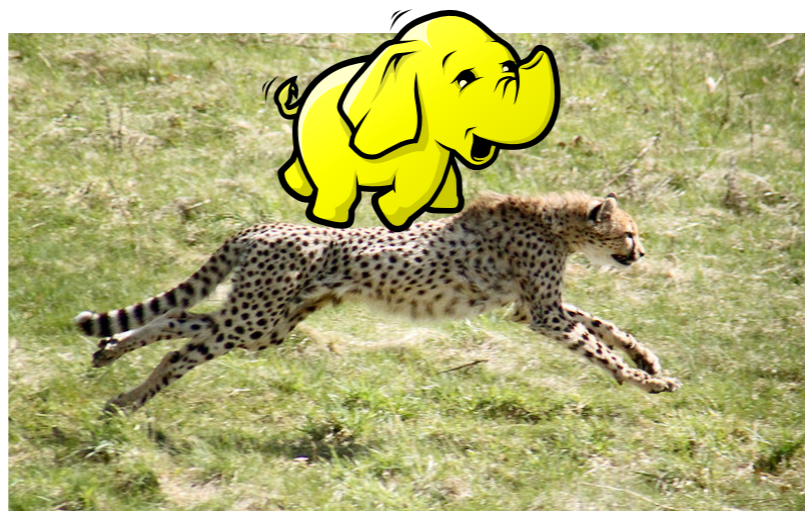# Hadoop++: Making a Yellow Elephant Run Like a Cheetah
## (Without It Even Noticing)

**Jens Dittrich[1]**     **Jorge-Arnulfo Quiané-Ruiz[1]**     **Alekh Jindal[1,2]**

**Yagiz Kargin[2]**     **Vinay Setty[2]**     **Jörg Schad[1]**

[1]Information Systems Group,
Saarland University
http://infosys.cs.uni-saarland.de

[2]International Max Planck
Research School for Computer Science
http://www.imprs-cs.de/

# The Parallel DBMS vs MapReduce Debate

|  | Parallel DBMS | MapReduce |
|---|---|---|
| licensing costs | usually high | none |
| administration | difficult | easy |
| upfront schema | must have | not required |
| user | advanced | beginner |
| scalability | 10-100es of nodes | >10,000 nodes |
| failover, large clusters | suboptimal | very good |
| performance | very good | suboptimal |

- see also [Pavlo etal, SIGMOD 2009] comparison
  - benchmark to compare Parallel DBMS with MapReduce
  - showed superiority of Parallel DBMS over MapReduce

# MapReduce ≠ MapReduce ≠ MapReduce

- but, MapReduce is **three different** things:

(1) a **programming paradigm**:
- it allows users to specify analytical tasks
- need to provide two functions only: `map()` and `reduce()`

(2) a description of a **processing pipeline and system**:
- that system computes the result to a MapReduce-job
- MapReduce-job: `map()`, `reduce()`, and some input data
- scales to very large clusters, > 10,000 nodes

(3) several implementations of (2):
- Google's proprietary MapReduce, Hadoop, ...

# Related Work

| | | (1) Programming Paradigm | | |
|---|---|---|---|---|
| | | MapReduce | SQL | Hybrid |
| (2) Processing pipeline and system | MapReduce | Hadoop | Hive | |
| | PDBMS | Greenplum Vertica | | |
| | Hybrid | | | HadoopDB |

# Related Work

| | (1) Programming Paradigm | | |
| --- | --- | --- | --- |
| | **MapReduce** | **SQL** | **Hybrid** |
| **(2) Processing pipeline and system** — **MapReduce** | Hadoop | Hive | |
| **PDBMS** | Greenplum Vertica | back to initial interface hurdle | |
| **Hybrid** | | | HadoopDB |

# Related Work

| | | (1) Programming Paradigm | | |
|---|---|---|---|---|
| | | **MapReduce** | **SQL** | **Hybrid** |
| **(2) Processing pipeline and system** | **MapReduce** | Hadoop | Hive | back to initial interface hurdle |
| | **PDBMS** | Greenplum proprietary, expensive Vertica | | |
| | **Hybrid** | | | HadoopDB |

# Related Work

| | | (1) Programming Paradigm | | |
|---|---|---|---|---|
| | | **MapReduce** | **SQL** | **Hybrid** |
| **(2) Processing pipeline and system** | **MapReduce** | Hadoop | Hive | |
| | **PDBMS** | Greenplum proprietary, expensive Vertica | back to initial interface hurdle | |
| | **Hybrid** | | admin costs? | HadoopDB |

# Related Work

| | | (1) Programming Paradigm | | |
|---|---|---|---|---|
| | | **MapReduce** | **SQL** | **Hybrid** |

**Research Challenge:**

Can we invent a system that:

(1) keeps the MapReduce programming paradigm **and** the MapReduce execution engine?

(2) approaches Parallel DBMSs in performance?

# Related Work

| | | (1) Programming Paradigm | | |
| --- | --- | --- | --- | --- |
| | | **MapReduce** | **SQL** | **Hybrid** |
| **(2) Processing pipeline and system** | **MapReduce** | Hadoop++ <br><br> Research Challenge: <br> Can we invent a system that <br> (1) keeps the MapReduce programming paradigm and the MapReduce execution engine? <br> (2) approaches Parallel DBMSs in performance? | Hive | back to initial interface hurdle |
| | **PDBMS** | Greenplum proprietary, expensive Vertica | admin costs? | |
| | **Hybrid** | | | HadoopDB |

# Hadoop++ System Vision

**(1) MapReduce programming paradigm**

↓ map(), reduce()

**MapReduce program analysis**
e.g. [Cafarella and Ré, WebDB2010]
[Iu and Zwaenepoel, EuroSys 2010]

↓ logical plan

**Optimization**
e.g. cost models [Morton et.al. SIGMOD 2010]

↓ optimized plan

**MapReduce program generation**
this paper, Hadoop++

↓ map'(), reduce'()

**(2) MapReduce processing pipeline and system**

# Features of Hadoop++

**(1)** we **do not change** the existing Hadoop framework at all

      advantage: no need to maintain and test Hadoop code changes

      advantage: future improvements of Hadoop orthogonal to Hadoop++

**(2)** **inject** our technology inside Hadoop, hide it

      advantage: clear layering

      advantage: no extra operators, no pipeline changes

**(3)** **do not change** the MapReduce programming paradigm

      advantage: nothing changes from the user-side

**(4)** still trick Hadoop into using **more efficient plans**

      advantage: improve runtime performance considerably

**How do we do this?**

Well, let's first better understand the existing Hadoop processing pipeline....

# Analysis: The Hadoop Plan



figure shows example with 4 mappers and 2 reducers

- partition data into blocks
- replicate data to nodes
- store data

- scan input data blocks
- form splits
- send data to processing nodes
- break data into records
- call `map()` for each record
- pregroup and preaggregate output
- store output locally

- redistribute data over processing nodes
- merge subsets belonging to same reducer into single file

- perform final grouping
- call `reduce()` for each group
- store output

# Observations on The Hadoop Plan

- again: no real operators, all hard-coded

- large distributed external merge sort

- sort in order to do a sort-based grouping

- full scan access at all times

- not only two functions, i.e. `map` and `reduce`,

- but...
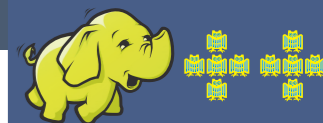
# Ten User-Defined Functions



- **The Hadoop Plan has ten user-defined functions (UDFs):**

  ```
  block
  split
  itemize
  mem
  map
  sh
  cmp
  grp
  combine
  reduce
  ```

figure shows example with 4 mappers and 2 reducers

# Hadoop++ Approach: Trojan Techniques

- **Trojan Index:**
    - at data load time:     create index
    - at query time:         use index access plan



- **Trojan Join:**
    - at data load time:     create co-partitions
    - at query time:         compute all join results locally

# Trojan Index Creation

**Desired layout:**

e.g. 8MB of index for 1GB of data



- **Index Creation Algorithm**:
  - read input split
  - add small clustered Trojan index (we use a CSS-tree)
  - add some metadata

- **Implementation**:
  - a MapReduce program

# Trojan Index Creation

$$\text{map}(\text{key } k, \text{value } v) \mapsto$$
$$[(\text{getSplitID}() \oplus \text{prj}_a (k \oplus v), k \oplus v)]$$

