

# STXXL and Thrill

## (Distributed Big Data Batch Processing in C++)

Michael Axtmann, Timo Bingmann, Peter Sanders, Sebastian Schlag, and 6 Students | 2016-09-21

INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMICS



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

# Example $T = [\text{dbadc}b\text{cc}b\text{abd}c\text{c}\$]$

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



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

# Flavours of Big Data Frameworks

- High Performance Computing (Supercomputers)

MPI

- Batch Processing

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

- 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.

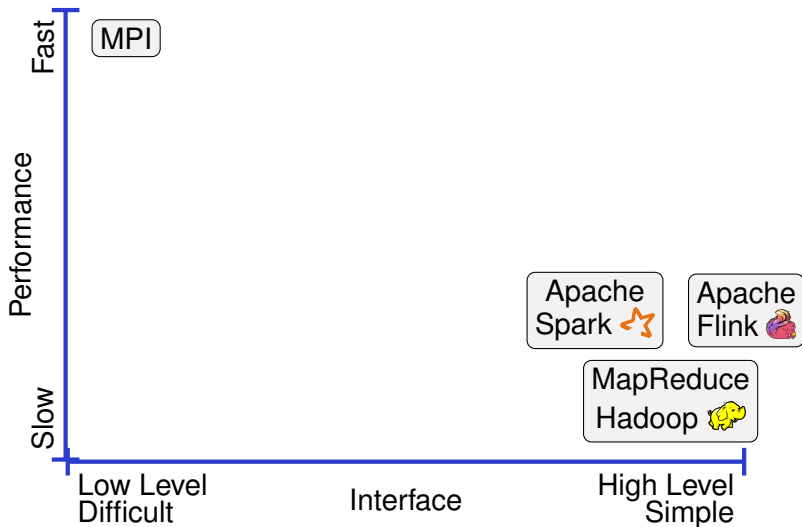
- Graph Processing

Google's Pregel, GraphLab 🐙, Giraph 🧱, GraphChi.

- Time-based Distributed Processing

Microsoft's Dryad, Microsoft's Naiad.

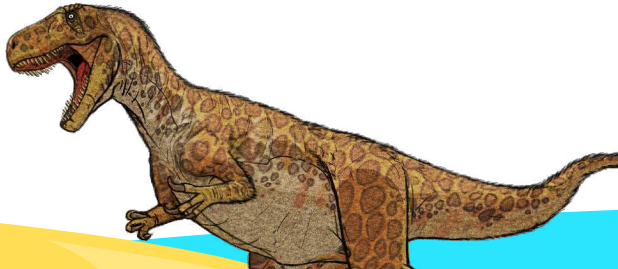
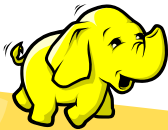
# Big Data Batch Processing



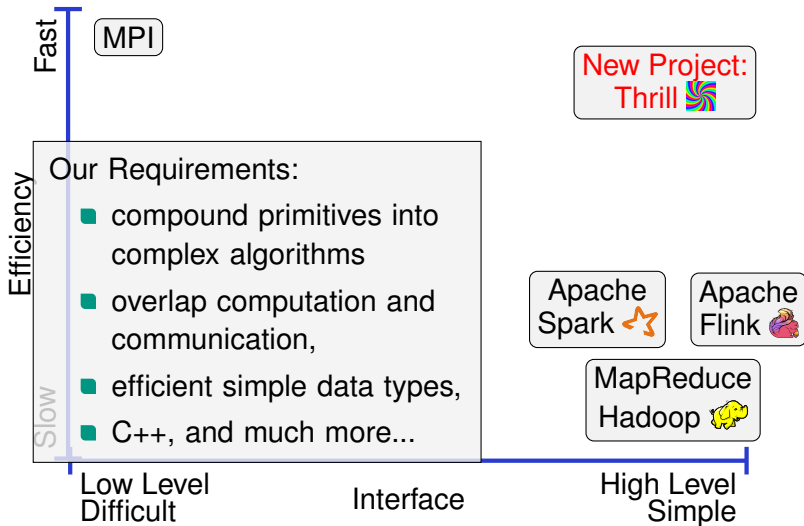
# Projektpraktikum: Verteilte Datenverarbeitung mit MapReduce

