS supports intra-query parallelism.

Andy is not aware of any new DBMS from last 15 years that doesn't use native OS threads unless they are Redis or Postgres forks.

INTER- VS. INTRA-QUERY PARALLELISM

Inter-Query: Execute multiple disparate queries simultaneously.

→ Increases throughput & reduces latency.

Intra-Query: Execute the operations of a single query in parallel.

→ Decreases latency for long-running queries, especially for OLAP queries.

INTER-QUERY PARALLELISM

Improve overall performance by allowing multiple queries to execute simultaneously.

If queries are read-only, then this requires almost no explicit coordination between queries.

→ Buffer pool can handle most of the sharing if necessary

Lecture #15

If multiple queries are updating the database at the same time, then this is hard to do correctly...

INTRA-QUERY PARALLELISM

Improve the performance of a single query by executing its operators in parallel.

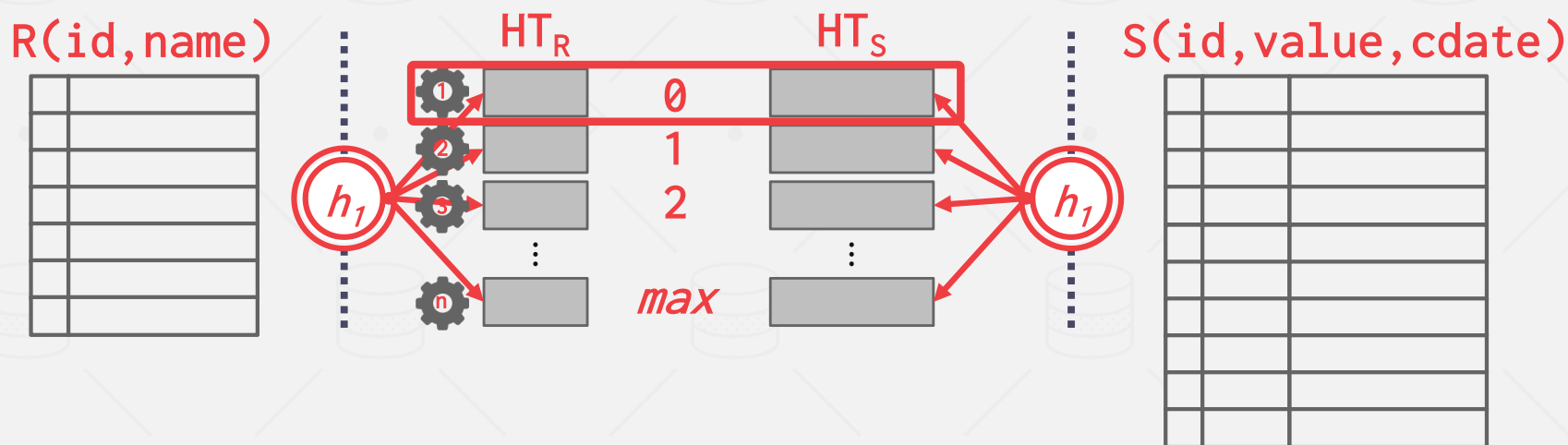
Think of organization of operators in terms of a *producer/consumer* paradigm.

There are parallel versions of every operator.

→ Can either have multiple threads access centralized data structures or use partitioning to divide work up.

e.g., PARALLEL GRACE HASH JOIN

Use a separate worker to perform the join for each level of buckets for **R** and **S** after partitioning.



INTRA-QUERY PARALLELISM

Approach #1: Intra-Operator (Horizontal)

Approach #2: Inter-Operator (Vertical)

Approach #3: Bushy

INTRA-OPERATOR PARALLELISM

Approach #1: Intra-Operator (Horizontal)

