Ph.D.

# ONE SIZE DOES NOT FIT ALL

Streaming

OLAP

OLTP

Archiving

Log-processing

Web-search

Scan-oriented

amazon.com

ebaY

Streaming

OLAP

OLTP

Archiving

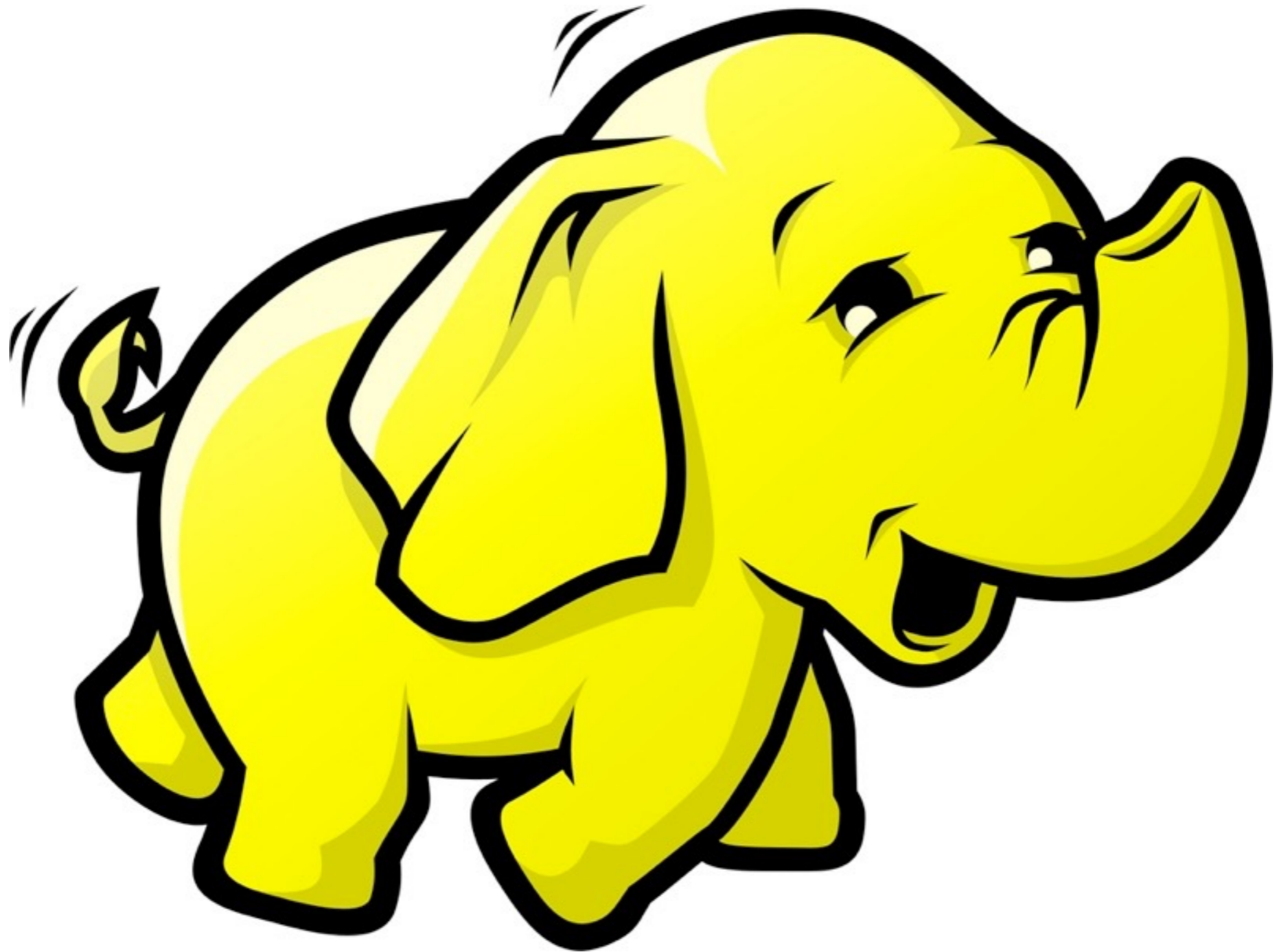Log-processing

Web-search

Scan-oriented

4

OLTP

OLAP

# Log-processing



5

Streaming

OLAP

OLTP

Archiving

Log-processing

Web-search

Scan-oriented

Indexes

Column

Row

Raw files

Row+Column

# Storage Views

| | | | |
|---|---|---|---|
| 1 | abc | 56 | 887.9 |
| 2 | fdg | 89 | 445.35 |
| 3 | poe | 67 | 234.67 |
| 4 | lkj | 12 | 385.92 |
| 5 | yui | 17 | 612.13 |
| 6 | omg | 90 | 148.9 |

