# The Metadata Triumvirate

Social Annotations, Anchor Texts and Search Queries

Michael G. Noll, Christoph Meinel

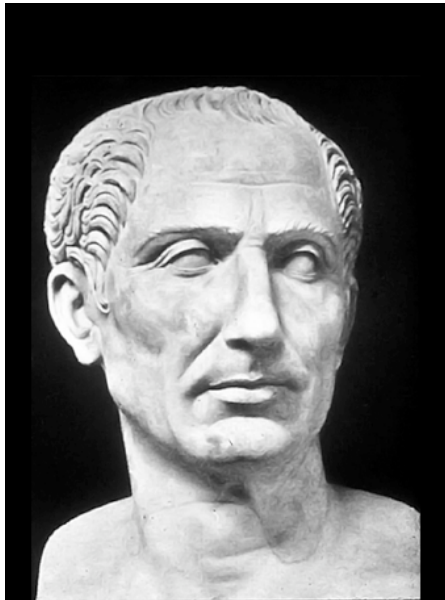# In this talk

**Metadata – "data about data"**

- Here: Web metadata, i.e. data about **WWW documents**

- **Variety of uses** for such metadata in Web information retrieval: indexing, ranking, filtering, …

- **Different types** of Web metadata:
  In this talk, we study and compare 3 very popular ones with the goal to improve our understanding of these metadata types, thereby helping us to improve existing IR algorithms or come up with new ones.
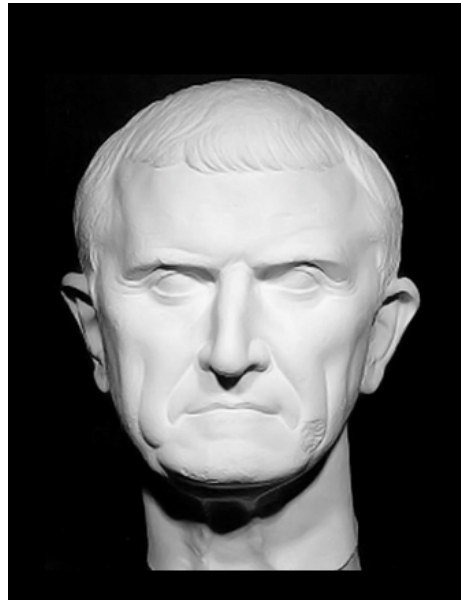
# The history of Triumvirates
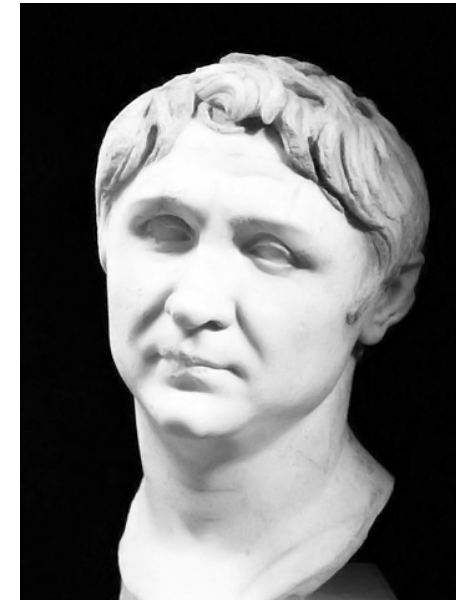
# The Metadata Triumvirate

**Triumvirate 1.0, 60 BC – "Conquer the World!"**
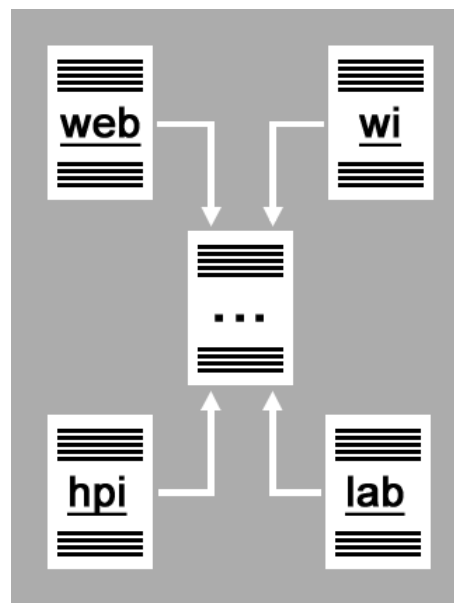


**Caesar**          **Crassus**          **Pompeius**

**Triumvirate 2.0, 2008 AD – "Conquer the World Wide Web?"**



**Social Annotations**         **Anchor Texts**         **Search Queries**

# Social Annotations



- **Definition of a social annotation:** list of "tags" (words) with which a social bookmark has been annotated

- Derived from **user-provided metadata**

- What does the social annotation **"web, conference, sydney, 2008"** tell about the user and the annotated document?

