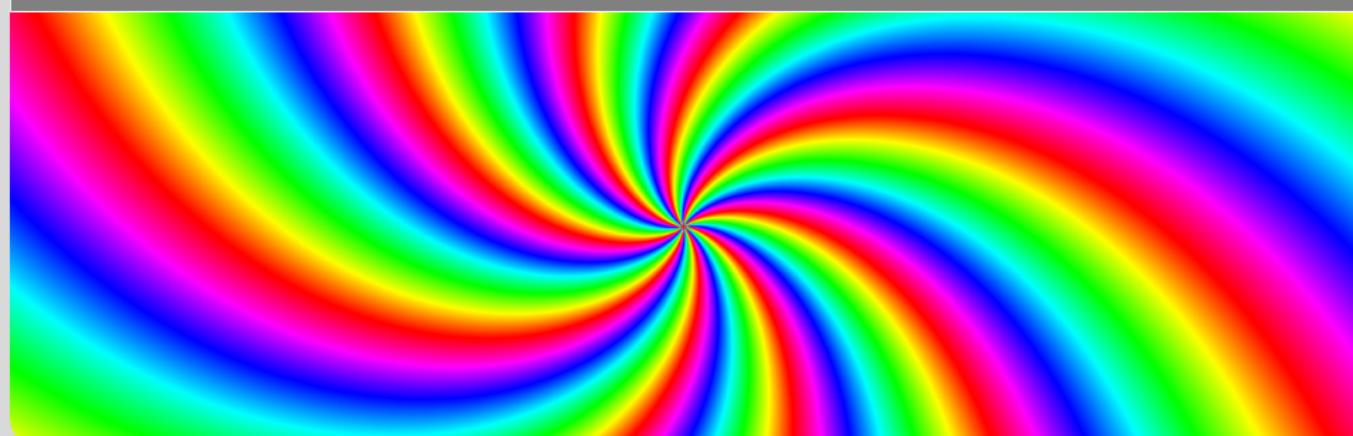


Thrill : High-Performance Algorithmic Distributed Batch Data Processing with C++

Timo Bingmann, Michael Axtmann, Peter Sanders, Sebastian Schlag, and 6 Students | 2017-09-01

INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMIC



Weak-Scaling Benchmarks

WordCountCC – $h \cdot 49$ GiB

- Reduce text files from CommonCrawl web corpus.

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

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

TeraSort – $h \cdot 16$ GiB

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

K-Means – $h \cdot 8.8$ GiB

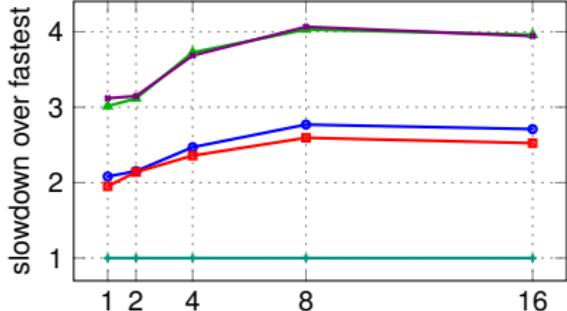
- Calculate K-Means clustering with 10 iterations.

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

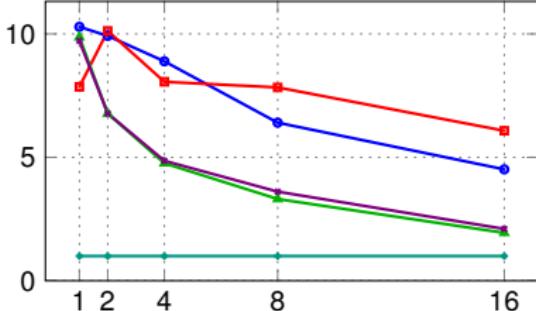
- 32 cores, Intel Xeon E5-2670v2, 2.5 GHz clock, 244 GiB RAM, 2 x 320 GB local SSD disk, ≈ 400 MiB/s read/write Ethernet network ≈ 1000 MiB/s throughput, Ubuntu 16.04.

Experimental Results: Slowdowns

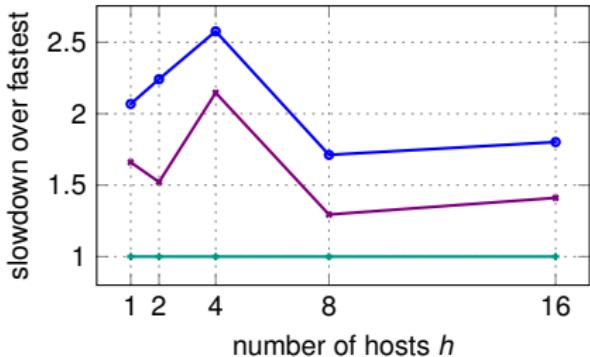
WordCountCC – $h \cdot 49$ GiB



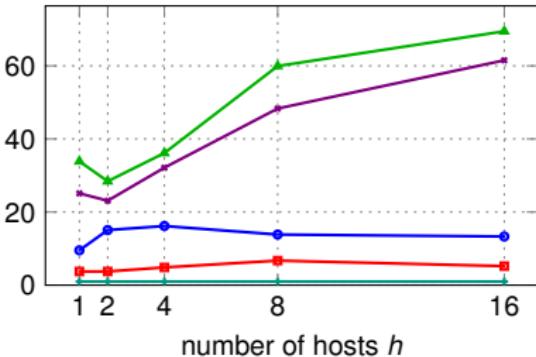
PageRank – $h \cdot 2.7$ GiB



TeraSort – $h \cdot 16$ GiB



KMeans – $h \cdot 8.8$ GiB



number of hosts h

number of hosts h

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

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

SA_i	LCP_i	$T_{SA_i \dots n}$
14	-	\$
9	0	a b d c c \$
2	1	a d c b c c b a b d c c \$
8	0	b a b d c c \$
1	2	b a d c b c c b a b d c c \$
5	1	b c c b a b d c c \$
10	1	b d c c \$
13	0	c \$
7	1	c b a b d c c \$
4	2	c b c c b a b d c c \$
12	1	c c \$
6	2	c c b a b d c c \$
0	0	d b a d c b c c b a b d c c \$
3	1	d c b c c b a b d c c \$
11	2	d c c \$

Suffix Sorting with DC3: Example

0 1 2 3 4 5 6 7 8 9 10

$$T = [\text{d } \boxed{\text{b}} \text{ a } \boxed{\text{c}} \text{ } \boxed{\text{b}} \text{ a } \boxed{\text{c}} \text{ } \boxed{\text{b}} \text{ d } \$ \$] = [t_i]_{i=0, \dots, n-1}$$

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

$$R = \boxed{1} \boxed{1} \boxed{2} \boxed{0} \boxed{0} \boxed{3} \$ \quad r_1 \quad r_4 \quad r_7 \quad r_2 \quad r_5 \quad r_8$$

$$\text{SA}_R = 3 \ 4 \ 0 \ 1 \ 2 \ 5 \ \$ \quad \text{ISA}_R = \boxed{2} \ \boxed{3} \ \boxed{4} \ \boxed{0} \ \boxed{1} \ \boxed{5} \ \$$$

