



Lecture #18

Carnegie Mellon University

ADVANCED DATABASE SYSTEMS

Query Execution & Scheduling

@Andy_Pavlo // 15-721 // Spring 2018

TODAY'S AGENDA

Process Models

Query Parallelization

Data Placement

Scheduling

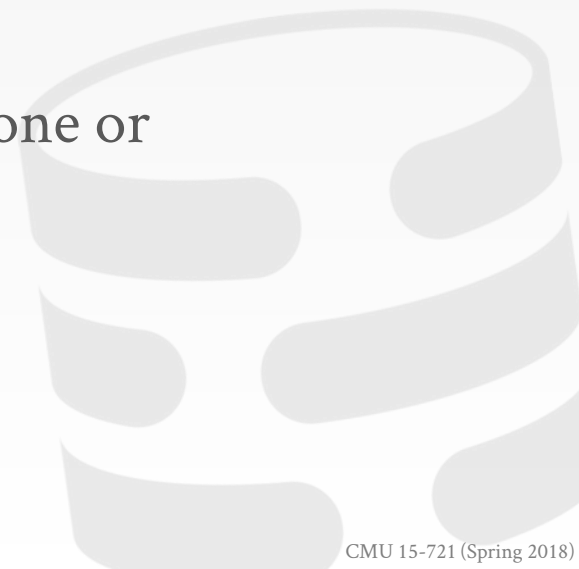


QUERY EXECUTION

A query plan is comprised of **operators**.

An **operator instance** is an invocation of an operator on some segment of data.

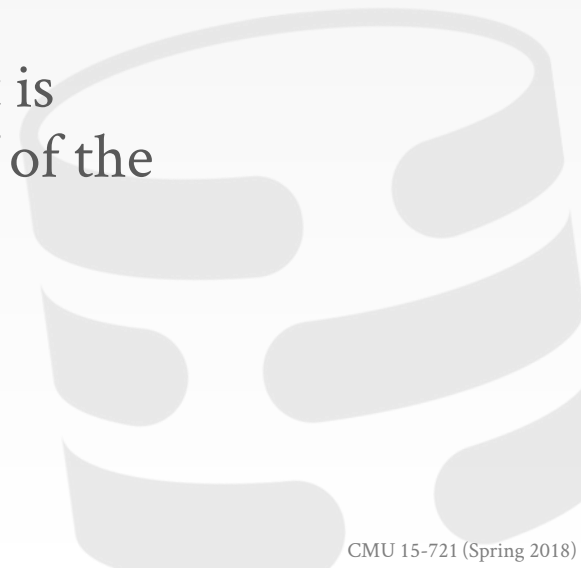
A **task** is the execution of a sequence of one or more operator instances.



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 MODELS

Approach #1: Process per DBMS Worker

Approach #2: Process Pool

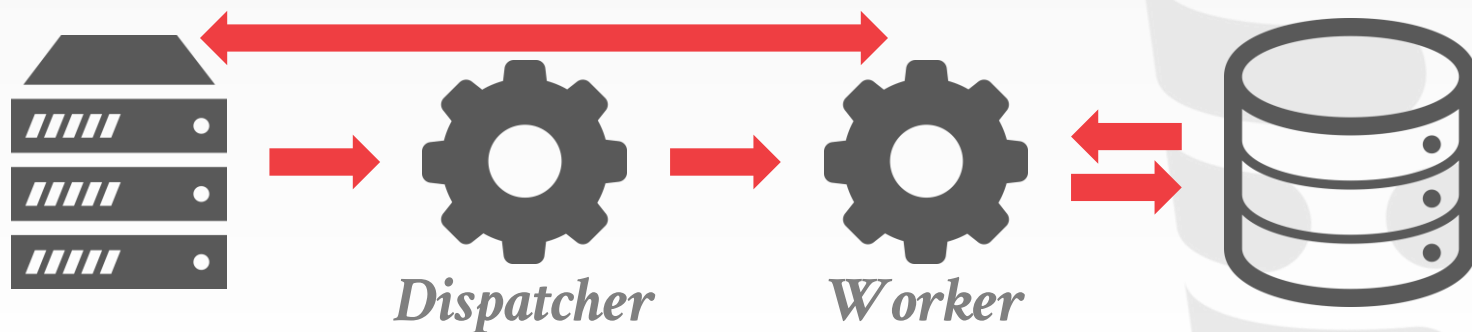
Approach #3: Thread per DBMS Worker



PROCESS PER WORKER

Each worker is a separate OS process.

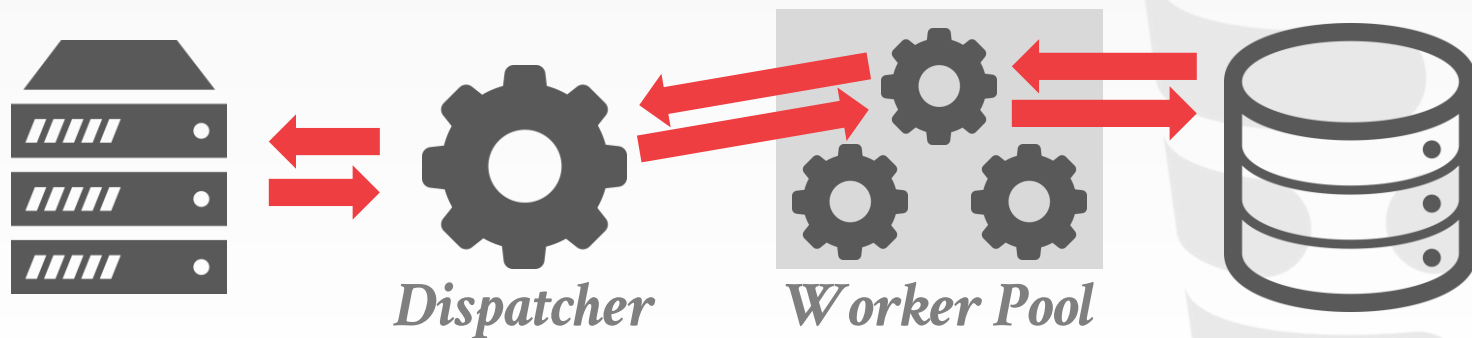
- Relies on OS scheduler.
- Use shared-memory for global data structures.
- A process crash doesn't take down entire system.
- Examples: IBM DB2, Postgres, Oracle



PROCESS POOL

A worker uses any process that is free in a pool

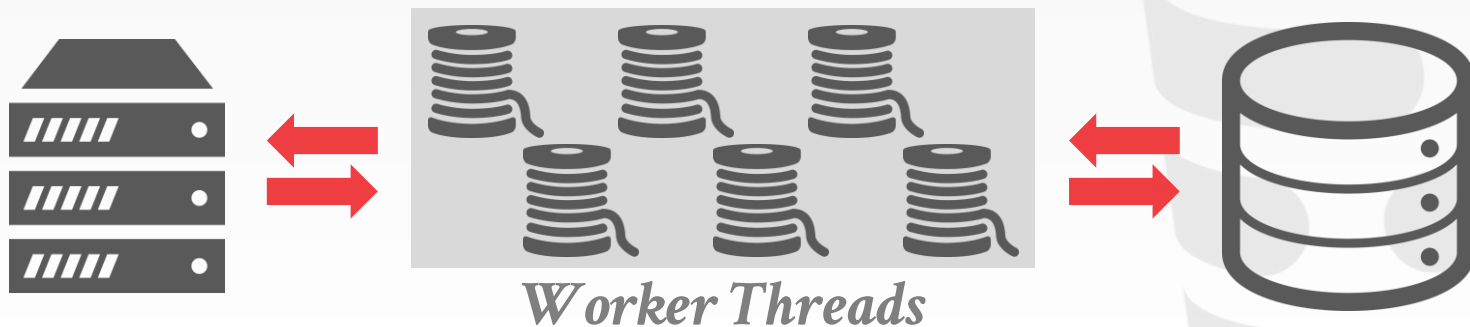
- Still relies on OS scheduler and shared memory.
- Bad for CPU cache locality.
- Examples: IBM DB2, Postgres (2015)



THREAD PER WORKER

Single process with multiple worker threads.

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



PROCESS MODELS

Using a multi-threaded architecture has several advantages:

- Less overhead per context switch.
- Don't have to manage shared memory.

The thread per worker model does not mean that you have intra-query parallelism.

I am not aware of any new DBMS built in the last 7-8 years that doesn't use threads.

SCHEDULING

For each query plan, the DBMS has to decide 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.



INTER-QUERY PARALLELISM

Improve overall performance by allowing multiple queries to execute simultaneously.

→ Provide the illusion of isolation through concurrency control scheme.

The difficulty of implementing a concurrency control scheme is not significantly affected by the DBMS's process model.

INTRA-QUERY PARALLELISM

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

Approach #1: Intra-Operator (Horizontal)

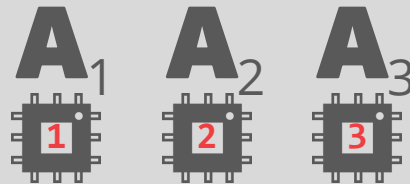
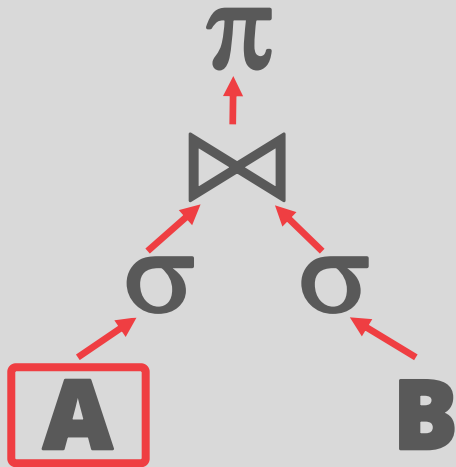
→ Operators are decomposed into independent instances that perform the same function on different subsets of data.

