# Multimedia Web Analysis Framework towards Development of Social Analysis Software

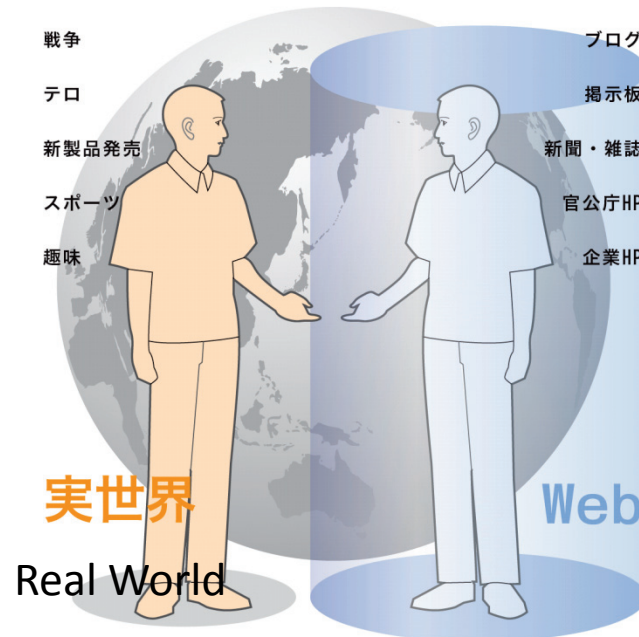**Masashi Toyoda**

IIS, University of Tokyo

**Shin'ichi Satoh**

National Institute of Informatics

# Web as a projection of the real world

- Web has been reflecting various social events in the real and virtual world
- Web could be used as social sensor

