

CStore: A Lock-free SIMD Skip List

Ziyi Liu (ziyil), Quan Quan (qqquan)

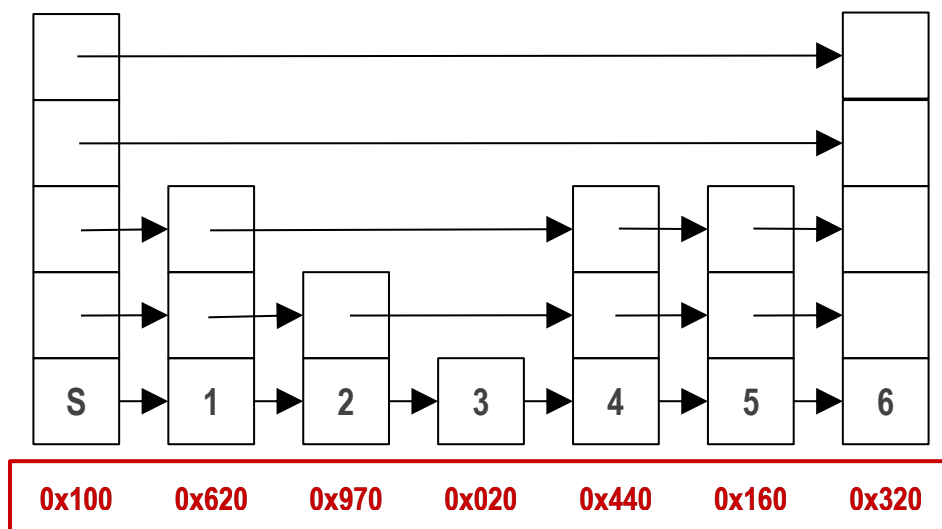
SUMMARY

In this project we implemented a lock-free SIMD skiplist which supports insert, delete and find. The project deliverable will be the library for this skiplist, and performance comparison with `std::map` and other lock-free non-SIMD skiplist on GHC machines.

BACKGROUND

According to Wikipedia, skiplist is a data structure that allows fast search within an ordered sequence of elements. On average it costs $O(\log N)$ to do searches.

However, since skiplist is an array of linked lists, and linked list has poor data locality (because it has to load memory again when going to a new node, as the example in the following graph shows), so it is actually slower than balance tree implementations like red-black tree.



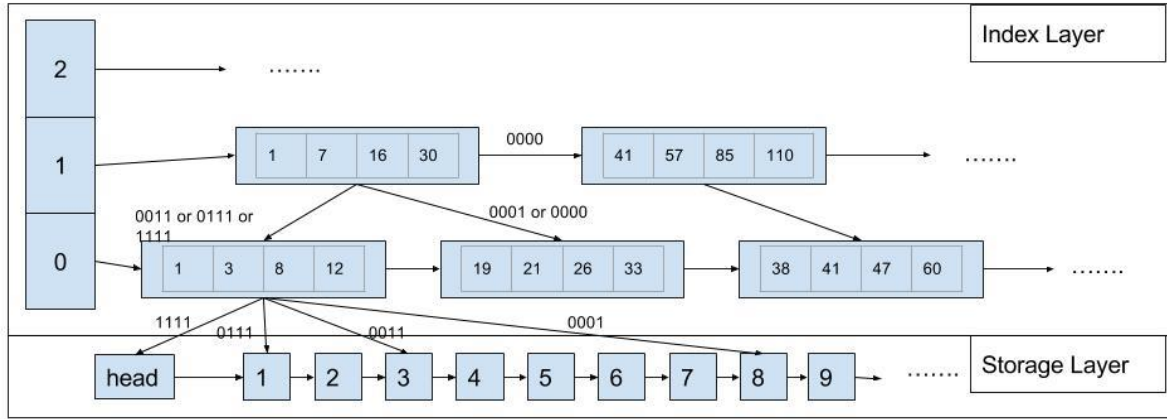
Graph 1. Skip List has poor data locality

One solution we apply in CStore is we group adjacent indexes together to improve data locality, and this also introduces the chance of applying SIMD to do comparisons and reduce the number of comparisons we have to do to search indexes horizontally.

