# Récolte des Connaissances de la Toile

## Gerhard Weikum

**Max Planck Institute for Informatics**
**Saarland University**
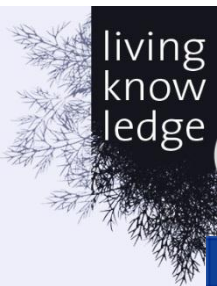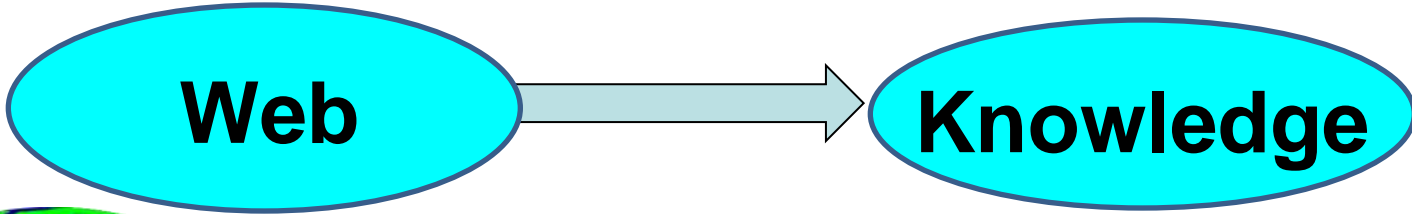**http://www.mpi-inf.mpg.de/~weikum/**

max planck institut
informatik

# Knowledge Harvesting From the Web

## Gerhard Weikum

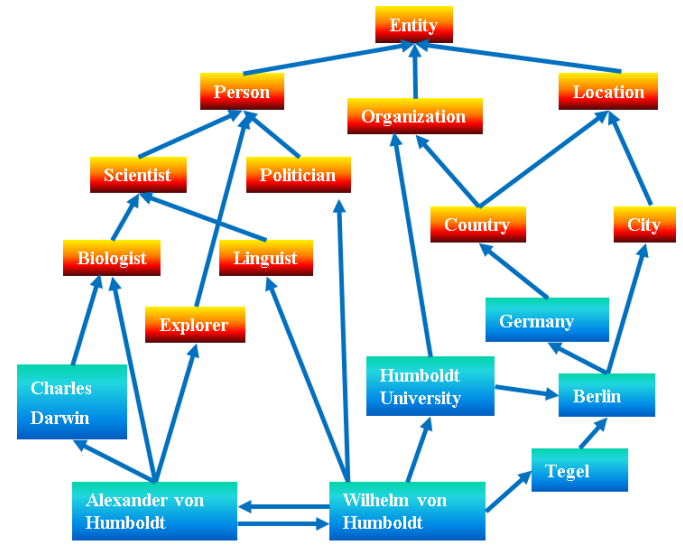**Max Planck Institute for Informatics**
**Saarland University**
**http://www.mpi-inf.mpg.de/~weikum/**

# Acknowledgements

# Knowledge from the Web



| Movie Title | Price | Rating |
|---|---|---|
| Sherlock Holmes 2 | 15.99 | 7.5 |
| Inglorious Basterds | 11.99 | 8.3 |
| OSS 117: Lost in Rio | 7.99 | 6.7 |
| | | |

# Knowledge from the Web

# Outline

★ **Knowledge for Machines**

★ **Construction of Knowledge Bases**

★ **KB Population from Text & Web Pages**

★ **Knowledge for Intelligent Applications**

★ **Opportunities & Challenges**

★ **Conclusions**

# Spectrum of Machine Knowledge

**factual:**
bornIn (SteveJobs, SanFrancisco), hasFounded (SteveJobs, Pixar),
hasWon (SteveJobs, NationalMedalOfTechnology), livedIn (SteveJobs, PaloAlto)

**taxonomic (ontology):**
instanceOf (SteveJobs, computerArchitects), instanceOf(SteveJobs, CEOs)
subclassOf (computerArchitects, engineers), subclassOf(CEOs, businesspeople)

**lexical (terminology):**
means ("Big Apple", NewYorkCity), means ("Apple", AppleComputerCorp)
means ("MS", Microsoft) , means ("MS", MultipleSclerosis)

**multi-lingual:**
meansInChinese („乔戈里峰", K2), meansInUrdu („کے ٹو", K2)
meansInFr („école", school (institution)), meansInFr („banc", school (of fish))
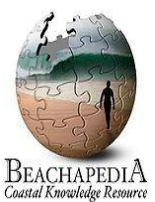
**temporal (fluents):**
hasWon (SteveJobs, NationalMedalOfTechnology)@1985
marriedTo (AlbertEinstein, MilevaMaric)@[6-Jan-1903, 14-Feb-1919]
presidentOf (NicolasSarkozy, France)@[16-May-2007, 15-May-2012]