# Storage Views

| | | | |
|---|---|---|---|
| Log | | 56 | 887.9 |
| | | 89 | 445.35 |
| | | 67 | 234.67 |
| 4 | lkj | 12 | 385.92 |
| 5 | yui | 17 | 612.13 |
| 6 | omg | 90 | 148.9 |

# Storage Views

| | | | |
|---|---|---|---|
| Log | | 56 | 887.9 |
| | | 89 | 445.35 |
| | | 67 | 234.67 |
| 4 | Row | | |
| 5 | yui | 17 | 612.13 |
| 6 | omg | 90 | 148.9 |

8

# Storage Views



| Log | | 56 | 887.9 |
|-----|-----|-----|-------|
| | | 89 | 445.35 |
| | | 67 | 234.67 |
| Column | Row | | |
| | yui | 17 | 612.13 |
| | omg | 90 | 148.9 |

# Storage Views

| | | |
|---|---|---|
| | 56 | 887.9 |
| Log | 89 | 445.35 |
| | 67 | 234.67 |

| | |
|---|---|
| | Row |
| Column | Column grouped |

8

# Storage Views

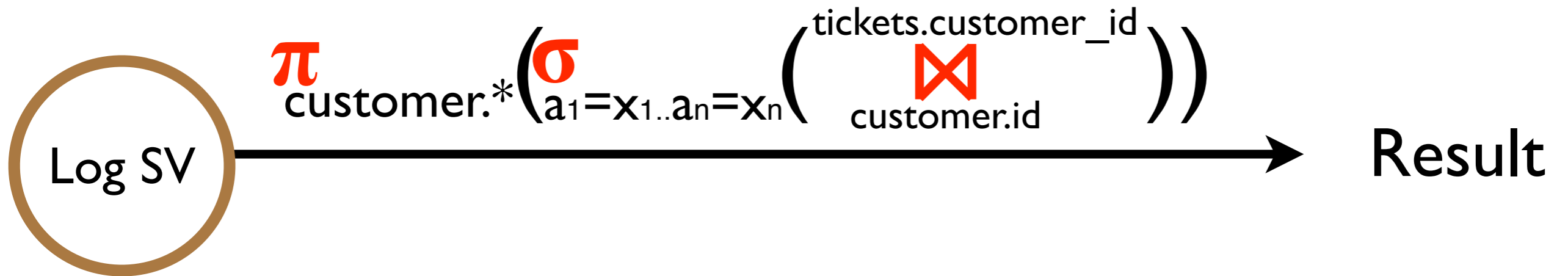# Storage Views

# Example: Flight Tickets



$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie^{tickets.customer\_id}_{customer.id}\right)\right)$$

Log SV → Result

$\sigma_{bag=tickets}$

ticketsLog: Log SV

$\sigma_{bag=customers}$

customersLog: Log SV

$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie^{tickets.customer\_id}_{customer.id}\right)\right)$$

Result

9

$$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=tickets}\right)\right)$$

$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$$

Log SV

Log SV

Log SV

Result

$$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=customers}\right)\right)$$

Log SV

$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$$

Log SV

Result

$$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=tickets}\right)\right)$$

Col SV  tickets

$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$$

Log SV

Result

$$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=customers}\right)\right)$$

ticketsLog

Row SV

$\sigma_{bag=tickets}$

Log SV

customers

$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$$

g SV

Result

$\sigma_{bag=customers}$

Log SV

customersLog

ticketsCold

Col SV

$\sigma_{time<now-7days}$

ticketsHot

Col SV

$$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$$

$\gamma_{\ldots}\left(\Gamma_{\ldots}\left(\sigma_{bag=tickets}\right)\right)$

9

# Example: Flight Tickets

$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=tickets}\right)\right)$

$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=customers}\right)\right)$

Log SV

$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{tickets.customer\_id = customer.id}\right)\right)$ → Result

$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{tickets.customer\_id}^{customer.id}\right)\right)$

Primary Log Store

$\sigma_{bag=customers |\ (bag=tickets\ \&\ tickets.class=E)}$

$\pi_{price,\ customer\_id}\left(\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=tickets}\right)\right)\right)$

$\pi_{price}$ → Result 1

$\pi_{id}\left(\sigma_{a_1=x_1}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$ → Result 2

$\pi_{name,\ email}\left(\sigma\left(\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=customers}\right)\right)\right)\right)$ → Result 3

$\pi_{name}\left(\sigma_{a_1=x_1}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$

$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie\right)\right)$

$\pi_{email}$ → Result 4

$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=tickets}\right)\right)$ ticketsLog

Row SV

→ Result

$\pi_{customers}$

$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$

$\sigma_{bag=tickets}$

$\sigma_{bag=customers}$

Col SV — Cold

Index SV

$\sigma_{time<now-7days}$

$\pi_{price,rid}$

$\pi_{customer.*}\left(\sigma_{a_1=x_1..a_n=x_n}\left(\bowtie_{customer.id}^{tickets.customer\_id}\right)\right)$

