

RIVERSIDE

SIGMOD 2025 Programming Contest – Team JobSeeking

Xiangyun Ding



Overview

• Our team

- Team name: **JobSeeking**, from University of California, Riverside
- Final result: **5'th place**!
- We presents a high-performance, multi-threaded query engine designed for the SIGMOD 2025 Programming Contest.

• Key techniques

• A **work-stealing scheduler** for shared-memory multi-core parallelism

Input/Output

- Input Data Format
 - ColumnarTable -> Column -> Page
 - A tree structure for join order

Letong Wang

- Parallel Column Read
 - **Parallel prefix sum**: we first perform a parallel scan over pages to sum up the number of rows, which gives each thread the precise starting offset in the array.
 - I/O cache: a column **cache**, keyed by a hash of the column's page content, stores pointers to already-read columns. This cache is crucial for performance.
- Highly parallelized I/O modules for rapid data ingestion and emission
- Skew-aware parallel hash join algorithm based on parallel integer-sort
- Custom memory allocators

Two-Phase Parallel Writing

- Counting Phase: Determine the number of pages required for each partition of the output.
- Writing Phase: Allocate memory for pages, then format and write data partitions into the appropriate pages in parallel.

Skew-Aware Parallel Hash Join

- Building the hash table: heavy/light Key partitioning via parallel integer-sort [1]
 - Heavy key identification: by sampling the keys and counting their frequencies, we identify heavy keys that appear more often than a dynamic threshold.
 - For heavy keys, we use a hash table to store them with their indices.
 - Light keys are distributed uniformly based on the high bits in parallel. We build one hash table for each light bucket.
 - For queries, we first look up the hash table for heavy keys. If the key is not found, we then look up the corresponding lightbucket hash table.
 - The hash table is implemented by linear-probing.

• Runtime optimizations

- Dynamic **join order selection**: we greedily choose the pair of tables whose product of row counts is minimized.
- Dynamic build: the smaller table is always chosen as the build side to construct the hash table.
- Hash table caching: completed hash tables for large tables can be **cached** and reused in subsequent joins on the same key, avoiding redundant build costs.

```
Step 1: Take samples (boxed), detect heavy keys,
assign bucket ids
```

Samples: 4 *3 2 *1 6 *2 9 *2 7 *1

4 light buckets: (highest two bits) 00, 01, 10, 11 3 heavy buckets: for 4, 6, 9

Step 2: Distribute records to corresponding buckets



Multi-Thread Parallelism

• Work-Steal Scheduler [2]

• Each worker maintains a **deque** for tasks

- Parallel primitives:
 - parallel-for
 - Simple application: parallelize any for-loop



• Each worker supports:

- **spaws**: pushes a task to the bottom of deque
- o get_task: if the deque is <u>not empty</u>, it pops from

the <u>bottom</u>

- if the deque is <u>empty</u>, it **steals** a task by
 - popping from the <u>top</u> of another deque
- Adaptive granularity: adaptively adjust task size for high-performance on general machines
- scan (prefix sum), interger sort, ...
- Parallel resources allocator
 - Allocate **thread pool** with size #cores/#cores*2
 - Allocate **memory pool** for intermediate results of parallel functions

References

[1] Dong, Xiaojun, Laxman Dhulipala, Yan Gu, and Yihan Sun. "Parallel integer sort: Theory and practice." In Proceedings of the 29th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming, pp. 301-315. 2024. [2] Blelloch, Guy E., Daniel Anderson, and Laxman Dhulipala. "ParlayLib-a toolkit for parallel algorithms on shared-memory multicore machines." In Proceedings of the 32nd ACM Symposium on Parallelism in Algorithms and Architectures, pp. 507-509. 2020.