



## **Web Mining**

# **Introduction and Course Organization**

**Prof. Dr. Christian Bizer  
Prof. Dr. Simone Ponzetto  
Alexander Brinkmann  
Ralph Peeters**

**FSS 2023**

- **Prof. Dr. Christian Bizer**
- **Professor for Information Systems V**
- **Research Interests:**
  - Information Extract from the Web
  - Large-Scale Data Integration
  - Data and Web Mining
- **Room: B6, 26 - B1.15**
- **Consultation: Wednesday 13:30-14:30**
- **eMail: [christian.bizer@uni-mannheim.de](mailto:christian.bizer@uni-mannheim.de)**
  
- **will teach the lectures on Web Usage Mining and Web Structure Mining**



# Hallo

- **Prof. Dr. Simone Ponzetto**
  - **Professor for Information Systems III**
  - **Research Interests:**
    - Natural Language Processing
    - Computational Social Science
  - **Room: B6 26, B 1.14**
  - **Consultation: Tuesday 13:30-14:30**
  - **eMail: [ponzetto@uni-mannheim.de](mailto:ponzetto@uni-mannheim.de)**
- 
- **will teach the lectures on Web Content Mining**



- **M. Sc. Wi-Inf. Alexander Brinkmann**
- **Graduate Research Associate**
- **Research Interests:**
  - Data Search using Deep Learning
  - Semantic Annotations in Web Pages
  - Product Data Categorization
- **Room: B6, 26, C 1.03**
- **eMail: [alexander.brinkmann@uni-mannheim.de](mailto:alexander.brinkmann@uni-mannheim.de)**
  
- **will teach the labs and will supervise student projects (IE684)**



- **M. Sc. Wi-Inf. Ralph Peeters**
- **Graduate Research Associate**
- **Research Interests:**
  - Entity Matching using Deep Learning
  - Product Data Integration
- **Room: B6, 26, C 1.04**
- **eMail: [ralph.peeters@uni-mannheim.de](mailto:ralph.peeters@uni-mannheim.de)**
  
- **will teach the labs and will supervise student projects (IE684)**



# Introduction and Course Organization

## 1. Course Organization

## 2. The World Wide Web

1. The Classic Document Web
2. The Web of Data
3. Web 2.0 Applications

## 3. What is Web Mining?

1. Web Usage Mining
2. Web Structure Mining
3. Web Content Mining

# 1. Course Organization

## ■ Lecture (IE 671, 3 ECTS)

- covers different types of Web Mining methods
- presents examples of Web Mining applications
- discusses how to evaluate the learned models

## ■ Labs

- students experiment with the methods using different Python libraries

## ■ Evaluation

- 60 min written exam

## ■ Student Projects (IE 684, 3 ECTS)

- teams of five students realize a Web Mining project
- teams may choose their own tasks and data sets  
(in addition, we will propose some suitable data sets and tasks)
- write a summary about the project, present the project results

## ■ Evaluation

- report + presentation + code

# Course Organization

## ■ Website

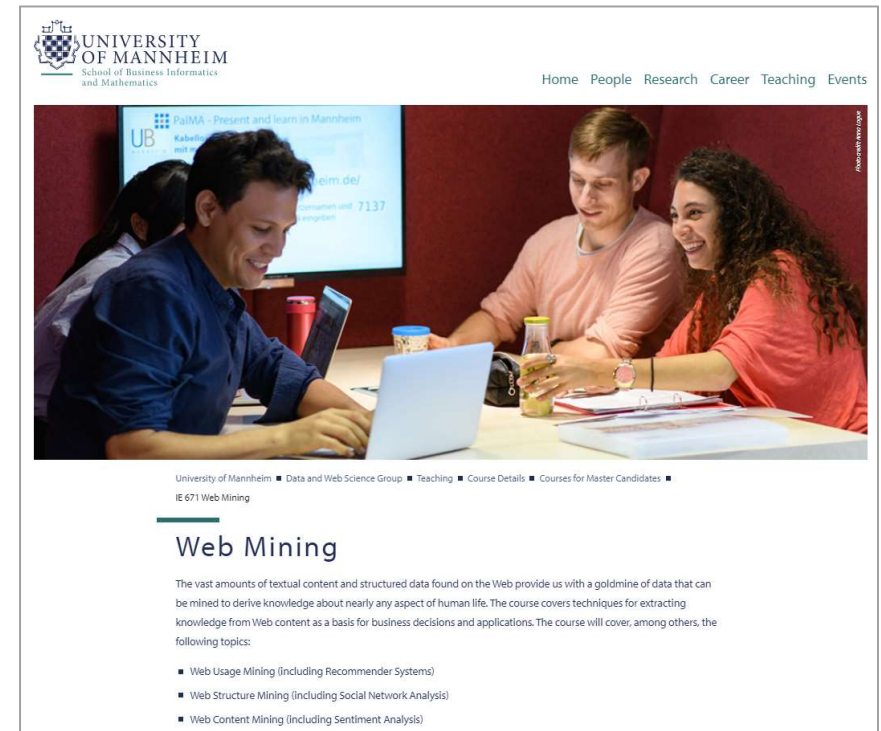
- organizational information
- lecture slides

## ■ Two IIAS Groups

- IE 671: Material for Labs + Mailing List (Web Mining [Ü])
- IE 684: Student Projects (Web Mining Project [PRO])

## ■ Time and Location

- Lecture: Tuesday, 10:15 to 11:45, B6 A101
- Labs: Thursday, 13:45 to 15:15, B6, A203



The screenshot shows the University of Mannheim website for the Web Mining course. The header includes the university logo and navigation links: Home, People, Research, Career, Teaching, Events. The main content area features a photograph of three students (two men and one woman) sitting around a table, working on laptops. Below the photo is a breadcrumb trail: University of Mannheim | Data and Web Science Group | Teaching | Course Details | Courses for Master Candidates | IE 671 Web Mining. The title 'Web Mining' is followed by a paragraph describing the course: 'The vast amounts of textual content and structured data found on the Web provide us with a goldmine of data that can be mined to derive knowledge about nearly any aspect of human life. The course covers techniques for extracting knowledge from Web content as a basis for business decisions and applications. The course will cover, among others, the following topics:'. A bulleted list of topics follows: Web Usage Mining (including Recommender Systems), Web Structure Mining (including Social Network Analysis), and Web Content Mining (including Sentiment Analysis).

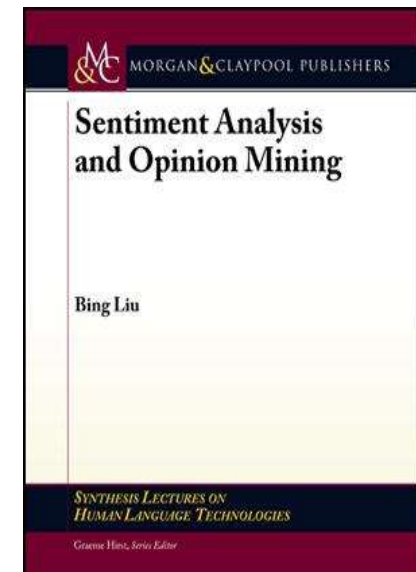
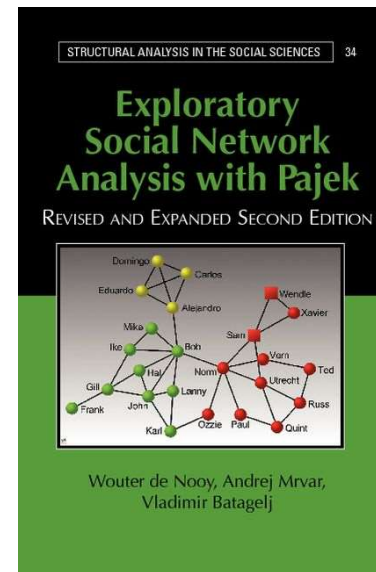
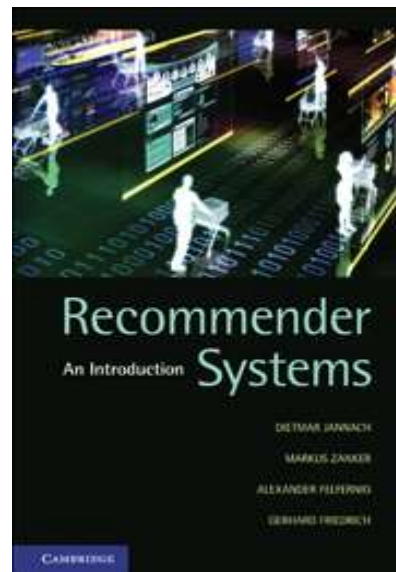


# Schedule

<b>Week</b>	<b>Lecture: Tuesday 10:15</b>	<b>Labs: Thursday 13:45</b>
<b>14.02.2023</b>	Introduction and Course Outline	- No Lab -
<b>21.02.2023</b>	Web Usage Mining (Bizer)	Web Usage Mining
<b>28.02.2023</b>	Web Usage Mining (Bizer)	Web Usage Mining
<b>07.03.2023</b>	Web Structure Mining (Bizer)	Web Structure Mining
<b>14.03.2023</b>	Web Structure Mining (Bizer)	Web Structure Mining
<b>21.03.2023</b>	Web Content Mining (Ponzetto)	Web Content Mining
<b>28.03.2023</b>	Web Content Mining (Ponzetto)	Web Content Mining
	- Easter break -	
<b>18.04.2023</b>	Introduction to the Student Projects	Preparation of project outlines
<b>25.04.2023</b>	Feedback on the projects outline	Project work
<b>02.05.2023</b>	Project work	Coaching
<b>09.05.2023</b>	Project work	Coaching
<b>16.05.2023</b>	Project work	Coaching
<b>23.05.2023</b>	Project work	Coaching
<b>30.05.2022</b>	Project presentations	Project presentations
<b>XX.06.2023</b>	Final exam	-


# Literature

1. Bing Liu: Web Data Mining. 2nd Edition, Springer.
2. Dietmar Jannach, et al.: Recommender Systems: An Introduction. Cambridge University Press.
3. Wouter de Nooy et al.: Exploratory Social Network Analysis with Pajek. Cambridge University Press.
4. Bing Liu: Sentiment Analysis and Opinion Mining. Morgan & Claypool.



# Software Libraries

- Within the Labs, we use the following libraries

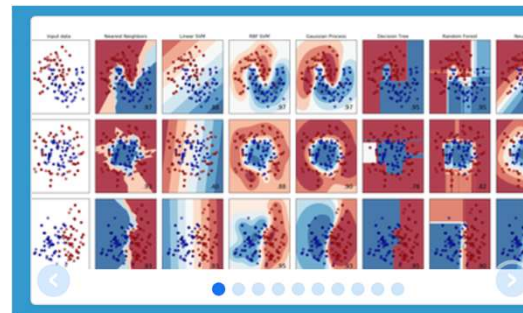


A Python scikit for recommender systems.

Home  
Documentation  
GitHub page

★ Star    🍴 Fork

Maintained by Nicolas Hug  
Page built with Jekyll and Hyde



## scikit-learn

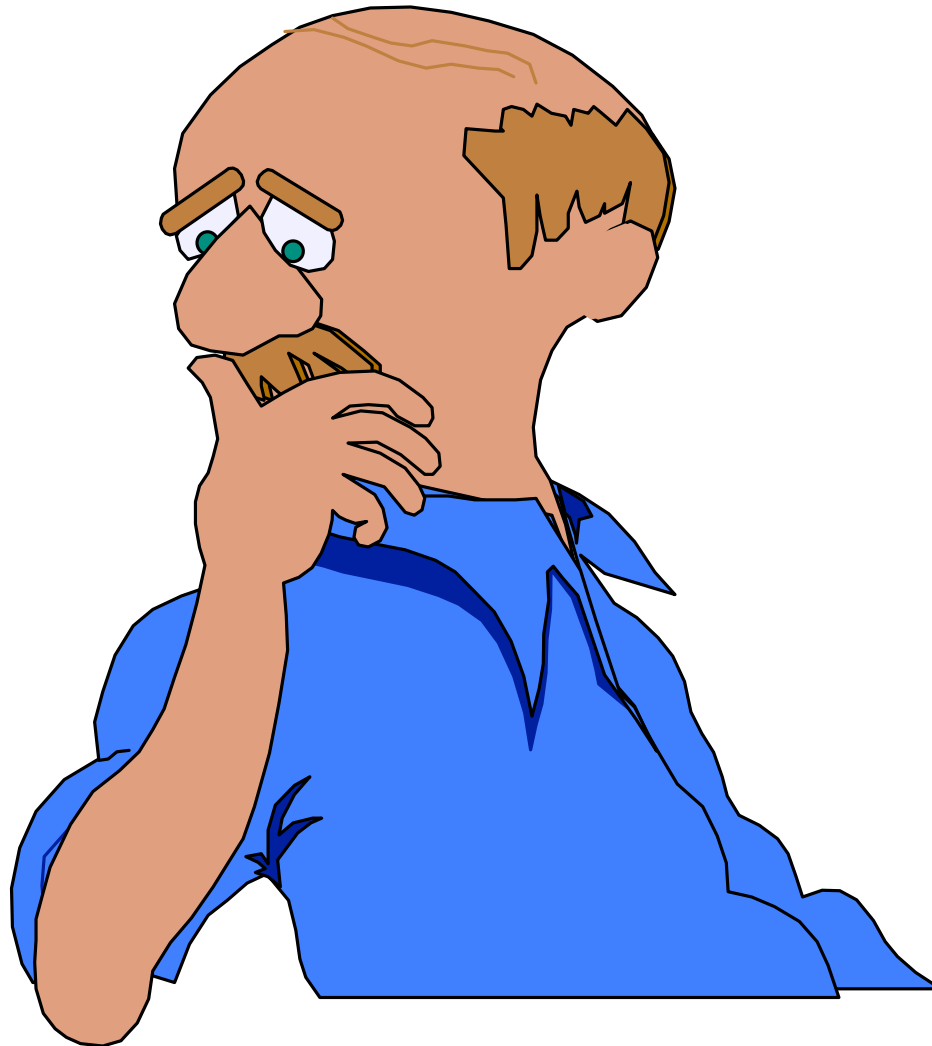
*Machine Learning in Python*

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license



- We assume some fluency in Python and Jupyter ...

# Questions?



## 2. The World Wide Web

The Web is a global decentralized information space build on a set of technical standards for the identification, retrieval and representation of content.

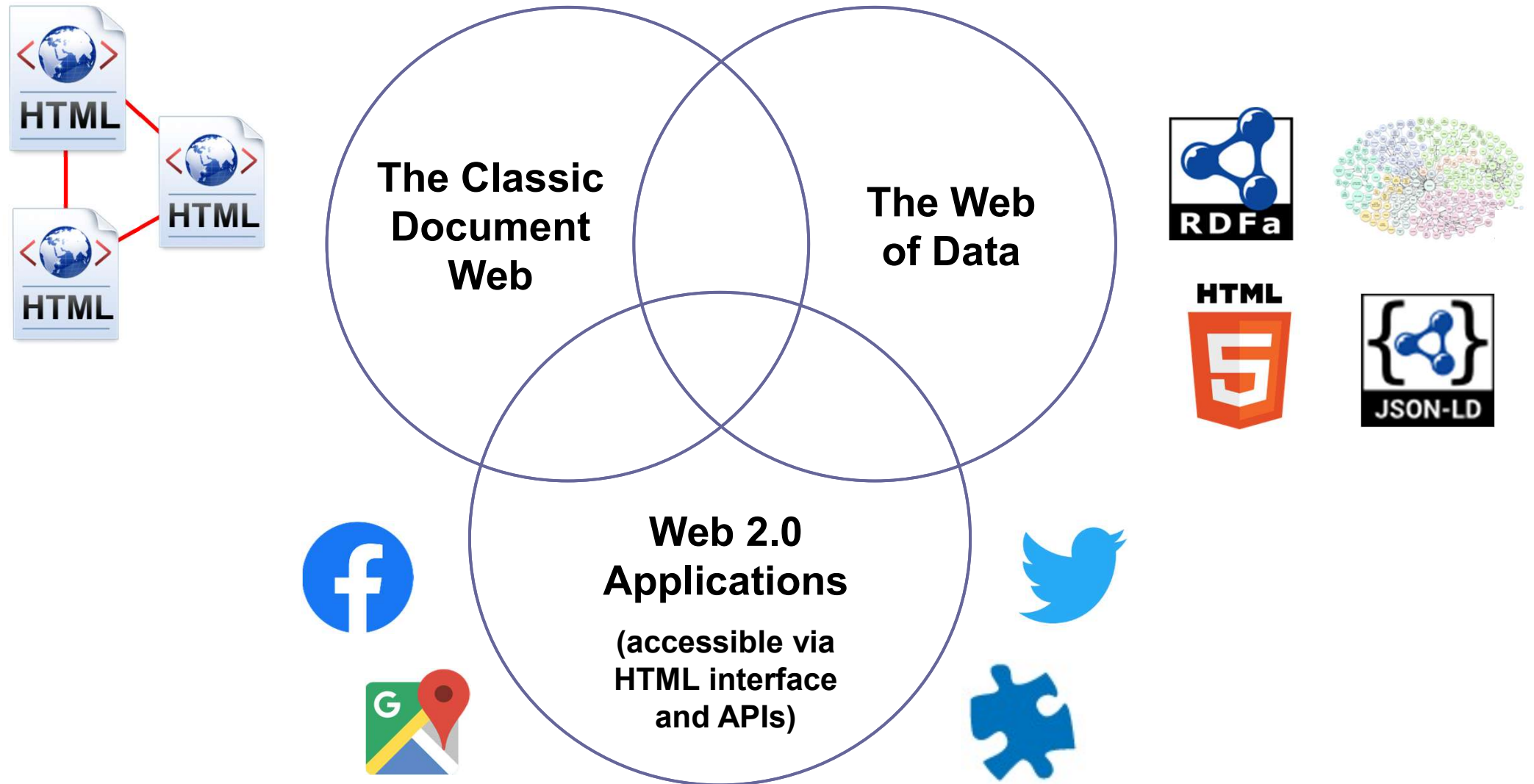


- **Uniform Resource Identifiers (URIs):** Globally unique identification of Web resources.
- **Hypertext Transfer Protokoll (HTTP):** Protocol for interacting with Web resources.
- **Content Formats:** HTML, XML, RDF, ...
- **The Web was invented in 1989 at CERN by Tim Berners-Lee**
- **Architectural Principles of the Web**



Architecture of the World Wide Web, Volume One,  
W3C Specification, 2004, <http://www.w3.org/TR/webarch/>

# Topology of the Web Today



# 2.1 The Classic Document Web

Global information space consisting of interlinked **resources** (HTML pages, images, multimedia).



