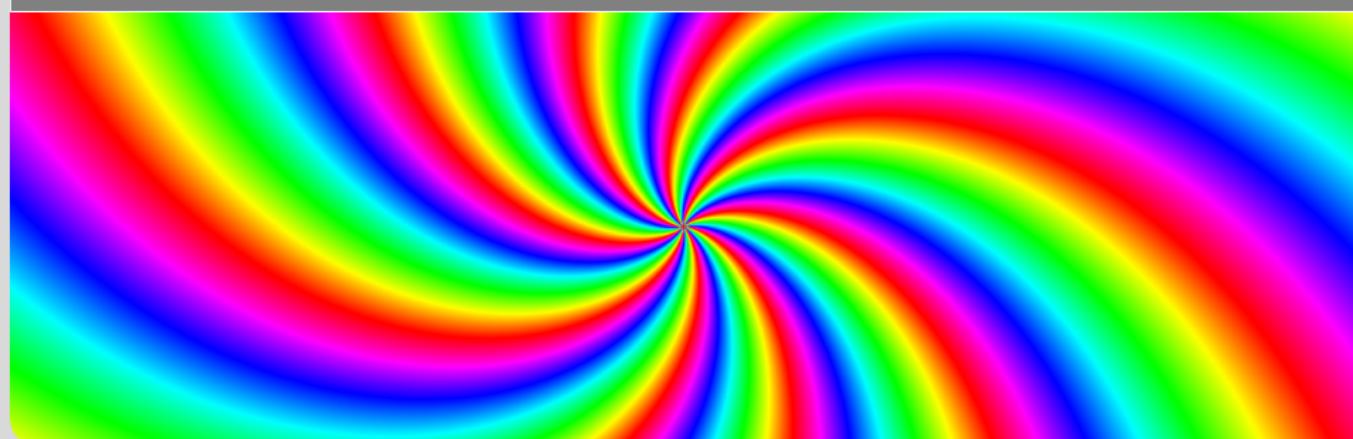


Thrill : High-Performance Algorithmic Distributed Batch Data Processing with C++

Timo Bingmann, Michael Axtmann, Peter Sanders, Sebastian Schlag, and 6 Students · 2019-08-28

INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMIC



Weak-Scaling Benchmarks

WordCountCC – $h \cdot 49$ GiB 222 lines

- Reduce text files from CommonCrawl web corpus.

PageRank – $h \cdot 2.7$ GiB, $|E| \approx h \cdot 158$ M 410 lines

- Calculate PageRank using join of current ranks with outgoing links and reduce by contributions. 10 iterations.

TeraSort – $h \cdot 16$ GiB 141 lines

- Distributed (external) sorting of 100 byte random records.

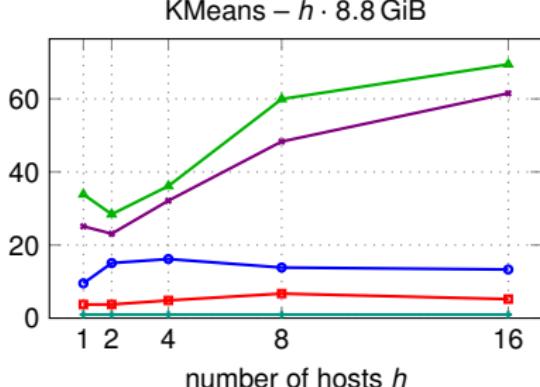
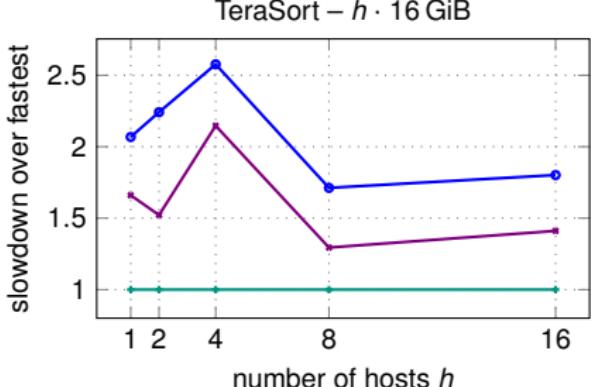
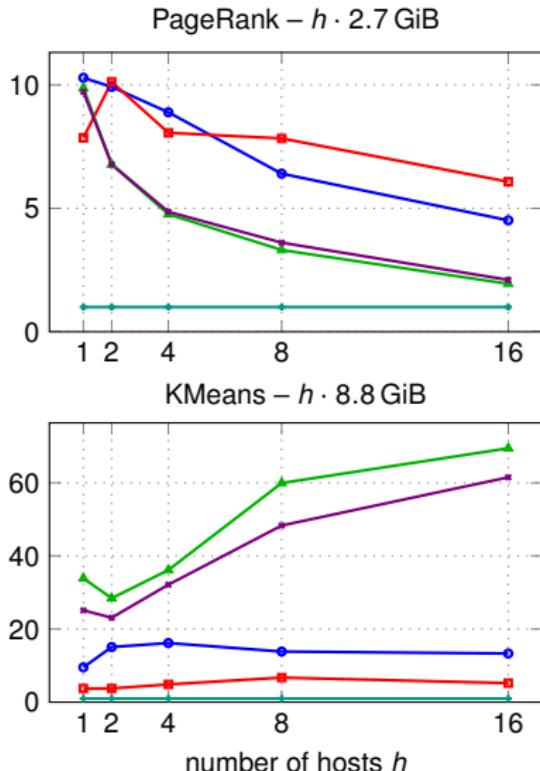
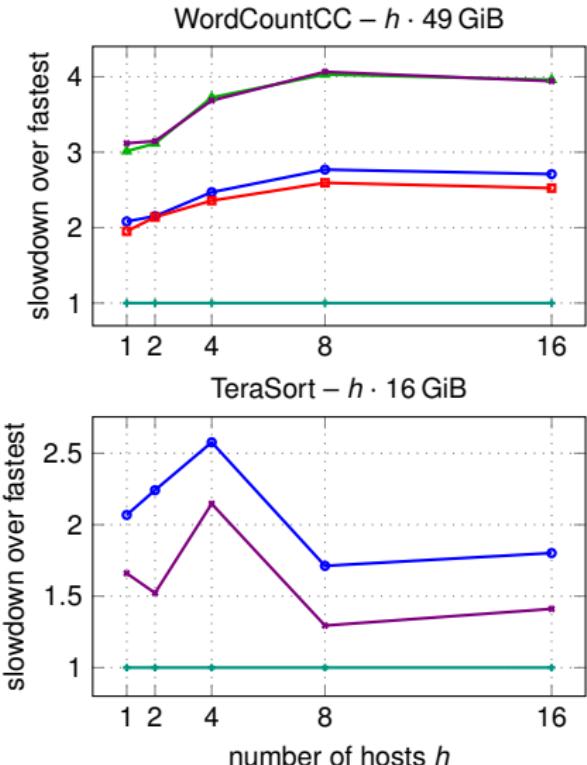
K-Means – $h \cdot 8.8$ GiB 357 lines

- Calculate K-Means clustering with 10 iterations.