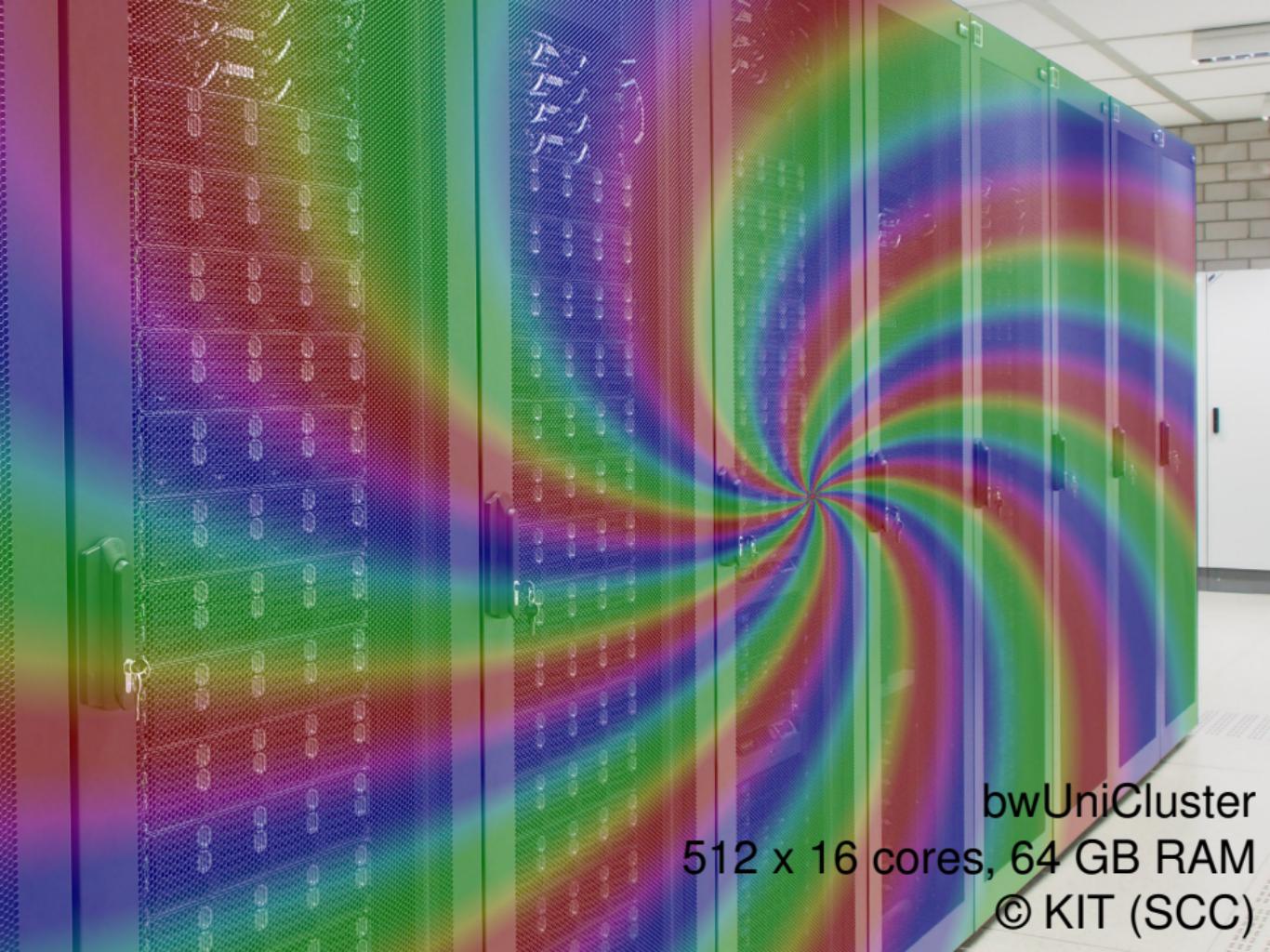
$$S_0 = [(d, b, \color{red}{2}, \color{blue}{0}, \color{green}{0}), (c, b, \color{red}{3}, \color{blue}{1}, \color{green}{3}), (c, b, \color{red}{4}, \color{blue}{5}, \color{green}{6})] \quad (t_i, t_{i+1}, r_{i+1}, r_{i+2}, i)$$

$$S_1 = [(2, b, \color{red}{0}, \color{blue}{1}), (3, b, \color{red}{1}, \color{blue}{4}), (4, b, \color{red}{5}, \color{blue}{7})] \quad (r_{i+1}, t_{i+1}, r_{i+2}, i+1)$$

$$S_2 = [(0, a, c, \color{red}{3}, \color{blue}{2}), (1, a, c, \color{red}{4}, \color{blue}{5}), (5, d, \$, \color{red}{6}, \color{blue}{8})] \quad (r_{i+2}, t_{i+2}, t'_{i+3}, r'_{i+4}, i+2)$$

$$\text{SA}_T = \text{Merge}(\text{Sort}(S_0), \text{Sort}(S_1), \text{Sort}(S_2))$$

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



bwUniCluster
512 x 16 cores, 64 GB RAM
© 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.

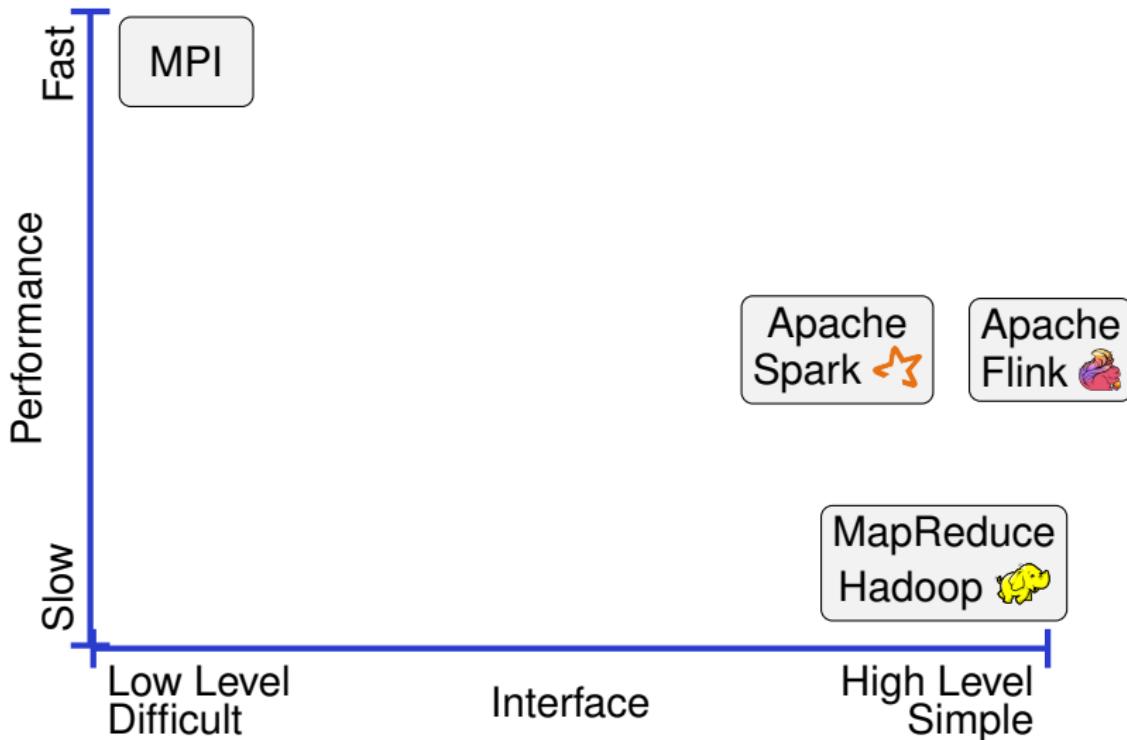
- Graph Processing

Google's Pregel, GraphLab , Giraph , GraphChi.

