# Web Data Management: Powering the New Web

**Raghu Ramakrishnan**

**Chief Scientist for Audience, Yahoo!**

**Research Fellow, Yahoo! Research**

**(On leave, Univ. of Wisconsin-Madison)**

# Outline

- **Trends in Search and Information Discovery**

  – Move towards task-centricity

  – Need to interpret content

- **Evolution of Online Communities**

  – Social Search

  – PeopleWeb

  – Community Information Management

- **Web Data Infrastructure**

  – Massively distributed computing

  – Hosted services

  – Heterogeneous content

# Further Reading

- Content, Metadata, and Behavioral Information: Directions for Yahoo! Research, *The Yahoo! Research Team*, IEEE Data Engineering Bulletin, Dec 2006 (Special Issue on Web-Scale Data, Systems, and Semantics)

- Systems, Communities, Community Systems, on the Web, *Community Systems Group at Yahoo! Research*, SIGMOD Record, Sept 2007

- Towards a PeopleWeb, R. Ramakrishnan and A. Tomkins, IEEE Computer, August 2007 (Special Issue on Web Search)

**Research**

# Community Systems Group @ Yahoo! Research

| | | |
|---|---|---|
| **Philip Bohannon** | Deepak Agrawal | Parag Agrawal |
| **Brian Cooper** | Sihem Amer-Yahia | Tyson Condie |
| **Nilesh Dalvi** | Ravi Kumar | Pedro DeRose |
| **Minos Garofalakis** | Cameron Marlow | Alban Galland |
| **Hans-Arno Jacobsen** | Srujana Merugu | Nitin Gupta |
| **Vinay Kakade** | Chris Olston | Ashwin Machanavajjhala |
| **Dan Kifer** | Bo Pang | Warren Shen |
| **Raghu Ramakrishnan** | Ben Reed | Julia Stoyanovich |
| **Adam Silberstein** | Keerthi Selvaraj | Fan Yang |
| **Utkarsh Srivastava** | Jai Shanmugasundaram | |
| **Ramana Yerneni** | Andrew Tomkins | |
| **Cong Yu** | | |

Trends in Search

# Search and Content Supply

- Premise:
  - People don't want to search
  - People want to get tasks done

I want to book a vacation in Tuscany.

Start                                                                    Finish

Google   YAHOO! FARE   TuscanDream *the ancient soul of Italian art, fashion and culture*   SiXT

# Y! Shortcuts

# Google Base

# Structure ➡ Intent

"seafood san francisco"

↳ Category: restaurant
Location: San Francisco

Reserve a table for two tonight at SF's best Sushi Bar and get a free sake, compliments of OpenTable!

Category: restaurant Location: San Francisco

Alamo Square **Seafood** Grill –  (415) 440-2828
803 Fillmore St, **San Francisco**, CA – 0.93mi – map

Category: restaurant Location: San Francisco

# Steps to Task-Centricity

- **Information integration**
  - Information extraction
  - Schema normalization

- **Structure**
  - Extract and exploit

hotel near leicester square

Search

**Welcome to The Savoy**
Located on The Strand in the heart of the West End theatre district,

# Semantic Structure – Not Easy

- Colorado, Texas

- Oregon, Alaska

- Peru, Bolivia

- Peru, Argentina

- Washington, Nevada County, California

- Bush, Cheney & Rice, WA state

**Research**

# How Do We Circumvent?

- Unleash community computing
  - Social structure
  - Incentive mechanisms
- More in a moment

Evolution of Online Communities

# Rate of content creation

- Estimated growth of content
  - Published content from traditional sources: 3-4 Gb/day
  - Professional web content: ~2 Gb/day
  - User-generated content: 8-10 Gb/day
  - Private text content: ~3 Tb/day (200x more)
  - Upper bound on typed content: ~700 Tb/day

