



Cloud Computing @Yahoo!

Raghu Ramakrishnan

**Yahoo! Fellow
Chief Scientist, Search and Cloud Platforms**

(Many slides courtesy of others at Yahoo!) 1

Cloud Services @Y!: Use Cases

My Yahoo! | May 12, 2009

Sign In | New here? Sign Up



Web | Images | Video | Local | Shopping | More

Web Search

MY FAVORITES + Add

View Yahoo! Sites

[Yahoo! Mail](#)

Autos

eBay

Finance (Down)

Flickr

Games

Messenger

Movies

Music

MySpace

omg!

Personals

Sports

TV

Weather

RECOMMENDED

Deal Of The Day

Buzz

Shine

Edit

+ Add

Yahoo! Mail

Check your Yahoo! Mail

Back to Yahoo! Home

YAHOO! MAIL

Machine Learning
(e.g. Spam filters)

Stay connected

Preview and keep up-to-date with your mail.

Attachment Storage

Yahoo! Mail

Back to Yahoo! Home

MAIL

Attachment Storage

McAfee

Chase C...

Dave Ramsey

Lottery Results - L.A. Times

White House officials say no decision has been... - L.A. Times

updated 10:32 am PDT

More: News Popular Buzz

Dow: 8,385.40 -0.39% Nasdaq: 1,700.55 -1.77%

Sponsored by: Scottrade Enter stock symbol Get Quotes

TOP SEARCHES

1. Cameron Diaz

2. Leonard Nimoy

3. Jason Bourne

4. Sarah Fawcett

5. Star Trek: The Next Generation

6. U.S. Stamps

7. 10 Reasons

Search Index

TOYOTA moving forward

PERFECT TIMING EVENT

Ads Optimization

ROUGH JUNE 1ST. See local offers

Toyota Sales Event - Ad Feedback

« Prev | Next »

Image/Video Storage & Delivery

YAHOO! NEWS

A baby black jaguar is named by its mom

During star races across space

More news photos »

700

What's Your Credit Score? A bad credit score is 600 or below. See yours in 2 easy steps for \$0. By Experian®



Yahoo! Data Scale

Massive user base and engagement

- 640M+ unique users, 11B page visits/month
- Hundreds of petabytes of storage
- Hundreds of billions of objects
- Hundreds of thousands of requests/sec, 200B events/day, 200 PB/day

Global

- Tens of globally distributed data centers
- Serving each region at low latencies

Challenging Users

- Rapidly extracting value from voluminous data
- Downtime is not an option (outages cost \$millions)
- Variable usage patterns

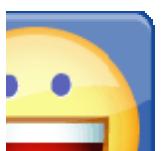




“Just look at our homepage, for example. Since we began pairing our **content optimization technology** with editorial expertise, we’ve seen **click-through rates in the Today module more than double**. And we’re making additional improvements to this technology that will make the user experience ever more personally relevant.”

Carol Bartz, Analyst Call, January 27, 2010

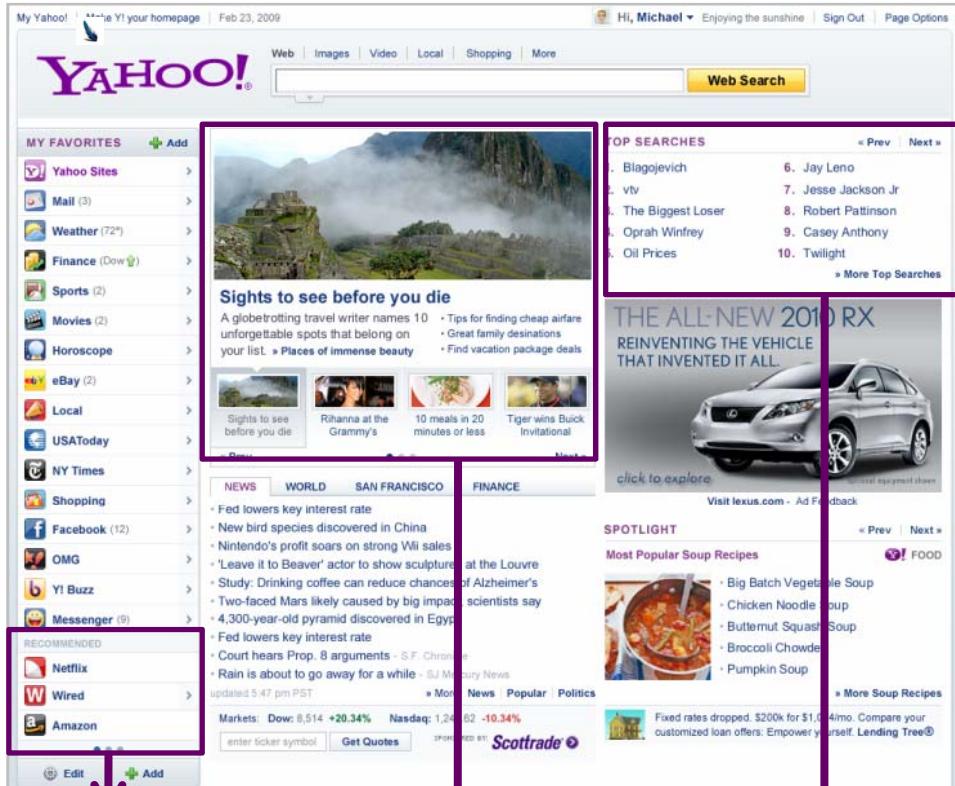
CONTENT OPTIMIZATION FOR PORTALS



Yahoo! Front Page

Product Objective

Prioritize small pool of editorially programmed packages to optimize engagement in real-time



Recommended links

News Interests

Top Searches

+79% clicks

vs. randomly selected

+160% clicks

vs. one size fits all

+43% clicks

vs. editor selected

Key Features

Package Ranker (CORE)

Ranks packages by expected CTR based on data collected every 5 minutes

Dashboard (CORE)

Provides real-time insights into performance by package, segment, and property

Mix Management (Property)

Ensures editorial voice is maintained and user gets a variety of content

Package rotation (Property)

Tracks which stories a user has seen and rotates them after user has seen them for a certain period of time

Key Performance Indicators

Lifts in quantitative metrics

Editorial Voice Preserved



Content Optimization & Cloud

Offline Modeling

- Exploratory data analysis
- Regression, feature selection, collaborative filtering (factorization)
- Seed online models & explore/exploit methods at good initial points
- Reduce the set of candidate items



Hadoop
Large amount of historical data (user event streams)



Online Learning

- Online regression models, time-series models
- Model the temporal dynamics
- Provide fast learning for per-item models

Near real-time user feedback

Explore/Exploit

- Multi-armed bandits
- Find the best way of collecting real-time user feedback (for new items)



Data Management in CORE



1) User click history logs
stored in HDFS

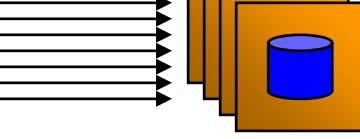


2) Hadoop job builds
models of user
preferences



HDFS

3) Hadoop reduce
writes models to
Sherpa user table



4) Models read from
Sherpa influence users'
frontpage content



Candidate
content





Matching Users to Content

We learn how user attributes correlates with engagement in each item

	Default	Male	Female	18-24	25-34	Heavy Sports
	7.1	-0.4	+0.4	+0.3	+0.1	-0.5
	6.8	+1.0	-1.0	+0.2	+0.3	+2.1
	6.5	-0.6	+0.6	+0.5	+0.3	-0.8
	6.2	0	0	-0.7	-0.5	-0.3
	5.9	-1.1	+1.1	-0.5	-0.2	-0.2

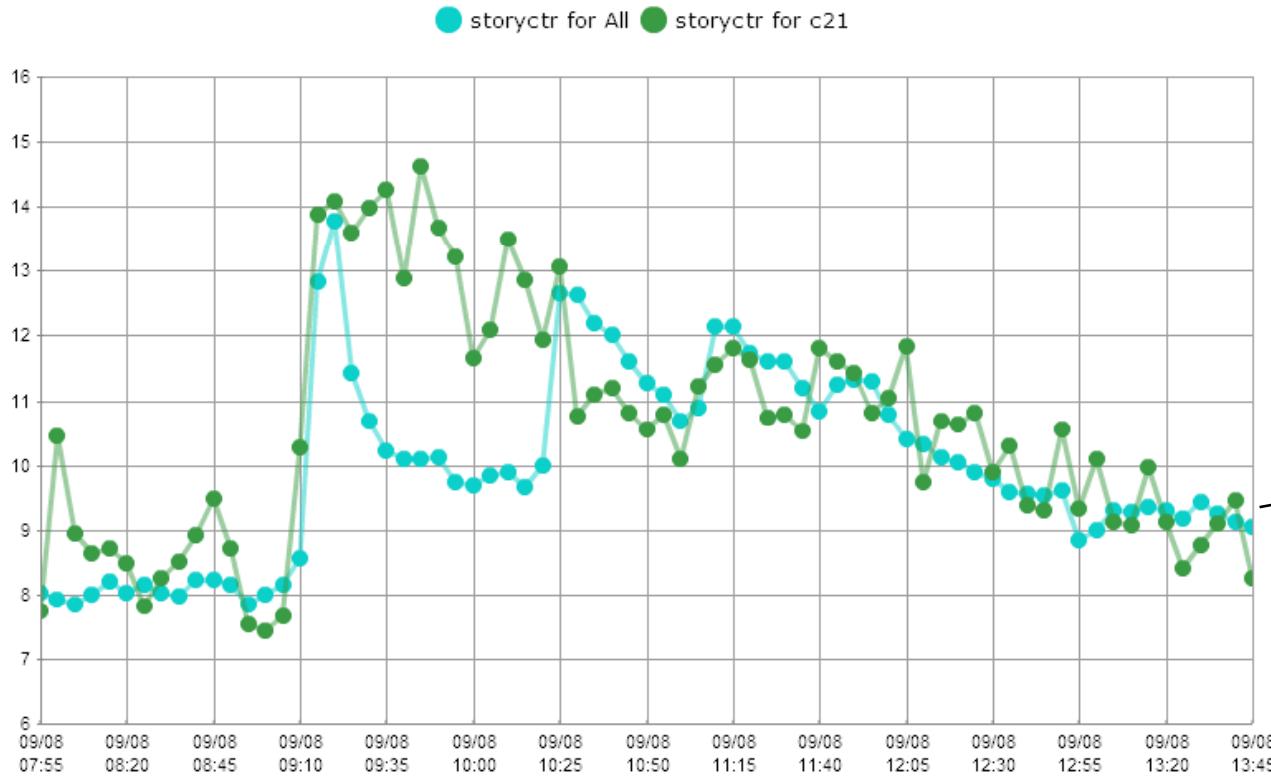


We compute rankings for each user based on his/her attributes



CORE Dashboard: Overall CTR

Compare performance of models and historical benchmarks



Compare
buckets and
models over
time



See which
content was
promoted most
across time



Compare
bucket
metrics

Bucket	Page Views	All Clicks	Story Clicks	Footer Clicks	Overall CTR	Story CTR	Footer CTR	Lift compared to
All	43,005,783	13,061,688	4,286,407	8,775,281	30.37	9.97	20.4	0
c21	281,821	84,853	29,777	55,076	30.11	10.57	19.54	6.02