- Used for Web search personalization, emerging semantics, content classification, expert identification, …

# Social Annotations



Sergey's social annotation of Facebook.com

- **Definition of an anchor text:**
  words within <a>…</a> HTML element

- Derived from **Web link structure**

- What do the anchor texts **"web", "wi", "hpi", "lab"** tell about the linked page?

- Used for gaining more information about the linked Web pages, for improving indexing and ranking techniques, …

# Search Queries



- **Definition of a search query:**
  search keywords of the user's query

- Derived from search query logs,
  i.e. **user interactions**

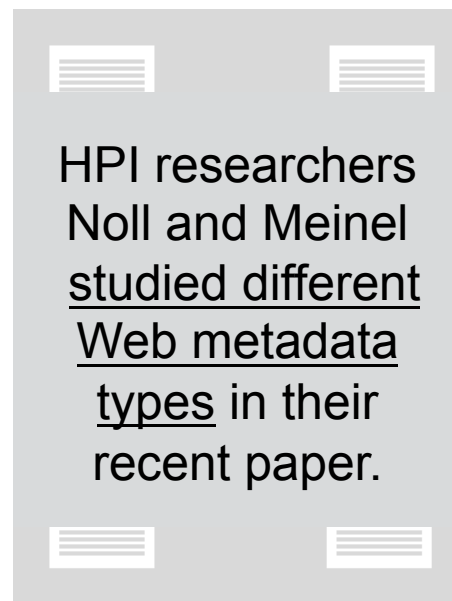- What does the search **"web wi 2008"**
  tell about the searcher or the clicked search
  result document?

- Used for query rewriting, user profiling,
  extracting semantics, …
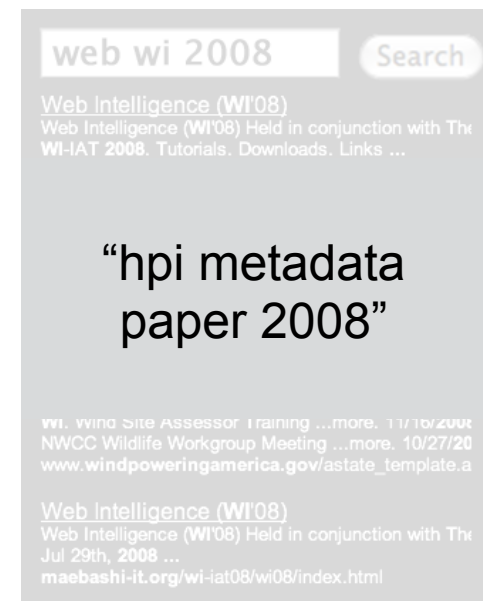
# Questions we want to answer

# How do these different types of metadata compare?

animals architecture art august austral
birthday black blackandwhite blue bosto
canon car cat chicago china chi

"metadata, paper, social web, hpi, research, 2008"

scotland sea seattle show sky snow s
taiwan texas thailand tokyo toronto tour
vacation vancouver washington wate

HPI researchers Noll and Meinel <u>studied different Web metadata types</u> in their recent paper.

web wi 2008          Search

Web Intelligence (WI'08)
Web Intelligence (WI'08) Held in conjunction with The
WI-IAT 2008. Tutorials. Downloads. Links ...

"hpi metadata paper 2008"

WI. Wind Site Assessor Training ...more. 11/16/2008
NWCC Wildlife Workgroup Meeting ...more. 10/27/20
www.windpoweringamerica.gov/astate_template.a

Web Intelligence (WI'08)
Web Intelligence (WI'08) Held in conjunction with The
Jul 29th, 2008 ...
maebashi-it.org/wi-iat08/wi08/index.html

**Social Annotations**          **Anchor Texts**          **Search Queries**

**Five questions**

- Q1: **Volume** of data per single metadata item?

- Q2: **New data** per metadata type?

- Q3: **Homogeneous** or **heterogeneous** metadata?

- Q4: **Similarity** between metadata types?

- Q5: Usefulness for **classification** of web documents?

