

Large Scale Text Analysis Using the Map/Reduce Hierarchy

David Buttler

This work is performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344

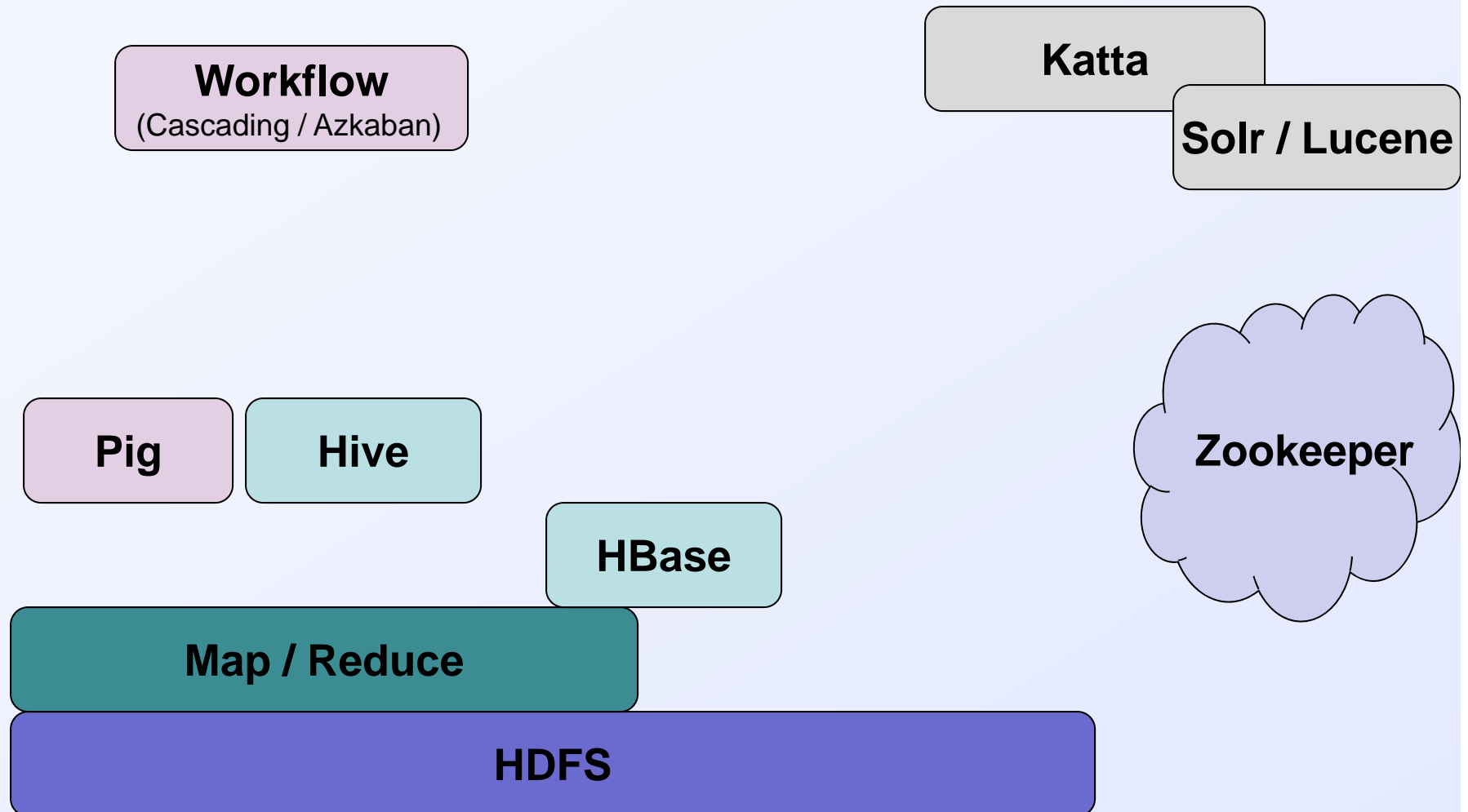
[Lawrence Livermore National Laboratory](#)

Large scale computing with commodity hardware

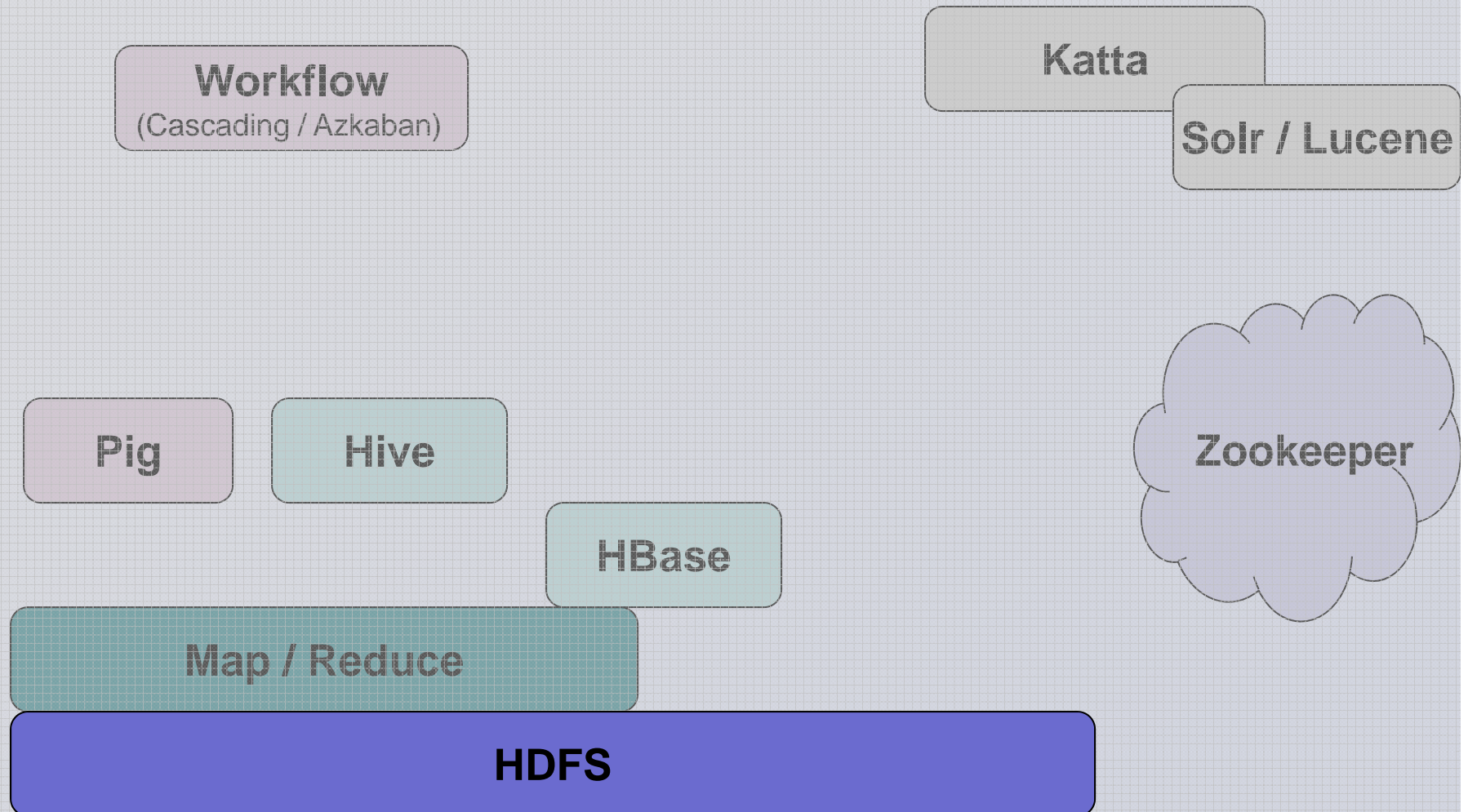
- Origins
 - Google GFS, Map/Reduce, BigTable
 - Microsoft Azure
 - **Hadoop**: Yahoo! / Open Source Software
- Why do we care: k-mer lexing
 - 10 hours on a single fat node
 - 1 hour on an old cluster