# Knowledge Harvesting from Web Sources

**(Automatic) Construction of Comprehensive Knowledge Bases**
**100 Mio's of entities, classes, relationships, common-sense properties**
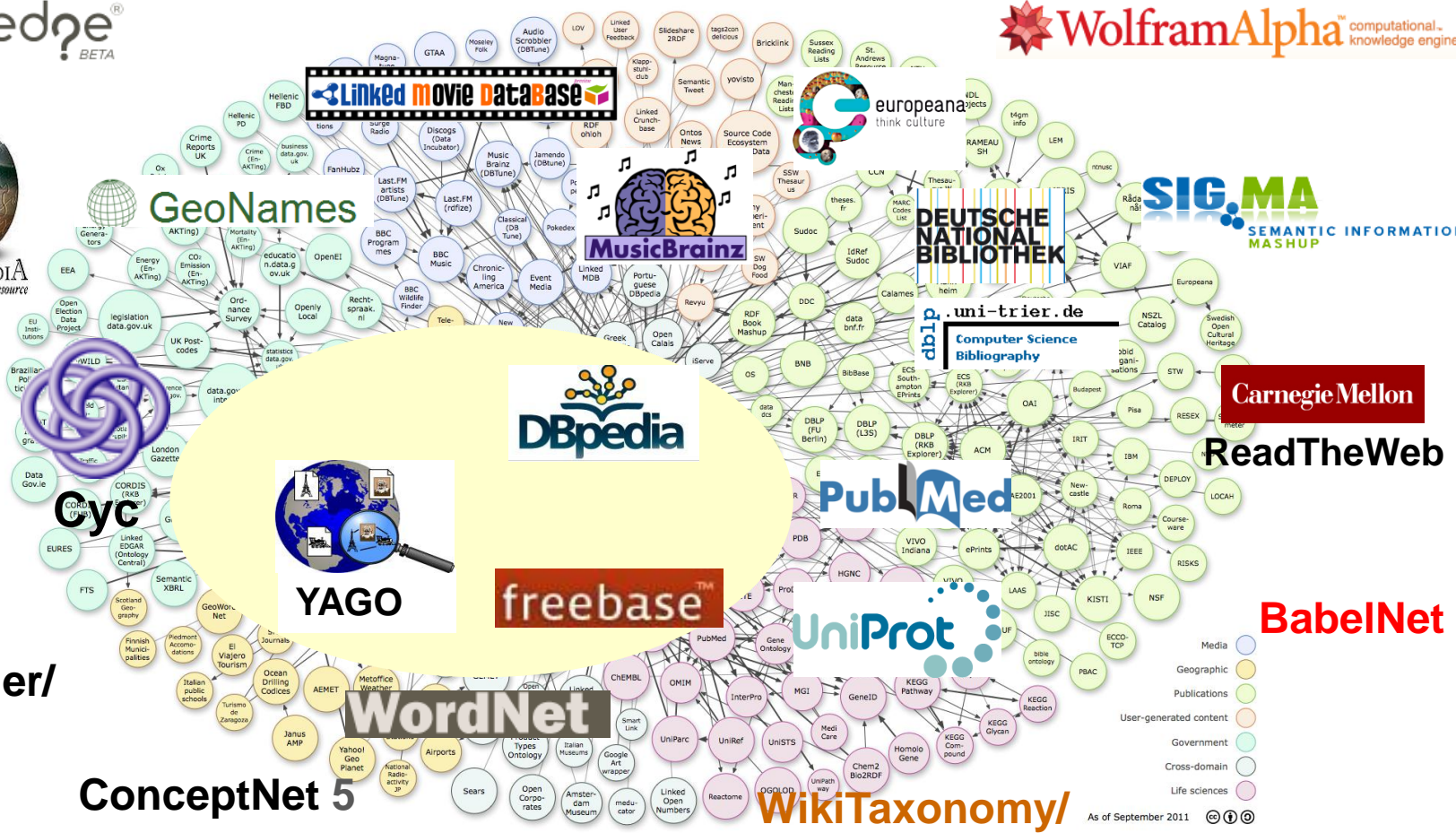
# Knowledge for Intelligence

**Enabling technology for:**
- ★ **disambiguation** in written & spoken natural language
- ★ **deep reasoning** (e.g. QA to win quiz game)
- ★ **machine reading** (e.g. to summarize book or corpus)
- ★ **semantic search** in terms of entities&relations (not keywords&pages)
- ★ **entity-level linkage** for the Web of Data

★ Politicians who are also scientists?

★ European composers who have won film music awards?

★ French professors who founded Internet companies?