COKE Dashboard: Segment Heat Map



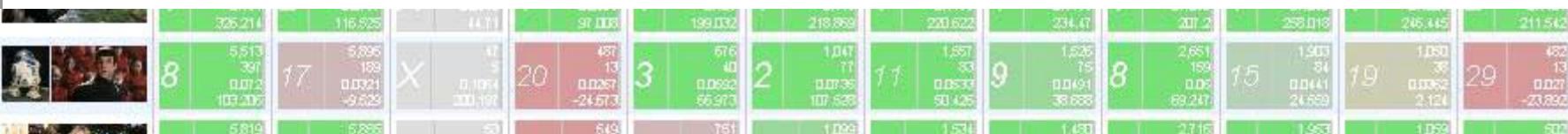
Package	male	female	OMG	BUAuto	BUEnt	BU Fin	Health	BUSport	NBA	BUTrav	ALL
	408,260 18,440 0.0452 8,417	390,404 14,449 0.037 -11,113	270,039 16,940 0.0621 50,661	121,080 1,389 0.061 46,564	270,038 16,940 0.0621 50,661	325,873 20,012 0.0614 47,488	195,796 12,763 0.0652 56,553	350,152 21,454 0.0613 47,152	132,916 9,457 0.0712 70,879	123,388 9,457 0.064 53,691	523,611 38,457 0.0416 0
1	8,067 852 0.1056 153,654	7,657 674 0.088 111,405	1 1 0.1405 231,406	5,125 2,382 0.1201 188,362	1 1 0.1405 231,406	5,125 2,382 0.1201 221,221	6,415 858 0.1337 299	3,769 532 0.1412 299	6,750 917 0.1359 226,272	2,585 385 0.1489 218,294	2,490 330 0.1326 130,143
2	9,968 644 0.0646 55,164	12,847 771 0.0605 45,256	3 2 0.1033 148,043	8,569 3,529 0.0924 121,86	4 2 0.1033 148,043	8,569 3,529 0.0946 121,252	9,744 922 0.0946 154,537	3 6,067 643 0.108 136,702	10,187 1,004 0.0986 164,088	3,820 420 0.1099 157,598	4,031 433 0.1073 48,798
3	3,326 249 0.0749 79,8	3,954 212 0.0536 28,169	5 2 0.0916 120,066	2,521 1,004 0.1016 143,995	2 5 0.0916 120,066	2,521 1,004 0.0915 119,782	3,016 276 0.0915 140,167	5 1,860 186 0.1 140,167	3,291 310 0.0942 126,229	1,141 136 0.1192 186,264	1,039 100 0.0962 131,152
4	2,952 133 0.0519 24,617	2,004 81 0.0404 -2,926	11 13 0.0576 134,403	1,250 122 0.0976 94,73	6 3 0.0811 134,403	1,250 122 0.0939 125,53	1,608 151 0.0939 125,53	2 919 103 0.1121 169,175	1,669 154 0.0923 121,604	655 74 0.1121 171,334	591 56 0.0931 123,506
5	2,881 206 0.0715 11,127	3,242 230 0.0709 70,384	2 4 0.0946 127,295	2,071 949 0.1001 140,42	4 3 0.0946 127,295	2,071 949 0.0972 133,368	2,614 254 0.0972 133,368	4 1,605 165 0.1028 146,901	2,740 239 0.0907 109,489	1,036 94 0.0907 117,912	958 78 0.0814 95,543
6	10,785 649 0.0602 44,523	12,768 742 0.0581 39,571	4 7 0.0809 94,261	8,580 694 0.0806 93,584	7 7 0.0809 94,261	8,580 694 0.0817 96,332	9,725 795 0.0817 96,332	6 6,138 590 0.0896 115,204	10,670 866 0.0812 94,905	3,785 321 0.0815 110,122	21,331 1,641 0.06 115,104
7	22,202 1,212 0.0546 31,106	23,328 1,200 0.0514 23,543	6 7 0.0547 98,535	15,593 533 0.0813 95,374	5 6 0.0813 95,374	15,593 533 0.0827 98,535	17,652 1,376 0.078 87,214	8 10,797 915 0.0847 103,532	19,050 8,052 0.0799 91,882	6,639 604 0.0911 118,498	6,435 552 0.0898 106,018
8	26,695 1,160 0.0435 4,401	35,405 1,530 0.0432 3,186	10 8 0.0793 90,371	19,832 552 0.0704 69,011	9 8 0.0793 69,011	19,832 552 0.0755 69,011	21,143 1,541 0.0755 81,26	7 13,721 1,167 0.0851 104,261	22,168 1,743 0.0785 88,836	8,249 788 0.0955 129,424	8,327 689 0.0827 98,721
9	7,745 518 0.0669 60,628	7,202 185 0.0257 -38,308	4 26 13 0.0657 57,889	4,898 322 0.0641 54,007	15 13 0.0657 57,889	4,898 322 0.0657 57,889	6,051 423 0.0699 67,891	19 3,652 295 0.0643 54,544	6,436 506 0.0786 88,882	2,562 308 0.1202 188,726	2,359 169 0.0716 72,057
10	7,699 460 0.0597 43,495	7,201 169 0.0235 -43,635	29 29 0.0707 67,239	4,809 340 0.0696 69,8	11 10 0.0707 67,239	4,809 340 0.0707 69,8	6,004 433 0.0721 73,205	14 3,544 243 0.0686 64,674	6,247 475 0.0716 82,615	2,482 291 0.1035 148,682	2,329 167 0.0711 122,211
11	7,688 393 0.0511 22,777	7,229 336 0.0465 11,528	9 9 0.0759 82,196	4,785 363 0.061 46,418	9 17 0.0759 82,196	4,785 363 0.0668 82,196	6,037 403 0.0668 60,324	13 3,501 245 0.0759 68,069	6,319 430 0.0668 63,431	2,397 182 0.0759 82,395	2,212 152 0.0667 15,809
12	2,138 126 0.0589 41,539	2,479 95 0.0383 -7,963	10 14 0.0751 81,757	1,546 117 0.0764 83,484	8 10 0.0751 81,757	1,546 117 0.0757 81,757	1,941 137 0.0706 69,515	9 1,206 91 0.0755 81,221	2,026 137 0.0676 62,403	723 58 0.0802 92,665	5,347 252 0.0886 106,019
13	28,932 1,160 0.0401	11,695 262 0.0224	15 31 0.0644	11,432 736 0.0652 81,757	24 24 0.0652	11,432 736 0.0644 81,757	14,567 871 0.0602	20 8,011 902 0.0621	16,818 1,085 0.0643	5,444 293 0.0643	45,041 1,523 0.0538



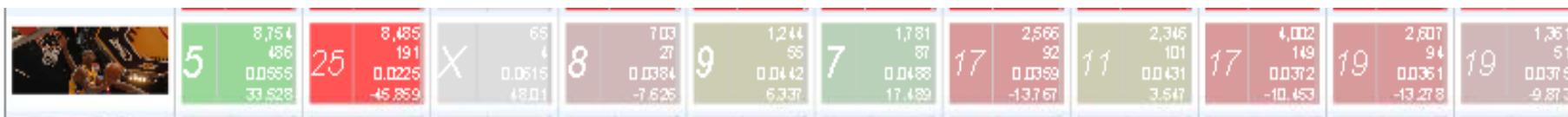
Examples



- **ACQUISITION:** A “Star Trek” package was #3 with 18-20 demo, #2 with 21-24 demo, but #9 overall. We can acquire younger audiences with targeted content like this.



- **ENGAGEMENT:** “Kobe’s astonishing shot” was #25 with women, but #5 with men. We can better engage men (or sports fans) by showing more like this, women by showing less.



- **REACH:** A package about a hair-pulling soccer player was just plain interesting to everyone (#1-3). We can maintain reach by programming content for the mass audience.





NEXT-GEN SEARCH



http://search.yahoo.com/search?p=julia+roberts&ei=UTF-8&fr=moz35

Julia Roberts - Yahoo! Search Res... julia roberts

YAHOO! julia roberts Search

julia roberts twins
julia roberts movies
lyle julia roberts
julia roberts babies
julia roberts henry daniel moder

Explore related concepts:
actor
Epis...
Search for julia roberts actor
Pretty woman
Ber...
julia roberts photos
Julia Roberts News
julia roberts biography
Julia Roberts

Search Pad
SearchScan - On

40,500,000 results for **julia roberts**

Related People

Scarlett Johansson
Emma Roberts
Hilary Swank
Lindsay Lohan
Tom Hanks
Halle Berry

Julia Roberts - Image Results

News & Photos Videos Twitter

[more Julia Roberts photos...](#)

Latest News:

[Hindus concerned about Julia Roberts' "Eat, Pray, Love"](#) - New Kerala - 6 hours ago
[Trailer For 'Eat Pray Love' Starring Julia Roberts](#) - KPBS San Diego - Mar 19 03:18pm
[Link Party: Julia Roberts' New Movie Will Teach You How to Live](#) - El Online - Mar 18 05:48pm
[more Julia Roberts news...](#)

Julia Roberts - Wikipedia

[Early life](#) | [Career](#) | [Influence](#) | [Personal life](#)

Julia Fiona Roberts is an American actress. She is known for starring in the romantic comedy Pretty Woman opposite Richard Gere, which grossed \$463 million worldwide. After receiving...

en.wikipedia.org/wiki/Julia_Roberts - 122k - [Cached](#)

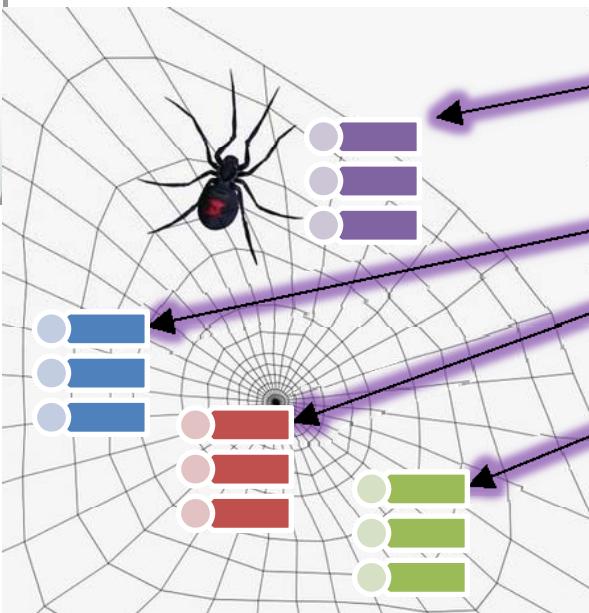
Information aggregation... Concept-centric



Web of Concepts



rich, aggregated data



Aggregated KB

concept

madonna

mumbai
restaurant

san jose

• • •

INDEX



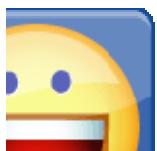
SERP



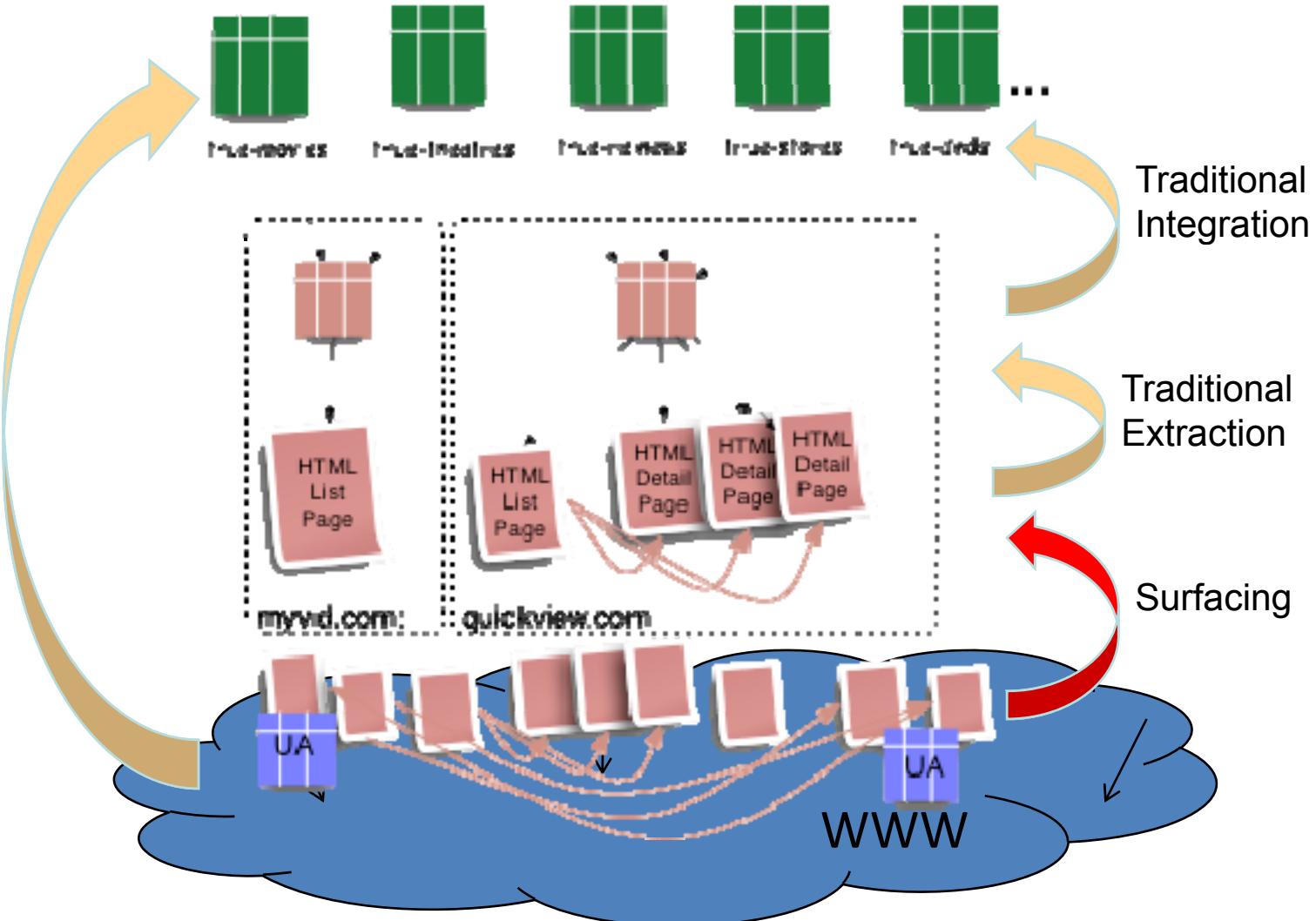
The “index” is keyed by concept instance, and organizes all relevant information (data describing the concept instance and its relationship to other instances), wherever it is drawn from, in semantically meaningful ways



End-to-
End



Web IE: Surfacing, Extraction, Integration





shower head santa clara

Search

Options ▾

Start typing to see suggestions.

