

Carnegie Mellon University ADVANCED DATABASE SYSTEMS

Ouery Scheduling

Andy Pavlo // 15-721 // Spring 2023

LAST CLASS

- We discussed query processing models.
- \rightarrow Vectorized model is best for OLAP.
- \rightarrow Top-to-bottom (pull) approach is probably better.



QUERY EXECUTION

A query plan is a DAG of **operators**.

An <u>operator instance</u> is an invocation of an operator on a unique segment of data.

A <u>task</u> is a sequence of one or more operator instances.

A <u>task set</u> is the collection of executable tasks for a logical pipeline.



ECMU·DB 15-721 (Spring 2023

QUERY EXECUTION

A query plan is a DAG of **operators**.

An <u>operator instance</u> is an invocation of an operator on a unique segment of data.

A <u>task</u> is a sequence of one or more operator instances.

A <u>task set</u> is the collection of executable tasks for a logical pipeline.

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



SCHEDULING

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

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

The DBMS *always* knows more than the OS.

SCHEDULING GOALS

Goal #1: Throughput → Maximize the # of completed queries.

Goal #2: Fairness

 \rightarrow Ensure that no query is starved for resources

Goal #3: Query Responsiveness

 \rightarrow Minimize tail latencies (especially for short queries)

Goal #4: Low Overhead

→ Workers should spend most of their time executing tasks not figuring out what task to run next.

ECMU·DB 15-721 (Spring 2023

TODAY'S AGENDA

Worker Allocation Data Placement Scheduler Implementations

PROCESS MODEL

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

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

We will assume that the DBMS is multi-threaded.



WORKER ALLOCATION

Approach #1: One Worker per Core

- \rightarrow Each core is assigned one thread that is pinned to that core in the OS.
- → See see setaffinity

Approach #2: Multiple Workers per Core

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

TASK ASSIGNMENT

Approach #1: Push

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

Approach #1: Pull

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

10

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 must be aware of its hardware memory layout.

 \rightarrow Uniform vs. Non-Uniform Memory Access

UNIFORM MEMORY ACCESS



ECMU·DB 15-721 (Spring 2023)

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 and numactl



MEMORY ALLOCATION

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

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

Almost nothing:

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

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



MEMORY ALLOCATION LOCATION

Approach #1: Interleaving

 \rightarrow Distribute allocated memory uniformly across CPUs.

Approach #2: First-Touch

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

The OS can try to move memory to another NUMA region from observed access patterns.



DATA PLACEMENT - OLAP

Sequential Scan on 10m tuples Processor: 8 sockets, 10 cores per node (2x HT)



Source: <u>Haibin Lin</u> SCMU-DB 15-721 (Spring 2023)

PARTITIONING VS. PLACEMENT

A **<u>partitioning</u>** scheme is used to split the database based on some policy.

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

A <u>placement</u> scheme then tells the DBMS where to put those partitions. \rightarrow Round-robin \rightarrow Interleave across cores



OBSERVATION

- We have the following so far:
- \rightarrow Task Assignment Model
- \rightarrow Data Placement Policy

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

- \rightarrow This is relatively easy for OLTP queries.
- \rightarrow 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.
- → Can still assign tasks to threads based on data location to maximize local data processing.

MORSEL-DRIVEN SCHEDULING

Dynamic scheduling of tasks that operate over horizontal partitions called "morsels" distributed across cores.

- \rightarrow One worker per core.
- \rightarrow One morsel per task.
- \rightarrow Pull-based task assignment.
- \rightarrow Round-robin data placement.

Supports parallel, NUMA-aware operator implementations.

HYPER - ARCHITECTURE

No separate dispatcher thread.

The workers 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.

 \rightarrow If there are no local tasks, then the worker just pulls the next task from the global work queue.

HYPER - DATA PARTITIONING



Data Table





ECMU·DB 15-721 (Spring 2023)



Global Task Queue





15-721 (Spring 2023)





15-721 (Spring 2023)





*----

Global Task Queue







SECMU-DB 15-721 (Spring 2023)



15-721 (Spring 2023)





15-721 (Spring 2023)





15-721 (Spring 2023)



MORSEL-DRIVEN SCHEDULING

Because there is only one worker per core and one morsel per task, HyPer must use work stealing because otherwise threads could sit idle waiting for stragglers.

The DBMS uses a lock-free hash table to maintain the global work queues.

26

OBSERVATION

Tasks can have different execution costs per tuple. \rightarrow Example: Simple Selection vs. String Matching

HyPer also has no notion of execution priorities.

- \rightarrow All query tasks are executed with the same
- → Short-running queries get blocked behind long-running queries.

UMBRA - MORSEL SCHEDULING 2.0

Tasks are not created statically at runtime. Each task may contain multiple morsels. Modern implementation of stride scheduling. Priority decay.





