

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

Timo Bingmann, Michael Axtmann, Peter Sanders, Sebastian Schlag, and 6 Students | 2016-12-06

INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMICS

Abstract

We present on-going work on a new distributed Big Data processing framework called Thrill. It is a C++ framework consisting of a set of basic scalable algorithmic primitives like mapping, reducing, sorting, merging, joining, and additional MPI-like collectives. This set of primitives goes beyond traditional Map/Reduce and can be combined into larger more complex algorithms, such as WordCount, PageRank, k-means clustering, and suffix sorting. These complex algorithms can then be run on very large inputs using a distributed computing cluster. Among the main design goals of Thrill is to lose very little performance when composing primitives such that small data types are well supported. Thrill thus raises the questions of a) how to design algorithms using the scalable primitives, b) whether additional primitives should be added, and c) if one can improve the existing ones using new ideas to reduce communication volume and latency.

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

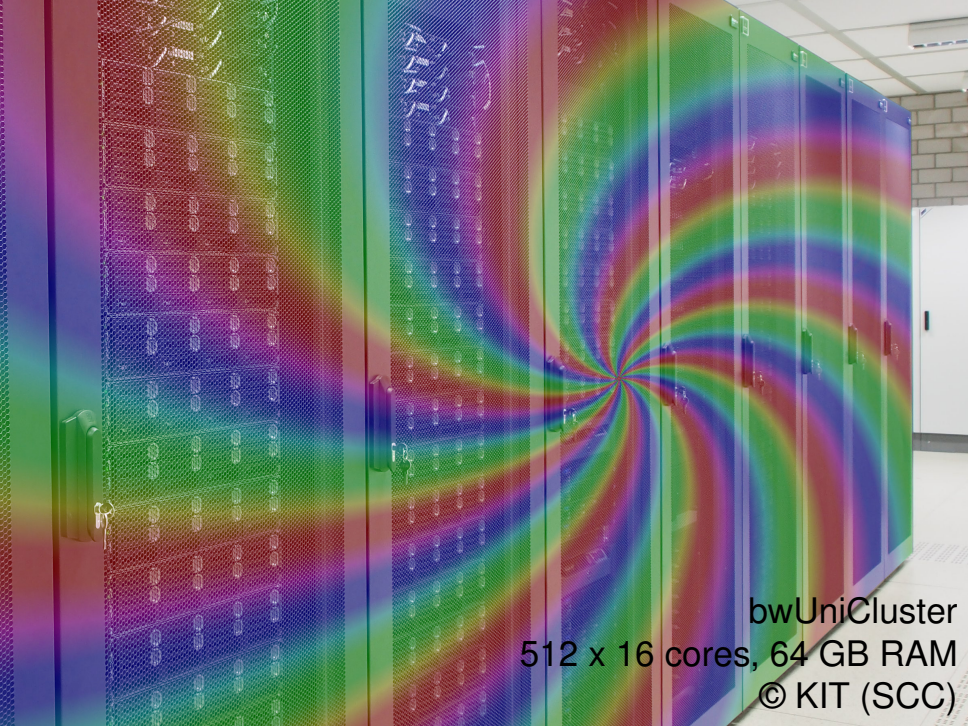
Example $T = [\text{dbadcbccbabdccc\$}]$

i	T_i
0	d b a d c b c c b a b d c c \$
1	b a d c b c c b a b d c c \$
2	a d c b c c b a b d c c \$
3	d c b c c b a b d c c \$
4	c b c c b a b d c c \$
5	b c c b a b d c c \$
6	c c b a b d c c \$
7	c b a b d c c \$
8	b a b d c c \$
9	a b d c c \$
10	b d c c \$
11	d c c \$
12	c c \$
13	c \$
14	\$

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14

Example $T = [\text{dbadc}bc\text{c}b\text{ab}d\text{c}c\$]$

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 \$



bwUniCluster
512 x 16 cores, 64 GB RAM
© KIT (SCC)