## The Size of the Web

- **overall size: 1 trillion pages on the Web at once**
  - announced by Google in 2008
  - <http://googleblog.blogspot.de/2008/07/we-knew-web-was-big.html>
- **Indexed Web: approx. 50 billion pages**
  - estimate based on search engine hit counts for popular words
  - Example: The word „the“ appears in 67% of all English pages and has 25.2 billion hits on Google
  - <http://www.worldwidewebsite.com/>





# Public Web Corpora

- **The non-profit Common Crawl project regularly crawls the Web and publishes large Web corpora**
  - size: 2.5 to 3.5 billion pages
  - pay-level-domains: 30 – 35 million
  - release cycle: 1 month
- **Public download from Amazon S3**
  - size: around 80 Terrabytes compressed
- **Public alternative to private corpora owned by Google and Microsoft**
- **Disadvantages:**
  - one order of magnitude smaller than private crawls
  - monthly releases compared to permanent updates



Sebastian Nagel





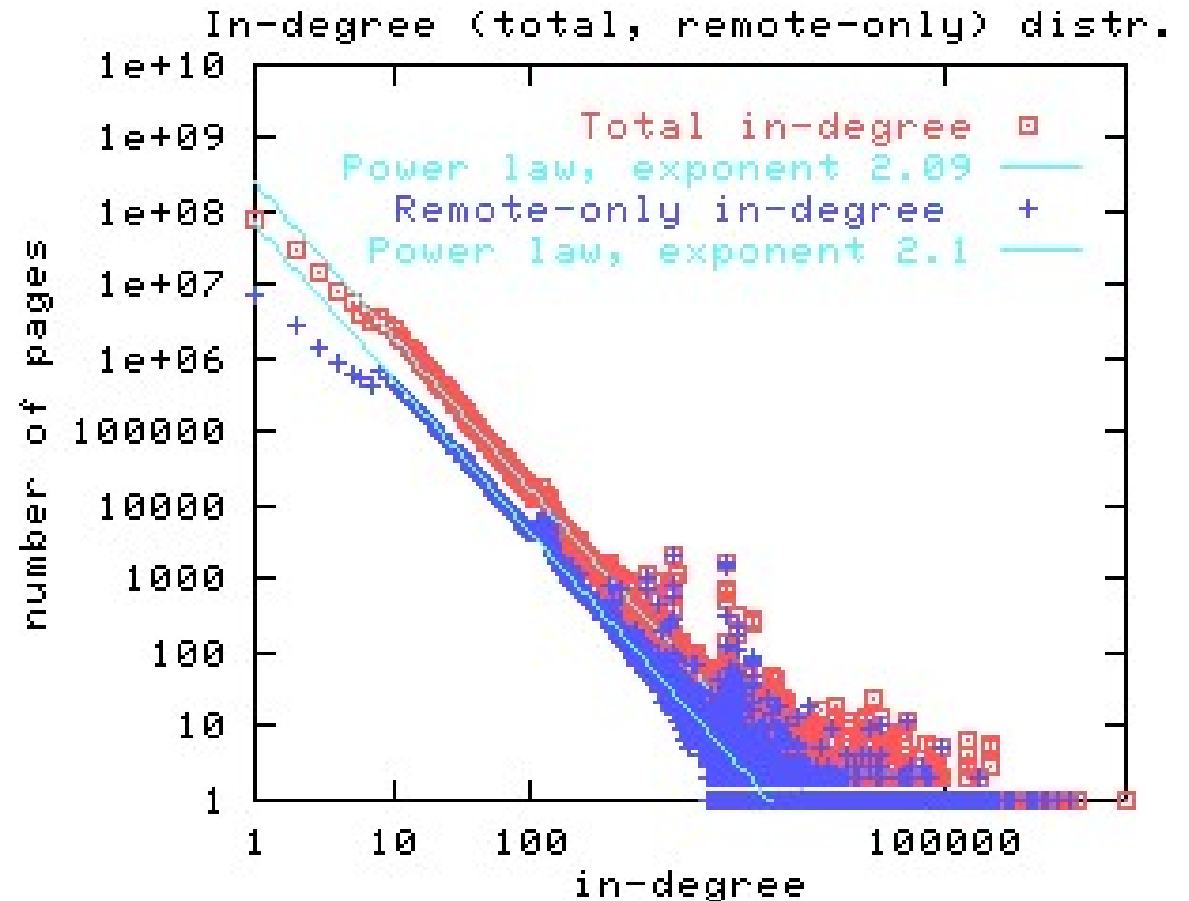
# Link Structure of the Web: In-Degree

## ■ The link distribution follows (kind of) a power law

- A small number of pages is target of many links
- A large number of pages is target of only a few or no links

## ■ Classic Paper:

- Broder et al.:  
Graph Structure  
in the Web. WWW2000
- AltaVista crawl with over  
200 million pages and  
1.5 billion links
- Conclusion: Log-log scale  
plot shows power-law

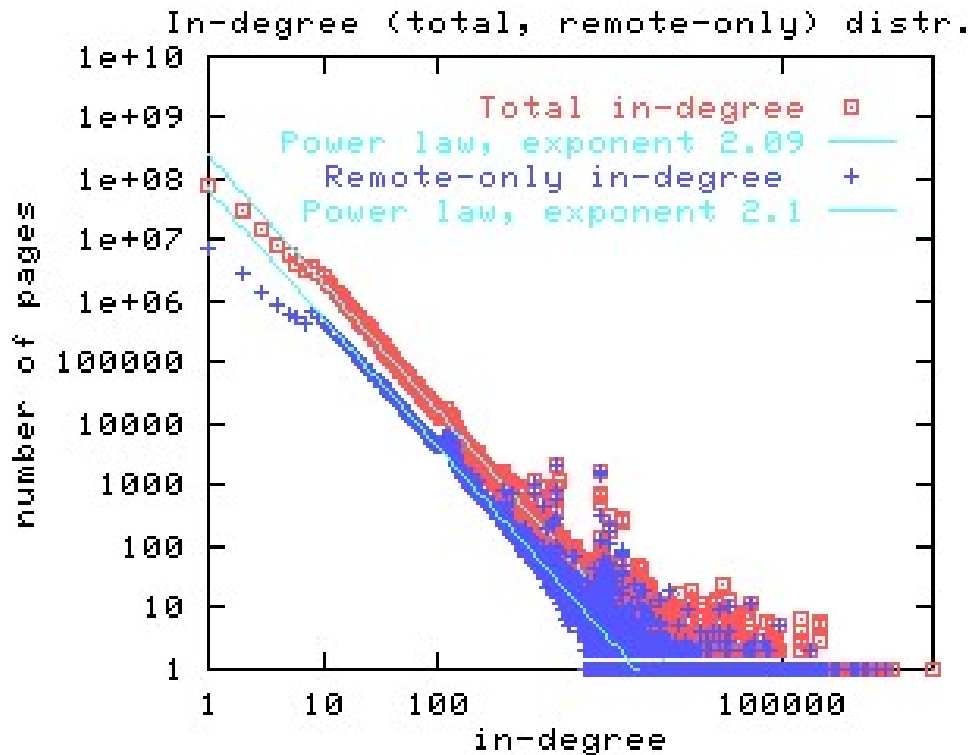


# In-Degree Distribution

Broder et al. (2000)

Power law with exponent **2.1**

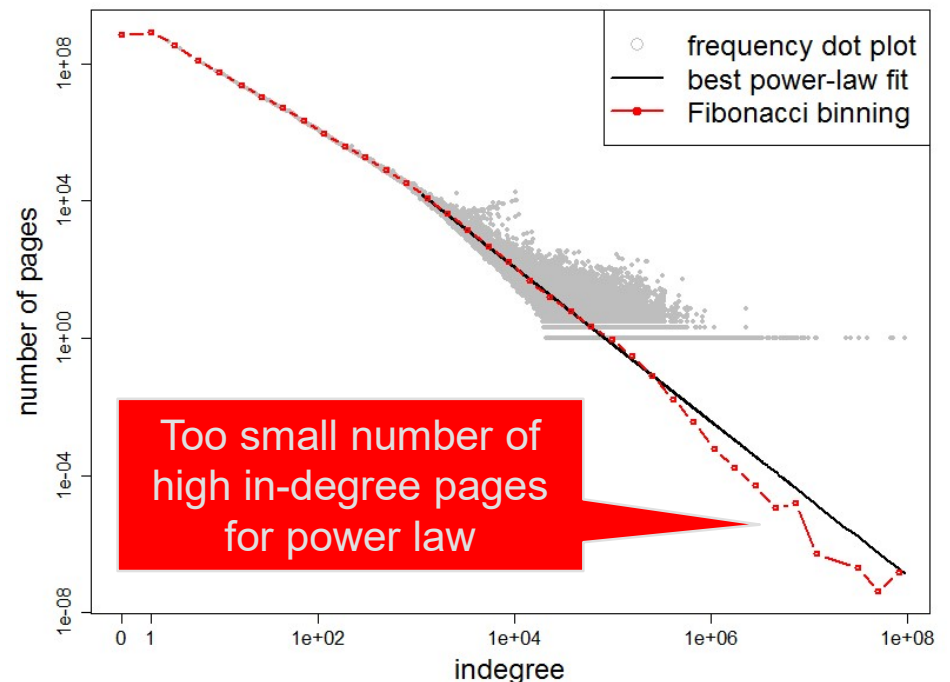
(200 million pages and 1.5 billion links from Altavista crawl 2000)



WDC Hyperlink Graph (2012)

Best power law exponent **2.24**

(3 billion pages and 128 billion links from Common Crawl 2012)



# Link Structure of the Web: Bow-Tie

## Four major components (Border et al., WWW2000)

### ■ Central Strongly Connected Component (SCC)

- pages that can reach one another along directed links
- about 30% of the Web (normal pages)

### ■ IN Group

- can reach SCC but cannot be reached from it
- about 20% (maybe new pages or boring ones)

### ■ OUT Group

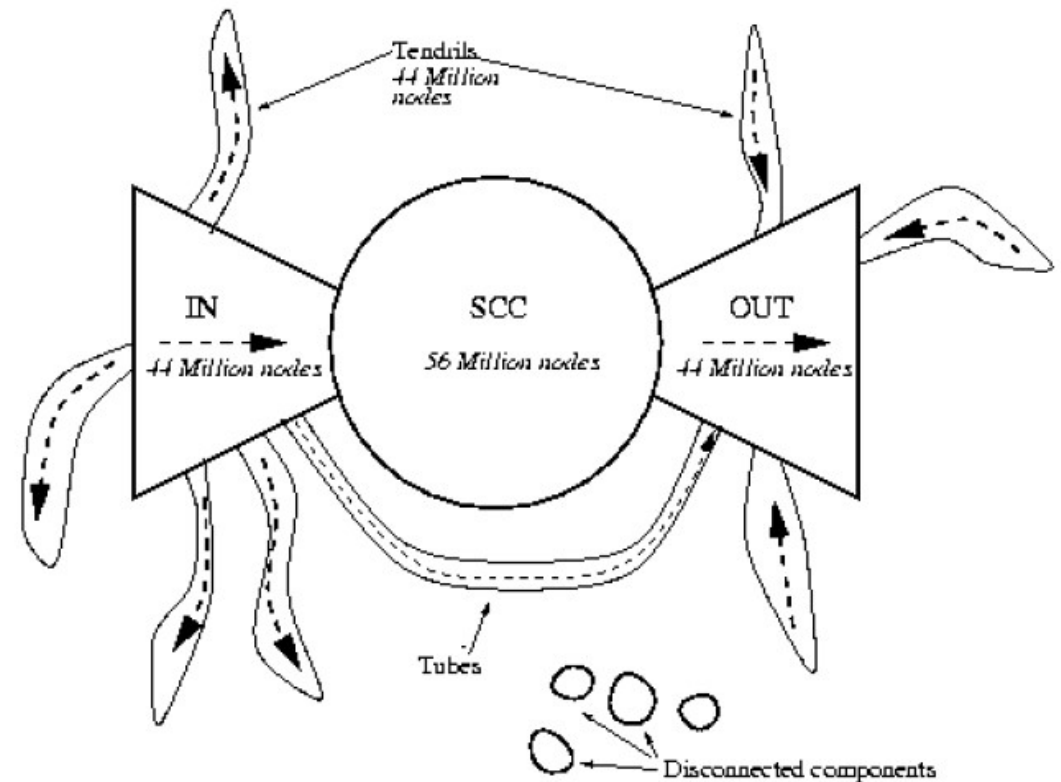
- can be reached from SCC but cannot reach it
- about 20% (maybe company pages that don't link)

### ■ Tendrils

- cannot reach SCC and cannot be reached by it
- about 20%

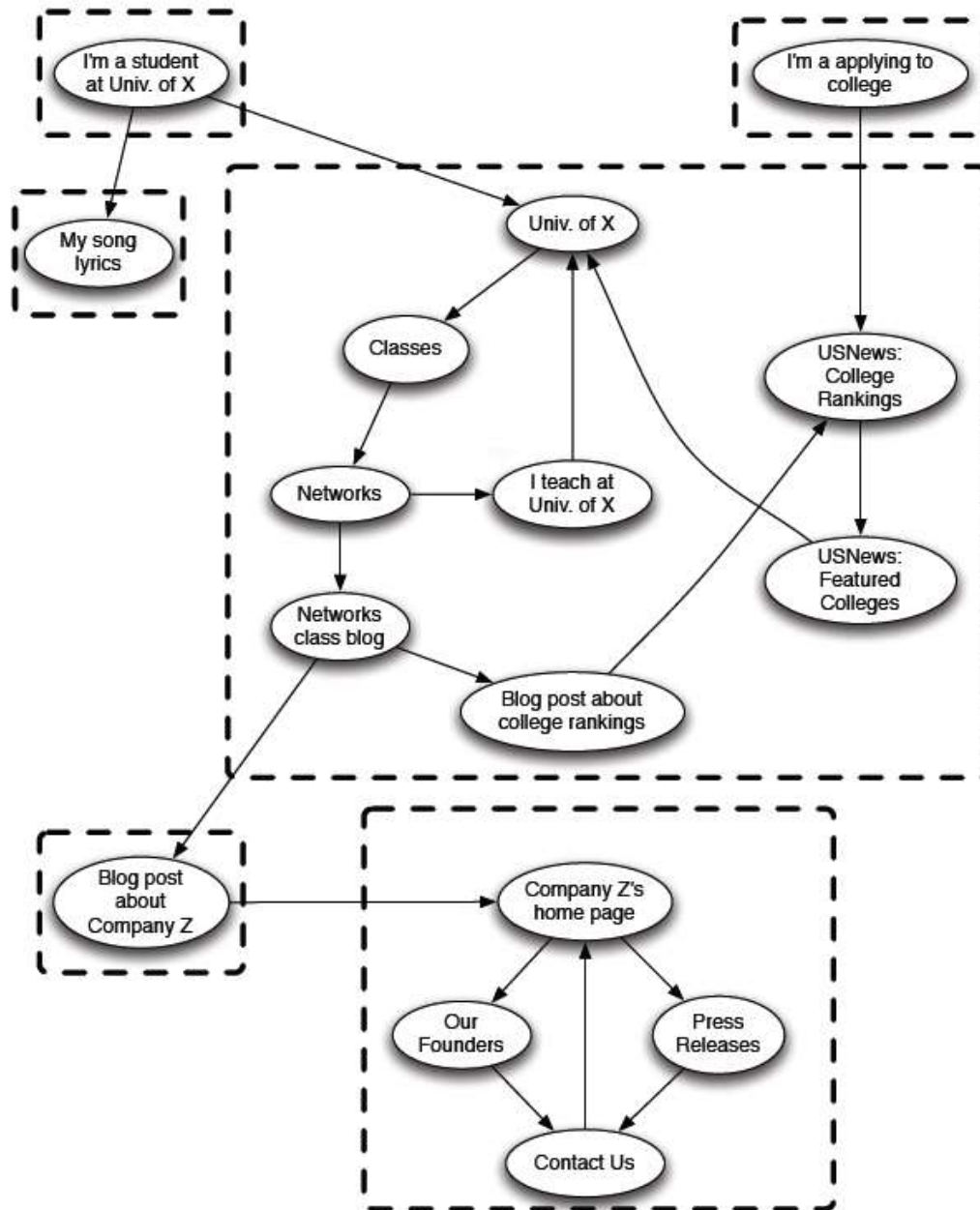
### ■ Unconnected

- about 10%



**Probability of path  
between nodes is 24%**

# Web Graph with Strongly Connected Components



A strongly connected component (SCC) in a directed graph is a subset of the nodes such that:

1. every node in the subset has a path to every other node
2. the subset is not part of some larger set with the property that every node can reach every other.