$\gamma_{recent}\left(\Gamma_{bag,key}\left(\sigma_{bag=tickets}\right)\right)$

$\sigma_{time>=now-7days}$ customersLog

$\gamma_{count(*)>=5}\left(\Gamma_{customer\_id}\right)$ ticketsCold

→ Result

$\pi_{id,rid}$

$\sigma_{time<now-7days}$

ticketsHot

$\bowtie_{customer.id}^{tickets.customer\_id}$

Index SV

Frequent Fliers (Adaptive Partial Index)

# Example: Flight Tickets

# Rodent Store

# What to store?



Data Files

copy 1    copy 2    copy 3

# How to store?

Data Files → **?**

$\left[\begin{array}{c} \\ \\ \\ \\ \end{array}\right.$ + $\triangle$a  $\triangle$b

# Where to store?



Data Files → **?**

# Example Use-cases

- WWHow! File System

- WWHow! RAID

- WWHow! Relational DBMS

- WWHow! Cloud

STORE '/Users/Bob/Conferences/Talks/*.*'
**WHAT** *.(pdf | ppt), *.pdf
**WHERE** vise4
**HOW** encryption(rsa) FOR *;

STORE '/Users/Bob/Conferences/Talks/*.*'

WHAT *.(pdf | ppt), *.pdf

HOW encryption(rsa) FOR *

**PREFERENCE** Availability='high';

# OctopusDB

- Cool Vision

- Tough to realize

# C-Store

# Trojan Columns

**Application**

User

Database

Query Processor

Relations

**UDF Storage Layer**

Physical Representation

| File 1 | File 2 | File 3 | .... | File n |

23

# Trojan Columns

## Relation

| Customer | | |
|---|---|---|
| name | phone | market_segment |
| smith | 2134 | automobile |
| john | 3425 | household |
| kim | 6756 | furniture |
| joe | 9878 | building |
| mark | 4312 | building |
| steve | 2435 | automobile |
| jim | 5766 | household |
| ian | 8789 | household |

## Physical Table

| Customer_trojan | | |
|---|---|---|
| segment_ID | attribute_ID | blob_data |
| 1 | name | smith, john, kim, joe |
| 1 | phone | 2134, 3425, 6756, 9878 |
| 1 | market_segment | automobile, household, furniture, building |
| 2 | name | mark, steve, jim, ian |
| 2 | phone | 4312, 2435, 5766, 8789 |
| 2 | market_segment | building, automobile, household, household |

# Trojan Columns

## Relation

| Customer | | |
|---|---|---|
| name | phone | market_segment |
| smith | 2134 | automobile |
| john | 3425 | household |
| kim | 6756 | furniture |
| joe | 9878 | building |
| mark | 4312 | building |
| steve | 2435 | automobile |
| jim | 5766 | household |
| ian | 8789 | household |

write-UDF

Tuple Iterator

(a) Convert row tuples into blobs

(c) Get next row data

Data Parser

(b) Store blob data

Data Accessor

## Physical Table

| Customer_trojan | | |
|---|---|---|
| segment_ID | attribute_ID | blob_data |
| 1 | name | smith, john, kim, joe |
| 1 | phone | 2134, 3425, 6756, 9878 |
| 1 | market_segment | automobile, household, furniture, building |
| 2 | name | mark, steve, jim, ian |
| 2 | phone | 4312, 2435, 5766, 8789 |
| 2 | market_segment | building, automobile, household, household |

# Trojan Columns



## Relation

| Customer | | |
|---|---|---|
| name | phone | market_segment |
| smith | 2134 | automobile |
| john | 3425 | household |
| kim | 6756 | furniture |
| joe | 9878 | building |
| mark | 4312 | building |
| steve | 2435 | automobile |
| jim | 5766 | household |
| ian | 8789 | household |

## Physical Table

| Customer_trojan | | |
|---|---|---|
| segment_ID | attribute_ID | blob_data |
| 1 | name | smith, john, kim, joe |
| 1 | phone | 2134, 3425, 6756, 9878 |
| 1 | market_segment | automobile, household, furniture, building |
| 2 | name | mark, steve, jim, ian |
| 2 | phone | 4312, 2435, 5766, 8789 |
| 2 | market_segment | building, automobile, household, household |

25

# Example: TPC-H Query 6

**Result**

↑

$\gamma_{agg}$ (extendedprice * discount)

↑

$\sigma$ shipdate BETWEEN
'1994-01-01' AND '1995-01-01'
AND discount BETWEEN
0.05 AND 0.07
AND quantity < 24

↑

$\pi$ quantity, discount
extendedprice, shipdate

↑

**Scan**

↑

lineitem

# Example: TPC-H Query 6

**Result**

↑

$\gamma_{agg}$ (extendedprice * discount)

↑

$\sigma$ shipdate BETWEEN
'1994-01-01' AND '1995-01-01'
AND discount BETWEEN
0.05 AND 0.07
AND quantity < 24