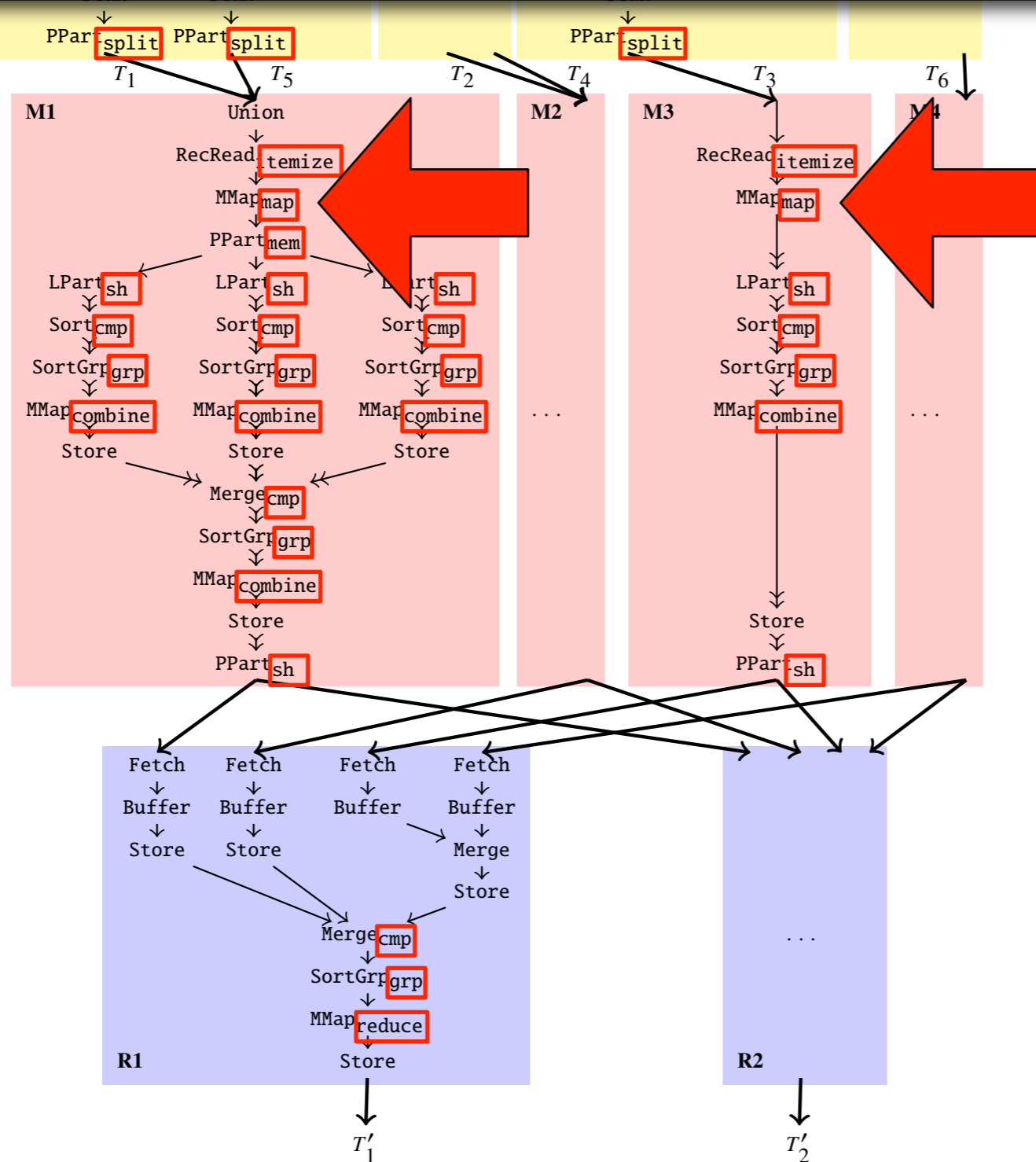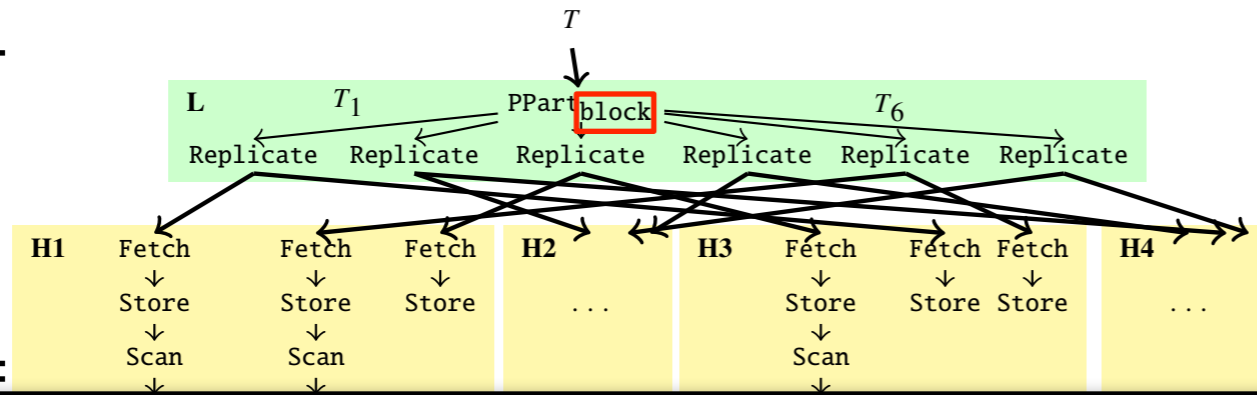form intermediate key with splitID and index key $a$



**Map Phase**

**Shuffle Phase**

**Reduce Phase**
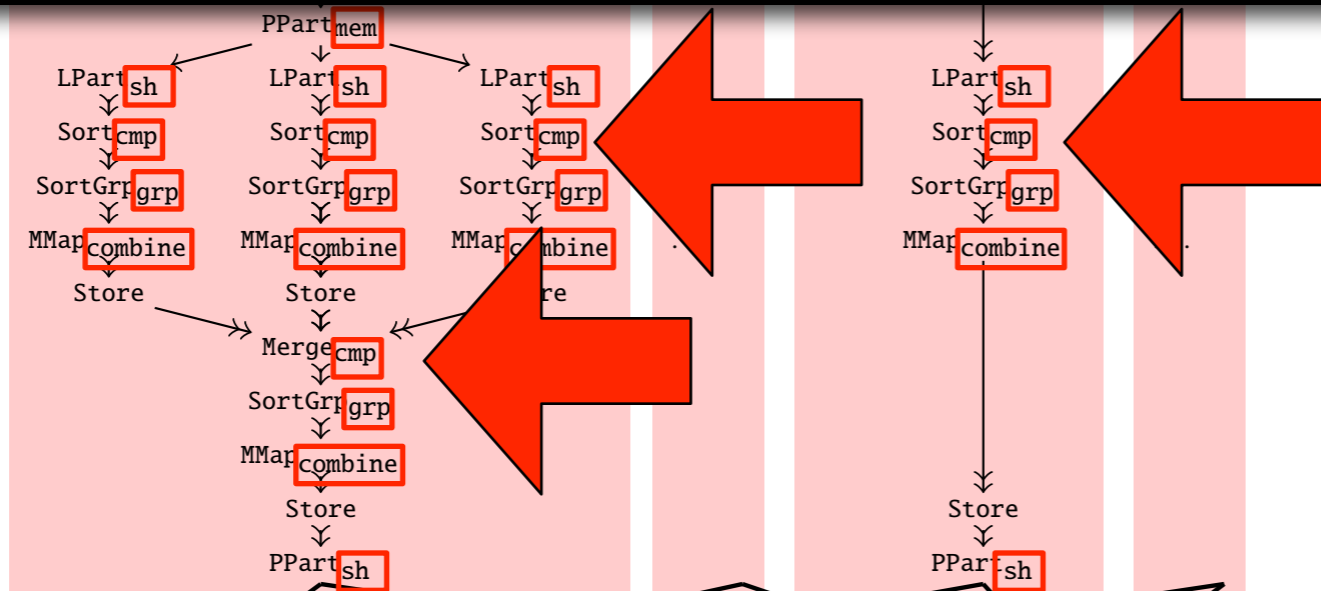
figure shows example with 4 mappers and 2 reducers

$\oplus$: concatenate schemas

12

# Trojan Index Creation

$$\mathrm{cmp}(\mathrm{key}\ k1, \mathrm{key}\ k2) \mapsto \mathrm{compare}(k1.a\ ,\ k2.a)$$