★ Relationships between
Alexander Pushkin, Evariste Galois, Johnny Ringo, and Hamlet?

★ Enzymes that inhibit HIV?
Influenza drugs for teens with high blood pressure?
...

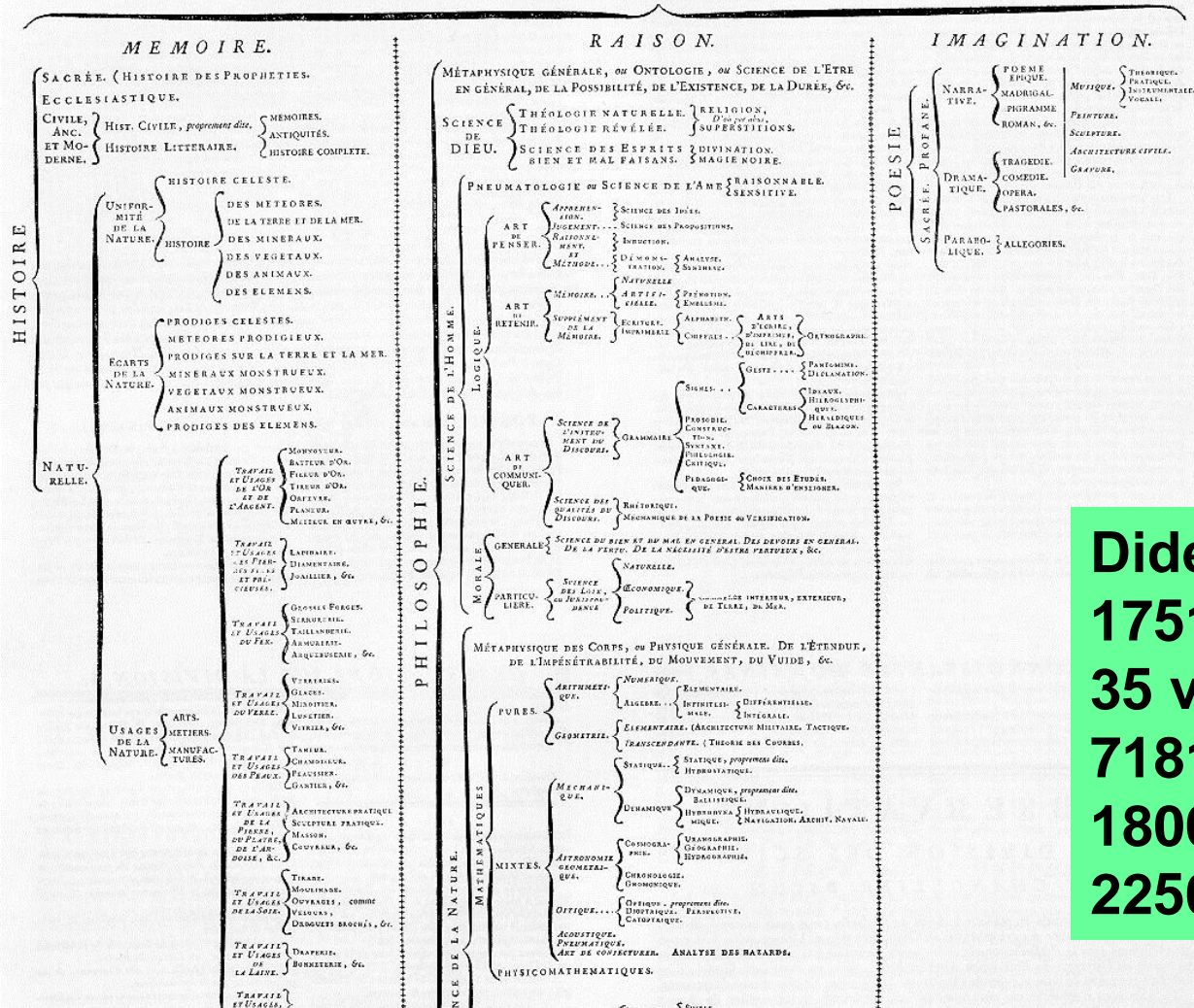# Outline

✓ **Knowledge for Machines**

★ **Construction of Knowledge Bases**

★ **KB Population from Text & Web Pages**

★ **Knowledge for Intelligent Applications**

★ **Opportunities & Challenges**

★ **Conclusions**
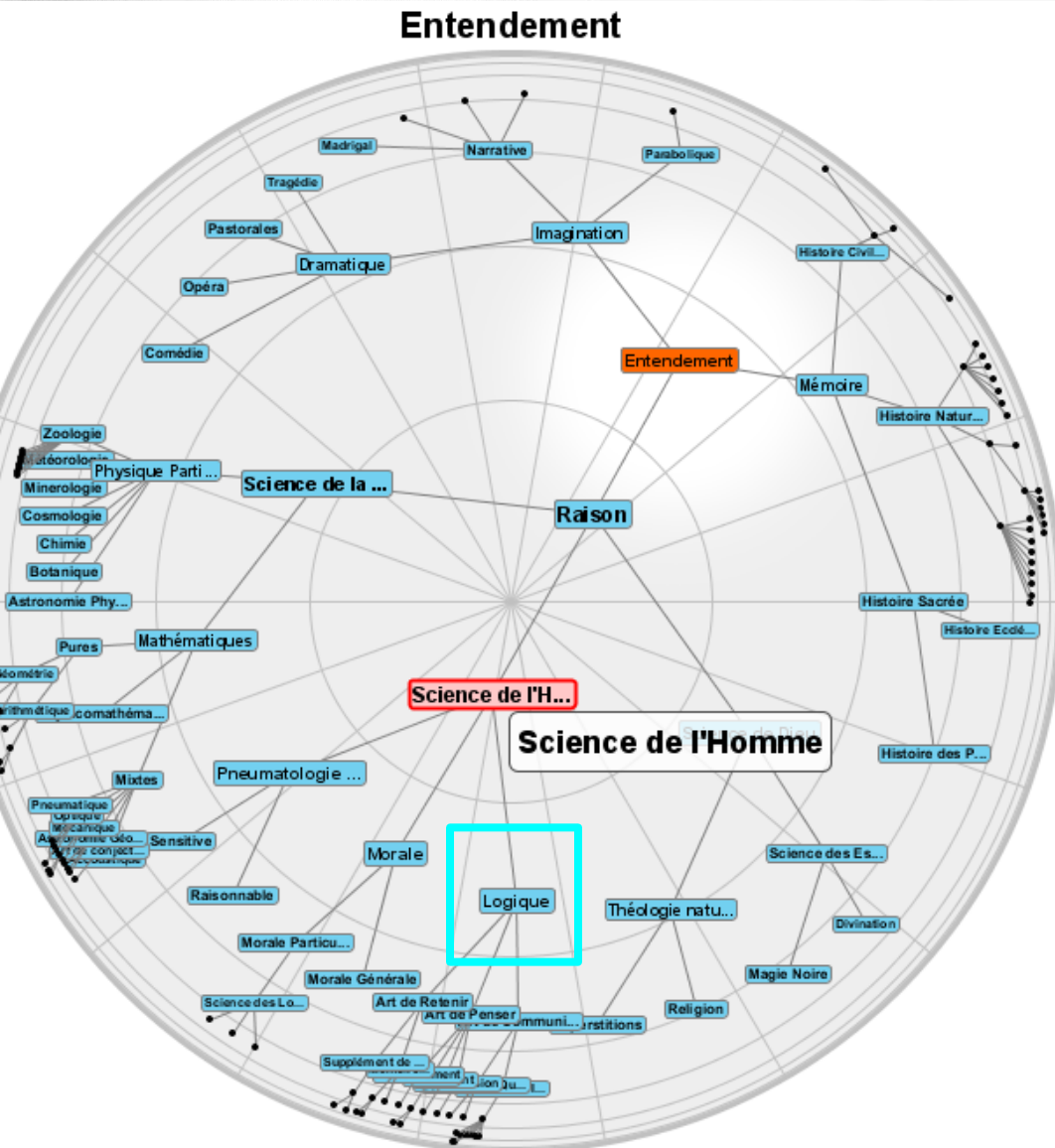
# The World's First Knowledge Base

Denis
Diderot
(1713-1784)

Jean le Rond
d'Alembert
(1717-1783)

**Diderot & d'Alembert (Eds.)**
**1751-1780**
**35 volumes, 4250 copies**
**71818 articles, 3129 figures**
**18000 pages, 20 Mio. words**
**2250 contributors**

# The World's First Knowledge Base



Entendement

Science de l'Homme

# Modern Knowledge Bases

**WordNet** project
(1985-now)

**George Miller**

**Christiane Fellbaum**

## WordNet Search - 3.1
- WordNet home page - Glossary - Help

### Noun

- S: (n) **enterprise**, endeavor, endeavour (a purposeful or industrious undertaking (especially one that requires effort or boldness)) *"he had doubts about the whole enterprise"*
- S: (n) **enterprise** (an organization created for business ventures) *"a growing enterprise must have a bold leader"*
- S: (n) **enterprise**, enterprisingness, initiative, go-ahead (readiness to embark on bold new ventures)