Approach #2: Inter-Operator (Vertical)

→ Operations are overlapped in order to pipeline data from one stage to the next without materialization.

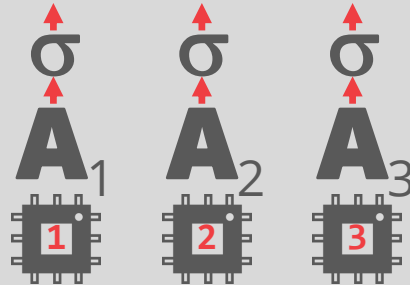
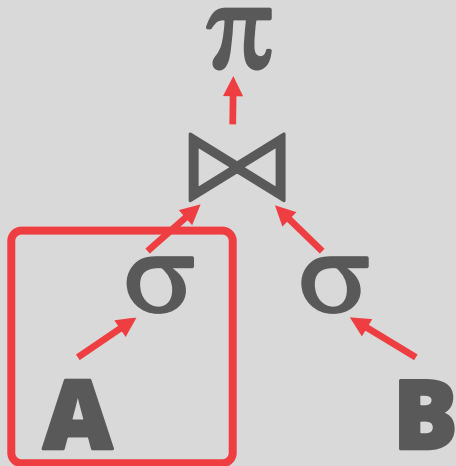
INTRA-OPERATOR PARALLELISM

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



INTRA-OPERATOR PARALLELISM

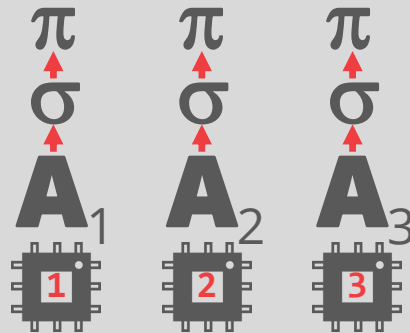
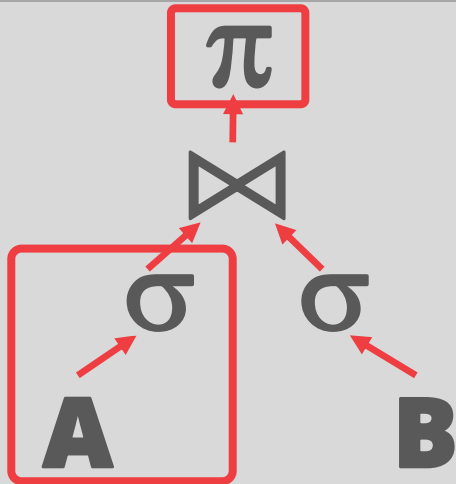
```
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
AND A.value < 99
AND B.value > 100
```



INTRA-OPERATOR PARALLELISM

```

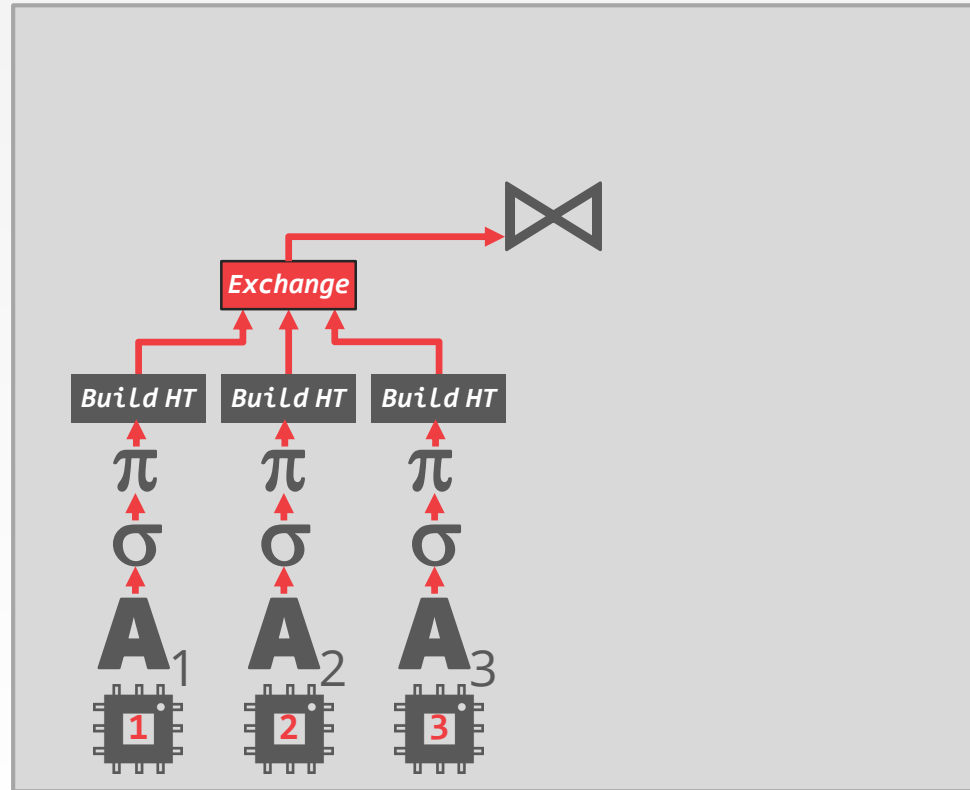
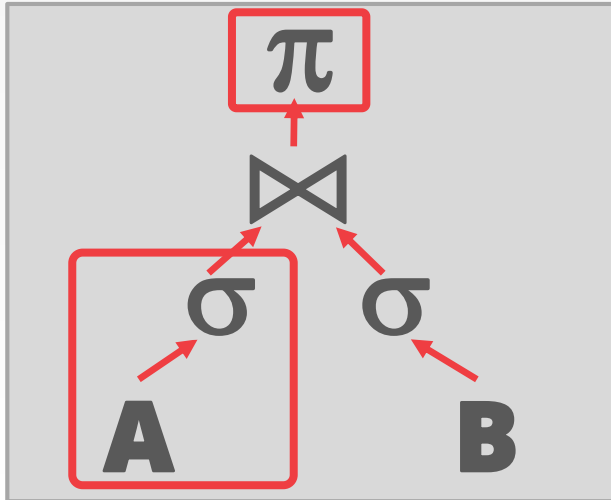
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
      AND A.value < 99
      AND B.value > 100
  
```



INTRA-OPERATOR PARALLELISM

```

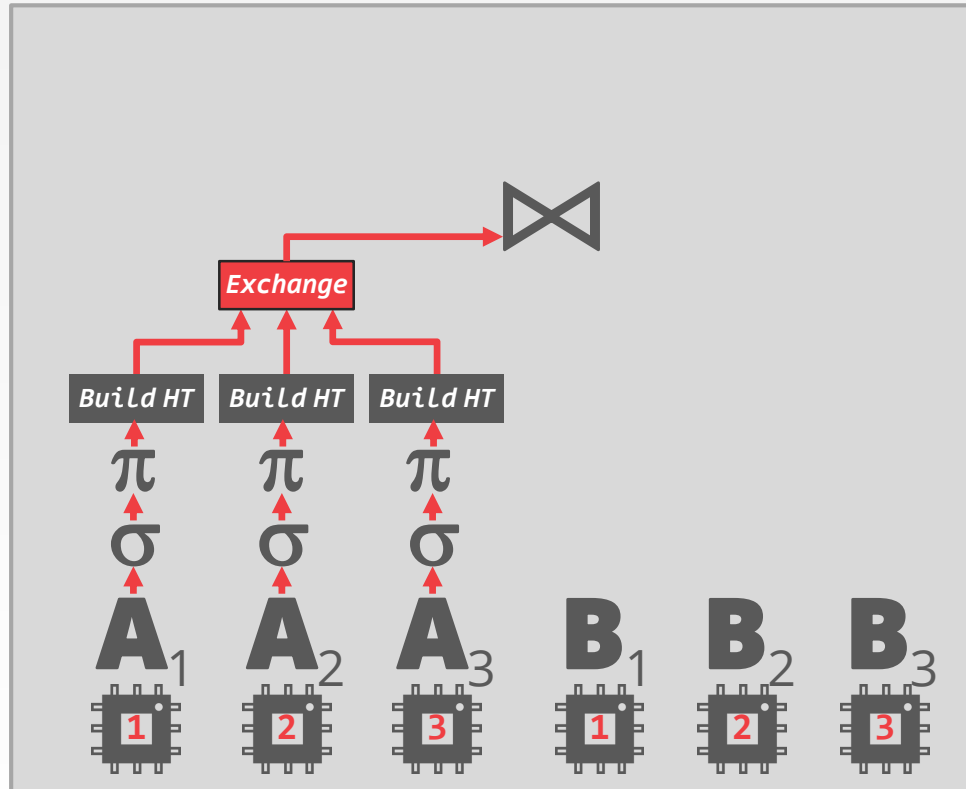
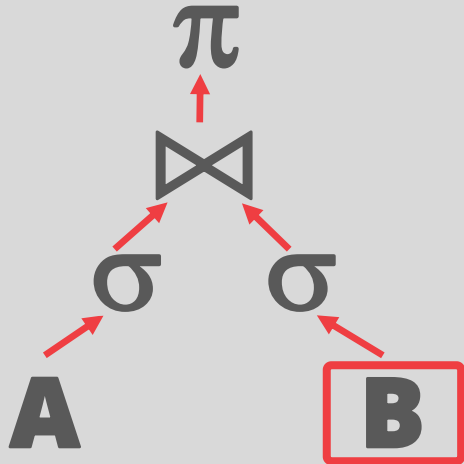
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
AND A.value < 99
AND B.value > 100
  