sort on index key only

**Data Load Phase**

**Map Phase**

**Shuffle Phase**

**Reduce Phase**

figure shows example with 4 mappers and 2 reducers

$\oplus$: concatenate schemas

12

# Trojan Index Creation



figure shows example with 4 mappers and 2 reducers

⊕: concatenate schemas    12

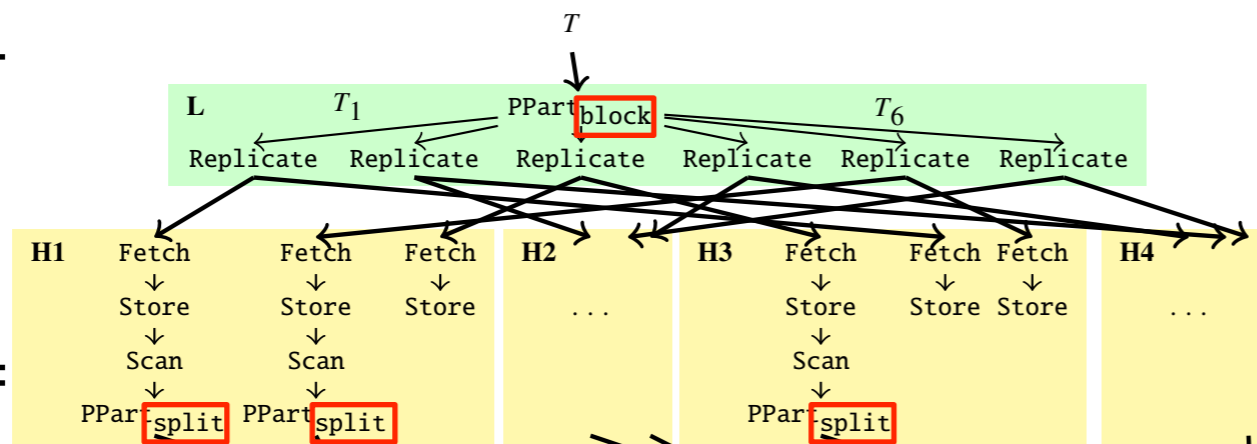# Trojan Index Creation

partition load map reduce

**Data Load Phase**

**Map Phase**

**Reduce Phase**

$$\text{grp}(\text{key } k1, \text{key } k2) \mapsto \text{compare}(k1.splitID, k2.splitID)$$
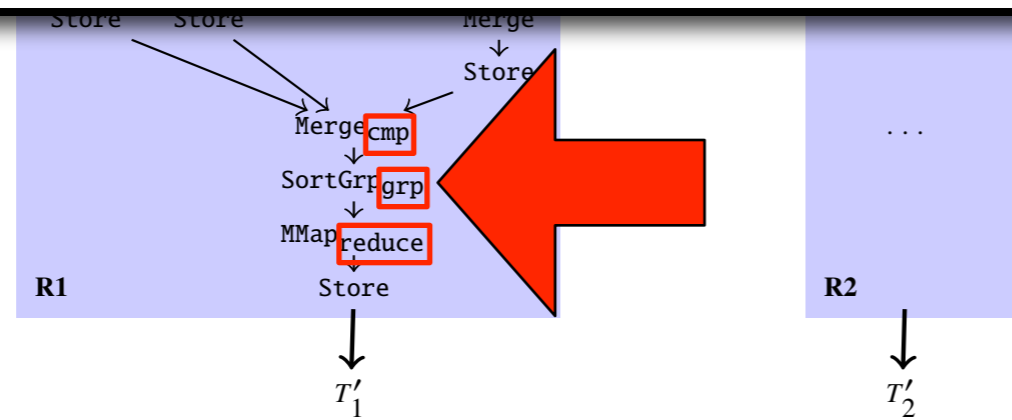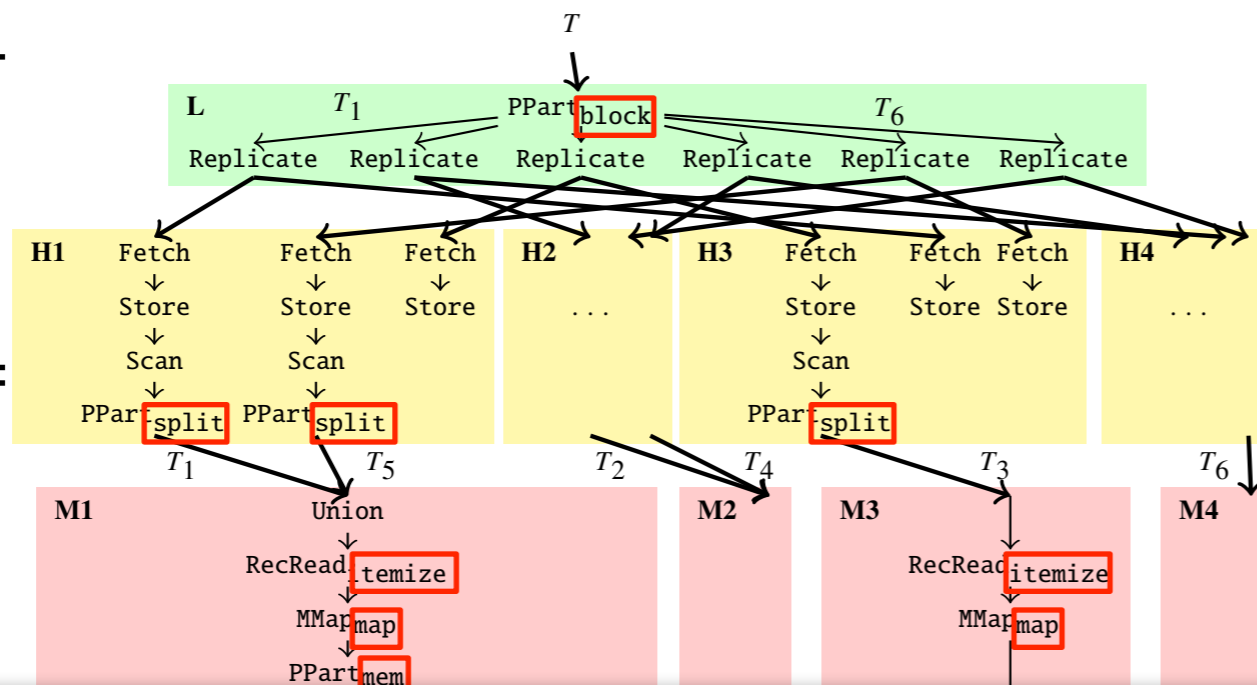
build groups on splitID only

figure shows example with 4 mappers and 2 reducers

$\oplus$: concatenate schemas

12

# Trojan Index Creation



figure shows example with 4 mappers and 2 reducers

$\oplus$: concatenate schemas

12

# Trojan Index Query Processing



DataSet · Trojan Index · SData T · Indexed Split *i* · Header · Footer

- **Query Algorithm**:
  - for each split:
    - read footer to obtain split size
    - read header to obtain [key$_{min}$, key$_{max}$]-range of index
    - if search key overlaps [key$_{min}$, key$_{max}$]-range:
      - read CSS-tree into main memory
      - read only records qualifying for search predicate
      - only pass those records to map()
    - else
      - skip this split

- **Implementation**:
  - a MapReduce program
  - provide `split` and `itemize` UDF
  - everything else unchanged

# Trojan Join Co-Partitioning

join T.a=S.b

**Desired layout:**



- **Co-Partition Creation Algorithm:**
  - read input data
  - create co-partitioned data based on join keys of two relations
  - add some metadata
- **Implementation:**
  - a MapReduce program

# Trojan Join Co-Partitioning Details



legend: partition · load · map · reduce

$$\text{map}(\text{key } k, \text{value } v) \mapsto$$

$$\begin{cases} [(\text{prj}_a\ (k \oplus v), k \oplus v)] & \text{if input}(k \oplus v) = T, \\ [(\text{prj}_b\ (k \oplus v), k \oplus v)] & \text{if input}(k \oplus v) = S . \end{cases}$$

form intermediate key with join key a from T and b from S

Map Phase

Shuffle Phase

Reduce Phase

figure shows example with 4 mappers and 2 reducers

join T.a=S.b    $\oplus$: concatenate schemas

15

# Trojan Join Co-Partitioning Details



partition　load　map　reduce

cmp(key $k1$, key $k2$)　use default ← sort on join key only

Map Phase

Shuffle Phase

Reduce Phase

figure shows example with 4 mappers and 2 reducers

join T.a=S.b　$\oplus$: concatenate schemas

# Trojan Join Co-Partitioning Details



figure shows example with 4 mappers and 2 reducers

join T.a=S.b      $\oplus$: concatenate schemas

# Trojan Join Co-Partitioning Details



**Legend:** partition, load, map, reduce

**Data Load Phase**

T → L: $T_1$ ... PPart block ... $T_6$ — Replicate, Replicate, Replicate, Replicate, Replicate, Replicate

**H1 / H2 / H3 / H4:** Fetch → Store → Scan → PPart split → $T_1$, $T_5$, $T_2$, $T_4$, $T_3$, $T_6$

**Map Phase**

**M1:** Union → RecRead itemize → MMap map → PPart mem → LPart sh → Sort cmp → SortGrp grp → MMap combine → Store → Merge cmp → SortGrp grp → MMap combine → Store