Suffix Sorting with DC3: Example

	0	1	2	3	4	5	6	7	8	9	10						
$T =$	d	b	a	c	b	a	c	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					
$R =$		1	1	2	0	0	3		\$		r_1	r_4	r_7	r_2	r_5	r_8	
$SA_R =$		3	4	0	1	2	5		\$	$ISA_R =$	2	3	4	0	1	5	\$
$S_0 =$		[(d, b,	2, 0, 0),	(c, b,	3, 1, 3),	(c, b,	4, 5, 6)]										$(t_i, t_{i+1}, r_{i+1}, r_{i+2}, i)$
$S_1 =$		[(2, b,	0, 1),	(3, b,	1, 4),	(4, b,	5, 7)]										$(r_{i+1}, t_{i+1}, r_{i+2}, i+1)$
$S_2 =$		[(0, a, c,	3, 2),	(1, a, c,	4, 5),	(5, d, \$,	6, 8)]										$(r_{i+2}, t_{i+2}, t'_{i+3}, r'_{i+4}, i+2)$
$SA_T =$		Merge(Sort(S_0),	Sort(S_1),	Sort(S_2))													

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

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.

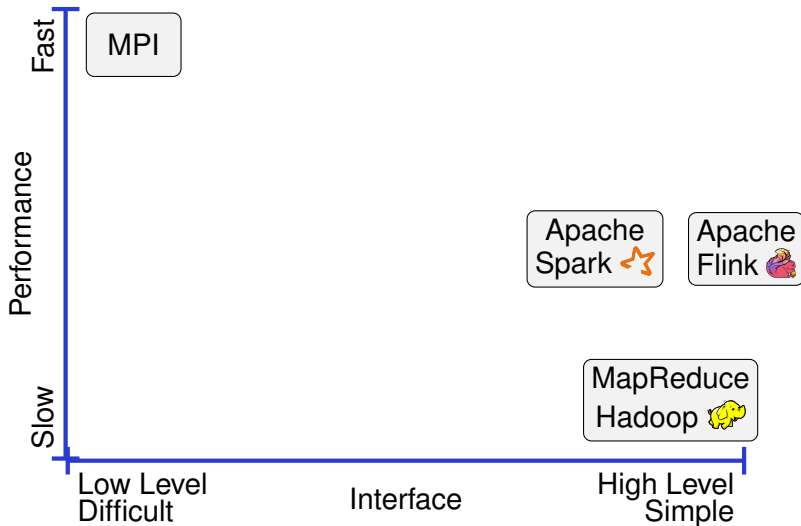
■ Graph Processing

Google's Pregel, GraphLab , Giraph , GraphChi.

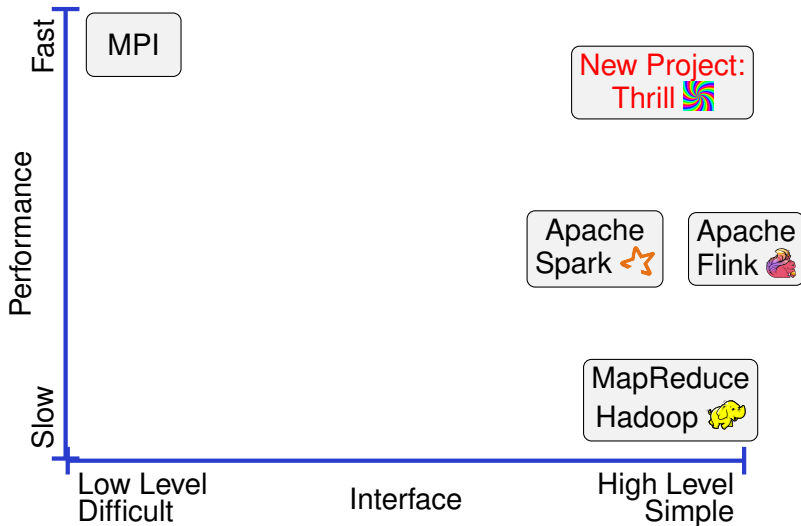
■ Time-based Distributed Processing

Microsoft's Dryad, Microsoft's Naiad.