# Metadata

- Estimated growth of metadata
    - Anchortext: 100Mb/day
    - Tags: 40Mb/day
    - Pageviews: 100-200Gb/day
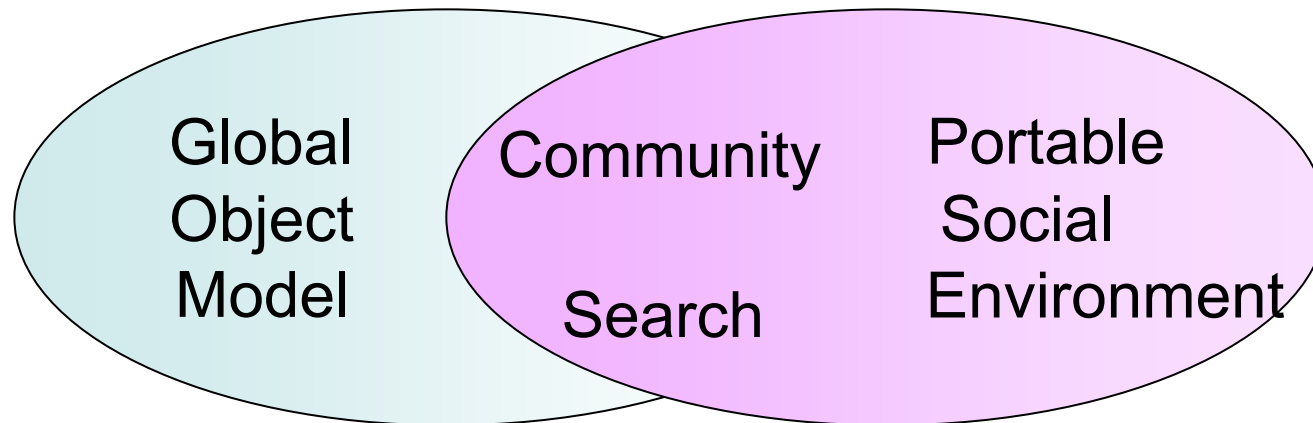    - Reviews: Around 10Mb/day
    - Ratings: <small>

Drove most advances in search from 1996-present

Increasingly rich and available, but not yet useful in search

This is in spite of the fact that interactions on the web are currently limited by the fact that each site is essentially a silo

# PeopleWeb: Site-Centric ⟹ People-Centric

**Global Object Model**

**Community**

**Search**

**Portable Social Environment**

- **Common web-wide id for objects (incl. users)**
  - Even common attributes? (e.g., *pixels* for camera objects)

- As users move across sites, their personas and social networks will be carried along

- **Increased semantics on the web through community activity (another path to the goals of the Semantic Web)**

**(Towards a PeopleWeb, Ramakrishnan & Tomkins, IEEE Computer, August 2007)**

# Social Search

- **Improve web search by**

  – Learning from shared community interactions, and leveraging community interactions to create and refine content

    - Enhance and amplify user interactions

  – Expanding search results to include sources of information (e.g., experts, sub-communities of shared interest)

Reputation, Quality, Trust, Privacy

Photos: **Explore Flickr** · **Learn More**

flickr BETA

# Tags / jaguar / clusters

jaguar   **SEARCH**

(Or, try an advanced search.)

**car**, **cars**, **auto**, etype, automobile, classic, vintage, autoshow, red, show

➡ **See more in this cluster...**

**zoo**, **animal**, **cat**, animals, bigcat, seattle, woodlandparkzoo, sleep, edinburgh, caged

➡ **See more in this cluster...**

**guitar**, **fender**

➡ **See more in this cluster...**

**aircraft**, **raf**

➡ **See more in this cluster...**

These are the *most recent* photos tagged with **jaguar**. See more...

# Web Search Results for "Lisa"



**Web** | Images | Video | Directory | Local | News | Shopping

**YAHOO! SEARCH**  [ Lisa ]  [ Search the Web ]  [ My Web (beta) ]

**My Web** BETA    My Search History OFF | On          Search Services    Adv

**Search Results**                    Results **1 – 10** of about **129,000,000** for <u>Lisa</u> -

Also try: <u>lisa lynnette clark</u>, <u>lisa loeb</u>, <u>lisa raye</u>, <u>mona lisa</u>  More...

Y! News Results for **Lisa**
**Lisa** Lynn Sargeson (Olsen) Marsolek - Independent Record - 4 minutes ago
UCLA's **Lisa** Willis Named Women's Basketball Pac-10 Player of the Week - Pac
By **LISA** MEYER TRIGG Editor - Banner Graphic - Nov 23 11:37 AM
Yahoo! Shortcut - About

*Latest news results for "Lisa". Mostly about people because Lisa is a popular name*

Y! My Web Results for **Lisa** (41)

*41 results from My Web!*

1. The Localization Industry Standards Association : home page
   Remember me. Quick Links. Welcome to **LISA**. Becoming a global enterprise is one of the most important challenges that your organization will ever face. There is no one right way to do it, but you should not have to reinvent the wheel. ... **LISA** is the leading international forum for organizations doing business globally ...
   Category: Software > Translation
   RSS: View as XML - Add to My Yahoo!
   www.**lisa**.org - More from this site - Save - Block