↑

$\pi_{quantity, discount}$
    extendedprice, shipdate

↑

**Scan**

↑

lineitem

**Result**

↑

$\gamma_{agg}$ (extendedprice * discount)

↑

$\sigma$ shipdate BETWEEN
'1994-01-01' AND '1995-01-01'
AND discount BETWEEN
0.05 AND 0.07
AND quantity < 24

↑

$\pi_{quantity, discount}$
    extendedprice, shipdate

↑

**Scan**

↑

lineitem

scanUDF

# Example: TPC-H Query 6

**Result**

$\gamma_{agg}$ (extendedprice * discount)

$\sigma$ shipdate BETWEEN
'1994-01-01' AND '1995-01-01'
AND discount BETWEEN
0.05 AND 0.07
AND quantity < 24

$\pi$ quantity, discount
extendedprice, shipdate

**Scan**

lineitem

ndedprice * discount)

END) / SUM(extended

⋈

26

| | | | |
|---|---|---|---|
| partsupp | | | |
| supplier | 10.144428 | *** | |
| nation | 5.5904746 | *** | |
| region | 5.8037666 | *** | |

# Benchmark Results *

| | Standard Row | Trojan Columns | Trojan Columns (SP) | Standard Row | Troja |
|---|---|---|---|---|---|
| Q1 | 76.730296 | 19.293983 | 24.208052774 | 230.19089 | 57.8 |
| Q6 | 77.589034 | 8.6532381 | 11.235220175 | 232.7671 | 25.9 |
| Q12 | 92.486038 | 37.331905 | 40.598335758 | 277.45811 | 111. |
| Q14 | 81.207649 | 30.788114 | 59.597473787 | 243.62295 | 92.3 |
| Q3 | 111.88261 | 809.38127 | | 335.64782 | 2428 |
| Q5 | 99.729039 | 169.34457 | | 299.18712 | 508. |
| Q10 | 110.93664 | 119.46429 | | 332.80993 | 358. |
| Q19 | 79.140857 | 43.115296 | | 237.42257 | 129. |

3 72.41696 74.32579 98.73742 322.3892 2491.637

# Benchmark Results *

5x

| | Standard Row | Trojan Columns | Trojan Columns (SP) | Standard Row | Troja... |
|---|---|---|---|---|---|
| Q1 | 76.730296 | 19.293983 | 24.208052774 | 230.19089 | 57.8 |
| Q6 | 77.589034 | 8.6532381 | 11.235220175 | 232.7671 | 25.9 |
| Q12 | 92.486038 | 37.331905 | 40.598335758 | 277.45811 | 111. |
| Q14 | 81.207649 | 30.788114 | 59.597473787 | 243.62295 | 92.3 |
| Q3 | 111.88261 | 809.38127 | | 335.64782 | 2428 |
| Q5 | 99.729039 | 169.34457 | | 299.18712 | 508. |
| Q10 | 110.93664 | 119.46429 | | 332.80993 | 358. |
| Q19 | 79.140857 | 43.115296 | | 237.42257 | 129. |
| | 72.41696 | 74.32579 | 98.73742 | 322.3892 | 2491.637 |

1980s

# 1990s

# 2000s

# HYRISE

## 2010s

# 7 Vertical Partitioning Algorithms

- Brute Force

- Navathe's Algorithm

- HillClimb

- AutoPart

- HYRISE

- O$_2$P

- Trojan



sv-timemachine.net

# Four Comparison Metrics

- How Fast?

- How Good?

- How fragile?

- Where does it makes sense?

# Optimization Runtime

ance from Column Layouts

| | TPC-H | SSB |
|---|---|---|
| AutoPart | 3.71 | 5.29 |
| HillClimb | 3.71 | 5.29 |
| HYRISE | 1.58 | 5.27 |
| Navathe | -21.47 | 1.64 |
| O2P | -27.74 | 1.64 |
| Trojan | 3.71 | 0.05 |
| BruteForce | 3.71 | 5.29 |

32

# Effect of Buffer Size



33

# Comparison's Paper: Hadoop Vs PDBMS

# Comparison's Paper: Hadoop Vs PDBMS

# Analytical Query Performance

## Selection Task

## Join Task

# Trojan Index Creation



Indexed Split *i*

**Data Load Phase**

L, $T_1$, PPart_block, $T_6$
Replicate, Replicate, Replicate, Replicate, Replicate, Replicate

H1: Fetch → Store → Scan → PPart_split
Fetch → Store → Scan → PPart_split
Fetch → Store
H2: ...
H3: Fetch → Store → Scan → PPart_split
Fetch → Store
Fetch → Store
H4: ...

**Map Phase**

M1: Union → RecRead_itemize → MMap_map → PPart_mem
LPart_sh → Sort_cmp → SortGrp_grp → MMap_combine → Store
LPart_sh → Sort_cmp → SortGrp_grp → MMap_combine → Store
LPart_sh → Sort_cmp → SortGrp_grp → MMap_combine → Store
Merge_cmp → SortGrp_grp → MMap_combine → Store → PPart_sh

M2: ...