Another case is we have to handle the multi-threading scenario. There has been some lock-free solutions for skiplist, but since we are changing the structure of skiplist, they cannot apply in our case. We make use of the lock-free linked list solution and come up with some adaptations to suit our design and maintain the lock-free attribute.

APPROACH

Overall structure design:



Graph 2. Overall structure of CStore

Above is a representation of the architecture of our skiplist. We divide the skiplist into index layer and storage layer. In the index layer, we group multiple indexes of a node together in one index node, so we can make use of cache locality and apply SIMD comparisons. The index layer will not change until some upper limit of the number of inserts and deletes are triggered, in which case we will build a new index layer to replace the old one. **Our assumption is the requests are uniformly distributed on the key space, so making the index layer unchanged for some time will not hurt performance much.** In the storage layer, we only have an ordered singly-linked lock-free skiplist.

Index layer design:

Each index node holds VECTOR_SIZE number of indexes. On receiving a query, each index node will launch a gang of instances, each instance comparing the key to an index and report the result in the corresponding slot in the output array. Then we use the pre-built routing-table to find where to go next. VECTOR_SIZE is a configurable parameter, on GHC machine we can use avx instruction set which can compile 8-width vector, so a VECTOR_SIZE greater or equal to 8 is preferable. But a large vector may lead to excessive computation, since in skip lists most of the time we suppose the target key is in the middle of the indexes. After doing test with various data size, a VECTOR_SIZE of 8 proves to have the best performance.

Node Count	vector_size=8	vector_size=16	vector_size=24
5,000,000	0.311	0.310	0.325
10,000,000	0.342	0.343	0.358
25,000,000	0.37	0.379	0.395
50,000,000	0.396	0.402	0.43
100,000,000	0.427	0.445	0.473

Table 1. Query time with different vector sizes
(thread number = 16, total requests = 32M, read-only test)

There is no parallelism in building the index layer, because anyway we have to go through the entire storage layer and the storage layer is a linked list which is hard to do data

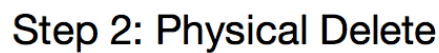
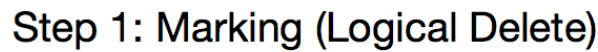
parallelization. We have a configurable parameter REBUILD_THRESHOLD to denote the maximum number of inserts and deletes before a rebuild of index layer. This parameter should be decided based on the number of nodes stored. In our testing, it is definitely not suggested to do frequent rebuilding, because going through the whole list to rebuild consumes a lot of resources. But scarcely rebuilding the index layer may also lead to slow lookups, so it is a tradeoff and this parameter needs to be carefully tuned. In our testing, based on a test setting of 500k initial nodes, 32 million requests with 60% inserts, 20% reads and 20% deletes, and 5 billion key space, 1/100, 1/5 and 2 times of the number of nodes all prove to be a good threshold.

rebuild threshold	time spent
1,000	7.5
2,500	7.5
5,000	7.5
7,500	7.6
10,000	7.4
20,000	8.1
1,000,000	7.2
2,000,000	8.1
5,000,000	9.1
∞	9.3

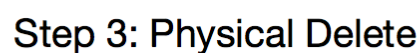
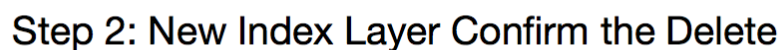
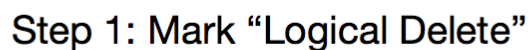
Table 2. Query time with different rebuild threshold

Storage layer design:

Our storage layer is a linked list, so we first turn to lock-free linked list solution suggested by Harris^[1]. The idea is to add a `is_delete` mark to each data node's next pointer, so we avoid the insert-delete conflict problem.



Later we find out that the index layer may have pointers pointing to a data node, so a data node should not be easily deleted or else the index layer may point to freed memory. We introduce a new `confirm_delete` mark to the node to help address this problem. `confirm_delete` will only be marked when rebuilding the index layer, so the new index layer will not build indexes on the nodes marked `is_deleted`, and add `confirm_delete` mark to this node so later operations can safely do physical delete.