# Size of Central Strongly Connected Component

## Largest SCC

■ Broder, 2000:

**27.7%**

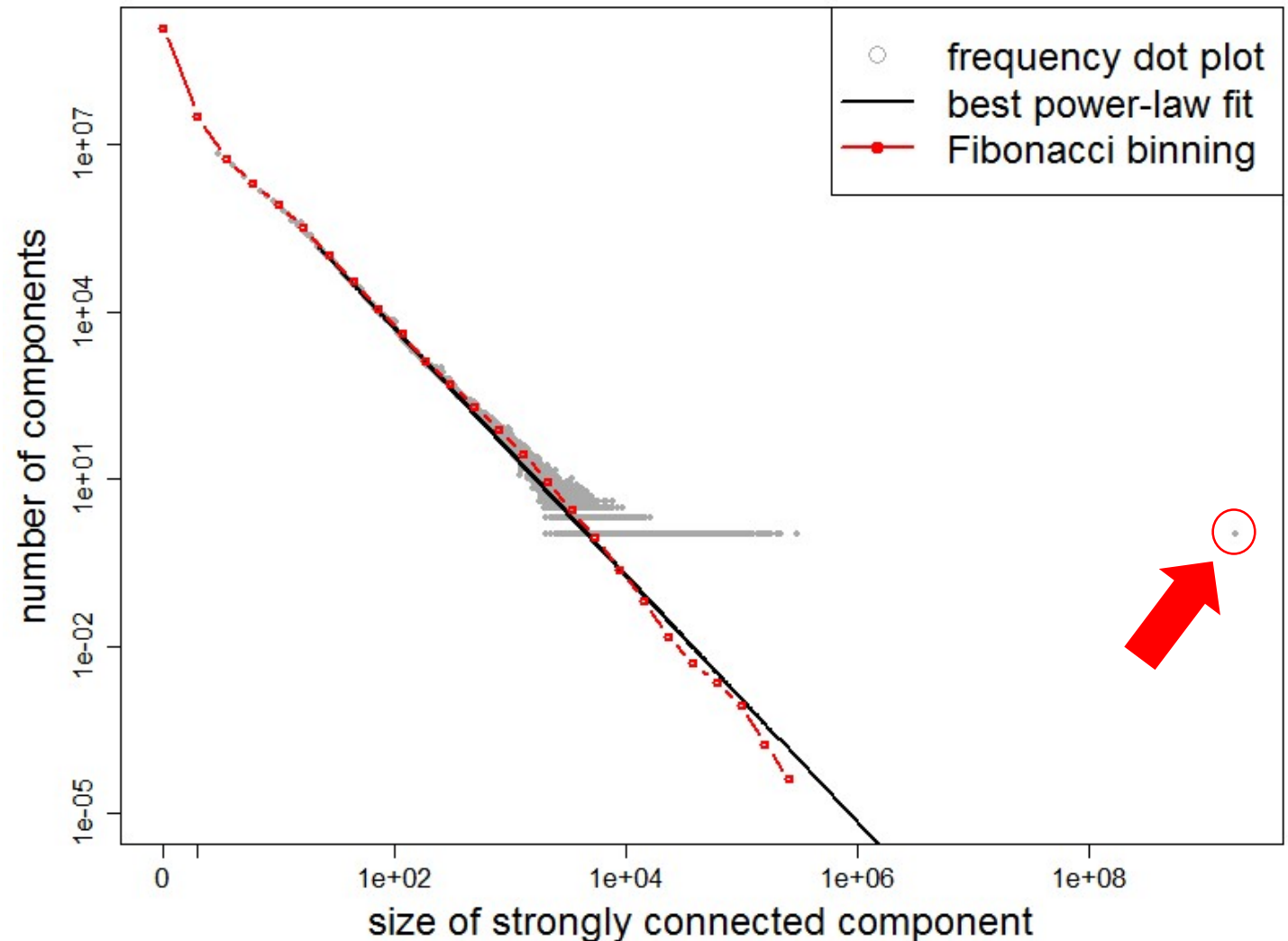
■ WDC, 2012:

**51.3 %**

➔ Factor 1.8 larger

➔ Also, factor 4.9 more links/page

➔ The Web has become denser



## 2.2. The Web of Data

### More and more Websites

- semantically **markup the content** of their HTML pages
- **publish structured data** in addition to HTML

### Markup Formats:

**Microdata**



**RDFa**

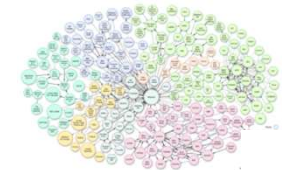


### Structured Data Formats:

**JSON-LD**



**Linked Data**



**Data is crawl-able using generic crawlers in contrast to Web APIs**

# Microdata

- format for annotating structured data within webpages
- proposed in 2009 by WHATWG as part of HTML5 work



```
<div itemtype="http://schema.org/Hotel">
  <span itemprop="name">Vienna Marriott Hotel</span>
  <span itemprop="address" itemscope="" itemtype="http://schema.org/PostalAddress">
    <span itemprop="streetAddress">Parkring 12a</span>
    <span itemprop="addressLocality">Vienna</span>
  </span>
  <div itemprop="aggregateRating" itemscope itemtype="http://schema.org/AggregateRating">
    <span itemprop="ratingValue"> 4 </span> stars-based on
    <span itemprop="reviewCount"> 250 </span> reviews.
  </div>
```

- used for embedding data into the HEAD of HTML pages
- putting data in HEAD is recommended by Google as it is empirically less error prone than annotations in BODY



```
<script type="application/ld+json">
{  "@context": "http://schema.org",
   "@type": "Product",
   "description": "Has six preset cooking ....",
   "name": "Kenmore White 17\" Microwave",
   "offers": {
     "@type": "Offer",
     "availability": "http://schema.org/InStock",
     "price": "55.00",
     "priceCurrency": "USD"
   },
}
</script>
```



- ask site owners since 2011 to annotate data for enriching search results
- 675 Types: Event, Local Business, Product, Review, Job Offer
- Encoding: Microdata, RDFa, JSON-LD

schema.org

Home Schemas Documentation

**Thing > Organization > LocalBusiness**

A particular physical business or branch of an organization. Examples of LocalBusiness include a restaurant, a particular branch of a restaurant chain, a branch of a bank, a medical practice, a club, a bowling alley, etc.


Property	Expected Type	Description
<b>Properties from Thing</b>		
description	Text	A short description of the item.
image	URL	URL of an image of the item.
name	Text	The name of the item.
url	URL	URL of the item.
<b>Properties from Place</b>		
address	<a href="#">PostalAddress</a>	Physical address of the item.
aggregateRating	<a href="#">AggregateRating</a>	The overall rating, based on a collection of reviews or ratings, of the item.
containedIn	<a href="#">Place</a>	The basic containment relation between places.



# Usage of Schema.org Data @ Google

**Gramercy Tavern - Flatiron - New York, NY | Yelp**  
www.yelp.com › Restaurants › American (New) ▾  
★★★★★ Rating: 4.5 - 1,288 reviews - Price range: \$\$\$\$  
Jeff C and I were in **New York** for vacation, and I wanted to treat him to a nice dinner for ..... **Gramercy Tavern** is certainly a legendary NY dining establishment.

**Gramercy Tavern Restaurant - New York, NY | OpenTable**  
www.opentable.com › ... › Gramercy restaurants ▾  
★★★★★ Rating: 4.7 - 508 reviews - Price range: \$50 and over  
Book now at **Gramercy Tavern** in **New York**, explore menu, see photos and read 508 reviews: "The menu was so limited but it was worth trying, food was deli..."



## The Black Keys

Band

The Black Keys is an American rock duo formed in Akron, Ohio in 2001. The group consists of Dan Auerbach and Patrick Carney. [Wikipedia](#)

**Origin:** Akron, Ohio, United States

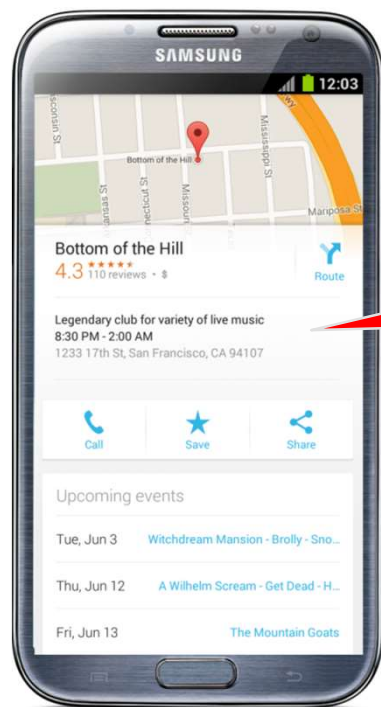
**Members:** Dan Auerbach, Patrick Carney

**Record labels:** Fat Possum Records, Nonesuch Records, V2 Records, Alive Natural Sound Records

**Awards:** Grammy Award for Best Rock Album, more

### Upcoming events

Jun 20 Fri	The Black Keys Neuhausen ob Eck (near you)
May 16 Fri	The Black Keys Gulf Shores, AL
Jun 22 Sun	The Black Keys Schneeßel