**M3:** RecRead itemize → MMap map → LPart sh → Sort cmp → SortGrp grp → MMap combine → Store

$$\mathrm{grp}(\text{key } k1, \text{key } k2)$$

use default  ←  build groups on join key only

**Reduce Phase**

**R1:** Store, Store, Merge → Store → Merge cmp → SortGrp grp → MMap reduce → Store → $T'_1$
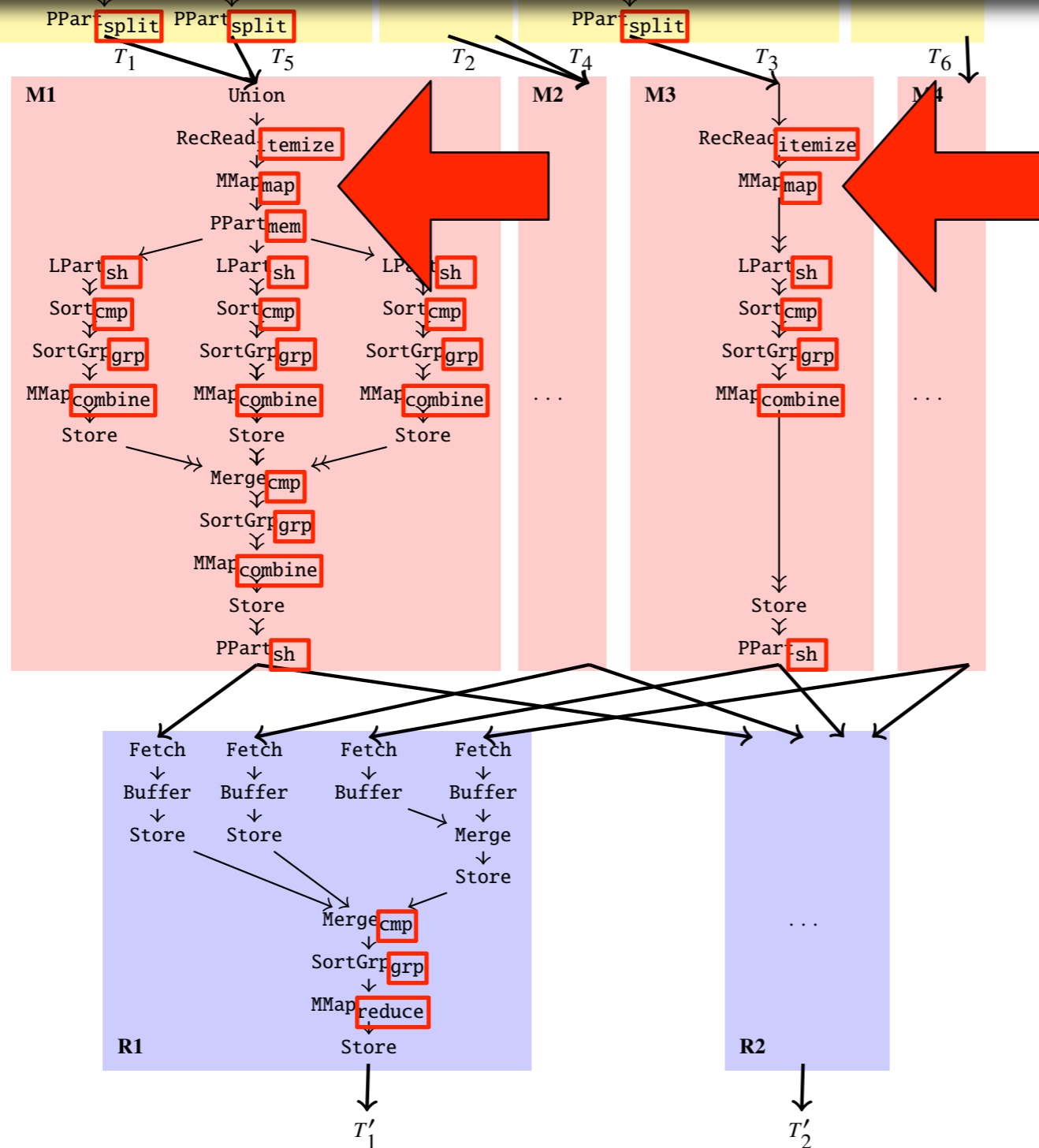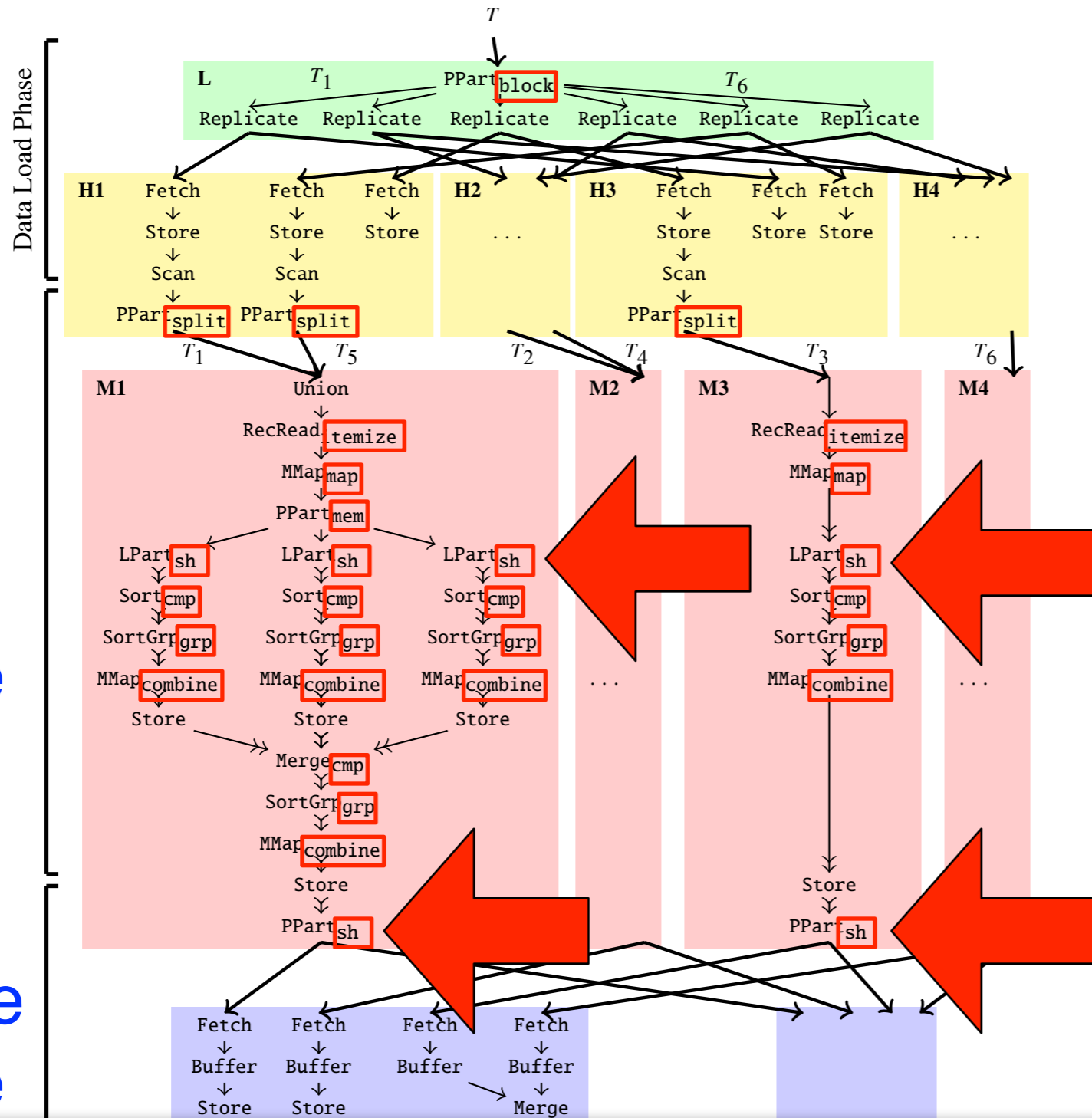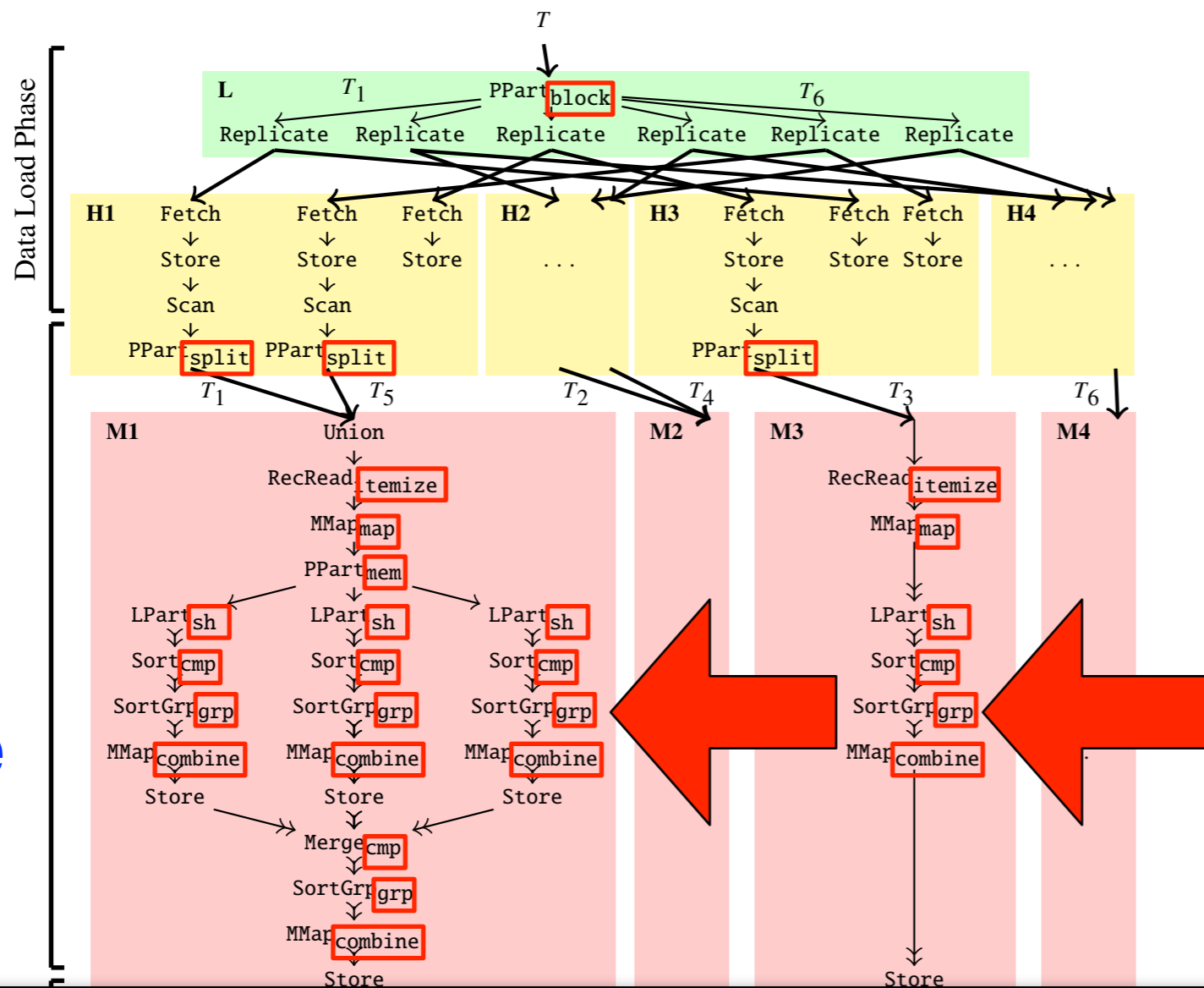
**R2:** ... → $T'_2$

