# Preservation Explorer and Vault for Web 2.0 User-Generated Content

http://han.ddns.comp.nus.edu.sg/prev/

Anqi Cui[1], Liner Yang[1], Dejun Hou[2],
Min-Yen Kan[2], Yiqun Liu[1], Min Zhang[1], Shaoping Ma[1]
[1]Tsinghua University, National University of Singapore[2]

# PrEV and NExT

- *Preserve* the past and today's Web 2.0 User-Generated Content (UGC) as a *Vault*, to help future researchers –
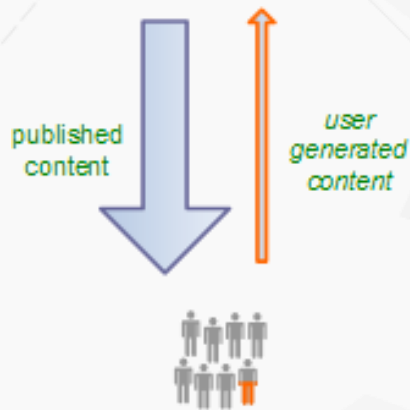- – *Explore* what was happening in our lives

# Motivation: The Web



Web 1.0
"the mostly read-only Web"
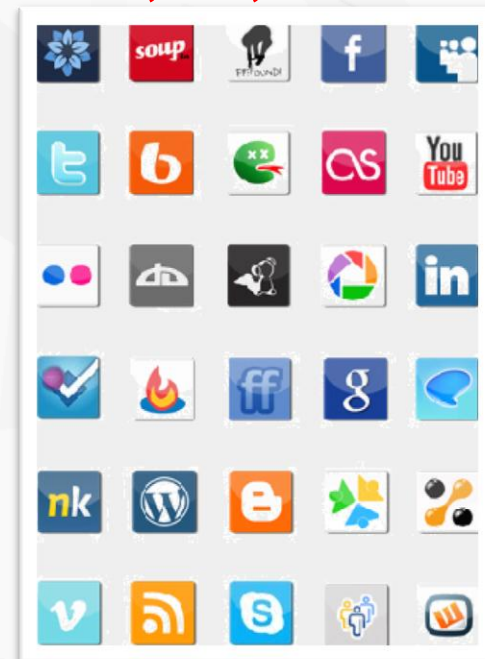250,000 sites

Web 2.0
"the wildly read-write Web"
80,000,000 sites

collective intelligence

> 600,000,000 sites

published content

user generated content

published content

user generated content

45 million global users
1996

1 billion+ global users
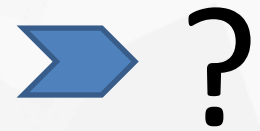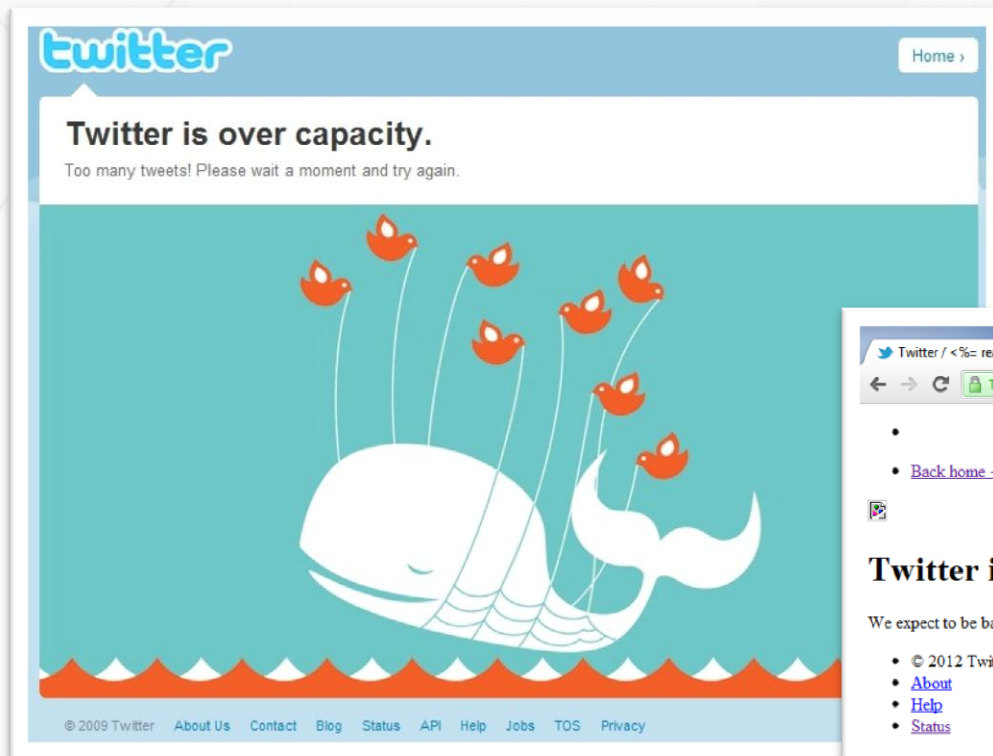2006