Big Data Batch Processing



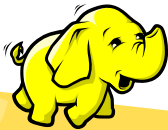
Big Data Batch Processing



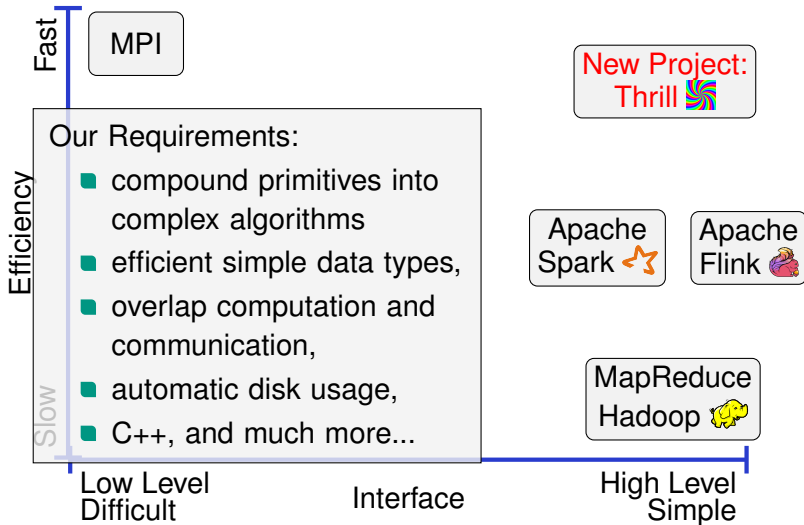
Projektpraktikum: Verteilte Datenverarbeitung mit MapReduce

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

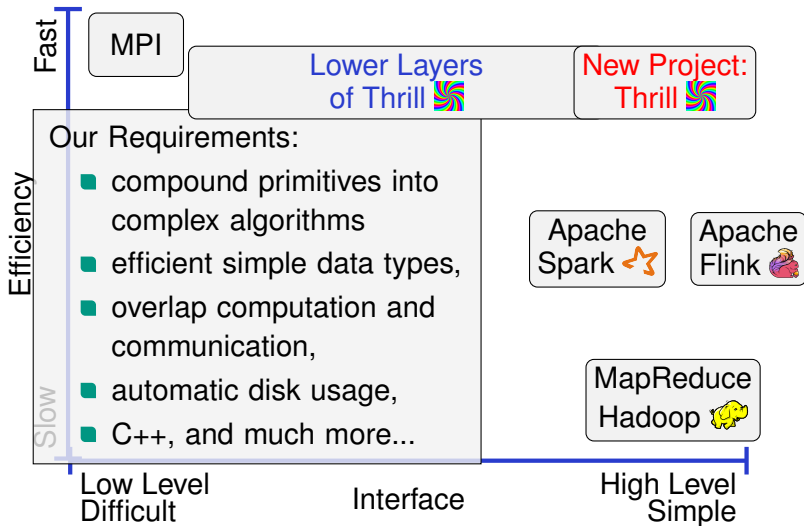
INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMICS



Big Data Batch Processing



Big Data Batch Processing



Why another Big Data Framework?

Sorting Records	Hadoop	Spark	Spark
Data Size	102.5 TB	100 TB	1000 TB
Elapsed Time	72 mins	23 mins	234 mins
# Hosts	2100	206	190
# Cores	50400	6592	6080
# Reducers	10 000	29 000	250 000
Rate	1.42 TB/min	4.27 TB/min	4.27 TB/min
Rate/host	11.2 MB/sec	345 MB/sec	375 MB/sec
Daytona Rules	Yes	Yes	No
Environment	dedicated	EC2 (i2.8xlarge)	

source: <http://databricks.com/blog/2014/10/10/spark-petabyte-sort.html>

EC2 (i2.8xlarge): 32 core each, 244 GiB RAM,
8 x 800 GiB SSD ($\approx 8 \times 400$ MB/s), 10 GBit Ethernet (≈ 800 MB/s).

Thrill's Design Goals

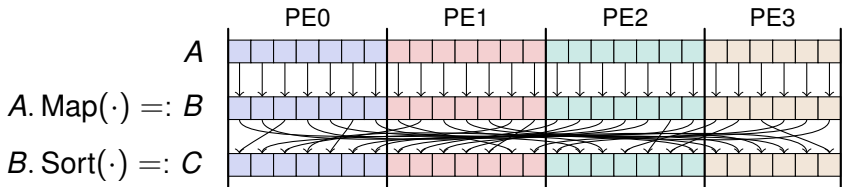
- **An easy way** to program **distributed** algorithms in C++.
- Distributed **arrays** of small items (characters or integers).
- **High-performance, parallelized C++** operations.
- **Locality-aware, in-memory** computation.
- **Transparently use disk** if needed
⇒ external memory or cache-oblivious 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>.