2. Laser Interferometer Space Antenna
   The Laser Interferometer Space Antenna is a mission that will detect and study gr
   sources involving massive black holes and galactic binaries. ... Download new LIS
   PDF file) **LISA** is a joint mission between the European Space Agency and NASA
   (Structure and Evolution ...
   **lisa**.jpl.nasa.gov - 19k - Cached - More from this site - Save - Block

*Web search results are very diversified, covering pages about organizations, projects, people, events, etc.*

3. ESA Science & Technology: **LISA**
   ... THE MISSION: **LISA** is an ESA-NASA mission involving three spacecraft flying approximately 5 million kilometres apart in ... Letter of Intent to Participate in **LISA** data processing study ...
   sci.esa.int/science-e/www/area/index.cfm?fareaid=27 - 31k - Cached - More from this site - Save - Block

4. a modern girl

# My Web 2.0 Search Results for "Lisa"

# Google Co-Op

## Query-based direct-display, programmed by Contributor

This query matches a pattern provided by Contributor…

…so SERP displays (query-specific) links programmed by Contributor.

Users "opts-in" by "subscribing" to them

Y! Research

# Challenges in Tag-Based Search

- How do we use these tags better?

- How do you cope with spam?

- What's the ratings and reputation system?


- The bigger challenge: where else can you exploit the power of the people?

- What are the incentive mechanisms?

    – Luis von Ahn (CMU): The ESP Game

Yahoo! Answers - Home - Mozilla Firefox

File   Edit   View   Go   Bookmarks   Yahoo!   Tools   Help

http://answers.yahoo.com/                                           Go   G

Y!   ▾   airbus 343              ▾   Search Web ▾   ⊕ ▾   ✉ Mail ▾   ➕ ▾   My My Yahoo!   Answers ▾   🧹 Games ▾   🎵 Music ▾   ❤ Personals ▾   Sign In ▾

Yahoo! My Yahoo! Mail

Search
the Web                    Search

# YAHOO! ANSWERS   [Sign In, My Account]

Home - Help - Forum
What's going on in Answers?

## ask.   ?

Ask a question on any topic and get answers from real people.

(you have **110** characters to work with)

**Post Question**

## answer.   ☺

Share what you know and you might make someone's day.
Featured Question

**Can I buy tile to match my 1950s-era kitchen countertop?**

## discover.   !

10 million answers and counting. Learn something new today.
Featured Topic

See what people are asking about in:

Travel

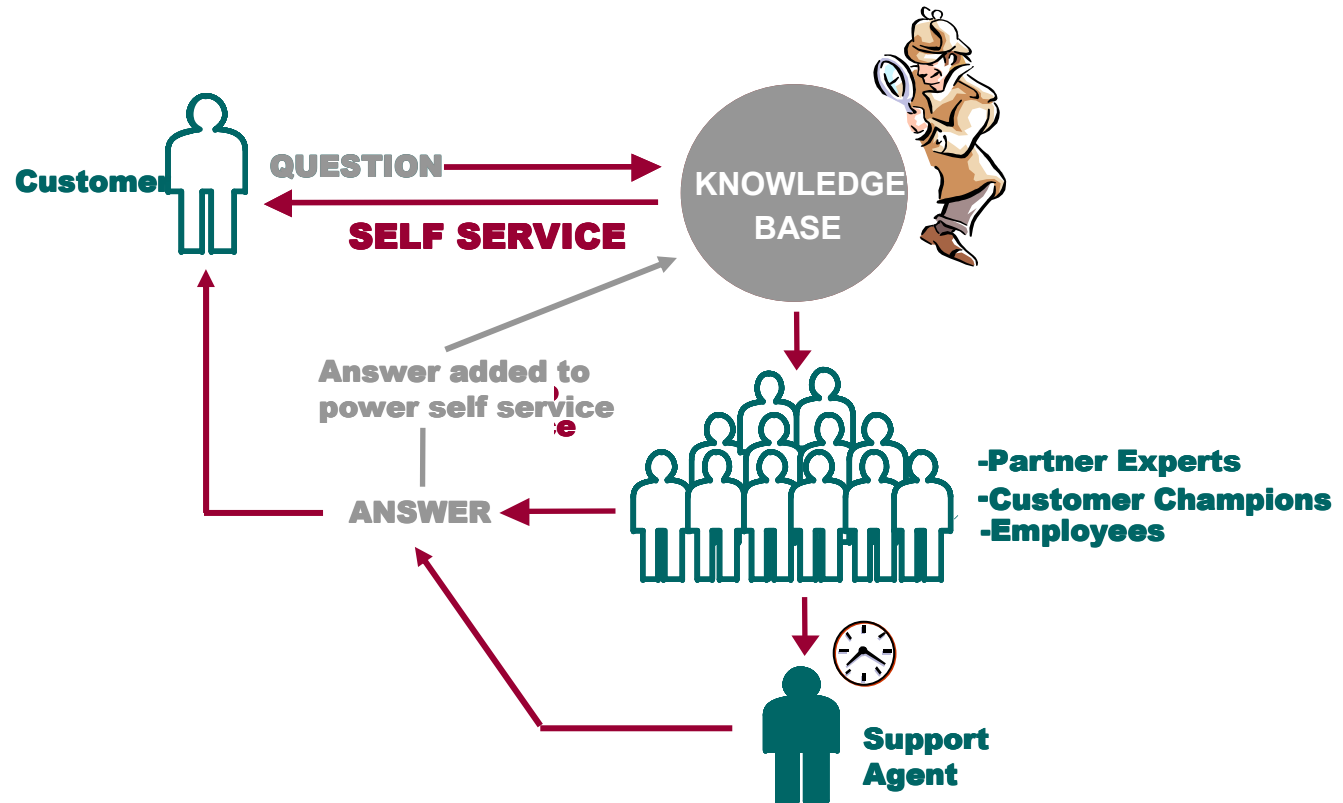**Search Yahoo! Answers:**               Search  Advanced                          **My Q&A**