Explore related concepts:

Settings

pool
Detached Single Family
Tub
faucets

tile
MLSListings
toilets
Real Estate MLS

Search Pad

SearchScan - On

577,000 results for **shower head santa cl...**

Show All

Google Sites

Yahoo! Local

Metacafe

Video Sites

Shower Head

High Performance Showerheads for Full Body Coverage. Learn More.
Moen.com/Showerheads

Sponsored Results

Sponsored Results

Shower Heads

Find **Shower Heads**. Find Out More at Guide2Faucets.

Guide2Faucets.com/Faucets

Shower Heads for Sale

Your Personal Guide to **shower heads** for sale

www.aashowerheads.info

[See your message here...](#)

Shower Head stores near Santa Clara, CA

Nearby City

All (26)	San Jose (8)	Mountain View (3)	Los Gatos (2)
Santa Clara (2)	Fremont (3)	Campbell (3)	Sunnyvale (1)

1 Conleff Plumbing Supply ★★★★★(5)

[conleff.com](#)
(408) 988-8005 - 2301 Lafayette St, **Santa Clara, CA**
[5 Reviews](#) ▾ | [Overview](#) ▾ | [1 Photo](#) ▾ | [Directions](#) ▾

2 Home Depot ★★★★★(9)

[homedepot.com](#)
(408) 492-9600 - 2435 Lafayette St, **Santa Clara, CA**
[5 Reviews](#) ▾ | [Overview](#) ▾ | [Directions](#) ▾

3 Kitchen & Bath Showplace ★★★★★(1)

[kbshowplace.com](#)
(408) 249-9880 - 1200 Campbell Ave, San Jose, CA
[Overview](#) ▾ | [1 Photo](#) ▾ | [Directions](#) ▾

[23 More Local Results...](#)



na baltimore - Yah...



Web Images Video Local Shopping News More ▾



eggplant parmigiana baltimore

Search

Options ▾

Start typing to see suggestions.

Explore related concepts:

Pizza

mushrooms

veal

tomato sauce

Chicken

Little Italy

Italian Restaurant

Crab Cake

Setting



30,300 results for
eggplant parmigiana ...:



Eggplant Parmigiana Restaurants near Baltimore, MD

Neighborhood

All (36)

Abell (1)

Canton (2)

Central Bal... (4)

Charles North (1)

Chinguapin ... (1)

Downtown (1)

Federal Hill (1)

- 1 Ciao Bella - Baltimore ★★★★★ (11)
[local.yahoo.com](#)

(410) 685-7733 - 236 S High St, **Baltimore**, MD

Menu: eggplant parmigiana

[4 Reviews](#) | [Overview](#) | [2 Photos](#) | [Directions](#)

- 2 Amicci's ★★★★★ (20)
[amiccis.com](#)

(410) 528-1096 - 231 S High St, **Baltimore**, MD

Menu: eggplant parmigiana

[14 Reviews](#) | [Overview](#) | [23 Photos](#) | [Directions](#)

- 3 Pasticcio ★★★★★ (8)
[local.yahoo.com](#)

(410) 522-7700 - 2400 Boston St, **Baltimore**, MD

Menu: eggplant parmigiana

[5 Reviews](#) | [Overview](#) | [3 Photos](#) | [Directions](#)

- 4 Caesar's Den ★★★★★ (7)
[caesarsden.com](#)

(410) 547-0820 - 223 S High St, **Baltimore**, MD

Menu: eggplant parmigiana

[4 Reviews](#) | [Overview](#) | [11 Photos](#) | [Directions](#)



© Yahoo! 2010. Data © NAVTEQ 2009



DATA MANAGEMENT IN THE CLOUD

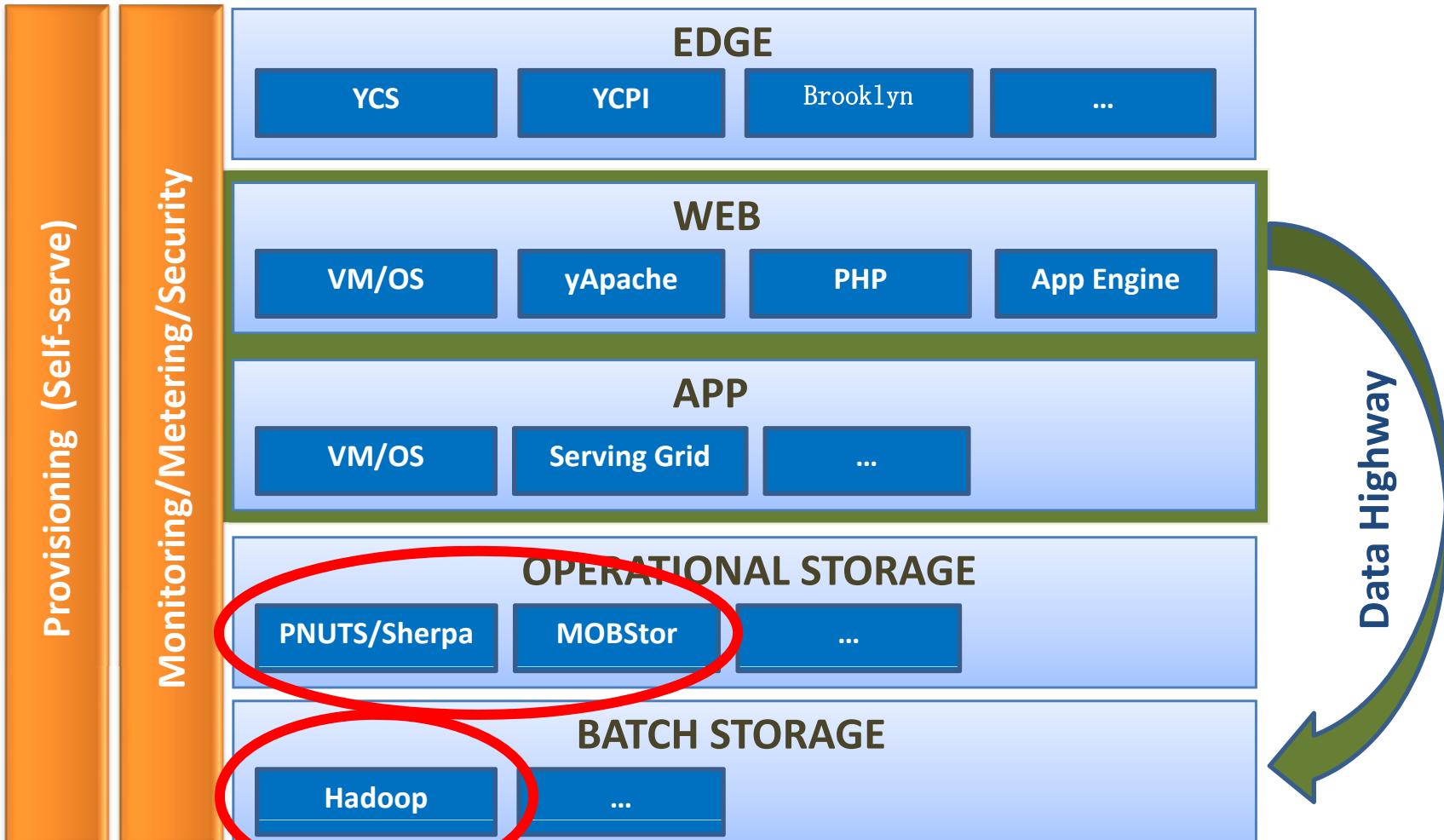


Requirements for Cloud Services

- **Multitenant.** A cloud service must support multiple, organizationally distant customers.
- **Elasticity.** Tenants should be able to negotiate and receive resources/QoS *on-demand* up to a large scale.
- **Resource Sharing.** Ideally, spare cloud resources should be transparently applied when a tenant's negotiated QoS is insufficient, e.g., due to spikes.
- **Horizontal scaling.** The cloud provider should be able to add cloud capacity in increments without affecting tenants of the service.
- **Metering.** A cloud service must support accounting that reasonably ascribes operational and capital expenditures to each of the tenants of the service.
- **Security.** A cloud service should be secure in that tenants are not made vulnerable because of loopholes in the cloud.
- **Availability.** A cloud service should be highly available.
- **Operability.** A cloud service should be easy to operate, with few operators. Operating costs should scale linearly or better with the capacity of the service.