M3: RecRead_itemize → MMap_map → LPart_sh → Sort_cmp → SortGrp_grp → PPart_sh

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

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

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

$$\text{sh}(\text{key } k, \text{value } v, \text{int } numPartitions) \mapsto k.splitID \% numPartitions$$

**Shuffle Phase**

Fetch → Buffer → Store
Fetch → Buffer → Store
Fetch → Buffer
Fetch → Buffer → Merge → Store

**Reduce Phase**

R1: Merge_cmp → SortGrp_grp → MMap_reduce → Store

$$\text{reduce}(\text{key } ik, \text{vset } ivs) \mapsto [(ivs \oplus \text{indexBuilder}_{a_i}(ivs))]$$

$T'_1$   $T'_2$

37

# Trojan Index Access

# Trojan Index Access

IndexScan   IndexScan

of Relation T   of Relation S   of Relation T

{offset, record}   {offset, re

**Map**

{a/b, record}   {splitID+a, rec

**Reduce**

Distributed
File System

Data Load Phase

DataSet

F

Indexed Split *i*

DataSet

F

Rearranged Co-Partitioned Split *i*

DataSet

F

Indexed Split *i*

DataSet

F

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

Co-Partitioned Split *i*

Map Phase

DataSet

F

Rearranged Co-Partitioned Split *i*

DataSet

F

Indexed Co-Partitioned Split *i*

**Algorithm 3**: Trojan Index itemize.next UDF

**Input**  : KeyType key, ValueType value
**Output**: has more records

1  **if** *offset < splitEnd* **then**
2      Record nextRecord = ReadNextRecord(*split*);
3      offset += nextRecord.size();
4      **if** *nextRecord.key < highKey* **then**
5          SetKeyValue(*key, value, nextRecord*);
6          **return** true;
7      **end**
8  **end**
9  **return** false;

Shuffle Phase

SortGrp_grp
↓
MMap_combine
↓
Store
↓
PPart_sh

Store
↓
PPart_sh

Reduce Phase

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

Merge_cmp
↓
SortGrp_grp
↓
MMap_reduce
↓
**R1**   Store

**R2**   . . .

$T'_1$   $T'_2$

# Selection Analytical Task *



* Pavlo et. al. A Comparison of Approaches to large-Scale Data Analysis. SIGMOD 2009

# Join Analytical Task *

* Pavlo et. al. A Comparison of Approaches to large-Scale Data Analysis. SIGMOD 2009

# Trojan Index

# Trojan Join

# Traditional Layouts

Row
(default)

Column*

PAX**

\* A. Floratou et al. Column-Oriented Storage Techniques for MapReduce. PVLDB, April, 2011

\*\*Y. He et al. RCFile: A fast and space-efficient data placement structure in MapReduce-based warehouse systems. ICDE, 2011

# Traditional Layouts

| | Row | Column | PAX |
|---|---|---|---|
| Non-required Reads | 🟥 | 🟩 | 🟩 |
| Network Costs | 🟩 | 🟥 | 🟩 |
| Data Block Placement | 🟩 | 🟥 | 🟩 |
| Tuple Reconstruction | 🟩 | 🟥 | 🟥 |

# Trojan Data Layouts

Replica 1

Replica 2

Replica 3

# Trojan Data Layouts

| | Row | Column | PAX | Trojan |
|---|---|---|---|---|
| Non-required Reads | 🟥 | 🟩 | 🟩 | 🟩 |
| Network Costs | 🟩 | 🟥 | 🟩 | 🟩 |
| Data Block Placement | 🟩 | 🟥 | 🟩 | 🟩 |
| Tuple Reconstruction | 🟩 | 🟥 | 🟥 | 🟩 |