War
Terro
Sports
Hobbies
Books
Movies

戦争
テロ
新製品発売
スポーツ
趣味

実世界
Real World

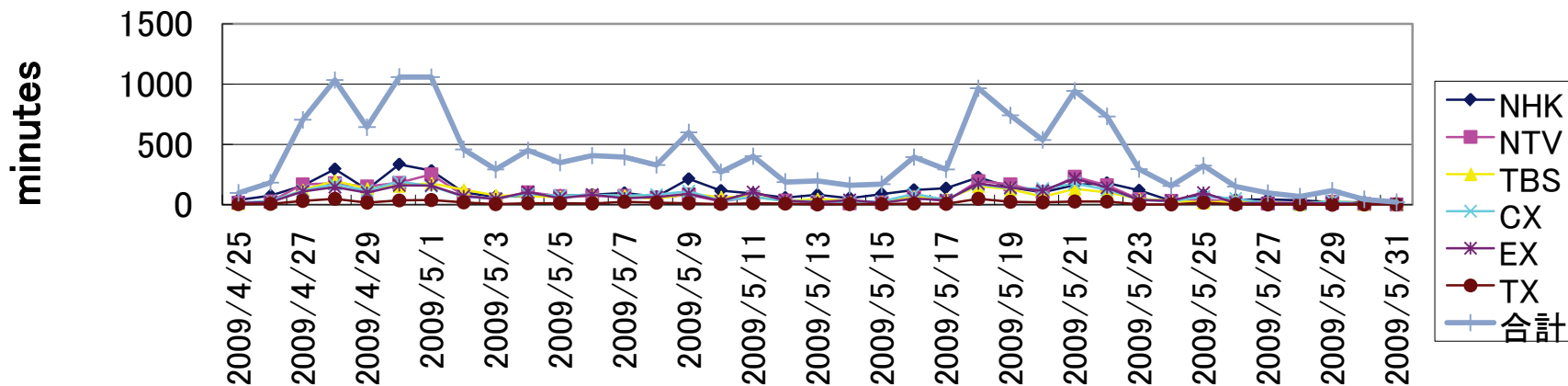ブログ
掲示板
新聞・雑誌
官公庁HP
企業HP

Web

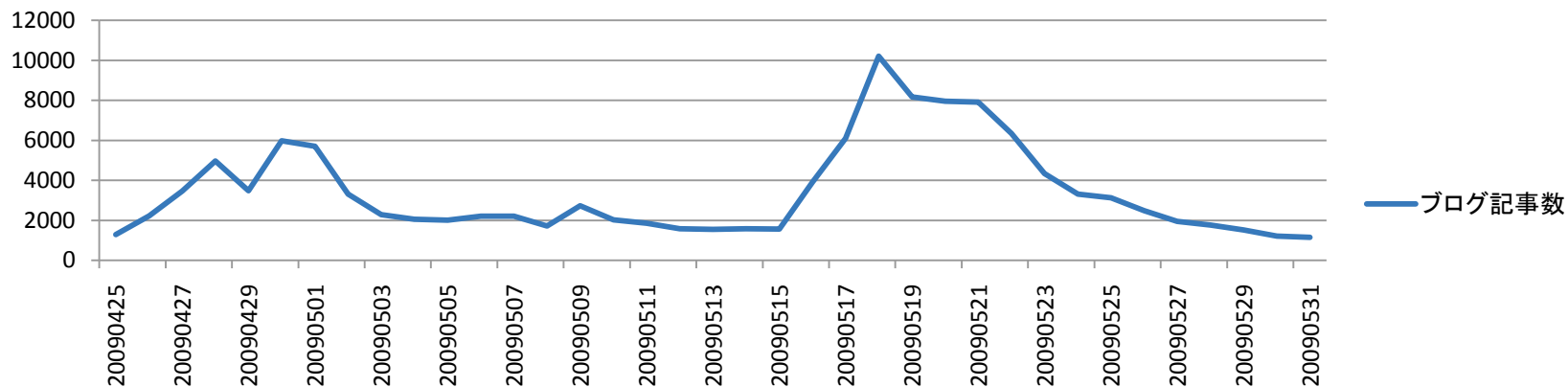Online news/magazines
Blogs
SNS
BBS
Youtube

# Correlation between Web and TV

- In many case, trends in Web and TV are similar
  - Example: broadcast time of TV news and #blog posts on swine-flu

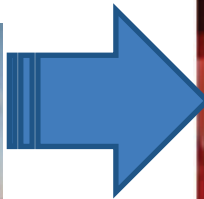Broadcast time of TV news on swine-flu（Meta-TV）



#blog posts on swine-flu

# Cross-media evolution of topics

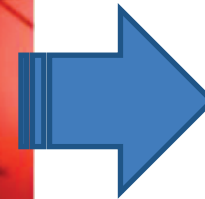- The first photo of an accident uploaded to twitpic, and used in TV News
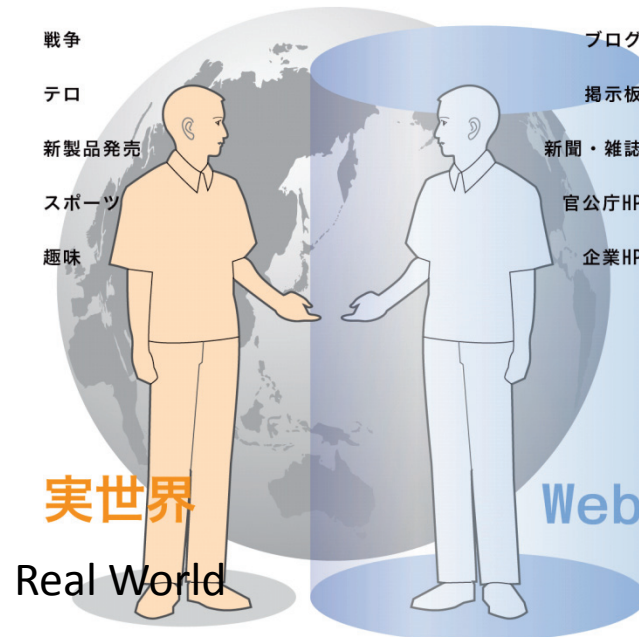


MSNBC   News (TV)

Topic evolution from a photo appeared in twitter to TV and news broadcasting

# Challenge: Tracking cross-media events

Web and media (TV, news papers, ...) dynamically interact reflecting the real world events