# Modern Knowledge Bases

- *direct hyponym* / *full hyponym*
  - **S:** (n) giant (an unusually large enterprise) *"Walton built a retail giant"*
  - **S:** (n) collective (members of a cooperative enterprise)
  - **S:** (n) business, concern, business concern, business organization, business organisation (a commercial or industrial enterprise and the people who constitute it) *"he bought his brother's business"; "a small mom-and-pop business"; "a racially integrated business concern"*
    - *direct hyponym* / **full hyponym**
      - **S:** (n) agency (a business or organization that provides a particular service, especially the mediation of transactions between two parties)
        - **S:** (n) advertising agency, ad agency (an agency that designs advertisement to call public attention to its clients)
        - **S:** (n) credit bureau (a private firm that maintains consumer credit data files and provides credit information to authorized users for a fee)
        - **S:** (n) detective agency (an agency that makes inquiries for its clients)
        - **S:** (n) employment agency, employment office (an agency that finds people to fill particular jobs or finds jobs for unemployed people)
        - **S:** (n) mercantile agency, commercial agency (an organization that provides businesses with credit ratings of other firms) *"Dun & Bradstreet is the largest mercantile agency in the United States"*
        - **S:** (n) news agency, press agency, wire service, press association, news organization, news organisation (an agency to collects news reports for newspapers and distributes it electronically)
          - **S:** (n) syndicate (a news agency that sells features or articles or photographs etc. to newspapers for simultaneous publication)
        - **S:** (n) service agency, service bureau, service firm (a business that makes its facilities available to others for a fee; achieves economy of scale)

  - **S:** (n) firm, house, business firm (the members of a business organization that owns or operates one or more establishments) *"he worked for a brokerage house"*
    - **S:** (n) corporation, corp (a business firm whose articles of incorporation have been approved in some state)
      - **S:** (n) conglomerate, empire (a group of diverse companies under common ownership and run as a single organization)
        - **S:** (n) publishing conglomerate, publishing empire (a conglomerate of publishing companies)
      - **S:** (n) large cap (a corporation with a large capitalization) *"he works for a large cap"*
      - **S:** (n) small cap (a corporation with a small capitalization) *"this annual conference is a showcase for ambitious small caps"*
      - **S:** (n) closed corporation, close corporation, private corporation, privately held corporation (a corporation owned by a few people; shares have no public market)
        - **S:** (n) family business (a corporation that is entirely owned by the members of a single family)
      - **S:** (n) closely held corporation (stock is publicly traded but most is held by a few shareholders who have no plans to sell)
      - **S:** (n) shell corporation, shell entity (a company that is incorporated but has no assets or operations)
      - **S:** (n) Federal Deposit Insurance Corporation, FDIC (a federally sponsored corporation that insures accounts in national banks and other

# Modern Knowledge Bases

S: (n) **enterprise** (an organization created for business ventures) *"a growing enterprise must have a bold leader"*

- *direct hyponym* / *full hyponym*
  - S: (n) giant (an unusually large enterprise) *"Walton built a retail giant"*
  - S: (n) collective (members of a cooperative enterprise)
  - S: (n) business, concern, business concern, business organization,
    - S: (n) **entrepreneur**, enterpriser (someone who organizes a business venture and assumes the risk for it)
      - *has instance*
        - S: (n) Gates, Bill Gates, William Henry Gates (United States computer entrepreneur whose software company made him the youngest multi-billionaire in the history of the United States (born in 1955))
        - S: (n) Sinclair, Clive Sinclair, Sir Clive Marles Sinclair (English electrical engineer who founded a company that introduced many innovative products (born in 1940))

**+  focus on classes and taxonomic structure**
**–  few or no instances (entities) of classes**

        - S: (n) capitalist (a person who invests capital in a business (especially a large business))
          - S: (n) person, individual, someone, somebody, mortal, soul (a human being) *"there was too much for one person to do"*
            - S: (n) organism, being (a living thing that has (or can

# Knowledge Communities & New Opportunities

Jimmy Wales

Larry Sanger

## Steve Jobs

From Wikipedia, the free encyclopedia

*For the biography, see Steve Jobs (biography).*

**Steven Paul Jobs** (/'dʒɒbz/; February 24, 1955 – October 5, 2011)[4][5] was an American businessman and inventor widely recognized as a charismatic pioneer of the personal computer revolution.[6][7] He was co-founder, chairman, and chief executive officer of Apple Inc. Jobs also co-founded and served as chief executive of Pixar Animation Studios; he became a member of the board of directors of The Walt Disney Company in 2006, following the acquisition of Pixar by Disney.