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

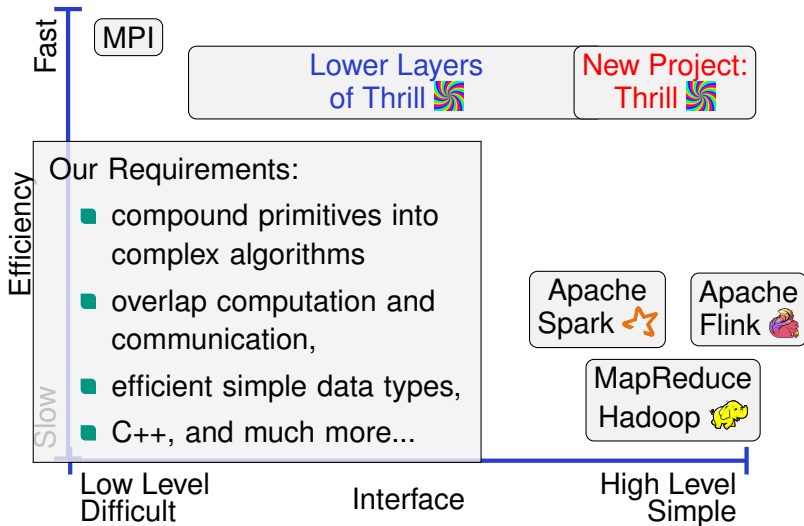
INSTITUTE OF THEORETICAL INFORMATICS – ALGORITHMICS



# Big Data Batch Processing



# Big Data Batch Processing





# Thrill's Design Goals

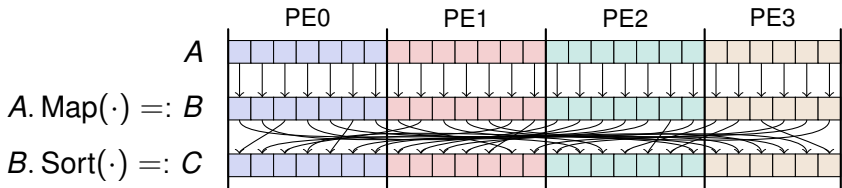
- A new and easier way to program distributed algorithms.
- Distributed arrays of small items (characters or integers).
- High-performance, parallelized C++ operations.
- Locality-aware, in-memory computation.
- Transparently use disk if needed  
⇒ external memory algorithms.
- Avoid all unnecessary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

## Current Status:

- Open-Source at <http://project-thrill.org> and Github.

# Distributed Immutable Array (DIA)

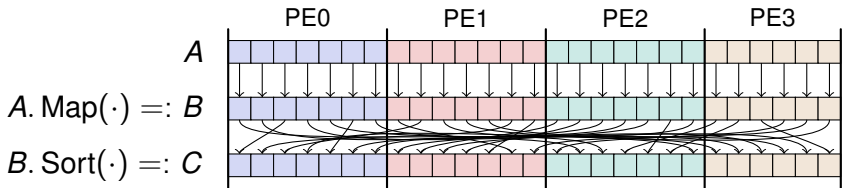
- User Programmer's View:
  - $\text{DIA}\langle T \rangle = \text{result}$  of an operation (local or distributed).
  - Model: **distributed array** of items  $T$  on the cluster
  - Cannot access items directly, instead use **transformations** and **actions**.



# Distributed Immutable Array (DIA)

## ■ User Programmer's View:

- $\text{DIA}\langle T \rangle$  = 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.