# Experimental Setup

# Experimental data

**We created our own experimental data set "CABS120k08" in 2008**

- Bootstrapped by an intersection of
  **AOL500k** and **Open Directory Project**

  + targeted **Web crawl**
  + scraping **Delicious**
  + retrieving **Google PageRank**

  = metadata for **120,000** web documents

| Overview of CABS120k08 |
|---|
| 120,000 web documents |
| 2,600,000 search queries |
| 85,000 categories |
| 2,200,000 anchor texts |
| 1,300,000 social annotations |
| 120,000 PageRank scores |

Data set (500 MB) is available for download at:
**http://www.michael-noll.com/cabs120k08/**

# Experimental Results

# Q1: Volume of data per single metadata item?

"Does a social annotation provide more data than an anchor text?"

or: "How much data do users provide when using a specific metadata type?"

**Approach**

- Measure size of a single metadata item by its "length"

- Definitions of length for…

  - Social annotation → number of **tags**

  - Anchor text → number of **words**

  - Search query → number of **search keywords**

# Experimental results

**Mean length**

- Social annotation:  2.49
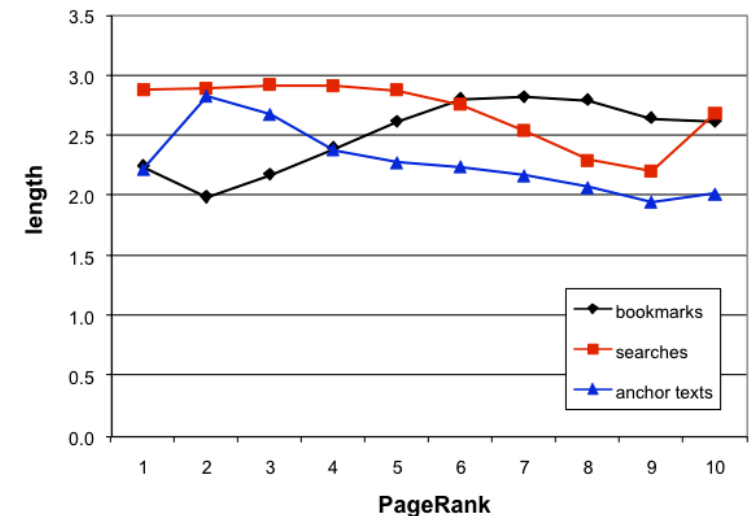
- Anchor text:  2.43

- Search query:  2.89

→ Surprisingly, **2.x** seems to be a "magic number" for user behavior across different problem domains (social bookmarking, hyperlink creation, Web search). Human psychology?

# Experimental results

Correlation of length with document popularity:

- **positively** for social annotations
- **negatively** for anchor texts and search queries



→ Anchor texts provide more metadata for less popular documents, whereas social annotations do so for popular ones

## Q2: New data per metadata type?

"How helpful is an analysis of a given metadata type
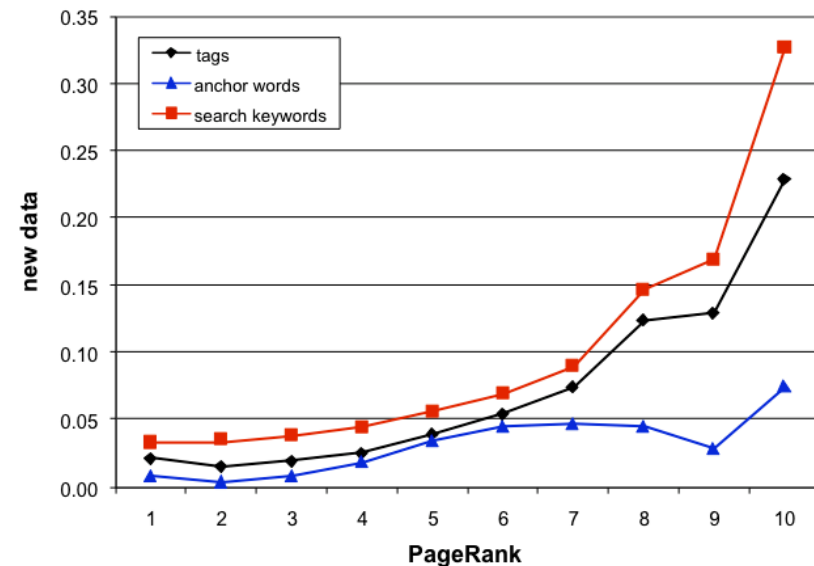for discovering new information?"

**Approach**

- Measure "novelty" of data provided by each metadata type

- Novelty is defined as the **percentage of unique terms** which are **new to a Web document**, i.e. terms that are not already present in the document's <TITLE>, <BODY>, plus selected HTML metadata

- For example, to retrieve a Web document in a search for "biology" even though the query term "biology" is not part of the document's HTML content.