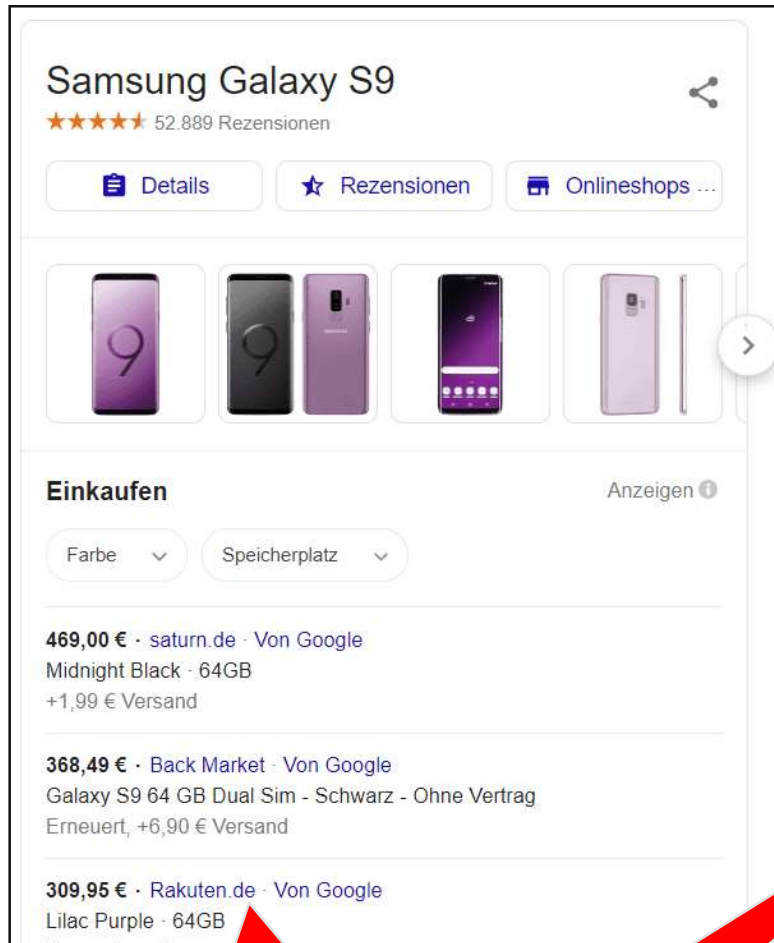


Data snippets within search results

Local businesses on maps

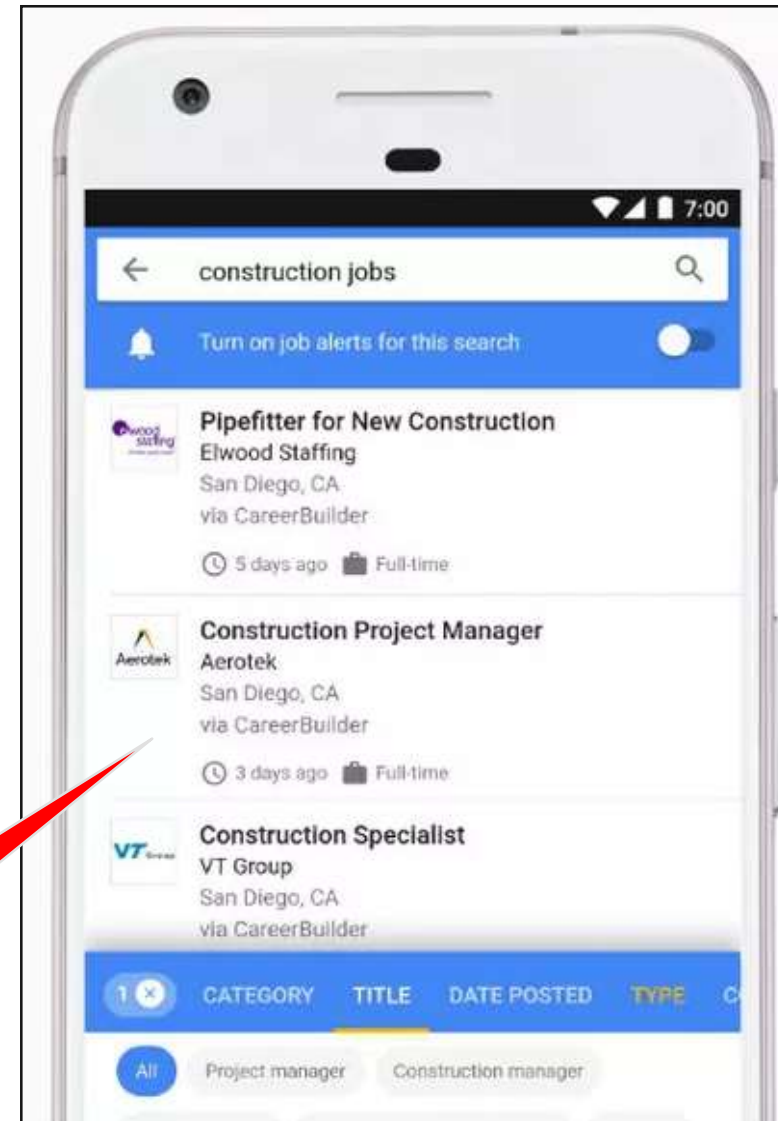
Data snippets within info boxes

# Usage of Schema.org Data @ Google



Product offers

Job offers



<https://developers.google.com/search/docs/guides/search-gallery>

# The Web Data Commons Project



Common Crawl



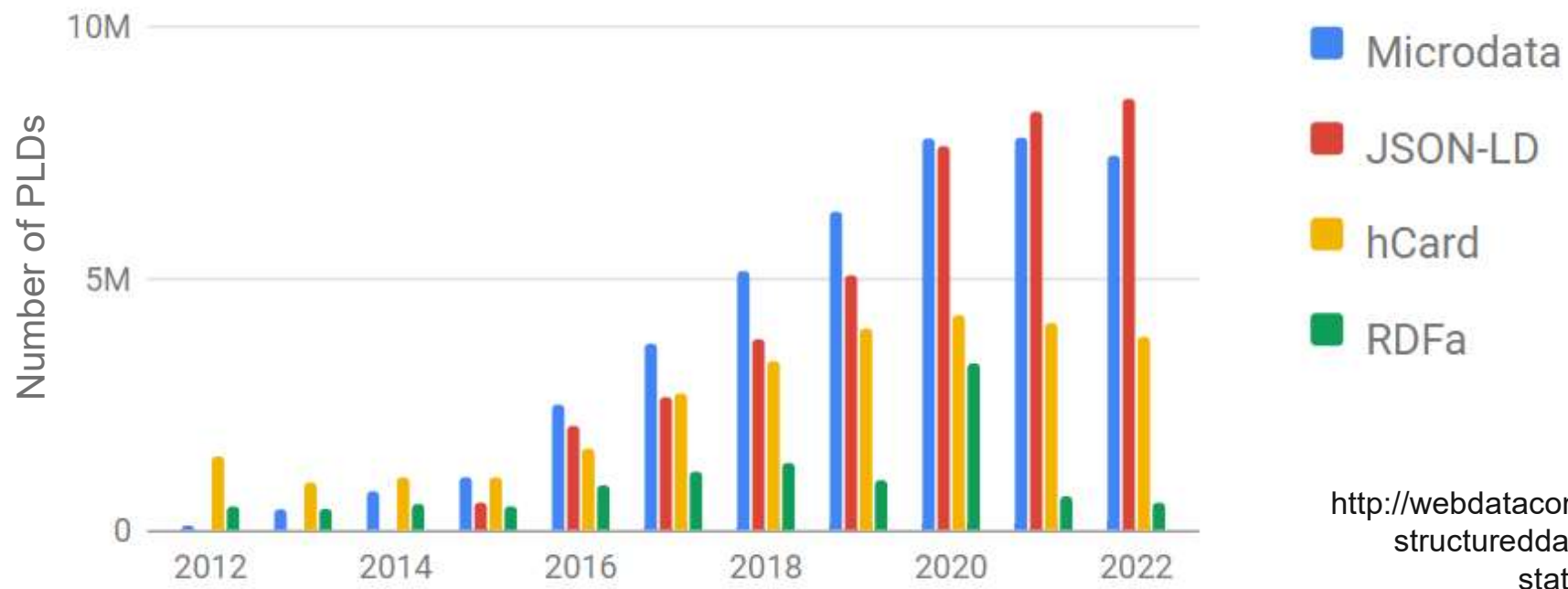
- extracts all Microformat, Microdata, RDFa, JSON-LD data from the **Common Crawl**
- analyzes and provides the extracted data for download
- statistics about some extraction runs
  - 2022 CC Corpus: **3.0 billion HTML pages** → **86.5 billion RDF triples**
  - 2019 CC Corpus: 2.4 billion HTML pages → 44.2 billion RDF triples
  - 2014 CC Corpus: 2.0 billion HTML pages → 20.4 billion RDF triples
  - 2010 CC Corpus: 2.8 billion HTML pages → 5.1 billion RDF triples
- uses **100 machines on Amazon EC2**
  - approx. 3500 machine/hours → 1000 Euro
- <http://webdatacommons.org/structureddata/>



# Overall Adoption 2022

1.5 billion HTML pages out of the 3.15 billion pages provide semantic annotations (46.9%).

14.2 million pay-level-domains (PLDs) out of the 33.8 million pay-level-domains covered by the crawl provide semantic annotations (42.0%).



# Frequently used Schema.org Classes

Class	# Websites (PLDs)	
	Microdata	JSON-LD
schema:WebPage	1,124,583	121,393
schema:Product	812,205	40,169
schema:Offer	676,899	57,756
schema:BreadcrumbList	621,344	205,971
schema:Article	612,361	57,082
schema:Organization	510,069	1,349,775
schema:PostalAddress	502,615	176,500
schema:ImageObject	360,875	111,946
schema:Blog	337,843	12,174
schema:Person	324,349	335,784
schema:LocalBusiness	294,390	249,017
schema:AggregateRating	258,078	23,105
schema:Review	124,022	6,622
schema:Place	92,127	66,396
schema:Event	88,130	63,605

<http://webdatacommons.org/structureddata/2018-12/>

# Adoption by Travel Websites

Top 15 Travel Websites	schema:Hotel
Booking.com	Yes
TripAdvisor	Yes
Expedia	Yes
Agoda	Yes
Hotels.com	Yes
Kayak	Yes
Priceline	Yes
Travelocity	Yes
Orbitz	Yes
ChoiceHotels	Yes
HolidayCheck	Yes
ChoiceHotels	Yes
InterContinental Hotels Group	Yes
Marriott International	Yes
Global Hyatt Corp.	No

Adoption:  
**93 %**

# Hands-on: How to get the Data?

- N-Quads: <http://www.webdatacommons.org/structureddata/>
- JSON: <http://webdatacommons.org/structureddata/schemaorgtables/>

## Class-Specific Subsets of the Schema.org Data

Class Name	Total Number of	Top Classes (Entity Count)	Total File Size	Quad File
<a href="http://schema.org/AdministrativeArea">http://schema.org/AdministrativeArea</a>	Quads: 1,724,857 URLs: 85,625 Hosts: 63	<a href="http://schema.org/AdministrativeArea">http://schema.org/AdministrativeArea</a> (100,671) <a href="http://schema.org/GeoCoordinates">http://schema.org/GeoCoordinates</a> (84,152) <a href="http://schema.org/Country">http://schema.org/Country</a> (83,851) <a href="http://schema.org/Continent">http://schema.org/Continent</a> (83,567)	23 MB	<a href="#">schemaorgAdministrativeArea.nq.gz (sample)</a>
<a href="http://schema.org/Hotel">http://schema.org/Hotel</a>	Quads: 148,211,253 URLs: 3,136,152 Hosts: 5,337	<a href="http://schema.org/Rating">http://schema.org/Rating</a> (7,007,590) <a href="http://schema.org/Hotel">http://schema.org/Hotel</a> (6,335,124) <a href="http://schema.org/Review">http://schema.org/Review</a> (4,408,551) <a href="http://schema.org/AggregateRating">http://schema.org/AggregateRating</a> (3,936,372)	2,994 MB	<a href="#">schemaorgHotel.nq.gz (sample)</a>
<a href="http://schema.org/JobPosting">http://schema.org/JobPosting</a>	Quads: 234,475,135 URLs: 2,011,332 Hosts: 3,962	<a href="http://schema.org/JobPosting">http://schema.org/JobPosting</a> (22,804,279) <a href="http://schema.org/Place">http://schema.org/Place</a> (16,321,339) <a href="http://schema.org/Organization">http://schema.org/Organization</a> (12,164,867) <a href="http://schema.org/Postaladdress">http://schema.org/Postaladdress</a> (7,516,387)	5,078 MB	<a href="#">schemaorgJobPosting.nq.gz (sample)</a>
<a href="http://schema.org/PostalAddress">http://schema.org/PostalAddress</a>	Quads: 776,573,609 URLs: 13,475,055 Hosts: 131,064	<a href="http://schema.org/PostalAddress">http://schema.org/PostalAddress</a> (48,086,763) <a href="http://schema.org/LocalBusiness">http://schema.org/LocalBusiness</a> (16,641,260) <a href="http://schema.org/GeoCoordinates">http://schema.org/GeoCoordinates</a> (12,345,942) <a href="http://schema.org/Place">http://schema.org/Place</a> (9,071,774)	14,364 MB	<a href="#">schemaorgPostalAddress.nq.gz (sample)</a>
<a href="http://schema.org/Product">http://schema.org/Product</a>	Quads: 2,829,523,589 URLs: 48,314,143 Hosts: 104,118	<a href="http://schema.org/Product">http://schema.org/Product</a> (287,815,069) <a href="http://schema.org/Offer">http://schema.org/Offer</a> (221,781,710) <a href="http://schema.org/AggregateRating">http://schema.org/AggregateRating</a> (38,398,548) <a href="http://schema.org/Review">http://schema.org/Review</a> (26,209,678)	62,179 MB	<a href="#">schemaorgProduct.nq.gz (sample)</a>

- Only tip of the iceberg, as each website is only partly crawled.

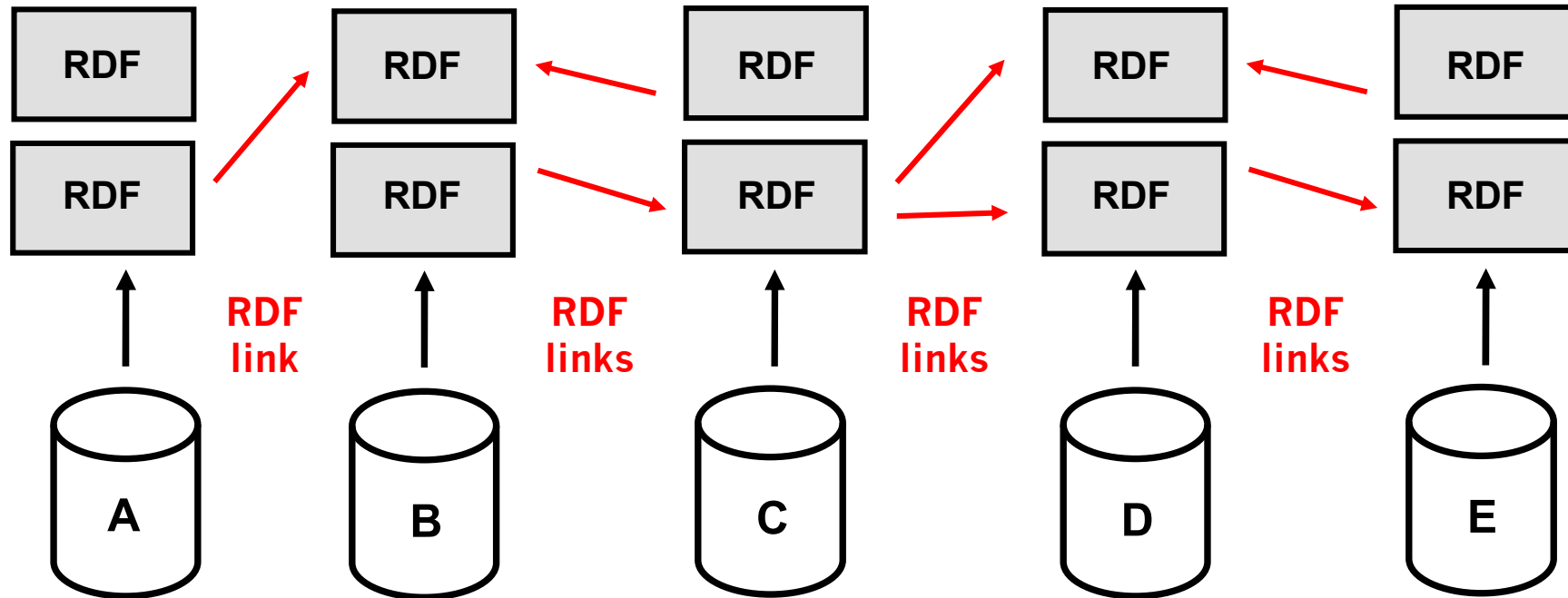


# 3. Alternative Approach: Linked Data

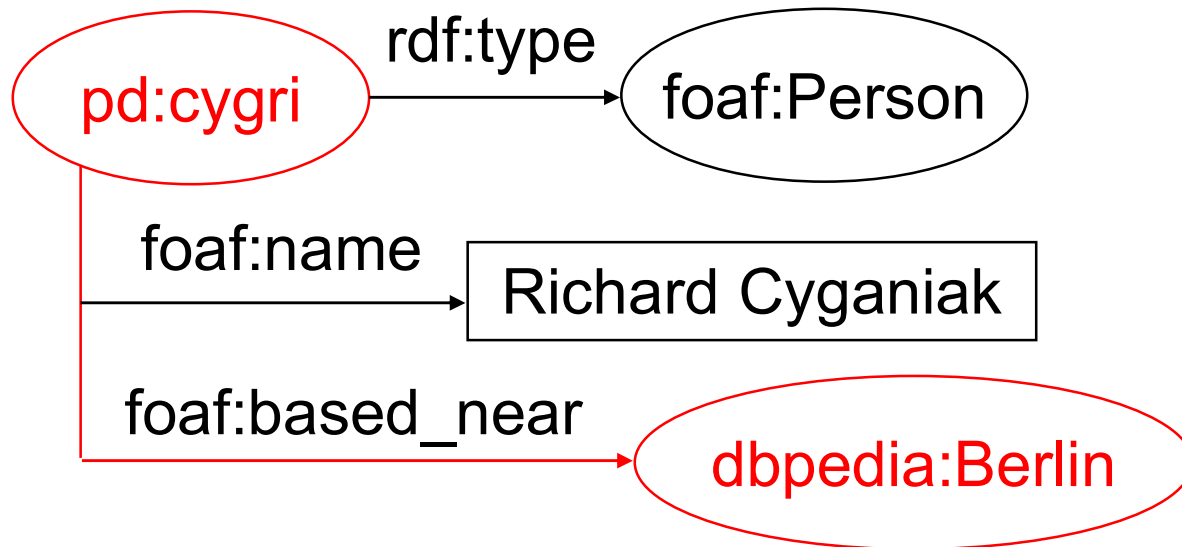


## Extend the Web with a single global data graph

- by using RDF to publish structured data on the Web
- by setting links between data items within different data sources



# Entities are identified with HTTP URIs

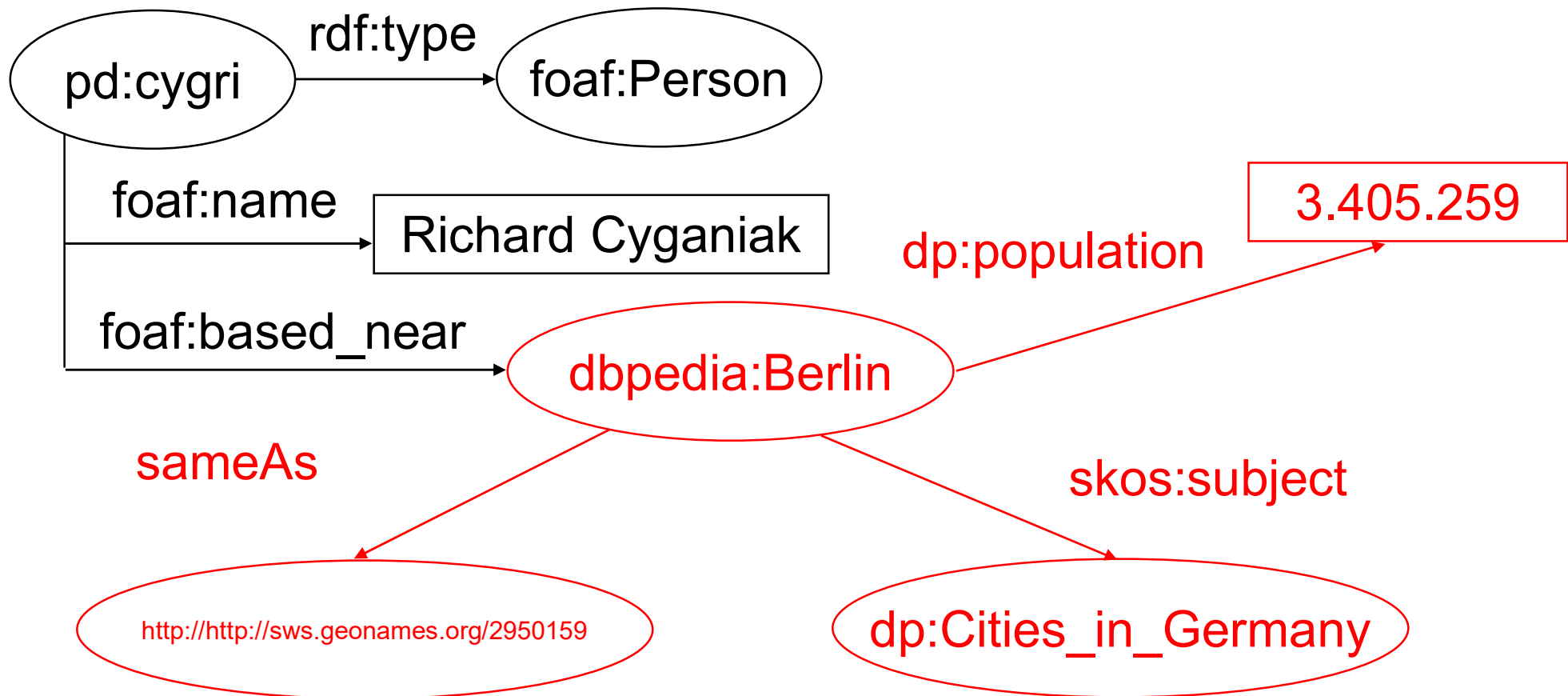


HTTP URIs take the role of global primary keys.