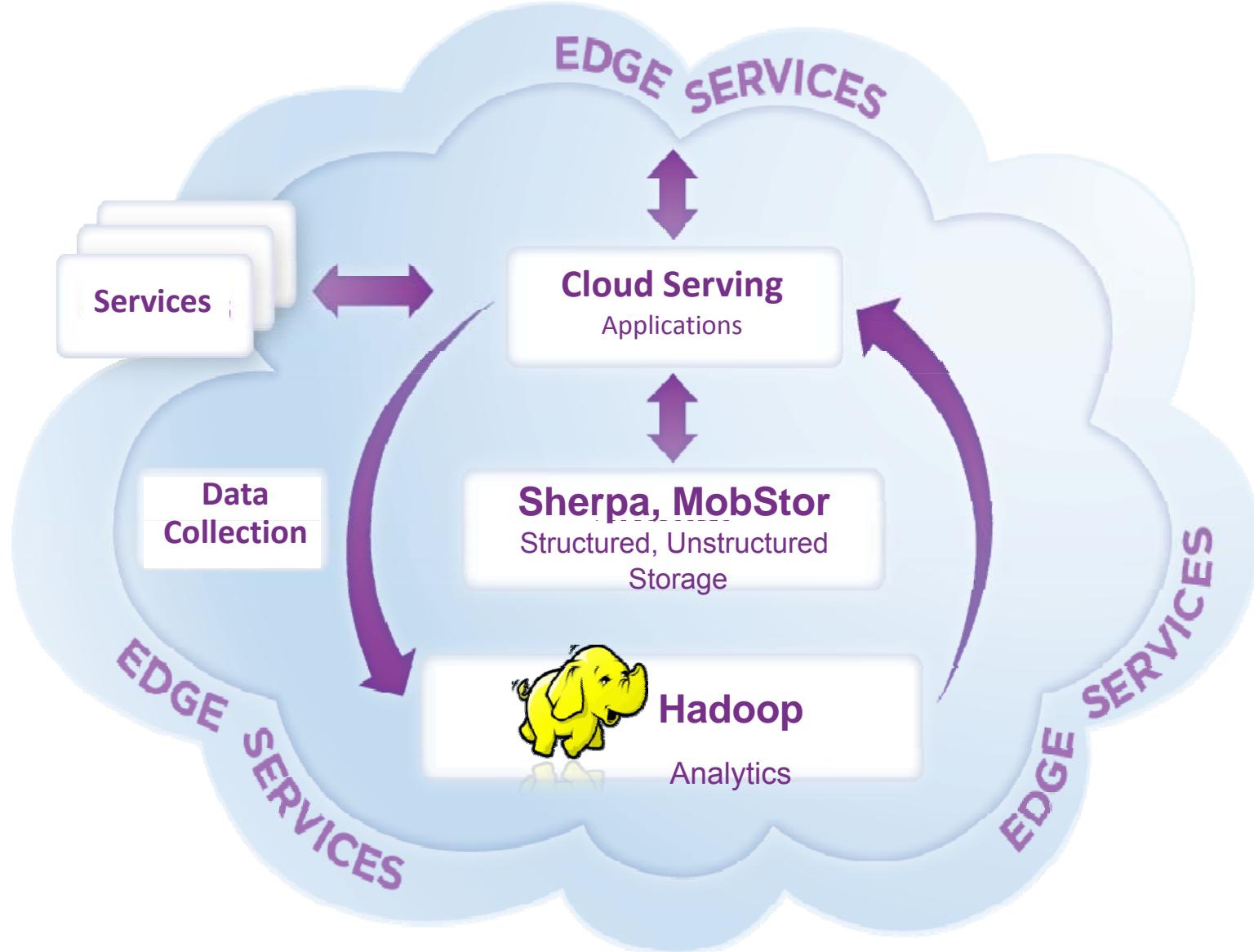


Yahoo! Cloud Stack



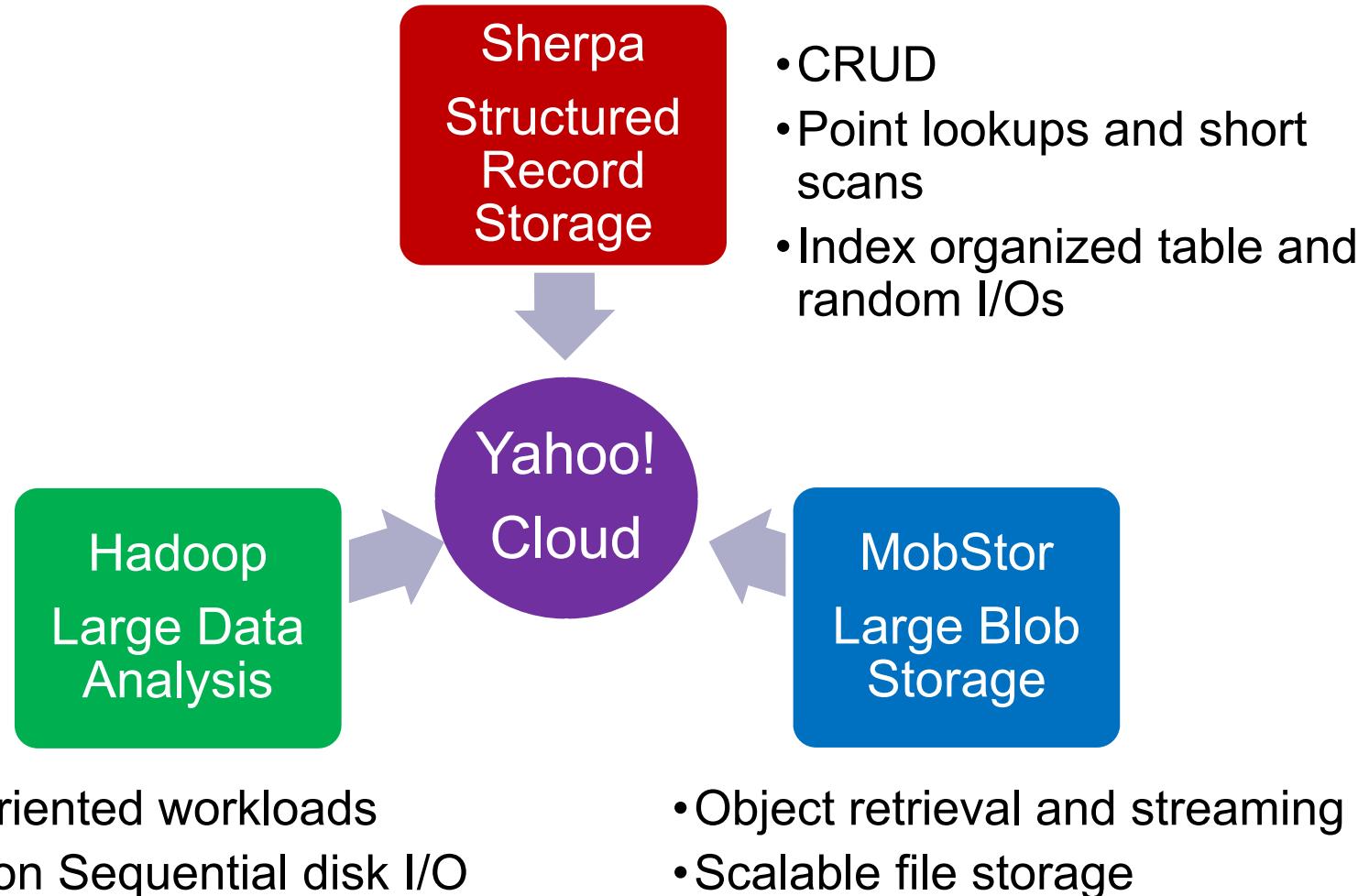


A Data-Centric View



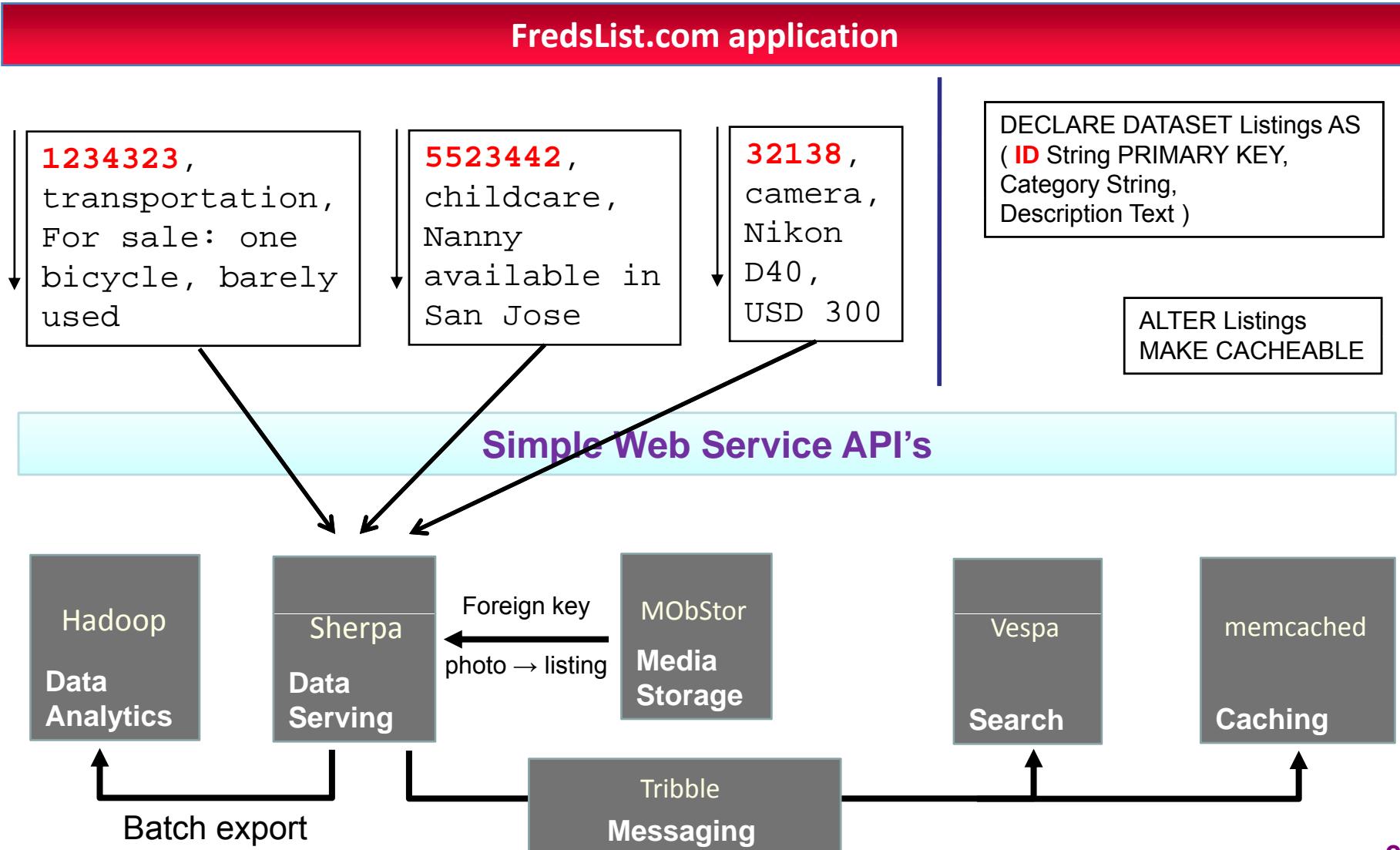


Cloud Data Management





The Yahoo! Data Cloud





Hadoop Core
(Core, Pig, Oozie,
Hive, Howl)

Ad BT and Inventory prediction, Content
Agility, UDA, COKE, Mail Spam, Search,
APG, Labs, Insights, Analytics

1+ million jobs per month
3.7 PB processed daily
90B events and 120 TB daily
70+ PB of Data

Map-Reduce and more ...

HADOOP: SCALABLE ANALYTICS





One Slide Hadoop Primer

Good for analyzing (scanning) huge files
Not great for serving (reading or writing individual objects)

Industry Challenge: Massive Data Sets Being Accumulated

"It's now the industrial revolution of data" – Joe Hellerstein, UC Berkeley

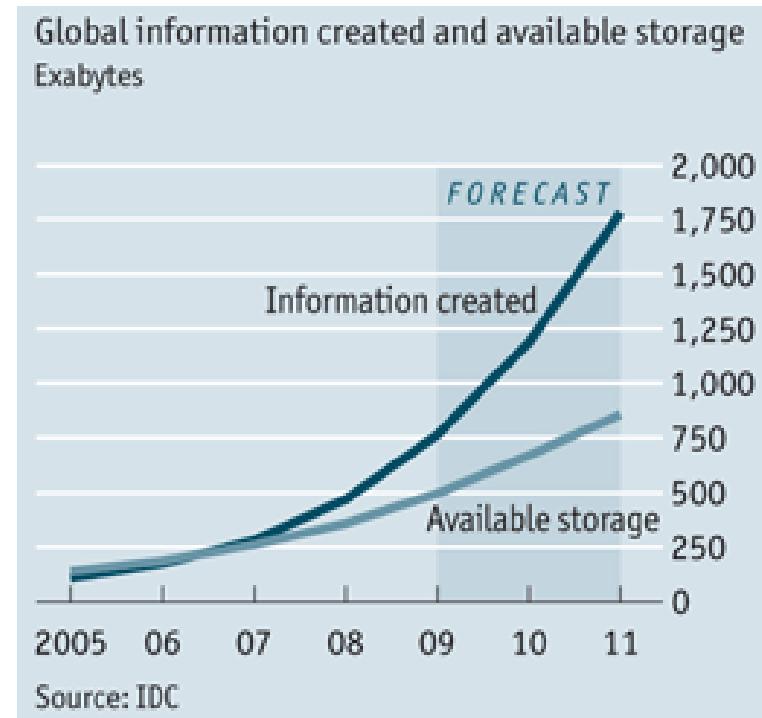
1 Massive datasets of highly-leveragable data amassed by relatively few players

"By 2020, the Digital Universe will be 44 times as large as it was in 2009" [IDC]

2 Deriving insights from the newly available mountains of data still challenging

3 Datasets become competitive differentiators

*Search query history,
Advertising/click-through data,
Surfing histories, Social graphs,
News, Finance, Sports and other
data feeds*