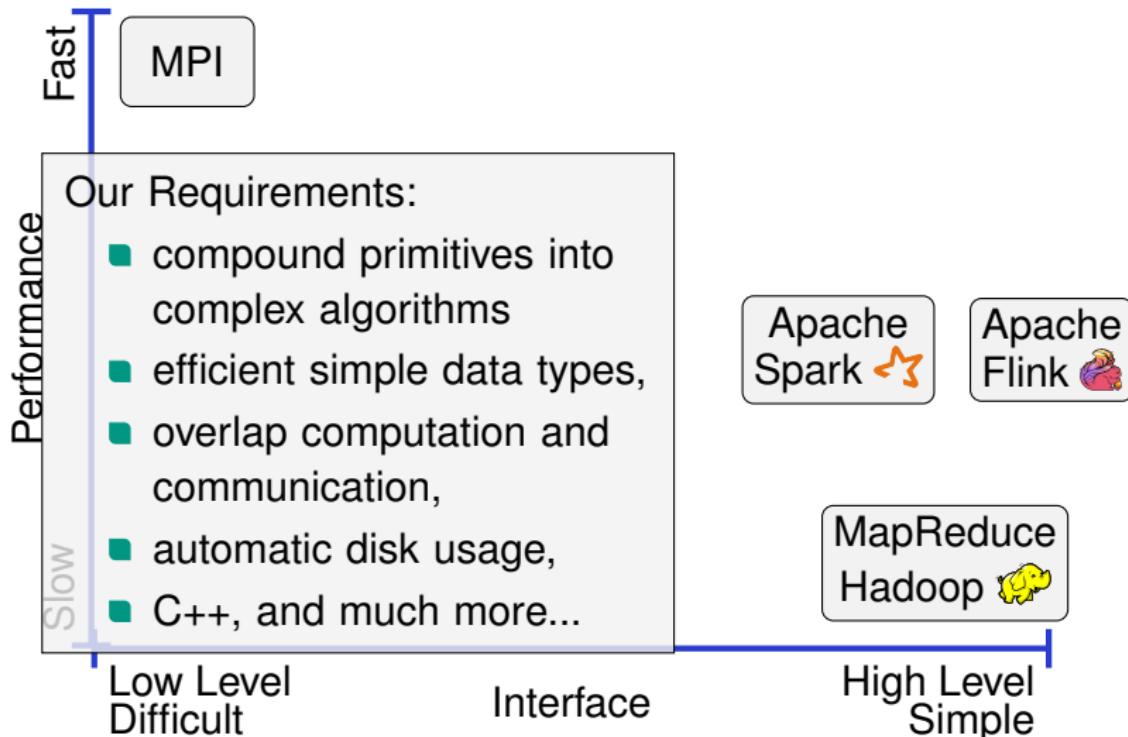
- Time-based Distributed Processing

Microsoft's Dryad, Microsoft's Naiad.

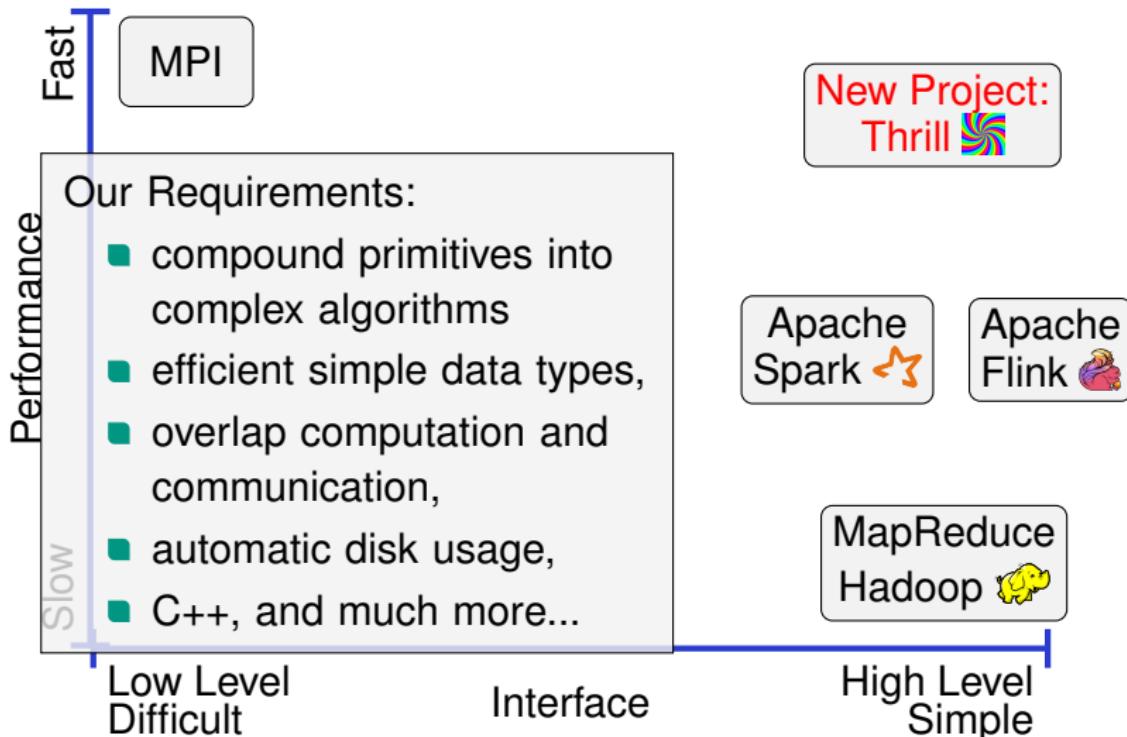
Big Data Batch Processing



Big Data Batch Processing



Big Data Batch Processing



Projektpraktikum: Verteilte Datenverarbeitung mit MapReduce

Timo Bingmann, Peter Sanders und Sebastian Schlag | 21. Oktober 2014 @ PdF Vorstellung

INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMIC



Thrill's Design Goals

- An easy way to program distributed algorithms in C++.
- Distributed arrays of small items (characters or integers).
- High-performance, parallelized C++ operations.
- Locality-aware, in-memory computation.
- Transparently use disk if needed
 - ⇒ external memory or cache-oblivious algorithms.
- Avoid all unnecessary round trips of data to memory (or disk).
- Optimize chaining of local operations.