2.2 billion global users
2012

# Reliable Services?
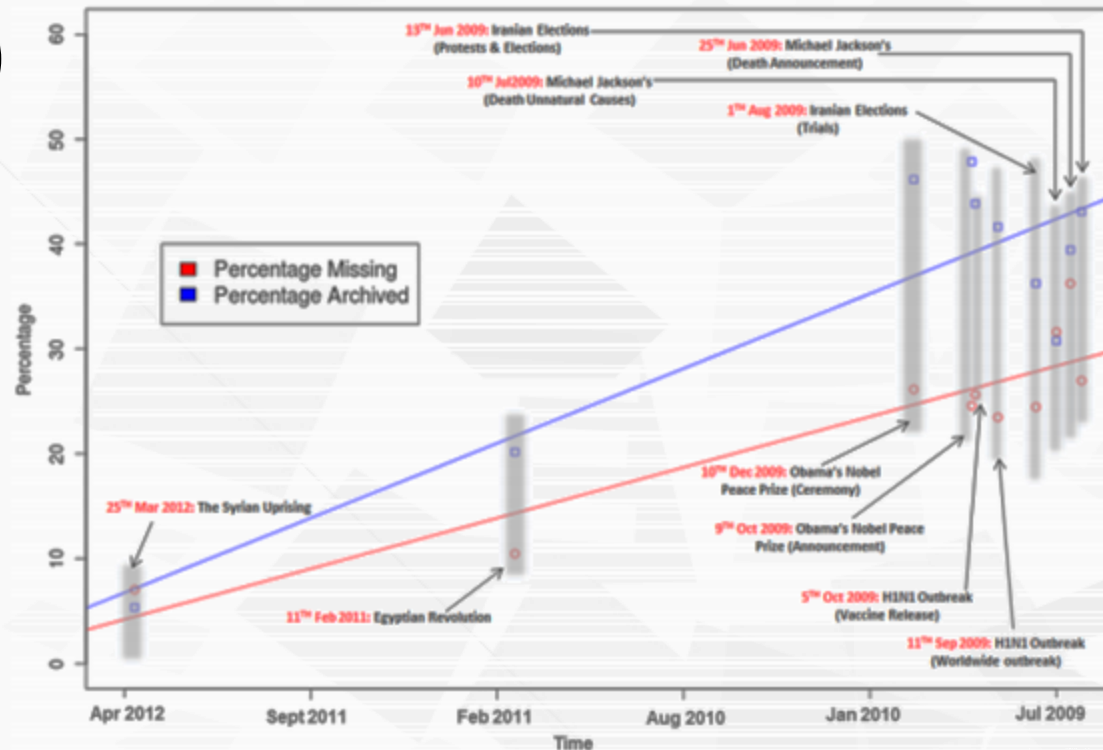
- The famous Twitter's Fail Whale:



Below: The Fail Whale failed on July 27th 4am UTC

# The Vanishing Web Contents

- Almost 30% of recorded history shared over social media has disappeared. (SalahEldeen 2012)

