

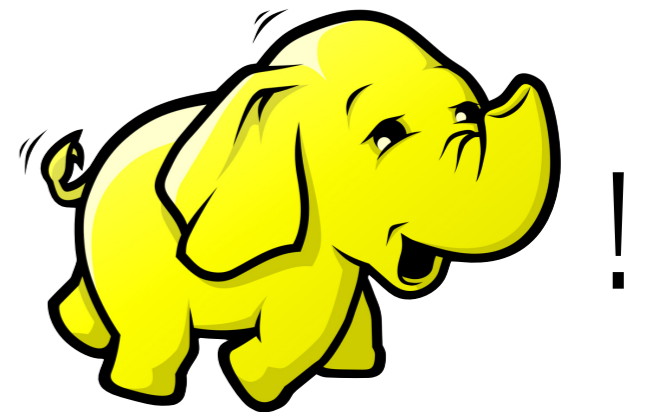
7

Things To Know

When Buying



for an



!

1

What Shoes? Why Shoes?

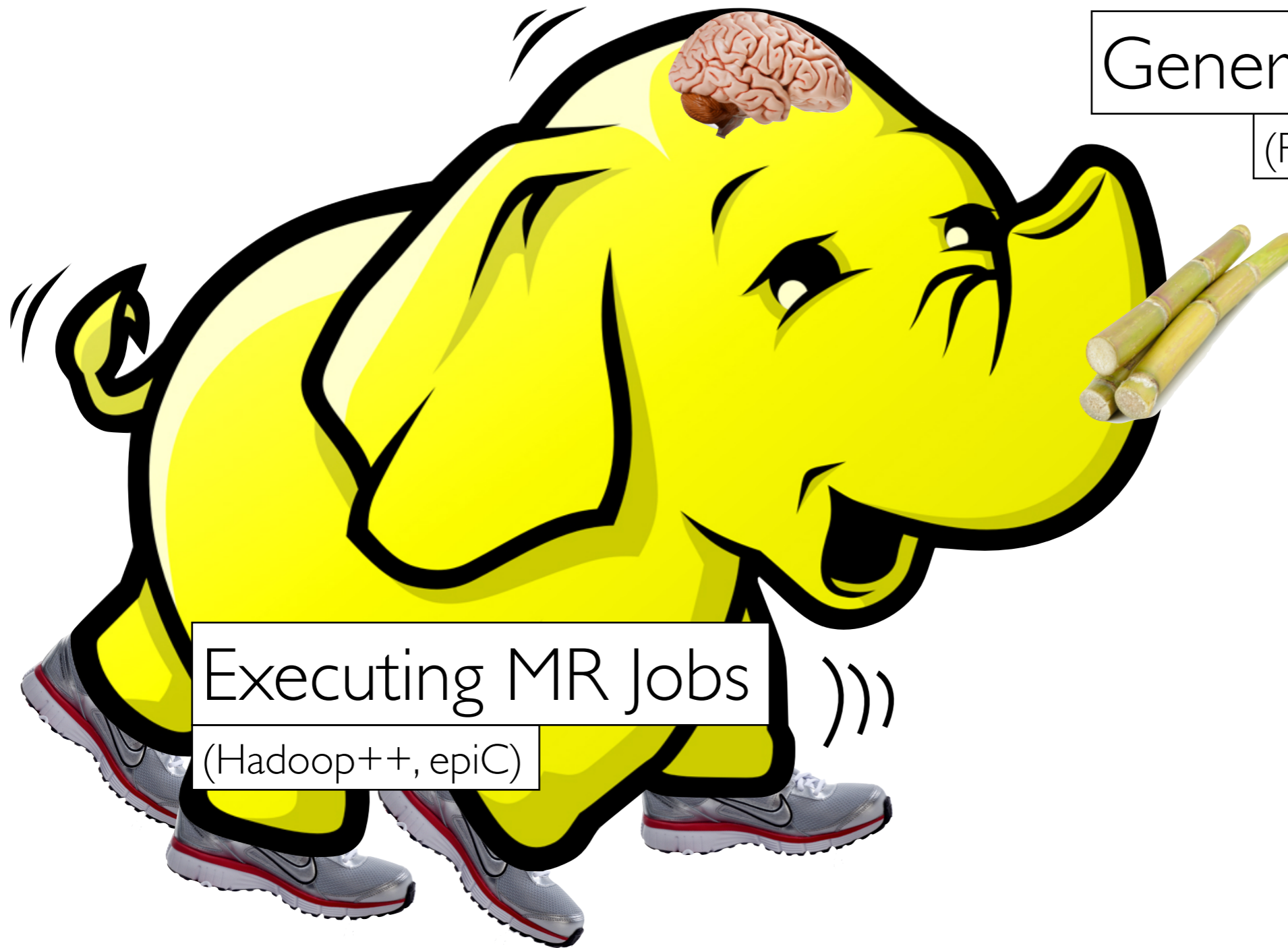


Analyzing MR Jobs

(HadoopToSQL, Manimal)

Generating MR Jobs

(PigLatin, Hive)



Executing MR Jobs

(Hadoop++, epiC)

Data Layouts & Access Paths !!

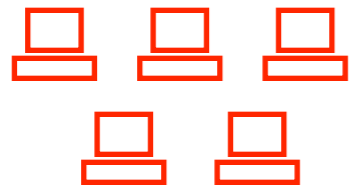
2

Why Elephant Needs Different Shoes?