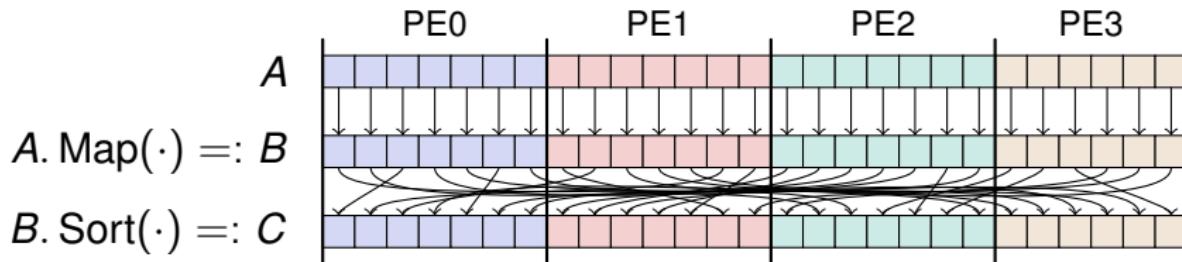
Current Status:

- open-source prototype at <http://github.com/thrill/thrill>.

Distributed Immutable Array (DIA)

- User Programmer's View:

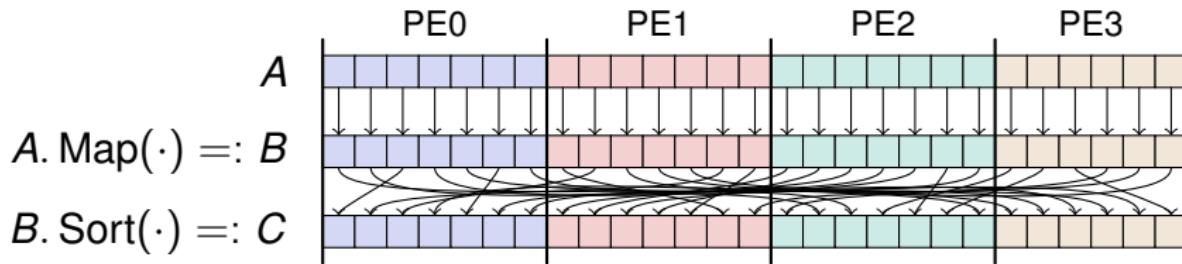
- $\text{DIA} < T >$ = result of an operation (local or distributed).
- Model: distributed array of items T on the cluster
- Cannot access items directly, instead use transformations and actions.



Distributed Immutable Array (DIA)

- User Programmer's View:

- $\text{DIA} < T >$ = 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.



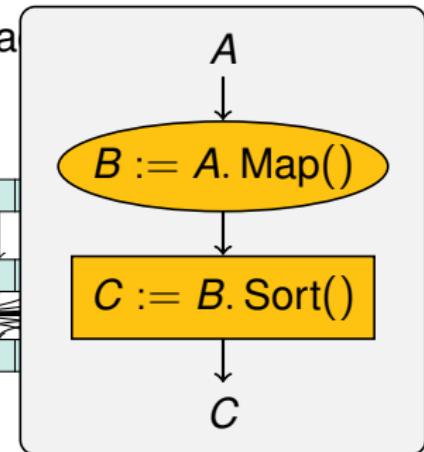
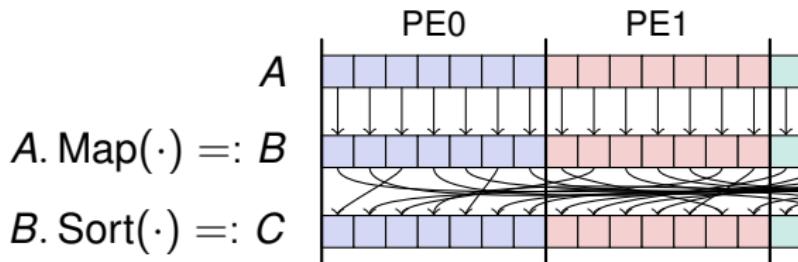
- Framework Designer's View:

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable. \implies build data-flow graph.
- $\text{DIA} < T >$ = 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 and actions.



■ Framework Designer's View:

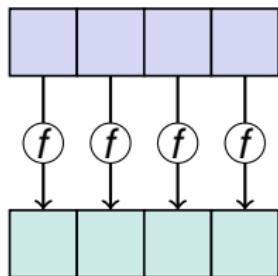
- Goals: distribute work, optimize execution on cluster, add redundancy where applicable. \implies build data-flow graph.
- DIA<T> = chain of computation items
- Let distributed operations choose “materialization”.

List of Primitives (Excerpt)

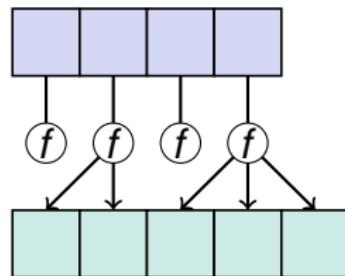
- Local Operations (**LOp**): input is **one item**, output ≥ 0 items.
Map(), **Filter()**, **FlatMap()**.
- Distributed Operations (**DOP**): input is a **DIA**, output is a **DIA**.
 - Sort()** Sort a DIA using comparisons.
 - ReduceBy()** Shuffle with Key Extractor, Hasher, and associative Reducer.
 - 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)

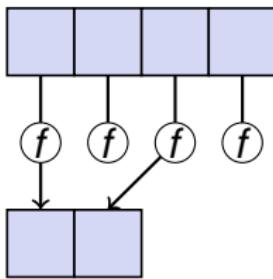
Map(f) : $\langle A \rangle \rightarrow \langle B \rangle$
 $f : A \rightarrow B$



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



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



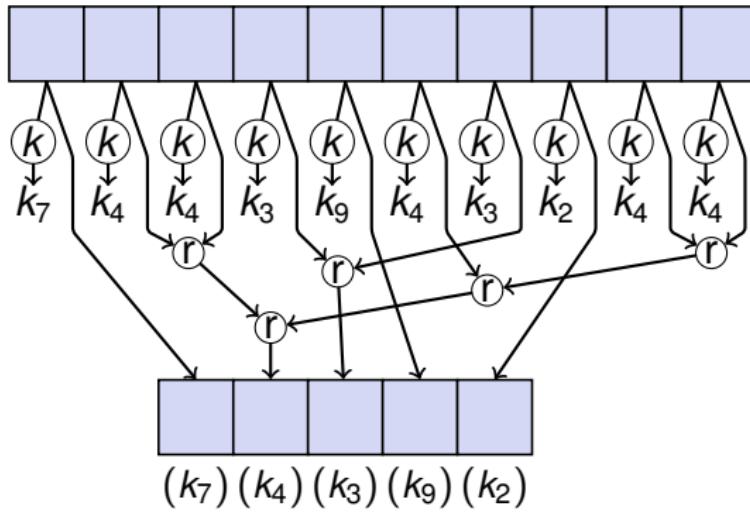
Currently: no rebalancing during LOps.