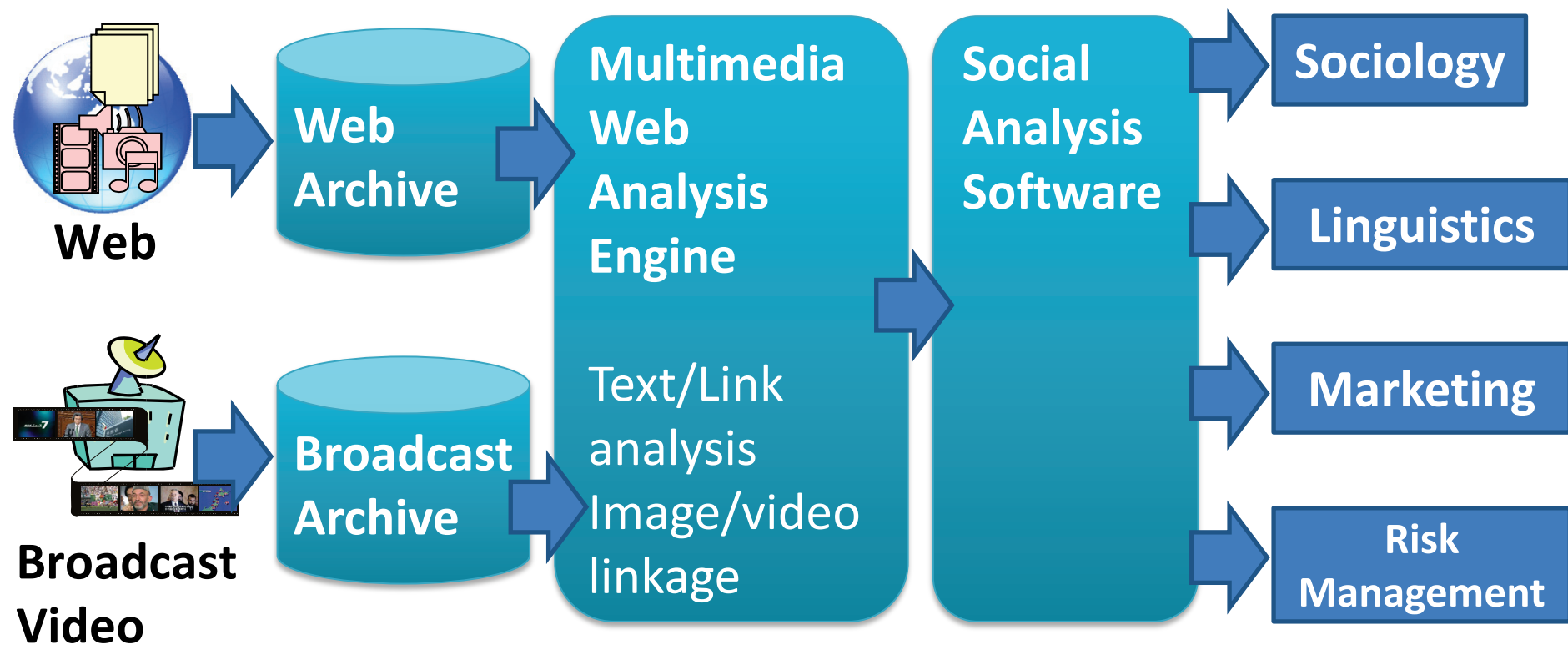Can we track the evolution of topics over multiple media?