```



INTRA-OPERATOR PARALLELISM

```

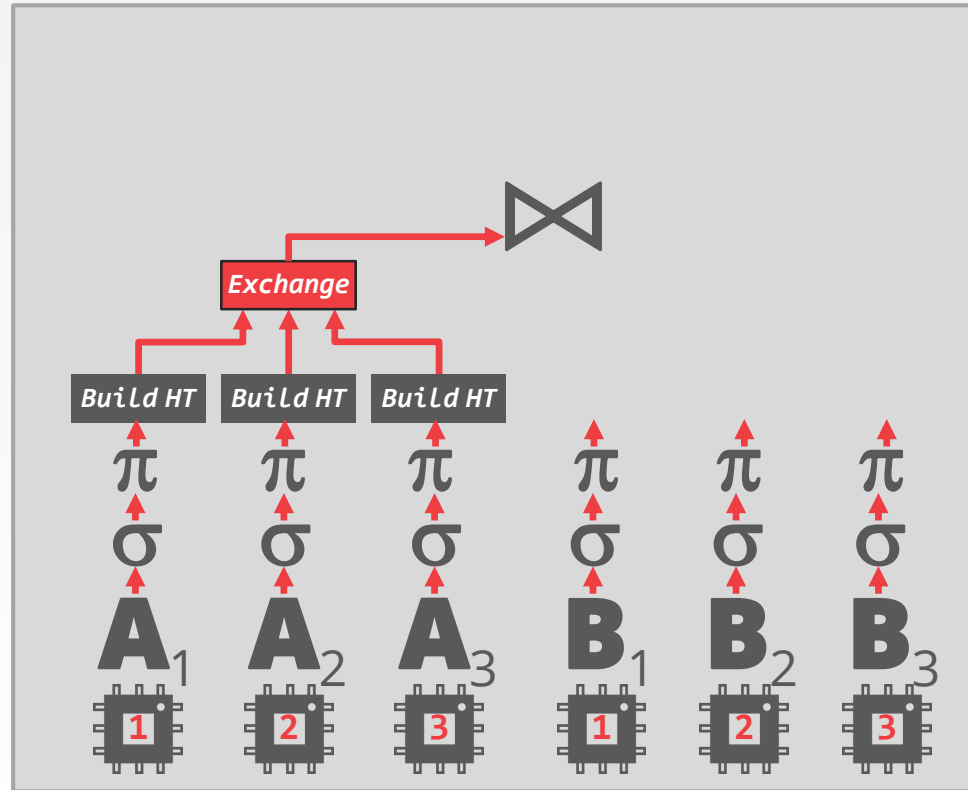
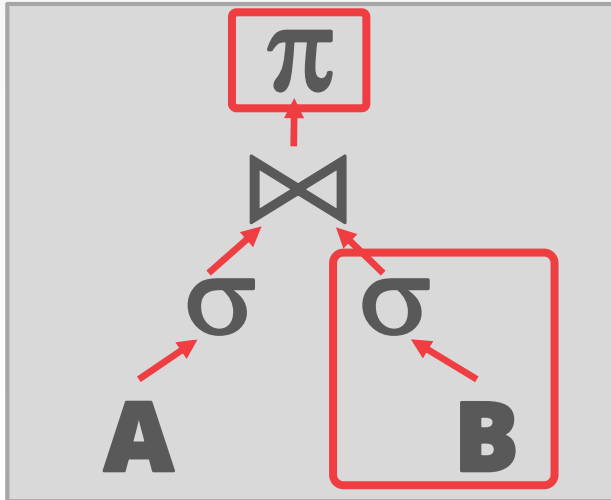
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
AND A.value < 99
AND B.value > 100
  
```



INTRA-OPERATOR PARALLELISM

```

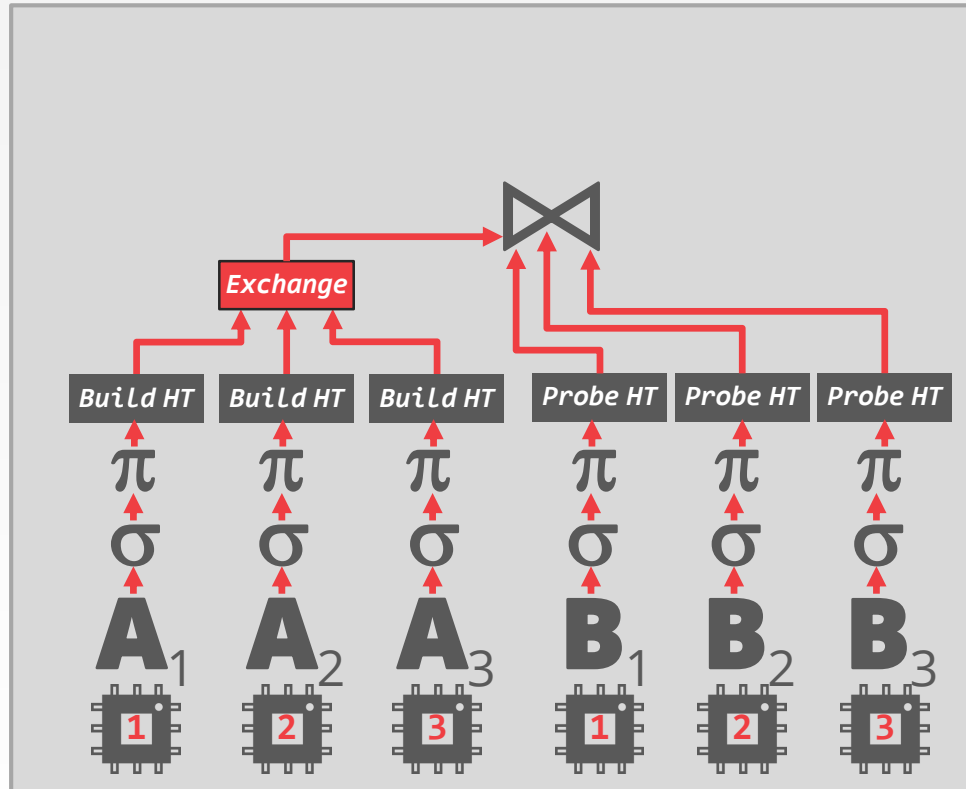
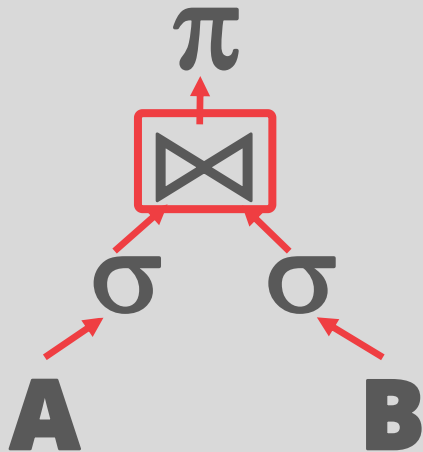
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
      AND A.value < 99
      AND B.value > 100
  
```



INTRA-OPERATOR PARALLELISM

```

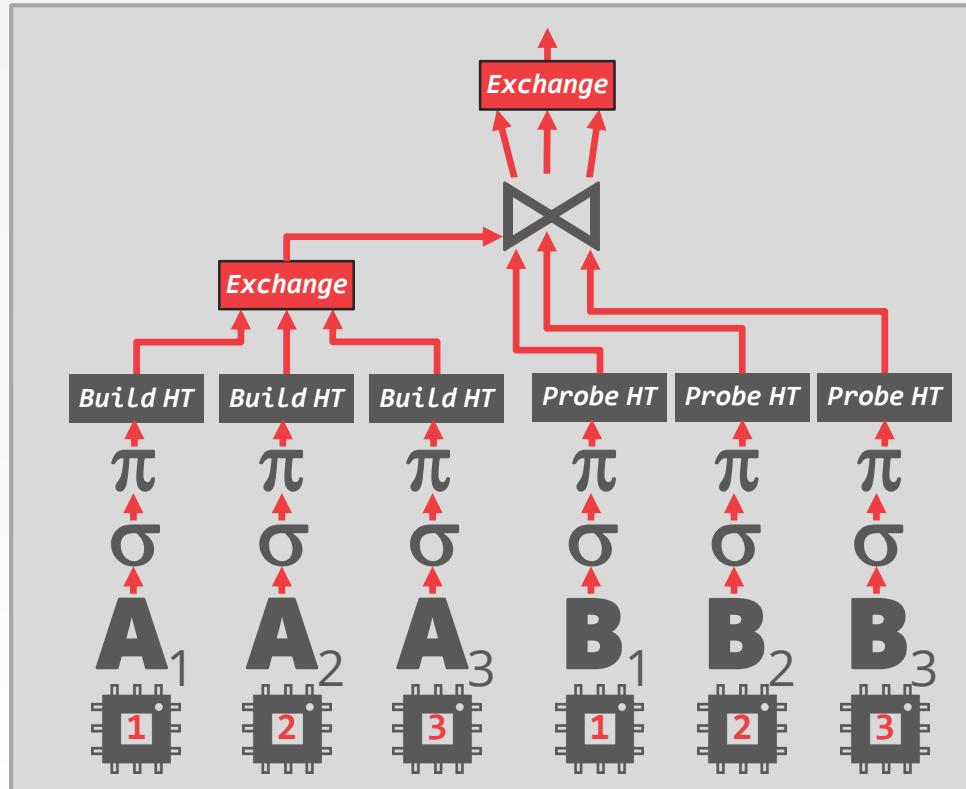
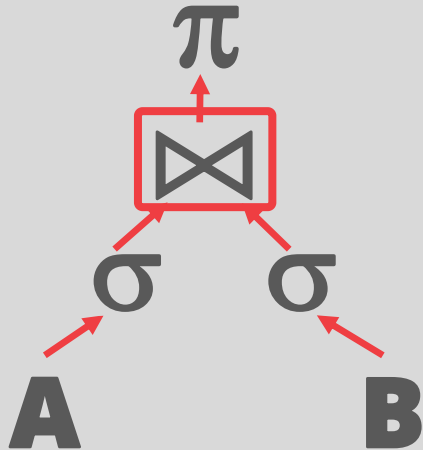
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
AND A.value < 99
AND B.value > 100
  