10 million answers. Thanks to all the world's Answerers.

☺ Ready to Participate?
Get Started!

### Categories

→ Arts & Humanities
→ Business & Finance
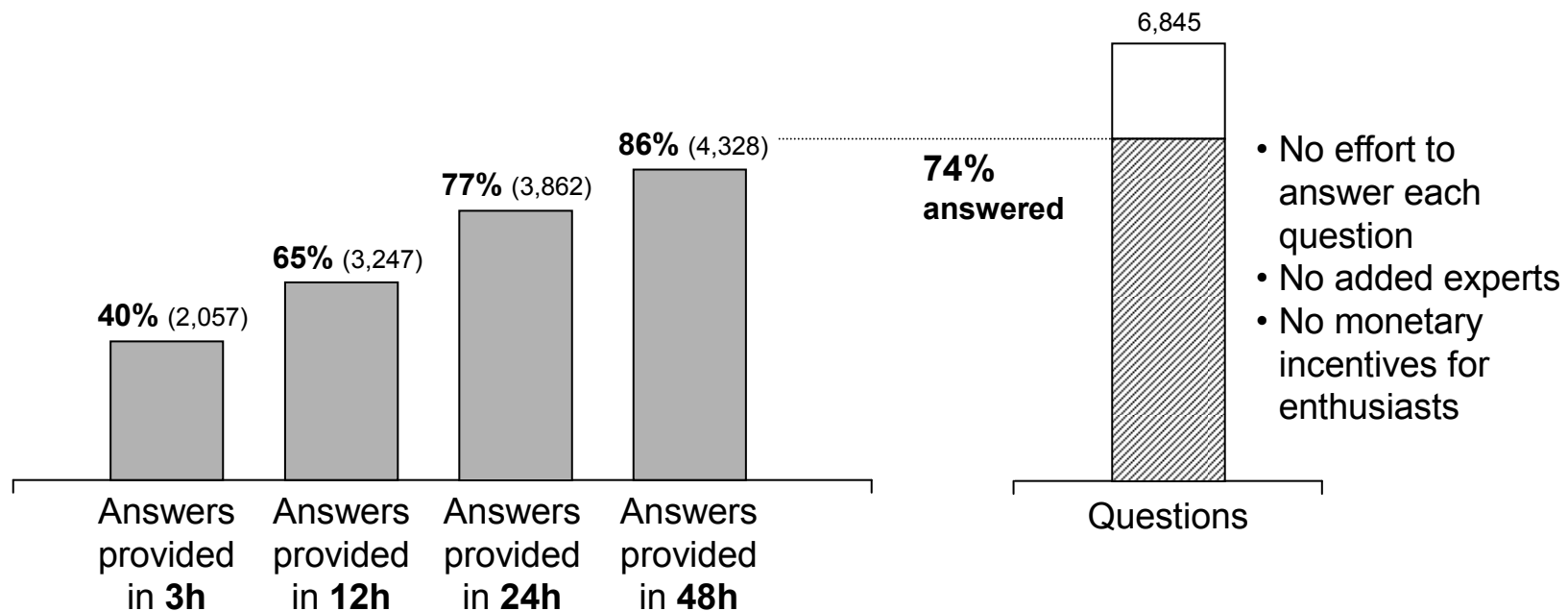→ Cars & Transportation
→ Computers & Internet
→ Consumer Electronics