In the late 1970s, Apple co-founder Steve Wozniak engineered one of the first commercially successful lines of personal computers, the Apple II series. Jobs directed its aesthetic design and marketing along with A.C. "Mike" Markkula, Jr. and others. In the early 1980s, Jobs was among the first to see the commercial potential of Xerox PARC's mouse-driven graphical user interface, which led to the creation of the Apple Lisa (engineered by Ken Rothmuller and John Couch) and, one year later, creation of Apple employee Jef Raskin's Macintosh.

After losing a power struggle with the board of directors in 1985, Jobs left Apple and founded NeXT, a computer platform development company specializing in the higher-education and business markets. NeXT was eventually acquired by Apple in 1996, which brought Jobs back to the company he co-founded, and provided Apple with the NeXTSTEP codebase, from which the Mac OS X was developed."[8] Jobs was named Apple advisor in 1996, interim CEO in 1997, and CEO from 2000 until his resignation. He oversaw the development of the iMac, iTunes, iPod, iPhone, and iPad and the company's Apple Retail Stores.[9] In 1986, he acquired the computer graphics division of Lucasfilm Ltd, which was spun off as Pixar Animation Studios.[10] He was credited in *Toy Story* (1995) as an executive producer. He remained CEO and majority shareholder at 50.1 percent until its acquisition by The Walt Disney Company in 2006,[11] making Jobs Disney's largest individual shareholder at seven percent and a member of Disney's Board of Directors.[12][13]

In 2003, Jobs was diagnosed with a pancreas neuroendocrine tumor. Though it was initially treated, he reported a hormone imbalance, underwent a liver transplant in 2009, and appeared progressively thinner as his health declined.[14] On medical leave for most of 2011, Jobs resigned as Apple CEO in August that year and was elected Chairman of the Board. On October 5, 2011, Jobs died of respiratory arrest related to his metastatic tumor. He

### Steve Jobs

Jobs holding a white iPhone 4 at Worldwide Developers Conference 2010

| | |
|---|---|
| **Born** | Steven Paul Jobs February 24, 1955[1][2] San Francisco, California, U.S.[1][2] |
| **Died** | October 5, 2011 (aged 56)[2] Palo Alto, California, U.S. |
| **Nationality** | American |
| *Alma mater* | Reed College (dropped out) |

# Knowledge Commu[nity] & New Opportuni[ties]

## Steve Jobs

From Wikipedia, the free encyclopedia

*For the biography, see Steve Jobs (biography).*

**Steven Paul Jobs** (/ˈdʒɒbz/; February 24, 1955 – October 5, 2011)[4][5] was an American busine[ss] inventor widely recognized as a charismatic pioneer of the personal computer revolution.[6][7] He w[as] chairman, and chief executive officer of Apple Inc. Jobs also co-founded and served as chief exec[utive] Animation Studios; he became a member of the board of directors of The Walt Disney Company the acquisition of Pixar by Disney.

In the late 1970s, Apple co-founder Steve Wozniak engineered one of the first commercially succ[essful] personal computers, the Apple II series. Jobs directed its aesthetic design and marketing along [with] Markkula, Jr. and others. In the early 1980s, Jobs was among the first to see the commercial pot[ential] PARC's mouse-driven graphical user interface, which led to the creation of the Apple Lisa (engine[ered]

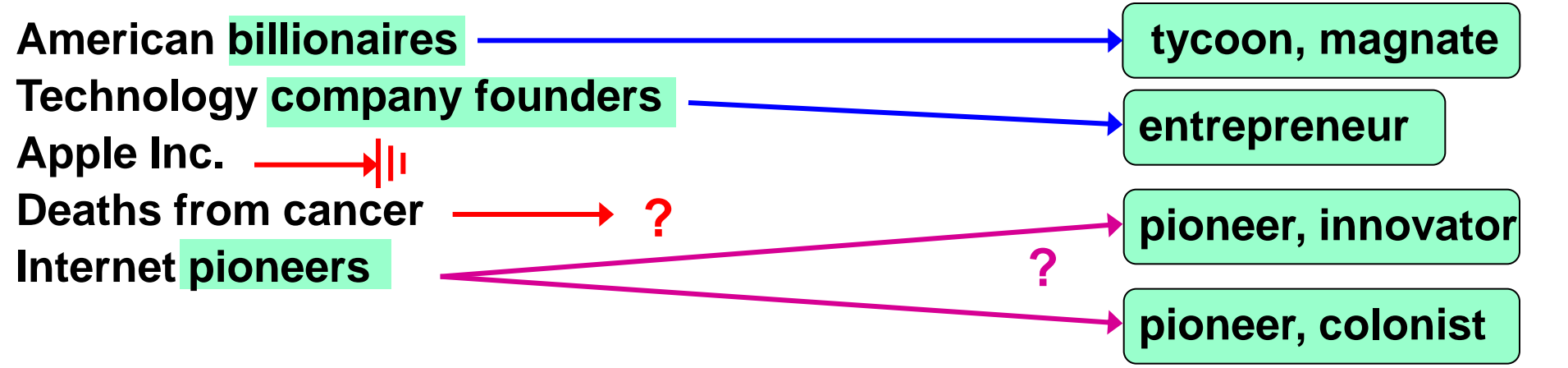| Born | Steven Paul Jobs February 24, 1955[1][2] San Francisco, California, U.S.[1][2] |
|---|---|
| Died | October 5, 2011 (aged 56)[2] Palo Alto, California, U.S. |
| Nationality | American |
| *Alma mater* | Reed College (dropped out) |
| Occupation | Co-founder, Chairman and CEO, Apple Inc., Co-founder and CEO, Pixar, Co-founder and CEO, NeXT Inc. |
| Years active | 1974–2011 |