```



INTRA-OPERATOR PARALLELISM

```

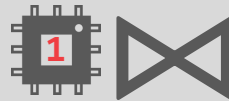
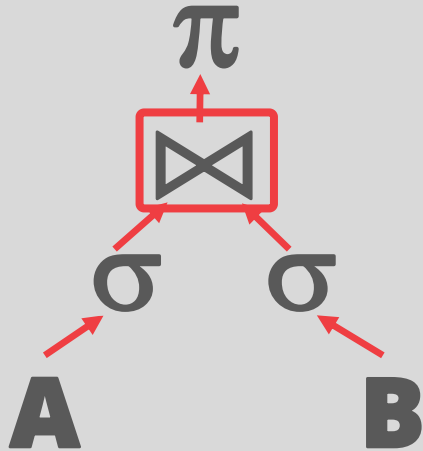
SELECT A.id, B.value
FROM A, B
WHERE A.id = B.id
      AND A.value < 99
      AND B.value > 100
  
```



INTER-OPERATOR PARALLELISM

```

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



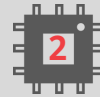
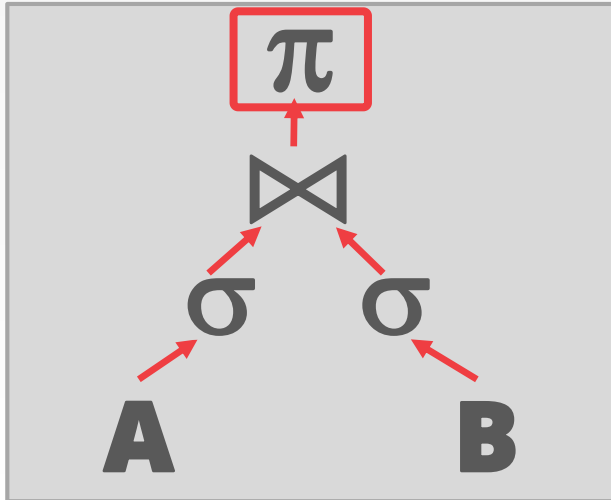
```

for  $r_1 \in$  outer:
  for  $r_2 \in$  inner:
    emit( $r_1 \bowtie r_2$ )
  
```

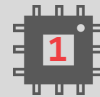
INTER-OPERATOR PARALLELISM

```

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


 π

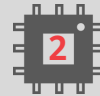
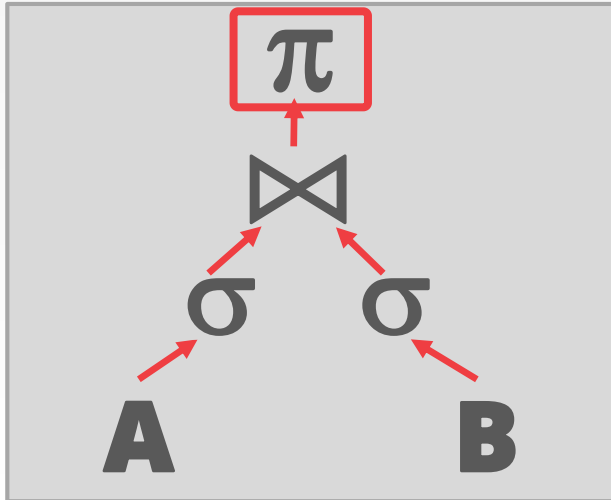
for $r \in$ incoming:
emit(πr)


 \bowtie

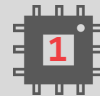
for $r_1 \in$ outer:
for $r_2 \in$ inner:
emit($r_1 \bowtie r_2$)

INTER-OPERATOR PARALLELISM

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


 π

for $r \in$ incoming:
emit(πr)


 \bowtie

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

OBSERVATION

Coming up with the right number of workers to use for a query plan depends on the number of CPU cores, the size of the data, and functionality of the operators.

WORKER ALLOCATION

Approach #1: One Worker per Core

- Each core is assigned one thread that is pinned to that core in the OS.
- See [sched_setaffinity](#)

Approach #2: Multiple Workers per Core

- Use a pool of workers per core (or per socket).
- Allows CPU cores to be fully utilized in case one worker at a core blocks.

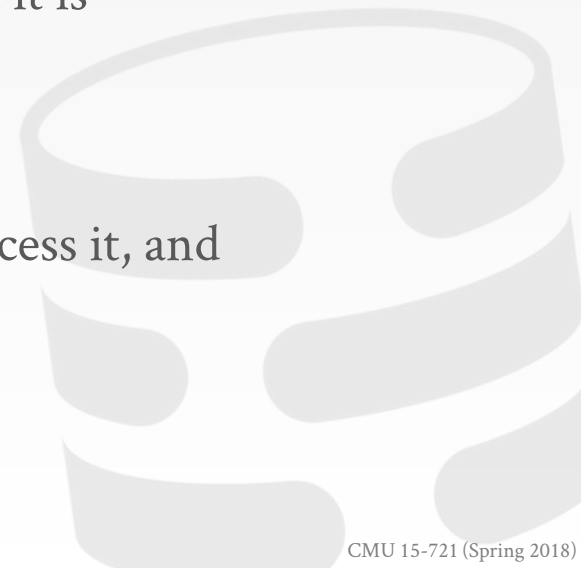
TASK ASSIGNMENT

Approach #1: Push

- A centralized dispatcher assigns tasks to workers and monitors their progress.
- When the worker notifies the dispatcher that it is finished, it is given a new task.

Approach #1: Pull

- Workers pull the next task from a queue, process it, and then return to get the next task.



OBSERVATION

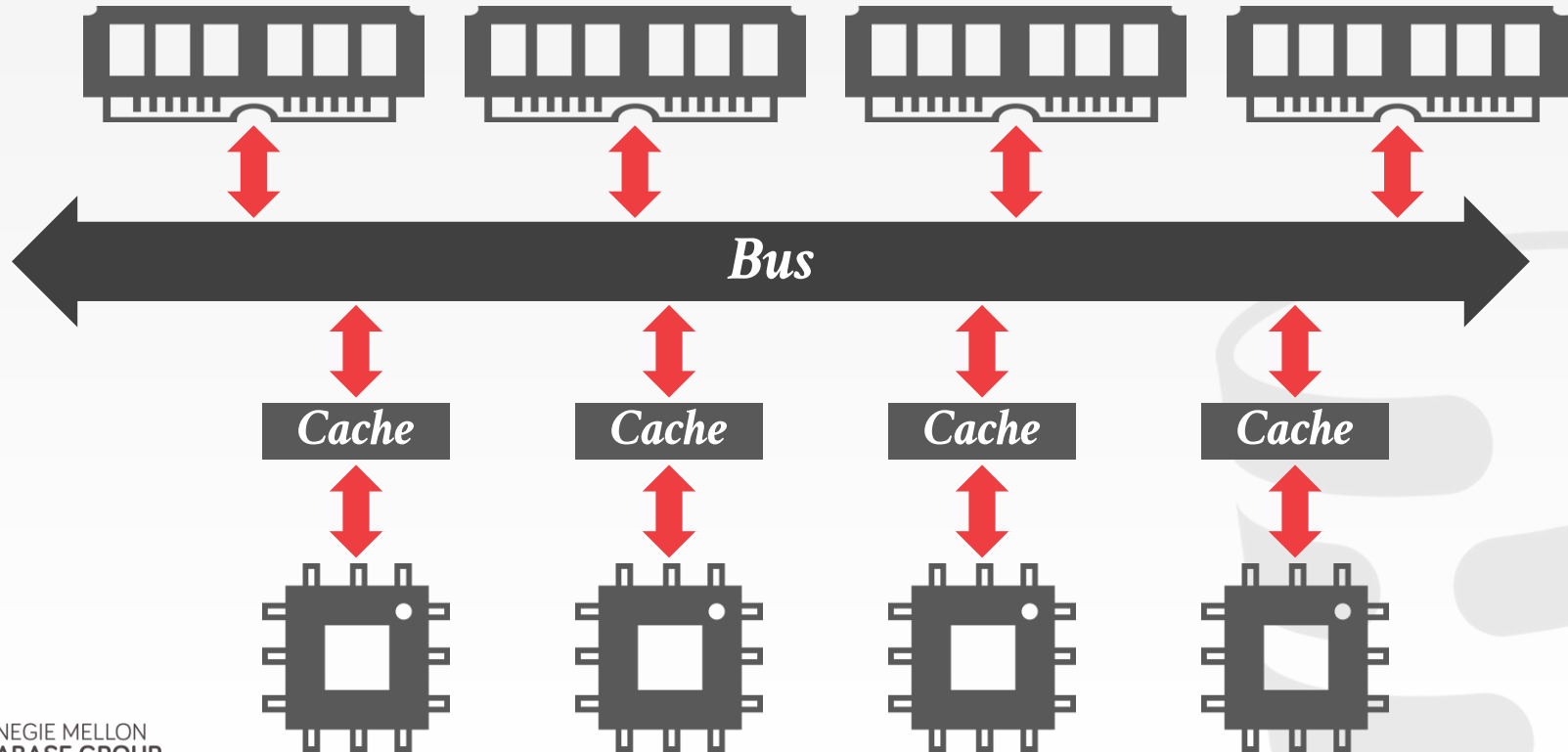
Regardless of what worker allocation or task assignment policy the DBMS uses, it's important that workers operate on local data.