Platform: $h \times$ r3.8xlarge systems on Amazon EC2 Cloud

- 32 cores, Intel Xeon E5-2670v2, 2.5 GHz clock, 244 GiB RAM, 2 x 320 GB local SSD disk, ≈ 400 MiB/s read/write Ethernet network ≈ 1000 MiB/s throughput, Ubuntu 16.04.

Experimental Results: Slowdowns



—●— Spark (Java) —■— Spark (Scala) —▲— Flink (Java) —●— Flink (Scala) —◆— Thrill

Example $T = [\text{tobeornottobe\$}]$

SA_i	LCP_i	$T_{SA_i \dots n}$
13	-	\$
11	0	b e \$
2	2	b e o r n o t t o b e \$
12	0	e \$
3	1	e o r n o t t o b e \$
6	0	n o t t o b e \$
10	0	o b e \$
1	3	o b e o r n o t t o b e \$
4	1	o r n o t t o b e \$
7	1	o t t o b e \$
5	0	r n o t t o b e \$
9	4	t o b e \$
0	1	t o b e o r n o t t o b e \$
8	0	t t o b e \$

Suffix Sorting with DC3: Example

	0	1	2	3	4	5	6	7	8	9	10	
$T = [d \boxed{b} a c \boxed{b} a c \boxed{b} d \$ \$]$												$= [t_i]_{i=0, \dots, n-1}$
triples	(bac,1),	(bac,4),	(bd\$,7),	(acb,2)	(acb,5),	(d\$\$,8)						
sorted	(acb,2)	(acb,5),	(bac,1),	(bac,4),	(bd\$,7),	(d\$\$,8)						
equal 0/1	0	0	1	0	1	1						
prefix sum	0	0	1	1	2	3						
$R = \boxed{1 \ 1 \ 2} \ \boxed{0 \ 0 \ 3} \ \$$							$r_1 \ r_4 \ r_7 \ r_2 \ r_5 \ r_8$					
$\text{SA}_R = 3 \ 4 \ 0 \ 1 \ 2 \ 5 \ \$$							$\text{ISA}_R = \boxed{2 \ 3 \ 4} \ \boxed{0 \ 1 \ 5} \ \$$					
$S_0 = [(d, b, \color{red}{2}, \color{blue}{0}, \color{teal}{0}), (c, b, \color{red}{3}, \color{blue}{1}, \color{teal}{3}), (c, b, \color{red}{4}, \color{blue}{5}, \color{teal}{6})]$							$(t_i, t_{i+1}, \color{red}{r_{i+1}}, \color{blue}{r_{i+2}}, \color{teal}{i})$					
$S_1 = [(\color{red}{2}, b, \color{blue}{0}, \color{teal}{1}), (\color{red}{3}, b, \color{blue}{1}, \color{teal}{4}), (\color{red}{4}, b, \color{blue}{5}, \color{teal}{7})]$							$(\color{red}{r_{i+1}}, t_{i+1}, \color{blue}{r_{i+2}}, \color{teal}{i+1})$					
$S_2 = [(\color{blue}{0}, a, c, \color{red}{3}, \color{teal}{2}), (\color{blue}{1}, a, c, \color{red}{4}, \color{teal}{5}), (\color{blue}{5}, d, \$, \color{red}{6}, \color{teal}{8})]$							$(\color{blue}{r_{i+2}}, t_{i+2}, t'_{i+3}, \color{red}{r'_{i+4}}, \color{teal}{i+2})$					
$\text{SA}_T = \text{Merge}(\text{Sort}(S_0), \text{Sort}(S_1), \text{Sort}(S_2))$												

$\Theta(\text{sort}(n))$



Google Cloud Platform

bwUniCluster
KIT (SCC)

Flavours of Big Data Frameworks

- Batch Processing

Google's MapReduce, Hadoop MapReduce , Apache Spark ,
Apache Flink  (Stratosphere), Google's FlumeJava.

- High Performance Computing (Supercomputers)

MPI

- Real-time Stream Processing

Apache Storm , Apache Spark Streaming, Google's MillWheel.

- Interactive Cached Queries

Google's Dremel, Powerdrill and BigQuery, Apache Drill .

- Sharded (NoSQL) Databases and Data Warehouses

MongoDB , Apache Cassandra, Apache Hive, Google BigTable,
Hypertable, Amazon RedShift, FoundationDB.

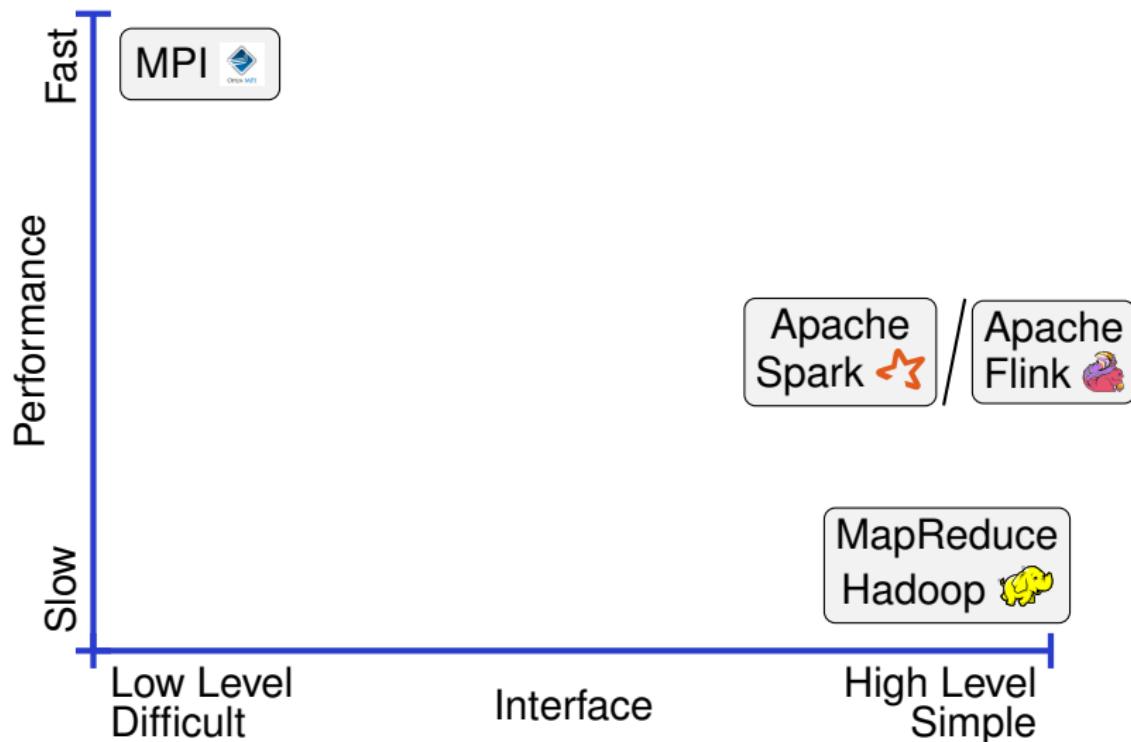
- Graph Processing

Google's Pregel, GraphLab , Giraph , GraphChi.