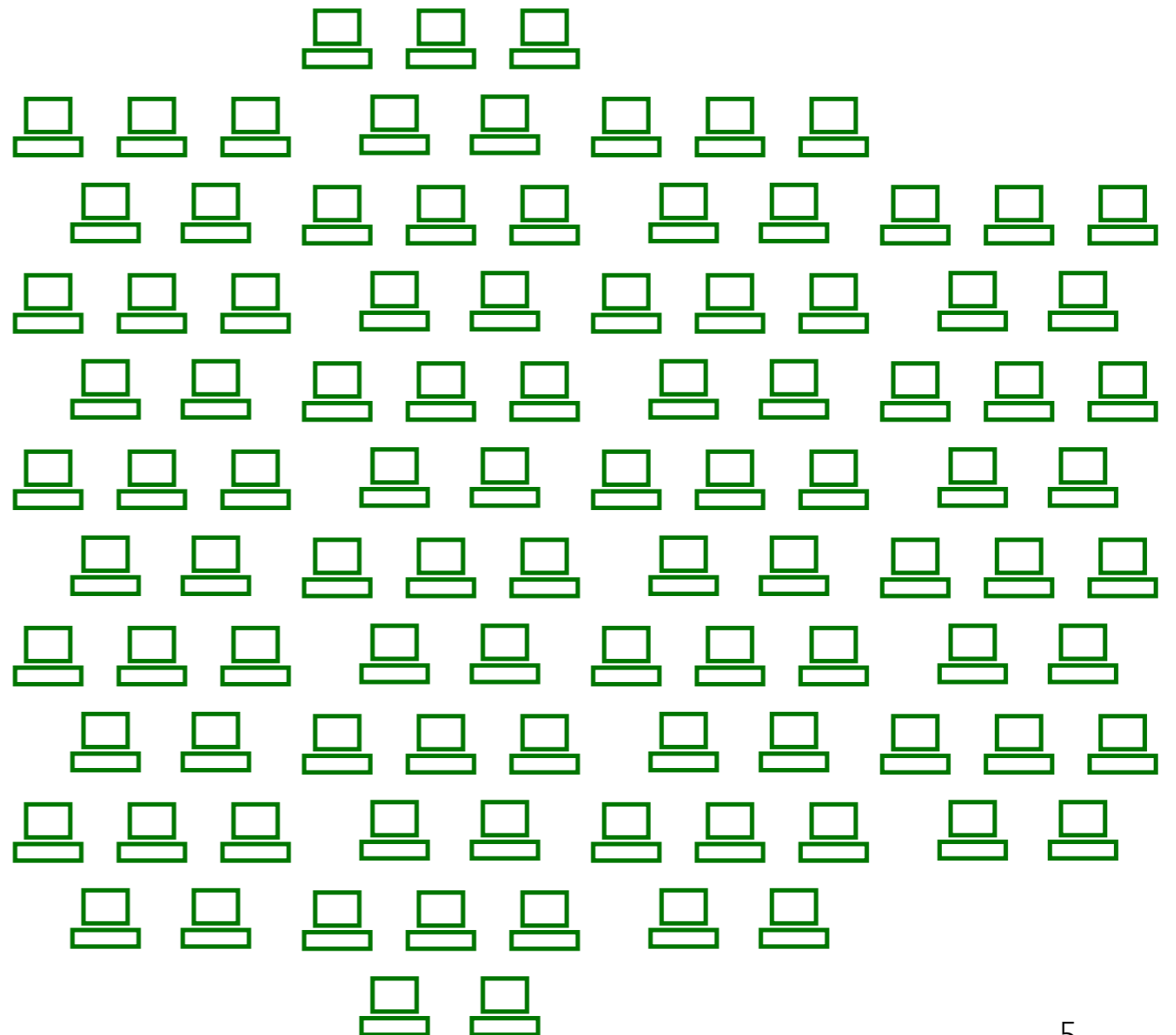


Very Large Scale Storage & Execution

DBMS



MapReduce



Large Data Block Sizes

DBMS

8 KB



MapReduce

1 GB



Block Level Data Replication

DBMS

001 alex bsc
002 tim msc
003 mat bsc
004 joel bsc
005 phil msc
006 ron msc
007 neo bsc
008 jack msc
009 jens bsc
010 tom msc

MapReduce

001 alex bsc
002 tim msc

003 mat bsc
004 joel bsc

005 phil msc
006 ron msc

007 neo bsc
008 jack msc

009 jens bsc
010 tom msc

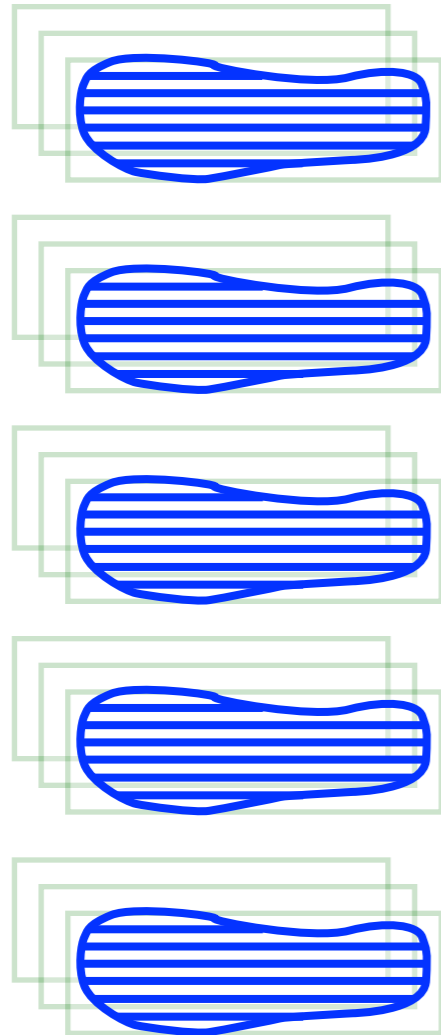
3

What's Wrong with Old Shoes?

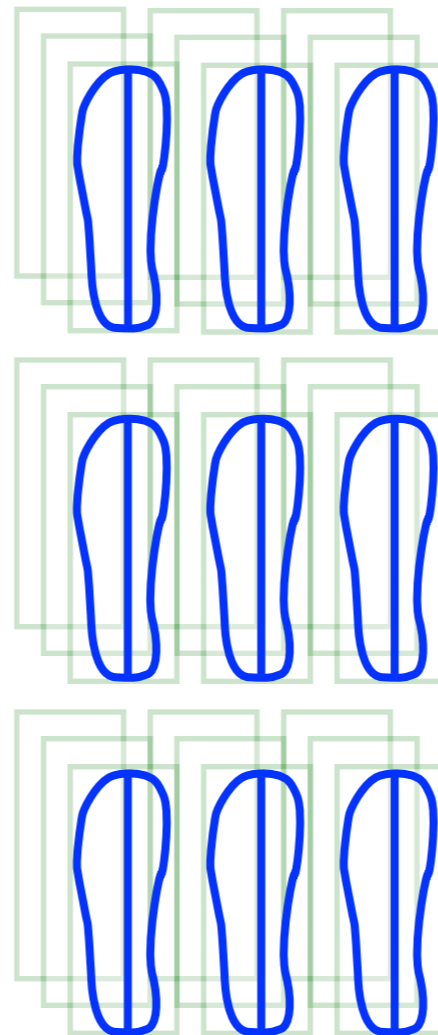


Current Data Layouts in Hadoop

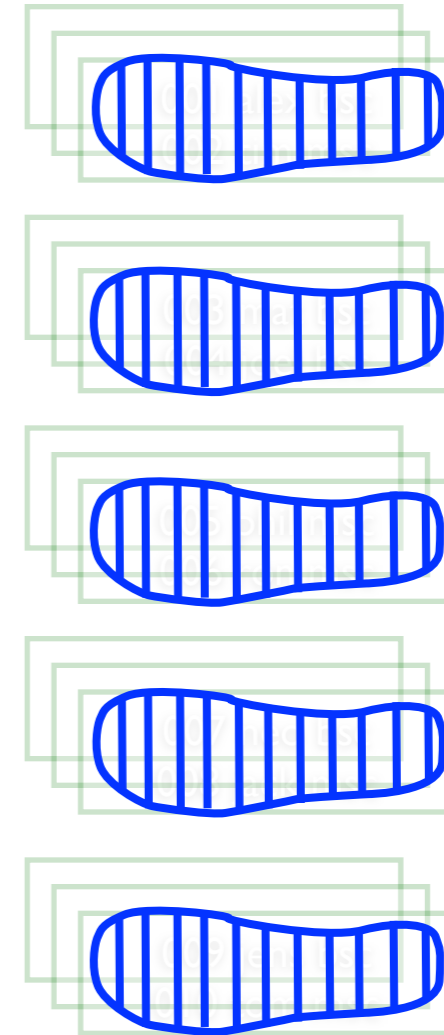
Row
(default)