The DBMS's scheduler has to be aware of its underlying hardware's memory layout.

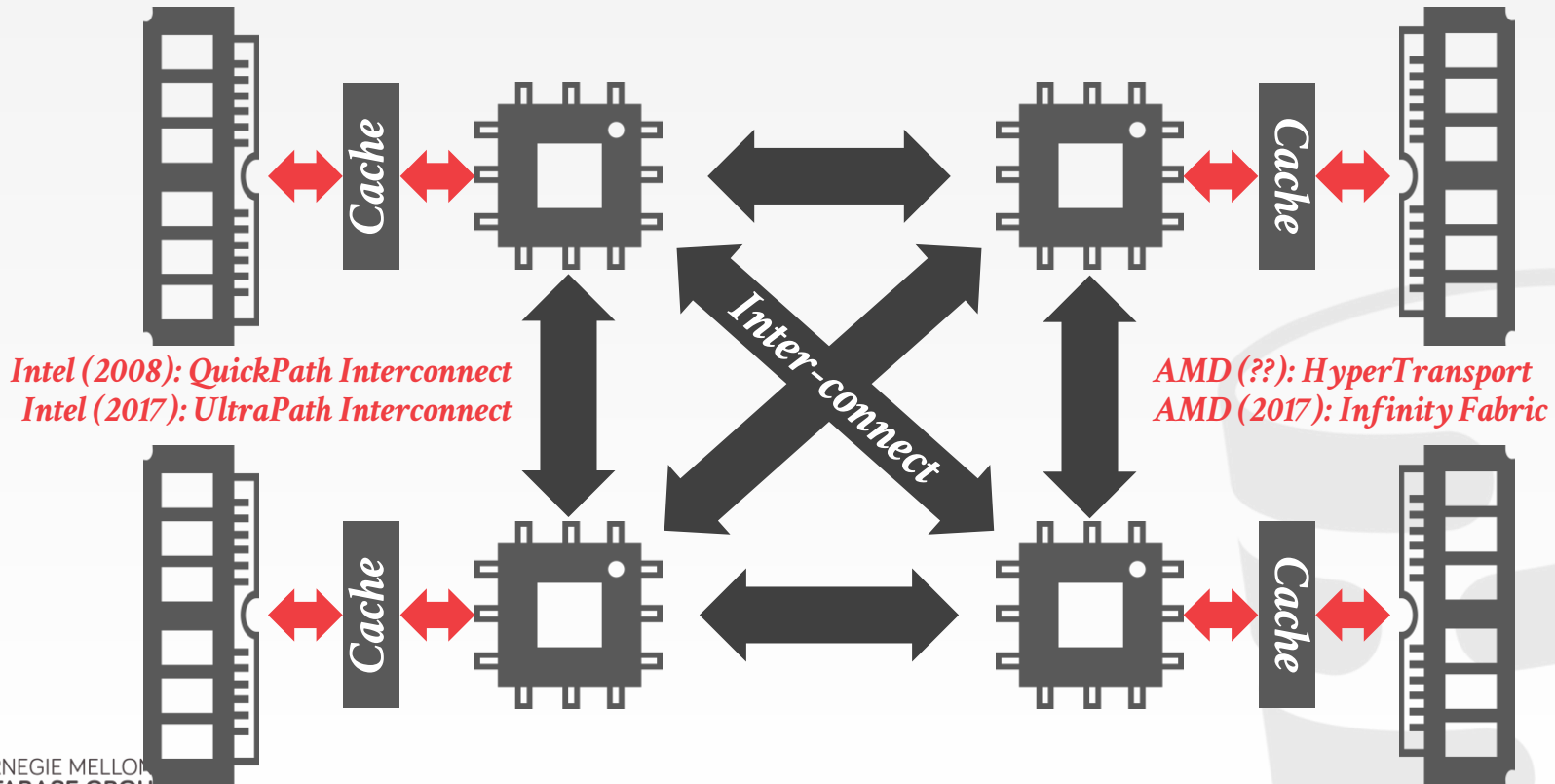
→ Uniform vs. Non-Uniform Memory Access



UNIFORM MEMORY ACCESS



NON-UNIFORM MEMORY ACCESS



DATA PLACEMENT

The DBMS can partition memory for a database and assign each partition to a CPU.

By controlling and tracking the location of partitions, it can schedule operators to execute on workers at the closest CPU core.

See Linux's [move_pages](#)



MEMORY ALLOCATION

What happens when the DBMS calls **malloc**?

→ Assume that the allocator doesn't already have a chunk of memory that it can give out.

Actually, almost nothing:

- The allocator will extend the process' data segment.
- But this new virtual memory is not immediately backed by physical memory.
- The OS only allocates physical memory when there is a page fault.

MEMORY ALLOCATION LOCATION

Now after a page fault, where does the OS allocate physical memory in a NUMA system?

Approach #1: Interleaving

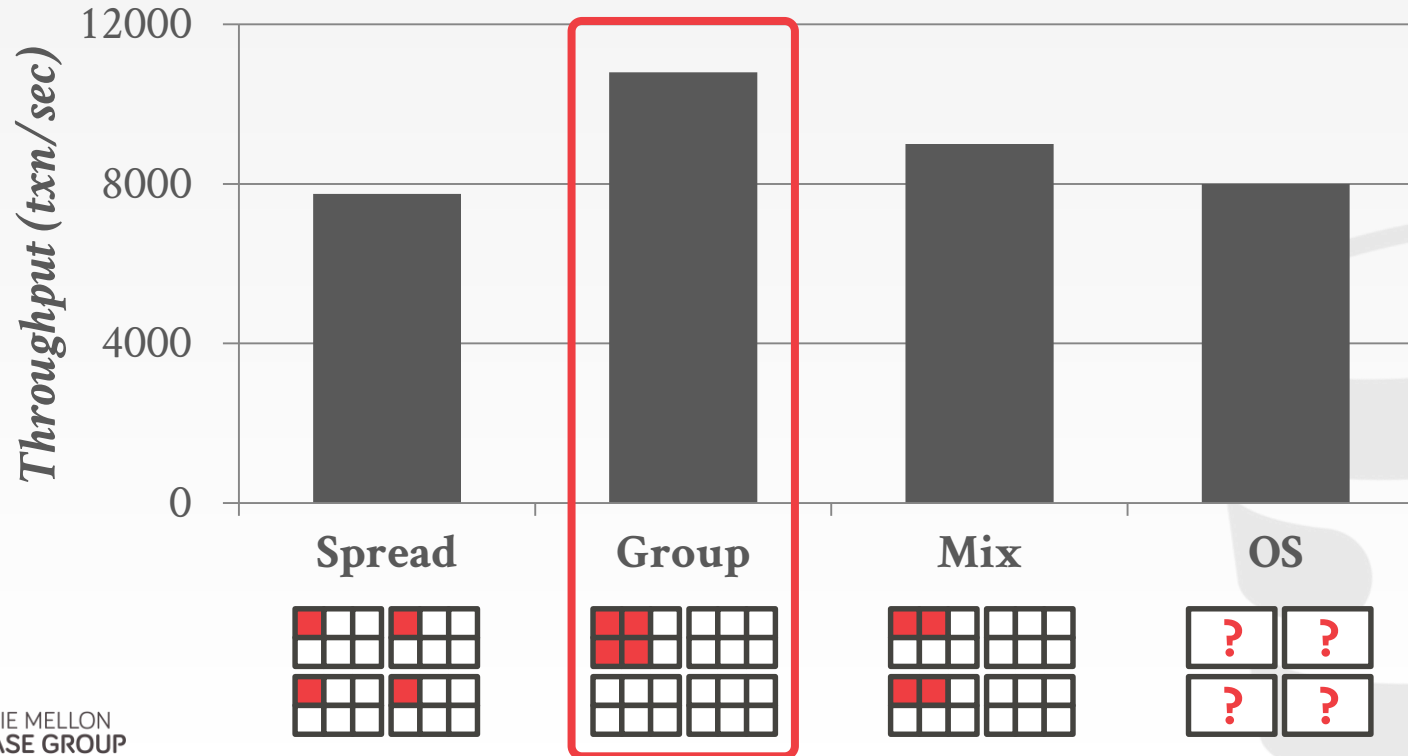
→ Distribute allocated memory uniformly across CPUs.