figure shows example with 4 mappers and 2 reducers

join T.a=S.b    ⊕: concatenate schemas

15

# Trojan Join Co-Partitioning Details

partition　load　map　reduce

**Data Load Phase**

$T$

**L**　$T_1$　PPart `block`　$T_6$

Replicate　Replicate　Replicate　Replicate　Replicate　Replicate

**H1** Fetch　Fetch　Fetch　**H2**　**H3** Fetch　Fetch Fetch　**H4**
↓　↓　↓　↓　↓ ↓
Store　Store　Store　...　Store　Store Store　...
↓　↓　↓
Scan　Scan　Scan
↓　↓　↓
PPart`split`　PPart`split`　PPart`split`
$T_1$　$T_5$　$T_2$　$T_4$　$T_3$　$T_6$

**M1**　Union　**M2**　**M3**　**M4**
↓
RecRead `temize`　RecRead `itemize`
↓　↓
MMap`map`　MMap`map`
↓　↓
PPart`mem`

$$\mathrm{reduce}(\mathrm{key}\ ik, \mathrm{vset}\ ivs) \mapsto [(\{ik\} \times ivs)]$$

build one co-group for each join value in a split

SortGrp`grp`
↓
MMap`combine`
↓
Store　Store
↓　↓
PPart`sh`　PPart`sh`

**Shuffle Phase**

Fetch　Fetch　Fetch　Fetch
↓　↓　↓　↓
Buffer　Buffer　Buffer　Buffer
↓　↓　↓
Store　Store　Merge
↓
Store

...

**Reduce Phase**

Merge`cmp`
↓
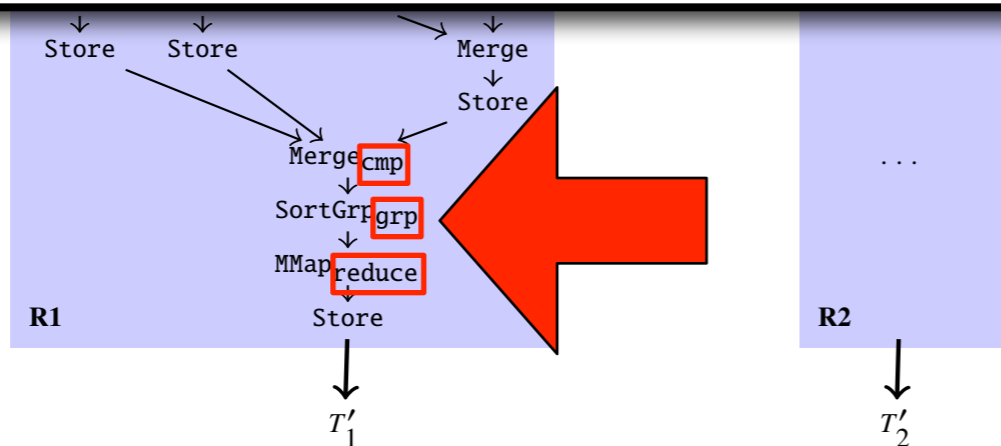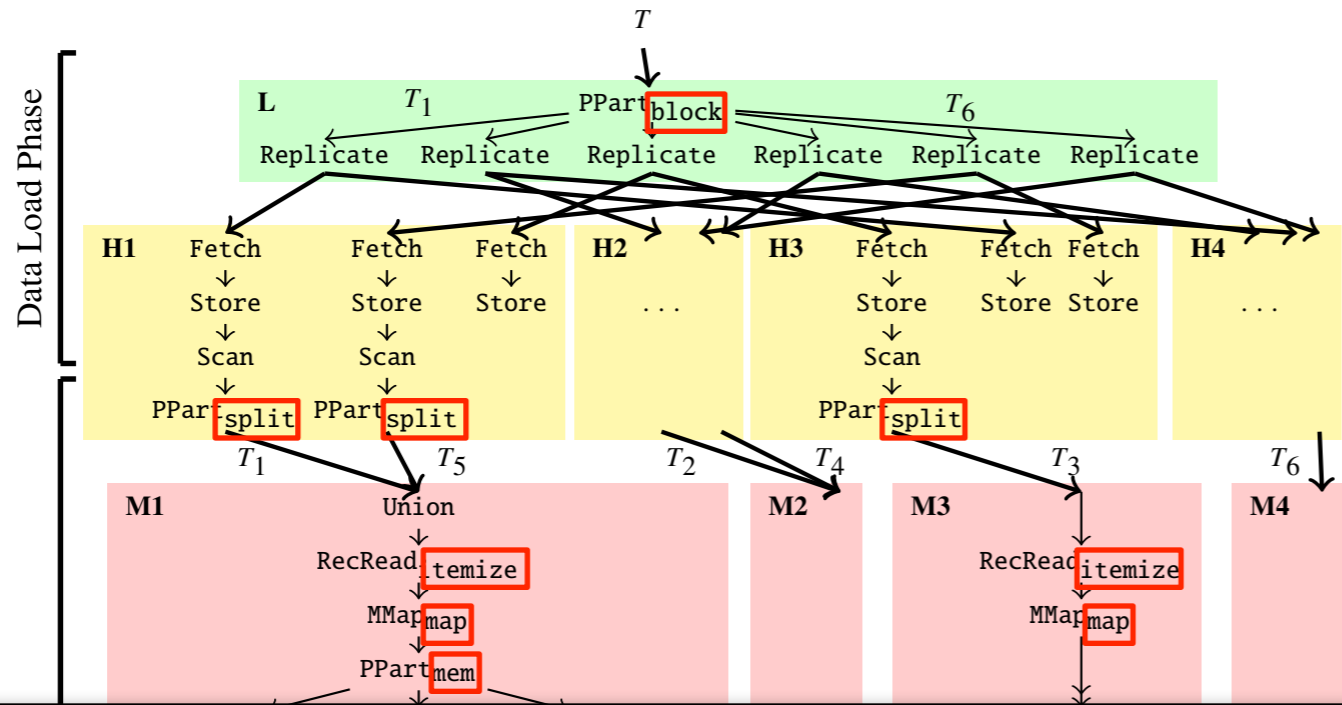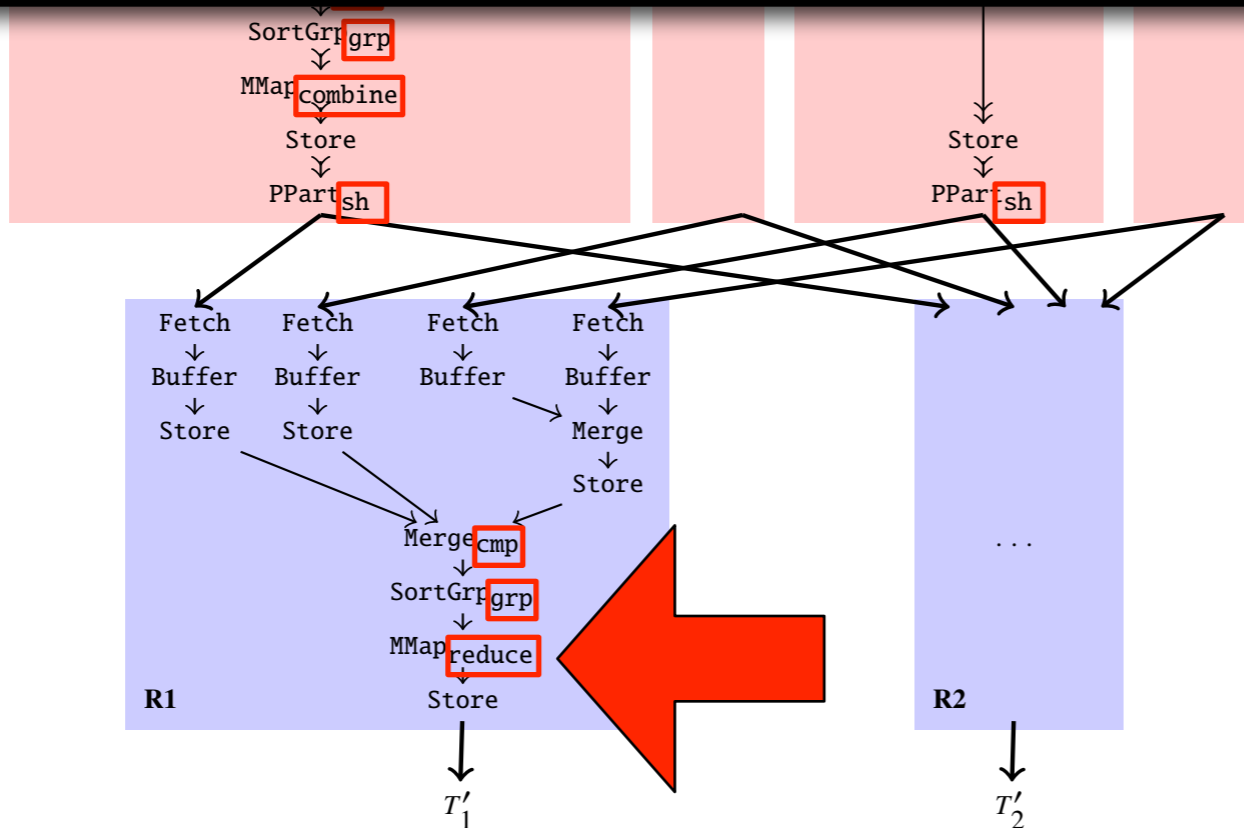SortGrp`grp`
↓
MMap`reduce`
**R1**　Store　**R2**
↓　↓
$T_1'$　$T_2'$