Column*















PAX**



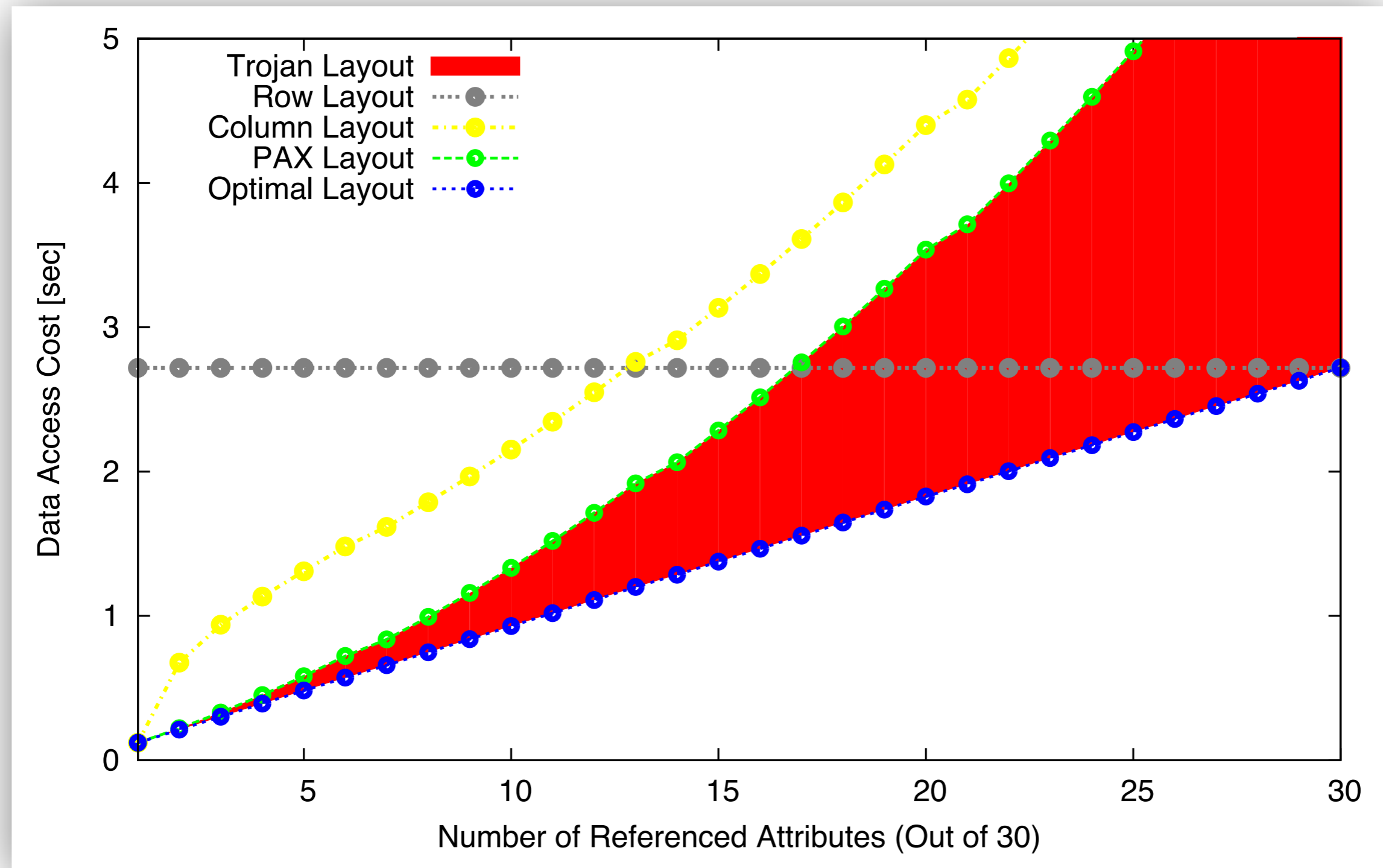
* A. Floratou et al. Column-Oriented Storage Techniques for MapReduce. PVLDB, April, 2011

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

Current Data Layouts in Hadoop

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

Current Data Layouts in Hadoop



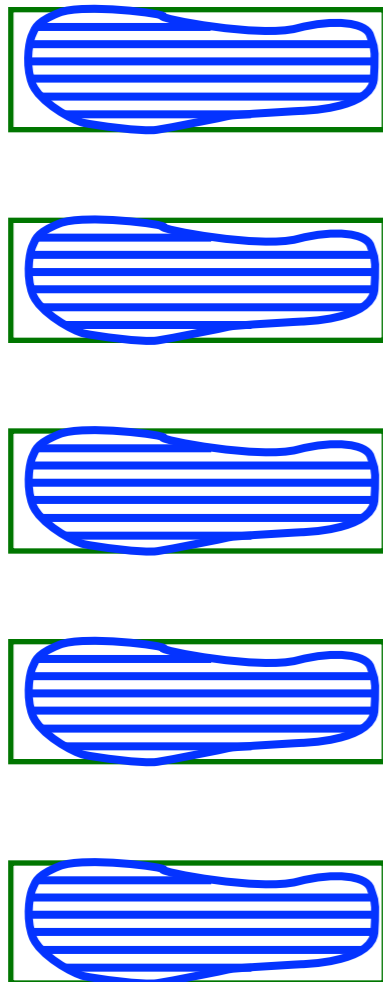
4

What Shoes do We Propose?