The M/R stack of open source software

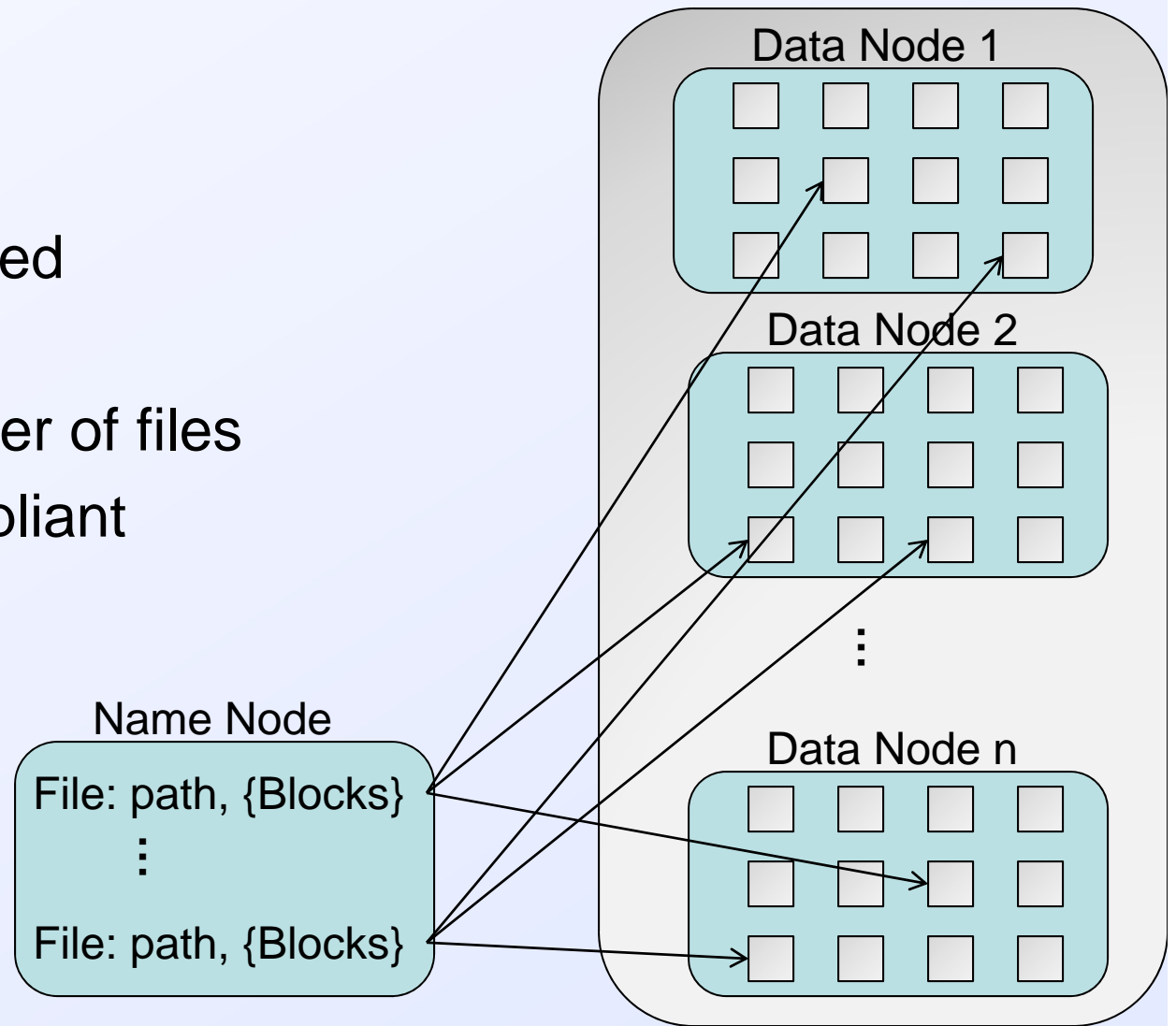


The M/R stack of open source software – HDFS

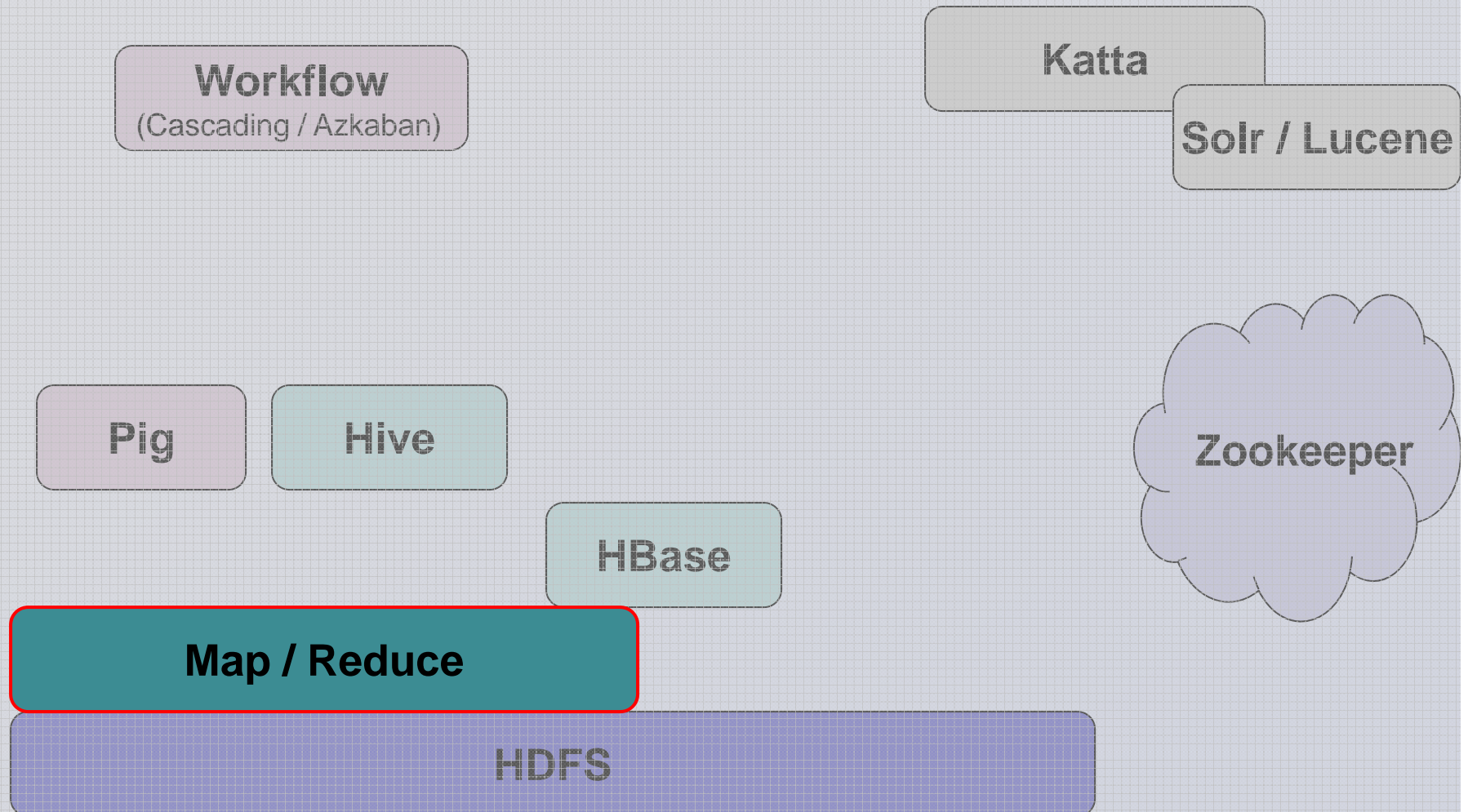


HDFS

- Replicated
- Distributed
- Centrally managed
 - SPOF
 - Limited number of files
- Not POSIX compliant
- Rack-aware

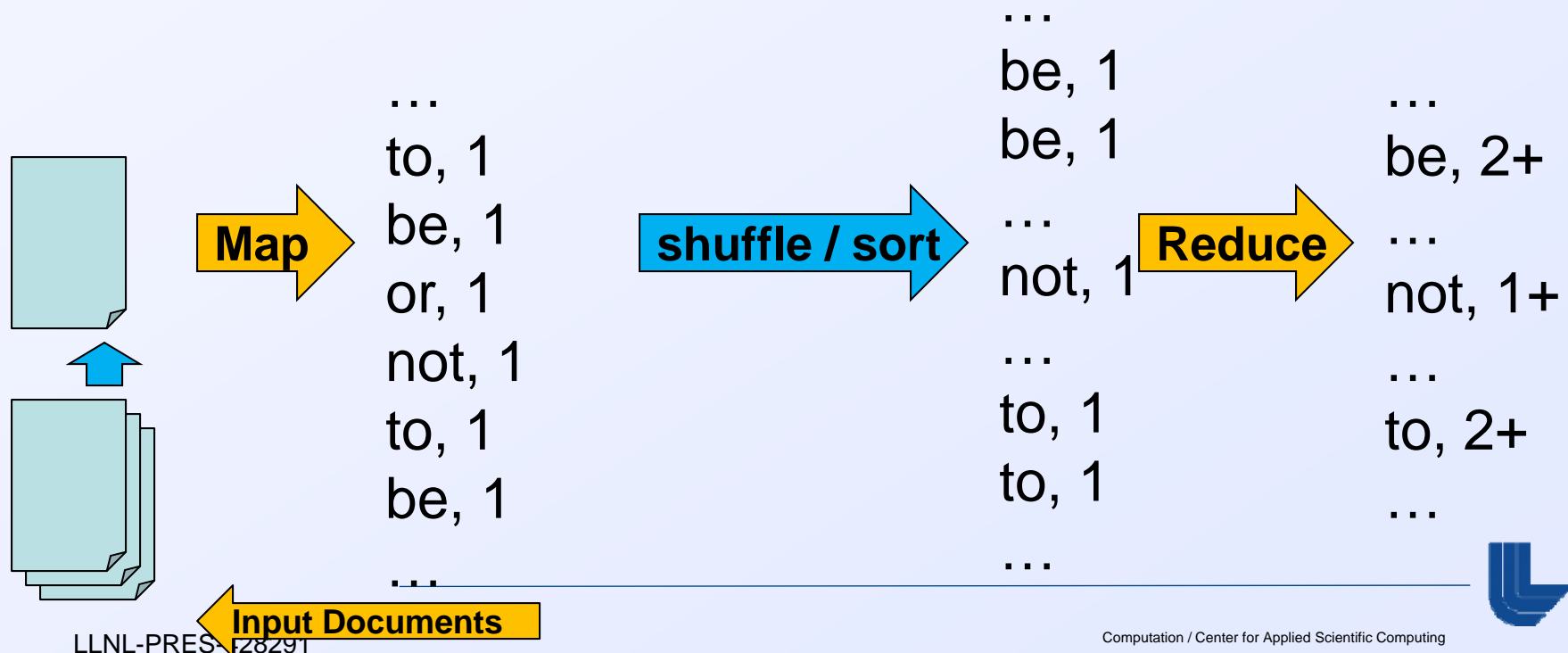


The M/R stack of open source software – Map / Reduce



Map/Reduce is functional programming distributed over a cluster

- Distributed computation
- Two phase computation
- Built-in shuffle/sort between phases
- Canonical example: word frequency count for the web