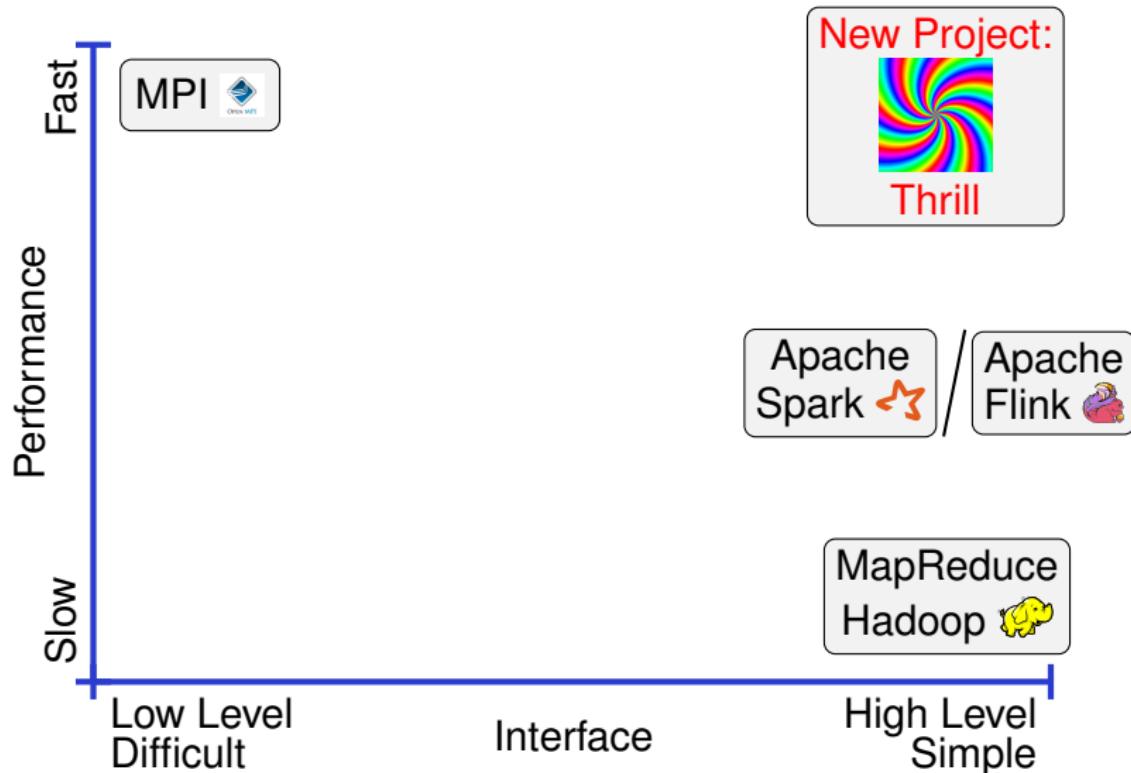
- Time-based Distributed Processing

Microsoft's Dryad, Microsoft's Naiad.