DOps: ReduceByKey

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

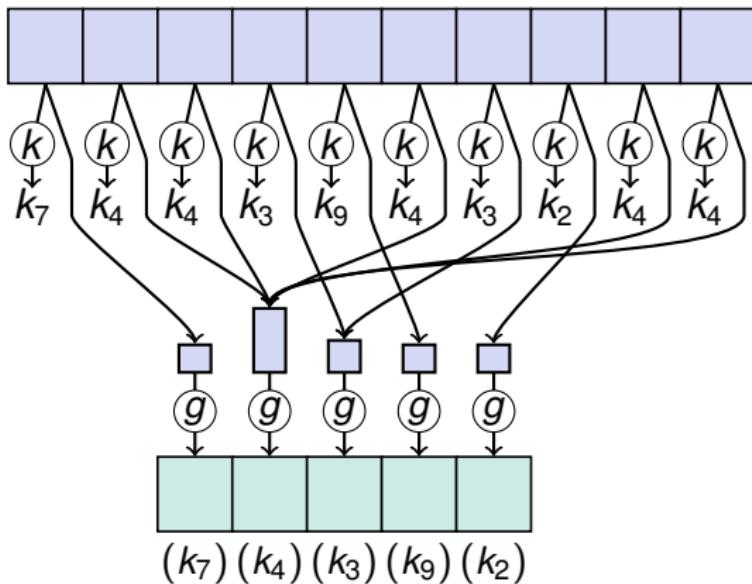
$k : A \rightarrow K$ key extractor

$r : A \times A \rightarrow A$ 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



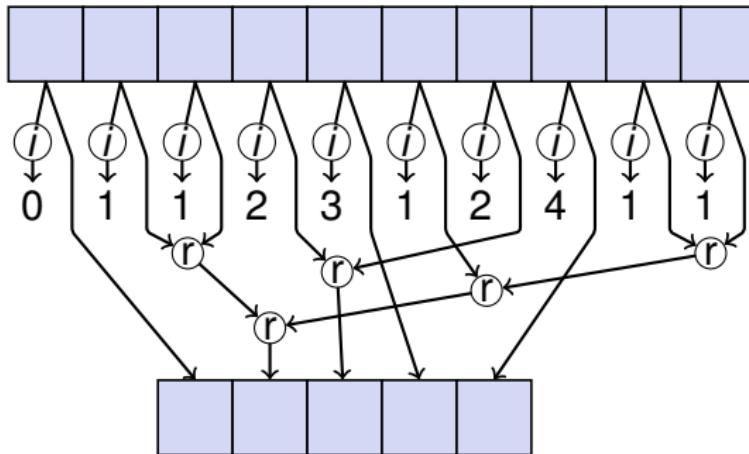
DOps: ReduceToIndex

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



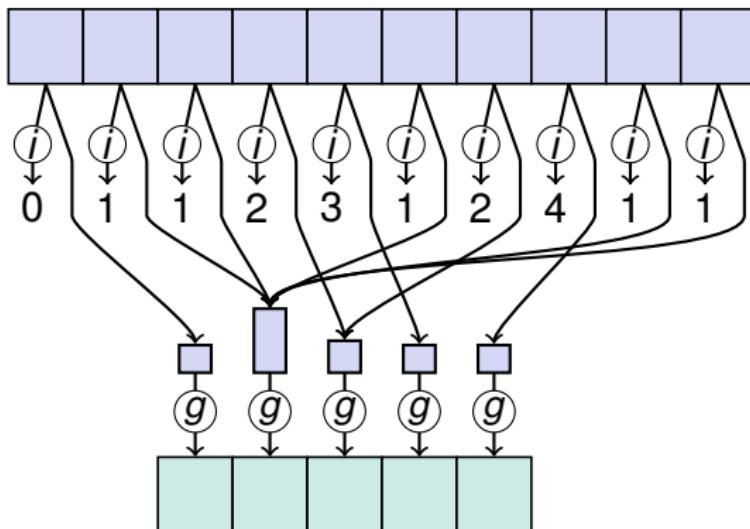
DOps: GroupToIndex

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

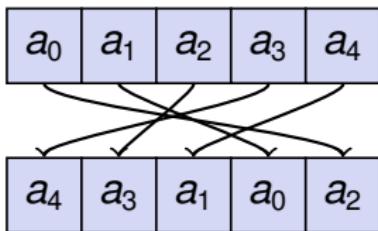
$g : \text{iterable}(A) \rightarrow B$ group function



DOps: Sort and Merge

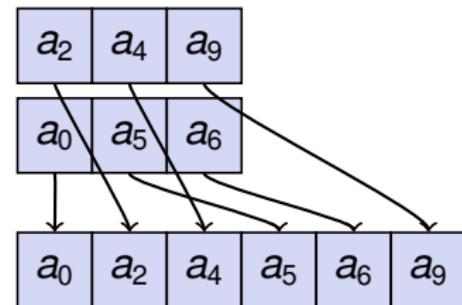
Sort(o) : $\langle A \rangle \rightarrow \langle A \rangle$

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



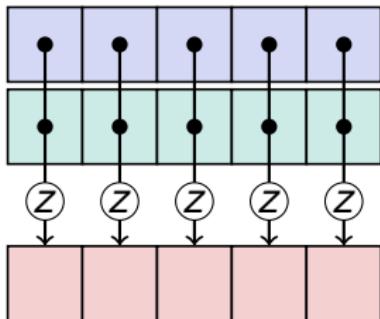
Merge(o) : $\langle A \rangle \times \langle A \rangle \cdots \rightarrow \langle A \rangle$