Approach #2: First-Touch

→ At the CPU of the thread that accessed the memory location that caused the page fault.

DATA PLACEMENT – OLTP

Workload: TPC-C Payment using 4 Workers
Processor: NUMA with 4 sockets (6 cores each)

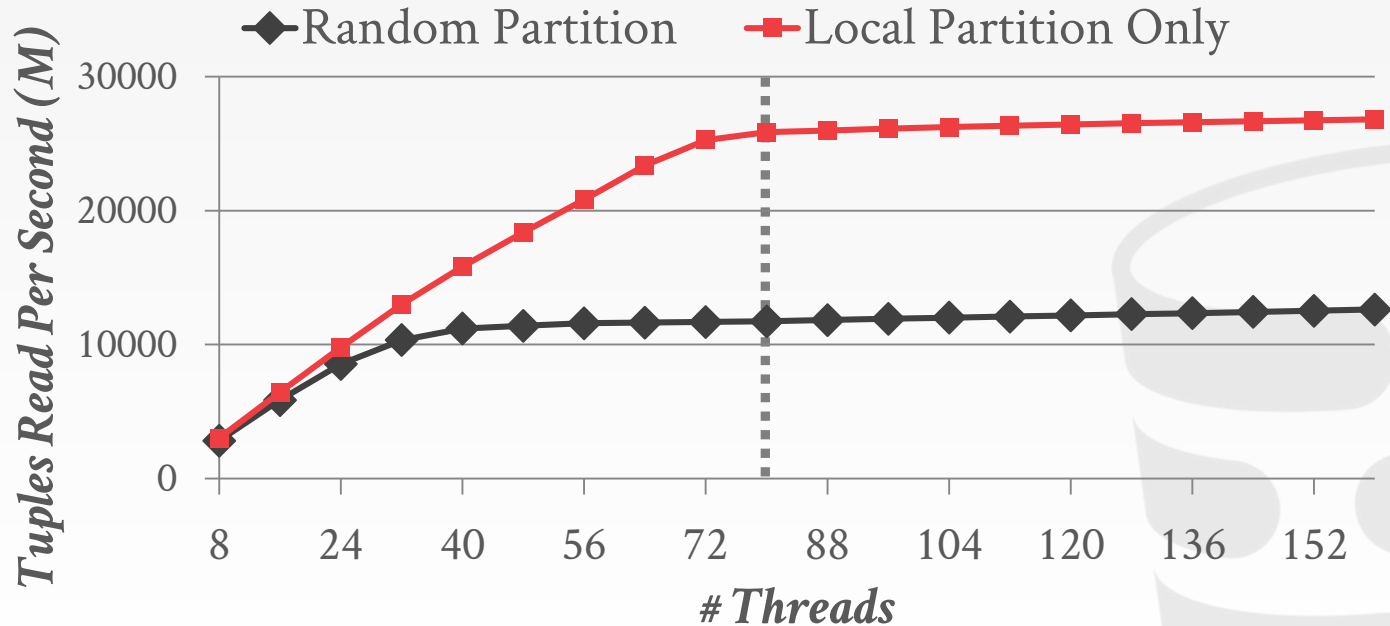


DATA PLACEMENT – OLAP

Database: 10 million tuples

Workload: Sequential Scan

Processor: 8 sockets, 10 cores per node (2x HT)



PARTITIONING VS. PLACEMENT

A **partitioning** scheme is used to split the database based on some policy.

- Round-robin
- Attribute Ranges
- Hashing
- Partial/Full Replication

A **placement** scheme then tells the DBMS where to put those partitions.

- Round-robin
- Interleave across cores



OBSERVATION

We have the following so far:

- Process Model
- Worker Allocation Model
- Task Assignment Model
- Data Placement Policy

But how do we decide how to create a set of tasks from a logical query plan?

- This is relatively easy for OLTP queries.
- Much harder for OLAP queries...



STATIC SCHEDULING

The DBMS decides how many threads to use to execute the query when it generates the plan.

It does **not** change while the query executes.

→ The easiest approach is to just use the same # of tasks as the # of cores.



MORSEL-DRIVEN SCHEDULING

Dynamic scheduling of tasks that operate over horizontal partitions called “morsels” that are distributed across cores.

- One worker per core
- Pull-based task assignment
- Round-robin data placement

Supports parallel, NUMA-aware operator implementations.



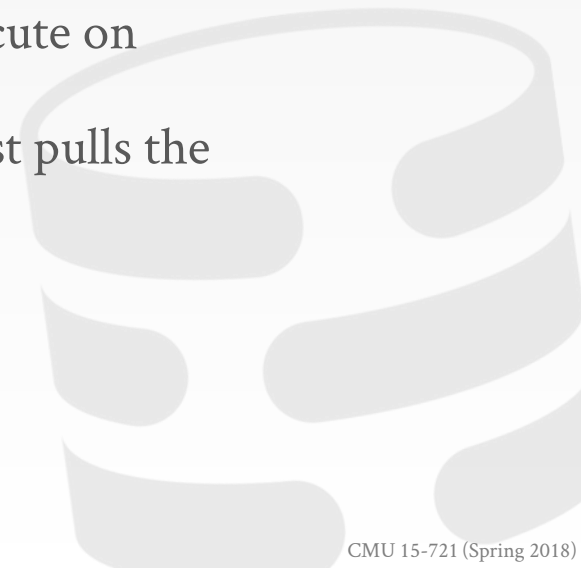
MORSEL-DRIVEN PARALLELISM: A NUMA-AWARE QUERY EVALUATION FRAMEWORK FOR THE MANY-CORE AGE
SIGMOD 2014

HYPER – ARCHITECTURE

No separate dispatcher thread.

The threads perform cooperative scheduling for each query plan using a single task queue.

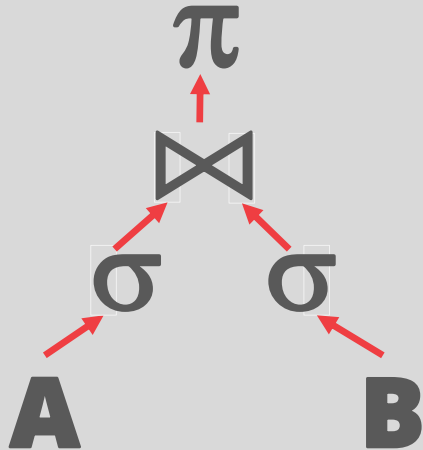
- Each worker tries to select tasks that will execute on morsels that are local to it.
- If there are no local tasks, then the worker just pulls the next task from the global work queue.



HYPER – DATA PARTITIONING

```

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



Data Table

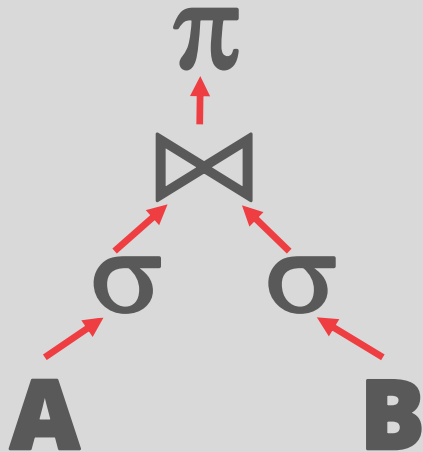
Morsels

	id	a1	a2	a3	
A ₁					}
A ₂					
A ₃					
					1
					2
					3

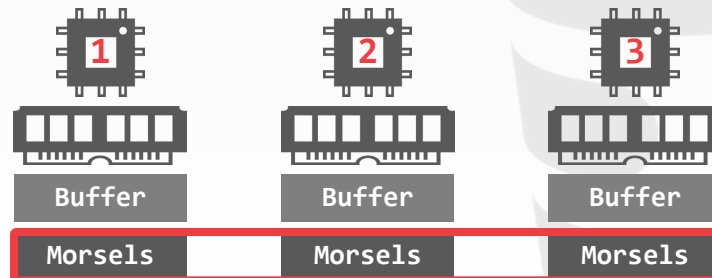
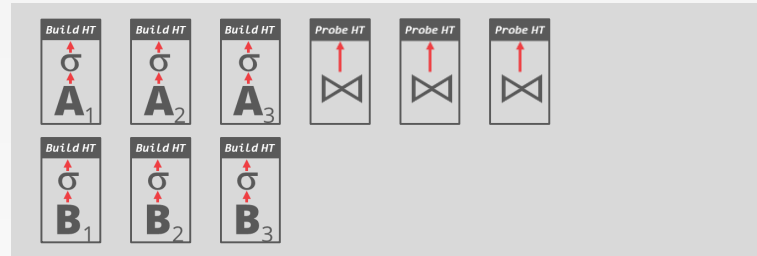
The table is divided into three horizontal sections labeled A₁, A₂, and A₃. Each section has four columns: id, a1, a2, and a3. To the right of the table, three red curly braces group the rows into three sections, each associated with a server icon labeled 1, 2, and 3 respectively.

HYPER – EXECUTION EXAMPLE

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



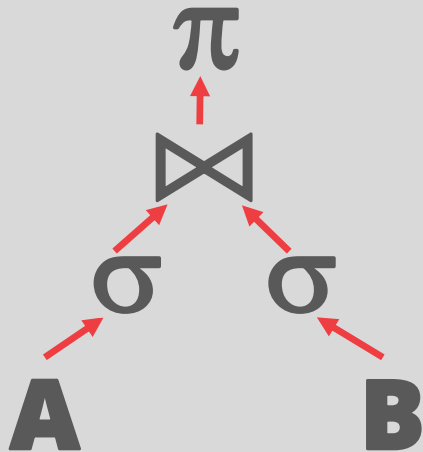
Global Task Queue



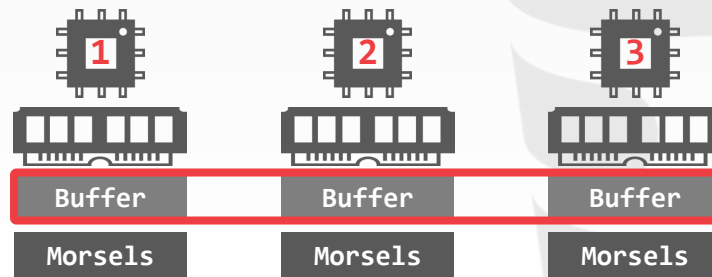
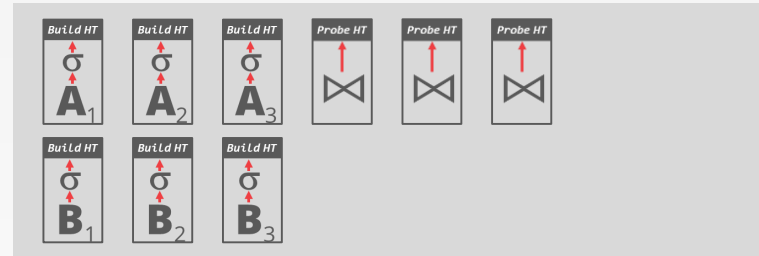
HYPER – EXECUTION EXAMPLE