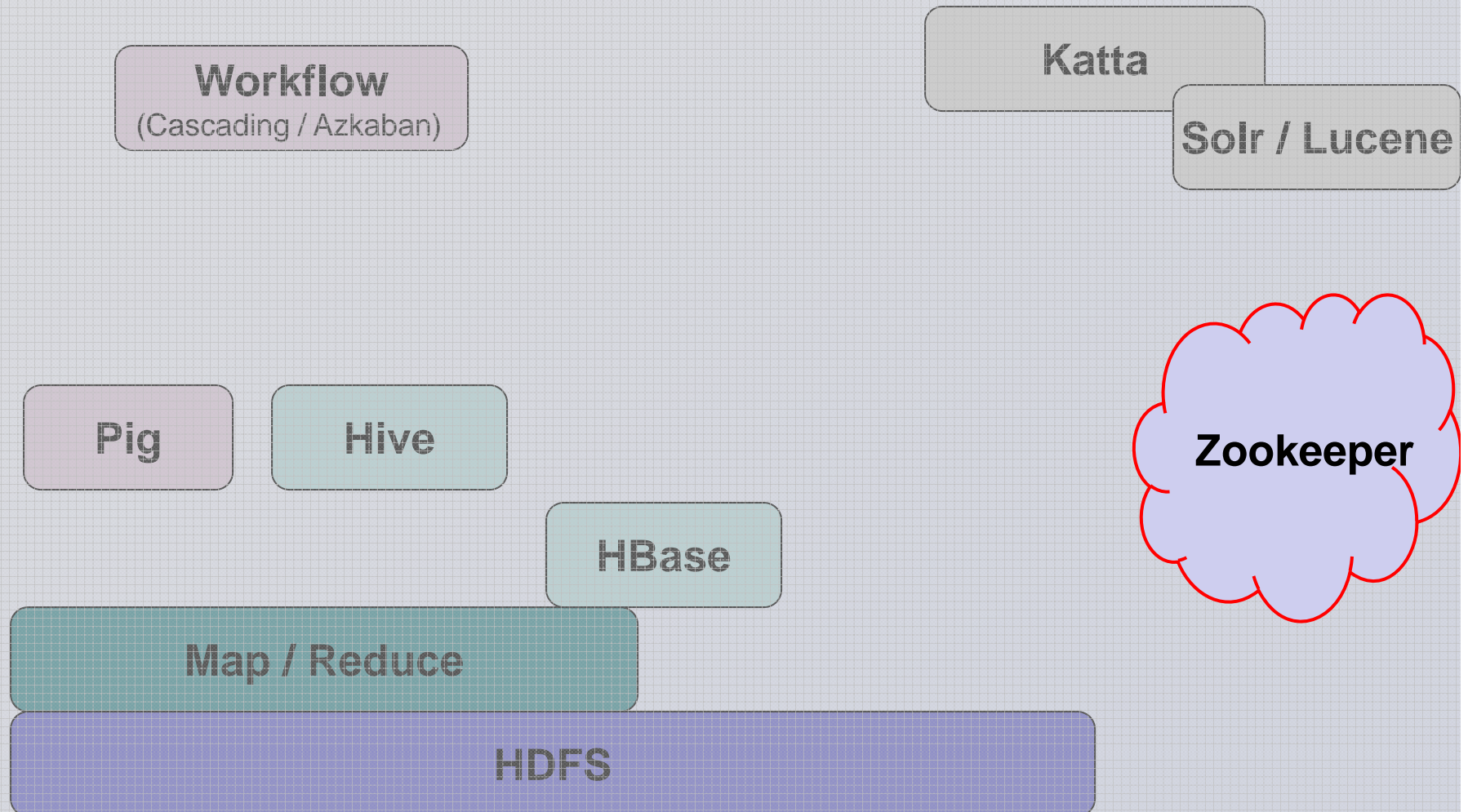


More interesting M/R examples

- Map input: document
- Map output: raw text
- Map input: text
- Map output: Named entity annotations



The M/R stack of open source software – Zookeeper



Zookeeper

- A highly available, scalable, distributed, configuration, consensus, group membership, leader election, naming, and coordination service
- Uses:
 - HBase: row locking; region key ranges; region server addresses
 - Katta: shard location information
 - Message queues
- Not: a large scale data store



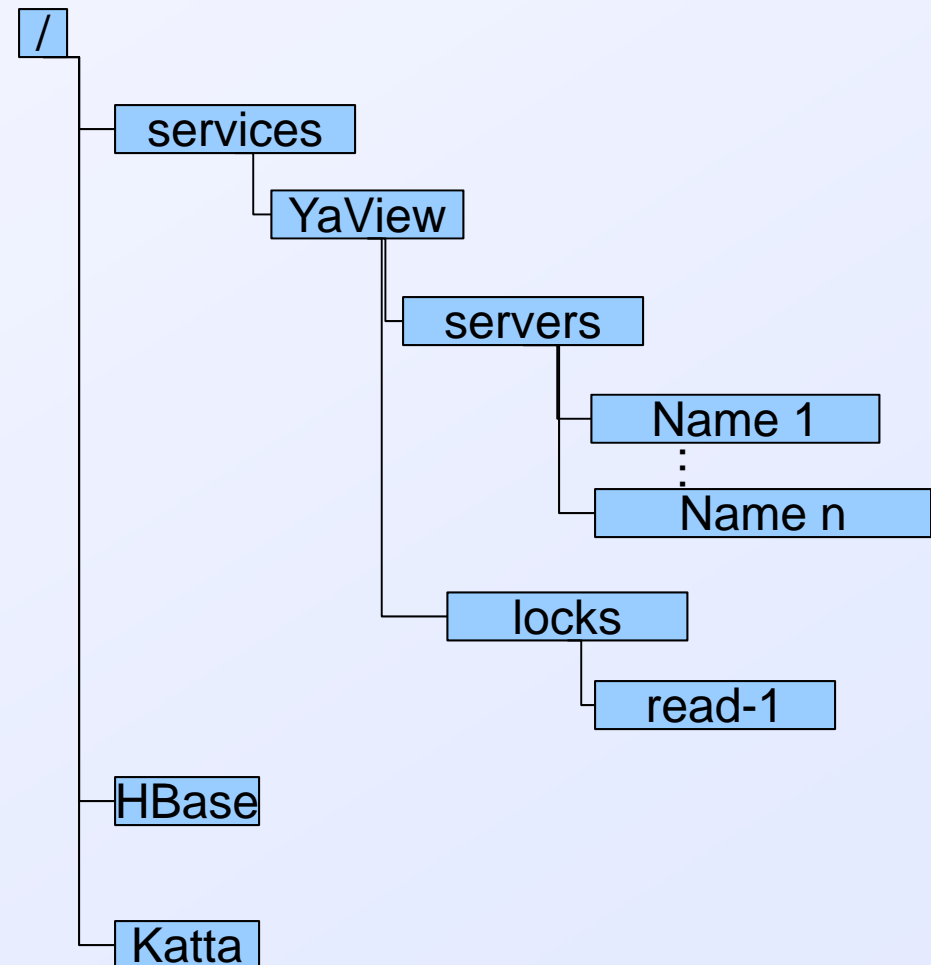
Zookeeper Guarantees

1. Clients will never detect old data.
2. Clients will get notified of a change to data they are watching within a bounded period of time.
3. All requests from a client will be processed in order.
4. All results received by a client will be consistent with results received by all other clients.

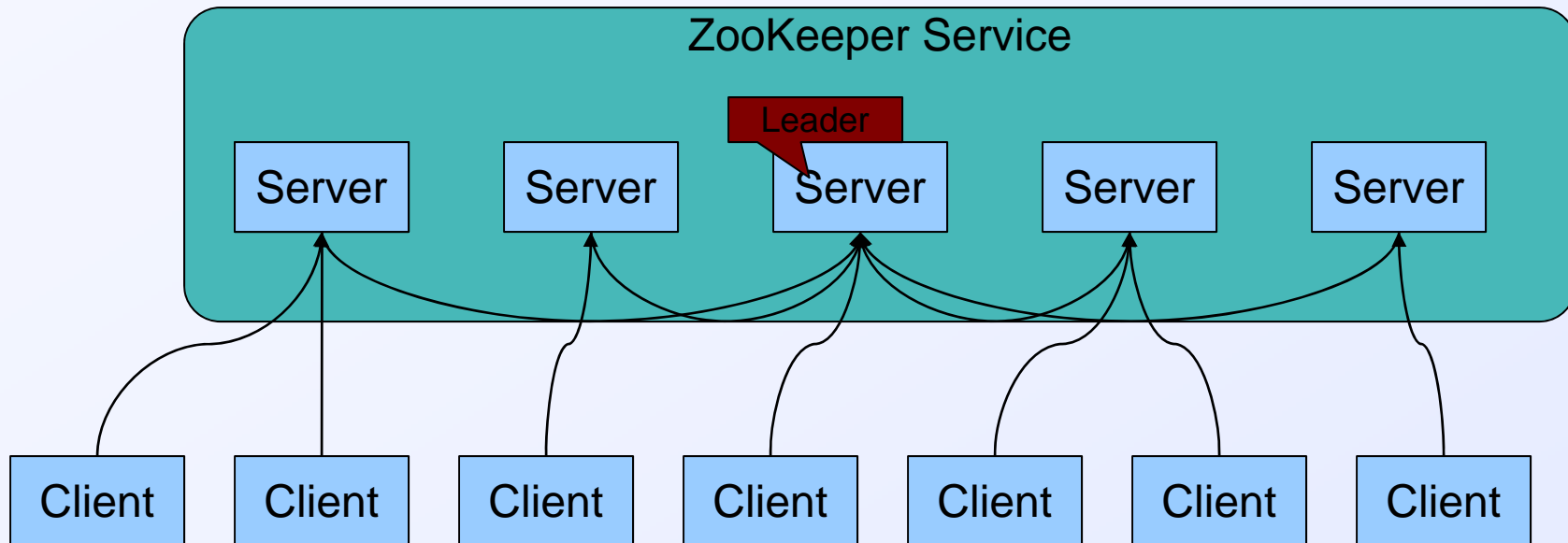


Zookeeper Data Model

- Hierarchical namespace
- Each znode has data and children
- data is read and written in its entirety
- Nodes store < 1MB data
- Writes go to all nodes