Yahoo!'s technology solution: Hadoop



Hadoop: Stability at Scale

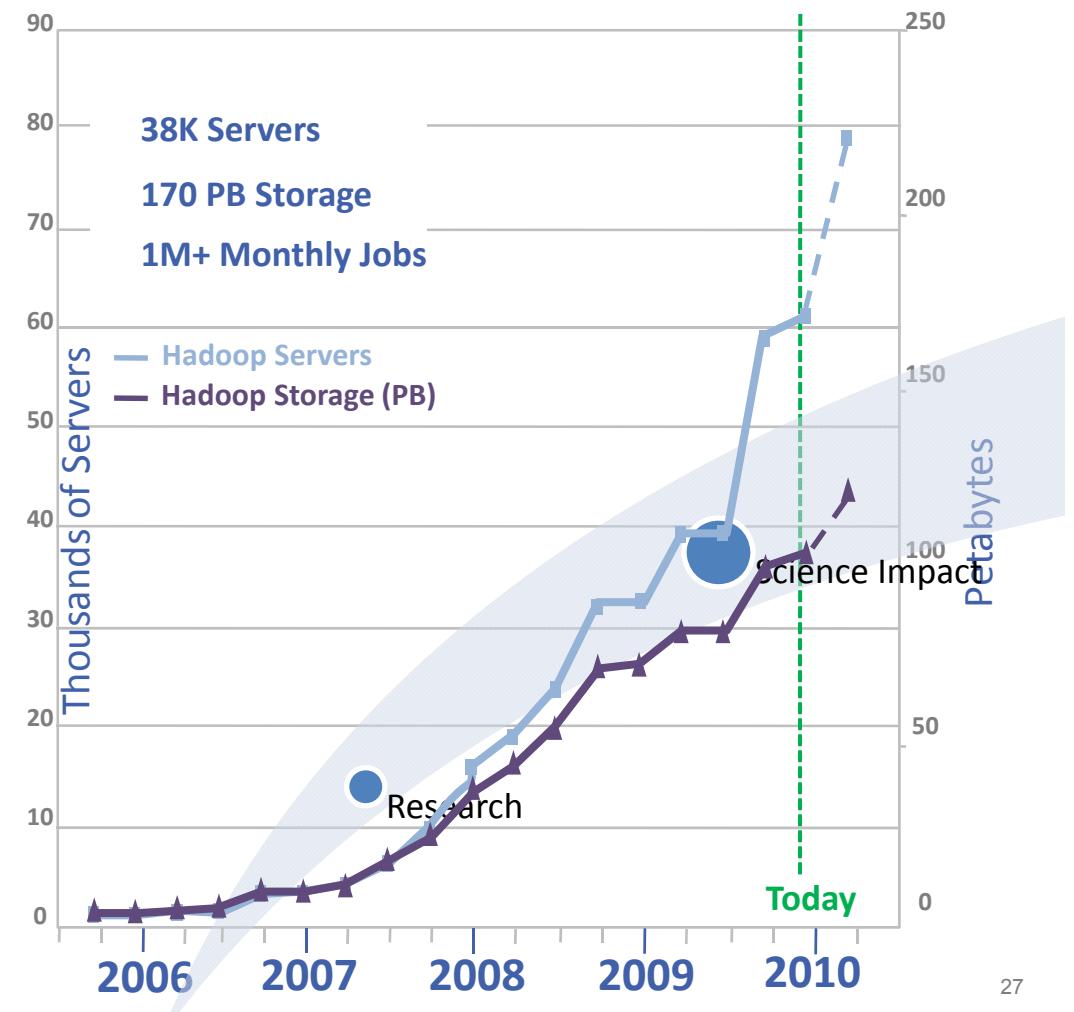
Hadoop powers the Yahoo! Network: must be rock-solid

We fix bugs before you see them

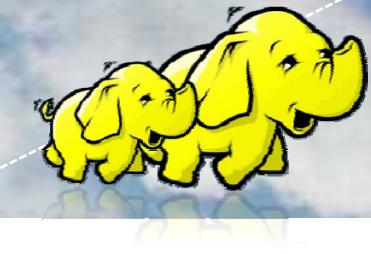
- We run very large clusters
- We have a large QA effort
- We run a huge variety of workloads

The Yahoo! Distribution of Hadoop

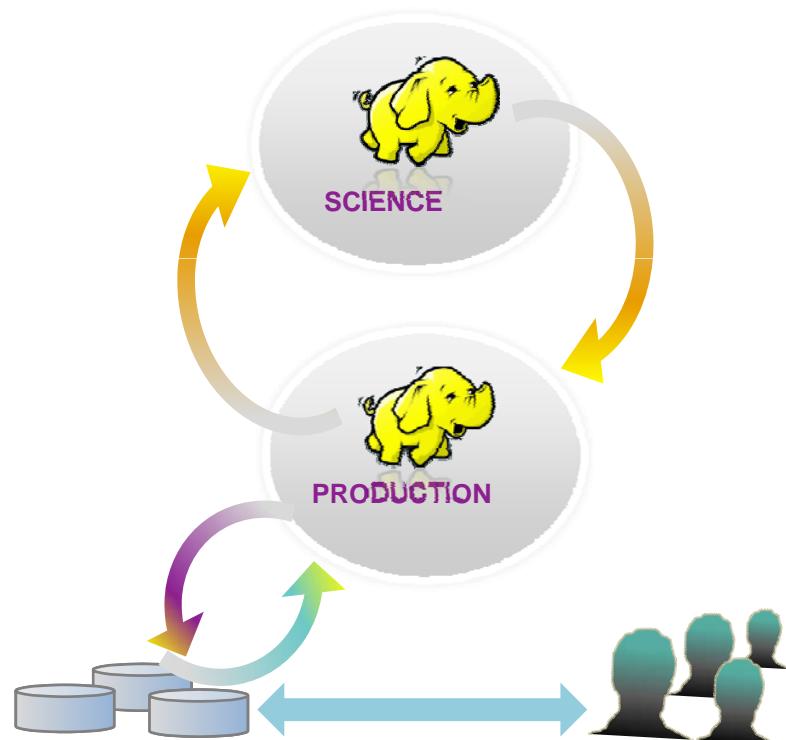
- We contribute our work to Apache
- We share the exact code we run
- We don't sell software or service



CASE STUDY YAHOO! MAIL



Enabling quick response in the spam arms race



- 450M mail boxes
- 5B+ deliveries/day
- Antispam models retrained every few hours on Hadoop

“ 40% less spam than Hotmail and 55% less spam than Gmail ”

Example: User Activity Modeling

Large dimensionality vector describing possible user activities
But a typical user has a sparse activity vector

Attribute	Possible Values	Typical values per user
Pages	~ MM	10 – 100
Queries	~ 100s of MM	Few
Ads	~ 100s of thousands	10s

Hadoop pipeline to model user interests from activities



Feature and Target Windows



Query

Visit Y! finance

T_0

Event of interest

Time

Moving Window



Feature Window



Target Window



User Modeling Pipeline

Component	Data Processed	Time
Data Acquisition	~ 1 Tb per time period	2 – 3 hours
Feature and Target Generation	~ 1 Tb * Size of feature window	4 - 6 hours
Model Training	~ 50 - 100 Gb	1 – 2 hours for 100's of models
Scoring	~ 500 Gb	1 hour

A Growing User Base

Year: 2007

YAHOO!



last.fm

Year: 2008

Google™

ablegrape

Cascading

INFORMATION SCIENCES INSTITUTE

zvents
Discover Things To Do

Cornell University Computing and Information Science

Vidalia

ImageShack online media hosting

A9

IBM®

ENORMO Every property. Everywhere.

The New York Times

Lookery Control freaks welcome

facebook

krugle

THE UNIVERSITY OF EDINBURGH

SECURITY ENHANCED DOMAIN NAME SYSTEM

NetSeer

News Corporation

LOTAME Locate, Target, & Message with Social Media

veoh™

parc® Palo Alto Research Center

Rackspace HOSTING

Year: 2009 - 2010

AOL

cloudera

cooliris

deepdyve BETA

TEXTMAP THE ENTITY SEARCH ENGINE

PSG College of Technology The Technology that looks the world

eyealike

iterend

tailsweep

hulu™ RapLeaf

Ning quantcast

USCIMS

VK SOLUTIONS Global Solutions Provider

NETFLIX

Microsoft

Terrier

acknowledge

LinkedIn

stampede beta

WorldLingo

TARAGANA Innovation • Quality • Simplicity

HOSTING HABITAT

HOLA SERVERS



Y!OS, COKE, LocDrop, Video, Media
Search history, Answers, Messenger,
BOSS, Image Search, Blog Search

15K requests per second
Over 1.5B records; 10sTB of data

ACID or BASE? Litmus tests are colorful, but the picture is cloudy

PNUTS: SCALABLE DATA SERVING





Requirements for Cloud Services

- **Multitenant.** A cloud service must support multiple, organizationally distant customers.
- **Elasticity.** Tenants should be able to negotiate and receive resources/QoS *on-demand* up to a large scale.
- **Resource Sharing.** Ideally, spare cloud resources should be transparently applied when a tenant's negotiated QoS is insufficient, e.g., due to spikes.
- **Horizontal scaling.** The cloud provider should be able to add cloud capacity in increments without affecting tenants of the service.
- **Metering.** A cloud service must support accounting that reasonably ascribes operational and capital expenditures to each of the tenants of the service.
- **Security.** A cloud service should be secure in that tenants are not made vulnerable because of loopholes in the cloud.
- **Availability.** A cloud service should be highly available.
- **Operability.** A cloud service should be easy to operate, with few operators. Operating costs should scale linearly or better with the capacity of the service.



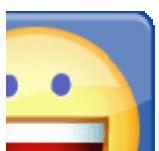
The World Has Changed

Web serving applications need:

- Scalability!
 - Elastic on demand, commodity boxes
- Flexible schemas
- Geographic distribution/replication
- High availability
- Low latency

Web serving applications willing to do without:

- Complex queries
- ACID transactions
 - But still benefit from support for data consistency





Typical Y! Applications



User logins and profiles

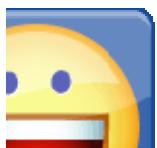
- Including changes that must not be lost!
 - But single-record “transactions” suffice

Events

- Alerts (e.g., news, price drops)
- Social network activity (e.g., user goes offline)
- Ad clicks, article clicks

Application-specific data

- Postings in message board
- Uploaded photos, tags
- Shopping carts



640M+ unique users, 11B pages/month
Hundreds of petabytes of storage
Hundreds of billions of objects
Hundred of thousands of requests/sec
Global, rapidly evolving workloads

What is PNUTS/Sherpa?

The diagram is enclosed in a red border and contains four main sections, each with a title and associated images:

- CREATE TABLE Parts (**
ID VARCHAR,
StockNumber INT,
Status VARCHAR
...
)
- Structured, flexible schema**
- Parallel database**
- Geographic replication**
- Hosted, managed infrastructure**

Each section includes a table with data and server rack icons. The tables have blue header rows and white data rows. The servers are shown in pairs, with one server in each pair having blue components and the other having red components.

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

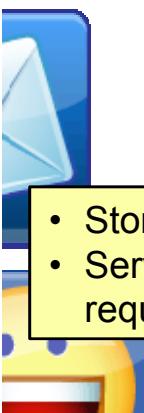
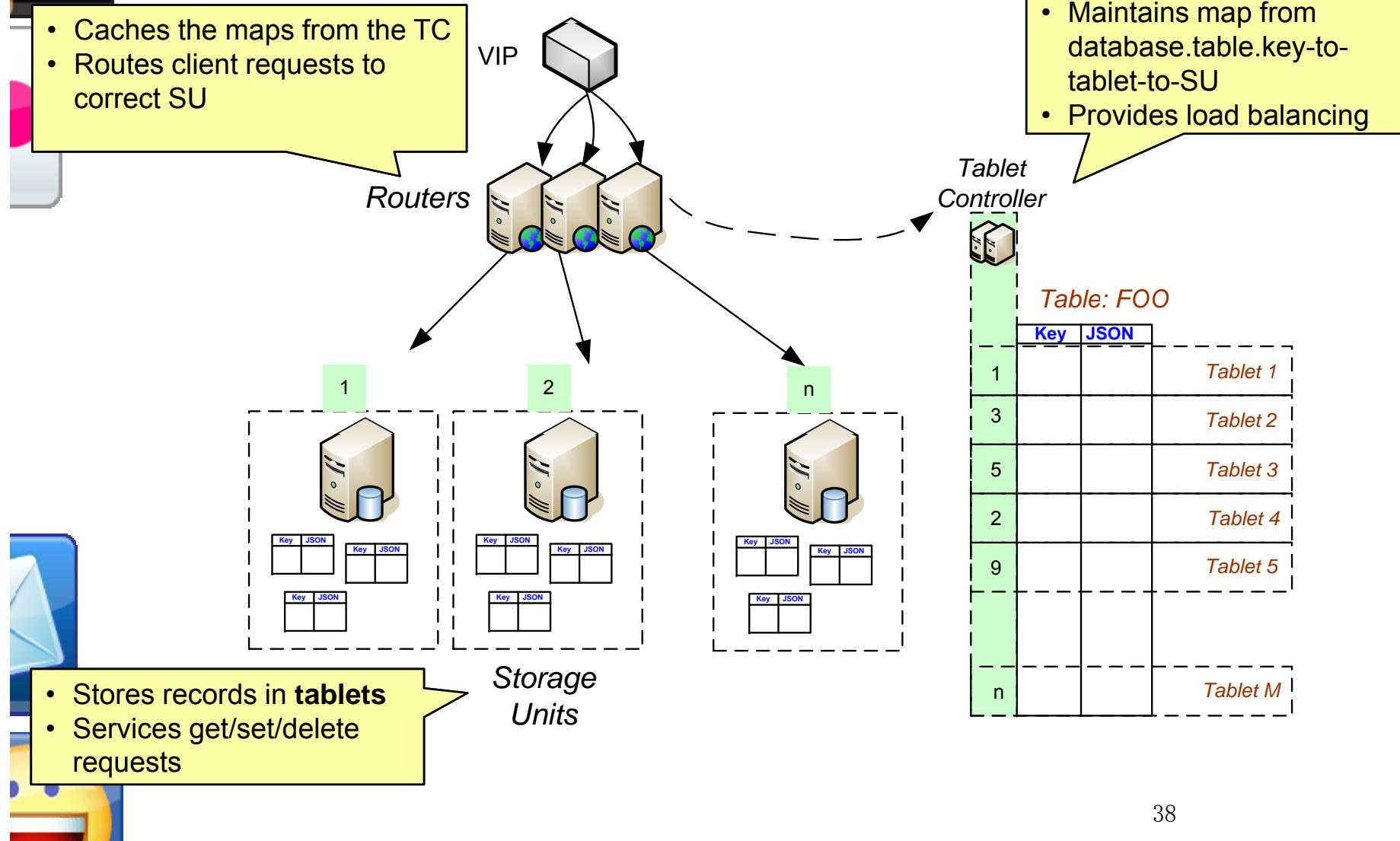
A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75656	C
F	15677	E