`pd:cygri` = <http://richard.cyganiak.de/foaf.rdf#cygri>

`dbpedia:Berlin` = <http://dbpedia.org/resource/Berlin>

# URIs can be linked to navigate the data space



## ■ By following RDF links applications can

- navigate the global data graph
- discover new data sources

# The Linked Open Data Cloud

**1,239** datasets connected by  
**16,147** sets of RDF links  
(as of March 2019)

## Legend

Cross Domain

Geography

Government

Life Sciences

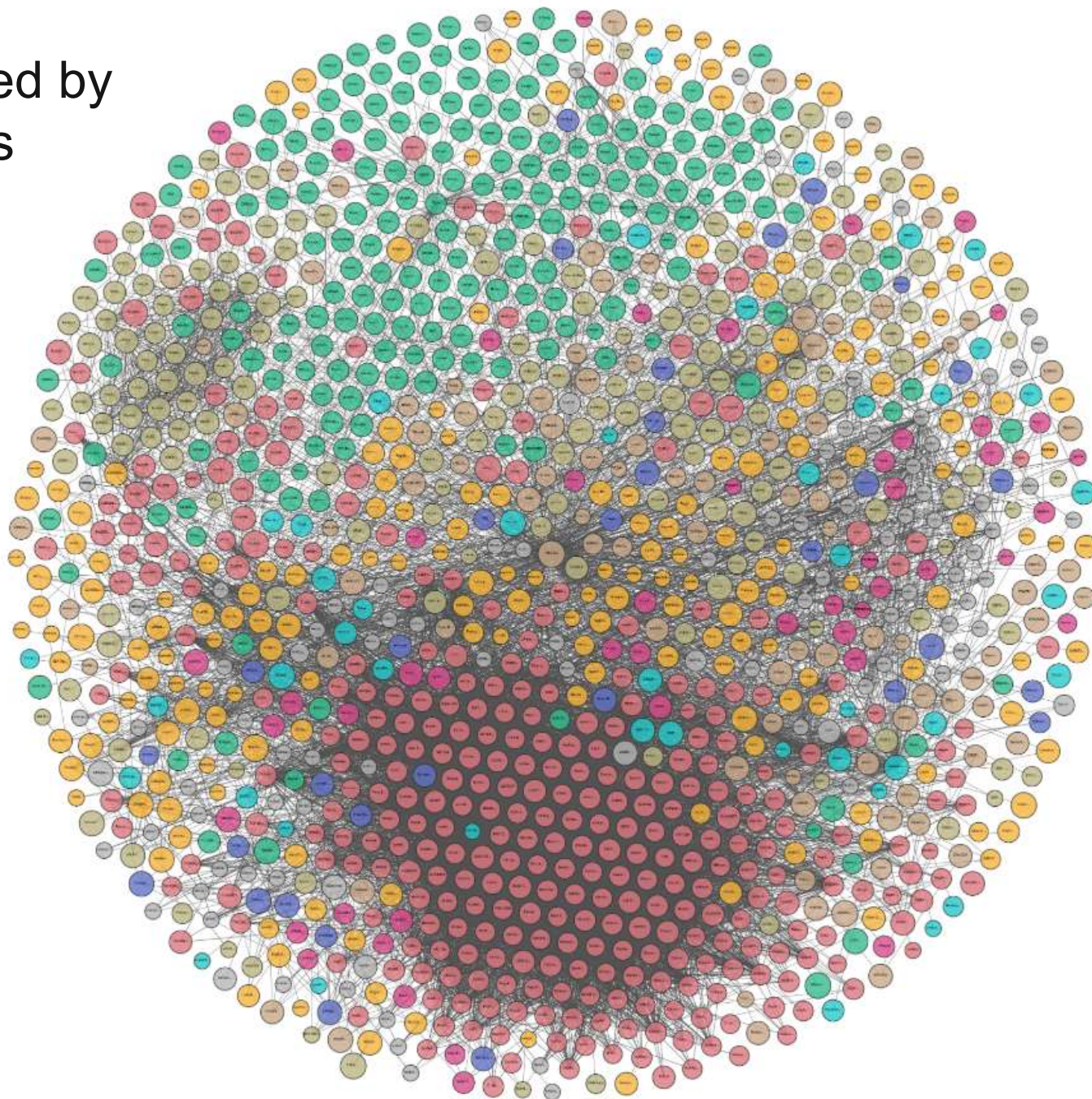
Linguistics

Media

Publications

Social Networking

User Generated



<https://lod-cloud.net/>



# Uptake in the Libraries Community

## ■ Institutions publishing Linked Data

- Library of Congress (subject headings and catalog)
- German National Library (PND dataset and subject headings)
- Swedish National Library (Libris catalog)
- Hungarian National Library (OPAC and digital library)
- Europeana Digital Library (catalog)
- Springer Nature (publications, researchers, projects)

## ■ Goals:

1. Interconnect resources between repositories (by topic, by location, by historical period, by ...)
2. Integrate library catalogs on global scale



LIBRARY OF  
CONGRESS



National Library  
of Sweden



europeana  
think culture



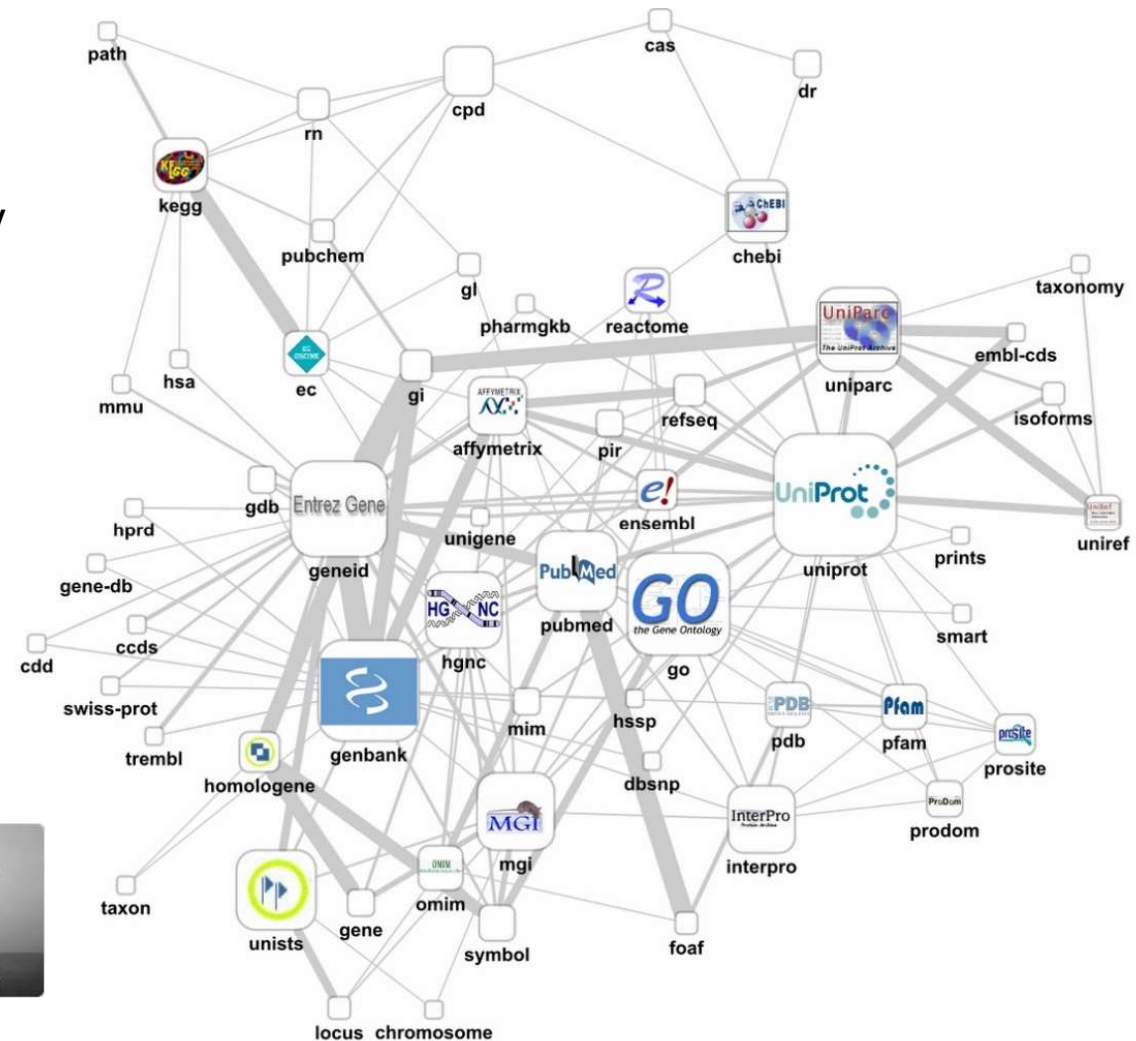
# Uptake in the Life Science Domain

## ■ Goals:

1. Connect life science datasets in order to support
  - biological knowledge discovery
  - drug development
2. Reuse results of previous integration efforts

## ■ Projects

**BIO** ↔ **RDF**



## 2. Web 2.0 Applications and Web APIs

- A multitude of **Web-based applications** has sprung up which enable users to share information.
- These applications
  - collect large **amounts of data** using proprietary schemata.
  - form **separate data spaces** that are only partly accessible from the Web via:
    1. HTML interfaces
    2. Web APIs



# Example: Facebook

## ■ Users (September 2018)

- 2,3 billion monthly active users
- including 1 billion mobile users

## ■ 740 billion friend connections

## ■ 4 million likes every minute

## ■ 250 billion photos uploaded

## ■ Data Volume

- 4 Petabyte of new data generated every day
- over 300 Petabyte in Facebook's data warehouse



<https://www.brandwatch.com/blog/facebook-statistics/>

<http://www.technologyreview.com/featuredstory/428150/what-facebook-knows/>



# Web APIs

- allow to access the data programmatically
- example of using the Twitter API:

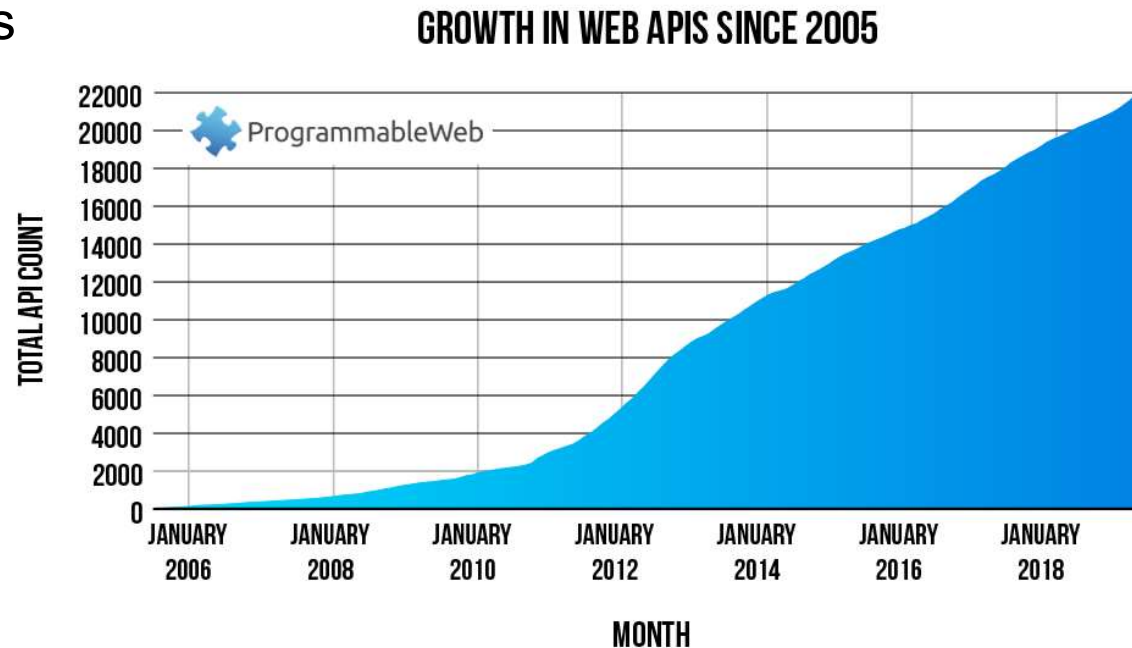
```
import pandas as pd

# Search tweets
dict_ = {'user': [], 'date': [], 'text': [], 'favorite_count': []}
for status in python_tweets.search(**query)['statuses']:
    dict_['user'].append(status['user']['screen_name'])
    dict_['date'].append(status['created_at'])
    dict_['text'].append(status['text'])
    dict_['favorite_count'].append(status['favorite_count'])

# Structure data in a pandas DataFrame for easier manipulation
df = pd.DataFrame(dict_)
df.sort_values(by='favorite_count', inplace=True, ascending=False)
df.head(5)
```

## ■ ProgrammableWeb API Catalog

- lists over 24,000 Web APIs
- lists over 6,800 mashups
- catalog maintained until 10-2022 (alternatives below\*)

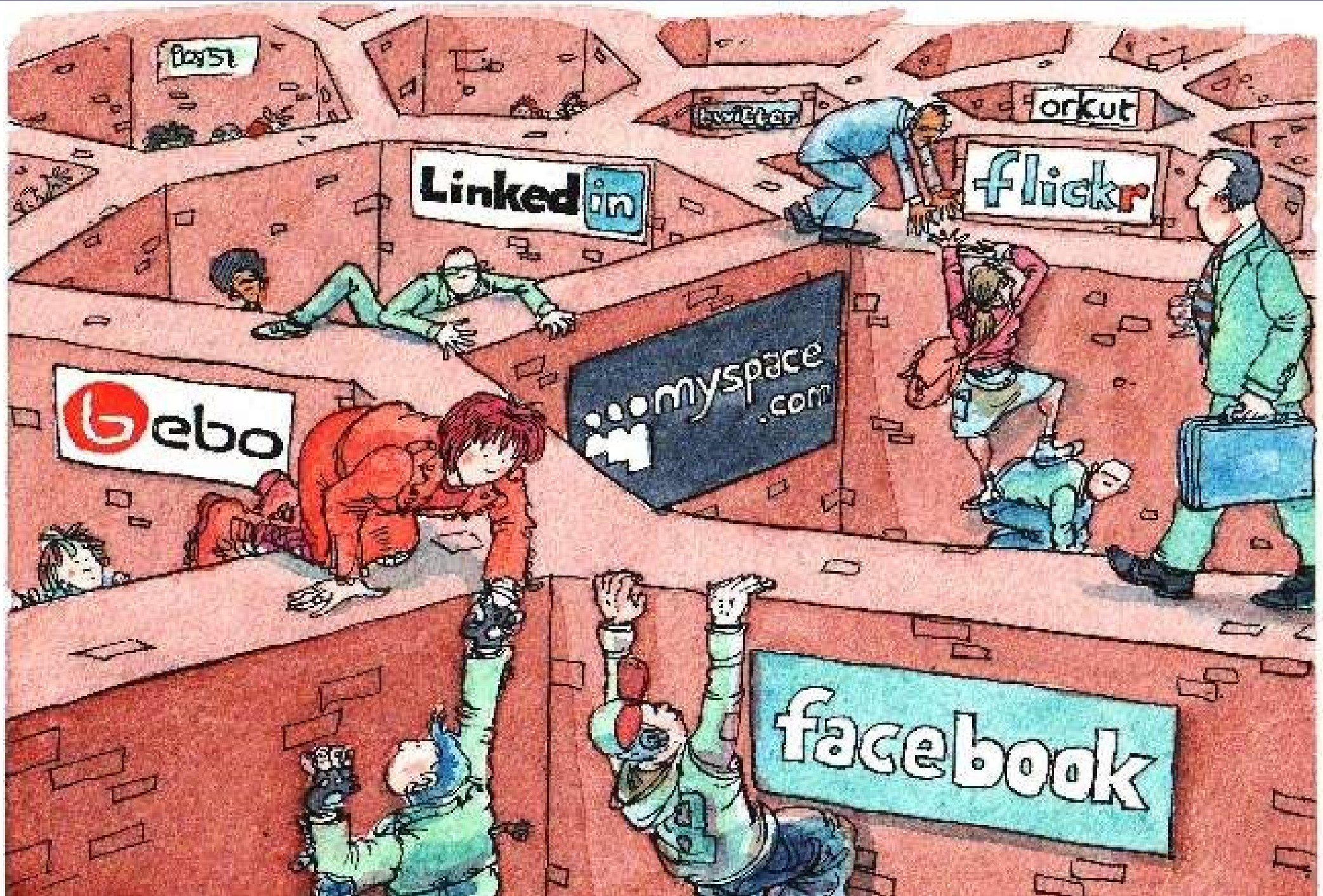


## ■ APIs usually provide only limited access

- restricted to specific queries (canned queries)
- restricted by amount of queries / number of results
- (try to) prevent crawling