# Preservation: Web 1.0 => 2.0


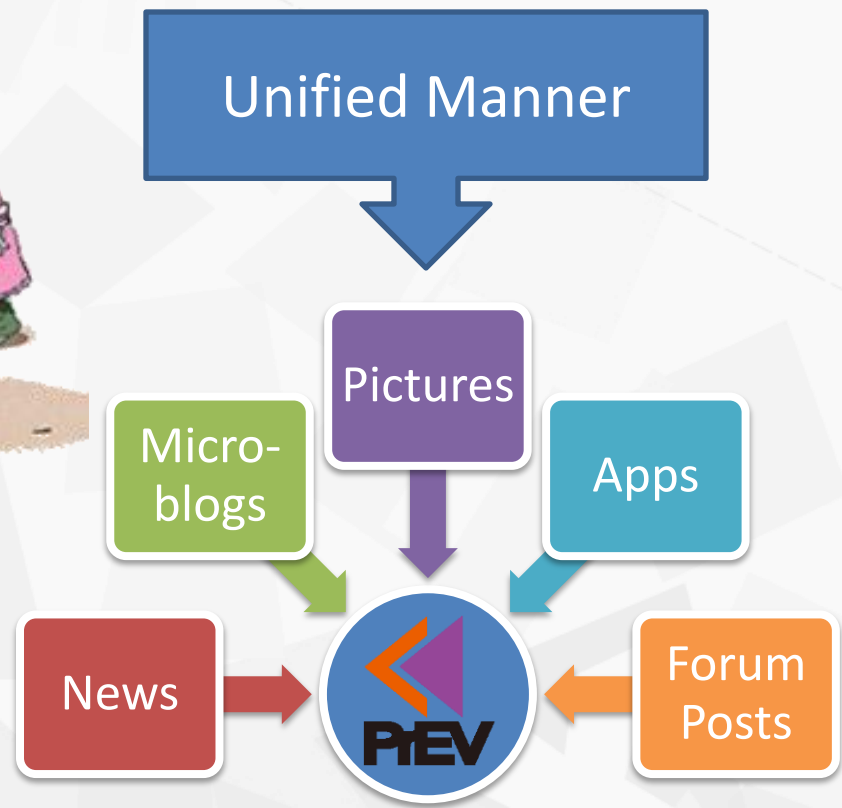(Kahle 1997)


(Albertsen 2003)


(Yan 2004)

Above: Internet Archive
Right: Country-wide:
  Norway's Paradigma
  China's Infomall

# Piecemeal => A Unified Manner



Piecemeal

(Campbell 2009, Hockx-Yu 2011)