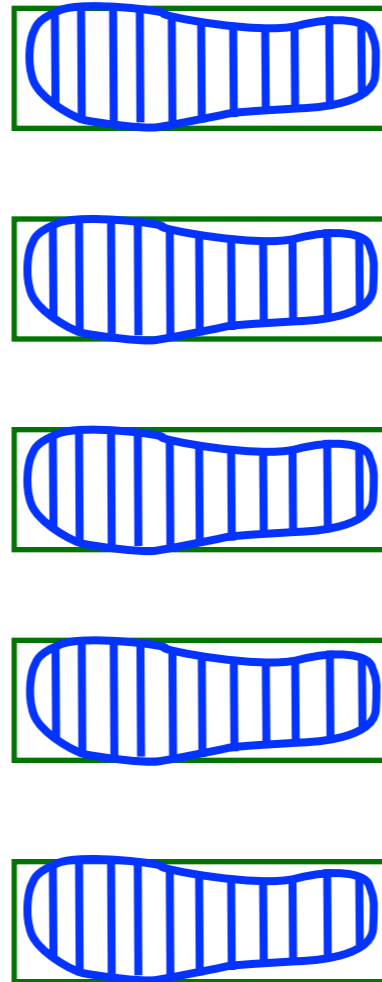


Trojan Data Layouts

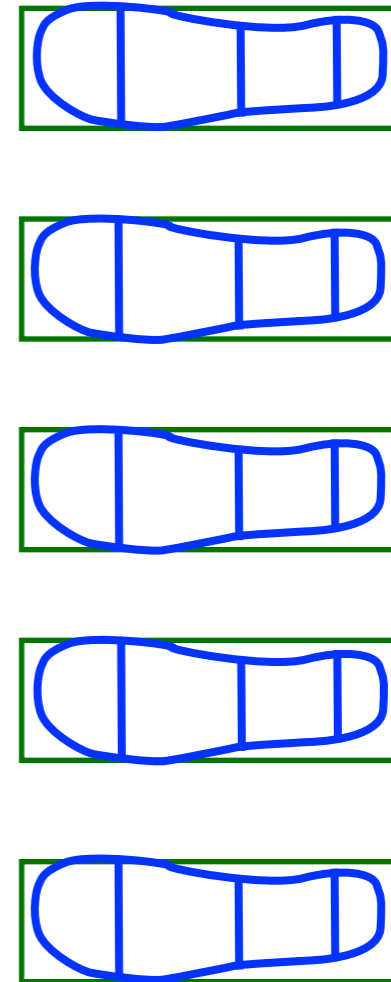
Replica 1



















Replica 2



Replica 3



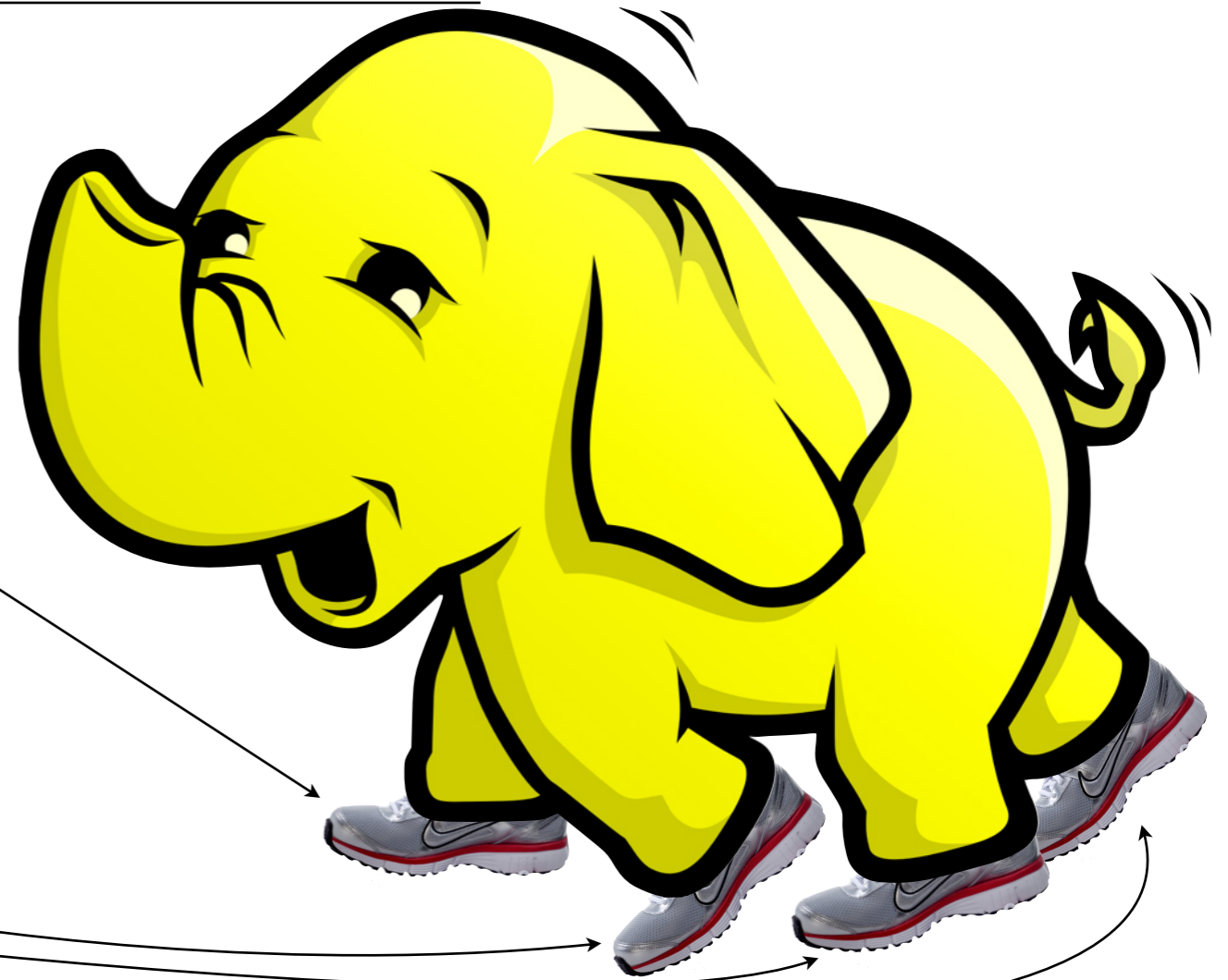
Trojan Data Layouts

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

Challenges in Trojan Data Layouts

How do we design shoe for one leg?

How do we design shoes for all legs?



How do we **make** the shoes from the design?