```

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

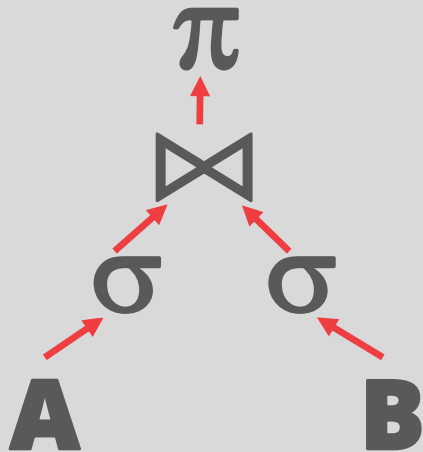


Global Task Queue

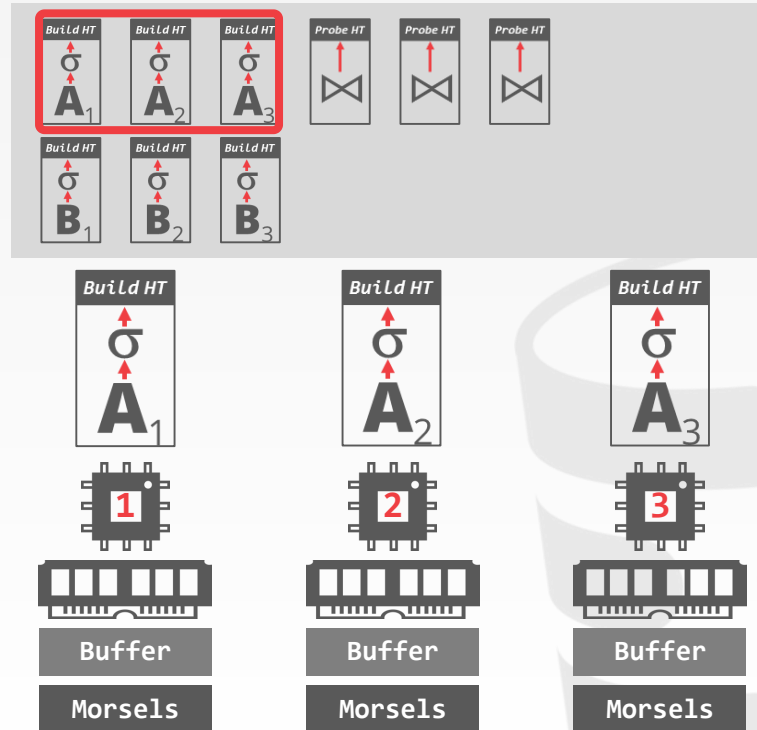


HYPER – EXECUTION EXAMPLE

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



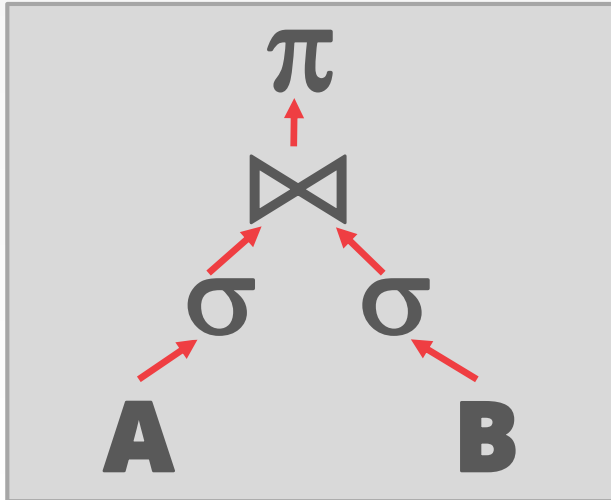
Global Task Queue



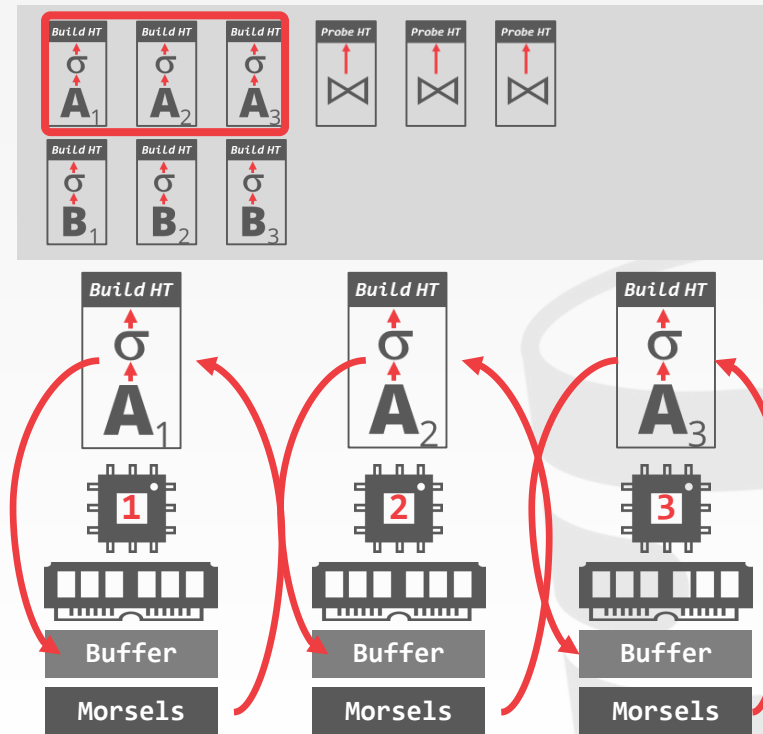
HYPER – EXECUTION EXAMPLE

```

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

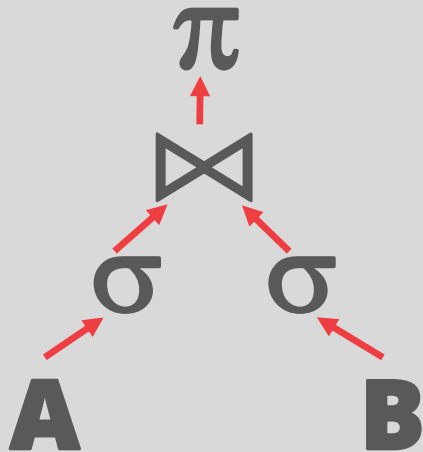


Global Task Queue

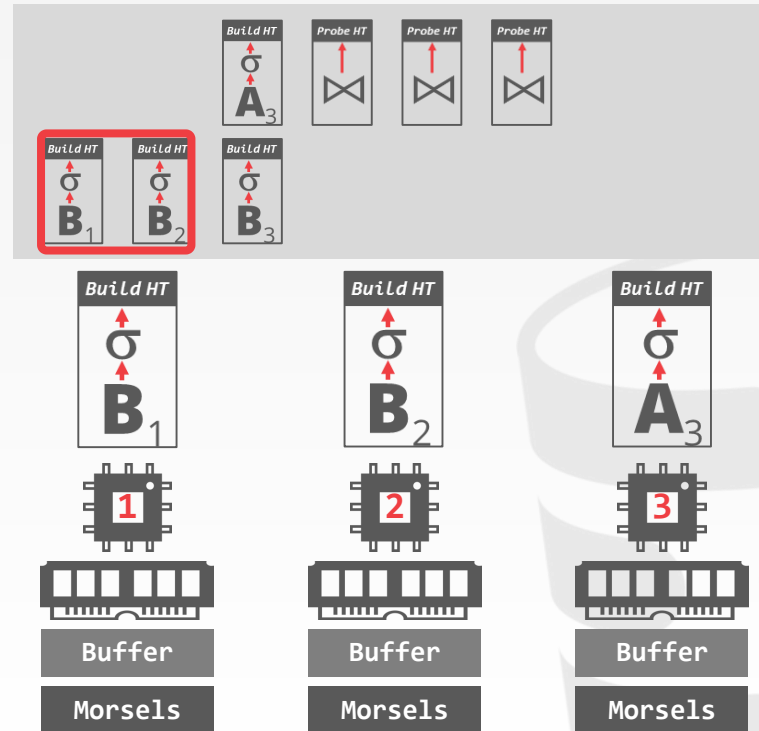


HYPER – EXECUTION EXAMPLE

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



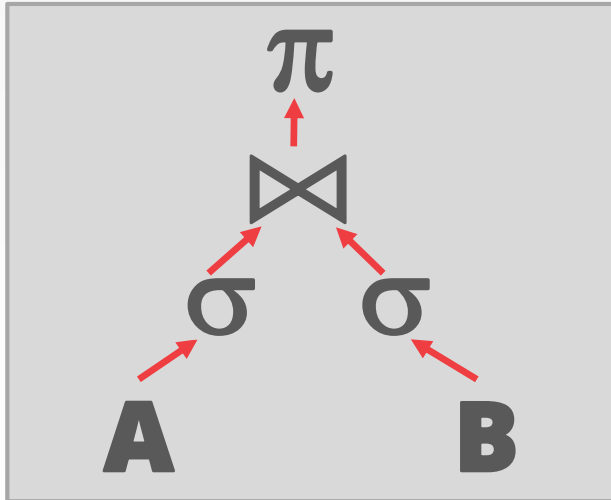
Global Task Queue



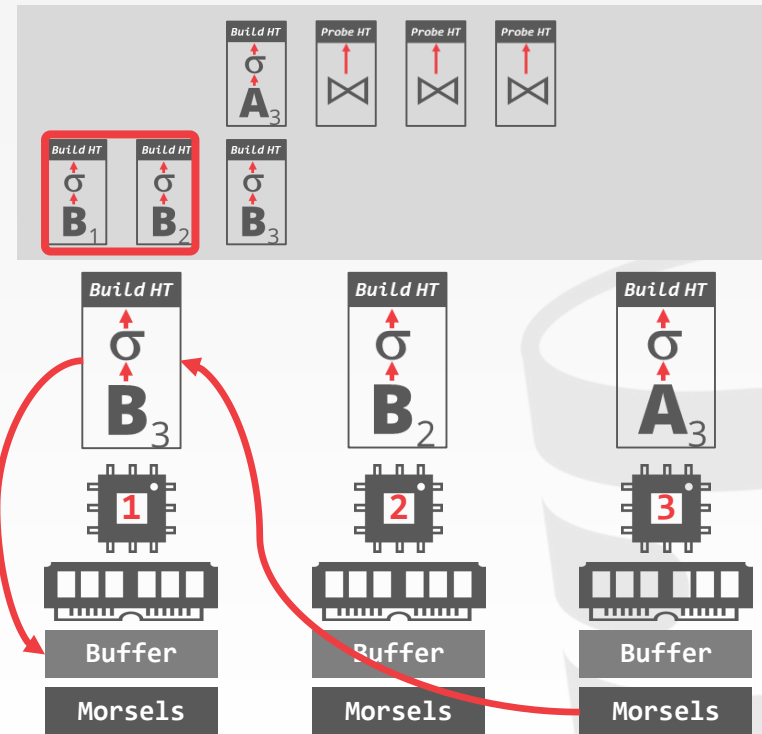
HYPER – EXECUTION EXAMPLE

