Custom 8-bit floating point value format for reducing shared memory bank conflict in approximate nearest neighbor search

Goal

- made.









interpret as FP32

While the representation accuracy of e4m4 is not significantly better than e5m3, representable range is only half. \Rightarrow We use e5m3.

The accuracy of ANNS

truth for the query.







Parameters

- of the search phase.
- num clusters : 100,000
- batch size : 10,000
- **PQ** bit : n = 8

Conclusion

- bank conflict in IVFPQ on GPU.
- The sign bit is omitted.
- little recall degradation.

Hiroyuki Ootomo¹, Akira Naruse² ¹Tokyo Institute of Technology ²NVIDIA

(*) num_probes : The number of clusters picked up in the first stage

We have developed custom 8-bit floating point formats for reducing

It can be converted from/to FP32 with a few operations. We have applied it to IVFPQ and improved the throughput with a