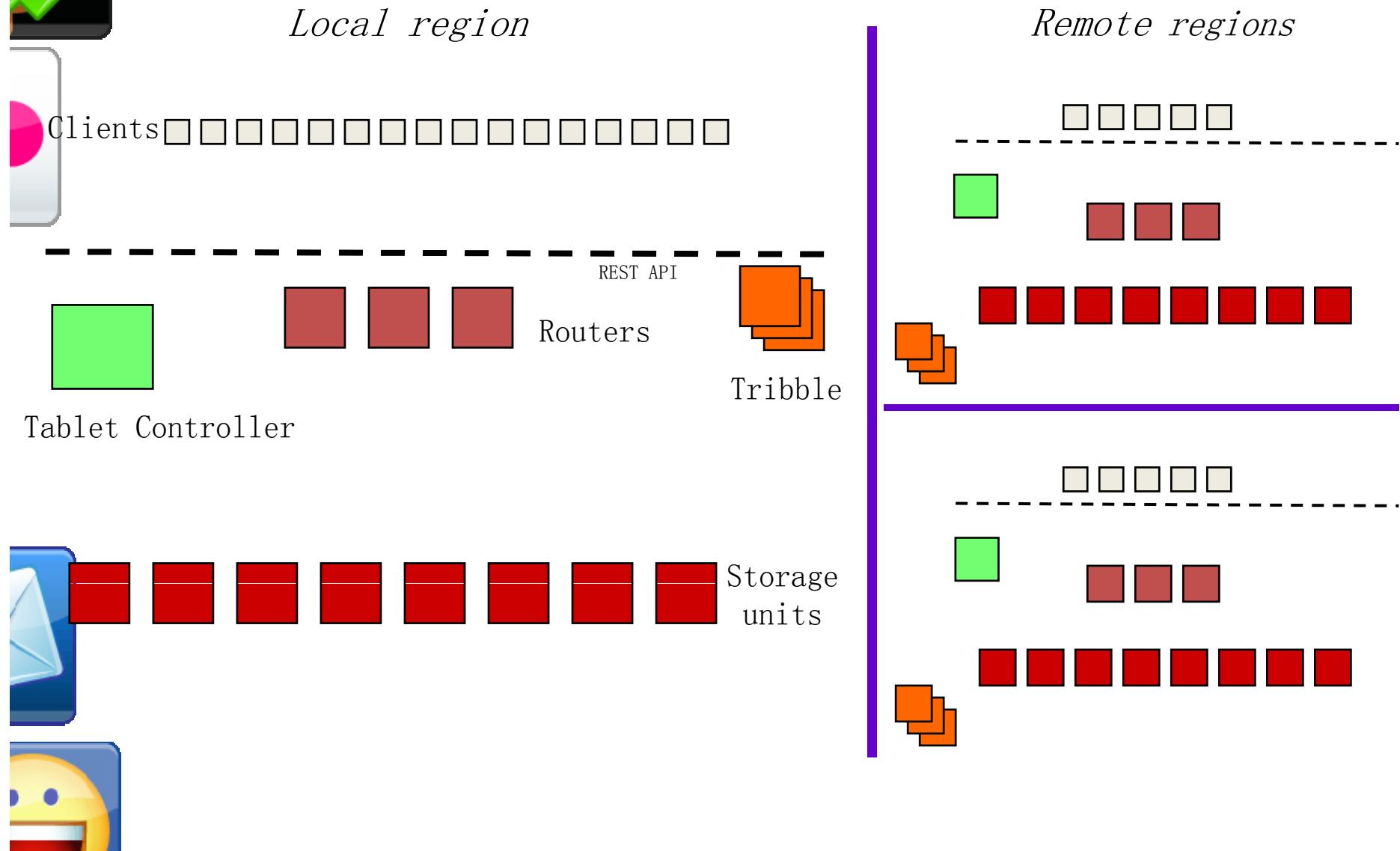
A	42342	E
B	42521	W
C	66354	W
D	12352	E
E	75	

PNUTS: Key Components





Architecture





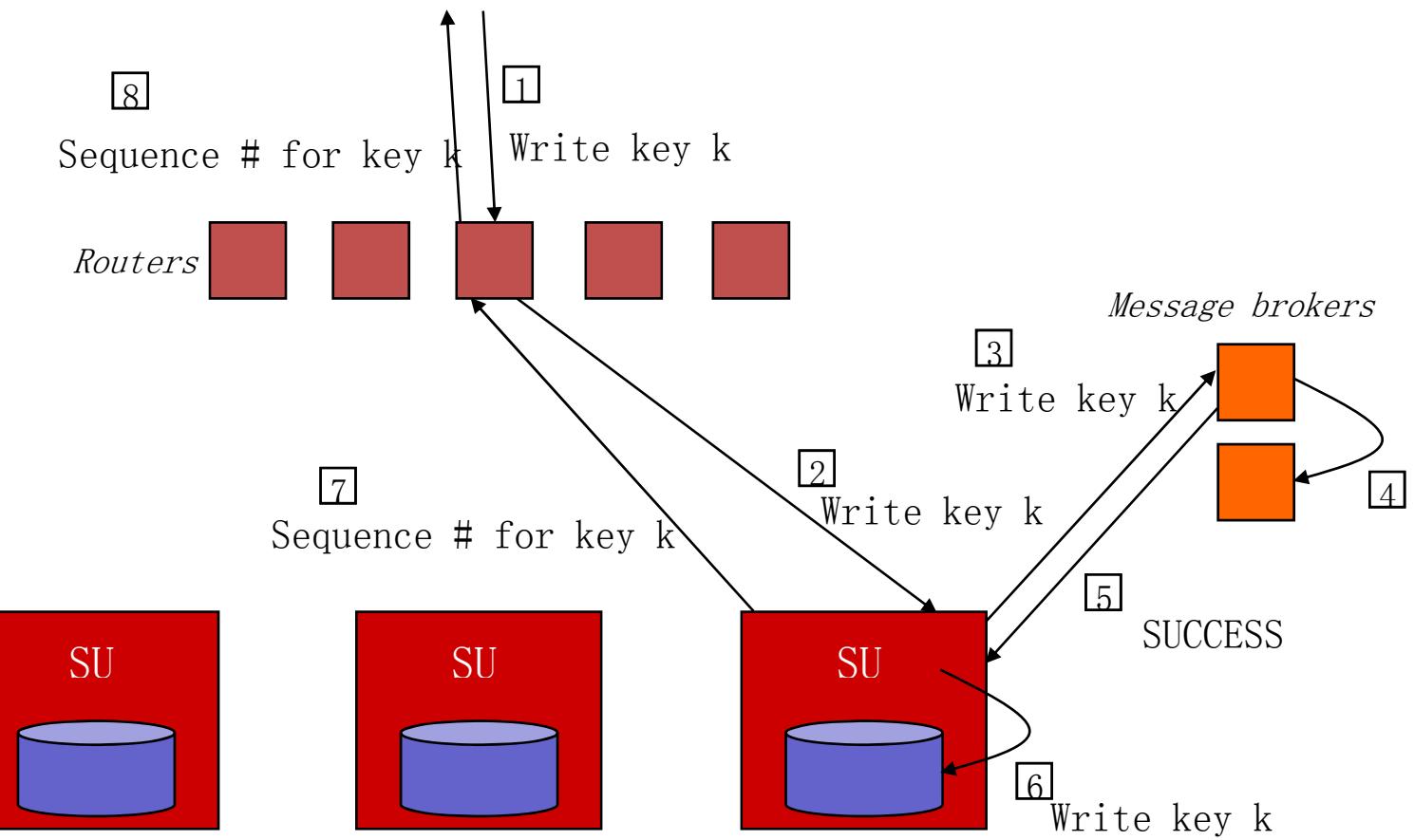
Flexible Schema



<i>Posted date</i>	<i>Listing id</i>	<i>Item</i>	<i>Price</i>	<i>Color</i>	<i>Condition</i>
6/1/07	424252	Couch	\$570		Good
6/1/07	763245	Bike	\$86		
6/3/07	211242	Car	\$1123	Red	Fair
6/5/07	421133	Lamp	\$15		



Updates

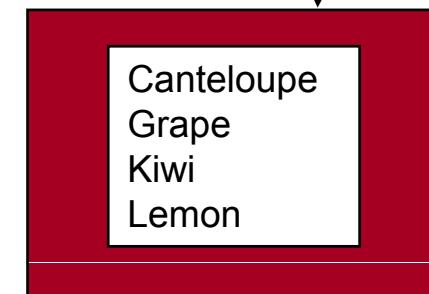
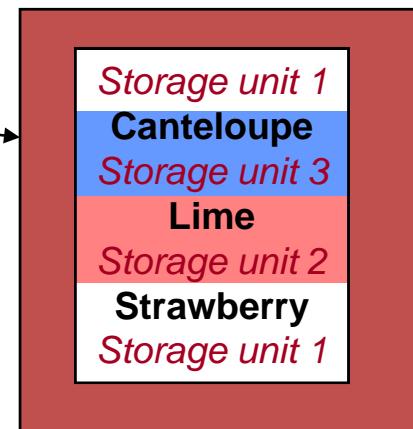
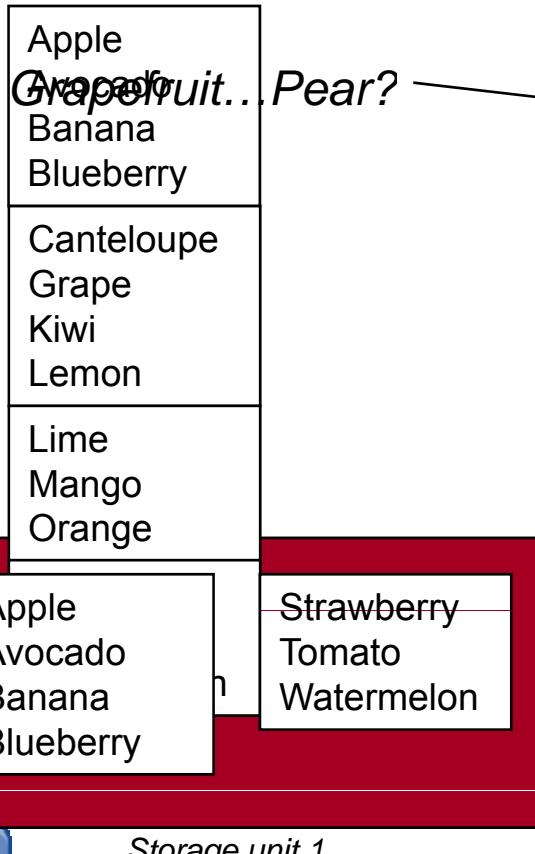
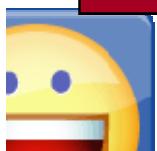


Tablets—Ordered Table

	<i>Name</i>	<i>Description</i>	<i>Price</i>
A	Apple	Apple is wisdom	\$1
	Avocado	But at what price?	\$3
	Banana	The perfect fruit	\$2
H	Grape	Grapes are good to eat	\$12
	Kiwi	New Zealand	\$8
	Lemon	How much did you pay for this lemon?	\$1
	Lime	Limes are green	\$9
Q	Orange	Arrgh! Don't get scurvy!	\$2
	Strawberry	Strawberry shortcake	\$900
Z	Tomato	Is this a vegetable?	\$14

Range Queries in YDOT

Clustered, ordered retrieval of records





ELASTICITY, OPERABILITY, HORIZONTAL SCALING





Distribution



6/1/07	Data shuffling for load balancing		
6/1/07	256623	Car	\$1123
6/2/07	636353	Bike	\$86
6/5/07	662113	Chair	\$10
6/7/07	121113	Lamp	\$19
6/9/07	887734	Bike	\$56
6/11/07	252111	Scooter	\$18
6/11/07	116458	Hammer	\$8000