Unified Manner

Pictures

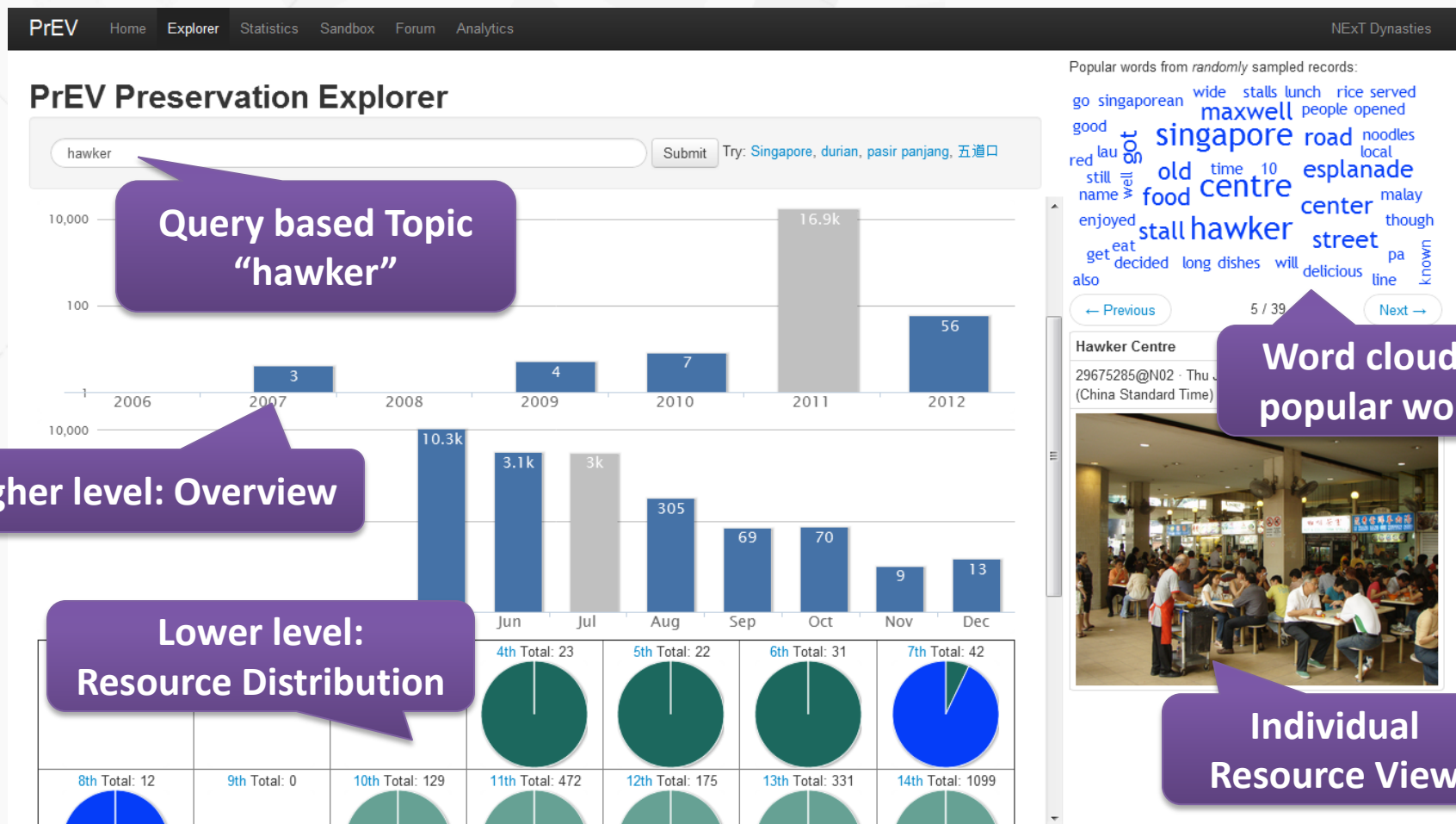Micro-blogs

Apps

News

PiEV

Forum Posts

# Scenario 1: Summarizing the Data

- Ryan's course project in 2022: Singaporean hawker center (food court) history
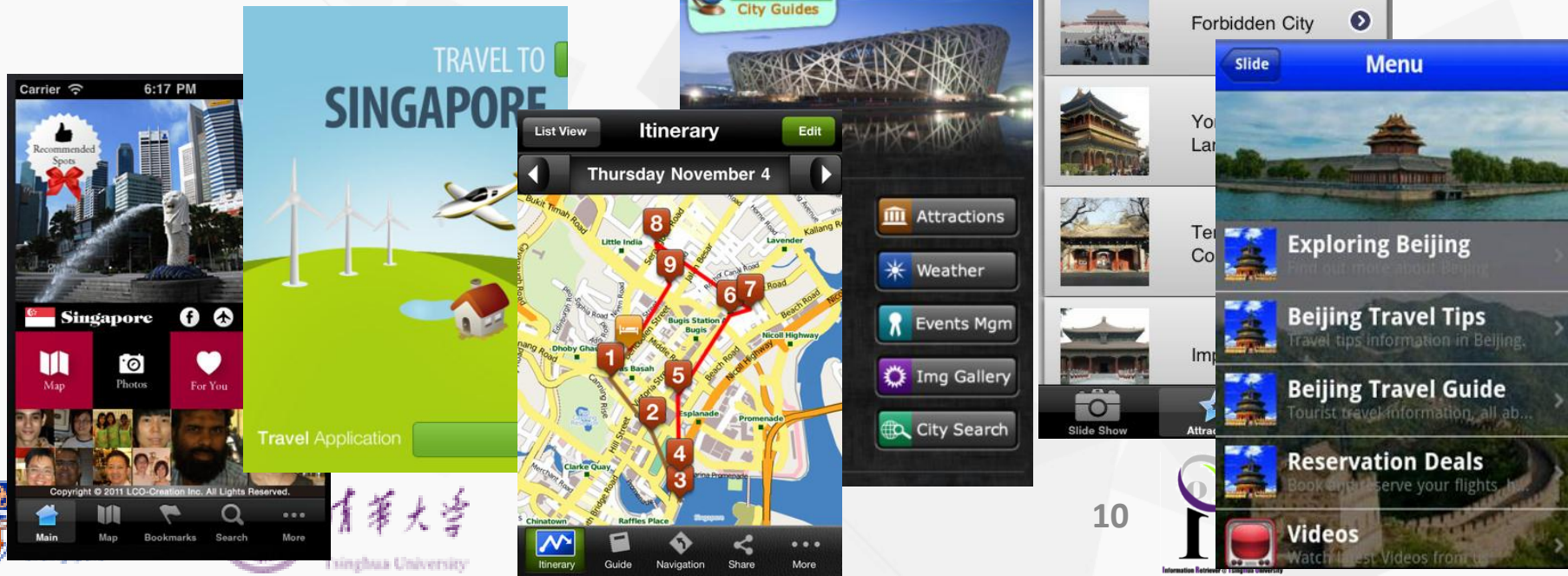
# Scenario 1: Query-based View