War
Terrorism
Sports
Hobbies

実世界
Real World

戦争
テロ
新製品発売
スポーツ
趣味

ブログ
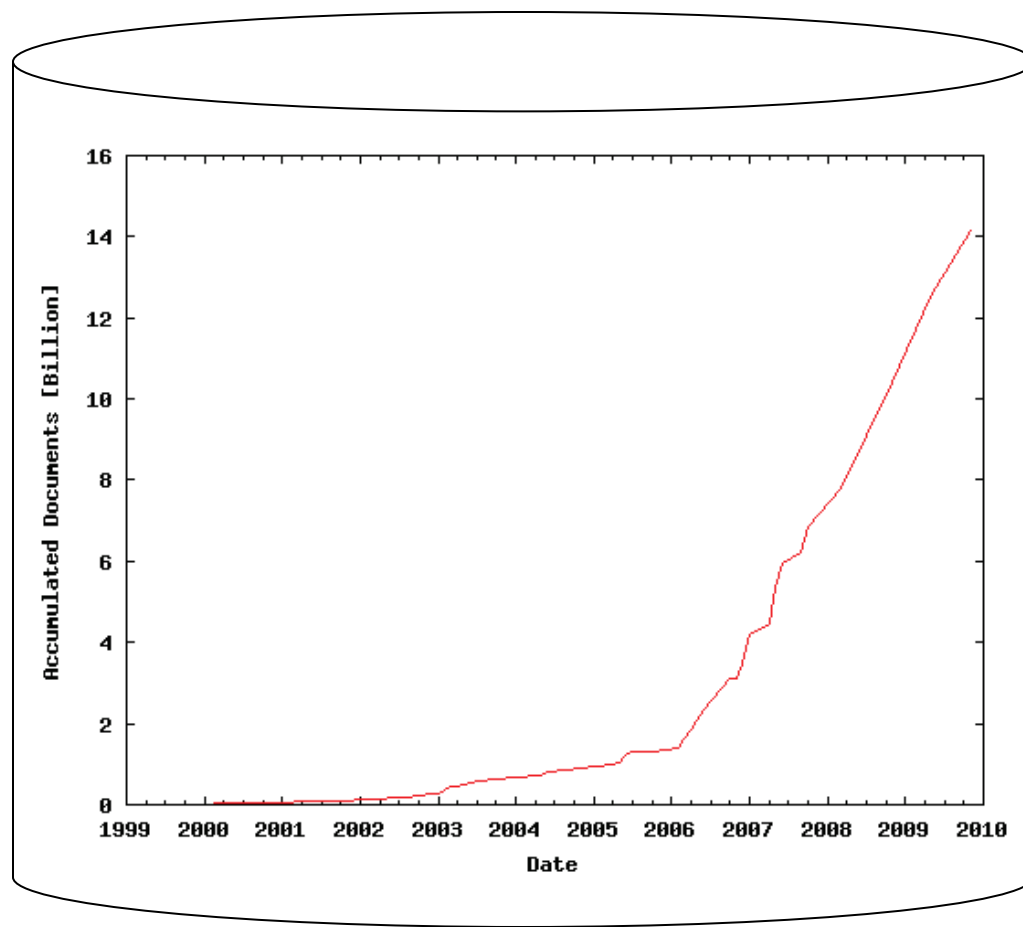掲示板
新聞・雑誌
官公庁HP
企業HP

Web

Online news/magazines
Blogs
SNS
BBS
Youtube

# Multimedia Web Analysis Framework for Social Analysis
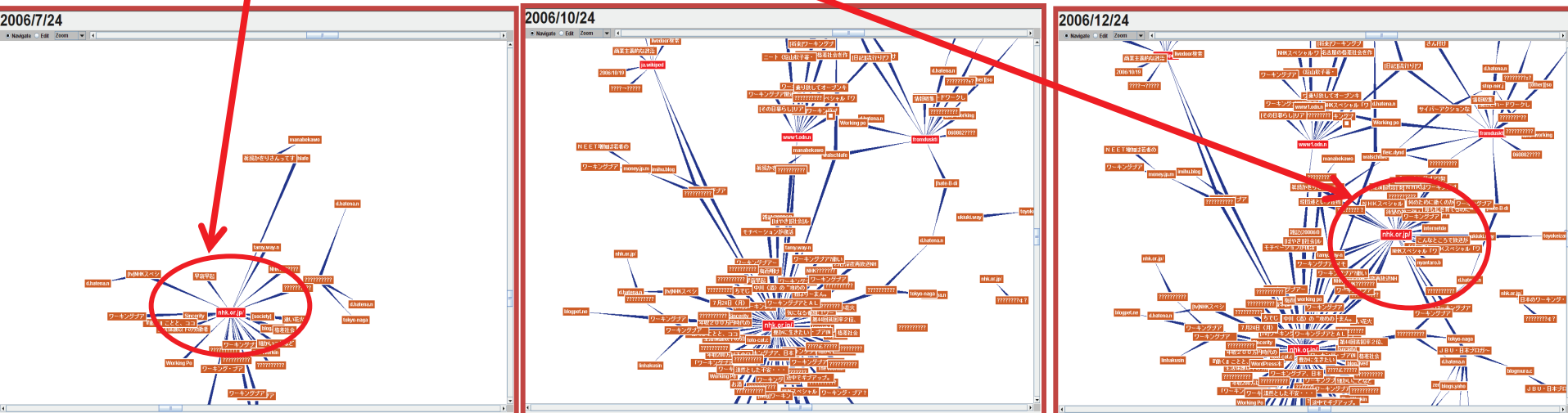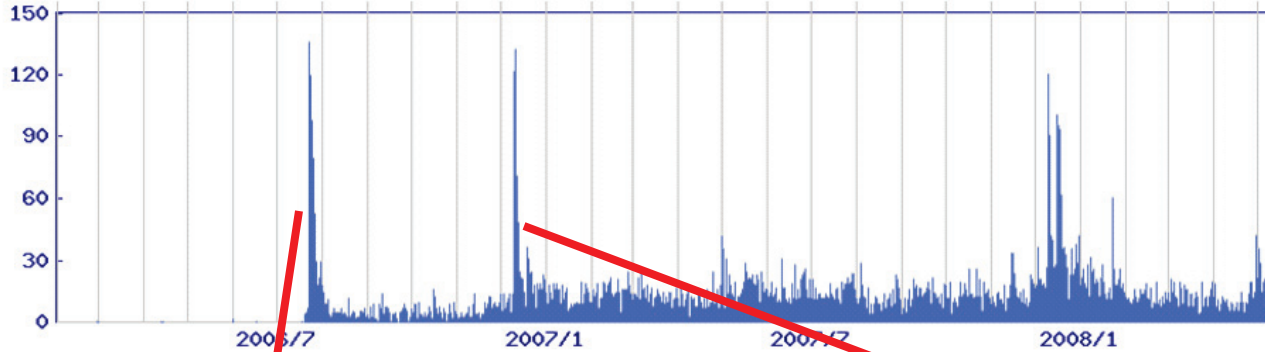
# 10 Year Web Archive in Univ. of Tokyo