Server 1



Server 2



Server 3



Server 4

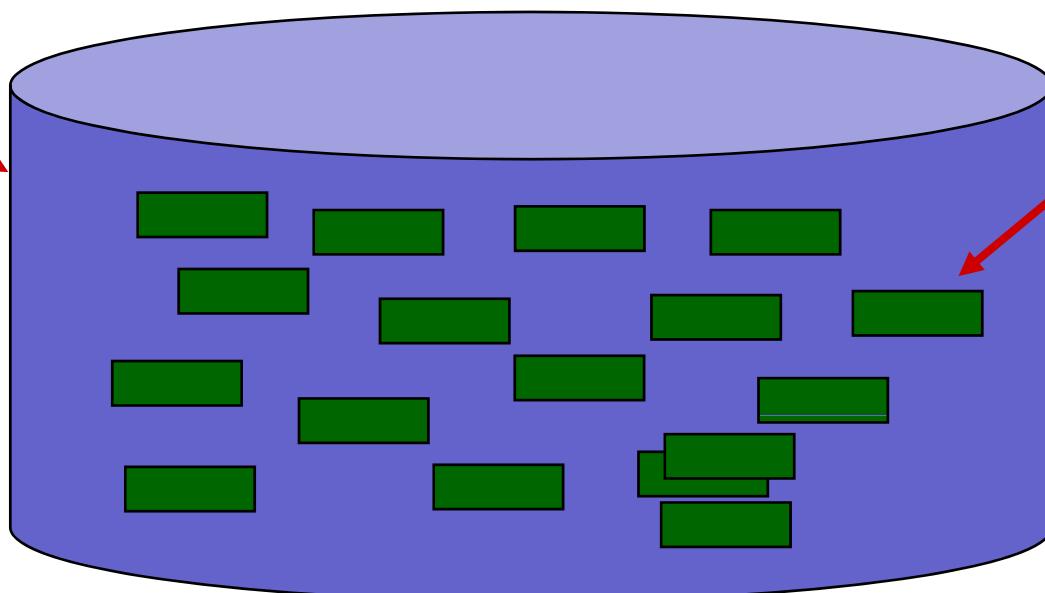


Tablet Splitting and Balancing

Each storage unit has many tablets (horizontal partitions of the table)
Storage unit may become a hotspot

Storage unit

Tablet



Overfull tablets split

Tablets may grow over time

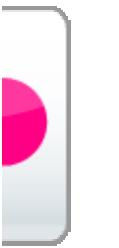
Shed load by moving tablets to other servers



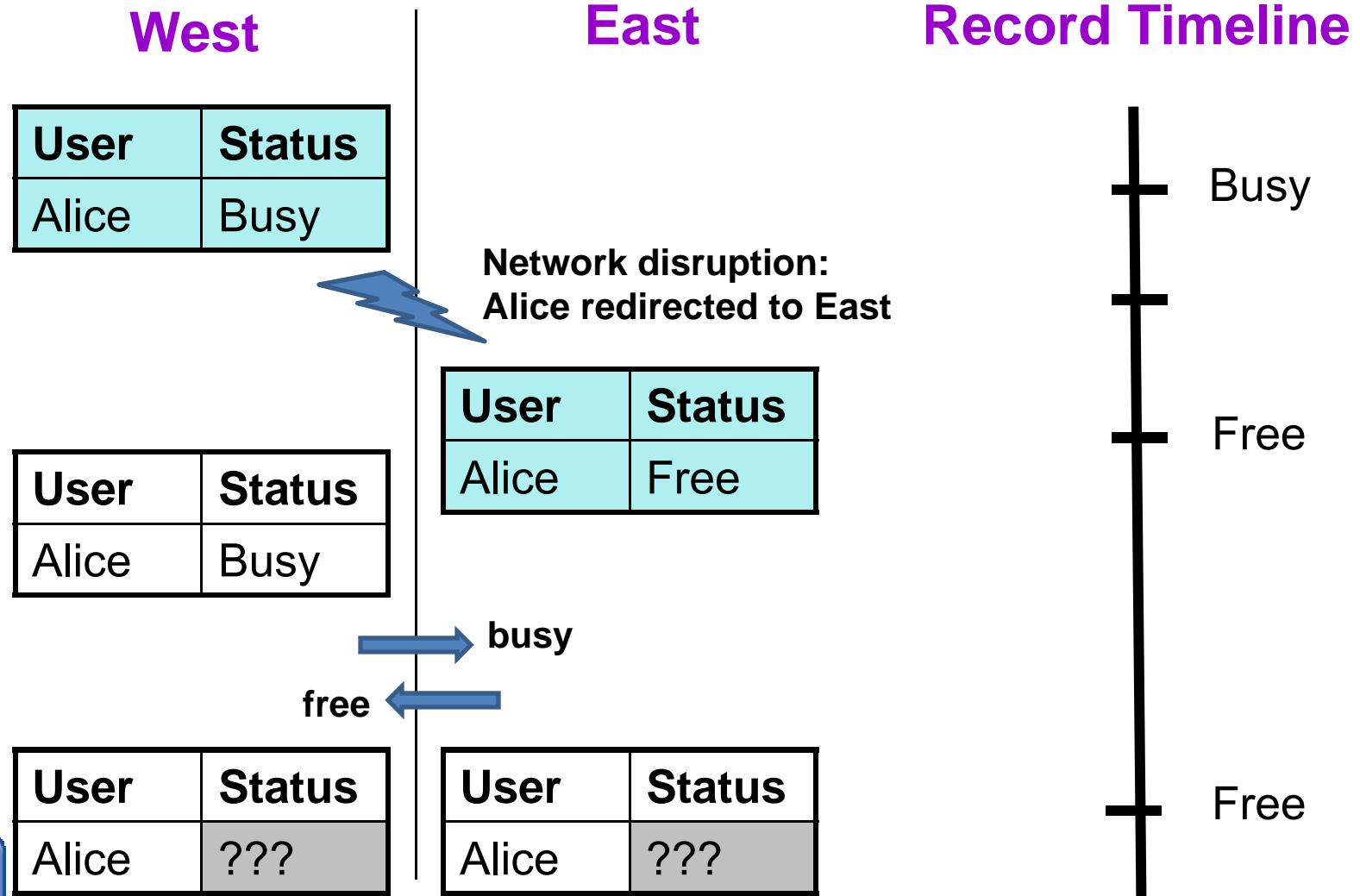
ASYNCHRONOUS REPLICATION AND CONSISTENCY



Asynchronous Replication

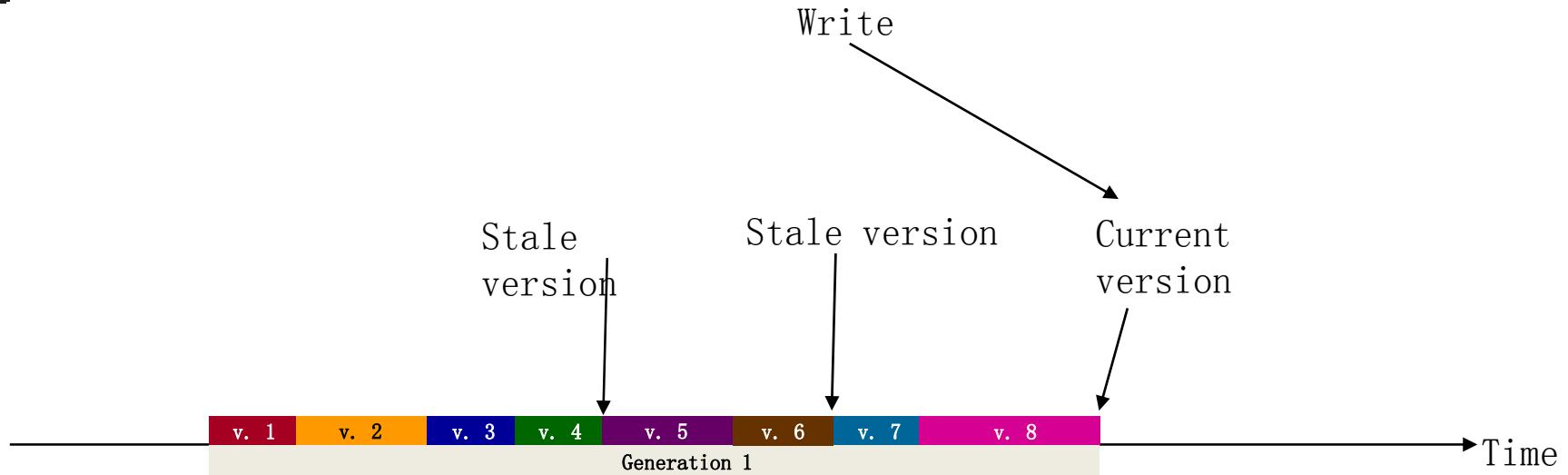


Consistency: Social Alice





PNUTS Consistency Model



Achieved via per-record primary copy protocol
(To maximize availability, record masterships automatically transferred if site fails)

Can be selectively weakened to eventual consistency
(local writes that are reconciled using version vectors)

Record Master

A	42342	
B	42521	
C	66354	
D	12352	
E	75656	
F	15677	



A	42342	
B	42521	
C	66354	
D	12352	
E	75656	
F	15677	

A	42342	
B	42521	
C	66354	
D	12352	
E	75656	
F	15677	



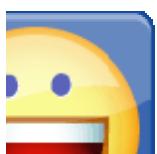
Consistency Techniques

Per-record mastering

- Each record is assigned a “master region”
 - May differ between records
- Updates to the record forwarded to the master region
- Ensures consistent ordering of updates

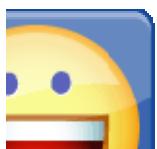
Tablet-level mastering

- Each tablet is assigned a “master region”
- Inserts and deletes of records forwarded to the master region
- Master region decides tablet splits



These details are hidden from the application

- Except for the latency impact!



Consistency Techniques

Availability
Consistency

Primary Key Constraint + Record Timeline

- Each tablet is assigned a “master region”
- Inserts of records forwarded to the master region
- Inserts and updates could fail during outages*

Record Timeline Consistency

- Each record is assigned a “master region”
- Updates to the record forwarded to the master region
- Inserts succeed, but updates could fail during outages*

Eventual Consistency

- Low latency updates and inserts done locally
- Per field timestamp used to merge updates

- ★ In case of SU or data center failure. We have failover tools!
- ★ Reads always will be sent to another region

Tablet Master

Region W

Key1	42342	E
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Key2: 42521



Region C

Tablet master

Key1	42342	E
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Region E

Key1	42342	E
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E



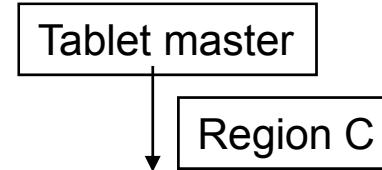
Tablet Mastership



Region W

Key1	42342	E
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Step 1: Forward
Req to Tablet Master



Key1	42342	E
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Region E

Key1	42342	E
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Step 2: Apply
Insert to Tablet Master



Key1	42342	E
Key2	42521	W
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Key1	42342	E
Key2	42521	W
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Step 4: Apply
Insert at Rec Master



Key1	42342	E
Key2	42521	W
Key3	66354	W
Key4	12352	E
Key5	75656	C
Key6	15677	E

Step 3: Replicate
Insert to Other Sites





Generalizing Record Timelines to Partition Timelines

Record → Partition of records with same key

- Tablet splits must respect partition boundaries
- Intra-partition ACID transactions can be done easily now
 - Single machine transactions!
 - With composite keys, this captures Azure and Google AE models
- Each partition is assigned a “master region”
 - May differ between partitions
- Updates to the partition forwarded to the master region
- Ensures consistent ordering of updates across nodes





AVAILABILITY



Possible Failure Modes

Failure type

Storage unit

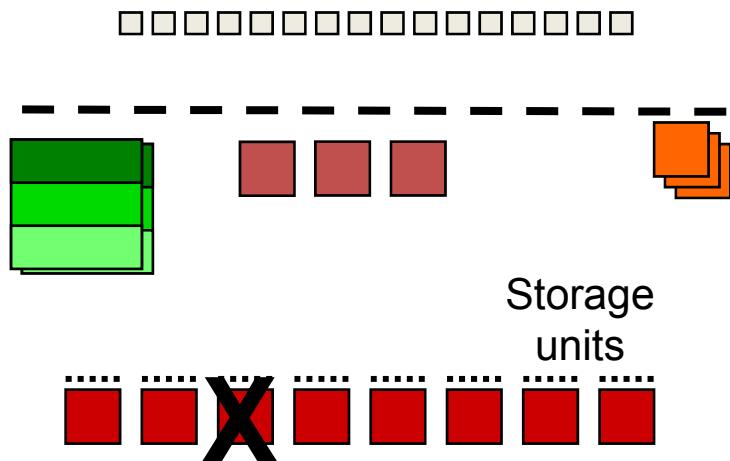
Consistency impact

None

Availability impact

Degraded service (forwards) for some data.