Categories: Steve Jobs │ 1955 births │ 2011 deaths │ American adoptees │ American billionaires │ American chief executives │ American computer businesspeople │ American industrial designers │ American inventors │ American people of German descent │ American people of Swiss descent │ American people of Syrian descent │ American technology company founders │ American Zen Buddhists │ Apple Inc. │ Apple Inc. employees │ Businesspeople from California │ Businesspeople in software │ Cancer deaths in California │ Computer designers │ Computer pioneers │ Deaths from pancreatic cancer │ Disney people │ Internet pioneers │ National Medal of Technology recipients │ NeXT │ Organ transplant recipients │ People from the San Francisco Bay Area │ Pescetarians │ Reed College alumni

# Automatic Knowledge Base Construction

map 300K Wikipedia categories onto 150K WordNet classes

American **billionaires** → **tycoon, magnate**

Technology **company founders** → **entrepreneur**

Apple Inc. —⊣‖

Deaths from cancer → **?**

Internet **pioneers** → **pioneer, innovator**

**?** → **pioneer, colonist**

Integrating entities & facts from Wikipedia with semantic classes in WordNet

→ **YAGO knowledge base**:

10 Mio. entities, 350 000 classes,
180 Mio. facts, 100 relations,
100 languages, 2 Bio. triples,
95% accuracy

**YAGO**

# Large-Scale Universal Knowledge Bases

**Yago:** **10 Mio. entities, 350 000 classes,**
**180 Mio. facts, 100 properties, 100 languages**
**high accuracy, no redundancy, limited coverage**
http://yago-knowledge.org

**Dbpedia:** **4 Mio. entities, 250 classes,**
**500 Mio. facts, 6000 properties**
**high coverage, live updates**
http://dbpedia.org

**Freebase:** **25 Mio. entities, 2000 topics,**
**100 Mio. facts, 4000 properties**
**interesting relations (e.g., romantic affairs)**
http://freebase.com

**NELL:** **300 000 entity names, 300 classes, 500 properties,**
**1 Mio. beliefs, 15 Mio. low-confidence beliefs**
**learned rules**
http://rtw.ml.cmu.edu/rtw/

**and more …**                    **plus Linked Data**

# Outline

✓ **Knowledge for Machines**

✓ **Construction of Knowledge Bases**

★ **KB Population from Text & Web Pages**

★ **Knowledge for Intelligent Applications**

★ **Opportunities & Challenges**

★ **Conclusions**

# Knowledge Base Population

**hasAdvisor (<person>, <person>)**    **AlmaMater (<person>, <university>)**

*Jim Gray* **was one of** *Harrison***'s** *best students*
*Barbara Liskov* **wrote her thesis** *under the guidance of McCarthy*
*Pierre Senellart* **was one of** *Serge***'s best students**
*Serge Abiteboul* **was grateful to** *his advisor Ginsburg*
*Seymour Ginsburg* **was** *influenced by Leibniz*
*Seymour* **and** *his advisor Ben Dushnik* **never** *co-authored* **any paper**
*Serge* **and** *Victor* *co-authored* **many papers**

➡ **hasAdvisor (JimGray, MikeHarrison)**
**hasAdvisor (BarbaraLiskov, JohnMcCarthy)**
**…**
**AlmaMater (JimGray, Berkeley)**
**AlmaMater (BarbaraLiskov, Stanford)**
**…**

**Pattern-based Gathering (statistical evidence)** **+** **Constraint-aware Reasoning (logical consistency)**

# Pattern-based Gathering of Fact Candidates
## (Linguistic Patterns and Statistical Evidence)

**Facts & _Fact Candidates_**

**(JimGray, MikeHarrison)**

**(BarbaraLiskov, JohnMcCarthy)**

_(Serge, Seymour)_
_(Pierre, Serge)_
_(Nicoleta, Ioana)_
_(Nicoleta, Fabian)_

_(Serge, Victor)_
_(Victor, Serge)_
_(Seymour, Leibniz)_
_(Seymour, Schwarzenegger)_
_(Nicoleta, Versailles)_