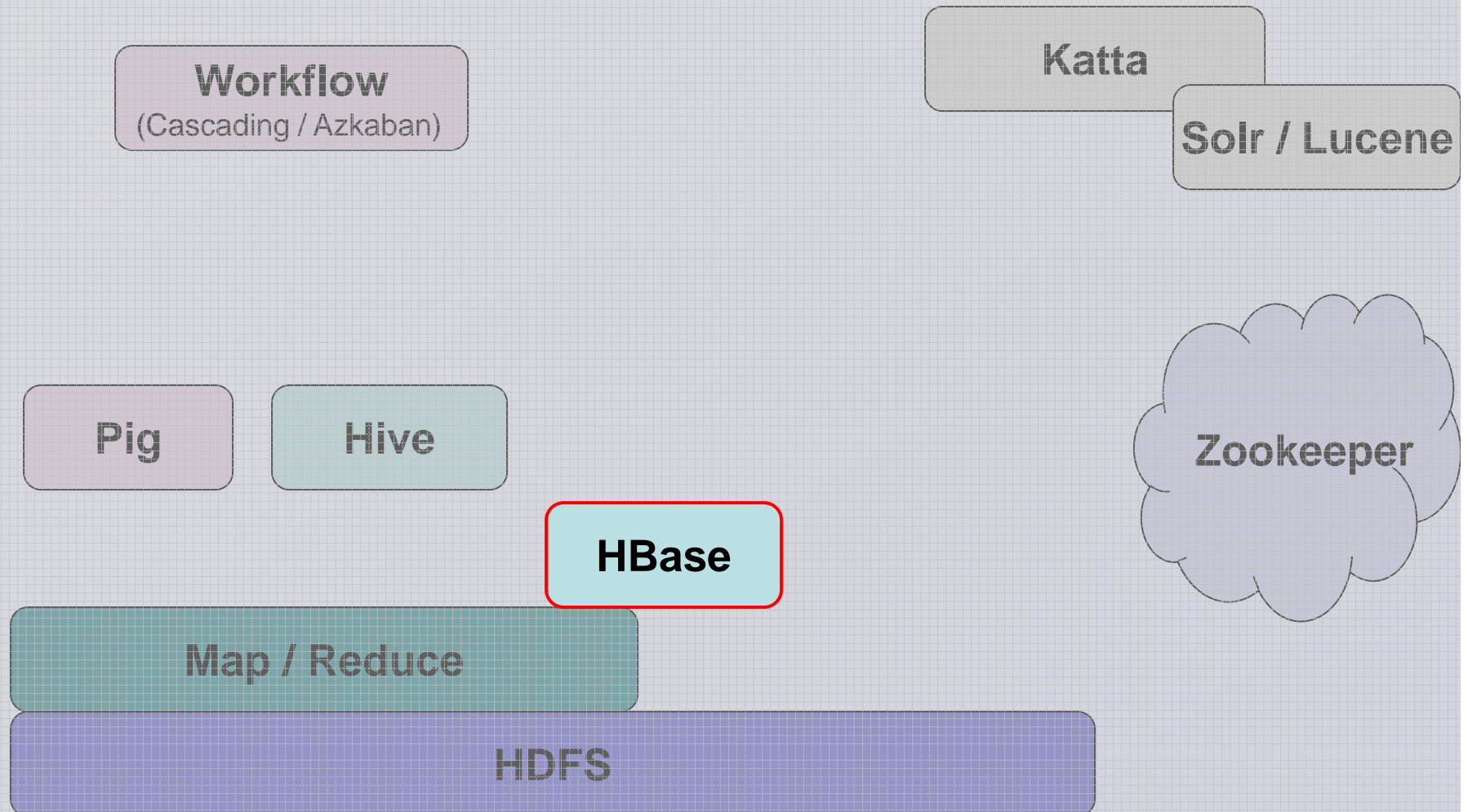
ZooKeeper Service



- All servers store a copy of the data (in memory)
- A leader is elected at startup
- Followers service clients, all updates go through leader
- Update responses are sent when a majority of servers have persisted the change



The M/R stack of open source software – HBase



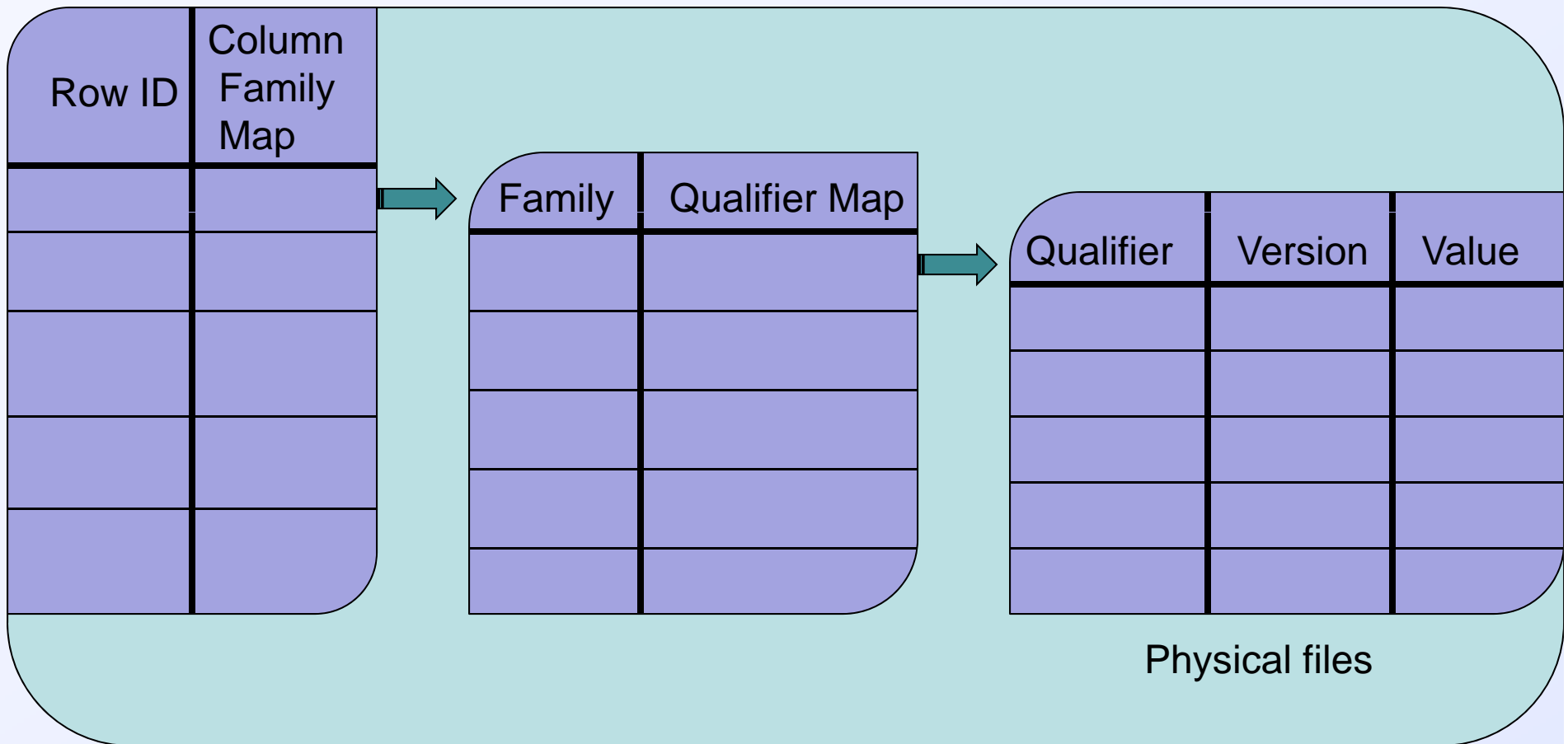
HBase

- Distributed column oriented data store
 - Only supports one data type
 - Tables are broken into regions
 - Regions are automatically split and redistributed
 - All data is local
- Scales to > 1M row / second insert rate (20 node cluster)
- Tightly integrated with Hadoop -> rows can be input/output for map/reduce tasks



HBase Data model

Table



HBase Data model (simplified)