# Layout Quality

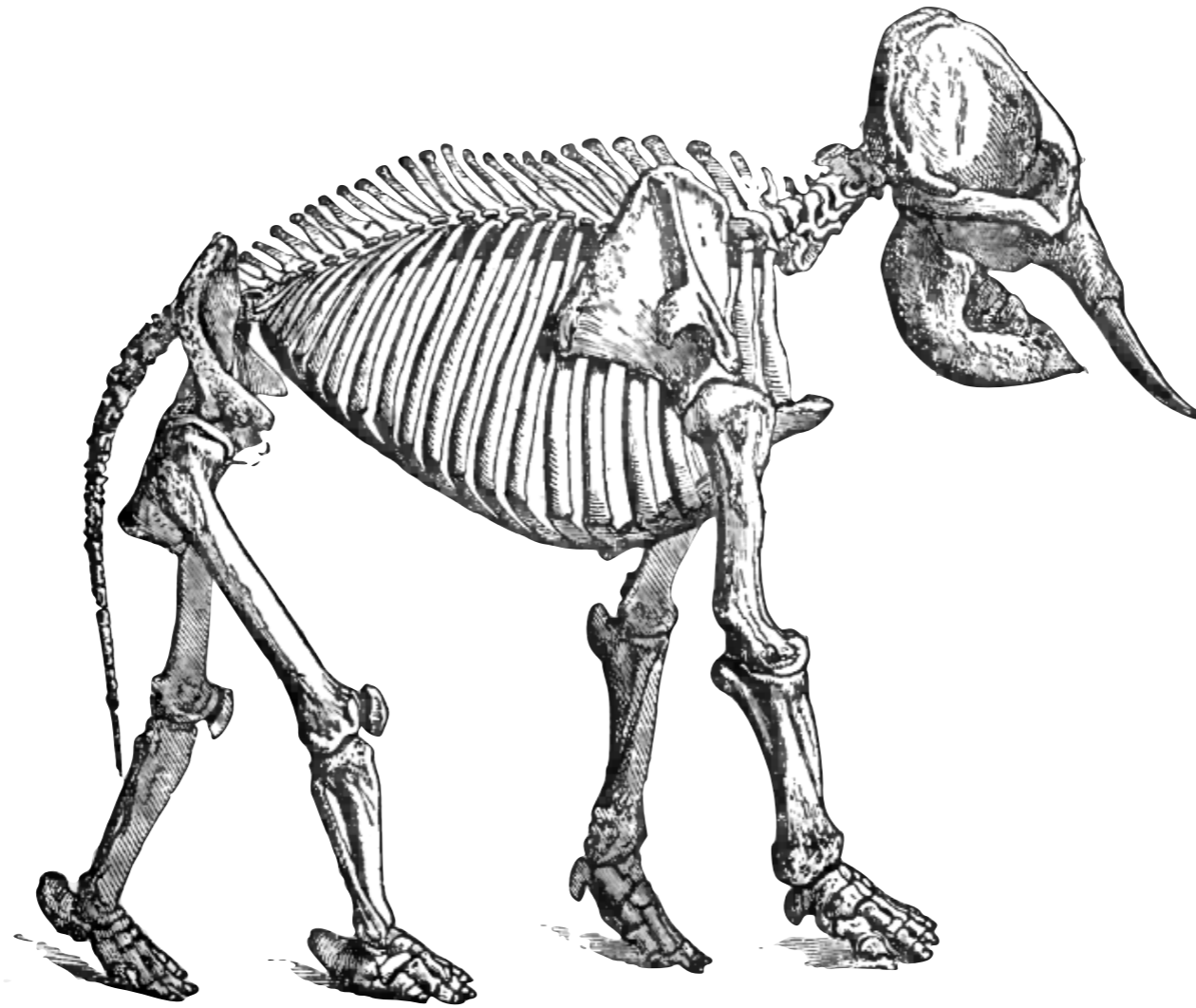| | #Non-required Attributes Read | #Joins in Tuple Reconstruction |
|---|---|---|
| HADOOP-ROW | 525 | 0 |
| HADOOP-PAX | 0 | 139 |
| Trojan Layout | 14 | 20 |