# Scenario 2: Historical Data Archive

- Blueberry Inc.: Develops "Follow Me" iPhone travel guide app for BJ & SG

- Review existing competitor apps' descriptions and reviews

# Scenario 2: API for Large-scale Access

```
{ "total":86, "count":10, "totalpage":9, "page":0, "data":[
  { "crawlresource":"twitter", "encoding":"en", "tweetcreatedat":"Sun Oct 02 19:56:51 +0800 2011",
    "url":"http://twitter.com/#!/ChristianLeeVO/status/120467276104339457",
    "maincontent":"Check this video I shot of Pulau Ubin for our Travel Now Singapore webseries and iphone
app http://t.co/mvfgJDqP via @youtube" },
  { "crawlresource":"weibo", "encoding":"zh", "weibocreatedat":"Sat Jan 14 19:13:27 +0800 2012",
    "url":"http://www.weibo.com/1910529591/y0Lf5mFAk",
    "maincontent":"#App推荐# 旅程规划: Routes. Planning your journeys【出行必备】，iPad/ iPhone通用。这款应用
可让你规划旅游景点，像是到某景点去拜访或是去某家大卖场购物等等。它会算出需要多远的距离以及所需的时间，就像...
http://t.cn/z0gtfEr （分享自 @App每日推送)" },
  { "crawlresource":"sgbjapps", "encoding": "others", "crawltime":"Wed Dec 23 00:00:00 +0800 2009",
    "maincontent":"Do You Love Travel ? If Yes，You Should Not Miss This App. Updated For Now! Download
this app to your iPhone to enjoy these beautiful scenery anywhere you go! These pictures are HD Photo You
can download the image to your iPhone or iPod and make it to wallpaper. No Ad No Wifi!",
    "name":"A Tourist Paradise <Singapore>" }, ... ] }
```