# Experimental results

- Generally, the amount of new information is **relatively low**

- ≤ 6% for 90% of documents

- Search queries >> social annotations >> anchor texts



→ Compared to anchor texts, social annotations are a better source of new data

→ However, similarity between social annotations and anchor texts (as we see later) is rather low = they provide **different** data, so both are useful!

# Q3: Homogeneous or heterogeneous metadata?

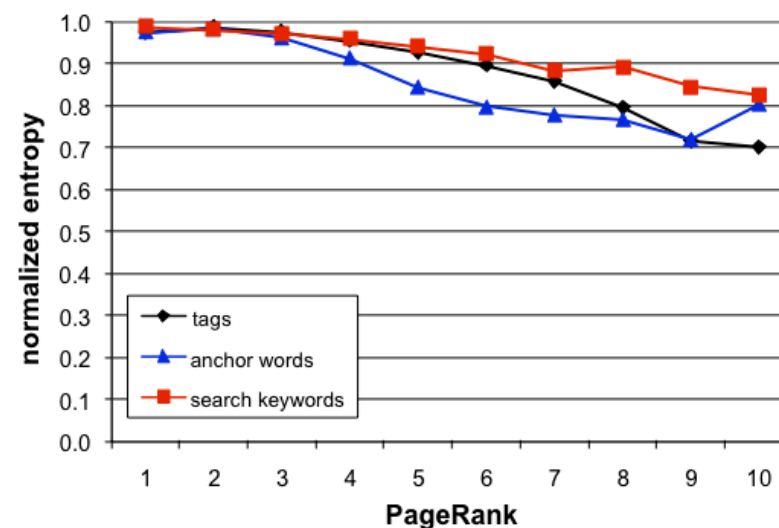"Is the data of each metadata type consistent/diverse/chaotic…?"

**Approach**

- Measure "diversity" of data **within** a given metadata type

- **Entropy** is used to measure diversity based on terms and term counts

- Note: Scoring a **high diversity** can indicate both **positive** (capturing different perceptions/meanings of content) and **negative** results (no consensus, noise).

# Experimental results

- Strong **negative correlation** with document popularity for all types: With increasing popularity, diversity of information decreases.

- Highest diversity for search queries: most "random" task, formulating good queries, spelling corrections ?

- Social annotations more diverse than anchor texts

→ Potential advantage for social annotations as they might capture information and meanings that anchor texts miss (cf. Bao et al. WWW 2007).

# Q4: Similarity between metadata types?

"How similar is the data provided by these metadata types?"

**Approach**

- Study the **interrelations** between metadata types

- **Pairwise cosine similarity** is used to measure similarity

- Preprocessing of terms: splitting ("new_york"), stemming, stop words

# Experimental results

|  | Social annotat. | Anchor texts | Search queries | Categories |
|---|---|---|---|---|
| **Social annotat.** | x | 0.126 | 0.126 | **0.189** |
| **Anchor texts** | 0.126 | x | **0.193** | 0.103 |
| **Search queries** | 0.126 | **0.193** | x | 0.102 |
| **Categories** | **0.189** | 0.103 | 0.102 | x |

Highest similarities for two pairs:

- sim(**social annotations, categories**) = 0.189 → "better" for classification?

- sim(**anchor texts, search queries**)　 = 0.193 → "better" for Web search?

**Q5: Usefulness for classification of web documents?**

"How helpful are these metadata types for classification tasks?"

**Approach:**
Matching data of each metadata type against a document's categorization trees from Open Directory Project

category depth

art

0.0

crafts

textiles

weaving

1.0

# Experimental results

- Strong **negative correlation** with document popularity for all types: With increasing popularity, broader classification scores are achieved.

- Social annotations are "used" for broader classification than anchor texts and search queries



→ Of all three, social annotations seem to be the best at classification tasks

# Conclusions

# The Metadata Triumvirate

**Triumvirate 1.0**



**Caesar**          **Crassus**          **Pompeius**

**Worked out quite well…**

Roman Empire, 44 BC

# The Metadata Triumvirate

**…however…**



| **Caesar** | **Crassus** | **Pompeius** |

**Metadata Triumvirate – no casualties (yet)!**



**Social Annotations**   **Anchor Texts**   **Search Queries**

# Conclusions

- First study to compare social annotations, anchor texts and search queries directly on a large volume of real-world data

- Starting point for future research

- Research data set CABS120k08, available for free download: http://www.michael-noll.com/cabs120k08/

Contact info:
Michael G. Noll
michael.noll@hpi.uni-potsdam.de
Hasso Plattner Institute
Potsdam, Germany