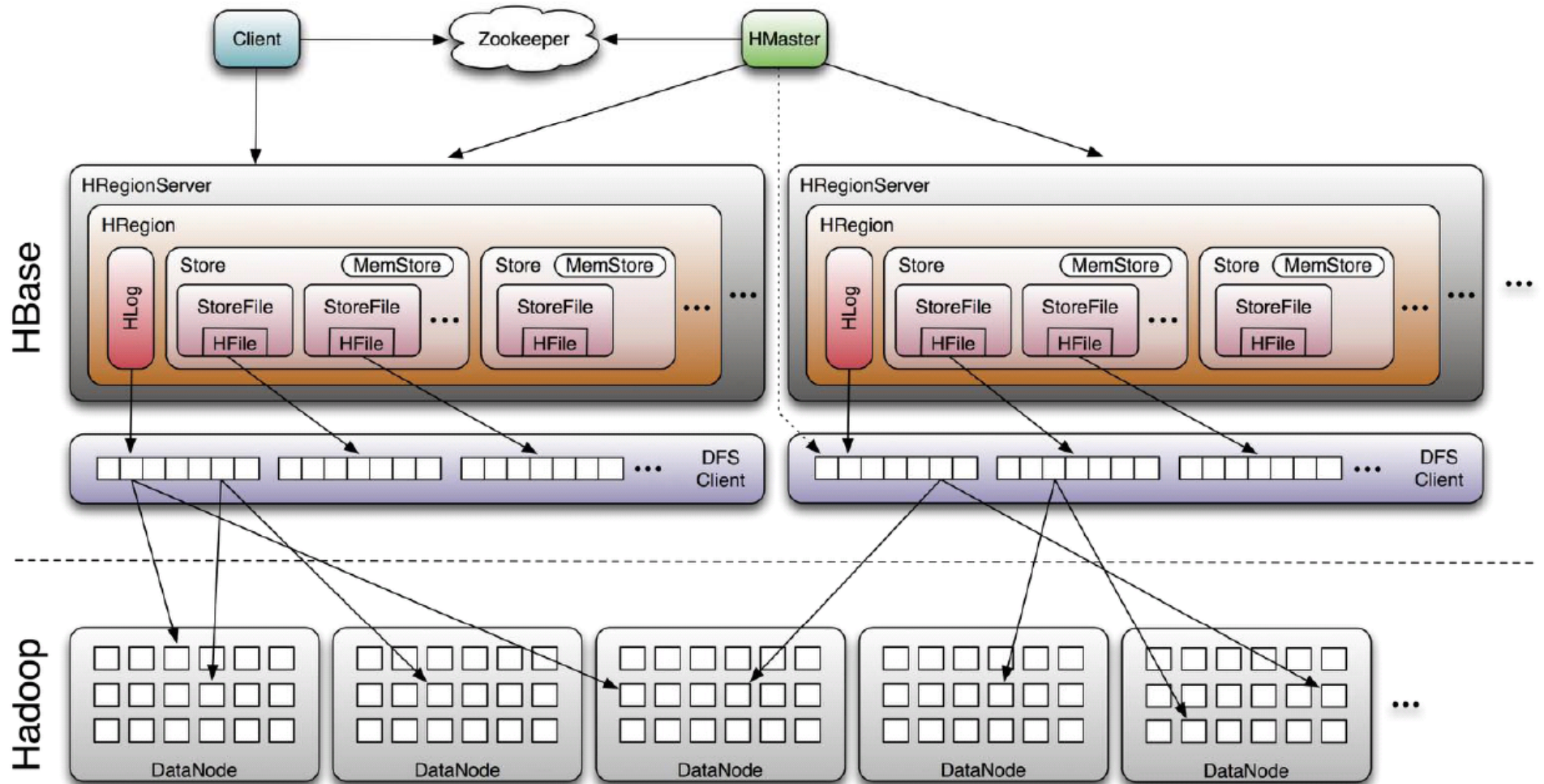
Table

Row ID	Family: Column: Version
Key	Value
⋮	⋮

Regions partitioned on row key



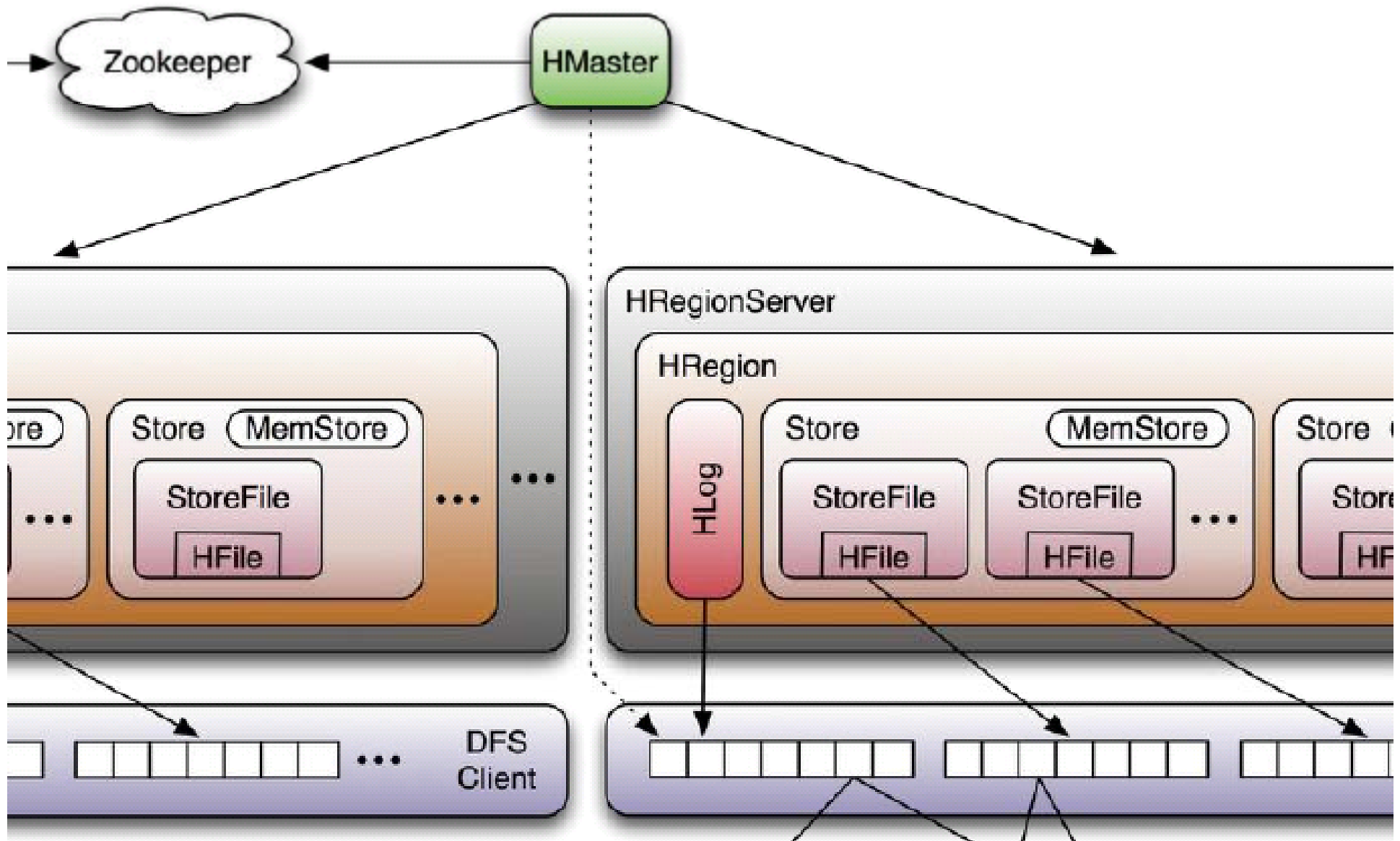
HBase System Architecture



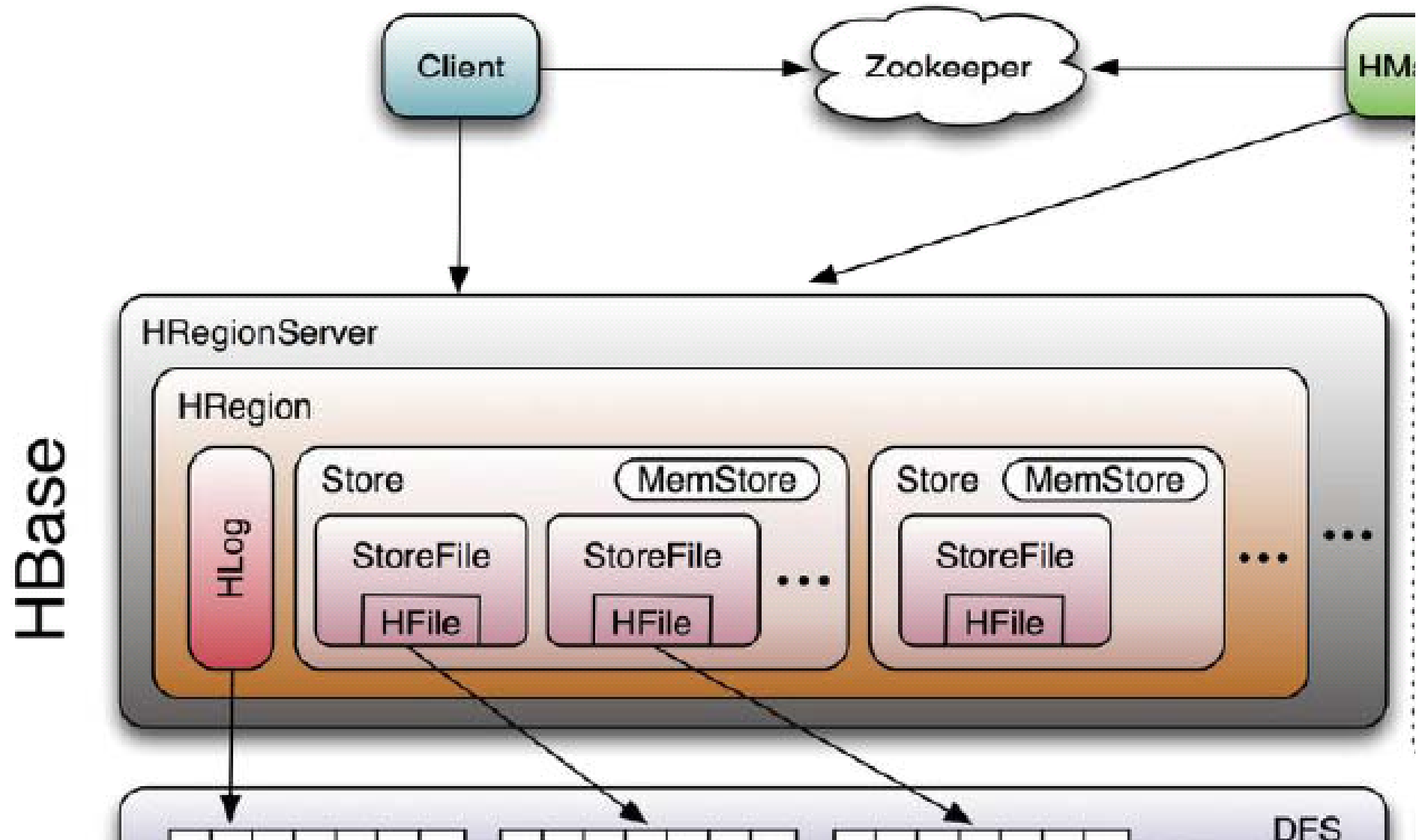
From <http://www.larsgeorge.com/2009/10/hbase-architecture-101-storage.html>