\* <https://nordicapis.com/13-api-directories-to-help-you-discover-apis/>

# Web APIs versus Linked Data or HTML Embedded Data



# 3. What is Web Mining?

## ■ Definition

**Non-trivial extraction of implicit, previously unknown and potentially useful information from**

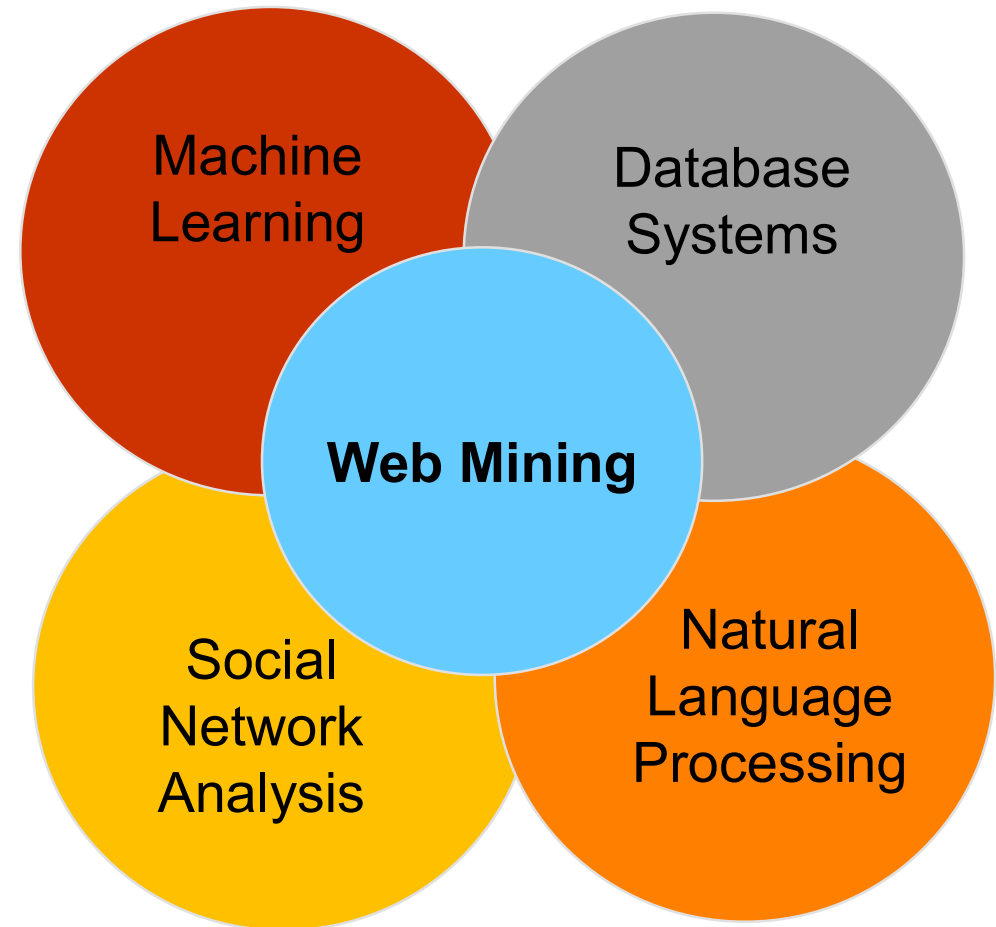
- **Web content,**
- **Web structure**
- **Web usage data.**

## ■ Recurring Challenges

1. huge amount of available data → requires sampling or cloud computing
2. semi-structured nature of data → mix of data and text mining techniques
3. heterogeneity of data → data integration and cleansing is a challenge
4. distributed nature of data → often requires large-scale crawling or relying on pre-crawled web corpora

# Web Mining is a Multi-Disciplinary Field

■ Draws ideas and techniques from



■ **Sub-Fields**

1. **Web Usage Mining**
2. **Web Structure Mining**
3. **Web Content Mining**

# 3.1 Web Usage Mining

## ■ Definition

**Discovery of patterns in click-streams and associated data collected or generated as a result of user interactions with one or more web sites or Web 2.0 applications.**

## ■ Typical Sources of Data

1. web server access logs
2. e-commerce and product-oriented user events (e.g., shopping cart changes, ad or product click-throughs, purchases)
3. user events on social network sites (e.g., likes, posts, comments)



## ■ Associated Data

1. page attributes, page content, site structure
2. additional domain knowledge and demographic data
3. user profiles or user ratings



# Web Usage Data: An Endless Sea

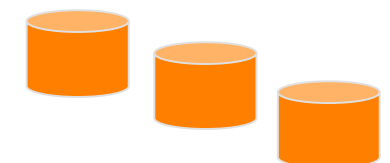
## 2021 This Is What Happens In An Internet Minute



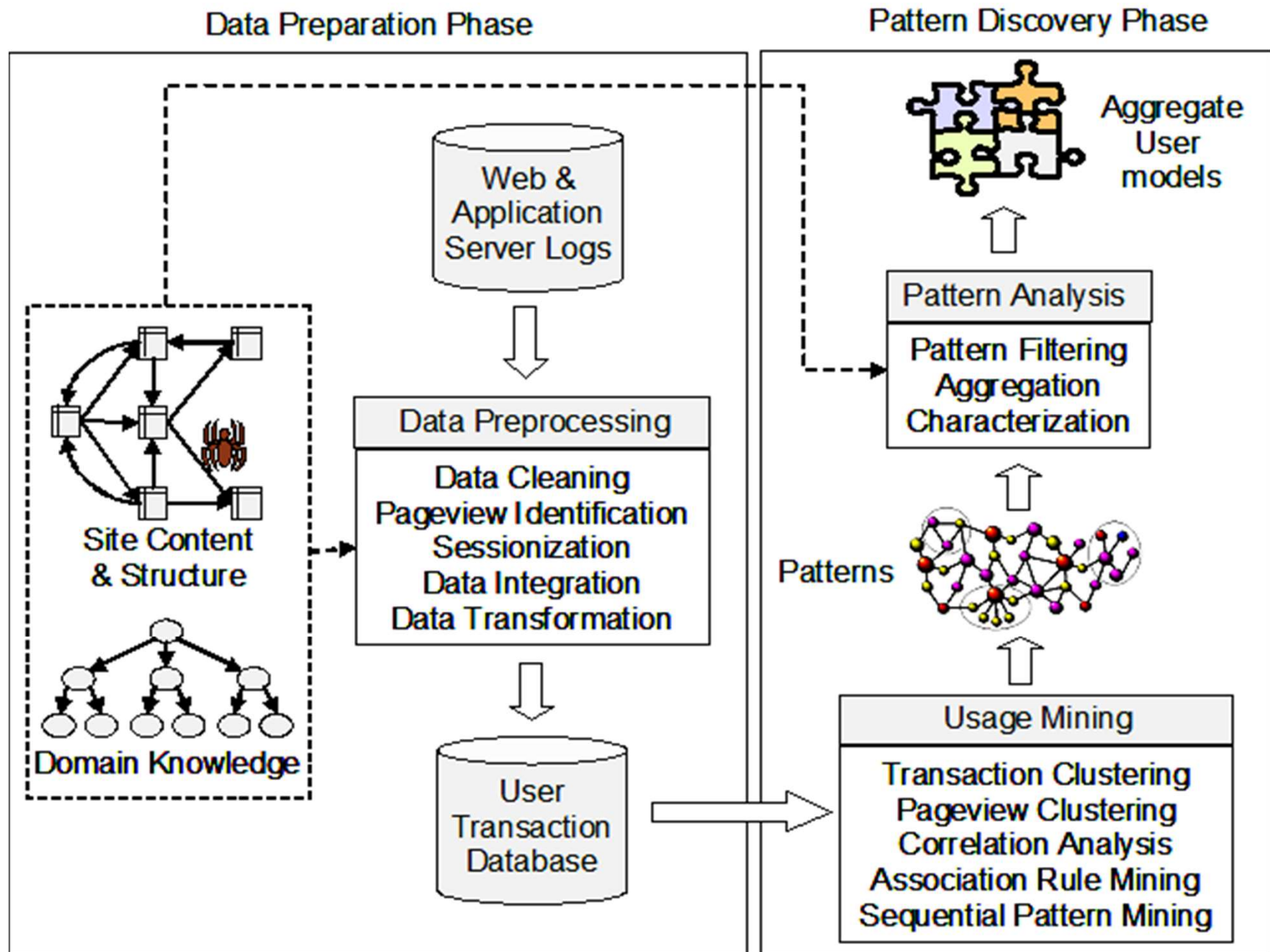
Get Traffic Analysis



Provide Access to



# The Web Usage Mining Process





# Example Application: Product Recommendation

amazon.de   DE Hallo, Simone [Konto und Listen](#) [Warenrücksendungen und Bestellungen](#) Mein Prime

Lieferung an Simone 69115 Heidelberg Erneut kaufen Simones Amazon Angebote Gutscheine Verkaufen Hilfe Küche, Haushalt & Wohnen

Mein Amazon Ihre besuchten Seiten **Ihre Empfehlungen** Verbessern Sie Ihre Empfehlungen Gutscheine Mein öffentliches Profil Mehr dazu

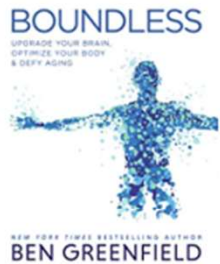
## Für Sie empfohlen

Simones Amazon >

### Bücher

Warum empfohlen?

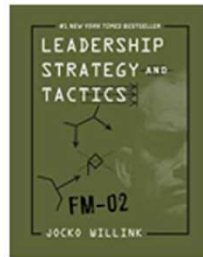
Artikel entfernen



**Boundless: Upgrade Your Brain, Optimize Your Body & Defy Aging**  
Ben Greenfield  
★★★★★ 4  
66,99 € ✓prime

Ähnliches

In den Einkaufswagen



**Leadership Strategy and Tactics: Field Manual**  
Jocko Willink  
★★★★★ 5  
18,99 € ✓prime

Ähnliches

In den Einkaufswagen



**Running Rewired: Reinvent Your Run for Stability, Strength, and...**  
Jay Dicharry  
★★★★★ 6  
23,59 € ✓prime

Ähnliches

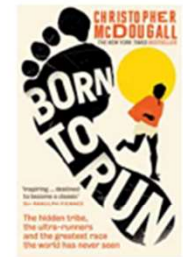
In den Einkaufswagen



**Training for the Uphill Athlete: A Manual for Mountain Runners and...**  
Steve House  
★★★★★ 8  
22,99 € ✓prime

Ähnliches

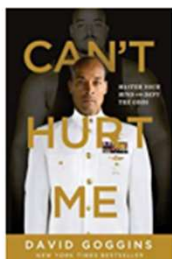
In den Einkaufswagen



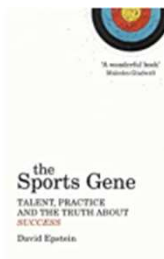
**Born to Run: The Hidden Tribe, the Ultra-Runners, and the Greatest...**  
Christopher McDougall  
★★★★★ 64  
10,29 € ✓prime

Ähnliches

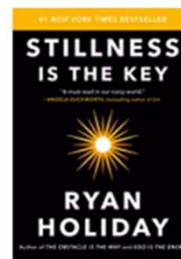
In den Einkaufswagen



**Can't Hurt Me: Master Your Mind and Defy the Odds**  
David Goggins



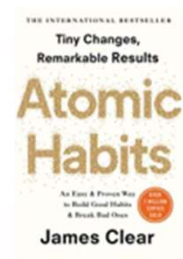
**The Sports Gene: Talent, Practice and the Truth About Success**  
David Epstein



**Stillness Is the Key**  
Ryan Holiday  
★★★★★ 31



**Das große Kettlebell-Trainingsbuch**  
Dr. Till Sukopp  
★★★★★ 60



**Atomic Habits: The life-changing million copy bestseller**  
James Clear

# Example Application: Personalized Search



kafka



All Images Videos Maps News More Settings Tools

About 51.700.000 results (0,58 seconds)

Ad · www.confluent.io/ ▾

## Confluent | Download Apache Kafka® Today | confluent.io

Monitor & Manage Data in Real-Time Using Our Scalable, Reliable, & Flexible Platform. Confluent Platform Includes the Latest Version of Apache Kafka, Plus Enterprise Features. Reduce Ops Burden. Streaming Data Service. Deployable in Minutes. 24/7 Support.

### What is Kafka?

More Info. About Kafka & Confluent Platform.

### Kafka Definitive Guide

Learn All About Kafka From its Original Developers in this eBook.

kafka.apache.org ▾

## Apache Kafka

Kafka® is used for building real-time data pipelines and streaming apps. It is horizontally scalable, fault-tolerant, wicked fast, and runs in production in thousands ...

[Introduction](#) · [Kafka Streams](#) · [Documentation](#) · [Apache Kafka Documentation](#)

kafka.apache.org > intro ▾

## Introduction - Apache Kafka - Apache Software

In Kafka the communication between the clients and the servers is done with a simple, high-performance, language agnostic TCP protocol. This protocol is ...


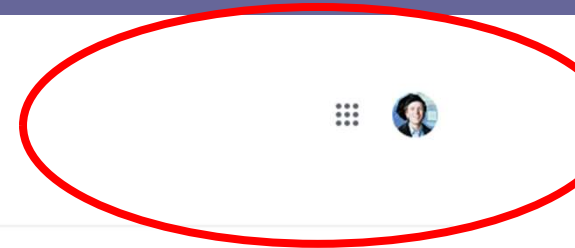
### People also ask

What is Kafka used for? ▾

What is meant by Kafka? ▾

What is Kafka and how it works? ▾

Is Kafka free? ▾



## Apache Kafka

Software

Apache Kafka is an open-source stream-processing software platform developed by LinkedIn and donated to the Apache Software Foundation, written in Scala and Java. The project aims to provide a unified, high-throughput, low-latency platform for handling real-time data feeds. [Wikipedia](#)

**License:** [Apache License 2.0](#)


**Developer(s):** [Apache Software Foundation](#)

**Initial release date:** January 2011


**Stable release:** 2.4.0 / [December 16, 2019](#); 48 days ago

**Written in:** [Scala](#), [Java](#)


People also search for View 15+ more




Apache Spark




RabbitMQ



Apache ZooKeeper



Apache Cassandra



Apache Hadoop

# Example Application: Personalized Search



kafka



Anmelden

Alle Bilder Videos News Bücher Mehr Einstellungen Suchfilter

Ungefähr 50.100.000 Ergebnisse (0,66 Sekunden)



Hinweise zum Datenschutz bei Google

SPÄTER ERINNERN

ANSEHEN

de.wikipedia.org › wiki › Franz\_Kafka

## Franz Kafka – Wikipedia

Franz Kafka (tschechisch gelegentlich František Kafka, jüdischer Name: אנשיל Anschel; \* 3. Juli 1883 in Prag, Österreich-Ungarn; † 3. Juni 1924 in Kierling, ...

Hermann Kafka · Kategorie:Franz Kafka · Kafkaesk · Die Verwandlung

kafka.apache.org › Diese Seite übersetzen

## Apache Kafka

Kafka® is used for building real-time data pipelines and streaming apps. It is horizontally scalable, fault-tolerant, wicked fast, and runs in production in thousands ...

Introduction · Kafka Streams · Apache Kafka Documentation · Documentation

www.franzkafka.de

## Franz Kafka

Dieses Portal zu Franz Kafka bietet zuverlässig und kurzweilig die wichtigsten Informationen über sein Werk.

www.franzkafka.de › franzkafka › das\_leben

## Das Leben - Franz Kafka

Juli wird Franz Kafka in Prag geboren. Er ist das erste Kind von Hermann Kafka ( 1852-1931) und seiner Frau Julie, geb. Löwy (1856-1934). Die jüdischen Eltern ...

www.inhaltsangabe.de › autoren › kafka

## Franz Kafka - Biografie und Inhaltsangaben - Inhaltsangabe.de



## Franz Kafka

Schriftsteller

Franz Kafka war ein deutschsprachiger Schriftsteller. Sein Hauptwerk bilden neben drei Romanfragmenten zahlreiche Erzählungen. [Wikipedia](#)

**Geboren:** 3. Juli 1883, Prag, Tschechien

**Gestorben:** 3. Juni 1924, Kierling, Klosterneuburg, Österreich

**Bestattet:** 11. Juni 1924, Neuer jüdischer Friedhof, Prag, Tschechien

**Kurzgeschichten:** Die Verwandlung, Das Urteil, Vor dem Gesetz, MEHR

**Beeinflusst von:** Fjodor Michailowitsch Dostojewski, MEHR

## Bücher

Über 45 weitere ansehen



Die



Das



Brief an



Das Urteil



Briefe an

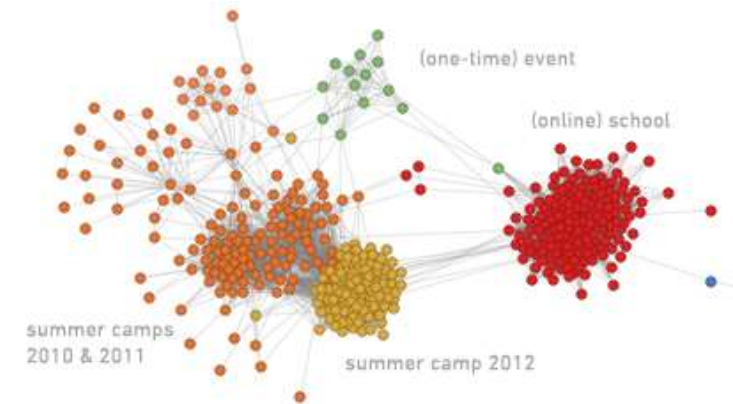


# 3.2 Web Structure Mining

## ■ Definition

**Discovery and interpretation of patterns in**

- 1. the hyperlink structure of the Web**
- 2. the social ties among actors that interact on the Web**



## ■ Typical sources of Web graphs

1. Web crawls including HTML pages and hyperlinks
2. social networks including explicit relations between actors (your Facebook friend network)
3. other types of community data (discussion forums, email conversations, ...)