1
2
3

1 

2 

3 

# System Architecture

- Three layers, loosely coupled
  - Preservation
  - Indexing
  - Interface

# 1. Preservation Layer

- Incoming Data Detection

- Data Format Recognition

- Record Storage

- Backup

# 1. Collected Resources (as of May 2012)

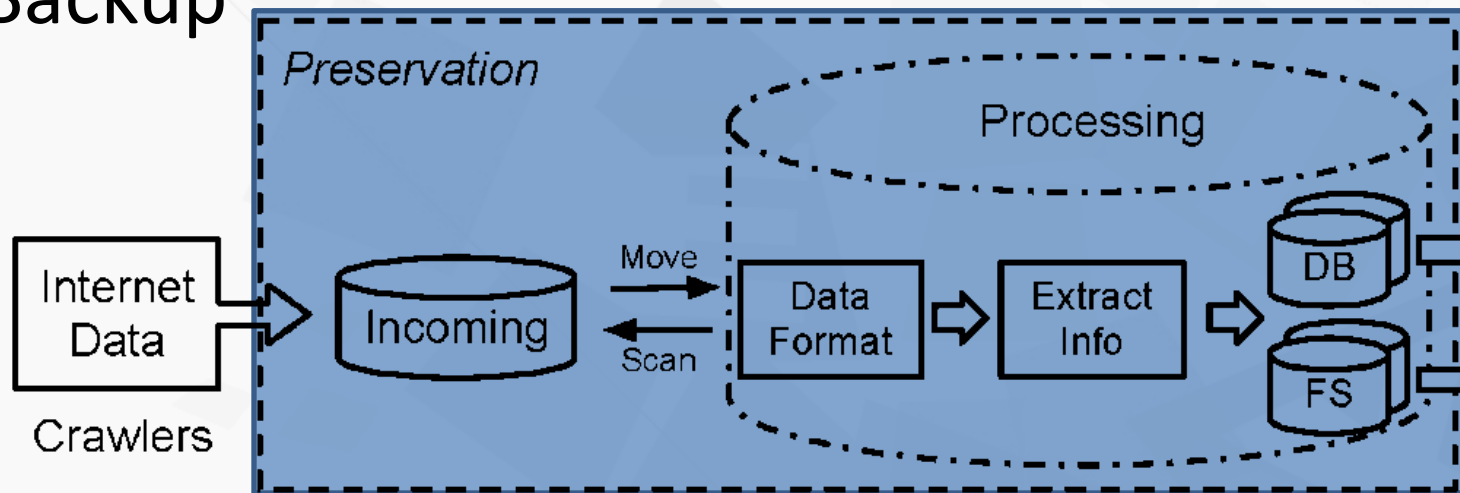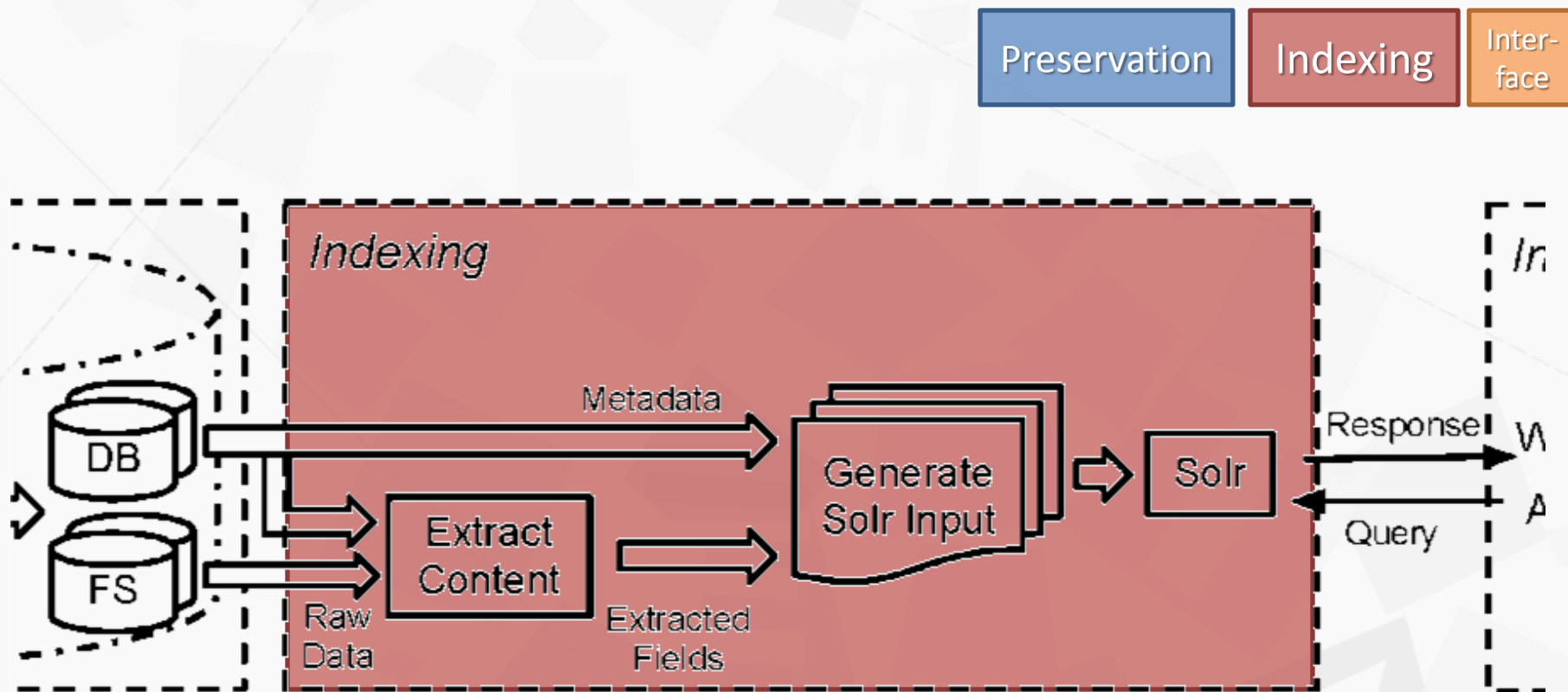| Data Type | Resource | No. of Records | |
|---|---|---|---|
| Microblog messages | Twitter | 229 M | 497M |
| | Weibo | 139 M | 191M |
| Photos with texts | Flickr, Panoramio | 2 M | |
| Food forum restaurants | 27 Singapore sites | 6 M (pages) | |
| | Fantong, Dianping (Chinese) | 78 K | |
| Public forum posts | 4 Chinese forums | 1 M (approx.) | |
| Product review products | 7 e-commerce sites | 70 K | |
| News articles | Sina News | 224 K | |
| | Guardian, Channel NewsAsia, Skysports, CNN, Economist, FoxNews, NewYorkTimes, StraitsTimes | 59 K | 3.3M |
| Wiki articles | Hudong (Chinese) | 1 M | |
| Traffic records | Singapore | 24 K | |
| | Beijing | 19 K | 30K |
| Question Answering articles | Baidu Zhidao (Chinese) | 33 K | |
| | Yahoo! Answers, WikiAnswers | 52 K | |
| Mobile Apps | US App Store | 617 K | |
| | Android Market | 345 K | |
| | Blackberry, Windows | 162 K | |