```

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



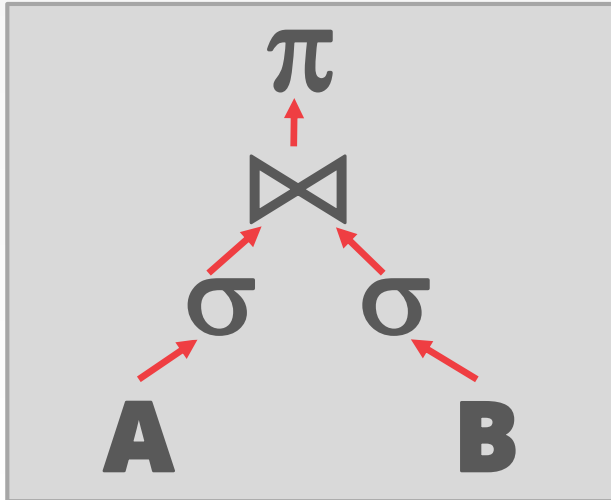
Global Task Queue



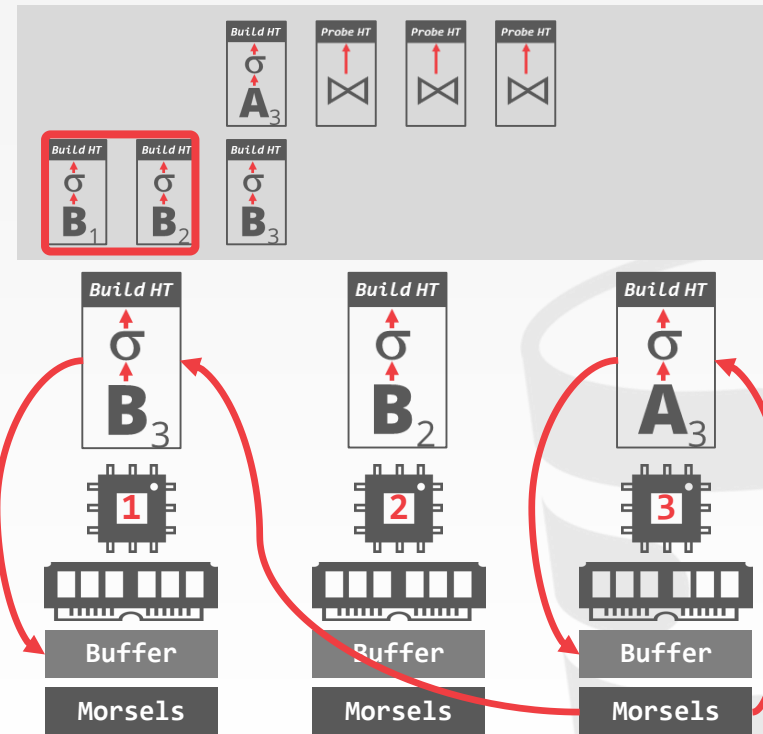
HYPER – EXECUTION EXAMPLE

```

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



Global Task Queue



MORSEL-DRIVEN SCHEDULING

Because there is only one worker per core, they have to use work stealing because otherwise threads could sit idle waiting for stragglers.

Uses a lock-free hash table to maintain the global work queues.

→ We will discuss hash tables next class...



SAP HANA – NUMA-AWARE SCHEDULER

Pull-based scheduling with multiple worker threads that are organized into groups (pools).

- Each CPU can have multiple groups.
- Each group has a soft and hard priority queue.

Uses a separate “watchdog” thread to check whether groups are saturated and can reassign tasks dynamically.



SCALING UP CONCURRENT MAIN-MEMORY COLUMN-STORE SCANS:
TOWARDS ADAPTIVE NUMA-AWARE DATA AND TASK PLACEMENT
VLDB 2015

SAP HANA – THREAD GROUPS

Each thread group has a soft and hard priority task queues.

→ Threads are allowed to steal tasks from other groups' soft queues.

Four different pools of thread per group:

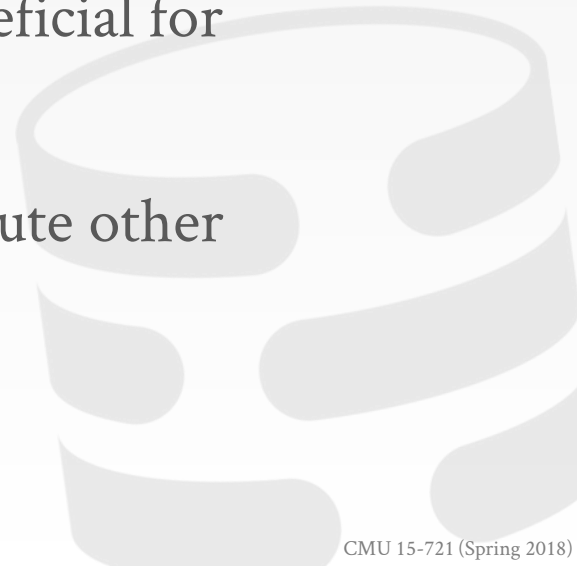
- **Working:** Actively executing a task.
- **Inactive:** Blocked inside of the kernel due to a latch.
- **Free:** Sleeps for a little, wake up to see whether there is a new task to execute.
- **Parked:** Like free but doesn't wake up on its own.

SAP HANA – NUMA-AWARE SCHEDULER

Can dynamically adjust thread pinning based on whether a task is CPU or memory bound.

Found that work stealing was not as beneficial for systems with a larger number of sockets.

Using thread groups allows cores to execute other tasks instead of just only queries.



PARTING THOUGHTS

A DBMS is a beautiful, strong-willed independent piece of software.

But it has to make sure that it uses its underlying hardware correctly.

- Data location is an important aspect of this.
- Tracking memory location in a single-node DBMS is the same as tracking shards in a distributed DBMS

Don't let the OS ruin your life.

NEXT CLASS

Hash Tables!

Hash Functions!

Hash Joins!