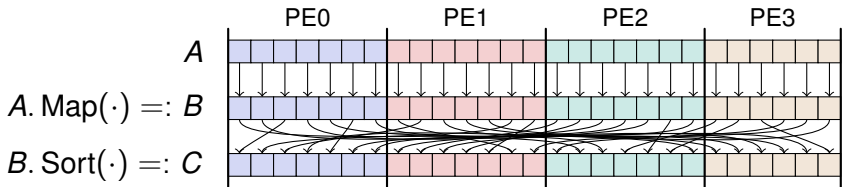
Distributed Immutable Array (DIA)

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

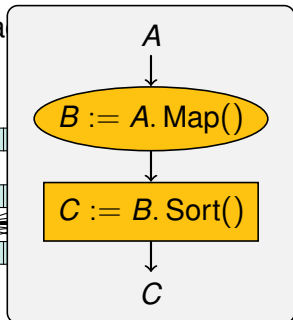
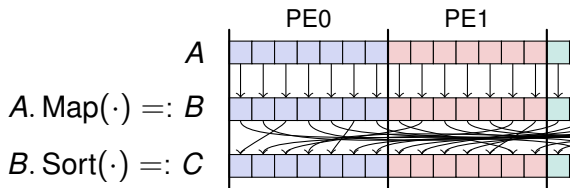


- Framework Designer's View:
 - Goals: distribute work, optimize execution on cluster, add redundancy where applicable. \implies **build data-flow graph**.
 - $\text{DIA}\langle T \rangle = \text{chain of computation items}$
 - Let distributed operations choose “materialization”.

Distributed Immutable Array (DIA)

User Programmer's View:

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



Framework Designer's View:

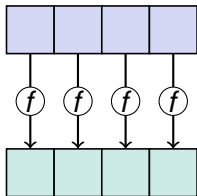
- Goals: distribute work, optimize execution on cluster, add redundancy where applicable. \implies **build data-flow graph**.
- $DIA\langle T \rangle$ = **chain of computation items**
- Let distributed operations choose “materialization”.

List of Primitives (Excerpt)

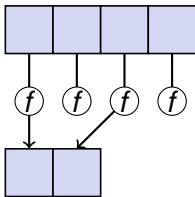
- Local Operations (LOp): input is **one item**, output ≥ 0 items.
`Map()`, `Filter()`, `FlatMap()`.
- Distributed Operations (DOp): input is a **DIA**, output is a **DIA**.
 - `Sort()` Sort a DIA using comparisons.
 - `ReduceBy()` Shuffle with Key Extractor, Hasher, and associative Reducer.
 - `GroupBy()`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

Local Operations (LOps)

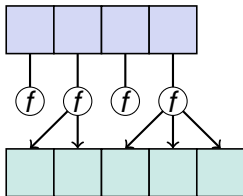
$$\text{Map}(f) : \langle A \rangle \rightarrow \langle B \rangle$$
$$f : A \rightarrow B$$



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



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



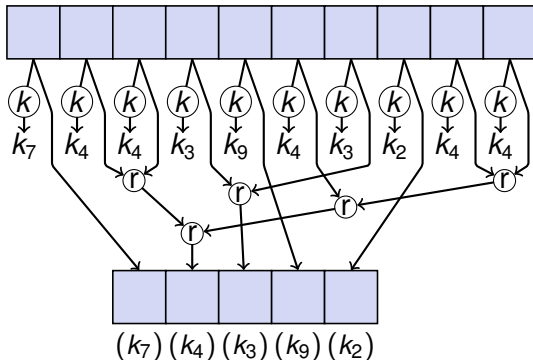
Currently: no rebalancing during LOps.