- One of the biggest Web archive in Asia
  - Focused on Japanese pages in any domain (1999〜)
- Contents of **14 billion URLs** are stored in total including HTML texts and images (Dec. 2009)

# Visual Analysis of Evolving Web

- The first page on a topic

- The most influential page

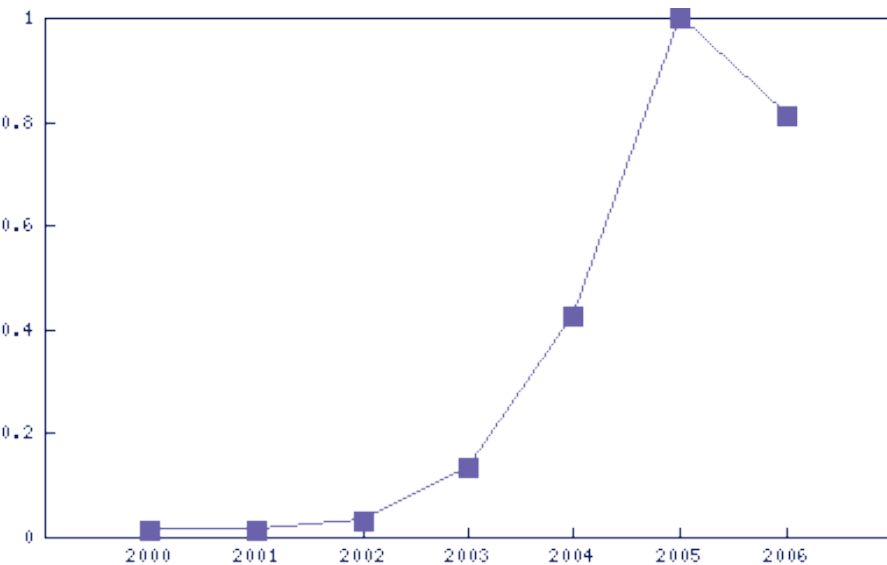Example: diffusion of the "working poor" problem in blogosphere