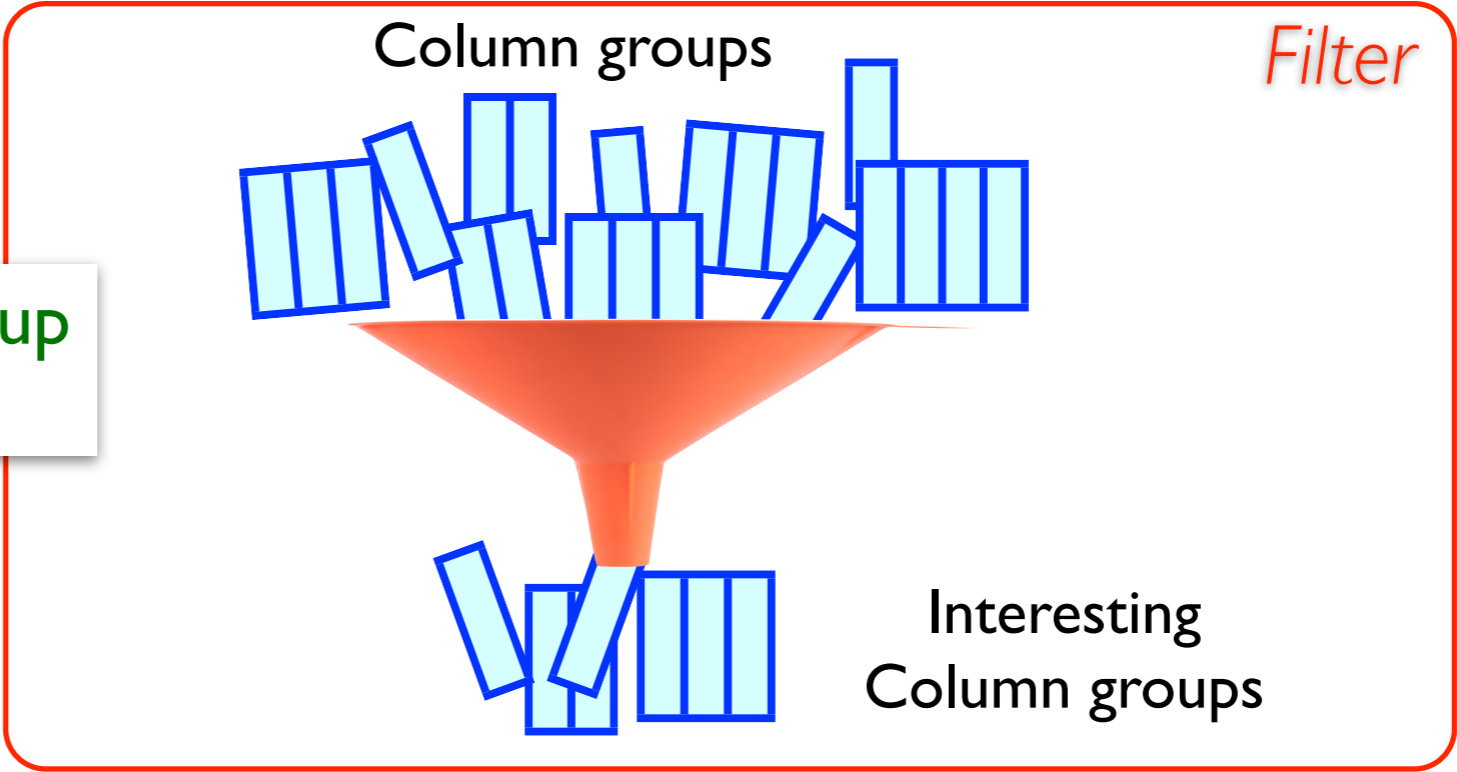
How Do We Design the Shoes?



