

Hasso Plattner Institute

Chair for Data Engineering Systems



Proposal Master Thesis

Hardware-Conscious SIMD-Accelerated Sort-Merge Joins in Multi Core In-Memory Database Systems

Finn Schöllkopf

Time frame: October 2024 - March 2024

Supervisor

Prof. Dr. Tilmann Rabl

Advisor

Florian Schmeller

1 Motivation

Traditional database systems have historically been designed for systems and architectures where I/O dominates performance. However, modern processors with multi-core architectures, advanced instruction sets, and other hardware accelerants like vector operations (SIMD) have significantly altered this landscape. Today’s in-memory database systems are no longer I/O bound and, therefore, need high intra-operator parallelism to fully utilize the multi-core architecture. To achieve maximum performance, cache locality, NUMA awareness, and using SIMD instructions for higher data parallelism should be considered. The join operator is a fundamental component of every database system. In recent years, the difference in performance between the sort-merge and radix-hash join has been the subject of ongoing debate. Kim et al. [10] projected that Sort-Merge Join would outperform hash-based alternatives with a factor of 1.35 – 1.65 with 512 bits, particularly for lower tuple counts with 256-bit SIMD registers. Albutiu et al. [1] reinforced this claim with recent results reporting that their NUMA-aware implementation of sort-merge join is superior to that of hash joins (without leveraging SIMD). Balakesen et al. [2] experimentally show contradicting results by implementing optimized versions for sort-merge and radix-hash join, showing that their implementation of radix-hash join is still superior. They use AVX2 in their implementation, allowing further work to explore wider SIMD registers (e.g., AVX512). With only very few public implementations, which are often experimental, there is still great value in an open-source, state-of-the-art sort-merge join implementation optimized for different architectures. Such an implementation can be a baseline for further research to measure improvements.

2 Goal of Thesis

This thesis aims to efficiently implement the sort-merge join algorithm, explicitly optimized for specific architectures and hardware components. While multiple papers exist about modern implementation approaches for sort-merge joins in in-memory database systems and SIMD sorting generally, only some have public implementations¹. Most SIMD sorting algorithms presented in the literature are not directly applicable to join operations as they usually use sorting keys of only 32 bits. We must track the row ID (rid) corresponding to the search key for a join. The current implementations of sort-merge join in literature use SSE and AVX2 intrinsics, but to our knowledge, there has yet to be an implementation using AVX512. Therefore, in the scope of this thesis, we want to integrate support for modern AVX512 sorting algorithms ([13], [14]) next to SSE and AVX2 into a complete sort-merge join operator. It would also be of value to see how existing and new approaches

¹Implementation of [2] published at <https://archive-systems.ethz.ch/node/334>

transfer to other CPU architectures like Arm with its Scalable Vector Extension (SVE). While some public implementations exist for modern and optimized sort-merge join, they have usually isolated implementations with a strong focus on the sorting step using randomly chosen input data, often already in the required data format. Also, they often skip the lookup of matching rows and the construction of the joined table. Hence, in this thesis, we want to integrate our implementation of the sort-merge join into Hyrise [8], a research in-memory database. Hyrise contains both a radix-based Hash-Join and sort-merge join. The sort-merge join uses radix cluster sorting, which uses pattern-defeating quicksort (boost) but no explicit SIMD instructions. It fundamentally differs from the modern approaches in the literature. These differences allow us to test our implementation against the existing sort-merge and hash-based join. Complete integration into an in-memory database allows us to benchmark sorting throughput in tuples per second and see how our sort-merge join operator compares to other join implementations in decision support benchmarks like TCP-H or TCP-DS with more realistic data. Benchmarking should also include measuring all important stages, like the initial data construction in the format of (key, rid) from the input relations, sorting and the final construction of the join table. We want to test our implementation on different architectures and hardware, taking advantage of differences in core count, cache size, SIMD registers with different widths, NUMA regions, and other hardware-specific properties.