# New Word Extraction for linguistics
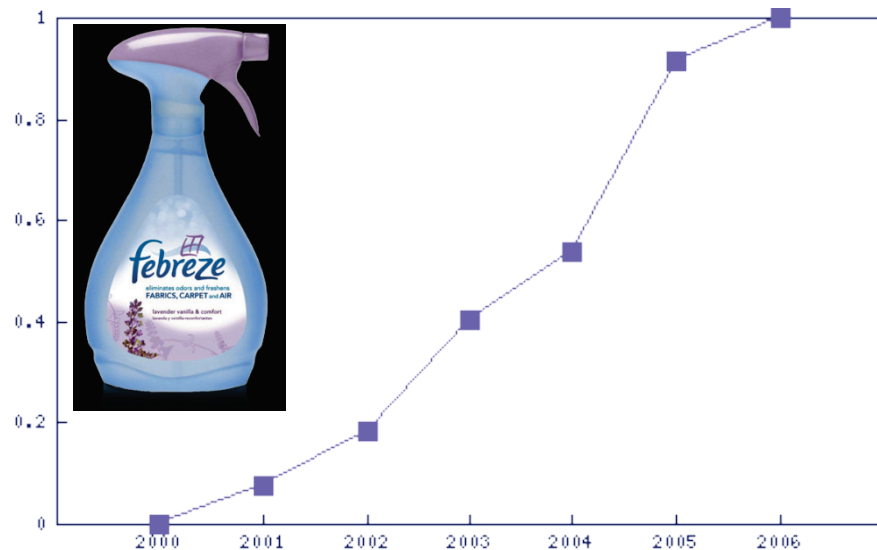
- Words not in dictionary but used as noun, verb,…
  - Learning the context of known words, and predict Part-of-Speech tags of unknown words
  - Filtering by the growth ratio of occurrences
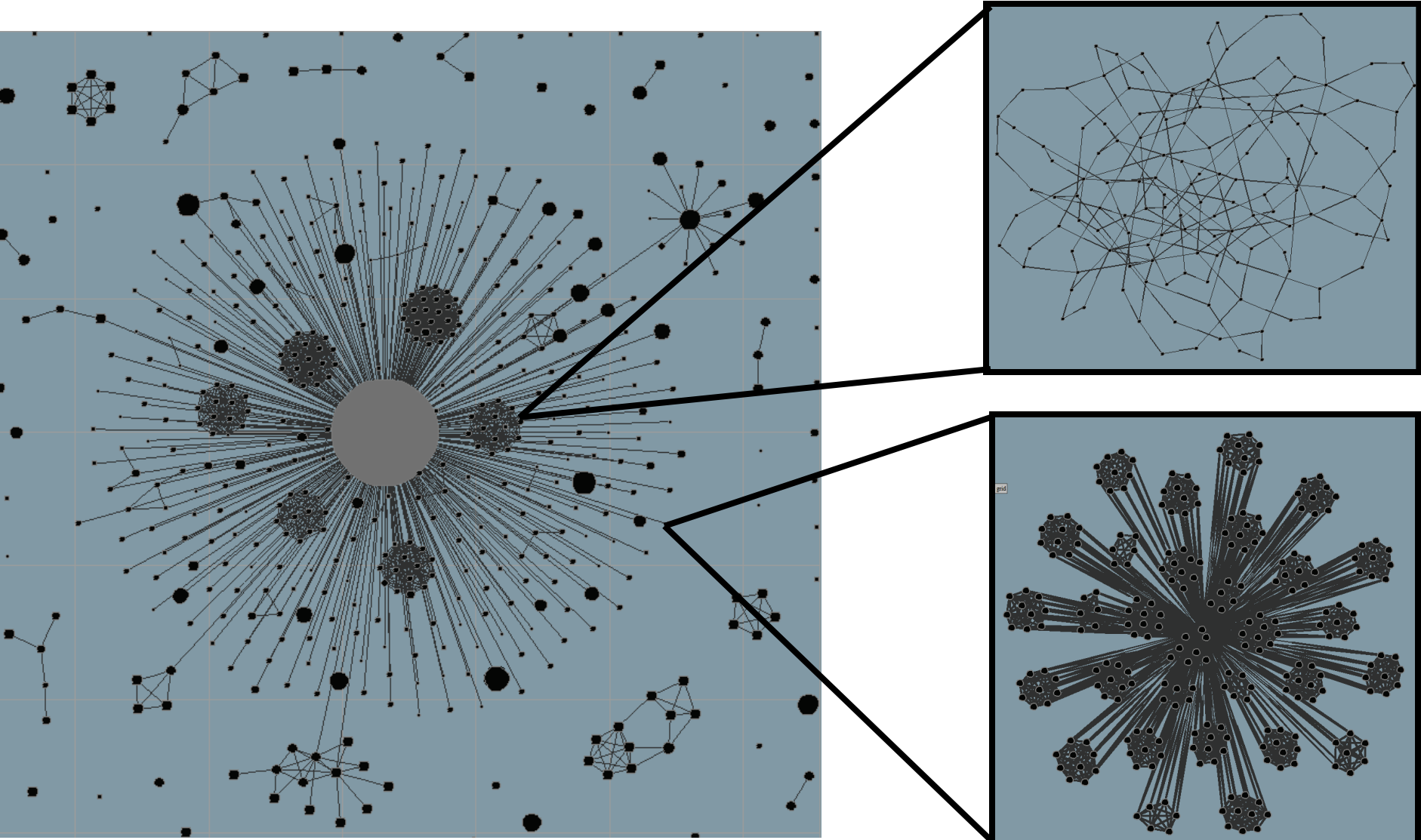
GUGU-RU (Verb: Search by Google)