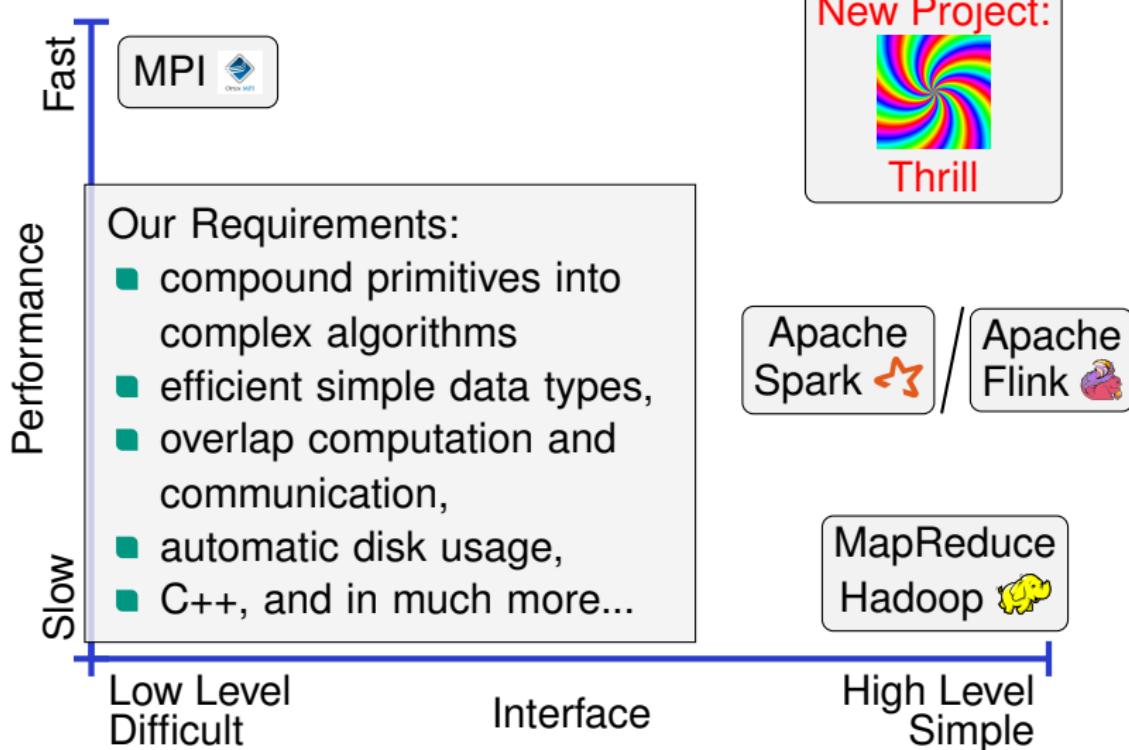
Big Data Batch Processing



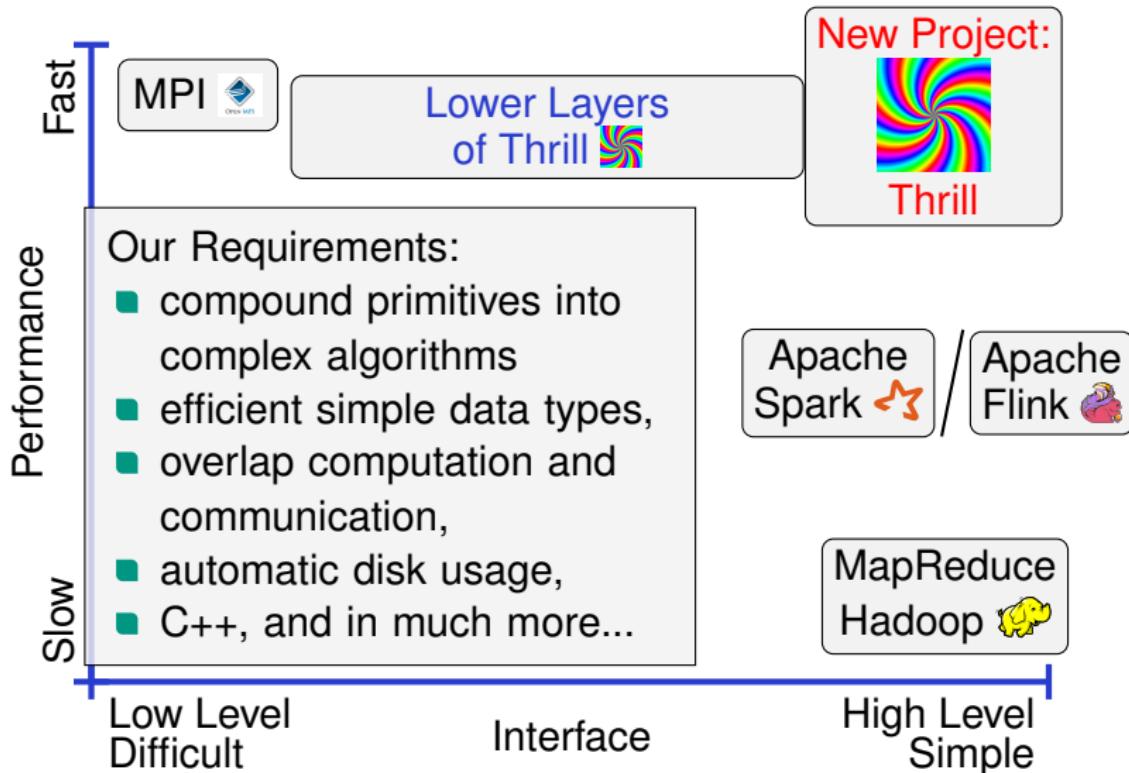
Big Data Batch Processing



Big Data Batch Processing



Big Data Batch Processing



Thrill's Goal and Current Status

An easy way to program distributed external algorithms in C++.

Current Status:

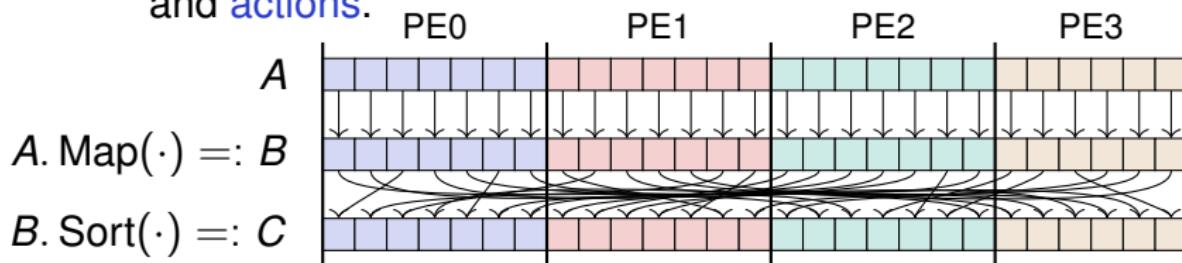
- Open-source prototype at <http://github.com/thrill/thrill>.
- $\approx 60\text{ K}$ lines of C++14 code, written by ≥ 12 contributors.
- Published at IEEE Conference on Big Data [B, et al. '16]
- Faster than Apache Spark and Flink on five micro benchmarks: WordCount1000/CC, PageRank, TeraSort, and K-Means.

Case Studies:

- Five suffix sorting algorithms [B, Gog, Kurpicz, BigData'18]
- Louvain graph clustering algorithm [Hamann et al. Euro-Par'18]
- Process scientific data on HPC (poster) [Karabin et al. SC'18]
- More: stochastic gradient descent, triangle counting, etc.
- Future: fault tolerance, scalability, predictability, and more.

Distributed Immutable Array (DIA)

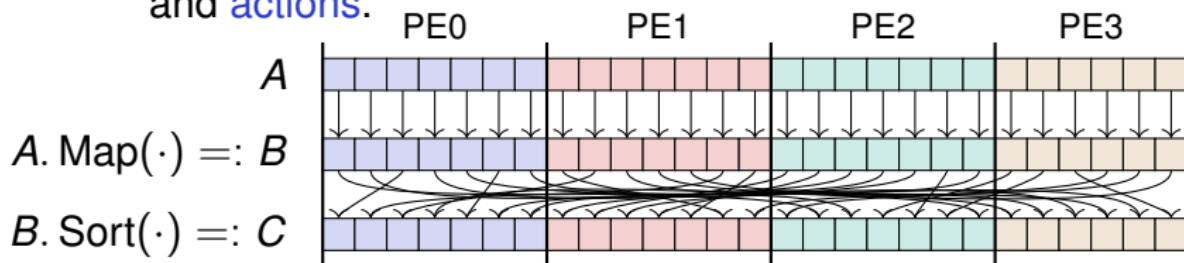
- User Programmer's View:
 - $\text{DIA}\langle T \rangle$ = result of an operation (local or distributed).
 - Model: distributed array of items T on the cluster
 - Cannot access items directly, instead use transformations and actions.



Distributed Immutable Array (DIA)

- User Programmer's View:

- $\text{DIA}\langle T \rangle$ = result of an operation (local or distributed).
- Model: distributed array of items T on the cluster
- Cannot access items directly, instead use transformations and actions.



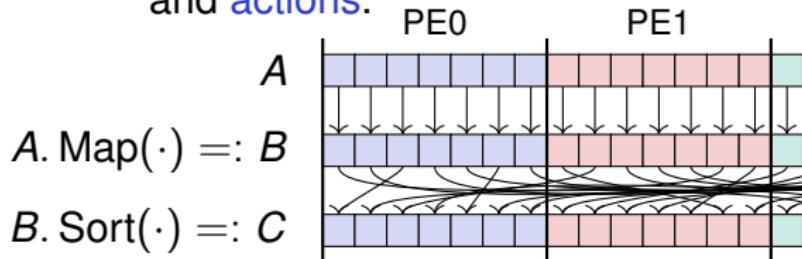
- Framework Designer's View:

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable. \Rightarrow build data-flow graph.
- $\text{DIA}\langle T \rangle$ = chain of computation items
- Let distributed operations choose “materialization”.

Distributed Immutable Array (DIA)

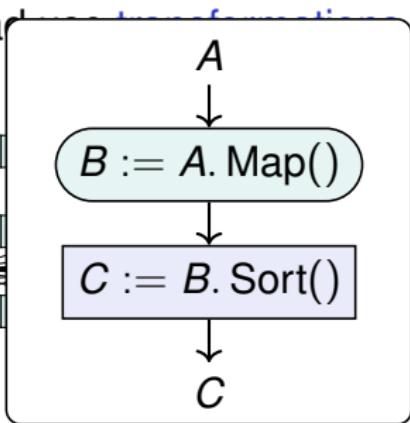
- User Programmer's View:

- DIA<T> = **result** of an operation (local or distributed).
- Model: **distributed array** of items T on the cluster
- Cannot access items directly, instead **operations** and **actions**.



- Framework Designer's View:

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable. \implies **build data-flow graph**.
- DIA<T> = **chain of computation items**
- Let distributed operations choose “materialization”.