3 Approach

The sort-merge join involves sorting both input relations. It is the most crucial and time-consuming part of the sort-merge join operation. Therefore, optimization efforts should primarily focus on this step, as it largely determines the runtime. Due to modern multi-core architectures, sorting should intensively utilize thread-level parallelism by multithreading. With the recent architectural trends of wider register widths for SIMD, sorting should also heavily use SIMD instructions to exploit SIMD data parallelism. In a multi-core context, merge sort is often preferred over quicksort, as the parallelization of the divide-and-conquer approach is straightforward, and it has other advantages over quicksort such as more predictable and cache-friendly memory access patterns and better load balancing through equal-sized partitioning. Sorting through SIMD registers can be achieved through sorting networks [4]. The sorting network compares elements in parallel in each step using SIMD min/max operations. A final transposition is needed, which requires additional SIMD shuffle instructions to complete the sorting. We can build sorting kernels for various input sizes² depending on the data type and register size. Merging can also benefit from SIMD acceleration. There are two standard merging networks: bitonic merge networks and odd-even merge networks ([4], SIMD accelerated [9]). Both

²https://bertdobbelaere.github.io/sorting_networks.html

scale poorly for bigger input sizes, with odd-even networks requiring slightly fewer comparisons but instead involving data movement and element masking. Therefore, we can use SIMD-accelerated merging networks as a kernel for small input sizes, e.g., by sequentially pulling already sorted data into SIMD registers and calling the merge kernel, which writes to the output and then fetches new data. We can merge different subparts of the data in different threads as long as we have enough sorted sublists. In the later round of the merge tree, with only a few sorted sublists remaining, it becomes increasingly more challenging to parallelize efficiently. However, even at this point, we can parallelize. Parallelization is made possible through the Merge Path [12]. This conceptual path allows us to parallelize a two-way merge by splitting it into non-overlapping segments that form disjoint sets of elements. We can then sequentially merge these segments in parallel. The sequential merging can again benefit from SIMD acceleration [13]. In the later stages, out-of-cache merging becomes necessary, quickly resulting in the memory bandwidth becoming the bottleneck of even a single-threaded merge routine. Therefore, multi-way merging [2], which consists of multiple two-way merge units (managed as tasks) connected via FIFO queues, is introduced. Only the leaves of the merge tree load data from memory. Blocking and task switching ensures that the combined FIFO queues fit into the CPU cache. This way, memory bandwidth can be reduced with a slight CPU overhead. Other hardware properties, such as NUMA-aware partitioning, can further assist in efficient sorting.

Before we can sort our input relations, the tuples need to be translated into a SIMD sortable format. Usually, a 64-bit pair (key, rid) (key & rid both 32-bit) is assumed. We, therefore, support a maximum relation size of 2^{32} . With most value types greater than 32 bits, we need to compress the values of the join columns to 32 bits. Methods like key-prefix [11] and XOR- and shift-based hash functions [6] have been used to represent keys using 32 bits.

After sorting both input relations, a final loop over both input relations suffices to find all join candidates. The sorted data is of the form (key, rid). Hence, we can use the row ID (rid) to find the respective tuples. As mentioned before some compression was used to generate the 32-bit search key from the join column value. Therefore, we might require additional validation and filtering in the final merge step.

4 Related Work

Many papers describe how sorting can be done efficiently in modern multi-core architectures for SIMD-accelerated sorting. Chhugani et al. [7] describe the concepts needed for efficient SIMD (SSE) sorting for both single—and multi-core execution, including sorting networks, bitonic—and odd-evenmerge networks, and how to deal with memory bandwidth limitations for large problem sizes through multiway merging. There are other ideas like MergePath [12] for merging only a few very large

sublists in parallel and SIMD accelerated. Kim et al. [10] implemented a sort-merge join using SSE intrinsics using these same concepts, projecting performance for wider SIMD widths that would outperform hash joins. Albutiu et al. [1] present MPSM, a sort-merge join implementation designed for modern multi-core and multi-socket NUMA processors using their custom sorting routine without SIMD. They did an experimental evaluation on a 32-core (4 socket) system, concluding that their sort-merge join implementation is faster than the respective hash join implementation of Blanas et al. [5]. Recent studies show that parallel radix-hash join has the best overall performance [3]. Therefore, Balkesen et al. [2] experimentally studied the performance of main-memory, parallel, multi-core join, and NUMA-aware algorithms, focusing on sort-merge and radix-hash join. They claim to provide the fastest in-memory join processing algorithms using sorting and hashing, and that sort-merge join gets more comparable in performance to radix-hash join with very large input sizes. Still, they conclude that the radix-hash join exceeds the sort-merge join for 256-bit SIMD. None of the papers mentioned above take advantage of 512-bit SIMD. There is research on SIMD sorting using 512-bit SIMD ([13],[14]). Still, to our knowledge, no research exists on implementing sort-merge join with the same optimizations and concepts like multiway merging to 512 SIMD registers.