DOps: ReduceByKey

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

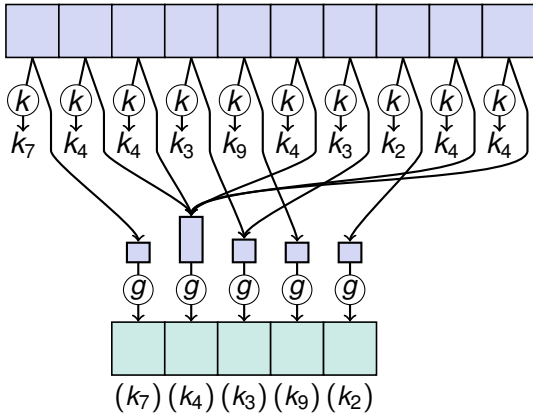
$k : A \rightarrow K$ key extractor

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



DOps: GroupByKey

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



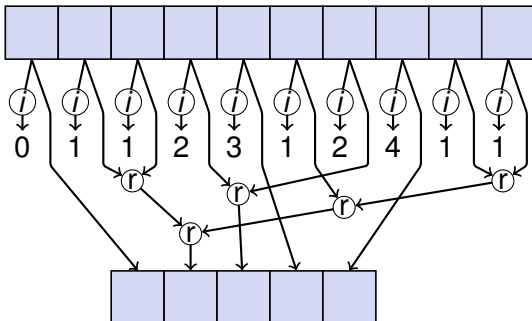
DOps: ReduceToIndex

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



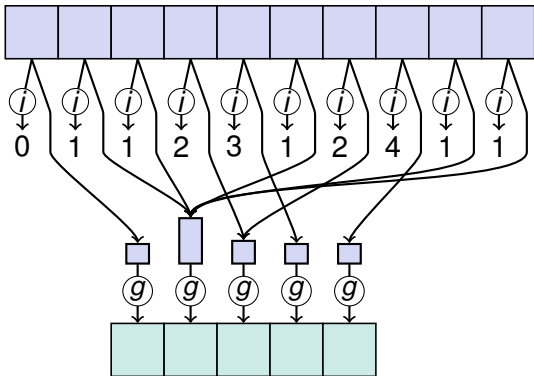
DOps: GroupToIndex

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

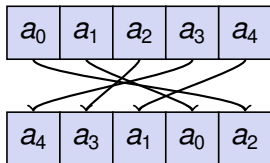
$g : \text{iterable}(A) \rightarrow B$ group function



DOps: Sort and Merge

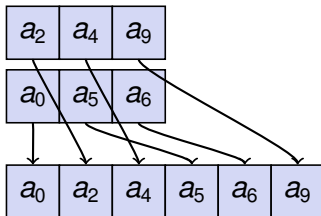
Sort(o) : $\langle A \rangle \rightarrow \langle A \rangle$

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



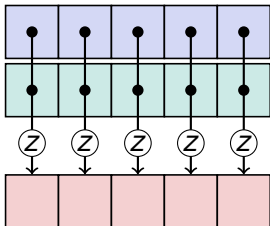
Merge(o) : $\langle A \rangle \times \langle A \rangle \cdots \rightarrow \langle A \rangle$

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

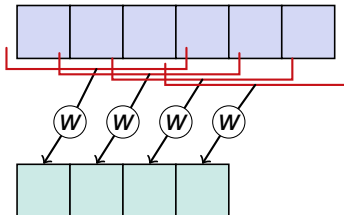


DOps: Zip and Window

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



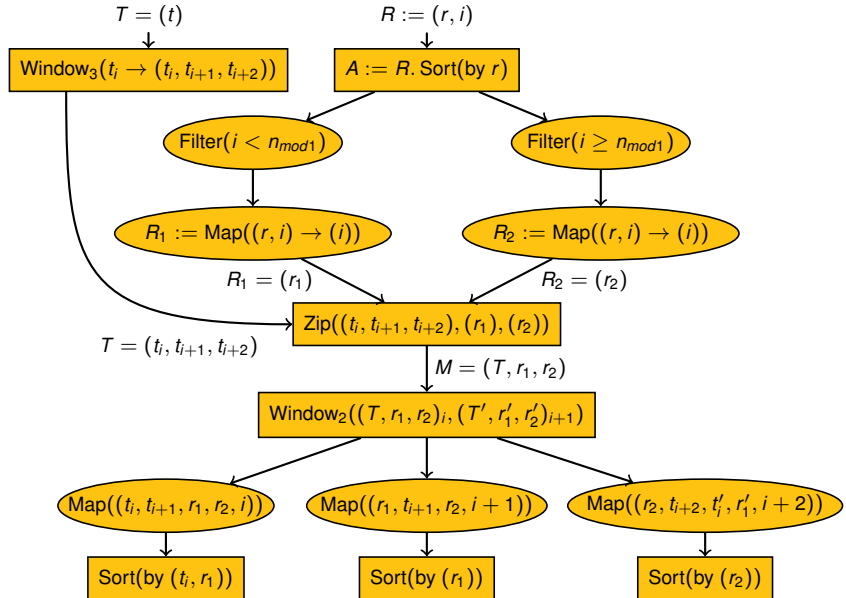
Window(k, w) : $\langle A \rangle \rightarrow \langle B \rangle$
 $k \in \mathbb{N}$ window size
 $w : A^k \rightarrow B$ window function