# Projection Analytical Task
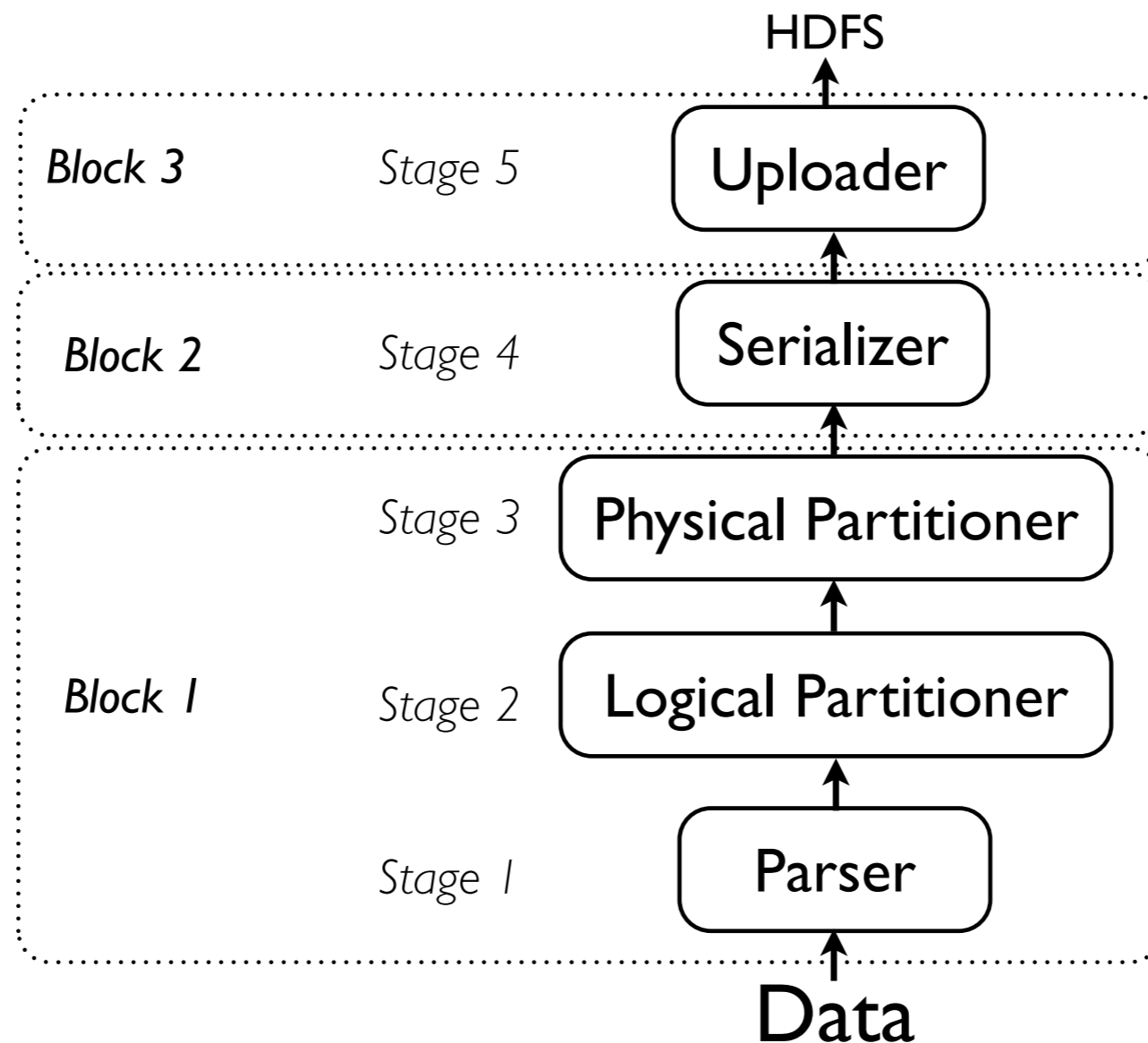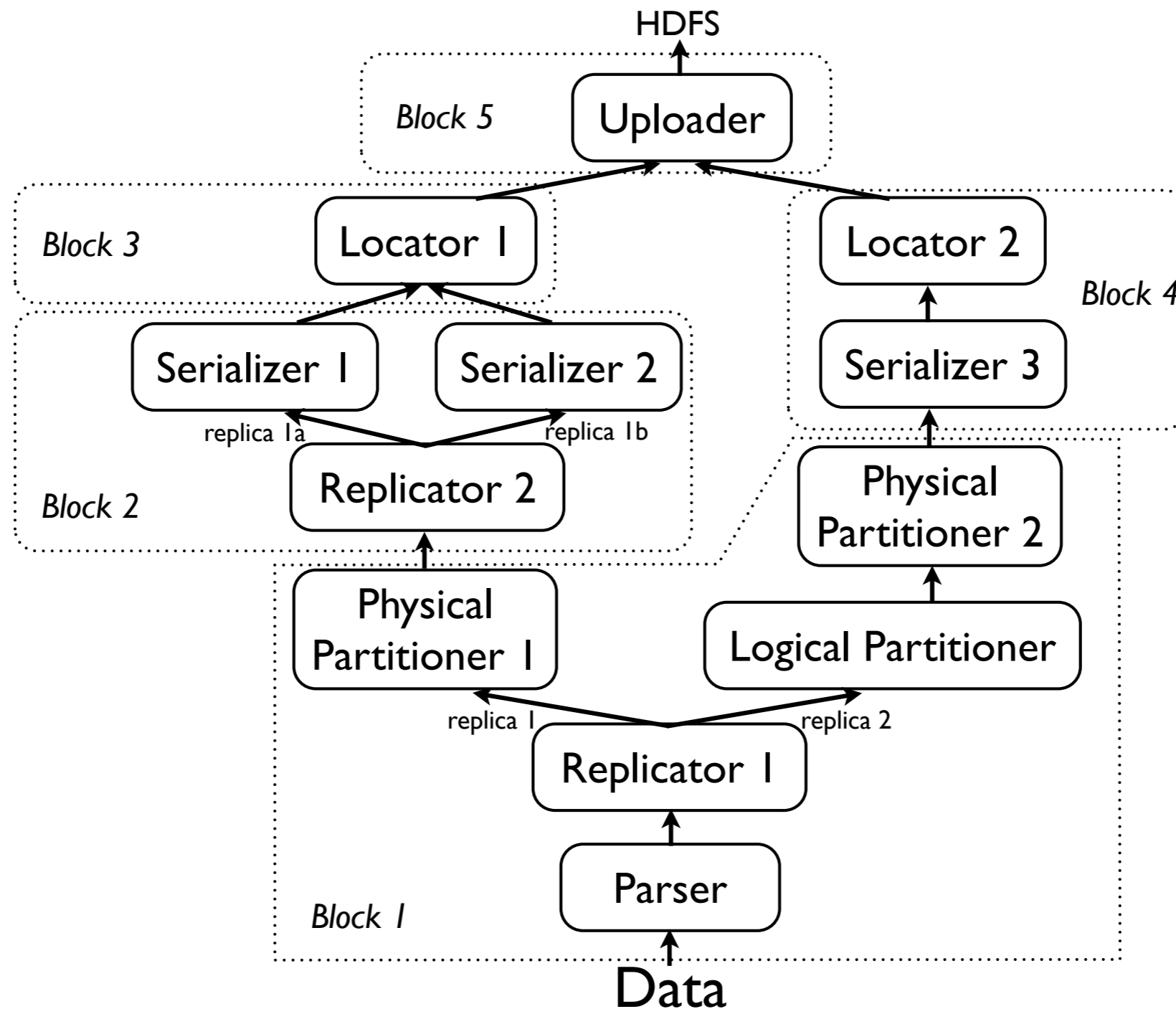
# Individual Jobs: Weblog, RecordReader