**Share what you know. Answer open questions.**

**How do you get a bleach spot out of your pants?**
Asked by girly_antagonist - Cleaning & Laundry - 1 second ago

**is there a home remedy forgetting rid of ants out side with out harming my dogs or plants?**
Asked by cheetarajade - Other - Home & Garden - 2 seconds ago

which poker site is the most profitable for a tournament player?
Asked by judas - Card Games - 12 seconds ago

what is the website where you can play all those games? there is a new TV commercial about it....?

Done

🟧

start      🖥 🔵 🐱 🔘 💿 🎮 😄 🎵      MS...   Ya...   3 F.. ▾   2 M. ▾   ST...   C:\   Ad...    10:20 PM

# How It Works

# Timely Answers

▶ **77% of answers provided within 24h**

6,845

**86%** (4,328)

**77%** (3,862)

**65%** (3,247)

**40%** (2,057)

**74% answered**

- No effort to answer each question
- No added experts
- No monetary incentives for enthusiasts

Answers provided in **3h**  Answers provided in **12h**  Answers provided in **24h**  Answers provided in **48h**

Questions

# Interesting Problems

- Question categorization

- Detecting undesirable questions & answers

- Identifying "trolls"

- Ranking results in Answers search

- Finding related questions

- Estimating question & answer quality

(Byron Dom: SIGIR talk)

# Supplying Search Content

- As information discovery and search become task-oriented, we need to find ways to create semantically rich summaries that address the user's information needs.  Three ways to do this:

  – Editorial, Extraction, **UGC**

    • Opportunity to focus creation of structured UGC feeds directly into this growing need!

Challenge: Design social interactions that lead to creation and maintenance of high-quality structured content

# Better Search via Information Extraction

- Extract, then exploit, structured data from raw text:

For years, **Microsoft Corporation CEO Bill Gates** was against open source. But today he appears to have changed his mind. "We can be open source. We love the concept of shared source," said **Bill Veghte**, a **Microsoft VP**. "That's a super-important shift for us in terms of code access." **Richard Stallman**, **founder** of the **Free Software Foundation**, countered saying…

Select  Name
From   PEOPLE
Where Organization = 'Microsoft'

PEOPLE

| Name | Title | Organization |
|------|-------|--------------|
| Bill Gates | CEO | Microsoft |
| Bill Veghte | VP | Microsoft |
| Richard Stallman | Founder | Free Soft.. |

Bill Gates

Bill Veghte

(from Cohen's IE tutorial, 2003)

# DBLife

- Integrated information about a (focused) real-world community

- Collaboratively built and maintained by the community

- **Semantic web via extraction & community**

# The DBLife Portal

- Faculty: AnHai Doan & Raghu Ramakrishnan
- Students: P. DeRose, W. Shen, F. Chen, R. McCann, Y. Lee, M. Sayyadian
- Prototype system up and running since early 2005
- Plan to release a public version of the system in Spring 2007
- 1164 sources, crawled daily, 11000+ pages / day
- 160+ MB, 121400+ people mentions, 5600+ persons
- See DE overview article, CIDR 2007 demo

# DBLife



Crawled daily, 11000+ pages = 160+ MB / day

# Entity Resolution



co-authors = A. Doan, Divesh Srivastava,

Raghu Ramakrishnan

# Resulting ER Graph

# Challenges

- **Extraction**
  - Domain-level vs. site-level
  - Blending extraction with other sources (feeds, wiki-style user edits)
  - Compositional, customizable approach to extraction planning
    - Cannot afford to implement extraction afresh in each application!

- **Maintenance of extracted information**
  - Managing information Extraction
  - Incremental maintenance of "extracted views" at large scales
  - Mass Collaboration—community-based maintenance

- **Exploitation**
  - Search/query over extracted structures in a community
  - Search across communities—semantic web through the back door!
  - Detect interesting events and changes

Right in the sweet spot of the "Data to Knowledge" thrust of NSF's CDI initiative!