## ■ Focuses on the structure, but can of course also be combined with content or usage mining techniques

# Identification of Prominent Nodes

Question: Who are the “most important” actors in a social network?

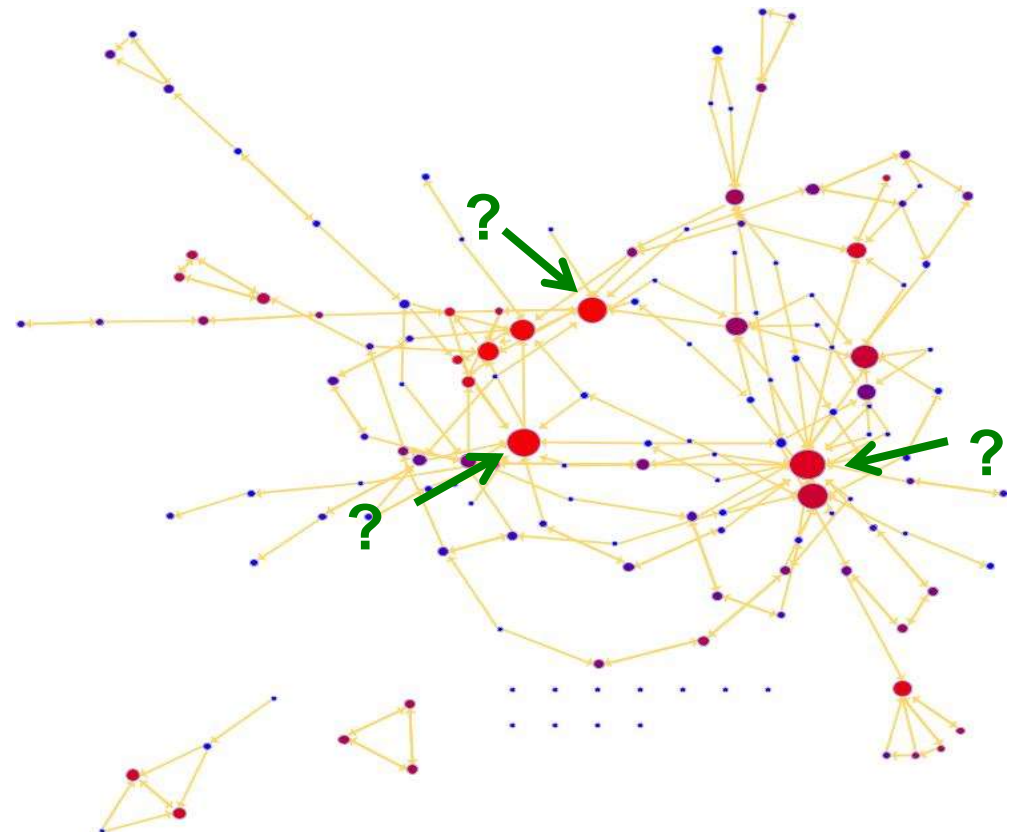


## Centrality

- A **central actor** is one involved in many edges.
- The direction of lines is not considered.

## Prestige

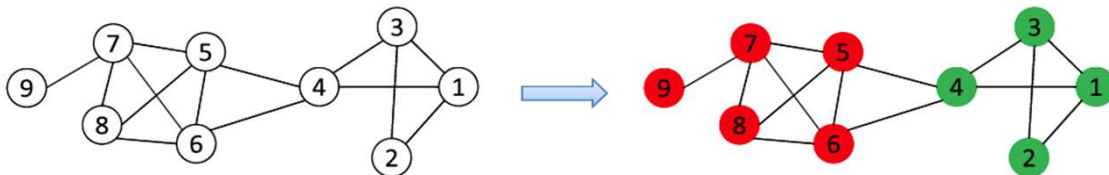
- A **prestigious actor** is one who is the target of many arcs.
- The direction of arcs is considered.



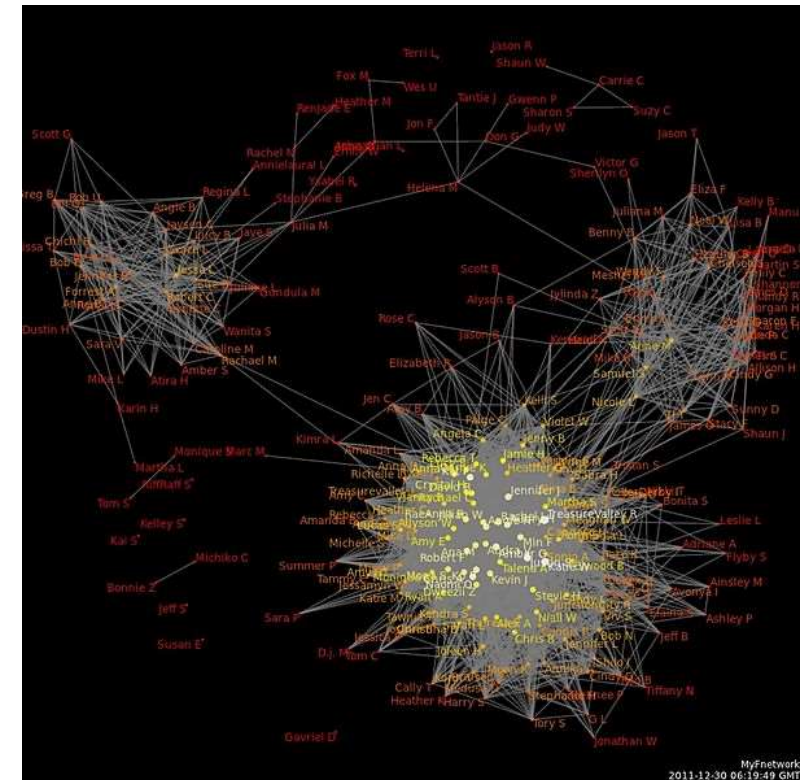
# Community Detection

A **community** is a set of actors between which interactions are (relatively) frequent.

- Finding a community in a social network is to identify a set of nodes such that they interact with each other more frequently than with those nodes outside the group.

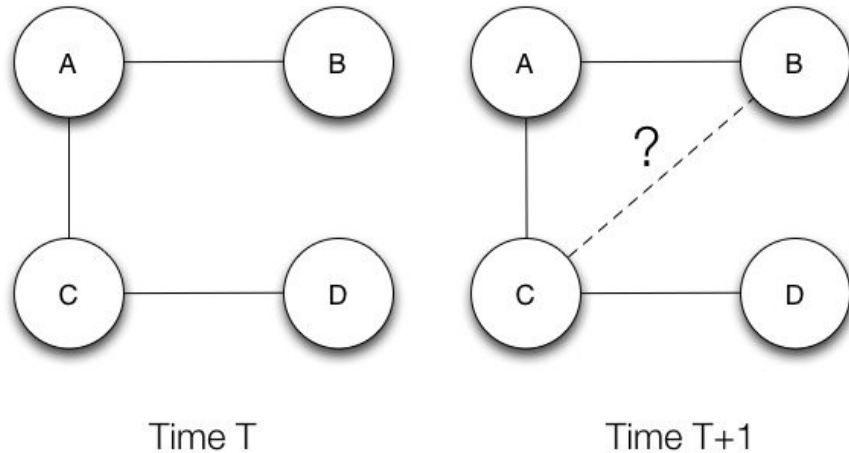


- Methods: Components, K-Cores, Islands, ...
- Applications: Recommendation based on communities, visualization of huge networks



# Machine Learning with Graphs

**Link Prediction:** Given a snapshot of a social network, can we infer which new interactions among its members are likely to occur in the near future? (Liben-Nowell & Kleinberg, 2007)



- Facebook: recommending possible friends
- Tinder: recommending potential matches

**Node Classification:** Predict node properties based on other node properties as well as neighboring nodes.

- Knowledge Graph: Type of a node: e.g., person, author, athlete



# 3.3 Web Content Mining

## ■ Definition

**Automatic extraction of useful information (facts, patterns) from Web content (text, images, multimedia).**

## ■ Content Mining Tasks

- Content Clustering
- Content Classification
- Sentiment Analysis
- Political Scaling
- Information Extraction

# Content Classification

- **Goal: Previously unseen documents/images should be assigned a class as accurately as possible.**
- **Applications**
  - News categorization
  - Product categorization
  - Spam detection
- **Classification methods**
  - Naive Bayes, Support Vector Machines, Deep Neural Nets, Transformers

# Content Clustering

- **Unsupervised Learning: Given a set of documents and a similarity measure among documents find clusters such that:**
  - documents in one cluster are more similar to one another
  - documents in separate clusters are less similar to one another
- **Applications**
  - Search result clustering
  - Topic discovery
- **Techniques**
  - Algorithms: K-Means, Hierarchical Clustering
  - Similarity measures: Cosine, Jaccard

# Mixture of Document Clustering and Classification

The screenshot shows the Google News interface. At the top, there is a search bar and navigation tabs for 'Headlines', 'Local', 'For You', and 'U.S.'. Below the navigation is a 'SECTIONS' sidebar on the left, which is highlighted with a red dashed box. The sidebar lists various news categories: Top Stories, World, U.S., Business, Technology, Entertainment, Sports, Science, Health, and a 'Manage sections' option. The main content area is titled 'Top Stories' and features a large article about Trump's budget. This article is also highlighted with a red dashed box. The article includes a photo of food cans, a headline, a source (Washington Post), and several related articles with their respective sources and timestamps. A 'View full coverage' link is at the bottom of the article.

Google News


Search

Headlines Local For You U.S. ▾

SECTIONS

- Top Stories
- World
- U.S.
- Business
- Technology
- Entertainment
- Sports
- Science
- Health
- Manage sections

## Top Stories



### Trump's budget hits poor Americans the hardest

Washington Post · 2h ago

RELATED COVERAGE

Legislative Outline for Rebuilding Infrastructure in America - WhiteHouse.gov  
**Most Referenced** · WhiteHouse.gov · 6h ago

MORE ABOUT

- Donald Trump
- White House
- President of the United States
- United States of America

Trump plan calls for \$1.5 trillion in 'investment' to fix nation's infrastructure  
ABC News · 1h ago

Trump budget wants billions more for border wall, immigration agents and judges  
USA TODAY · 1h ago

budget - WhiteHouse.gov  
**Most Referenced** · WhiteHouse.gov · 5m ago

Trump Budget Calls for Work Requirements for Housing Aid  
**Featured** · Governing · 1h ago

Is Trump Working On a Different Kind of 'Massacre'?  
**Opinion** · New York Times · 1h ago

View full coverage →

# Sentiment Analysis

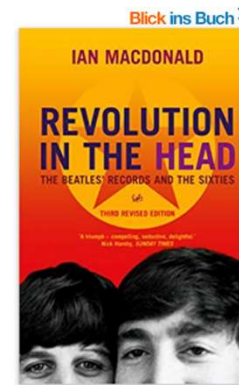
The basic task in sentiment analysis is classifying the polarity of a given text at the document, sentence, or feature/aspect level.

## ■ Polarity Values

- positive, neutral, negative
- stars

## ■ Applications

- Document-level: poll prediction from tweets
- Feature/Aspect-level: analysis of product reviews



**Revolution In The Head: The Beatles Records and the Sixties** (Englisch) Taschenbuch – 4. Dezember 2008

von Ian MacDonald (Autor)

★★★★☆ 175 Sternebewertungen

> Alle 3 Formate und Ausgaben anzeigen

Taschenbuch  
5,69 €

Lieferung 19.-24. Febr., wenn Sie Standardversand an der Kasse wählen. [Siehe Details.](#)

18 neu ab 5,69 € | 5 gebraucht ab 3,60 €

As dazzling as the decade they dominated, The Beatles almost single-handedly created pop music as we know it. Today, their songs are cited as seminal influences by stars like Oasis and Blur. Eloquently giving voice to their time, The Beatles quite simply changed the world.



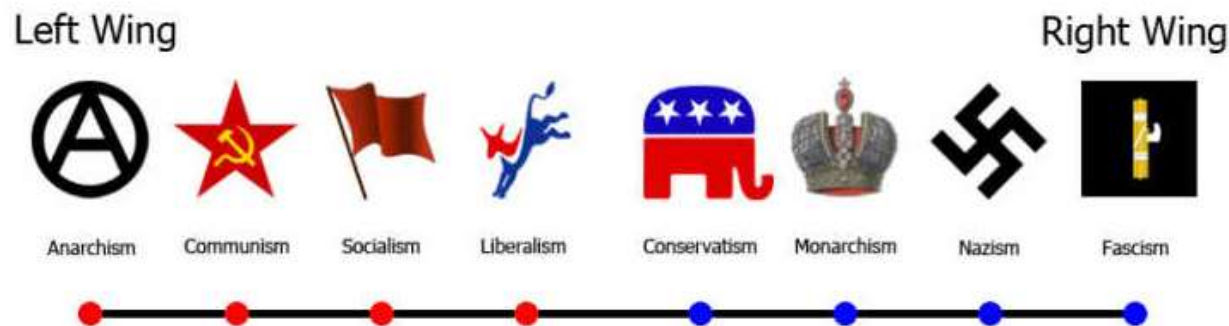
# Hate Speech Detection / Political Scaling

## ■ Hate Speech Detection in Social Media Content

Twitter	% posts	Whisper	% posts
I hate	70.5	I hate	66.4
I can't stand	7.7	I don't like	9.1
I don't like	7.2	I can't stand	7.4
I really hate	4.9	I really hate	3.1
I fucking hate	1.8	I fucking hate	3.0
I'm sick of	0.8	I'm sick of	1.4
I cannot stand	0.7	I'm so sick of	1.0
I fuckin hate	0.6	I just hate	0.9
I just hate	0.6	I really don't like	0.8
I'm so sick of	0.6	I secretly hate	0.7

Twitter		Whisper	
Hate target	% posts	Hate target	% posts
Nigga	31.11	Black people	10.10
White people	9.76	Fake people	9.77
Fake people	5.07	Fat people	8.46
Black people	4.91	Stupid people	7.84
Stupid people	2.62	Gay people	7.06
Rude people	2.60	White people	5.62
Negative people	2.53	Racist people	3.35
Ignorant people	2.13	Ignorant people	3.10
Nigger	1.84	Rude people	2.45
Ungrateful people	1.80	Old people	2.18

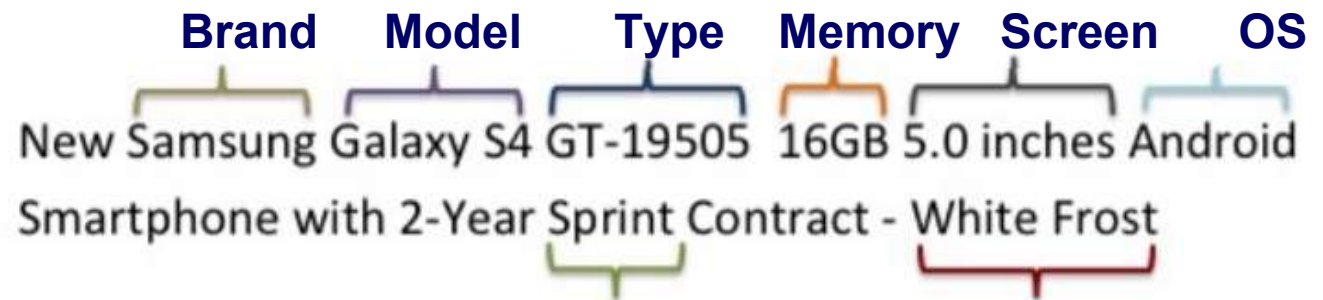
## ■ Political Scaling of Opinionated Texts