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

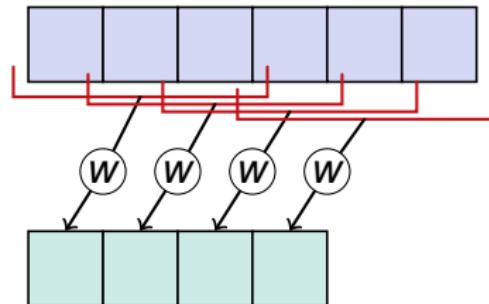


DOps: Zip and Window

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



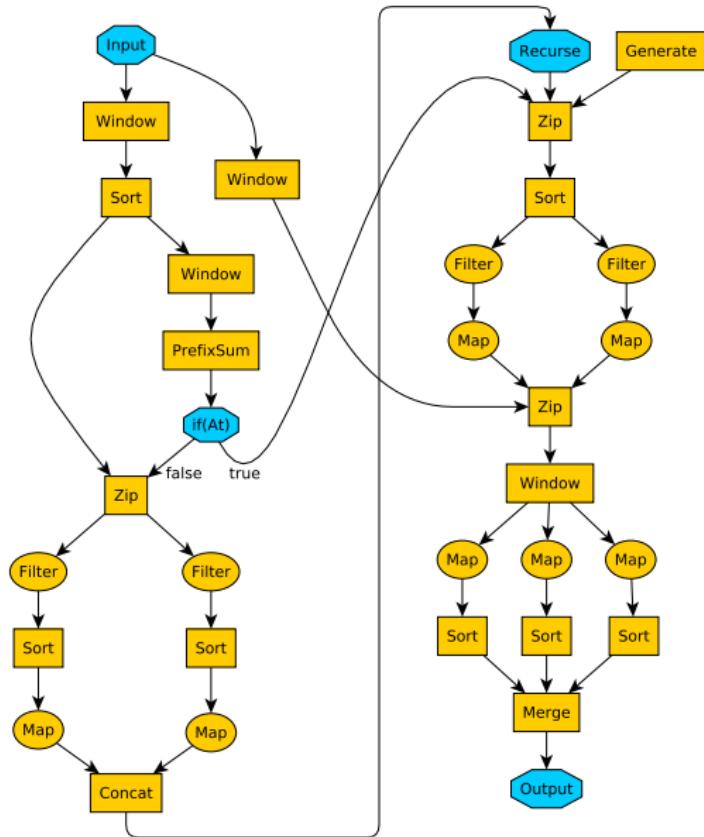
Window(k, w) : $\langle A \rangle \rightarrow \langle B \rangle$
 $k \in \mathbb{N}$ window size
 $w : A^k \rightarrow B$ window function



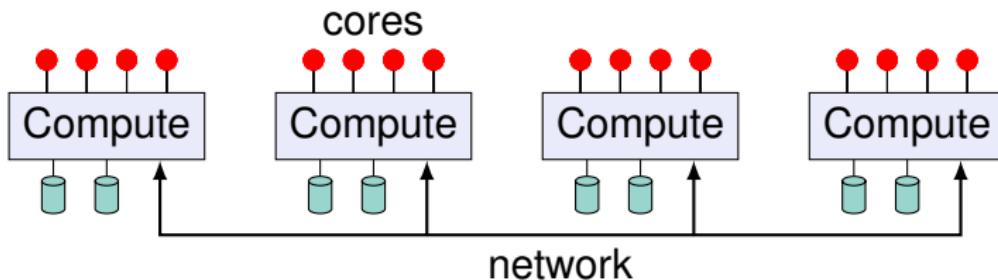
Example: WordCount in Thrill

```
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             Split(line, ' ', [&](std::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 }
```

A Suffix Sorting Algorithm: DC3



Execution on Cluster

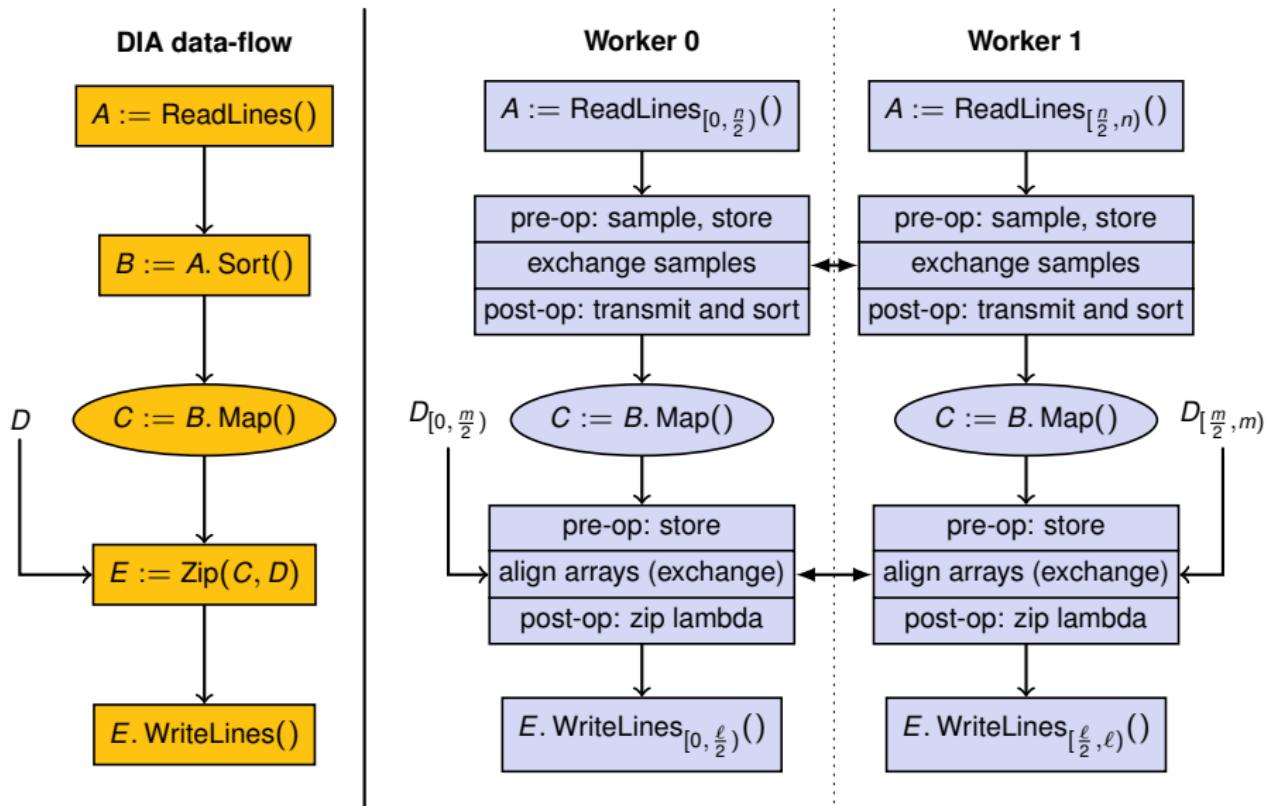


- Compile program into **one binary**, running on all hosts.
- **Collective coordination** of work on compute hosts, like MPI.
- **Control flow** is decided on by using C++ statements.
- Runs on MPI HPC clusters and on Amazon's EC2 cloud.