# Web Data Management: Massively Distributed Hosted Systems

# Two Key Subsystems

- **Serving system**
  - Takes queries and returns results

- **Content system**
  - Gathers input of various kinds (including crawling)
  - Generates the data sets used by serving system

- Both highly parallel

Goal: scaleup.
Hardware increments support larger loads.

Users ↔ **Serving System** ↔ Data sets

Logs

Content System ↔ Web sites

Data updates

Goal: speedup.
Hardware increments speed computations.

# A Case for Hosted Infrastructure

- What does it take to get the Next Great Thing off the ground?

- **Now:**
    - Set up multiple replicas of a clustered data store
    - Set up a system for indexing
    - Set up a system for caching
    - Set up auxiliary DBMS instances for reporting, etc.
    - Set up the feeds and messaging between them
    - Write the application logic
    - Fairly complex system at first line of new code

- **Our vision:**
    - Write the application logic
    - Use a hosted infrastructure to store and query your data

    - Or, as Joshua Shachter puts it: "The next cool thing shouldn't take a team of 30, it should be three guys, PHP and a long weekend"

# The PNUTS Project



| 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 |

CREATE TABLE Parts (
  ID VARCHAR,
  StockNumber INT,
  Status VARCHAR
  …
)

| D | 12352 | E |
|---|-------|---|
| E | 75656 | C |
| F | 15677 | E |

# Asynchronous replication

# Basic consistency model

- Record lifecycle
  1. Record inserted with a given primary key
  2. Record's non-primary key attributes updated
     - Primary key cannot be updated
  3. Record deleted
  4. Another record with the same primary key may subsequently be inserted

- What happens to a record with primary key "Brian"?

Record inserted  Update Update Delete

| v. 1 | v. 2 | v. 3 |

**Generation 1**

Record inserted  Update Update  Update  Delete

| v. 1 | v. 2 | v. 3 | v. 4 |

**Generation 2**

Record inserted  Delete

| v. 1 |

**Generation 3**

Time

**Research**

# The Big Picture: Sherpa Data Services

Applications

**YCA:** Authorization

**PNUTS Services**
- Query planning and execution
- Index maintenance

**Distributed infrastructure for tabular data**
- Data partitioning
- Update consistency
- Replication

**YDOT FS**
- Ordered tables

**YDHT FS**
- Hash tables

**YMB**
- Pub/sub messaging

**Zookeeper**
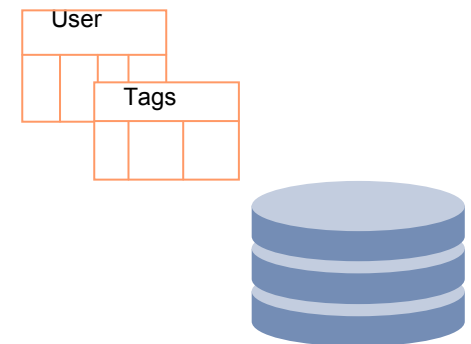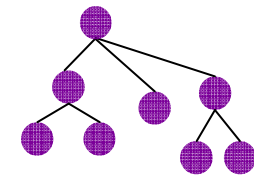- Consistency service

# Data Analysis Platforms

- Understanding online communities, and provisioning their data needs

  - Exploratory analysis over massive data sets

    - Challenges: Analyze shared, evolving social networks of users, content, and interactions to learn models of individual preferences and characteristics; community structure and dynamics; and to develop robust frameworks for evolution of authority and trust; extracting and exploiting structure from web content …

- **Examples:**
  - **Bigtable, Map-Reduce, Hadoop, PIG**

# The Bigger Picture

- **Software-as-a-service**

  – E.g., Salesforce.com

- **Hosted data systems**

  – Amazon's S3 and EC2

- **Web application development**

  – Ning, Ruby-on-rails

- **Change tracking**

  – Stream management

**Research**

# Grand Challenge

- **How to maintain and leverage structured, integrated views of web content**
  - Web meets DB … and neither is ready!
    - Interpreting and integrating information
      - Result pages that combine information from many sites
    - **Scalable serving of data/relationships**
      - **Multi-tenancy, QoS, auto-admin, performance**
  - Beyond search—web as app-delivery channel
    - Data-driven services, not DBMS software
      - Customizable hosted apps!
    - Desktop ➡ Web-top