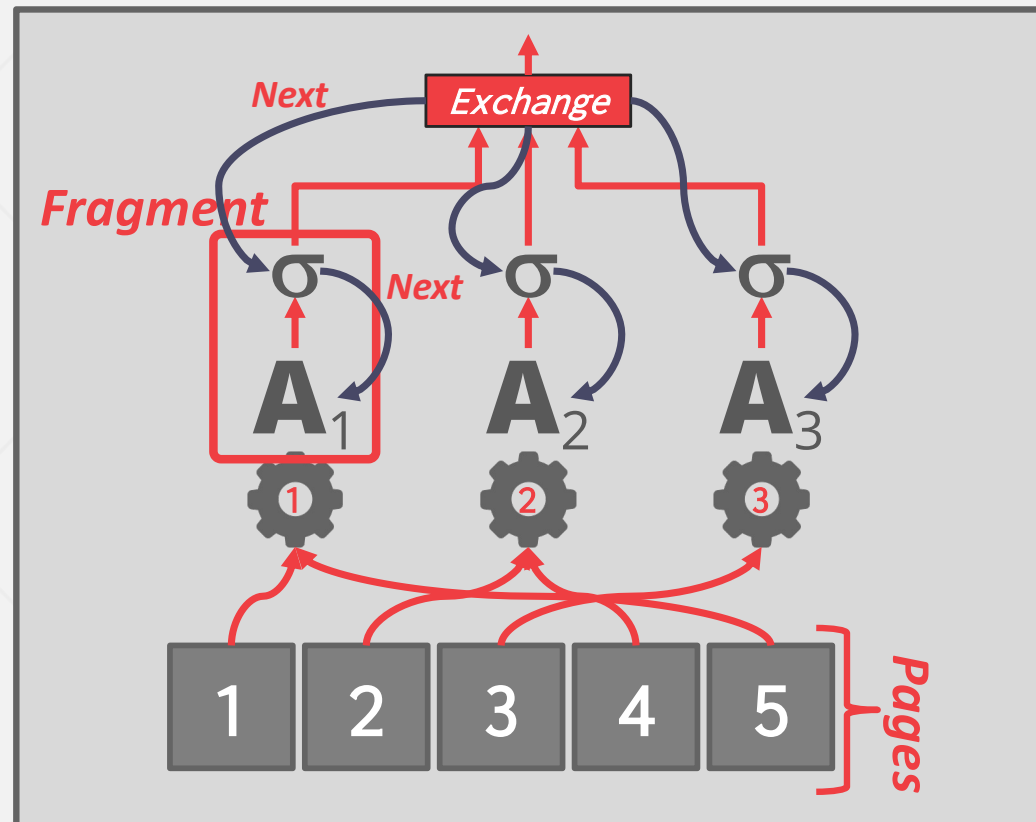
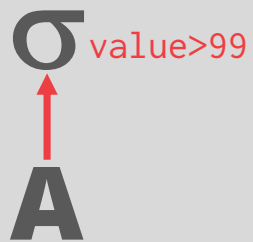
→ Decompose operators into independent fragments that perform the same function on different subsets of data.

The DBMS inserts an exchange operator into the query plan to coalesce/split results from multiple children/parent operators.

→ Postgres calls this "gather"

INTRA-OPERATOR PARALLELISM

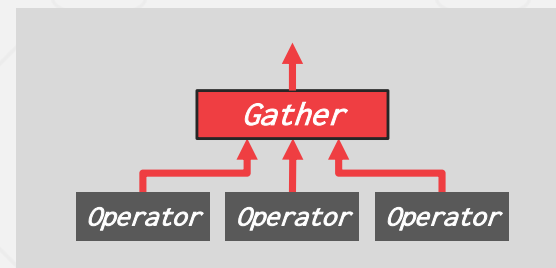
```
SELECT * FROM A  
WHERE A.val > 99
```



EXCHANGE OPERATOR

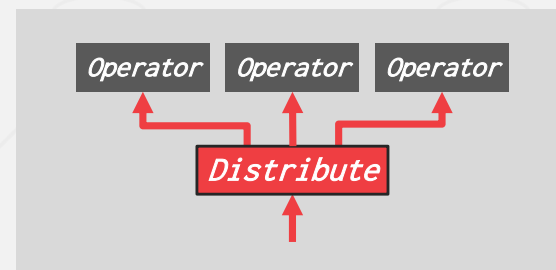
Exchange Type #1 – Gather

→ Combine the results from multiple workers into a single output stream.



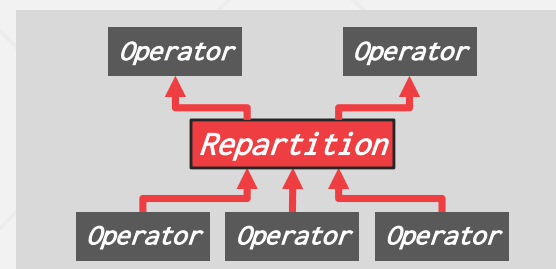
Exchange Type #2 – Distribute

→ Split a single input stream into multiple output streams.