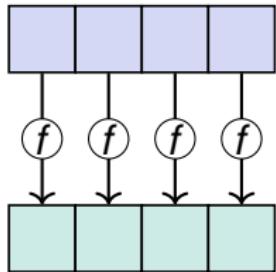


List of Primitives (Excerpt)

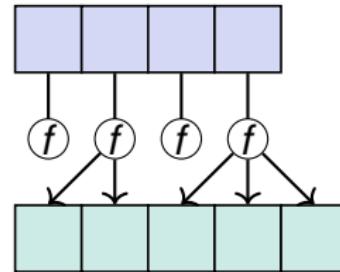
- Local Operations (**LOp**): input is **one item**, output ≥ 0 items.
Map(), **Filter()**, **FlatMap()**.
- Distributed Operations (**DOp**): input is a **DIA**, output is a **DIA**.
 - Sort()** Sort a DIA using comparisons.
 - ReduceBy()** Shuffle with Key Extractor, Hasher, and associative Reducer.
 - GroupBy()** Like ReduceBy, but with a general Reducer.
 - PrefixSum()** Compute (generalized) prefix sum on DIA.
 - Window_k()** Scan all k consecutive DIA items.
 - Zip()** Combine equal sized DIAs item-wise.
 - Union()** Combine equal typed DIAs in arbitrary order.
 - InnerJoin()** Join items from two DIAs by key.
- **Actions**: input is a **DIA**, output: ≥ 0 items **on every worker**.
Sum(), **Min()**, **ReadLines()**, **WriteLines()**, pretty much open.

Local Operations (LOps)

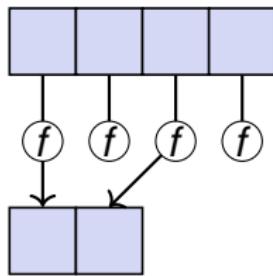
Map(f) : $\langle A \rangle \rightarrow \langle B \rangle$
 $f : A \rightarrow B$



FlatMap $\langle B \rangle(f) : \langle A \rangle \rightarrow \langle B \rangle$
 $f : A \rightarrow \text{array}(B)$



Filter(f) : $\langle A \rangle \rightarrow \langle A \rangle$
 $f : A \rightarrow \{\text{false}, \text{true}\}$

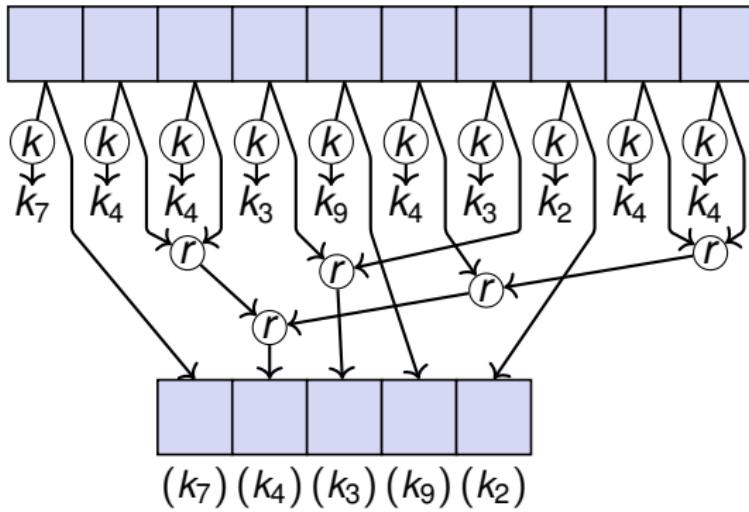


DOps: ReduceByKey

ReduceByKey(k, r) : $\langle A \rangle \rightarrow \langle A \rangle$

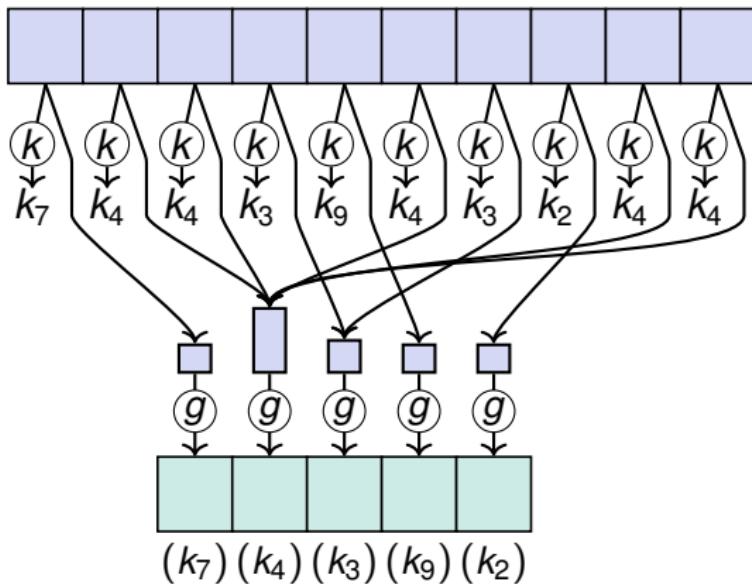
$k : A \rightarrow K$ key extractor

$r : A \times A \rightarrow A$ reduction



DOps: GroupByKey

GroupByKey(k, g) : $\langle A \rangle \rightarrow \langle B \rangle$
 $k : A \rightarrow K$ key extractor
 $g : \text{iterable}(A) \rightarrow B$ group function