Graph 3. Lock-free CStore storage layer delete steps

RESULTS

In order to give a comprehensive testing of CStore, we test the read-only scenarios and mix-requests scenarios separately. We will measure CStore's performance with time spent on the queries, and compare it against `std::map` (red-black tree) in read-only scenario and a non-SIMD lock-free skiplist implementation we found in this github repo^[2] (author: Jung-Sang Ahn jungsang.ahn@gmail.com). In both scenario, we generate requests with 16 threads launched by `pthread_create()` for each data structure, and each thread will send 200,000 requests and that's 3,200,000 in total. For each test run we measure the time spent as the time of the last thread exits minus the time of the first thread starts. The accepted time is the minimum time of three consecutive tests.

Read only test

First let's look at the read-only scenario. As setup, we give CStore and `std::map` the same number of initial key-value nodes, and the initial key-value data for both data structures are the same. Each thread will send the exact same requests to both data structures.

Node Count	CStore perf.	std::map perf.
5,000,000	0.308	0.274
10,000,000	0.34	0.318
25,000,000	0.375	0.37
50,000,000	0.406	0.421
100,000,000	0.445	0.46

Table 3. Read-only performance comparison with `std::map`



Graph 4. Read-only performance comparison with `std::map`

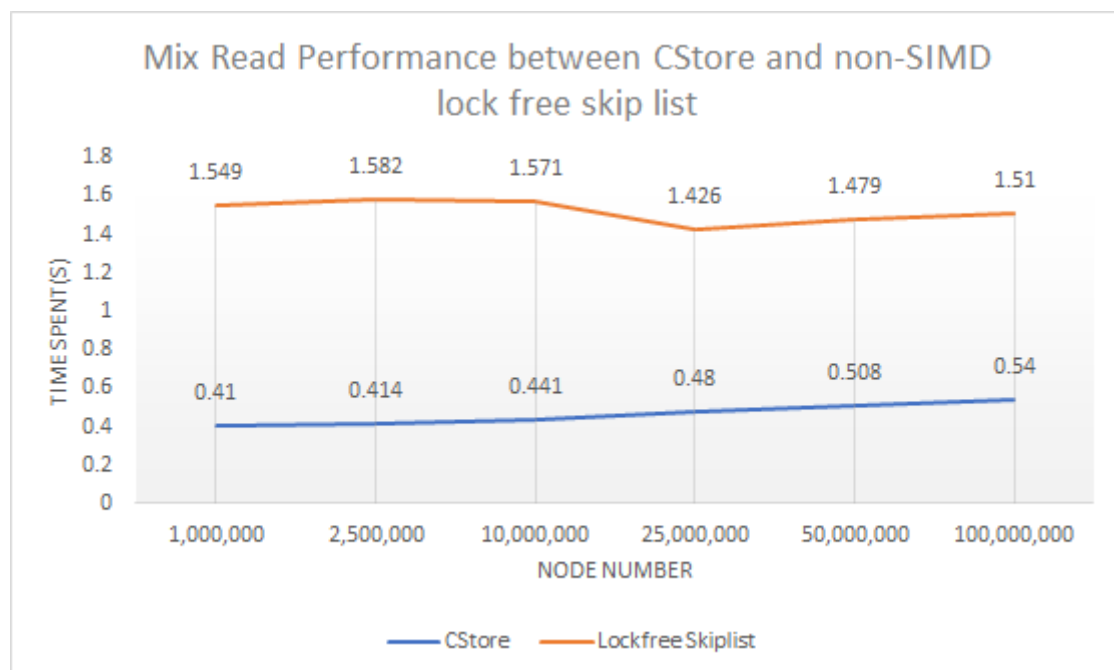
We can see from the result that when the number of nodes is not so huge, std::map has notable better performance than CStore, however with the number of nodes growing, CStore shows stable performance and outperforms std::map at last.