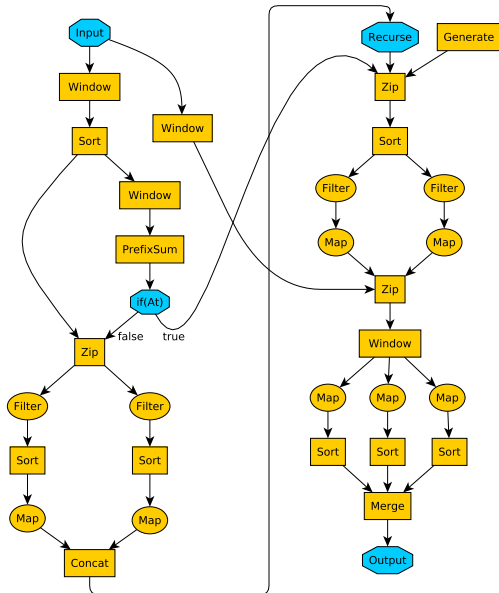
Example: WordCount in Thrill

```
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>(<br>
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     }); // DIA<Pair>
10 word_pairs.ReduceByKey(<br>
11     // key extractor: the word string
12     [](const Pair& p) { return p.first; },
13     // commutative reduction: add counters
14     [](const Pair& a, const Pair& b) {
15         return Pair(a.first, a.second + b.second);
16     }) // DIA<Pair>
17     .Map([](const Pair& p) {
18         return p.first + ": " + std::to_string(p.second); })
19     .WriteLines(output); // DIA<std::string>
20 }
```

Exert of DC3's Data-Flow Graph

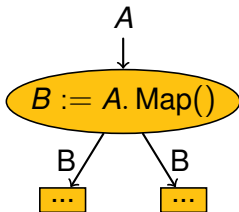


A Suffix Sorting Algorithm: DC3

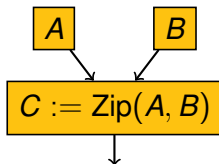
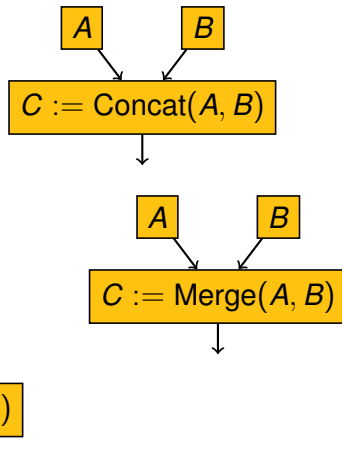


Structure of Data-Flow Graph

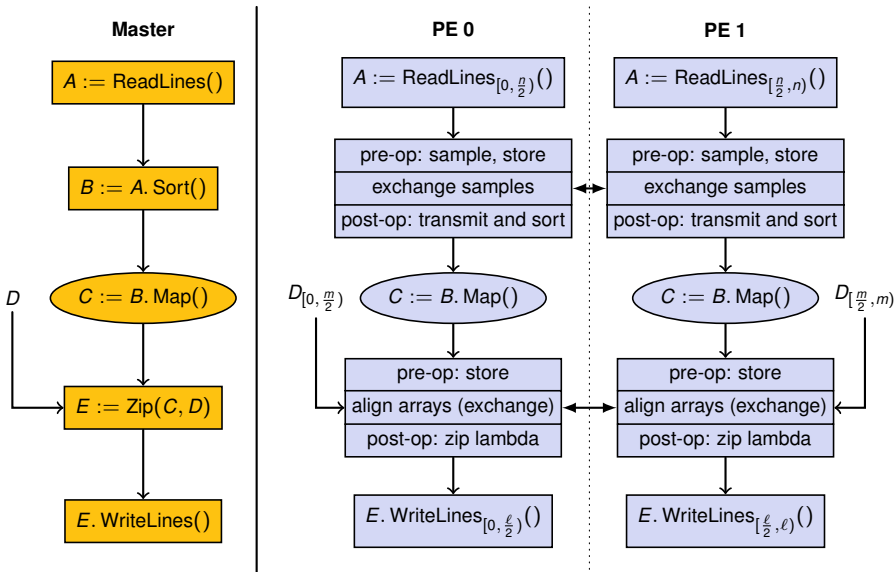
Any node can fork DIAs.