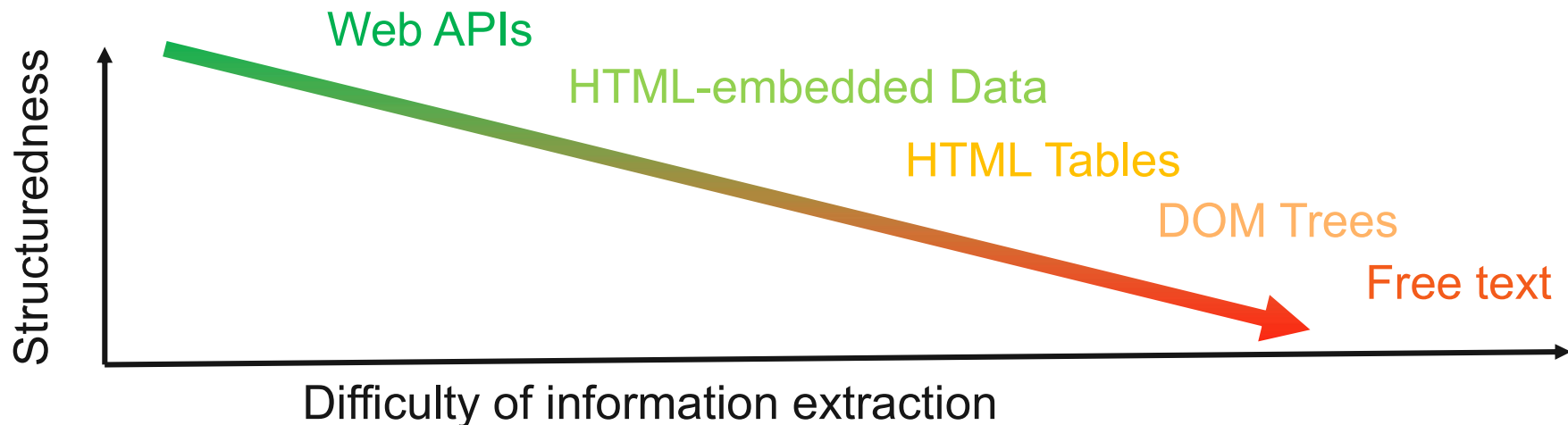
# Information Extraction

**Goal: Automatic extraction of structured information from unstructured or semi-structured Web content.**

- **Example of below 1NF data:**



- **The difficulty of the extraction depends on the structuredness**





# Examples: Information Extraction from the Web

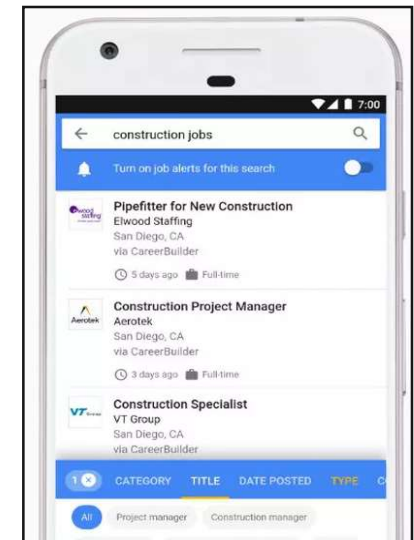
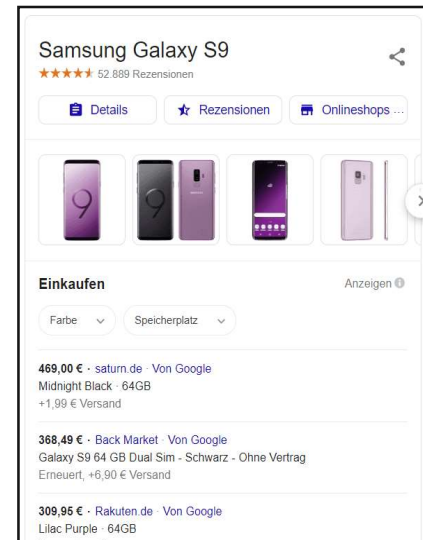


WIKIPEDIA  
The Free Encyclopedia

schema.org

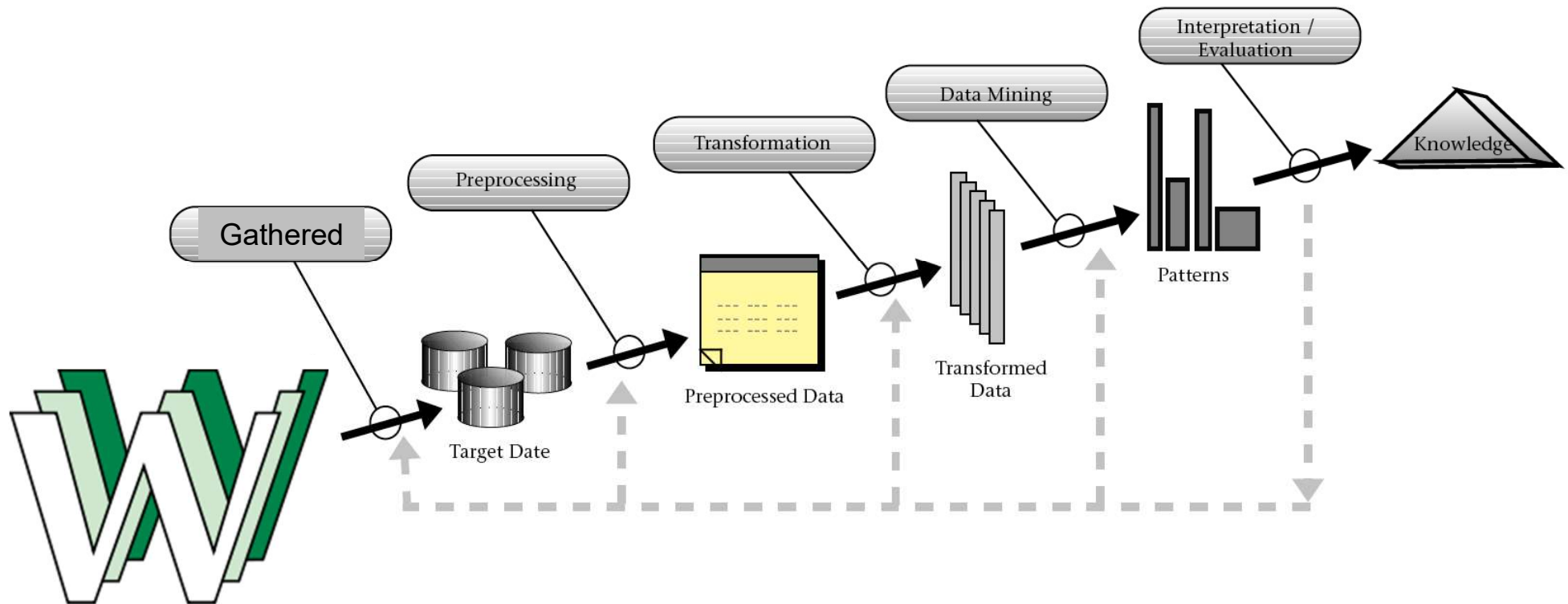


UNIVERSITÄT LEIPZIG



## 3.4. The Web Mining Process

Equal to the standard data mining process with the difference that data is gathered from the Web.



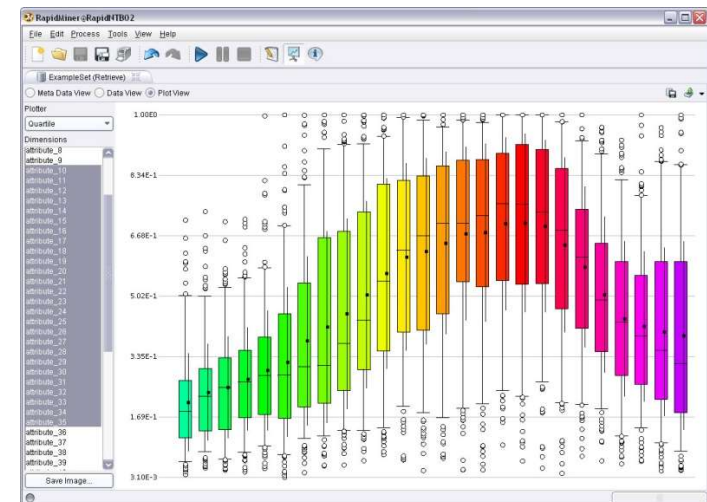
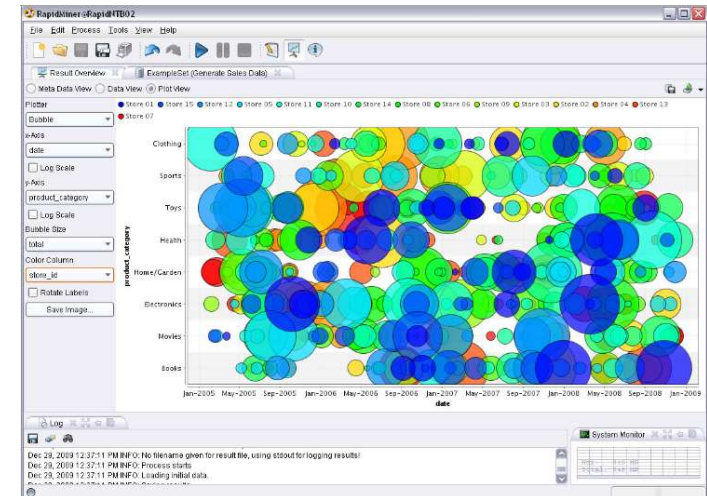
# Gathering and Exploration

## ■ Gathering of Web Data

- Crawl documents or data
- Retrieve data via Web API
- Download pre-gathered data sets

## ■ Exploration

- Get an initial understanding of the data
- Calculate basic summarization statistics
- Visualize the data
- Identify data problems such as outliers, missing values, duplicate records



# Preprocessing and Transformation

- **Transform data into a representation that is suitable for the chosen data mining methods**
  - amount of data (determines hardware requirements)
  - number of dimensions (represent relevant information using less attributes)
  - scales of attributes (nominal, ordinal, numeric)
- **Methods**
  - discretization and binarization
  - feature subset selection / dimensionality reduction
  - attribute transformation / text to term vector / embeddings
  - aggregation, sampling
  - integrate data from multiple sources
- **Good data preparation is key to producing valid and reliable models**
- **Data integration and preparation is estimated to take 70-80% of the time and effort of a data mining project**

# Actual Data Mining

- Input: Preprocessed Data

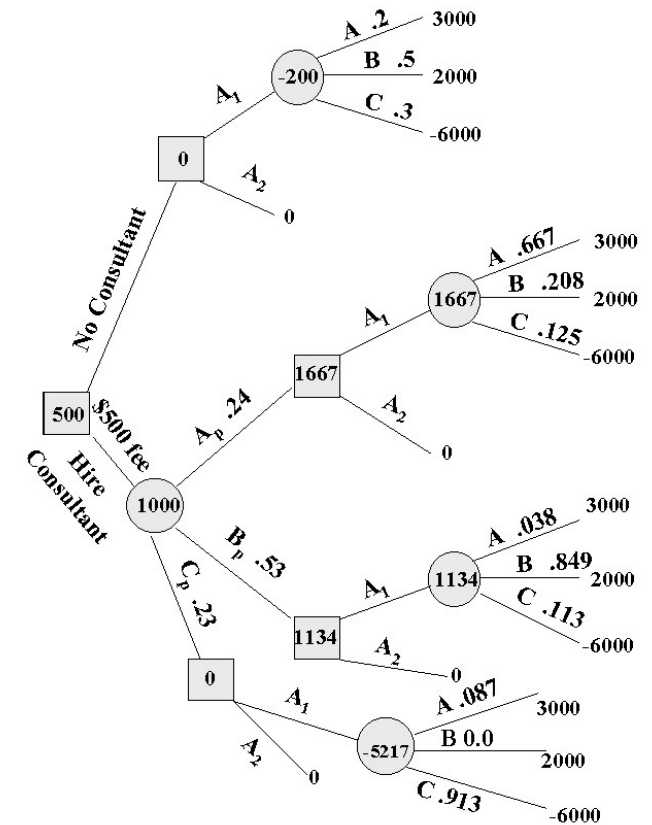
- Output: **Model / Patterns**

1. Apply data mining method

2. Evaluate resulting model / patterns

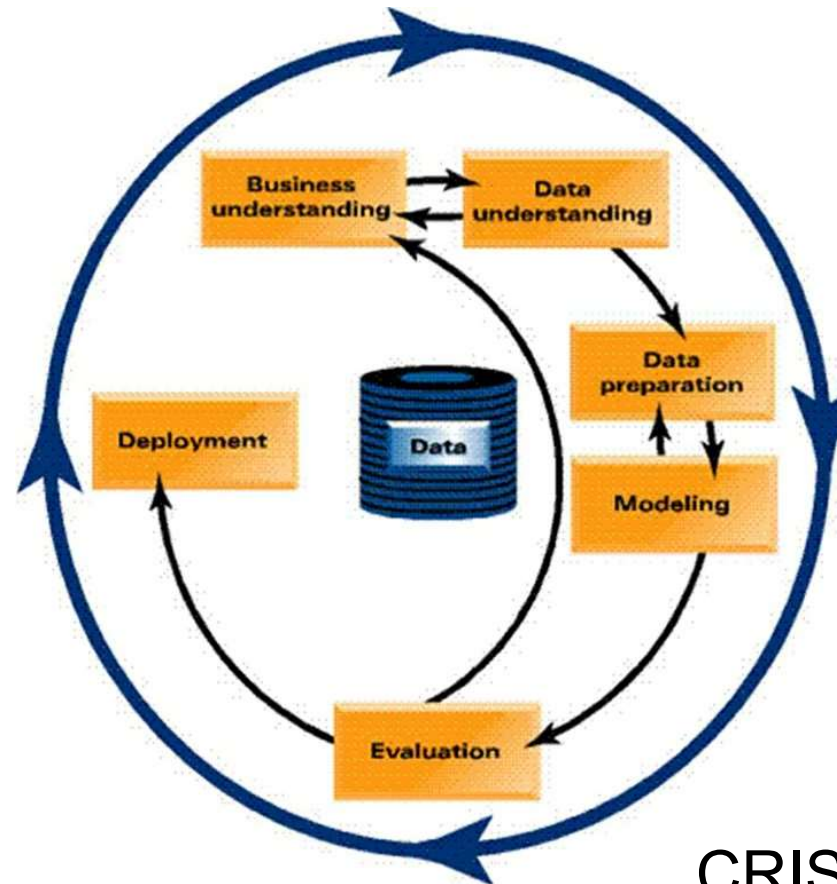
3. Iterate

- experiment with different parameter settings
- experiment with multiple alternative methods
- improve preprocessing and feature generation
- increase amount or quality of training data



# Deployment

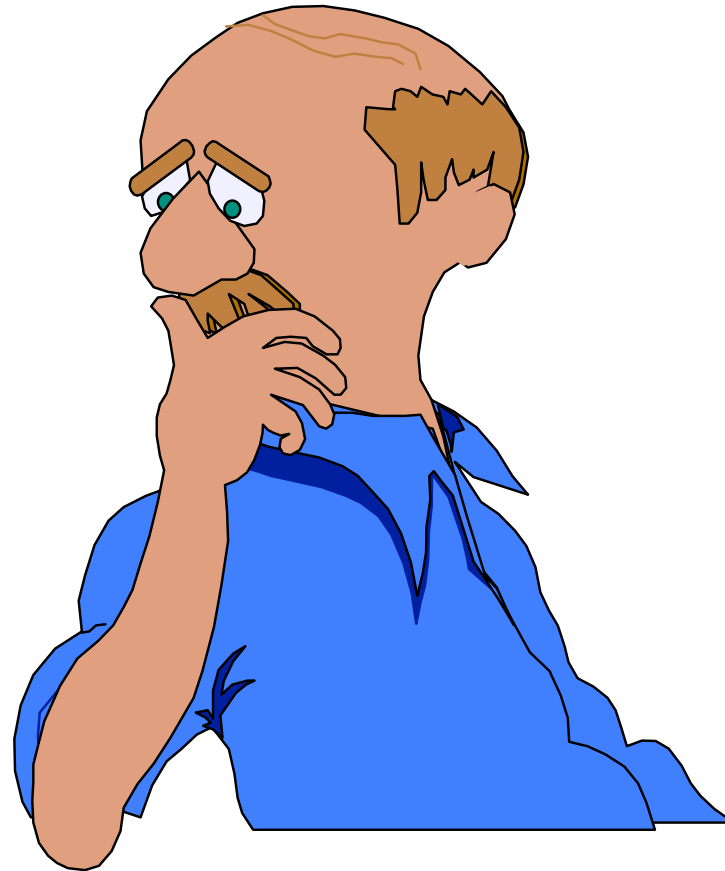
- Use model in the business context
- Keep iterating in order to maintain and improve model



CRISP-DM Process Model



# Questions?



- **This week: No lab!**
- Next week: Lecture and Lab: Web Usage Mining and Recommender Systems