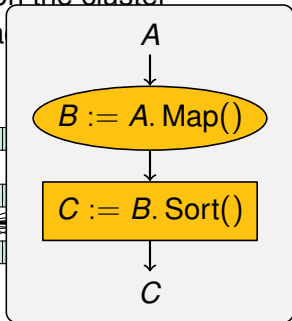
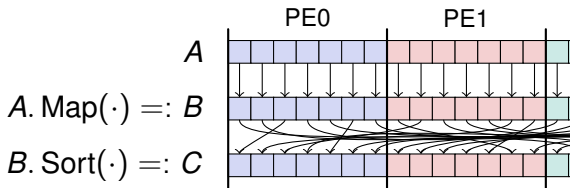
## ■ Framework Designer's View:

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable.  $\implies$  build data-flow graph.
- $\text{DIA}\langle T \rangle$  = chain of computation items

# 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 **actions**.



## ■ Framework Designer's View:

- Goals: distribute work, optimize execution on cluster, add redundancy where applicable.  $\implies$  **build data-flow graph**.
- $\text{DIA}\langle T \rangle$  = **chain of computation items**

# List of Primitives

- Local Operations (LOp): input is **one item**, output  $\geq 0$  items.  
`Map()`, `Filter()`, `FlatMap()`.
- Distributed Operations (DOp): input is a **DIA**, output is a **DIA**.
  - `Sort()` Sort a DIA using comparisons.
  - `ReduceByKey()` Shuffle with Key Extractor, Hasher, and associative Reducer.
  - `GroupByKey()` Like `ReduceByKey`, 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.
  - `Merge()` Merge equal typed DIAs using comparisons.
- **Actions**: input is a **DIA**, output:  $\geq 0$  items **on master**.  
`At()`, `Min()`, `Max()`, `Sum()`, `Sample()`, pretty much still open.

# Example: WordCount in Thrill

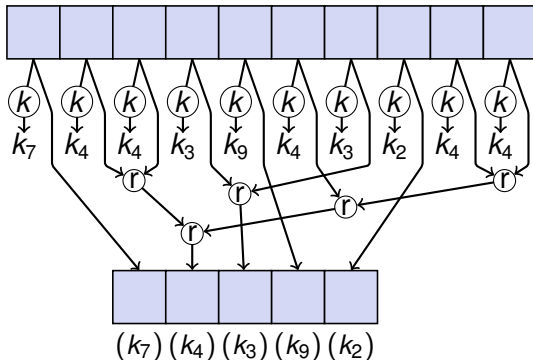
```
1 using Pair = std::pair<std::string, size_t>;
2 void WordCount(Context& ctx, std::string input, std::string output) {
3     auto word_pairs = ReadLines(ctx, input) // DIA<std::string>
4     .FlatMap<Pair>(
5         // flatmap lambda: split and emit each word
6         [](const std::string& line, auto emit) {
7             Split(line, ' ', [&](std::string_view sv) {
8                 emit(Pair(sv.to_string(), 1)); });
9     }); // DIA<Pair>
10 word_pairs.ReduceByKey(
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 }
```

# DOps: ReduceByKey

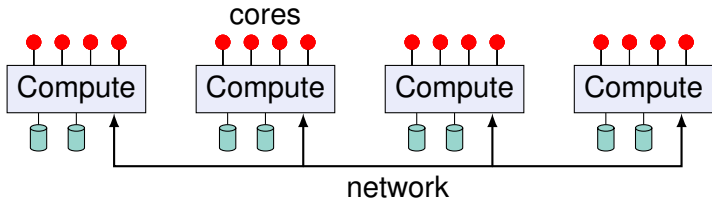
$\text{ReduceByKey}(k, r) : \langle A \rangle \rightarrow \langle A \rangle$