Exchange Type #3 – Repartition

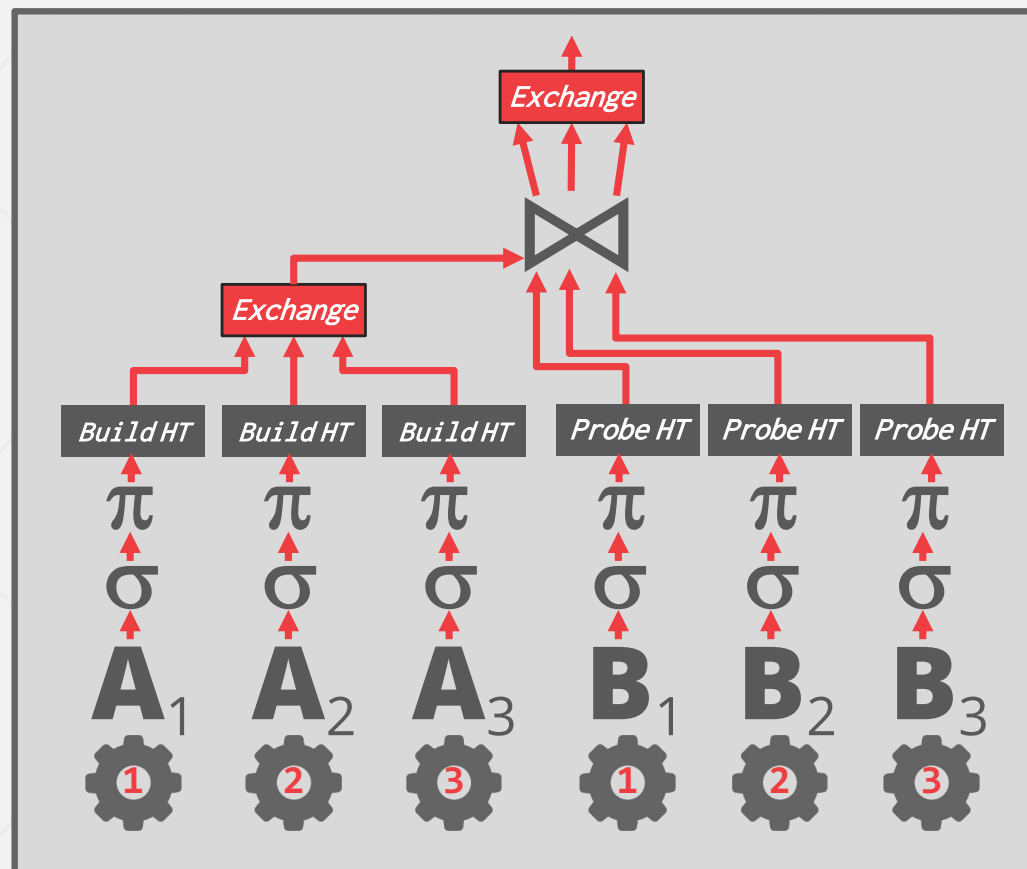
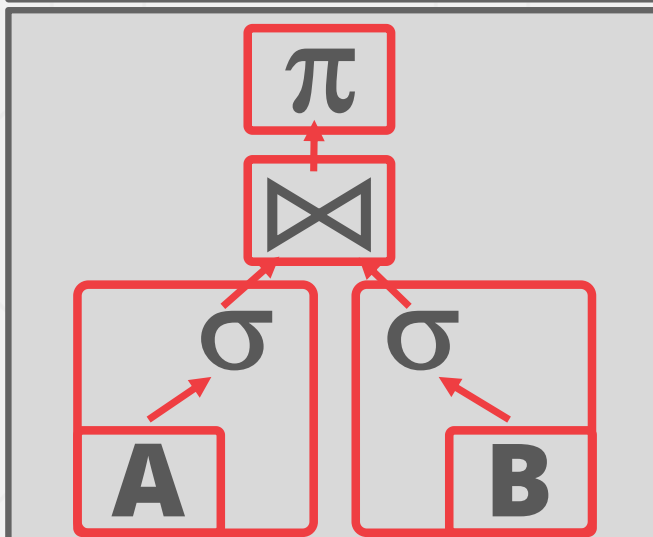
→ Shuffle multiple input streams across multiple output streams.