# Cartilage

# Cartilage

# Hadoop Stack

HBase

Hive

Pig

MapReduce

**Cartilage Query Engine**

HDFS

**Cartilage Upload Pipeline**

| Data File 1 | Data File 2 | ... | Data File n |

# Hadoop Stack

Queried Data

**Cartilage Query Engine**

HDFS

**Cartilage Upload Pipeline**

Input Data

# Upload Plans

# Upload Plans



HDFS

*Block 5* — Uploader

*Block 3* — Locator 1

Locator 2 — *Block 4*

Serializer 1    Serializer 2

Serializer 3

replica 1a    replica 1b

*Block 2* — Replicator 2

Physical Partitioner 2

Physical Partitioner 1    Logical Partitioner

replica 1    replica 2

Replicator 1

Parser

*Block 1*

Data

57

# Summary

ONE SIZE DOES NOT FIT ALL

CIDR 2011

WWHow! Layer

HYRISE

CIDR 2013
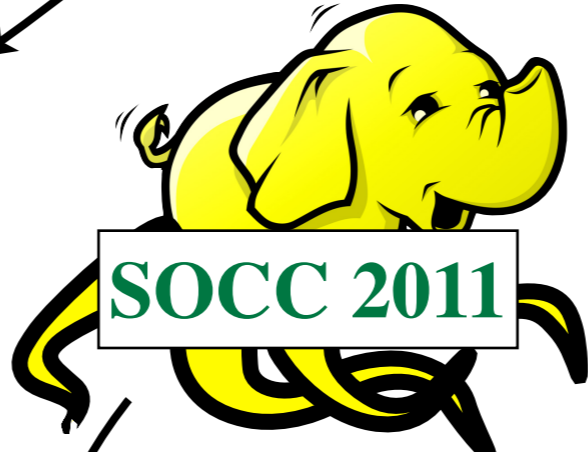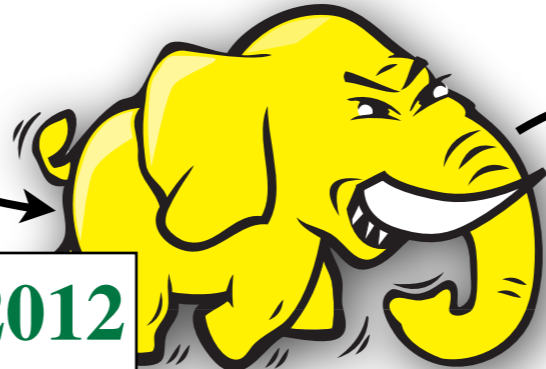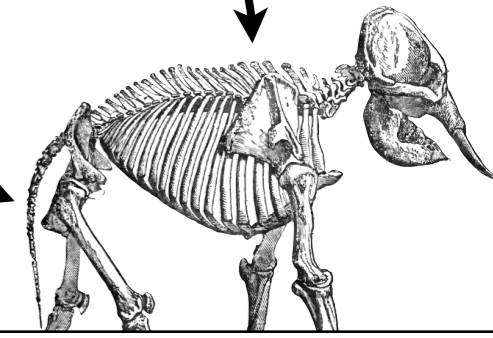
VLDB 2013

CIDR 2013

SOCC 2011

VLDB 2010

VLDB 2012

SIGMOD (demo)

# Acknowledgements

- Jens Dittrich

- Jorge Quiane

- Felix Martin Schuhknecht

- Endre Palatinus

- Karen Khachatryan

- Stefan Richter

- Alexander Bunte

- Sam Madden

- Stefan Richter

- Stefan Schuh

- Joerg Schad

- Yagiz Kargin

- Vinay Setty

- Vladimir Pavlov