$k : A \rightarrow K$       key extractor

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



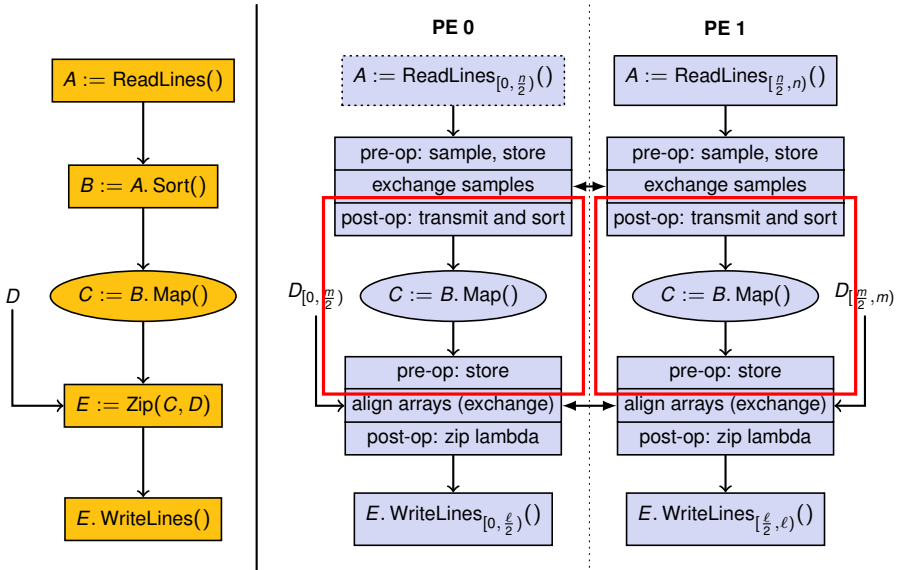
# Execution on Cluster



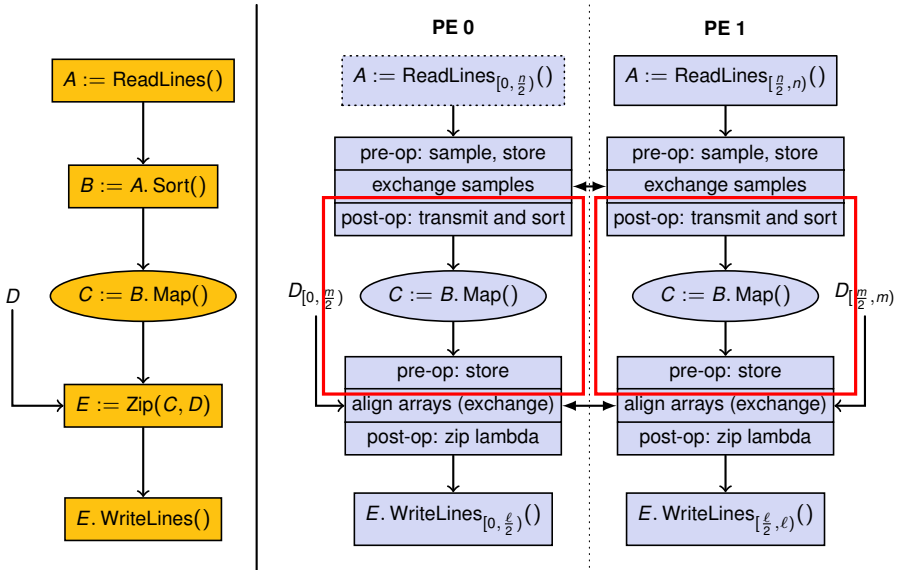
- Compile program into **one binary**, running on all hosts.
- **Collective** coordination of work on compute hosts, like MPI.
- **Control flow** is decided on by using C++ statements.
- Runs on MPI HPC clusters and on Amazon's EC2 cloud.



