## Grappa: enabling next-generation analytics tools via latency-tolerant distributed shared memory

Jacob Nelson, Brandon Myers, Brandon Holt, Preston Briggs, Luis Ceze, Simon Kahan, Mark Oskin {nelson, bholt, bdmyers, preston, luisceze, skahan, oskin}@cs.washington.edu University of Washington



More info, tech report, downloads: http://grappa.io



Grappa is a modern take on software distributed shared memory, tailored to exploit parallelism inherent in data-intensive applications to overcome their poor locality and input-dependent load distribution.

Grappa differs from traditional DSMs in three ways:

- Instead of minimizing per-operation latency for performance, Grappa *tolerates latency* with concurrency (latency-sensitive apps need not apply!)
- Grappa moves computation to data instead of caching data at computation
- Grappa operates at *byte granularity* rather than page granularity

## **Programming example**

Grappa's familiar single-system multithreaded C++ programming model enables easier development of analysis tools for terabyte-scale data. We provide sequential consistency for race-free programs using RPC-like atomic *delegate operations*, along with standard multithreaded synchronization primitives.

Here is an example of using Grappa's C++11 library to build a distributed parallel wordcount-like application with a simple hash table:

// distributed input array

GlobalAddress<char> chars;

// distributed hash table:

using Cell = map<char,int>;

Node 0 Node 2 Node 1 . . . Global Heap "p" "b" "g" | "h" "a" "o" "q" "d" "i" "C" "X" "C" hash("i") Cell[3] Cell[4] Cell[5] Cell[0] Cell[1] Cell[2] - - 🖌 "e"→1 |"a"→7 "b"→1 "i"→5 "f"→2 "c"→3 "|"→1

Latency tolerance has been applied successfully in hardware for nanosecond latencies (e.g., superscalar processors and GPUs). This project explores the application of this idea at distributed system scales with millisecond latencies.



Experiments run at Pacific Northwest National Laboratory



arks run on 31 nodes us

## **Key feature: Message Aggregation**

Commodity networks have a limited message injection rate, so building up larger packets from unrelated tasks is essential for smallmessage throughput (fine-grained random access to global memory).



handlers run messages via MPI/RDMA to dest. cores on home cores



locally per core

into buffer

Query on Seał with two different k values, compared with Spark using MEMORY\_ONLY fault tolerance

edge Friendster social network graph and 1.4B edge Twitter follower dataset, compared with GraphLab using two partitioning strategies

compared with Shark distributed query engine