INTRA-OPERATOR PARALLELISM

```

SELECT A.id, B.value
FROM A JOIN B
      ON A.id = B.id
WHERE A.value < 99
      AND B.value > 100
  
```



INTER-OPERATOR PARALLELISM

Approach #2: Inter-Operator (Vertical)

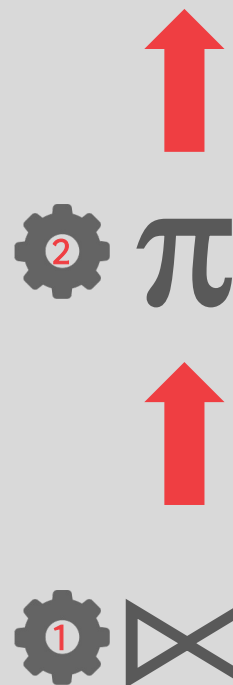
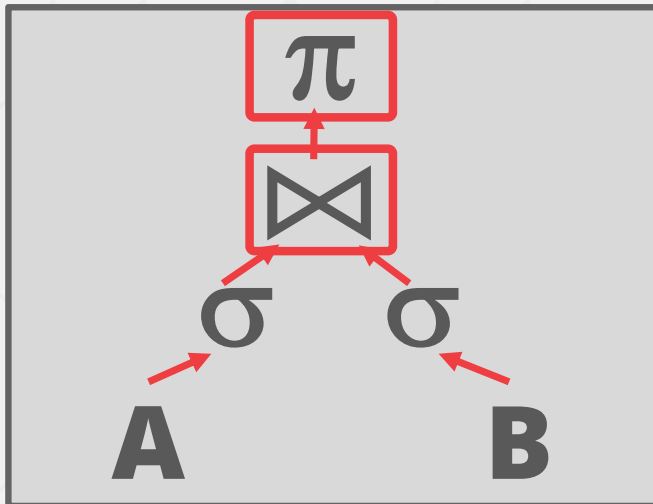
- Operations are overlapped in order to pipeline data from one stage to the next without materialization.
- Workers execute operators from different segments of a query plan at the same time.
- More common in streaming systems (continuous queries)

Also called pipeline parallelism.



INTER-OPERATOR PARALLELISM

```
SELECT A.id, B.value  
FROM A JOIN B  
ON A.id = B.id  
WHERE A.value < 99  
AND B.value > 100
```



for $r \in \text{incoming}$:
emit(πr)

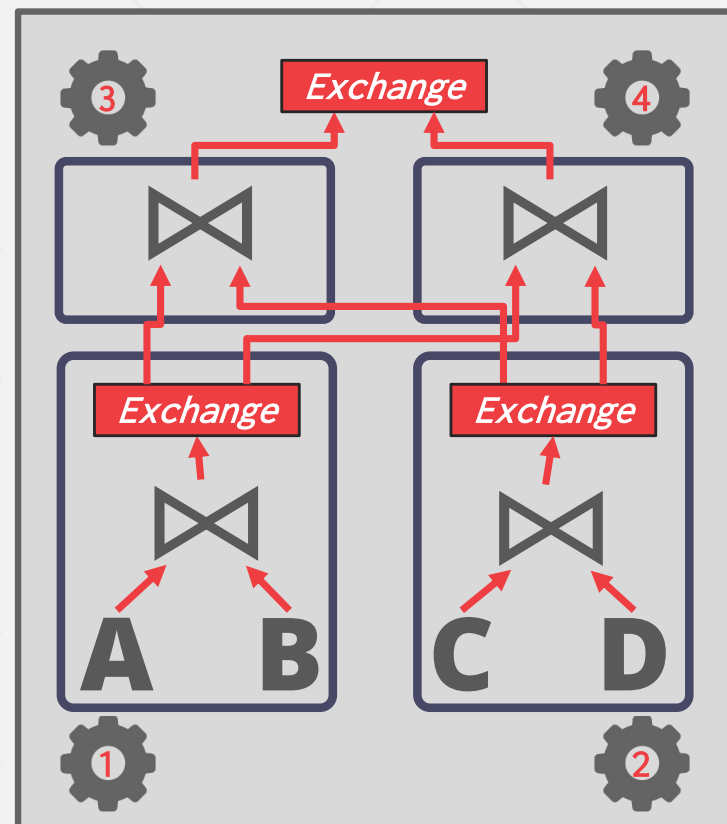
for $r_1 \in \text{outer}$:
for $r_2 \in \text{inner}$:
emit($r_1 \bowtie r_2$)