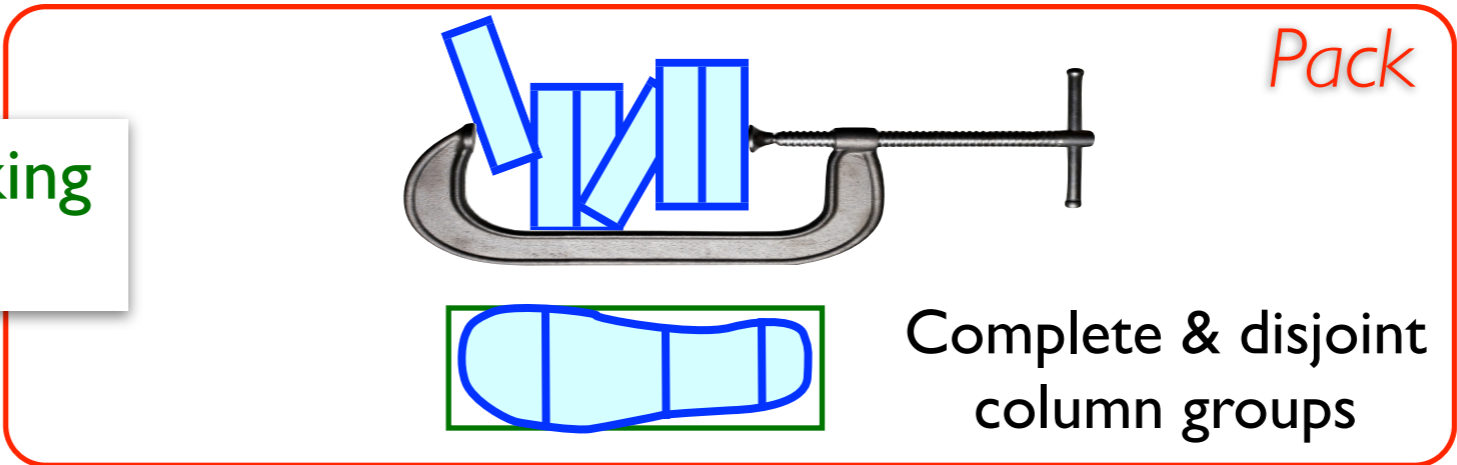
Single Replica



Novel Column Group Interestingness



Column Group Packing as 0-1 Knapsack



Multiple Replicas

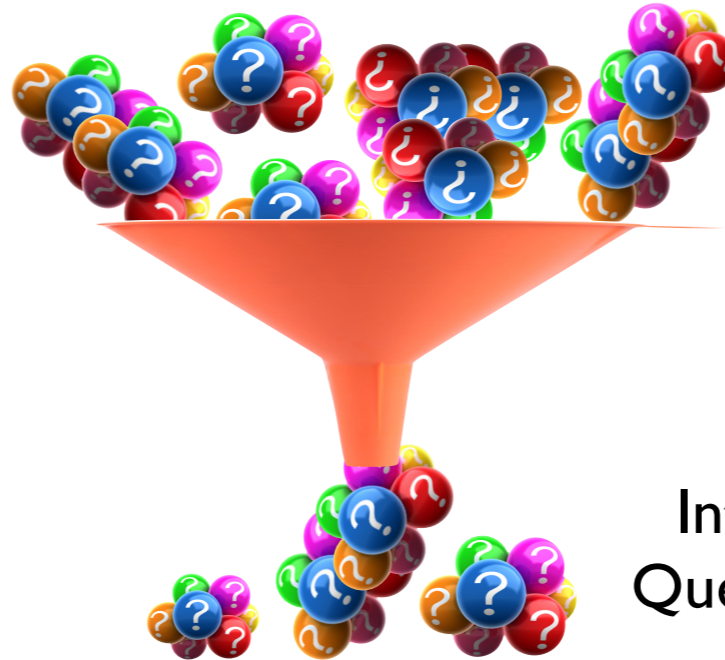


Queries



Query groups

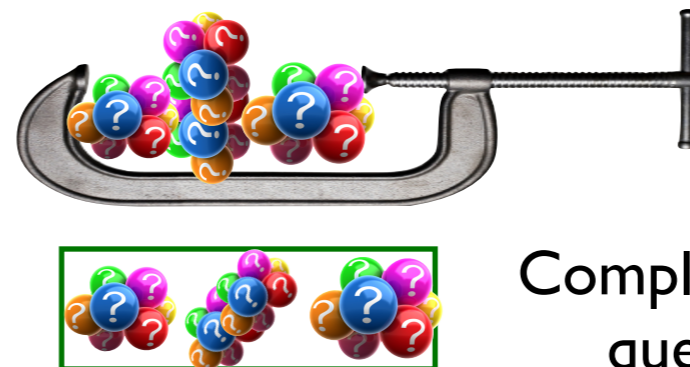
Filter



Interesting
Query groups

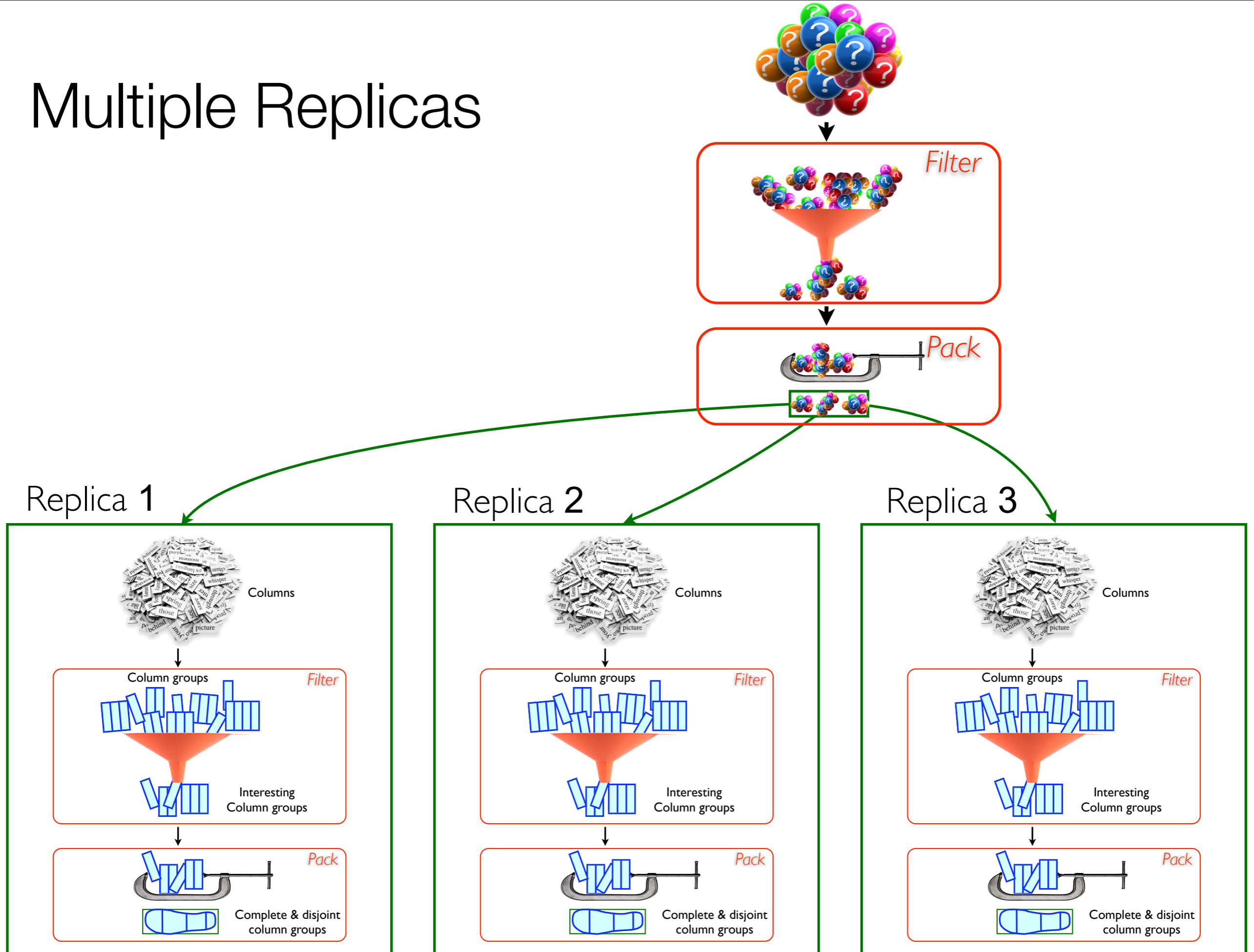


Pack



Complete & disjoint
query groups

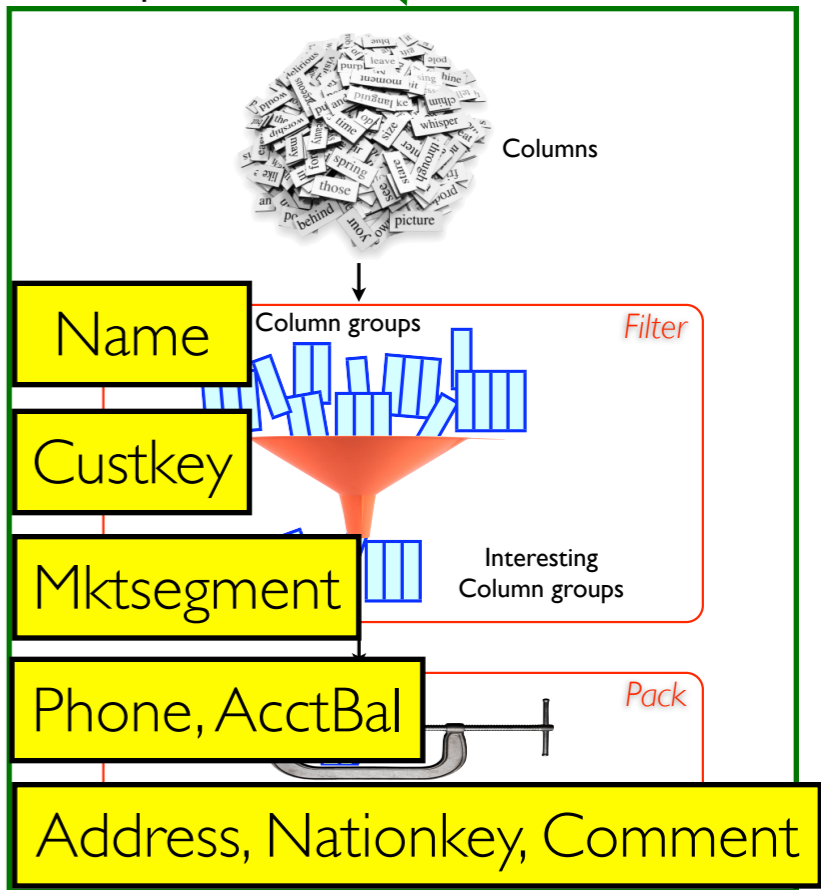
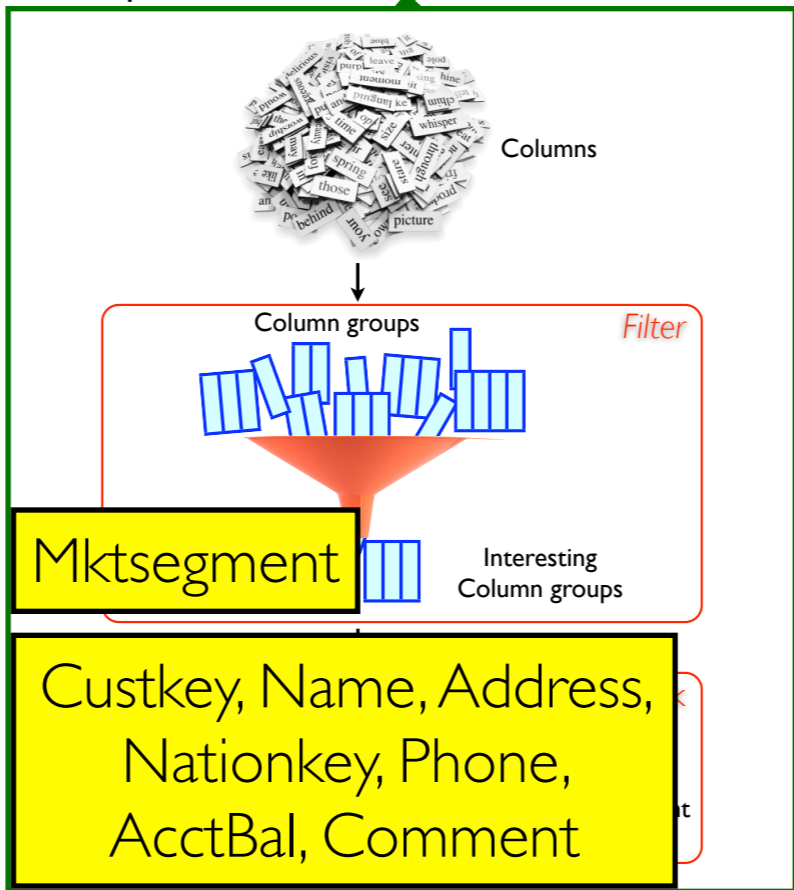
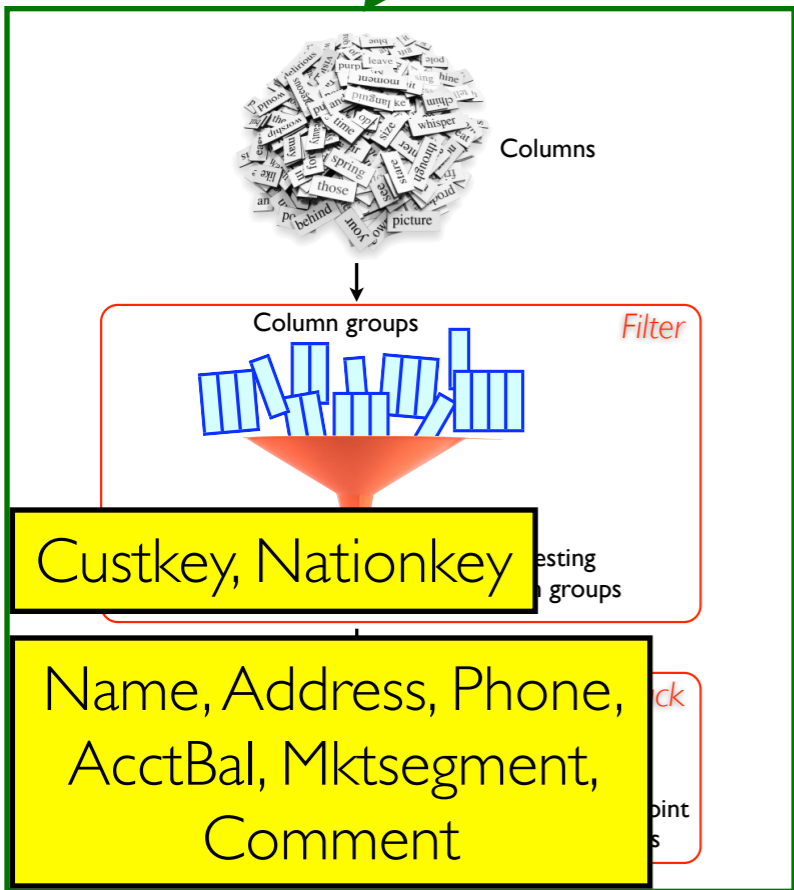
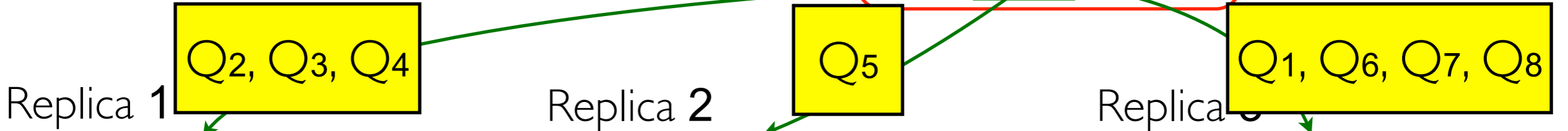
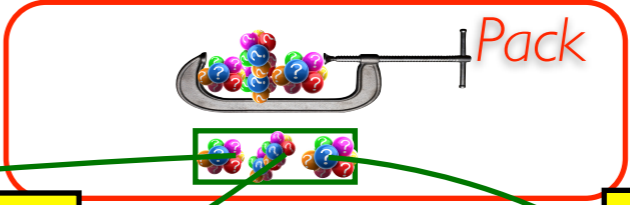
Multiple Replicas



Multiple Replicas

TPC-H Customer

Q1, Q2, Q3, Q4, Q5, Q6, Q7, Q8



Trojan Layout Advantages

- Multiple layouts for a given workload
- Default row layout still available
- Specialized replicas for different query sub-class
- Divide and conquer layout computation



How do We Ride the Elephant?



Putting It All Together



Load

Create trojan layout configuration file in HDFS
dataset layout-1 layout-2 layout-3



Query

Supply referenced attributes in JobConf
itemize UDF to transparently read the referenced attributes



Schedule
?

Three Optimization Options:

- data locality (default)
- best layout
- best layout & locality

7

How were the Field Trials?