# Pipelining Stages in Thrill



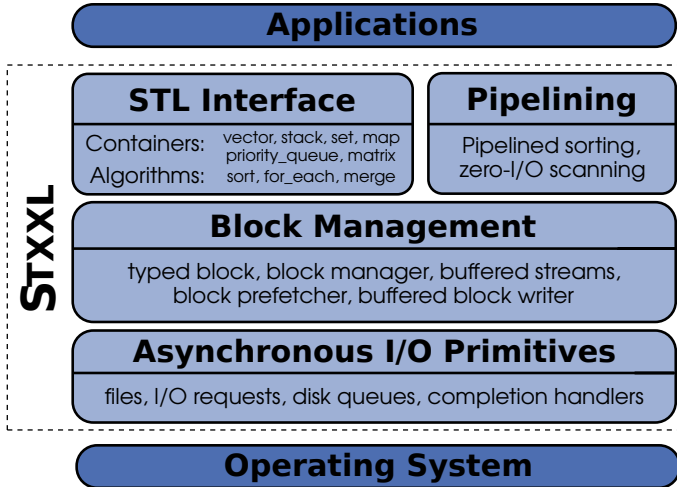
# Pipelining Stages in Thrill



# Layers of Thrill

<p><b>api: High-level User Interface</b> DIA&lt;T&gt;, Map, FlatMap, Filter, Reduce, Sort, Merge, ...</p>	
<p><b>core: Internal Algorithms</b> reducing hash tables (bucket and linear probing), multiway merge, stage executor</p>	
<p><b>data: Data Layer</b> Block, File, BlockQueue, Reader, Writer, Multiplexer, Streams, BlockPool (paging)</p>	<p><b>net: Network Layer</b> (Binomial Tree) Broadcast, Reduce, AllReduce, Async-Send/Recv, Dispatcher Backends: mock, TCP, MPI</p>
<p><b>io: Async File I/O</b> borrowed from STXXL</p>	
<p><b>common: Common Tools</b> Logger, Delegates, Math, ...</p>	<p><b>mem: Memory Limitation</b> Allocators, Counting</p>

# Layers of STXXL



# STXXL and Thrill

	STXXL	Thrill
Model	external	distributed external
Shared Memory	partially parallelized	inherently parallel
External Memory	explicit	via swapping
Items	fixed size	variable length
High-level API	containers, streams	DIA operations
Programming	mix of imperative and functional parts	
Pipelining via	nested template streams	functional templating, consume in pipeline
RAM Limit	manual (parameters)	automatic (stages)
Code Base	C++98	very modern C++14

# STXXL and Thrill

**STXXL API:** vector (**Paging**), sorter (**Sort**), sequence (**Scan**),  
map (**B-Tree**), unordered\_map (**Hash**), priority\_queue (**PQ**),  
parallel\_priority\_queue (**PPQ**), matrix (**Block Matrix**),  
Stream/Pipelining: stream::sort, stream::runs\_creator,  
stream::runs\_merger, ...

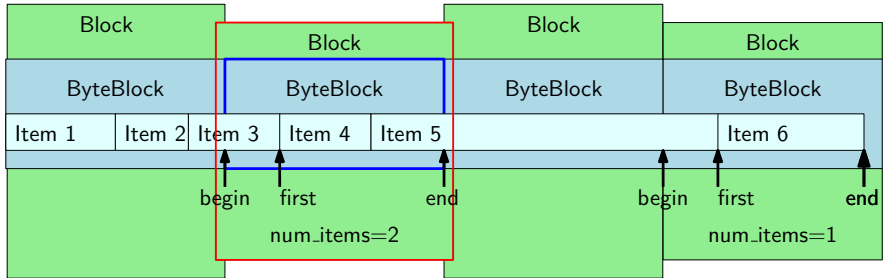
Lower Layers: **BID**, typed\_block, block\_manager,  
read\_write\_pool, buf\_istream, buf\_ostream, async I/O impl.

---

**Thrill API:** DIA.**Generate**, DIA.**Map**, DIA.**ReduceByKey**,  
DIA.**Window**, DIA.**Sort**, DIA.**Union**, DIA.**Zip**, DIA.**ReadBinary**, ...

Lower Layers: **Block**, ByteBlock, **File**, **BlockQueue**,  
BlockReader, BlockWriter, BlockPool (**Paging**), STXXL's  
BlockManager and async I/O implementations.

# File and Blocks in Thrill



# Future of STXXL and Thrill?

On STXXL:

- Pre-C++11 old-style code. But good architecture.
- Many simplifications possible with C++11.
- Shared memory **parallelism** is critically important.
- Missing some important convenience features:  
**variable length items**, **memory management**.

Thrill and STXXL:

- Share much in **common code layers**, esp. C++ tooling.
- Thrill is distributed and **inherently parallel**, but not the API itself.

Thank you for your attention!

Questions? ... Discussion?