BUSHY PARALLELISM

Approach #3: Bushy Parallelism

- Hybrid of intra- and inter-operator parallelism where workers execute multiple operators from different segments of a query plan at the same time.
- Still need exchange operators to combine intermediate results from segments.

```
SELECT *  
FROM A JOIN B JOIN C JOIN D
```



OBSERVATION

Using additional processes/threads to execute queries in parallel won't help if the disk is always the main bottleneck.

It can sometimes make the DBMS's performance worse if worker is accessing different segments of the disk at the same time.

I/O PARALLELISM

Split the DBMS across multiple storage devices to improve disk bandwidth latency.

Many different options that have trade-offs:

- Multiple Disks per Database
- One Database per Disk
- One Relation per Disk
- Split Relation across Multiple Disks

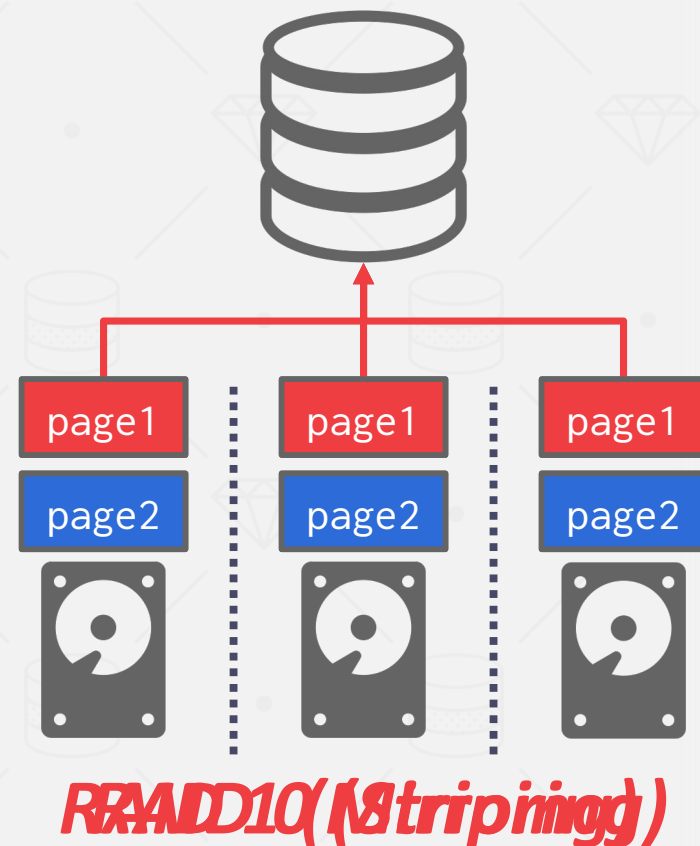
Some DBMSs support this natively. Others require admin to configure outside of DBMS.

MULTI-DISK PARALLELISM

Configure OS/hardware to store the DBMS's files across multiple storage devices.

- Storage Appliances
- RAID Configuration

This is transparent to the DBMS.



DATABASE PARTITIONING

Some DBMSs allow you to specify the disk location of each individual database.

→ The buffer pool manager maps a page to a disk location.

This is also easy to do at the filesystem level if the DBMS stores each database in a separate directory.

→ The DBMS recovery log file might still be shared if transactions can update multiple databases.

PARTITIONING

Split single logical table into disjoint physical segments that are stored/managed separately.

Partitioning should (ideally) be transparent to the application.

→ The application should only access logical tables and not have to worry about how things are physically stored.

We will cover this further when we talk about distributed databases.

CONCLUSION

Parallel execution is important, which is why (almost) every major DBMS supports it.

However, it is hard to get right.

- Coordination Overhead
- Scheduling
- Concurrency Issues
- Resource Contention