Setup

- Datasets

TPC-H Lineitem, TPC-H Customer, SSB LineOrder, SDSS PhotoObj

- Queries

First **8** queries from the respective benchmark for each table

- Methodology

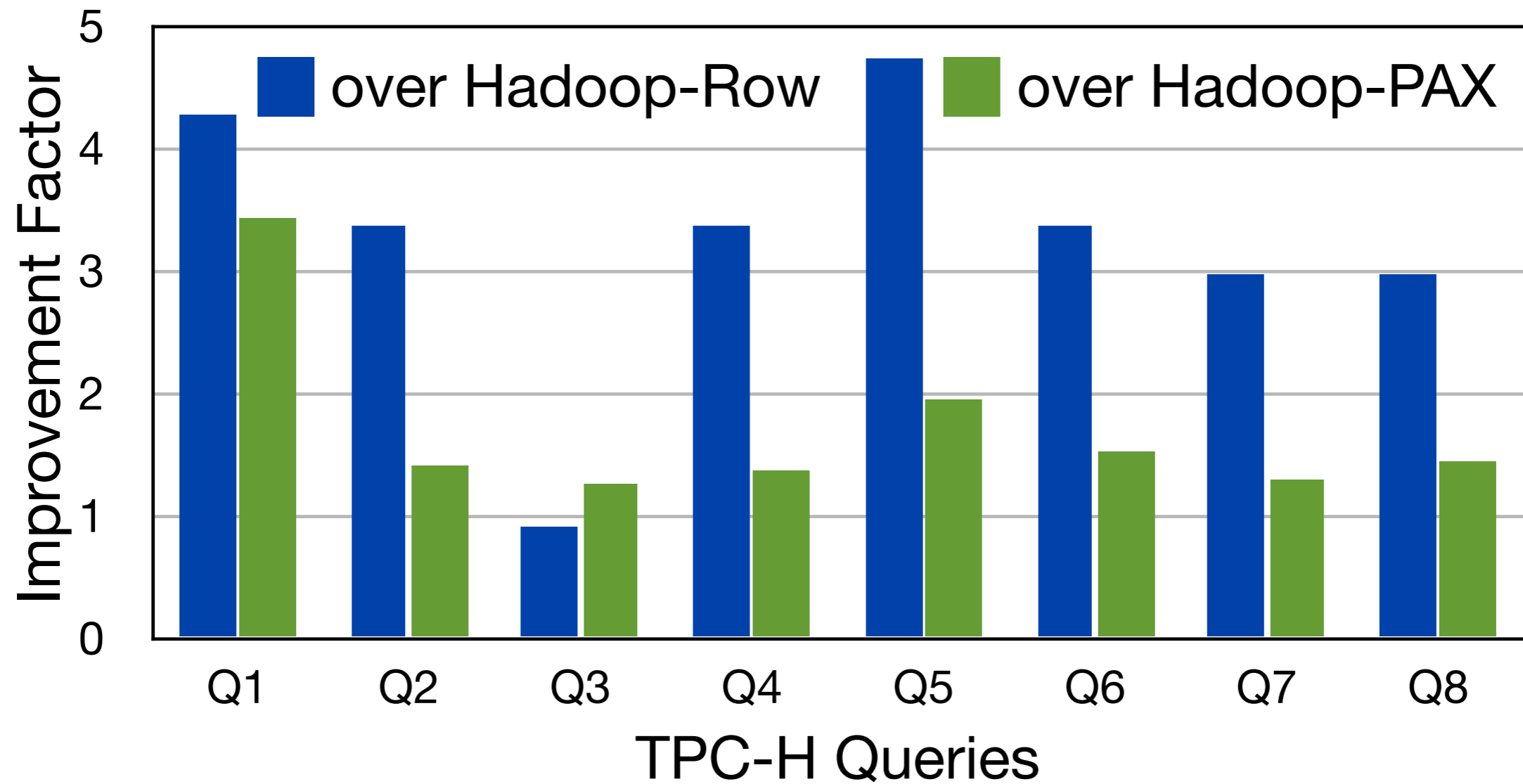
focus on scan and projection operators i.e. map-phase-only jobs
improvement: record reader time (I/O and tuple reconstruction)

- Hardware

50 virtual nodes in a **10** node cluster

Per-replica Trojan Layout Performance

TPC-H Lineitem



Layout Quality

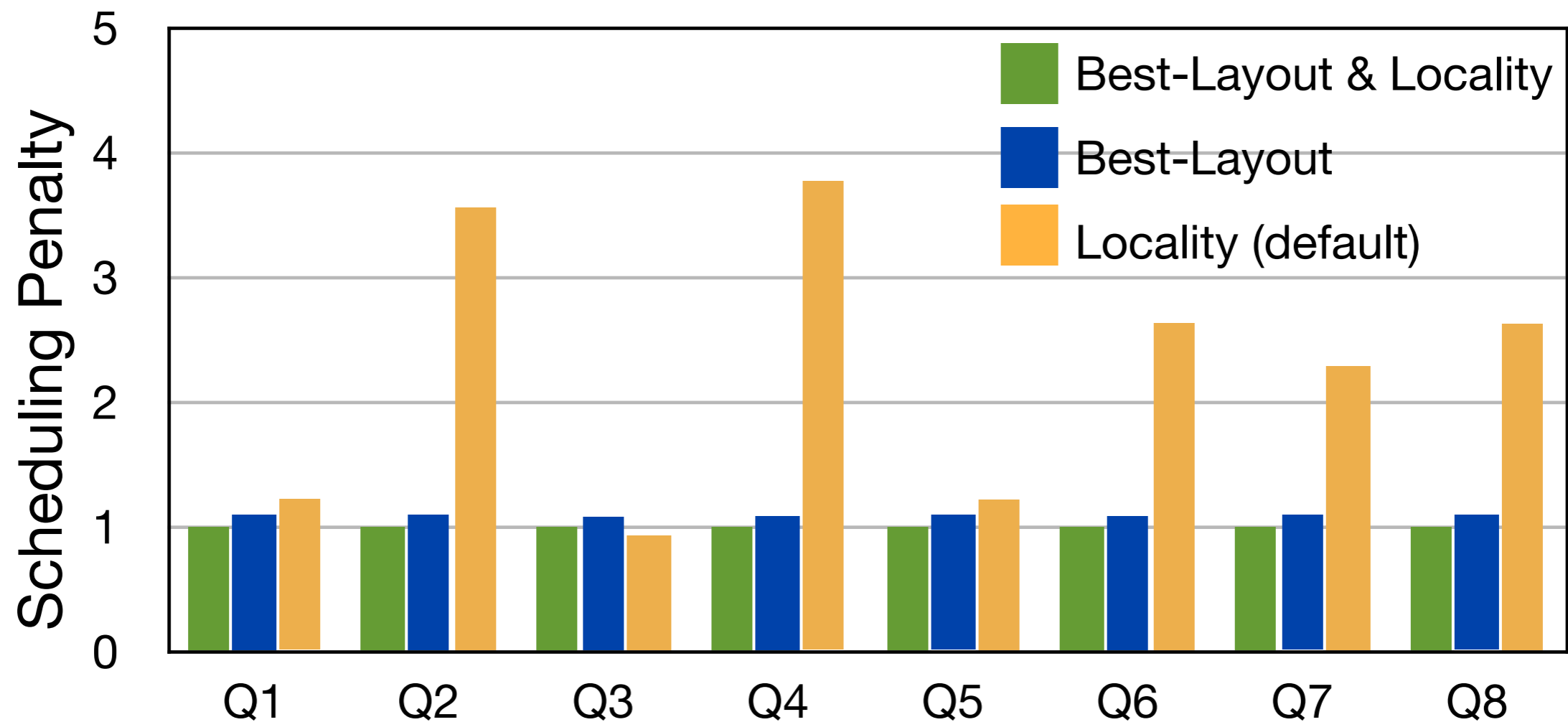
	#Non-required Attributes Read	#Joins in Tuple Reconstruction
HADOOP-ROW	525	0
HADOOP-PAX	0	139
HYRISE* Layout	2	64
Trojan Layout	14	20

>14% improvement over HYRISE

* M. Grund et al. HYRISE - A Main Memory Hybrid Storage Engine. PVLDB, November, 2010.

Scheduling Decisions

TPC-H Lineitem



Summary

- Data layouts crucial to MR job performance
- Exploit default data block replication in MR
- Novel algorithm to compute per-replica layouts
- Improvement: **4.8x** over Row, **3.5x** over PAX
- Better than HYRISE; **14%** improvement