*Updates are of Sep 2012. Please refer to the website for an up-to-date statistics.

# 2. Indexing Layer

Preservation   Indexing   Inter-face

# 3. Interface Layer

Preservation | Indexing | Inter-face
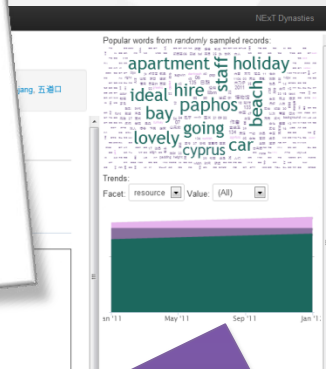
- Web frontend: Calendar view, word cloud, trends view, individual view

- API Frontend: Authentication, rate limiting (sandbox provided)



**Trends of resources of the query "paphos"**
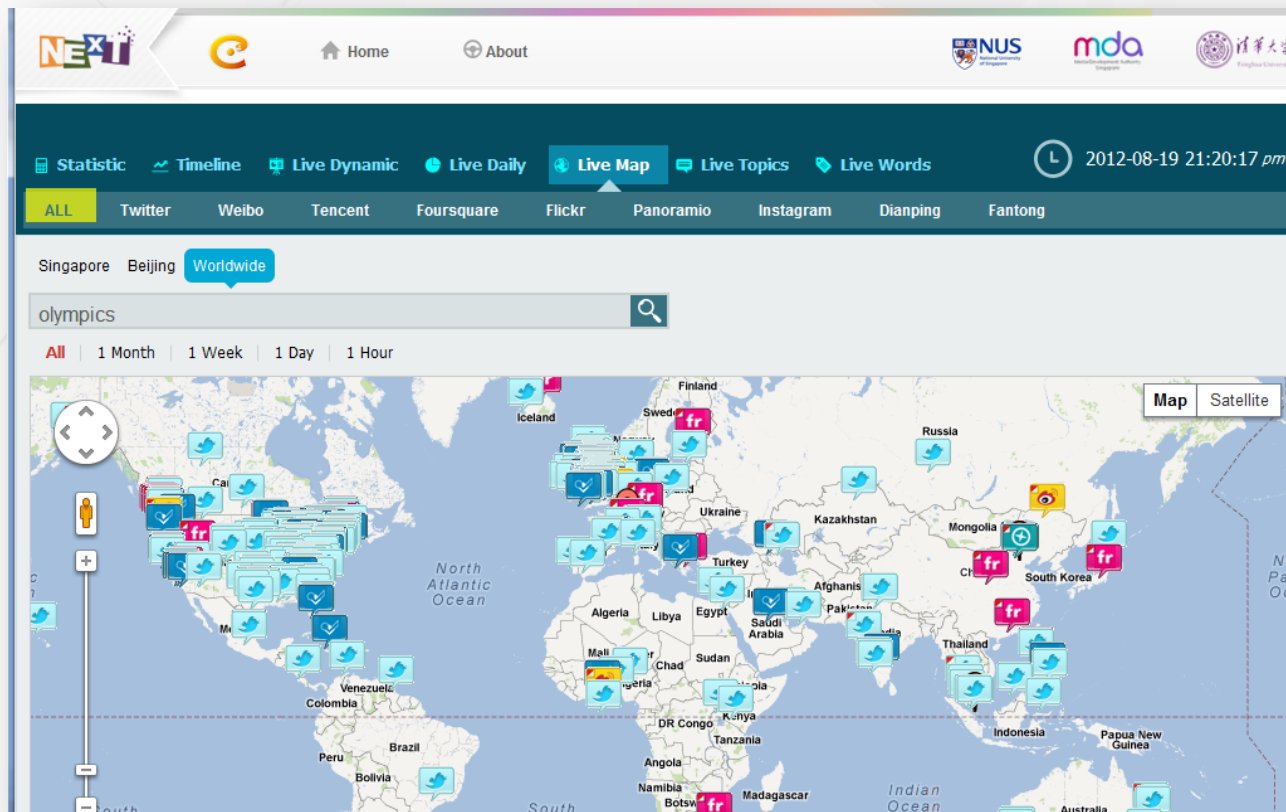
{ "total":86, "c...
    { "crawlresource ...                                     800 2011",
      "url":"http://twitter.com/w...
      "maincontent":"Check this video I shot o...              bseries and iphone
app http://t.co/mvfgJDqP via @youtube" },
    { "crawlresource":"weibo", "encoding":"zh", "weibocreatedat":"Sat Jan ...    0 2012",
      "url":"http://www.weibo.com/1910529591/y0Lf5mFAk",
      "maincontent":"#App推荐# 旅程规划: Routes. Planning your journeys【出行必备】，iPad/ iPhone通用。这款应用
可让你规划旅游景点，像是到某景点去玩功或是去某东大支场购物等等，它会算出需要多远的距离以及所需花时间，就像...
http://t.cn/z0gtfEr （分享自 @App每日推送）" },
    { "crawlresource":"sgbjapps", "encoding": "others", "crawltime":"Wed Dec 23 00:00:00 +0800 2009",
      "maincontent":"Do You Love Travel ? If Yes, You Should Not Miss This App. Updated For Now! Download
this app to your iPhone to enjoy these beautiful scenery anywhere you go! These pictures are HD Photo You
can download the image to your iPhone or iPod and make it to wallpaper. No Ad No Wifi!",
      "name":"A Tourist Paradise <Singapore>" }, ... ] }

# NExT Live UGC Web Portal

- http://137.132.145.151:8080/ugcp/

# Conclusion & Future Work

- What is PrEV?
  - PrEV: City-centric archiving system
  - Archive & unify multilingual Web 2.0 data
- Whom is it for?
  - Individual users: Discover the old good days
  - Enterprise-level use: Programmatically access a large amount of data for business and scientific research
- How is PrEV built?
  - Preservation layer: Collect data from different groups
  - Indexing layer: Faceted search
  - Interface layer: Flexible for different needs

- Future: System performance, user interfaces (visualization)
  - StrmWrd: The visualization tool https://github.com/THUNUS/StrmWrd

# Questions?

- Website:

  http://han.ddns.comp.nus.edu.sg/prev/ or
  http://tinyurl.com/prevweb
  (mobile devices supported)

- Contact:
  - Anqi Cui (Google+), @CAQ9 (Twitter)

# About NExT

- Crawling & mining UGCs in SG and BJ in:
  - Location-oriented: shared photos and check-in venues;
  - Topic-oriented: forums, question-answering, tweets;
  - Application-oriented: mobile applications and associated information and discussions; and
  - structured: factual, cultural and historical information.
- Carrying out research projects in the areas of: Extreme Database, Live Event Capturing and Sharing, Live Media Processing, **Live Text Search**, Live City
- http://next.comp.nus.edu.sg/

# Extreme Text Search Group

- Extreme search in text: real-time search + faceted search
  - Interesting research directions
  - Commercial and industrial applications for the next generation web
- Topics to be explored includes:
  - Twitter sentiment analysis
  - Mobile app ranking
  - Social and differential news analysis