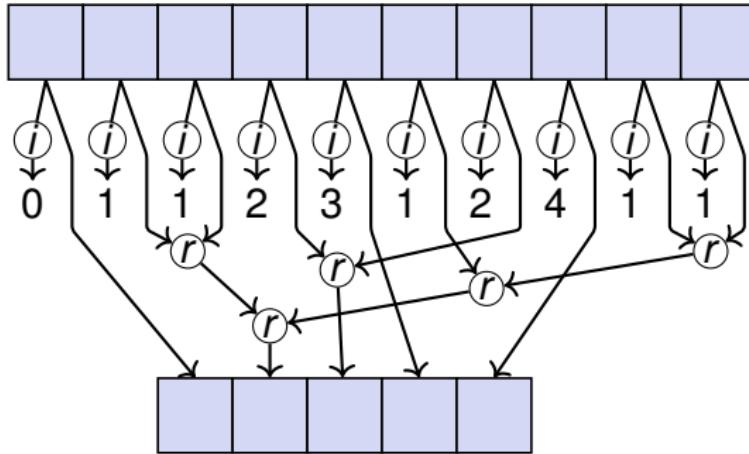
DOps: ReduceToIndex

ReduceToIndex(i, n, r) : $\langle A \rangle \rightarrow \langle A \rangle$

$i : A \rightarrow \{0..n - 1\}$ index extractor

$n \in \mathbb{N}_0$ result size

$r : A \times A \rightarrow A$ reduction



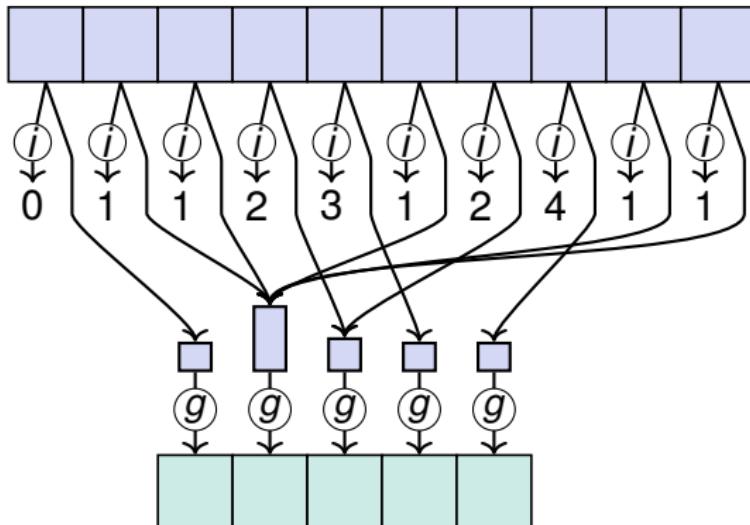
DOps: GroupToIndex

GroupToIndex(i , n , g) : $\langle A \rangle \rightarrow \langle B \rangle$

$i : A \rightarrow \{0..n - 1\}$ index extractor

$n \in \mathbb{N}_0$ result size

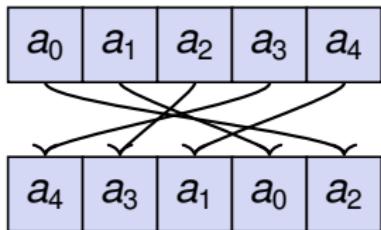
$g : \text{iterable}(A) \rightarrow B$ group function



DOps: Sort and Merge

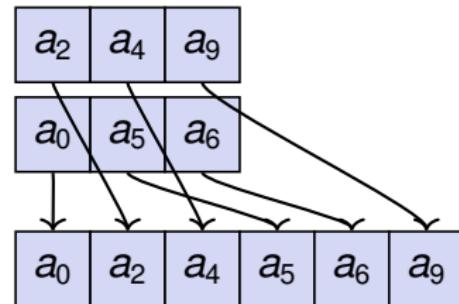
Sort(o) : $\langle A \rangle \rightarrow \langle A \rangle$

$o : A \times A \rightarrow \{ \text{false}, \text{true} \}$
(less) order relation



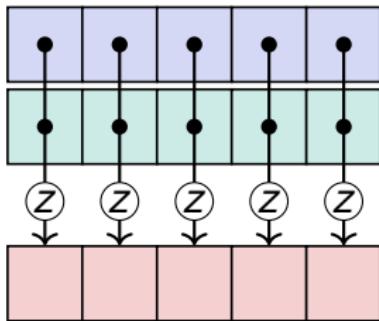
Merge(o) : $\langle A \rangle \times \langle A \rangle \cdots \rightarrow \langle A \rangle$

$o : A \times A \rightarrow \{ \text{false}, \text{true} \}$
(less) order relation

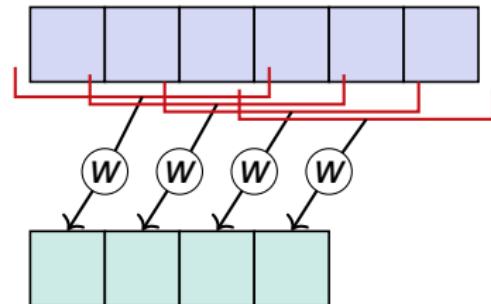


DOps: Zip and Window

Zip(z) : $\langle A \rangle \times \langle B \rangle \cdots \rightarrow \langle C \rangle$
 $z : A \times B \rightarrow C$
 zip function



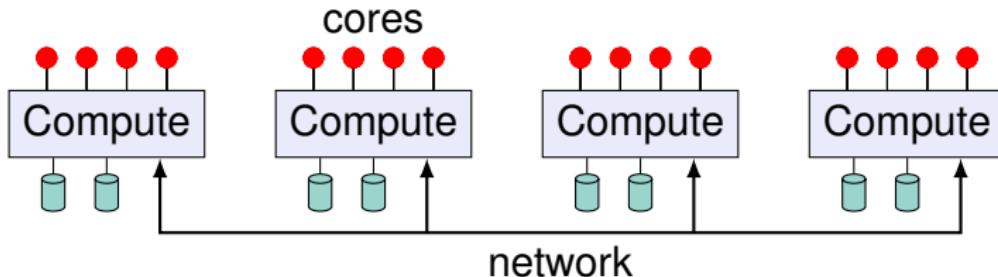
Window(k, w) : $\langle A \rangle \rightarrow \langle B \rangle$
 $k \in \mathbb{N}$ window size
 $w : A^k \rightarrow B$ window function



Example: WordCount in Thrill

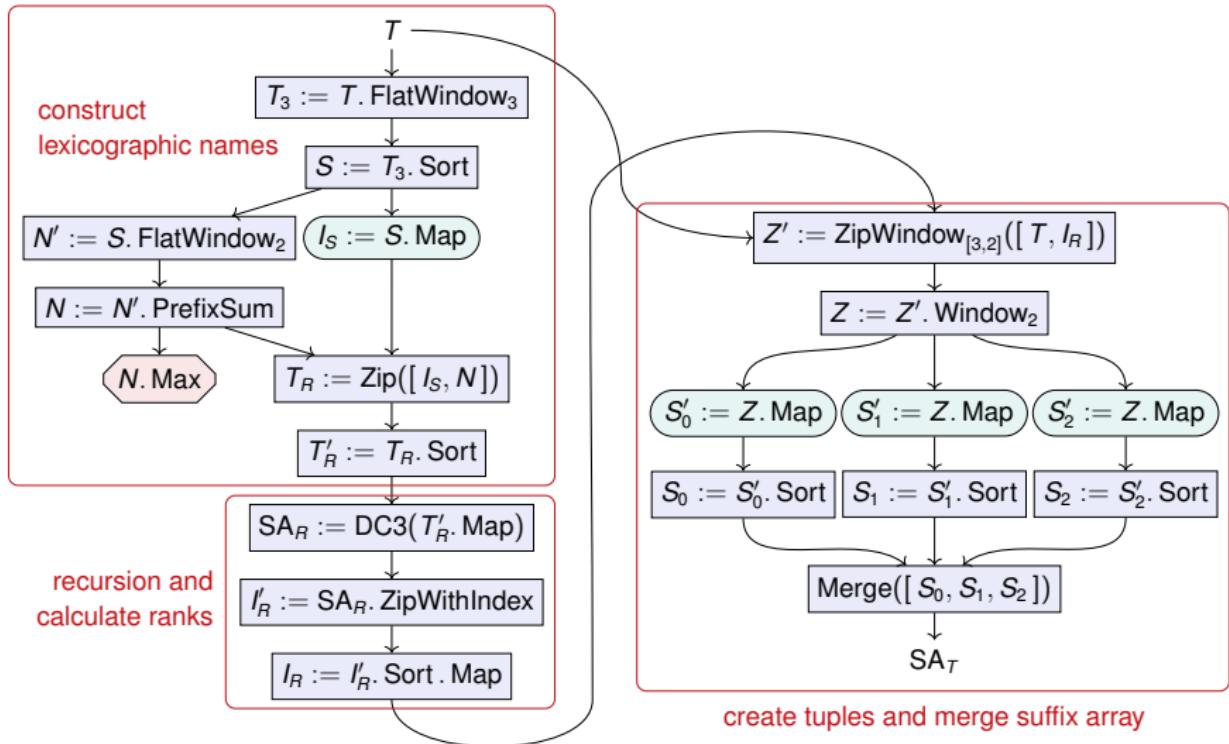
```
1 using Pair = std::pair<std::string, size_t>;
2 void WordCount(Context& ctx, std::string input, std::string output) {
3     auto word_pairs = ReadLines(ctx, input)      // DIA<std::string>
4     .FlatMap<Pair>(
5         // flatmap lambda: split and emit each word
6         [](const std::string& line, auto emit) {
7             tlx::split_view(' ', line, [&](tlx::string_view sv) {
8                 emit(Pair(sv.to_string(), 1)); });
9         });
10    word_pairs.ReduceByKey(
11        // key extractor: the word string
12        [](&Pair p) { return p.first; },
13        // commutative reduction: add counters
14        [](&Pair a, &Pair b) {
15            return Pair(a.first, a.second + b.second);
16        });
17        .Map([](&Pair p) {
18            return p.first + ":" + std::to_string(p.second); })
19        .WriteLines(output);
20 }
```