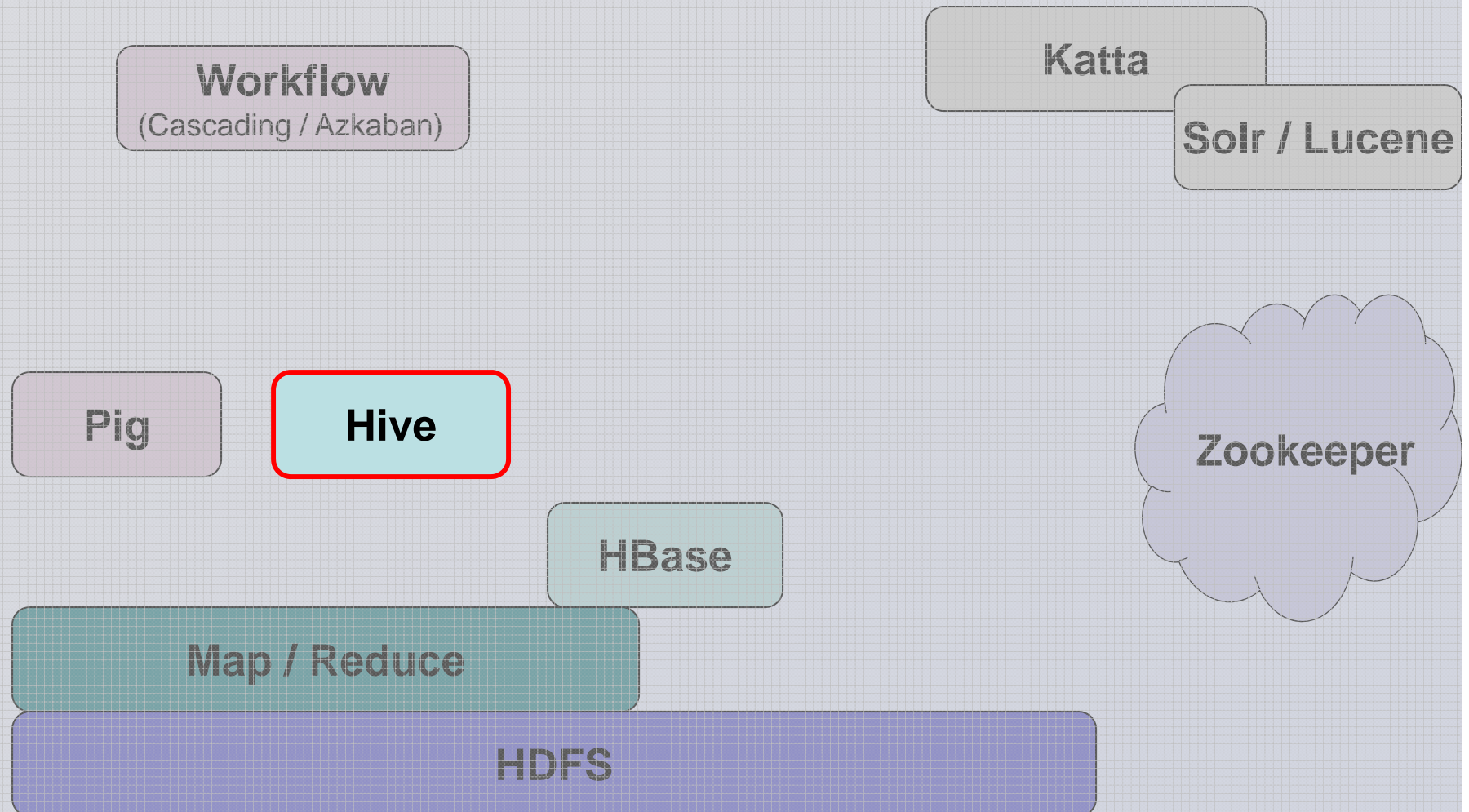
HBase Master manages region servers



Hbase Client directly access region servers for data



The M/R stack of open source software – Hive



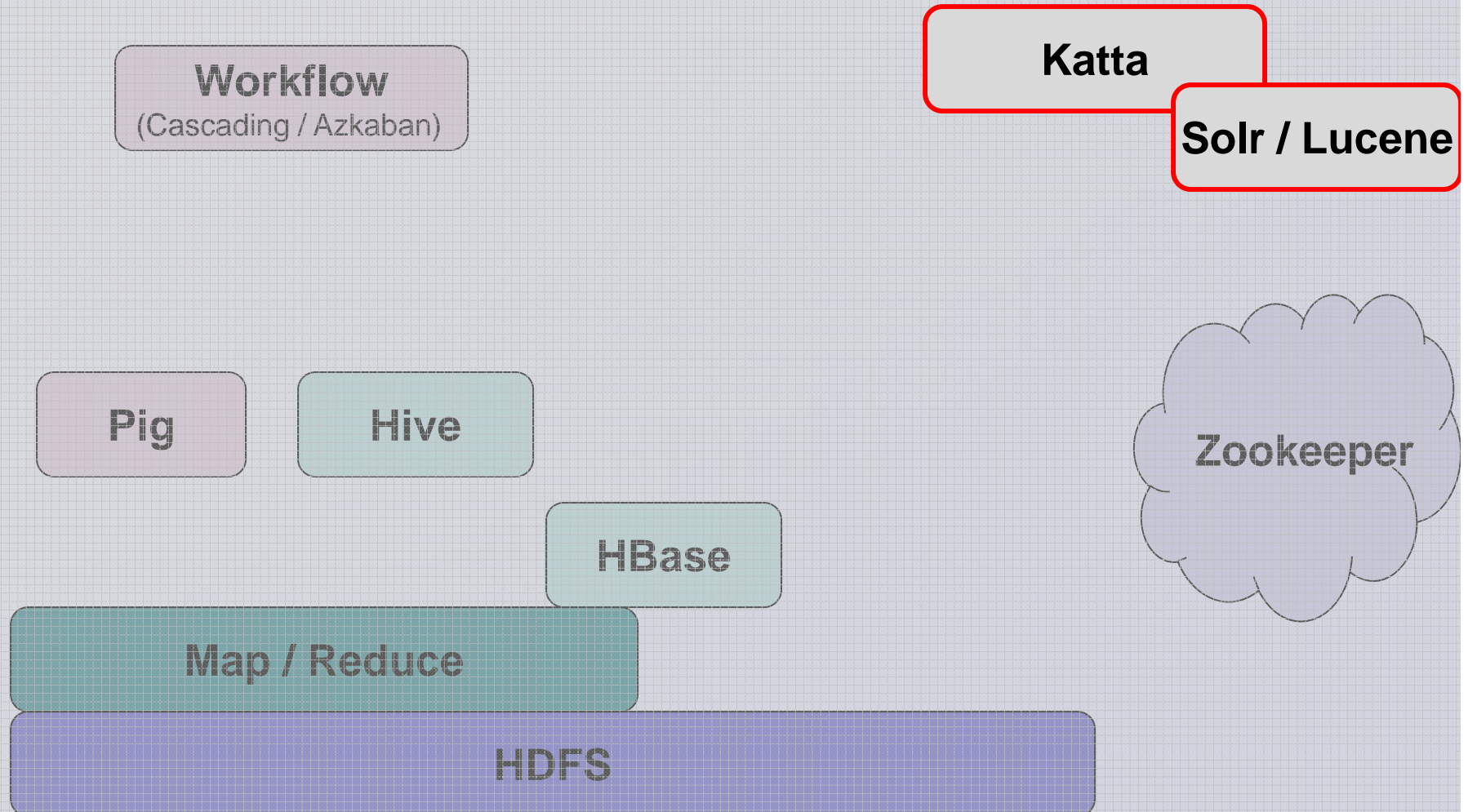
Hive provides and SQL-like interface to data

- Components
 - Shell: SQL-like command line; Web; JDBC
 - Driver: API interface
 - Compiler: parse, plan, optimize
 - Execution Engine: DAG of stages (M/R, HDFS, or metadata)
 - Metastore: schema, location in HDFS, SerDe

http://www.cloudera.com/videos/introduction_to_hive



The M/R stack of open source software – Solr / Katta



Solr

- Faceted text search interface built on top of Lucene
- Built as a native web app – drops into any web server



Faceted search is a foundational component for ad hoc document analysis

Category
< Any Category
Electronics
Audio & Video (82)
Computers & Add-Ons (5)
Accessories & Supplies (2)

Brand
< Any Brand
Samsung

Seller
Any Seller
ButterflyPhoto (23)
Amazon.com (22)
PRESTIGE CAMERA (16)
GENGLOBAL (16)
Crutchfield (15)
B&H Photo-Video (15)
Vanns (14)
> See more...

Eco-friendly
All Products
Energy Star (30)

Price
< Any Price
\$1000 to \$1999 GO


Shipping Option
Any Shipping Option
Can be shipped within one business day from Amazon.com (19)


"televisions" > Electronics > Samsung > \$1000-\$1999


Related Searches: [lcd televisions](#), [tv](#).


Showing 1 - 24 of 92 Results

< Previous | Page: 1 2 3

1.

Samsung LNT5265F 52-inch 1080p LCD HDTV
Buy new: ~~\$2,999.99~~ [Click to see price](#)
[11 Used & new](#) from \$1,574.06
In Stock
★★★★★ (183)