UMBRA - MORSEL SCHEDULING 2.0

Tasks are 1 Each task 1 Modern ir Priority de

FIREBOLT

How We Build Firebolt

Vaccination Database Tech Talks - Second Dose benjamin.wagner@firebolt.io





ECMU·DB 15-721 (Spring 2023)

Each worker maintains its own thread-local meta-data about the available tasks to execute.

- → Active Slots: Which entries in the global slot array have active task sets available.
- → **Change Mask:** Indicates when a new task set is added to the global slot array.
- → **Return Mask:** Indicates when a worker completes a task set.

Workers perform CaS updates to TLS meta-data to broadcast changes.











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

- \rightarrow Each CPU can have multiple groups.
- \rightarrow 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.



SAP HANA - THREAD GROUPS

DBMS maintains <u>soft</u> and <u>hard</u> priority task queues for each thread group.

 \rightarrow Threads can steal tasks from other groups' soft queues.

Four different pools of thread per group:

- → **Working**: Actively executing a task.
- \rightarrow **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.
- \rightarrow **Parked**: Waiting for a task (like a free thread) but blocked in the kernel until the watchdog thread wakes it up.

Dynamically adjust thread pinning based on whether a task is CPU or memory bound. \rightarrow Allow more cross-region stealing if DBMS is CPU-bound.

SAP found that work stealing was not as beneficial for systems with a larger number of sockets. \rightarrow HyPer (2-4 sockets) vs. HANA (64 sockets)

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













```
Tasks
         Build HT
```







15-721 (Spring 2023)

SCMUDB



15-721 (Spring 2023)

Tasks







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

- \rightarrow 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.



Startups SQLOS is a use Videos runs inside of Audio Newsletters hardware reso Extra Crunch Advertise \rightarrow Determines v Events \rightarrow Also manage More like logical da

Transportation Apple Non-preemp Tesla Security instrumented



SECMU-DB 15-721 (Spring 2023)

MICROSOFT SQL SERVER 2012 INTERNALS PEARSON 2013

How Microsoft brought SQL Server to Linux

Frederic Lardinois @frederic! / 12:00 pm EDT * July 17, 2017

Join Extra Crunch

Login

Search Q

Back in 2016, when Microsoft
announced that SQL Server would soon run on Linux, the news came as a major surprise to users and pundits alike. Over the course of the last year, Microsoft's support for Linux (and open source in general), has come into clearer focus and the company's mission now seems to be all about bringing its tools to wherever its users are. The company today launched the first release candidate of SQL Server 2017, which will be the first version to run on Windows, Linux and in Docker containers. The Docker container alone has already seen more than 1 million pulls, so there can be no doubt that there is a lot of interest in this new version. And while there are plenty of new features and speed improvements in this new version, the fact that SQL Server 2017 supports Linux remains one of the most interesting aspects of this release.

Ahead of today's announcement, I talked to Rohan Kumar, the general manager of Microsoft's Database Systems group, to get a bit more info about the history of this project and how his team managed to bring an extremely complex piece of software like SQL Server to Linux. Kumar, who has been at Microsoft for more than 18 years, noted that his team noticed many enterprises were starting to use SQL Server for their mission-critical workloads. But at the same time, they were also working in mixed environments that included both Windows Server and Linux. For many of these businesses, not being able to run their database of choice on Linux became a friction point.

"Talking to enterprises, it became clear that doing this was necessary," Kumar said. "We were forcing customers to use Windows as their platform of choice." In another incarnation of Microsoft, that probably would've been seen as something positive, but the company's strategy today is quite different.

X

Comment

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.

Other Examples: \rightarrow ScyllaDB \rightarrow FaunaDB \rightarrow CoroBase

15-721 (Spring 2023)

SELECT * FROM R WHERE R.val = ?

Approximate Plan

```
for t in R:
    if eval(predicate, tuple, params):
        emit(tuple)
```

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.

Other Examples: \rightarrow ScyllaDB \rightarrow FaunaDB \rightarrow CoroBase

15-721 (Spring 2023)

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)
```

SQLOS quantum i scheduler cannot e

DBMS developers explicit yield calls locations in the so

Other Examples: \rightarrow ScyllaDB \rightarrow FaunaDB \rightarrow CoroBase



ECMU-DB 15-721 (Spring 2023)

OBSERVATION

If requests arrive at the DBMS faster than it can execute them, then the system becomes overloaded.

The OS cannot help us here because it does not know what threads are doing:

- \rightarrow CPU Bound: Do nothing
- \rightarrow Memory Bound: OOM

Easiest DBMS Solution: Crash

38

FLOW CONTROL

Approach #1: Admission Control

 \rightarrow Abort new requests when the system believes that it will not have enough resources to execute that request.

Approach #2: Throttling

- \rightarrow Delay the responses to clients to increase the amount of time between requests.
- \rightarrow This assumes a synchronous submission scheme.

PARTING THOUGHTS

We ignored disk I/O scheduling.

A DBMS is a beautiful, strong-willed independent software. But it must use hardware correctly.

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

Do <u>not</u> let the OS ruin your life.

NEXT CLASS

Vectorized Query Execution