Execution on Cluster



- Compile program into **one binary**, running on all hosts.
- **Collective coordination** of work on compute hosts, like MPI.
- **Control flow** is decided on by using C++ statements.
- Runs on MPI HPC clusters and on Amazon's EC2 cloud.

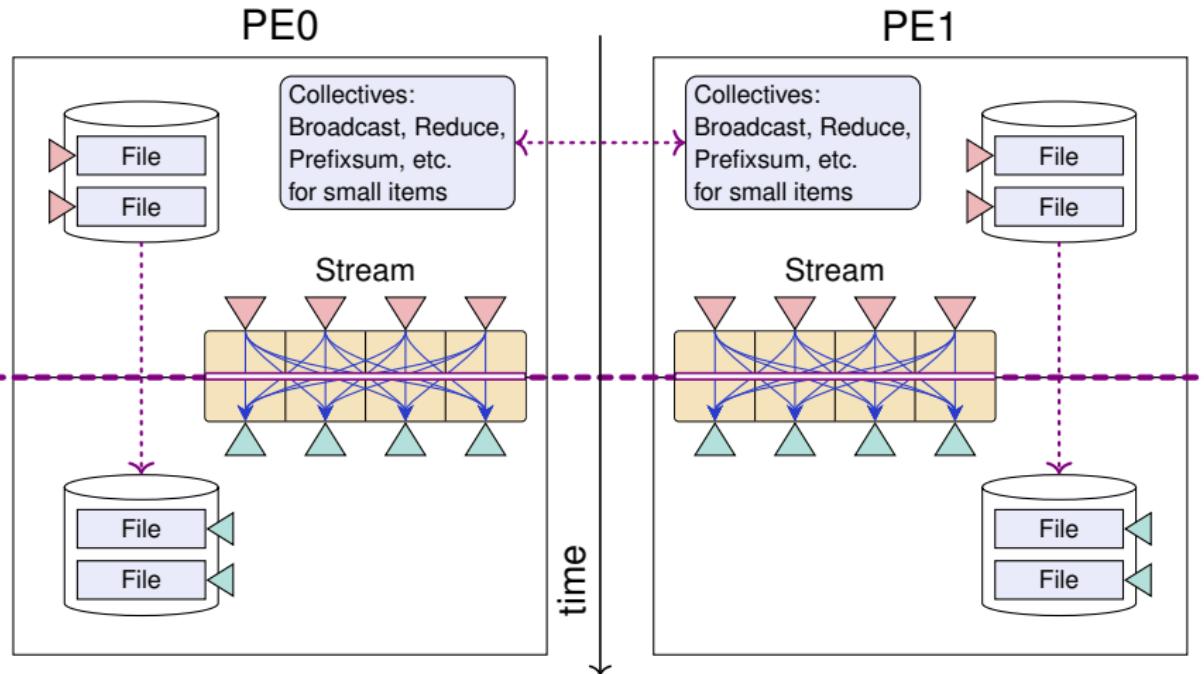
Data-Flow Graph of DC3



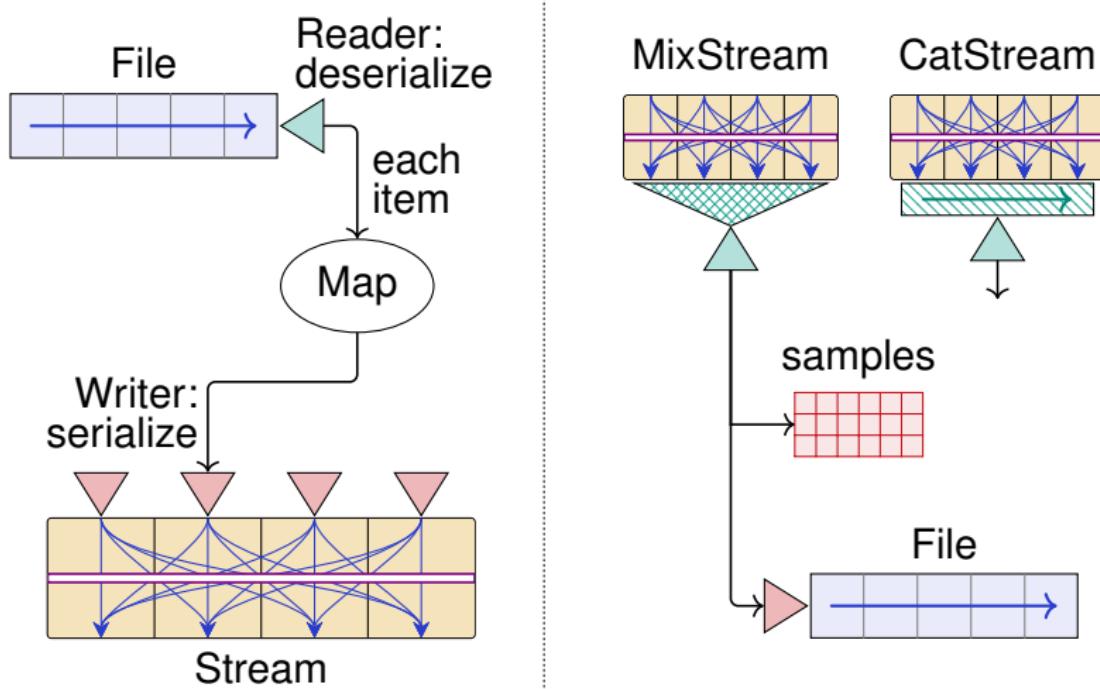
Layers of Thrill

api: High-level User Interface DIA<T>, Map, FlatMap, Filter, Reduce, Sort, Merge, ...				
core: Internal Algorithms reducing hash tables (bucket and linear probing), multiway merge, stage executor	vfs: Data FS local, S3, HDFS			
data: Data Layer Block, File, BlockQueue, Reader, Writer, Multiplexer, Streams, BlockPool (paging)	net: Network Layer (Binomial Tree) Broadcast, Reduce, AllReduce, Async-Send/Recv, Dispatcher Backends: <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td>mock</td> <td>tcp</td> <td>mpi</td> </tr> </table>	mock	tcp	mpi
mock	tcp	mpi		
io: Async File I/O borrowed from STXXL				
common: Common Tools Logger, Delegates, Math, ...	mem: Memory Limitation Allocators, Counting			

Thrift's Communication Abstraction



Thrill's Data Processing Pipelines



Weak-Scaling Benchmarks

WordCountCC – $h \cdot 49$ GiB 222 lines

- Reduce text files from CommonCrawl web corpus.