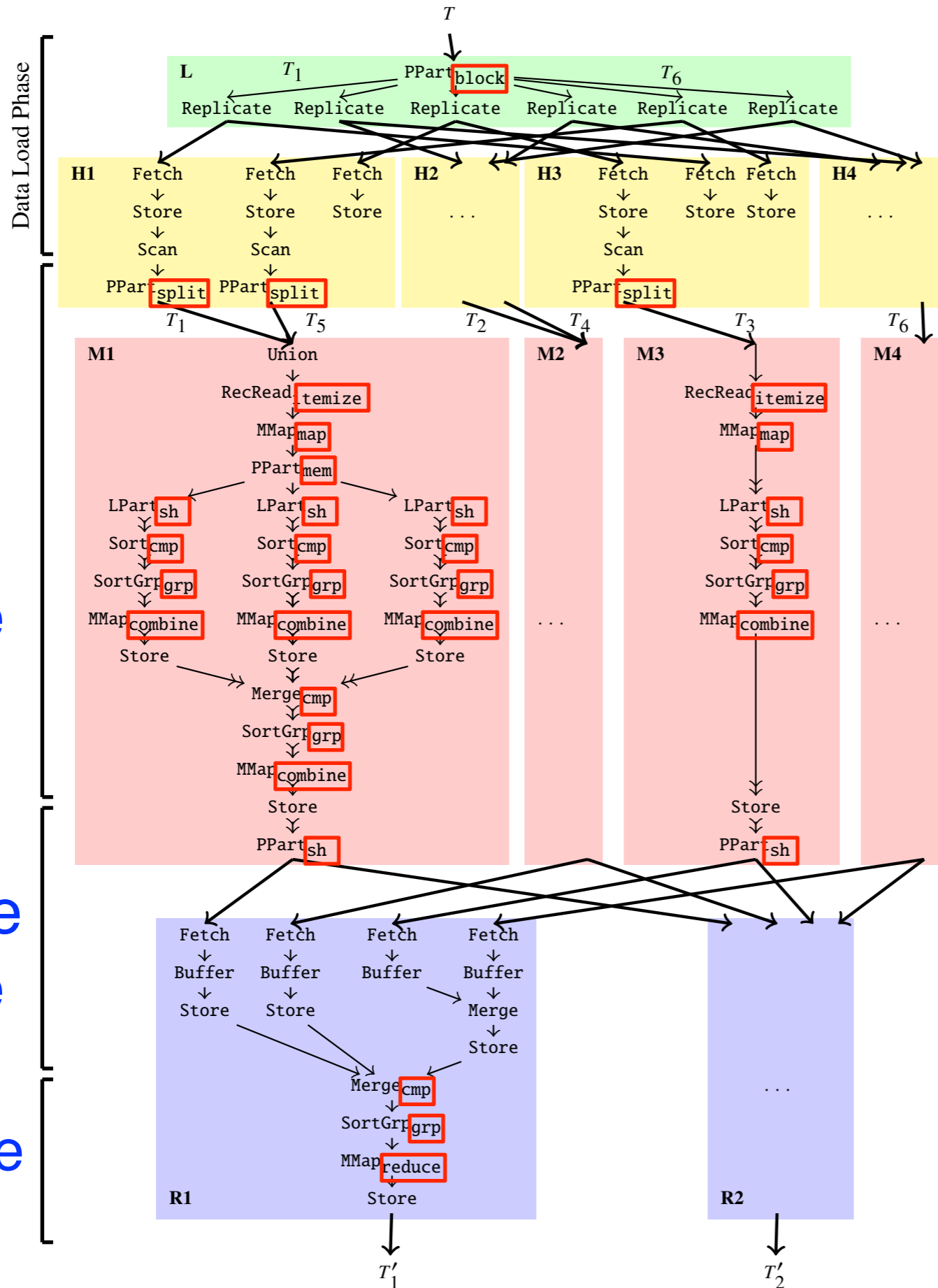figure shows example with 4 mappers and 2 reducers

join T.a=S.b　$\oplus$: concatenate schemas

15

# Trojan Join Co-Partitioning Details



figure shows example with 4 mappers and 2 reducers

join T.a=S.b      ⊕: concatenate schemas

**Notice.** Write-up of these UDFs in the CR has a small bug. See note on our website:

http://infosys.cs.uni-saarland.de/publications/DQJ+10CRv1correction.pdf

15

# Trojan Join Query Processing



- **Query Algorithm:**
  - read footer of each input split to determine split size
  - read records from each co-group in ascending order
  - build cross product for each co-group

- **Implementation:**
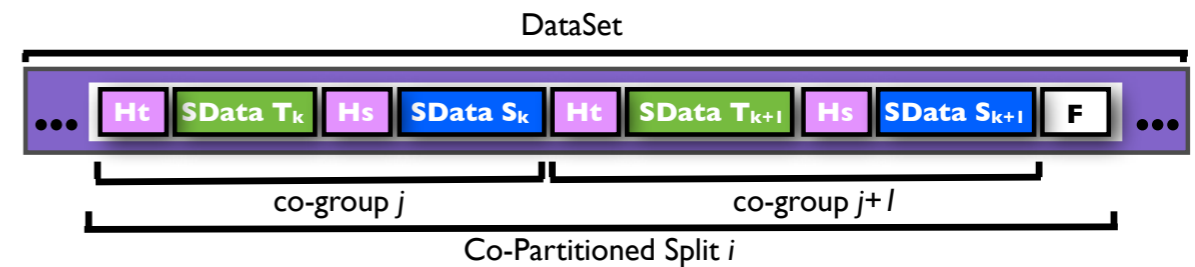  - a MapReduce program
  - provide `split` UDF

- **Option 1: map-side join**
  - trick: map function keeps some state
  - perform local join in `map()`
  - <span style="color:green">advantage</span>: no Reduce Phase (see paper)
  - <span style="color:red">drawback</span>: need to keep some state in `map()` for sort-based grouping

**Figure 1: The Hadoop Plan: Hadoop's processing pipeline expressed as a physical query execution plan**

# 3. HADOOP AS A PHYSICAL QUERY EXECUTION PLAN

as aggre-

As ex-

earch in-

me SQL

via map

and re-

tputs the

found by

ce func-

uzzword

appears,

a DBMS

work as

dvantage

perform

he overall

d DBMS

tasks are

-support,

and *(iii)*

ase, con-

umber of

# Option 2

■ **Algorithm**

- change
[joinkey,

- then map is being called with the data belonging to an entire co-group

- inside map: b ... les and co ... cross product

■ Advantages:

- no Reduce P

- but also: **no** ... map

**in fact:** we exploited an interesting order plus ... temize to **semantically reduce** data in map!

map phase

shuffle phase

reduce phase

LogPart$_{sh}$    LogPart$_{sh}$    LogPart$_{sh}$       LogPart$_{sh}$

Sort    Sort    Sort      Sort

SortGrp    SortGrp    SortGrp      SortGrp

MMap$_{combine}$   MMap$_{combine}$   MMap$_{combine}$  …    MMap$_{combine}$  …

Store    Store    Store

Merge

SortGrp

MMap$_{combine}$

Store

PhysPart$_{sh}$

DataSet

… | Ht | SData T$_k$ | Hs | SData S$_k$ | Ht | SData T$_{k+1}$ | Hs | SData S$_{k+1}$ | F | …

co-group *j*      co-group *j+1*

Co-Partitioned Split *i*

Store

PhysPart$_{sh}$

Fetch   Fetch   Fetch   Fetch

Buffer   Buffer   Buffer   Buffer

Store   Store     Merge

Store

Merge

SortGrp

MMap$_{reduce}$

Store

**R1**          …      **R2**

$T'_1$      $T'_2$

**Figure 1: The Hadoop Plan: Hadoop's processing pipeline expressed as a physical query execution plan**

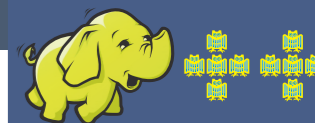## 3. HADOOP AS A PHYSICAL QUERY EXECUTION PLAN

# Trojan Index plus Trojan Join

- may combine both techniques

- may use index on join key

- may use index on different key

- may create multiple indexes inside the split

- in any case:
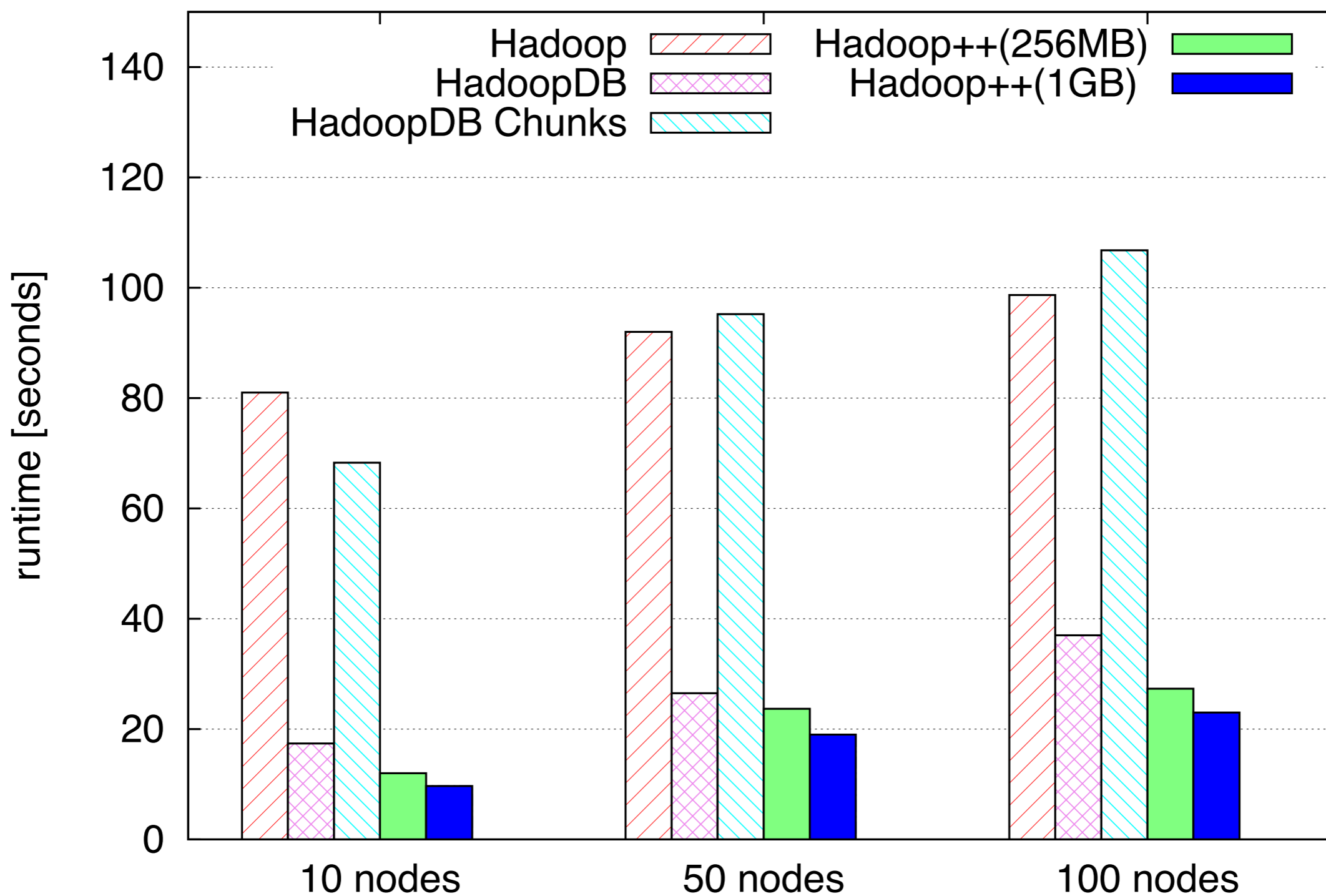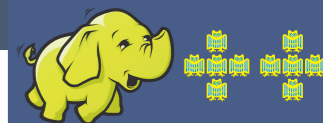  - **both** scan access **and** index access paths possible

# Experiments

- used benchmark as proposed in [Pavlo etal, SIGMOD 2009]

- benchmark defines several tasks

- two of them related to indexing and join processing
  - Selection Task
  - Join Task

- used up to 100 EC2 nodes as in HadoopDB-paper [Abouzeid etal, VLDB 2009]

- report average of three executions

- Some twist, see our paper:
  **Runtime Measurements in the Cloud: Observing, Analyzing, and Reducing Variance**
  Jörg Schad, Jens Dittrich, Jorge-Arnulfo Quiané-Ruiz
  **VLDB 2010**
  Research Session-14 : Experimental Analysis and Performance (i.e., yesterday)