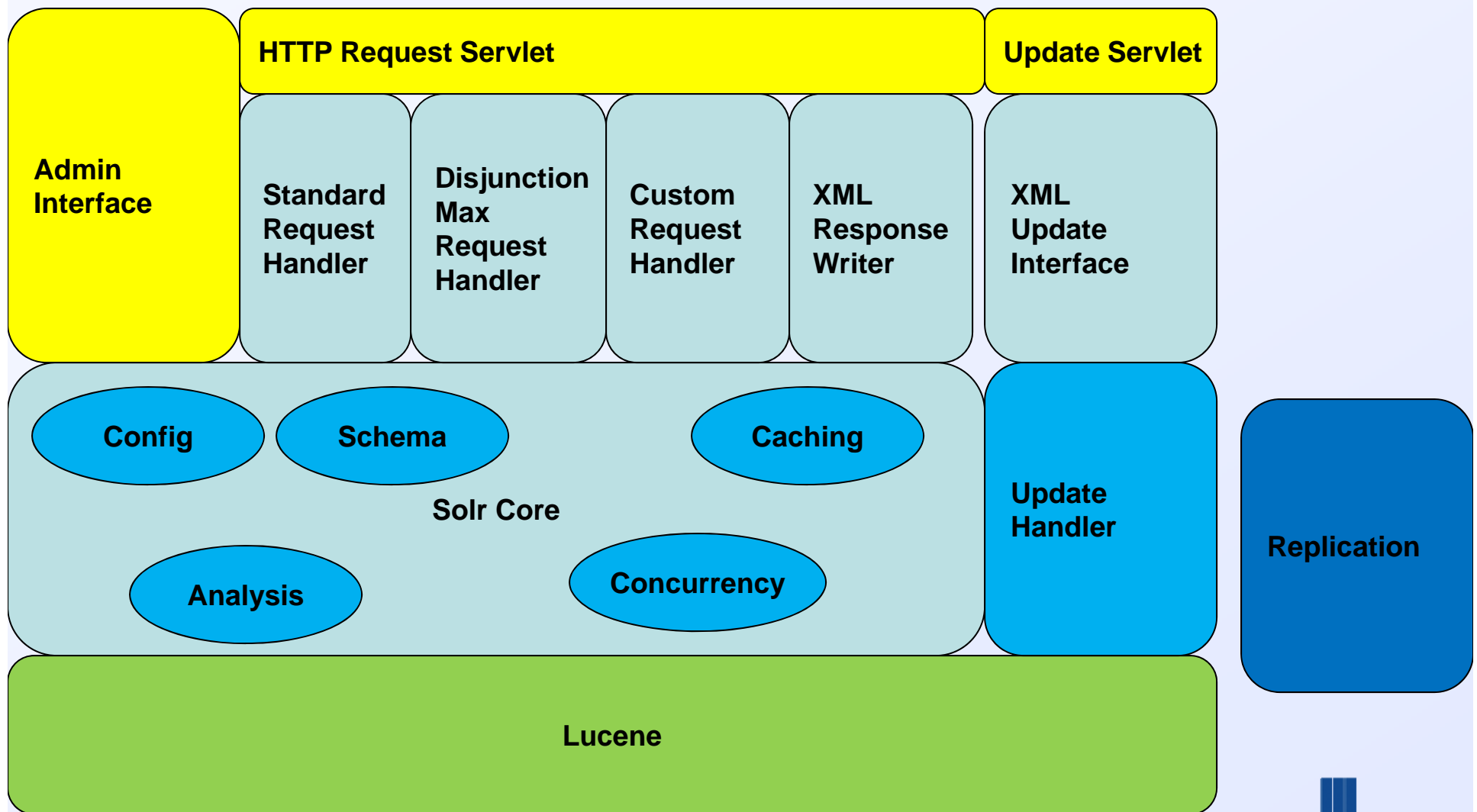
2.

Samsung HLT5687SAX 56-inch LCD HDTV
Buy new: ~~\$2,299.99~~ [Click to see price](#)
[5 Used & new](#) from \$1,190.05
In Stock
★★★★★ (151)

4.

Samsung LNT4671F 46-inch 1080p 120Hz LCD HDTV
Buy new: ~~\$3,199.99~~ **\$2,599.99**
[6 Used & new](#) from \$1,609.95
In Stock
★★★★★ (19)

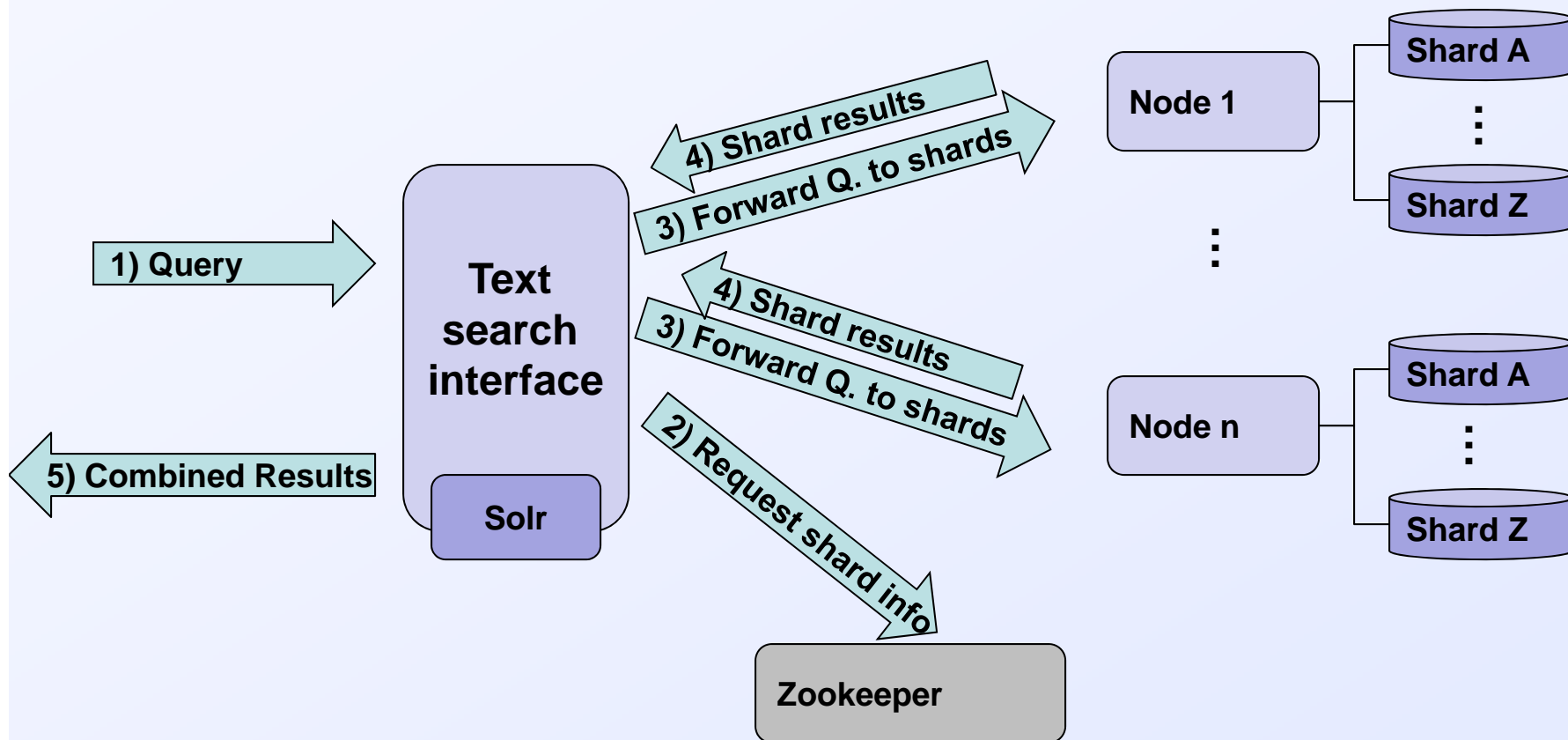
5.

Samsung HLT6176S 61-Inch LCD HDTV
Buy new: ~~\$2,299.99~~ [Click to see price](#)
[4 Used & new](#) from \$1,488.00
In Stock
★★★★★ (151)

5

Solr architecture



Katta provides vertical and horizontal scalability



Projects using Hadoop at LLNL

- Student projects
 - Bioinformatics [James Leek]
 - Continuous time LDA [Kurt Miller and Tina Elliasi-Rad]
- Advanced R&D projects
 - Network analytics
 - Keyword tagging and entity extraction
 - Faceted Search
- Research projects
 - READ LDRD
- Program deployments
 - BKMS



Bioinformatics (student project)

- KPATH: produce DNA signatures for detection of pathogens
 - k-mer lexing: produce set of unique DNA sequences of length 15-60
 - Sliding window
 - Discover k-mers that are unique between bacteria and viruses



K-mer parsing performance comparisons

- Lexing bacteria file 30 k-mer length [120 GB]
 - Optimized suffix tree [C implementation]
 - on single node, 256 GB RAM, 16 processor system
 - 10.5 hours
 - Custom hadoop implementation
 - 85 nodes, 8 GB RAM, dual processor [old]
 - ~1 hour



Unique K-mer grouping performance

- Group unique k-mers of length 15 [13 GB data]
 - Pig implementation using outer joins [10 LOC]
 - More than 9 hours
 - Customer hadoop implementation:
 - 3 hours 26 minutes