[12]

5 Project Plan

Sketch of a time line for the thesis with major milestones, e.g.

Time	Writing/Research	Prototype
------	------------------	-----------

Table 1: Planned Time Table

References

- [1] Martina-Cezara Albutiu, Alfons Kemper, and Thomas Neumann. Massively parallel sort-merge joins in main memory multi-core database systems. *Proc. VLDB Endow.*, 5(10):1064–1075, jun 2012.
- [2] Cagri Balkesen, Gustavo Alonso, Jens Teubner, and M. Tamer Özsu. Multi-core, main-memory joins: sort vs. hash revisited. *Proc. VLDB Endow.*, 7(1):85–96, sep 2013.
- [3] Cagri Balkesen, Jens Teubner, Gustavo Alonso, and M. Tamer Özsu. Main-memory hash joins on multi-core cpus: Tuning to the underlying hardware. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*, pages 362–373, 2013.
- [4] K. E. Batcher. Sorting networks and their applications. In *Proceedings of the April 30–May 2, 1968, Spring Joint Computer Conference, AFIPS '68 (Spring)*, page 307–314, New York, NY, USA, 1968. Association for Computing Machinery.
- [5] Spyros Blanas, Yinan Li, and Jignesh M. Patel. Design and evaluation of main memory hash join algorithms for multi-core cpus. In *Proceedings of the 2011 ACM SIGMOD International Conference on Management of Data, SIGMOD '11*, page 37–48, New York, NY, USA, 2011. Association for Computing Machinery.
- [6] S. Chen, A. Ailamaki, P.B. Gibbons, and T.C. Mowry. Improving hash join performance through prefetching. In *Proceedings. 20th International Conference on Data Engineering*, pages 116–127, 2004.
- [7] Jatin Chhugani, Anthony D. Nguyen, Victor W. Lee, William Macy, Mostafa Hagog, Yen-Kuang Chen, Akram Baransi, Sanjeev Kumar, and Pradeep Dubey. Efficient implementation of sorting on multi-core simd cpu architecture. *Proc. VLDB Endow.*, 1(2):1313–1324, aug 2008.
- [8] Markus Dreseler, Jan Kossmann, Martin Boissier, Stefan Klauck, Matthias Uflacker, and Hasso Plattner. Hyrise re-engineered: An extensible database system for research in relational in-memory data management. In Melanie Herschel, Helena Galhardas, Berthold Reinwald, Irini Fundulaki, Carsten Binnig, and Zoi Kaoudi, editors, *Advances in Database Technology - 22nd International Conference on Extending Database Technology, EDBT 2019, Lisbon, Portugal, March 26-29, 2019*, pages 313–324. OpenProceedings.org, 2019.
- [9] Hiroshi Inoue, Takao Moriyama, Hideaki Komatsu, and Toshio Nakatani. Aa-sort: A new parallel sorting algorithm for multi-core simd processors. In *16th International Conference on Parallel Architecture and Compilation Techniques (PACT 2007)*, pages 189–198, 2007.
- [10] Changkyu Kim, Tim Kaldewey, Victor W. Lee, Eric Sedlar, Anthony D. Nguyen, Nadathur Satish, Jatin Chhugani, Andrea Di Blas, and Pradeep Dubey. Sort vs. hash revisited: fast join implementation on modern multi-core cpus. *Proc. VLDB Endow.*, 2(2):1378–1389, aug 2009.
- [11] Chris Nyberg, Tom Barclay, Zarka Cvetanovic, Jim Gray, and Dave Lomet. Alpha-sort: a risc machine sort, 1994.
- [12] Saher Odeh, Oded Green, Zahi Mwassi, Oz Shmueli, and Yitzhak Birk. Merge path - parallel merging made simple. pages 1611–1618, 05 2012.
- [13] Alex Watkins and Oded Green. A fast and simple approach to merge and merge sort

using wide vector instructions. 11 2018.

- [14] Zekun Yin, Tianyu Zhang, André Müller, Hui Liu, Yanjie Wei, Bertil Schmidt, and Weiguo Liu. Efficient parallel sort on avx-512-based multi-core and many-core architectures. In *2019 IEEE 21st International Conference on High Performance Computing and Communications; IEEE 17th International Conference on Smart City; IEEE 5th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*, pages 168–176, 2019.