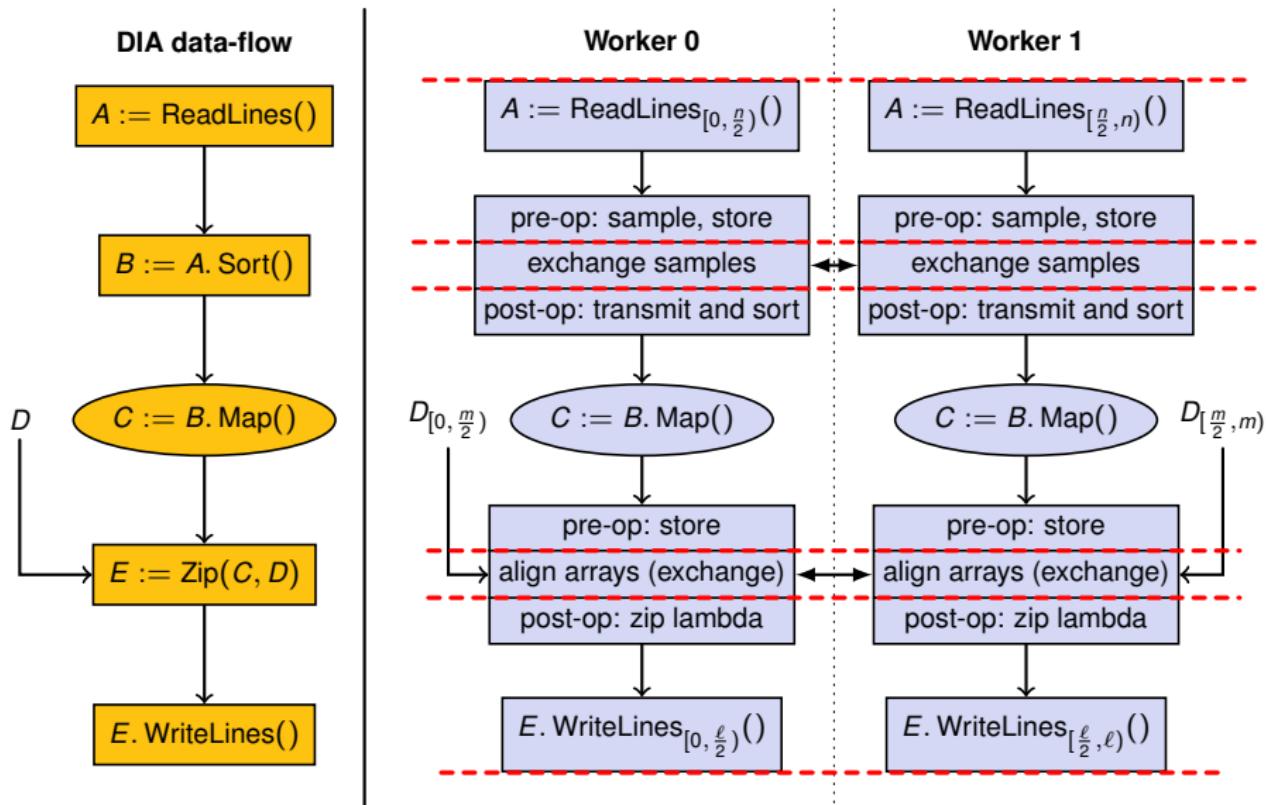
Layers of Thrill

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"><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			

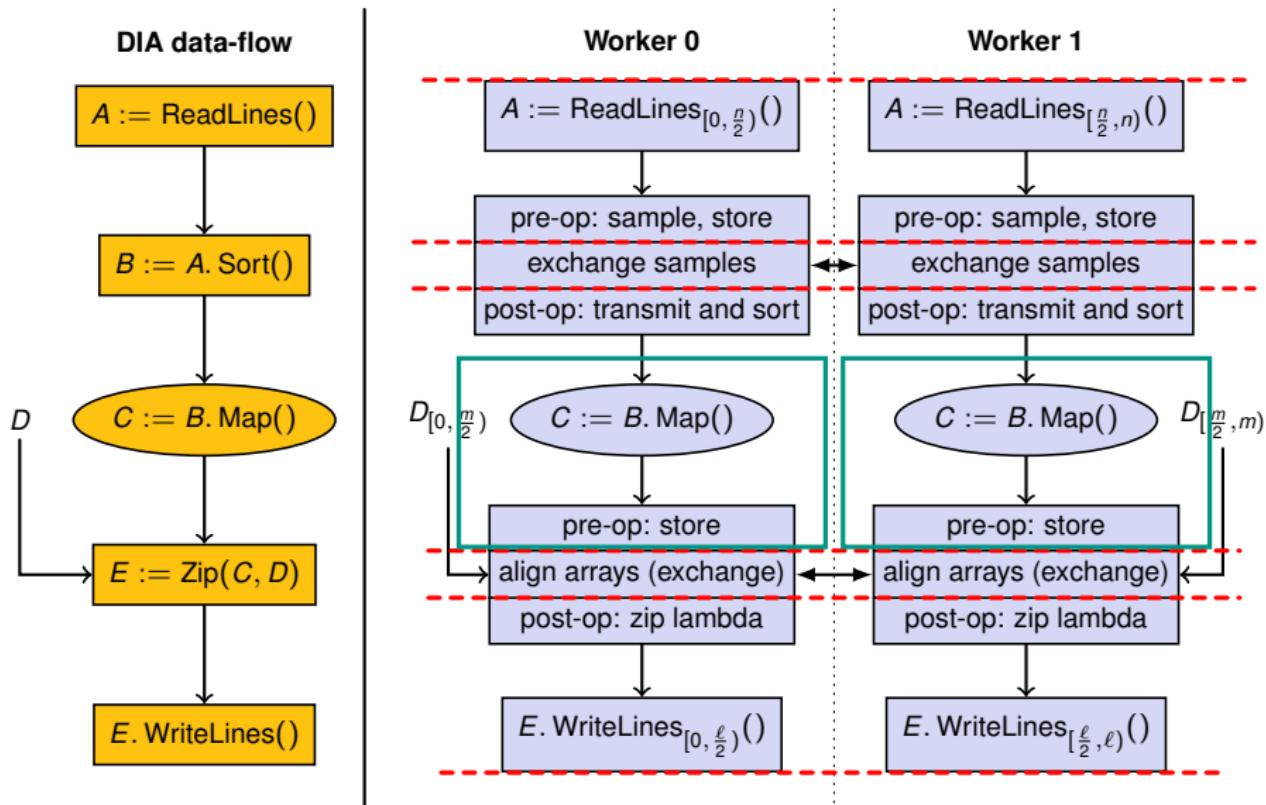
Mapping Data-Flow Nodes to Cluster



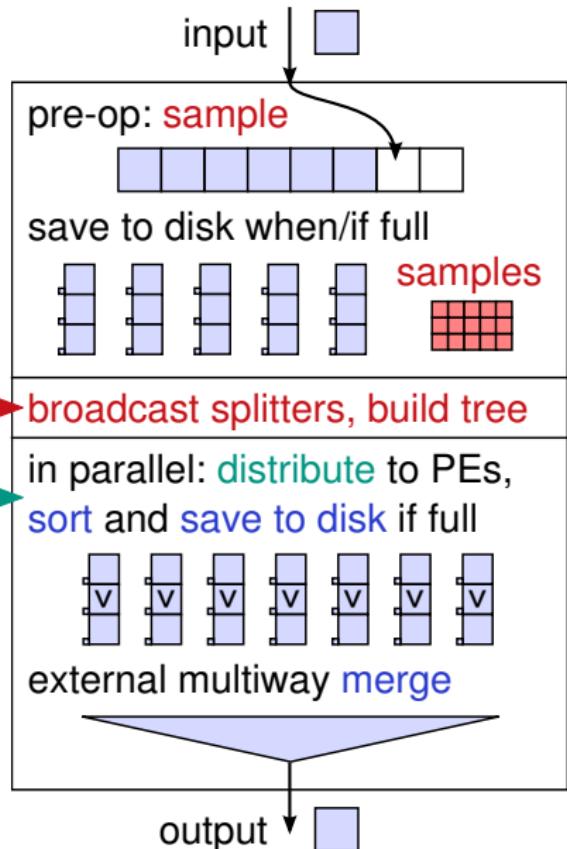
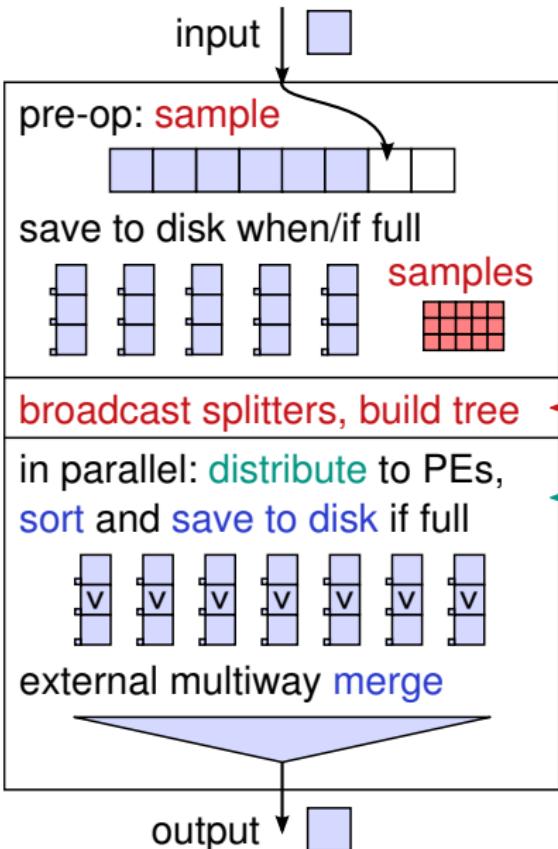
Mapping Data-Flow Nodes to Cluster