Updates and inserts fail for some records



Resolution

If data not lost: Reboot machine

If data lost: Copy lost tablets from a remote replica

Time to resolve

If data lost, hours or less (depending on tablet size and colo location). If no data lost, minutes.





Coping With Failures

A map of the United States with several data overlays. On the West Coast, there is a table with a large red X drawn over it. The table has rows labeled A through F and columns labeled W, E, C, and E. The values are:

A	2342		
B	42521		
C	66354		
D	12352		
E	75656		
F	15677		

In the center of the map, there is a small image of a woman sitting at a desk with a laptop. A white arrow points from the red X on the West Coast table to a table on the East Coast labeled "OVERRIDE W → E". This table has rows labeled A through F and columns labeled E, W, W, E, C, and E. The values are:

OVERRIDE W → E		
A	42342	
B	42521	
C	66354	
D	12352	
E	75656	
F	15677	

At the bottom of the map, there is another table with green cells. It has rows labeled A through F and columns labeled E, W, W, E, C, and E. The values are:

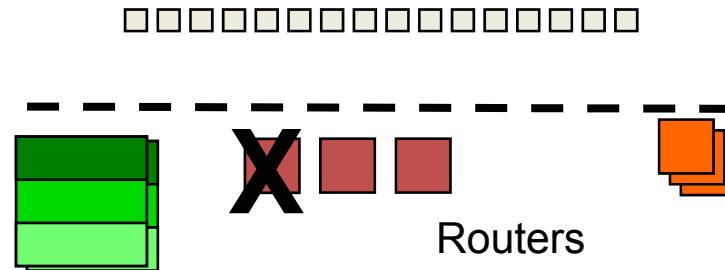
A	42342	
B	42521	
C	66354	
D	12352	
E	75656	
F	15677	



Possible Failure Modes

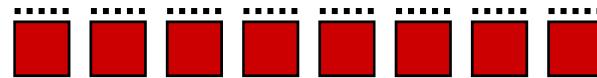
Failure type

Router



Consistency impact

None



Availability impact

None

Resolution

Boot router



Time to resolve

Minutes



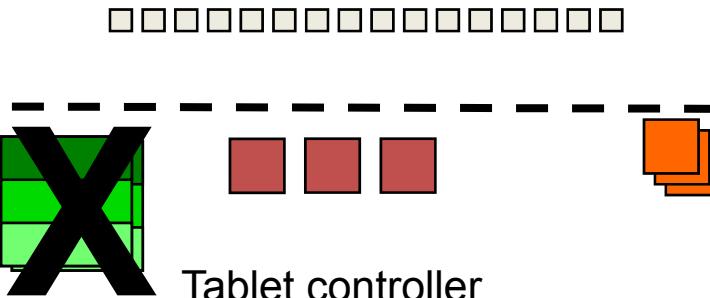


Possible Failure Modes



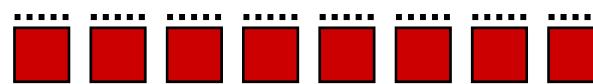
Failure

Tablet controller



Consistency impact

None



Availability impact

Some actions (e.g., tablet copy) will be blocked

Resolution



Start secondary controller



Time to resolve

Minutes

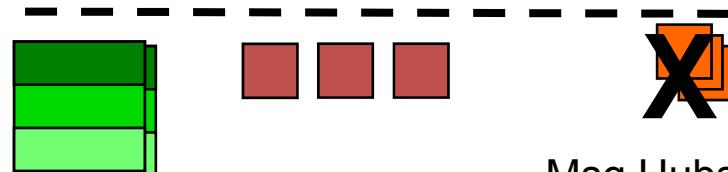


Possible Failure Modes



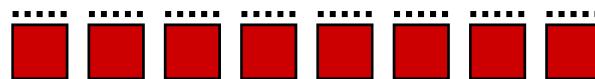
Failure

One msg hub node



Consistency impact

None



Availability impact

Writes fail for some records until a new secondary node takes over

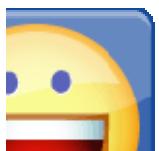
Resolution



Create new primary or secondary for lost topics

Time to resolve

Minutes



Possible Failure Modes



Failure

Colo power outage or partition



Consistency impact

Option to allow “relaxed consistency”
to improve availability

Availability impact

Some inserts, updates and
deletes cannot succeed

Some critical reads fail

Option to allow updates to proceed in
“relaxed consistency mode”



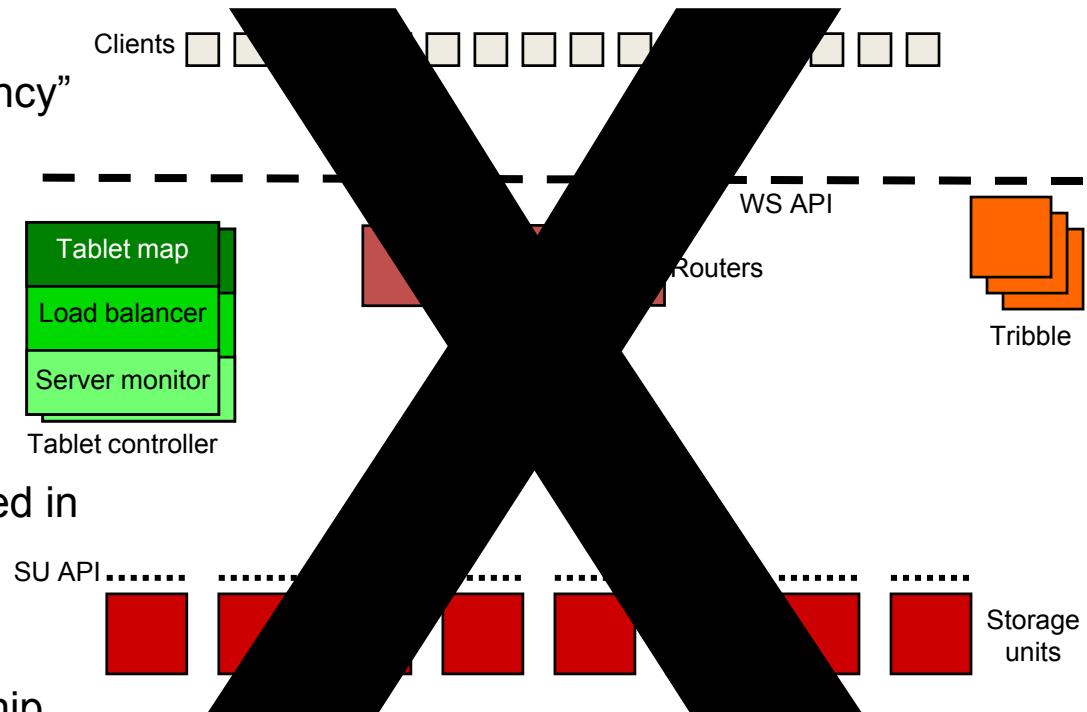
Resolution

Major overrides to force mastership
transfer; discard conflicting updates



Time to resolve

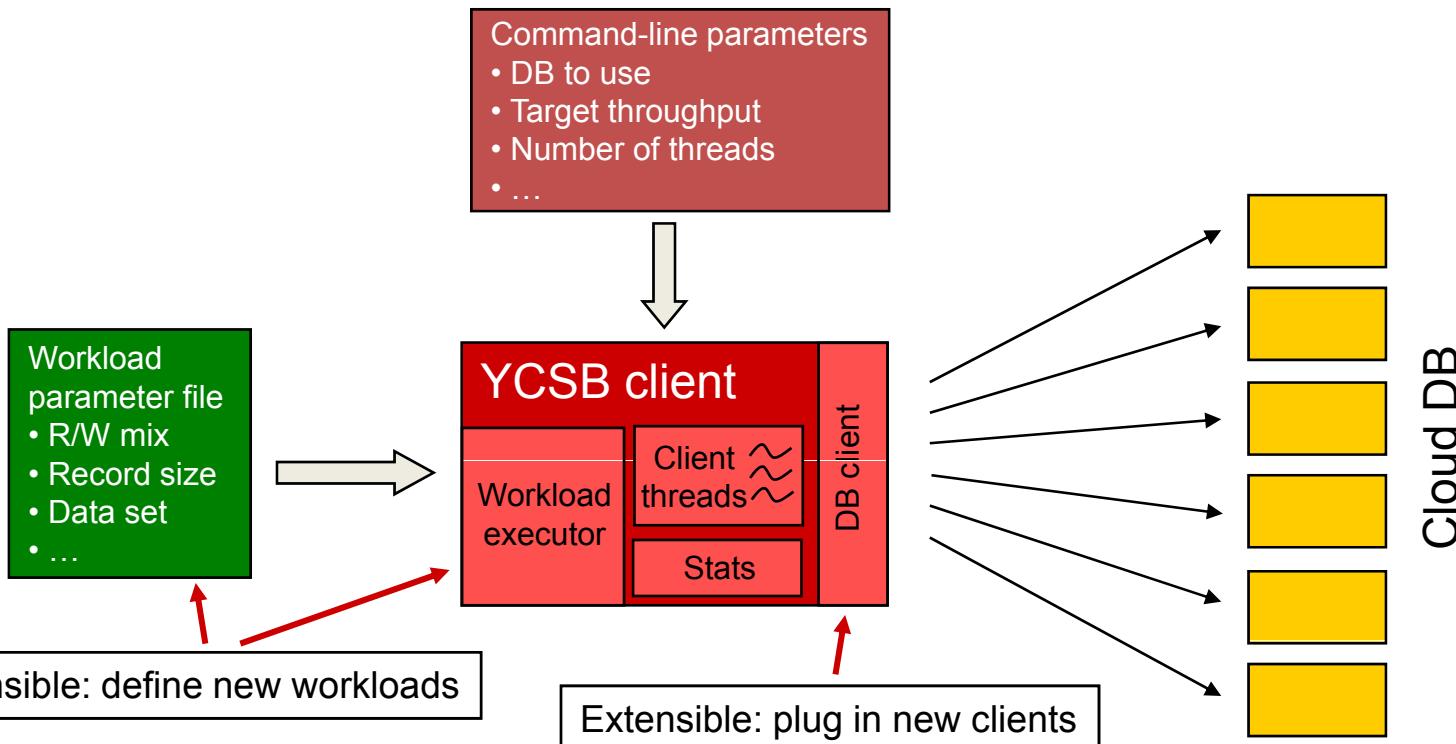
Hours



YCSB Benchmark Tool

Java application

- Many systems have Java APIs
- Other systems via HTTP/REST, JNI or some other solution





Further Reading on PNUTS



Efficient Bulk Insertion into a Distributed Ordered Table (SIGMOD 2008)
Adam Silberstein, Brian Cooper, Utkarsh Srivastava, Erik Vee,
Ramana Yerneni, Raghu Ramakrishnan

PNUTS: Yahoo!'s Hosted Data Serving Platform (VLDB 2008)
Brian Cooper, Raghu Ramakrishnan, Utkarsh Srivastava,
Adam Silberstein, Phil Bohannon, Hans-Arno Jacobsen,
Nick Puz, Daniel Weaver, Ramana Yerneni

Asynchronous View Maintenance for VLSD Databases (SIGMOD 2009)
Parag Agrawal, Adam Silberstein, Brian F. Cooper, Utkarsh Srivastava and
Raghu Ramakrishnan

Cloud Storage Design in a PNUTShell (**Beautiful Data, O'Reilly Media, 2009**)
Brian F. Cooper, Raghu Ramakrishnan, and Utkarsh Srivastava



Adaptively Parallelizing Distributed Range Queries (VLDB 2009)
Ymir Vigfusson, Adam Silberstein, Brian Cooper, Rodrigo Fonseca



**A Batch of PNUTS: Experiences Connecting Cloud Batch and
Serving Systems (SIGMOD 2011)**
Adam Silberstein et al.



Summary

- Yahoo! has an extensive cloud computing environment
 - Major open source contributor
 - Currently, focused on internal developers
 - Foundation layer for all Yahoo! products, e.g., Mail, Front Page, Search, Groups ...