- **therefore:** also executed scaled-down experiments on small local cluster to verify
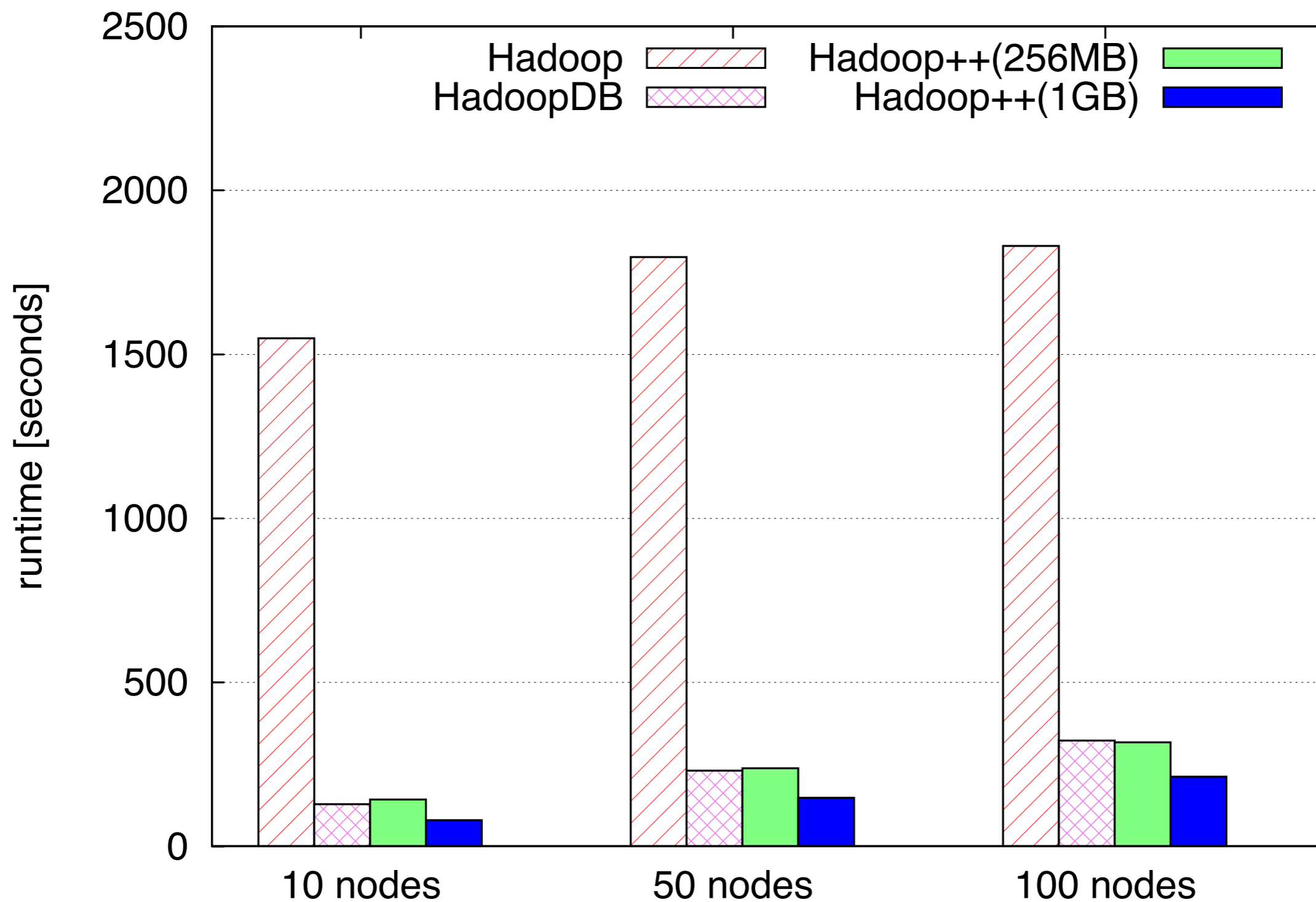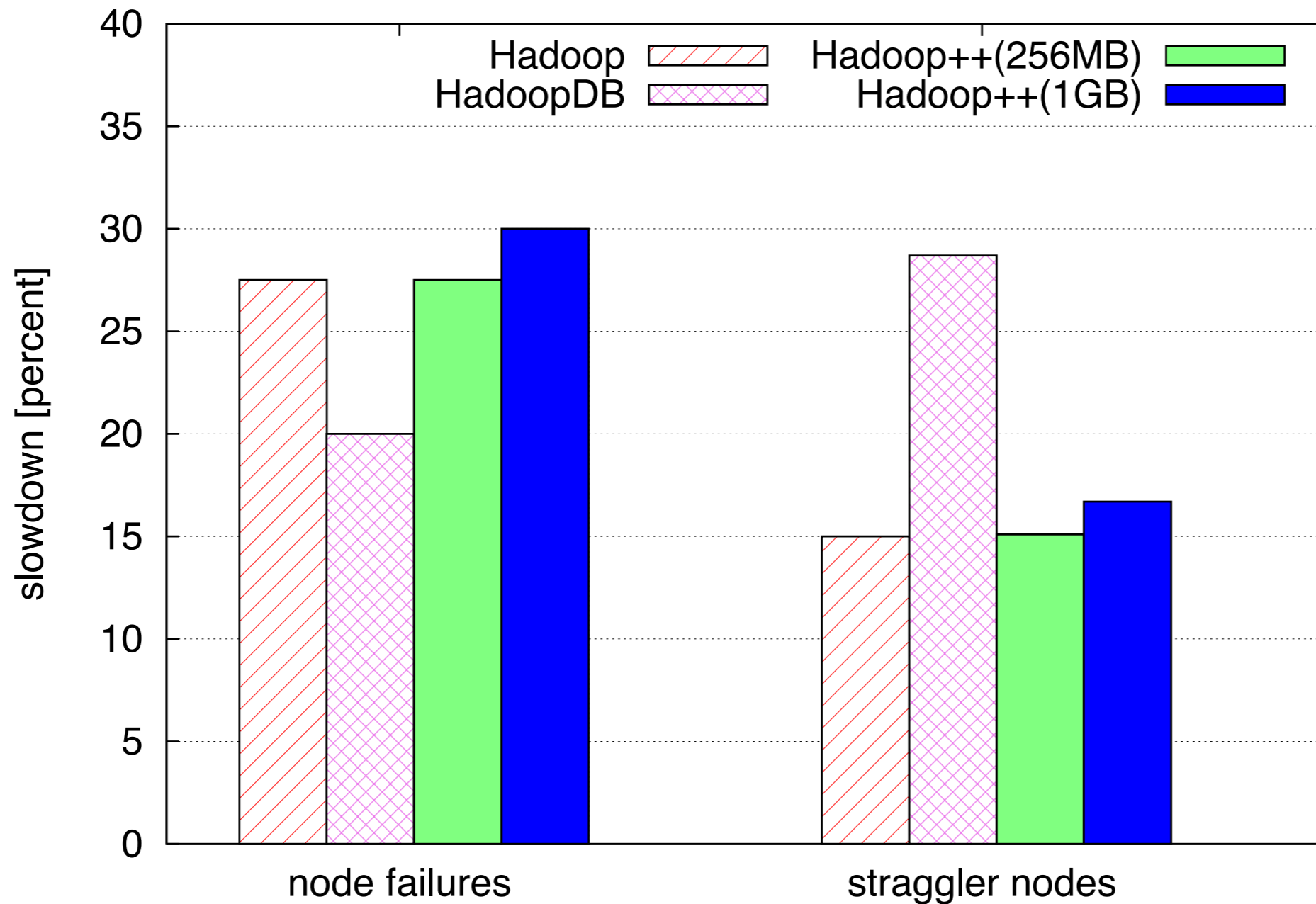
# Selection Task

# Join Task

# Failover



- we inherit fault tolerance from Hadoop!
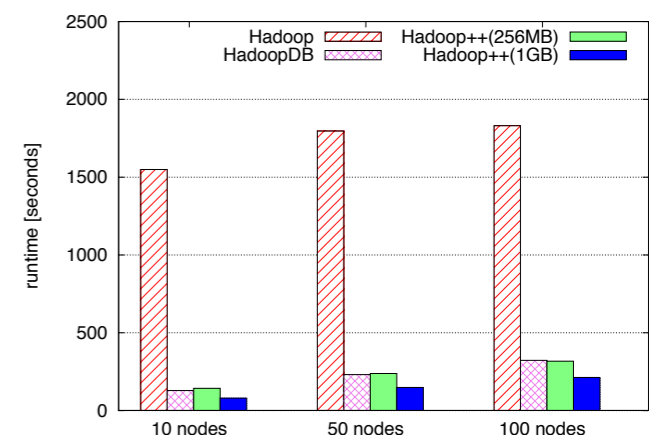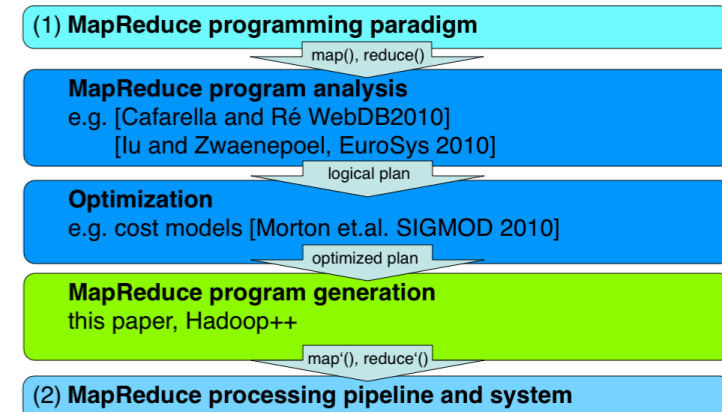
- the Trojan effect!

# Lessons Learned for our Community

- indexing, co-partitioning, preprocessing, etc....

- ...are **not exclusive** to database management systems

- all these techniques may be successfully used in **any** data processing system, not only DBMS

- just one thing matters:

- **"Do we know anything about the schema and the anticipated workload in advance?"**

- if **yes**, we may:
  - create appropriate indexes
  - create co-partitions
  - etc.

- this holds for **both**
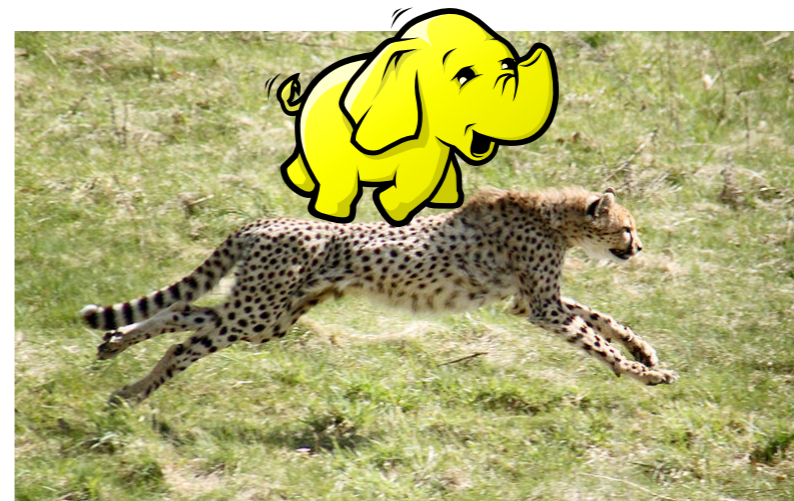  - DBMS
  - and MapReduce/Hadoop

# Conclusions

- we proposed Hadoop++

- a new approach to large scale data analysis

- keep the MapReduce interface
  **and** the MapReduce execution engine

- still: rewrite incoming MapReduce programs to more efficient ones

- inject code through **Trojan techniques**

- execute plans using existing MapReduce pipeline unchanged

- experimens with SIGMOD 2009 benchmark

- strong improvements in selection and join tasks

- up to a factor of 18 better than Hadoop

# Future Work



- other Trojan techniques

  ongoing

- research challenges when executing MapReduce on the Cloud

  **Flying Yellow Elephant: Predictable and Efficient MapReduce in the Cloud**
  Jörg Schad
  **VLDB PhD Workshop 2010 (see VLDB USB stick or online)**

- marry Hadoop++ with OctopusDB* one-size-fits-all DBMS

  **The Mimicking Octopus: Towards a one-size-fits-all Database Architecture**
  Alekh Jindal
  **VLDB PhD Workshop 2010 (see VLDB USB stick or online)**

*patent pending