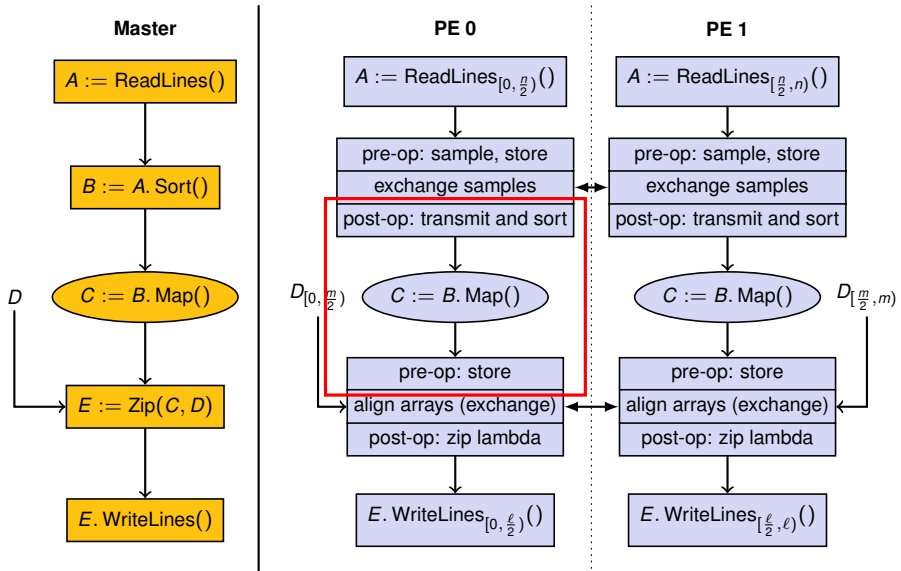
Combines are DOPs.



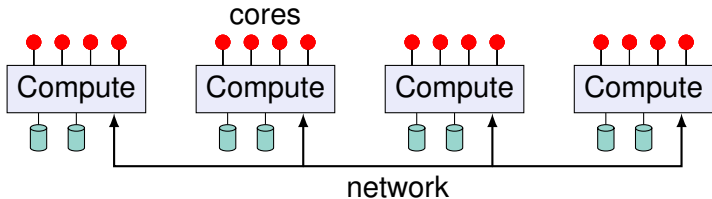
Mapping Data-Flow Nodes to Cluster



Mapping Data-Flow Nodes to Cluster

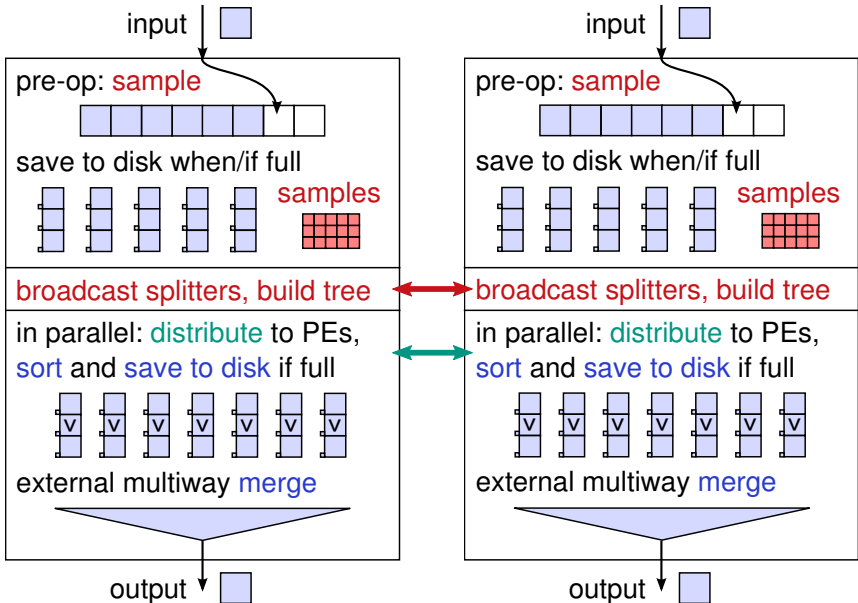


Execution on Cluster



- 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.

Sorting DOp



Benchmarks

WordCountCC

- Reduce text files from CommonCrawl web corpus.