FABU-RU (Verb: Use Febreze, a fabric & air refresher by P&G)

# Search Engine SPAM [Chung. 09, 10]

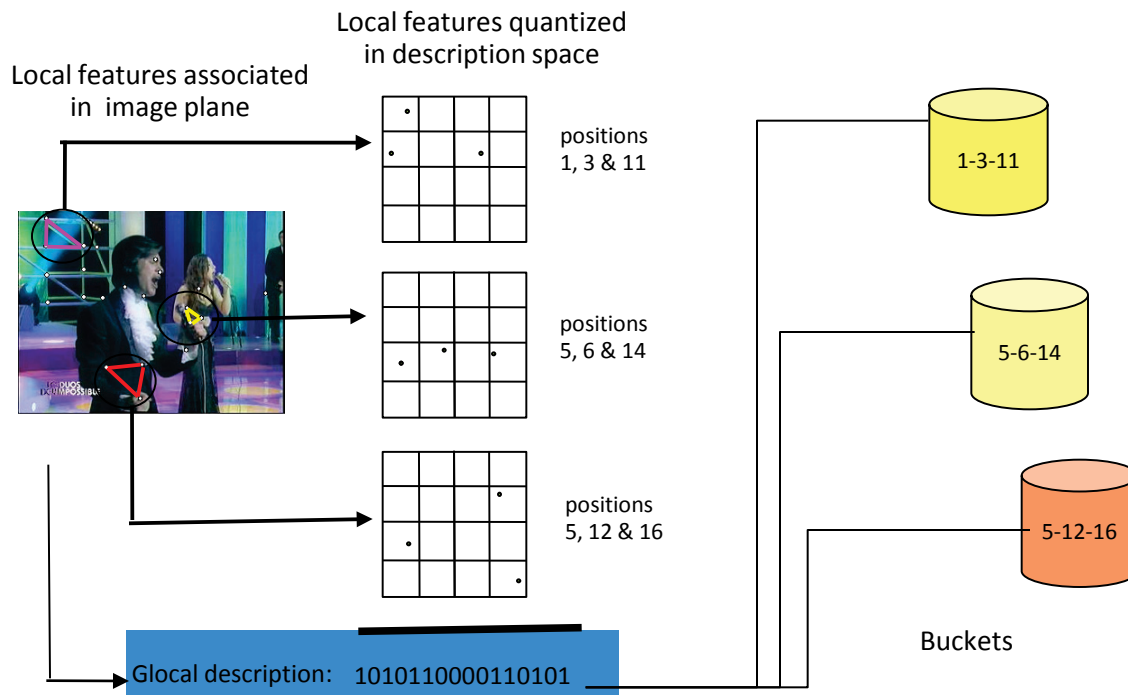- Many web sites intend to deceive SEs and boost their ranking

# Broadcast Archive in NII

- Capture and store Japanese TV broadcast
  - 7 channels
  - About 1 year (60,000 hours)

# Content-based image search in the video archives [Poullot et. al. 09]

1. Construct a database of the video archives (offline)
   - Extract local description in KeyFrames
   - Associate them in image place → define buckets in an index (spatial hashing)
   - Built a global description using the local ones and insert it in this buckets



Local features associated in image plane

Local features quantized in description space

positions 1, 3 & 11

positions 5, 6 & 14

positions 5, 12 & 16

1-3-11

5-6-14

5-12-16

Buckets

Glocal description:   1010110000110101

2. Search (online)
   - Same extraction process on query picture
   - Explore matching buckets
   - Rank result by similarity of global description

# Summary

- Tracking cross-media events is an important issue for social analysis

- Enabler: Large-scale Web and broadcast archive

- Analysis Engine:
  - Text/link based tracking techniques
  - Scalable video indexing technique

- Next step: Integration of Web/image analysis
  - Text/link/image based topic tracking
  - CF filter, face recognizer