Mix test

Next we will look at CStore's performance against a non-SIMD lock free skip list. The mix requests consists of ~20% reads, ~60% inserts and ~20% deletes.

Initial Node Count	CStore	Lockfree Skiplist
1,000,000	0.41	1.549
2,500,000	0.414	1.582
10,000,000	0.441	1.571
25,000,000	0.48	1.426
50,000,000	0.508	1.479
100,000,000	0.54	1.51

Table 4. Mix request performance comparison with non-SIMD lock-free skiplist



Graph 5. Mix request performance comparison with non-SIMD lock-free skiplist

In this test our solution is about 3 times faster than the non-SIMD lock free skiplist solution. Because of the fixed number of requests, with less nodes, the performance tends to be not consistent with that of many nodes because of high contention.

LIMITATION ANALYSIS

Because we are optimizing a linked-list based data structure, there are many techniques we learnt from this class cannot be applied here. Also, although combining indexes together improves data locality, because we are splitting the indexes for each data node into several

levels of index nodes, we actually created more linked-list nodes, and introduced constant more memory access (each time going down one level, need to load the memory once more).

In our implementation, we have different scenarios of read-heavy, insert heavy or read-only. We use perf to profile our implementation.

Insert 60%, read 20%, delete 20%

Samples: 2M of event 'cycles', Event count (approx.): 2033452723875

Overhead	Command	Shared Object	Symbol
78.58%	skiplist	skiplist	[.] LockfreeList::search_by_index
13.39%	skiplist	skiplist	[.] index_compare_ispc_sum
2.65%	skiplist	[kernel.kallsyms]	[k] read_hpet
1.44%	skiplist	skiplist	[.] testrun
0.70%	skiplist	skiplist	[.] std::uniform_int_distribution<int>::operator()

Insert 20%, read 60%, delete 20%

Samples: 1M of event 'cycles', Event count (approx.): 1616298617138

Overhead	Command	Shared Object	Symbol
79.25%	skiplist	skiplist	[.] LockfreeList::search_by_index
14.51%	skiplist	skiplist	[.] index_compare_ispc_sum
2.30%	skiplist	[kernel.kallsyms]	[k] read_hpet
0.97%	skiplist	skiplist	[.] std::uniform_int_distribution<int>::operator()
0.76%	skiplist	skiplist	[.] testrun

read only

Samples: 609K of event 'cycles', Event count (approx.): 506172951481

Overhead	Command	Shared Object	Symbol
63.19%	skiplist	skiplist	[.] LockfreeList::search_by_index
30.44%	skiplist	skiplist	[.] index_compare_ispc_sum
2.86%	skiplist	skiplist	[.] std::uniform_int_distribution<int>::operator()
2.04%	skiplist	[kernel.kallsyms]	[k] read_hpet
0.80%	skiplist	skiplist	[.] std::mersenne_twister_engine<unsigned long, 32

As we can see from the results, the main bottleneck is in the storage layer lock-free searching part, especially in insert heavy scenarios when conflict and backoffs are more common. This is hard to optimize since we already have an optimal lock-free solution. In read-only situation, the ispc SIMD compare takes a good amount of time. This is because each read/insert/delete operation has to go through the index layer and in each index node we have to run ispc SIMD compare to route.

Another problem is it seems SIMD instructions actually provides little speedup. The reason is if having proper max height and reasonable index height for each data node, the number of times to access index node for a lookup is small. For example, in our test above we generally have 10 million nodes. In this case a height of 23 (because 10 million is around 2^{23}) is sufficient for efficient lookup. Then, based on the probabilistic model, we merely need to traverse through the index layer horizontally and only need to access 23 index nodes to get to the storage layer. Comparing to the memory loading time, the speedup brought by reducing 4 or 5 cycles on comparison is minor.

Citations

- [1] Harris, Timothy L. "A pragmatic implementation of non-blocking linked-lists." *International Symposium on Distributed Computing*. Springer Berlin Heidelberg, 2001.
- [2] <https://github.com/greensky00/skiplist>

LIST OF WORK BY EACH STUDENT

Equal work was performed by both project members.