PageRank

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

TeraSort

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

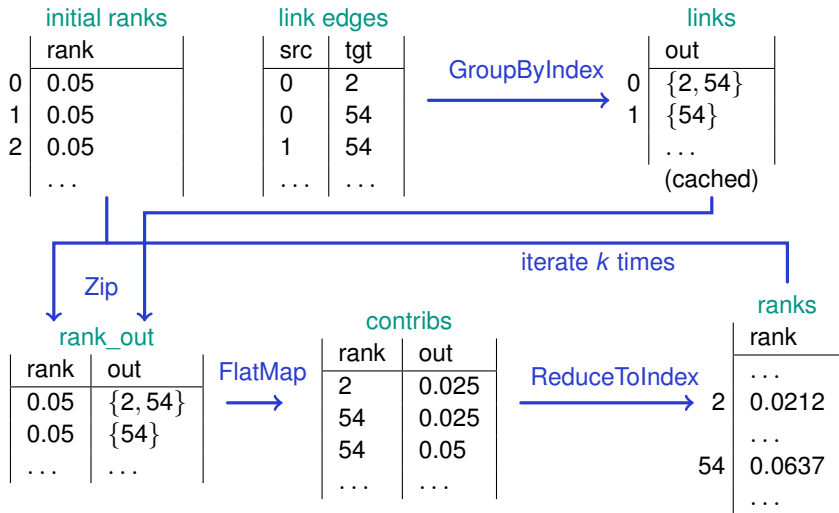
K-Means

- Calculate K-Means clustering with 10 iterations.

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

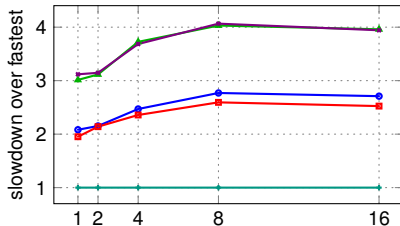
- 32 cores, Intel Xeon E5-2670v2, 2.5 GHz clock, 244 GiB RAM, 2 x 320 GB local SSD disk, ≈ 400 MiB/s bandwidth
Ethernet network ≈ 1000 MiB/s network, Ubuntu 16.04.

PageRank in Thrill

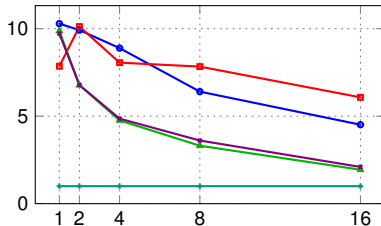


Experimental Results: Slowdowns

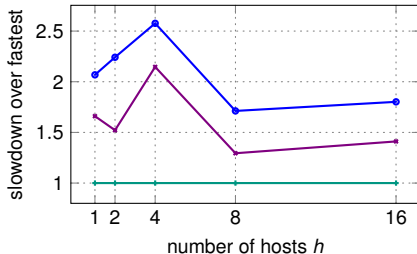
WordCountCC



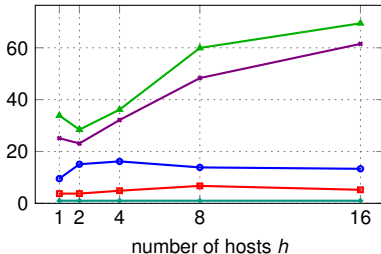
PageRank



TeraSort



KMeans



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

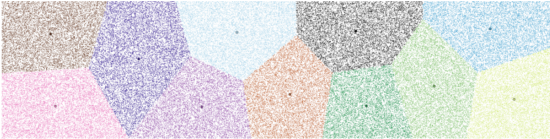
Main Page | **Related Pages** | Modules | Namespaces | Classes | Files | Examples | Search

- Thrill
- Thrill Documentation Overview
- Getting Started
- K-Means Tutorial
 - Step 1: Generate Random Points**
 - Step 2: Pick Random Centers and C
 - Step 3: ReduceByKey to Calculate f
 - Step 4: Iteration!
 - Step 5: Input and Output
 - Bonus Step 6: Boost Qi
- Thrill Layer Architecture
- Coding Style Guide
- Modules
- Namespaces
- Classes
- Files
- Examples

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 $\text{DIA}\langle T \rangle$? $\text{Graph}\langle V, E \rangle$? $\text{DenseMatrix}\langle T \rangle$?
- 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?