Mapping Data-Flow Nodes to Cluster



Sorting DOp



Weak-Scaling Benchmarks

WordCountCC – $h \cdot 49$ GiB

- Reduce text files from CommonCrawl web corpus.

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

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

TeraSort – $h \cdot 16$ GiB

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

K-Means – $h \cdot 8.8$ GiB

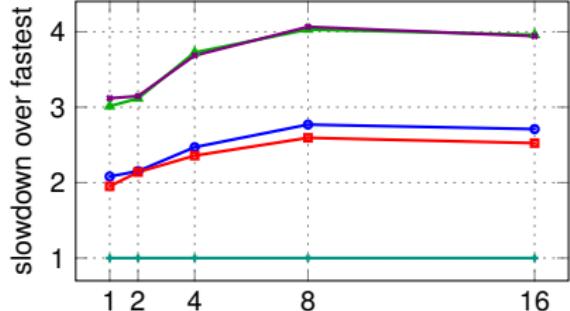
- Calculate K-Means clustering with 10 iterations.

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

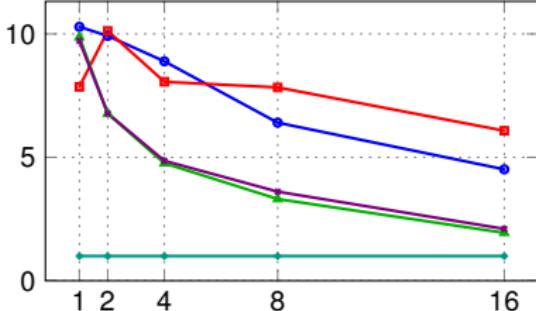
- 32 cores, Intel Xeon E5-2670v2, 2.5 GHz clock, 244 GiB RAM, 2 x 320 GB local SSD disk, ≈ 400 MiB/s read/write Ethernet network ≈ 1000 MiB/s throughput, Ubuntu 16.04.

Experimental Results: Slowdowns

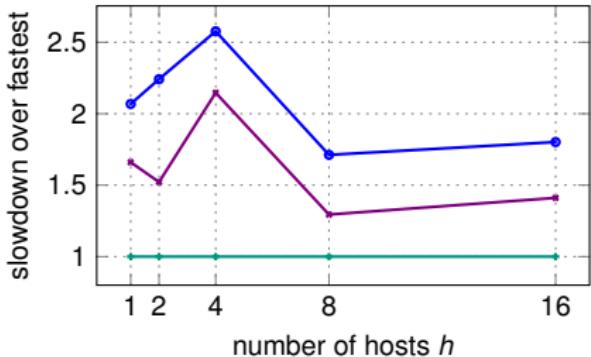
WordCountCC – $h \cdot 49$ GiB



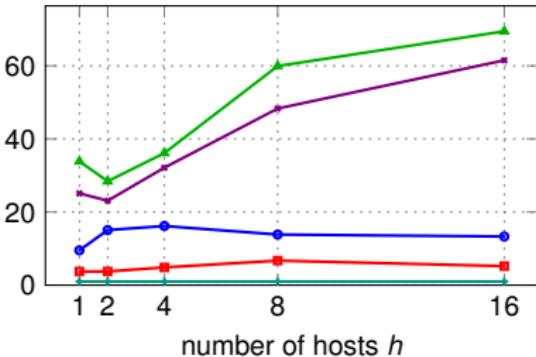
PageRank – $h \cdot 2.7$ GiB



TeraSort – $h \cdot 16$ GiB



KMeans – $h \cdot 8.8$ GiB



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

Thoughts on the Architecture

Thrill's Sweet Spot

- C++ toolkit for **experimenting** with distributed algorithms.
- 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 hosts **constantly running**, (the old MPI/HPC way of things, Hadoop/Spark do block-level scheduling).
- Memory management

K-Means Tutorial

Thrill 0.1

Main Page Related Pages Modules Namespaces Classes Files Examples Search

▼ Thrill

- ▼ Thrill Documentation Overview
 - ▶ Getting Started
 - ▼ K-Means Tutorial

Step 1: Generate Random Points

Welcome to the first step in the Thrill k-means tutorial. This tutorial will show how to implement the k-means clustering algorithm (Lloyd's algorithm) in Thrill.

The algorithm works as follows: Given a set of d-dimensional points, select k initial cluster center points at random. Then attempt to improve the centers by iteratively calculating new centers. This is done by classifying all points and associating them with their nearest center, and then taking the mean of all points associated to one cluster as the new center. This will be repeated a constant number of iterations.



We will implement this algorithm in Thrill, and only work with two-dimensional points for simplicity. Furthermore, we will hard-code many constants to make the code easier to understand.

In this step 1, let us start with generating random 2-dimensional points and outputting them for debugging.

We first need a Point class to represent the points. We may add some calculation functions to it later on.

```
//! A 2-dimensional point with double precision
struct Point {
    //! point coordinates
    double x, y;
};
```

For outputting the Point class, we need to add an operator `<<` for `std::ostream`, which is the standard way for

Thrill Documentation Overview K-Means Tutorial Generated on Tue Sep 20 2016 19:24:29 for Thrill by doxygen 1.8.5

Current and Future Work

- Open-Source at <http://project-thrill.org> and Github.
- High quality, **very modern C++14** code.

Ideas for Future Work:

- Distributed rank()/select() and wavelet tree construction.
- Beyond DIA<T>? Graph<V,E>? DenseMatrix<T>?
- Fault tolerance? Go from p to $p - 1$ workers?
- Communication efficient distributed operations for Thrill.
- Distributed functional programming language on top of Thrill.

Thank you for your attention!
Questions?