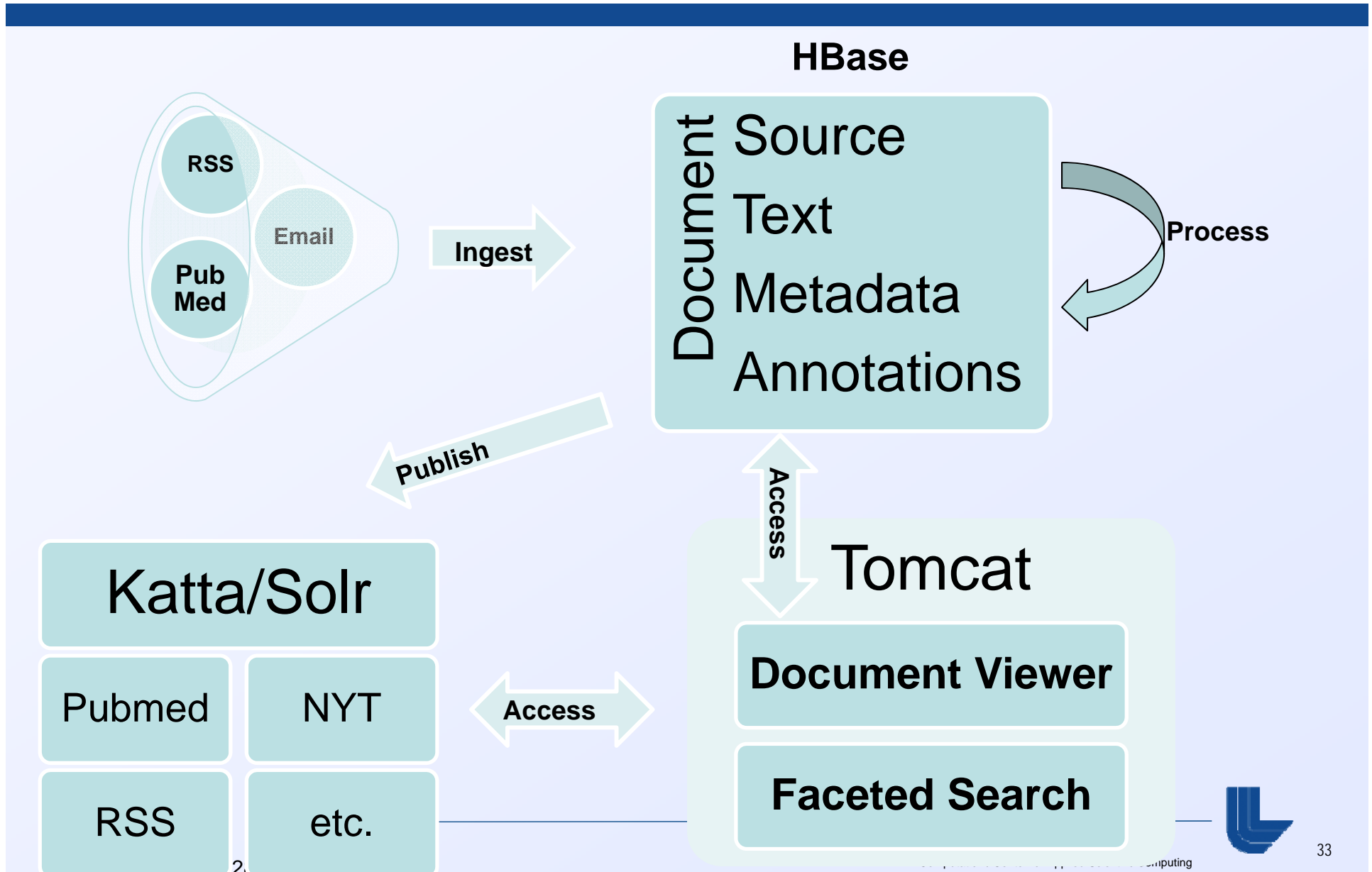


Network data

- HDFS provides storage layer for large repositories of network data
- Hive provides an SQL interface
- Performance on single query for 6 months of data:
 - Tuned Oracle DB: hours to days
 - Hive: minutes



Hadoop-based document management architecture



Example Data Flow

Initial Load

- Load original documents into document table

M/R

Parser

- Custom map code to extract text and meta data

M/R

NLP

- Named Entity Extraction (SNER)
- Parsing / Coreference

M/R

Topics

- Send corpus slices to LDA for topic modeling

Index

- Write specific HBase columns to faceted Solr index shards

M/R

Serve

- Manage indexes with Katta over HDFS



Keyword tagging & Entity extraction

- Keyword tagging
 - Large dictionaries (100K terms)
 - Finite state machine to store dictionary
- Named Entity Recognition
 - Stanford NER
 - CRF model [People, Organizations, Location]



Performance of Keyword tagging & Entity extraction

- 21M Pubmed entries + 1M news articles
 - 11M Pubmed abstracts
- 55K dictionary key phrases
- 6 node cluster [16 core, 96 GB RAM, 6 TB disk]
- Keyword Tagging
 - 8 minutes, 34 seconds
- Named Entity Annotation
 - 1 hr 58 minutes





Query NYTSolr:

☒ Extracted ☐ Source

☒ person ☐ organization ☐ location

☒ All ☐ Hierarchical

- | | | | |
|--|---|---|----------------------------------|
| <input type="checkbox"/> Abbas, Mahmoud | <input type="checkbox"/> Abbas, Mahmoud (Prime Min) | <input type="checkbox"/> Abt, Samuel | <input type="checkbox"/> Akhta |
| <input type="checkbox"/> Altman, Daniel | <input type="checkbox"/> Anderson, Dave | <input type="checkbox"/> Anderson, Michael P (Lt Col) | <input type="checkbox"/> Anna |
| <input type="checkbox"/> Araton, Harvey | <input type="checkbox"/> Archibold, Randal C | <input type="checkbox"/> Armstrong, Lance | <input type="checkbox"/> Aske |
| <input type="checkbox"/> Baker, Al | <input type="checkbox"/> Banerjee, Neela | <input type="checkbox"/> Barboza, David | <input type="checkbox"/> Barri |
| <input type="checkbox"/> Battista, Judy | <input type="checkbox"/> Beers, Christopher | <input type="checkbox"/> Bell, Jack | <input type="checkbox"/> Belsc |
| <input type="checkbox"/> Bennet, James | <input type="checkbox"/> Berlusconi, Silvio (Prime Min) | <input type="checkbox"/> Bernstein, Viv | <input type="checkbox"/> Blair, |
| <input type="checkbox"/> Bloomberg, Michael R | <input type="checkbox"/> Bloomberg, Michael R (Mayor) | <input type="checkbox"/> Bradsher, Keith | <input type="checkbox"/> Brant |
| <input type="checkbox"/> Bremer, L Paul III | <input type="checkbox"/> Brick, Michael | <input type="checkbox"/> Broder, John M | <input type="checkbox"/> Brow |
| <input type="checkbox"/> Bryant, Kobe | <input type="checkbox"/> Bush, George W | <input type="checkbox"/> Bush, George W (Pres) | <input type="checkbox"/> Cald |
| <input type="checkbox"/> Carr, David | <input type="checkbox"/> Chawla, Kalpana (Dr) | <input type="checkbox"/> Cheney, Dick (Vice Pres) | <input type="checkbox"/> Clare |
| <input type="checkbox"/> Clark, Laurel Salton (Dr) | <input type="checkbox"/> Clemens, Roger | <input type="checkbox"/> Clinton, Bill | <input type="checkbox"/> Coop |
| <input type="checkbox"/> Curry, Jack | <input type="checkbox"/> Curtis, Ben | <input type="checkbox"/> Davis, Gray | <input type="checkbox"/> Davis |
| <input type="checkbox"/> Davis, James E (Councilman) | <input type="checkbox"/> Dennehy, Patrick | <input type="checkbox"/> Dewan, Shaila K | <input type="checkbox"/> Dots |
| <input type="checkbox"/> Dunlap, David W | <input type="checkbox"/> Dunning, Jennifer | <input type="checkbox"/> Elliott, Andrea | <input type="checkbox"/> Elliott |
| <input type="checkbox"/> Feder, Barnaby J | <input type="checkbox"/> Feuer, Alan | <input type="checkbox"/> Filkins, Dexter | <input type="checkbox"/> Finley |
| <input type="checkbox"/> Fuerbringer, Jonathan | <input type="checkbox"/> Glater, Jonathan D | <input type="checkbox"/> Greenhouse, Steven | <input type="checkbox"/> Grim |
| <input type="checkbox"/> Hakim, Danny | <input type="checkbox"/> Hanley, Robert | <input type="checkbox"/> Hart, Ariel | <input type="checkbox"/> Healy |
| <input type="checkbox"/> Hermoso, Rafael | <input type="checkbox"/> Herszenhorn, David M | <input type="checkbox"/> Hicks, Jonathan P | <input type="checkbox"/> Hode |



Query NYTSolr:

☐ Extracted ☒ Source

☐ author ☒ topics ☐ taxonomy ☐ online_section ☐ print_section ☐ series_name

☒ All ☐ Hierarchical

☐ Accidents and Safety

☐ Airlines and Airplanes

☐ Archaeology and Anthropology

☐ Atomic Weapons

☐ Basketball

☐ Bombs and Explosives

☐ Cellular Telephones

☐ Classical Music

☐ Computer Software

☐ Courts

☐ Decisions and Verdicts

☐ Drug Abuse and Traffic

☐ Election Issues

☐ Ethics

☐ Fish and Other Marine Life

☐ Accounting and Accountants

☐ Alcoholic Beverages

☐ Architecture

☐ Automobiles

☐ Bicycles and Bicycling

☐ Books and Literature

☐ Child Abuse and Neglect

☐ Coaches and Managers

☐ Computers and the Internet

☐ Credit

☐ Demonstrations and Riots

☐ Drugs (Pharmaceuticals)

☐ Elections

☐ Executives and Management

☐ Food

☐ Acquired Immune Deficiency Syndrome

☐ All Star Games

☐ Area Planning and Renewal

☐ Bankruptcies

☐ Biographical Information

☐ British Open (Golf)

☐ Children and Youth

☐ College Athletics

☐ Consumer Protection

☐ Crime and Criminals

☐ Diet and Nutrition

☐ Economic Conditions and Trends

☐ Electric Light and Power

☐ Finances

☐ Football

☐ Adverse Forces

☐ Agec

☐ Appa

☐ Arms Forces

☐ Bank

☐ Biolo

☐ Budg

☐ City C

☐ Colle

☐ Cook

☐ Danc

☐ Discr

☐ Edito

☐ Elect






☐ Fines

☐ Fines

Extracted

 person	29
 organization	31
 location	29


Source

 author	37
 topics	202
 taxonomy	254
 online_section	52
 print_section	40

 desk	40
--	----

Name	Not	Count
Business/Financial Desk	(X)	31
Editorial Desk	(X)	4
Metropolitan Desk	(X)	3
Foreign Desk	(X)	2



 descriptor	273
--	-----

Results

iScore: 0.5 [Why?](#)

Id: [e0e92e1e7535ba6c2e8f8042aedef983b4a](#)

Solr: NYTSolr

Solr Score: 6.085178

ABSTRACT:

A Philippine antigraft court has dealt a blow to efforts conglomerate San Miguel to regain control of the share

Coming after nearly 17 years of legal battles, the court San Miguel chairman, Eduardo M. Cojuangco Jr., had a controlling stake in one of the country's largest banks the purchase and awarded ownership of the stake in

iScore: 0.5 [Why?](#)

Id: [c8e609ae732796fab759447646a5eb57c](#)

Solr: NYTSolr

Solr Score: 6.085178

Faceted Search Indexing Performance

- Creating 1 Solr index on 1M news articles:
 - 8hrs 16 min
 - Map: 37 min
 - Reduce: 8 hrs 14 min
- Creating 50 Solr indexes on 1M news articles:
 - 55 min
 - Map: 7 min
 - Reduce: 54 min



Open-sourced products and others in the open source pipeline

- iScore
 - Content-based personalization
 - [pre-Hadoop]
- Reconcile
 - Coreference resolution software built on open source tools [with Cornell and U. Utah]
 - Additional adaptation to Hadoop
- Dunk
 - An elegant java annotation system that allows you to have the fields of a java object serialized (deserialized) to (from) an HBase table
 - Simplifies queries, object construction, and map/reduce formulation



Questions?

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