PageRank – $h \cdot 2.7$ GiB, $|E| \approx h \cdot 158$ M 410 lines

- Calculate PageRank using join of current ranks with outgoing links and reduce by contributions. 10 iterations.

TeraSort – $h \cdot 16$ GiB 141 lines

- Distributed (external) sorting of 100 byte random records.

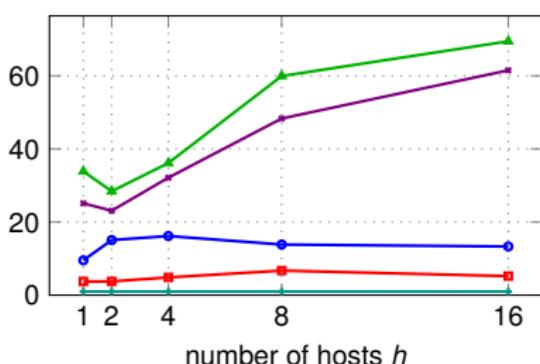
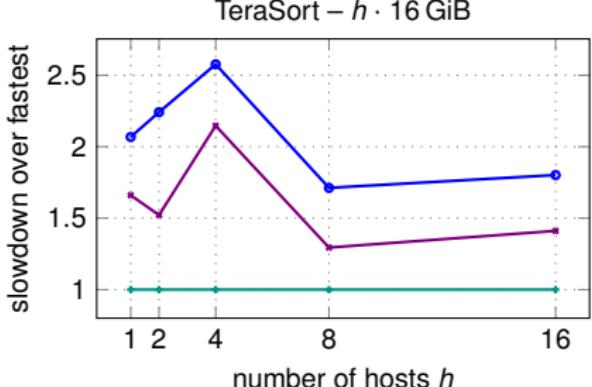
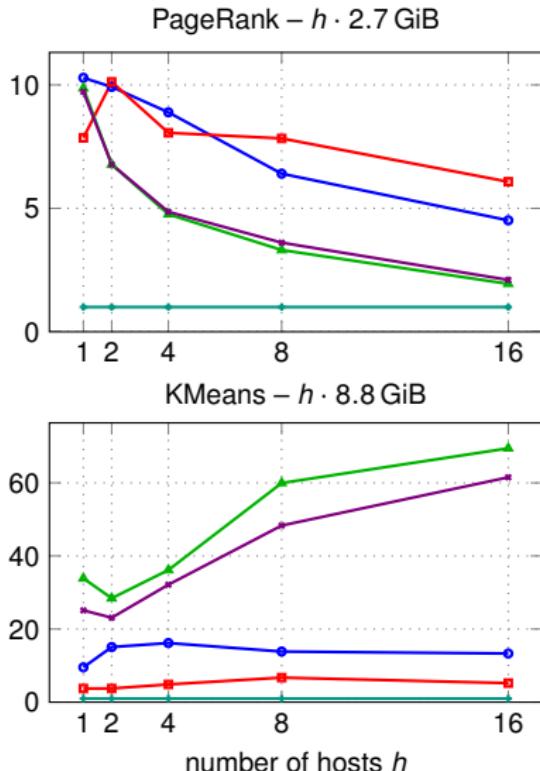
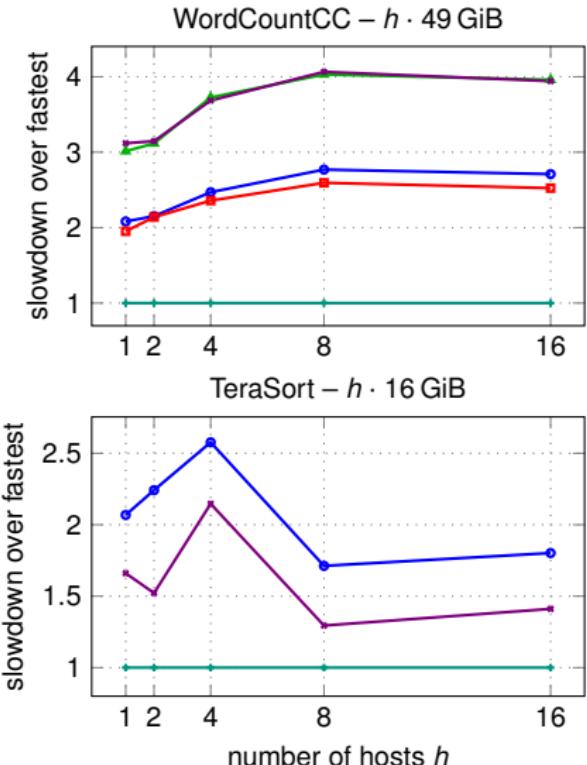
K-Means – $h \cdot 8.8$ GiB 357 lines

- Calculate K-Means clustering with 10 iterations.

Platform: $h \times$ r3.8xlarge systems on Amazon EC2 Cloud

- 32 cores, Intel Xeon E5-2670v2, 2.5 GHz clock, 244 GiB RAM, 2 x 320 GB local SSD disk, ≈ 400 MiB/s read/write Ethernet network ≈ 1000 MiB/s throughput, Ubuntu 16.04.

Experimental Results: Slowdowns



—●— Spark (Java) —■— Spark (Scala) —▲— Flink (Java) —●— Flink (Scala) —◆— Thrill

Thoughts on the Architecture

Thrill's Sweet Spot

- C++ toolkit for **implementing** distributed algorithms **quickly**.
- Platform to **engineer** and evaluate distributed primitives.
- Efficient processing of **small items** and **pipelining** of primitives.
- Platform for implementing on-the-fly compiled queries?

Open Questions

- **Compile-time optimization** only – no run-time algorithm selection or (statistical) knowledge about the data.
- Assumes h identical hosts **constantly running**, (the old MPI/HPC way, Hadoop/Spark do block-level scheduling).
- **Memory management**
- **Predictability** and **scalability** to 1 million cores

Current and Future Work

- Open-Source at <http://project-thrill.org> and Github.
- High quality, **very modern C++14** code.
- A **K-Mean tutorial** is available!



Ideas for Future Work:

- Distributed rank()/select() and succinct bit vectors?
- Beyond DIA<T>? Graph<V, E>? Matrix<T>?
- **Fault tolerance** in the algorithms and **scalability** to large clusters.
- Stream/Micro-batch or Query processing framework?

Thank you for your attention!
Questions?