**Patterns**

**X and his advisor Y**

**X under the guidance of Y**

**X and Y in their paper**

**X co-authored with Y**

**X was influenced by Y**

**…**

- **good for recall**
- **noisy, drifting**
- **not robust enough for high precision**

# Constrained Reasoning for Logical Consistency

Use **knowledge** (consistency contraints)
for joint reasoning on hypotheses
and pruning of false candidates

$\forall$ **x, y, z: hasA(x,y) $\land$ hasA(x,z) $\Rightarrow$ y=z**
$\forall$ **x, y: hasA(x,y) $\Rightarrow$ $\neg$ hasA(y,x)**
$\forall$ **x, y, p: occurs (x, y, p) $\land$ goodPattern(p, hasA)**
　　　　**$\Rightarrow$ hasA(x,y)**
$\forall$ **x, y, p: occurs (x, y, p) $\land$ hasA(x,y)**
　　　　**$\Rightarrow$ goodPattern(p, hasA)**
$\forall$ **x, y: hasA(x,y) $\Rightarrow$ type(x)=Person**
$\forall$ **x, y: hasA(x,y) $\Rightarrow$ type(y)=Scientist**
$\forall$ **x, y: hasA(x,y) $\Rightarrow$ hasStudent(y,x)**
$\forall$ **x, y: AlmaMater(x,y) $\Rightarrow$ type(y)=University**
$\forall$ **x,y,u: hasA(x,y) $\land$ AlmaMater(x,u) $\Rightarrow$ FacultyOf (y,u)**
$\forall$ **x,y,s,t: hasA(x,y) $\land$ gradYear(x,s) $\land$ gradYear(y,t) $\Rightarrow$ t < s**

# Constrained Reasoning for Logical Consistency

**Use knowledge (consistency contraints)
for joint reasoning on hypotheses
and pruning of false candidates**

**ground atoms:**

hasA(Serge, Seymour)
hasA(Nicoleta, Ioana)
~~hasA(Nicoleta, Fabian)~~
~~hasA(Serge, Victor)~~
~~hasA (Victor, Serge)~~
~~hasA (Nicoleta, Versailles)~~
~~hasA (Seymour, Leibniz)~~
…

$\forall$ x, y, z: hasA(x,y) $\wedge$ hasA(x,z) $\Rightarrow$ y=z
$\forall$ x, y: hasA(x,y) $\Rightarrow$ $\neg$ hasA(y,x)
$\forall$ x, y, p: occurs (x, y, p) $\wedge$ goodPattern(p, hasA)
            $\Rightarrow$ hasA(x,y)
$\forall$ x, y, p: occurs (x, y, p) $\wedge$ hasA(x,y)
            $\Rightarrow$ goodPattern(p, hasA)
$\forall$ x, y: hasA(x,y) $\Rightarrow$ type(x)=Person
$\forall$ x, y: hasA(x,y) $\Rightarrow$ type(y)=Scientist
$\forall$ x, y: hasA(x,y) $\Rightarrow$ studentOf(y,x)
$\forall$ x, y: AlmaMater(x,y) $\Rightarrow$ type(y)=University
$\forall$ x,y,u: hasA(x,y) $\wedge$ AlmaMater(x,u) $\Rightarrow$ FacultyOf (y,u)
$\forall$ x,y,s,t: hasA(x,y) $\wedge$ gradYear(x,s) $\wedge$ gradYear(y,t) $\Rightarrow$ s < t

**Find consistent subset(s) of atoms ("possible world(s)", "the truth")
    $\rightarrow$ customized Weighted MaxSat solver for set of clauses
    $\rightarrow$ max a posteriori for probabilistic factor graph**

# PROSPERA: Web-Scale Experiments

[N. Nakashole et al.: WSDM'11]

- on **ClueWeb'09** corpus: **500 Mio.** English Web pages
- with **Hadoop cluster** of 10x16 cores and 10x48 GB memory
- 10 seed examples, 5 counter examples
  for each of **15 relations** on sports and academia



more than **100,000 facts** acquired after 6 iterations
with overall **precision ≈ 90%**
and **99% precision@1000** for each of the relations

**www.mpi-inf.mpg.de/yago-naga/prospera/**

# Outline

✓ **Knowledge for Machines**

✓ **Construction of Knowledge Bases**

✓ **KB Population from Text & Web Pages**

★ **Knowledge for Intelligent Applications**

★ **Opportunities & Challenges**

★ **Conclusions**

# Application: Question Answering

William Wilkinson's "An Account of the Principalities of Wallachia and Moldavia" inspired this author's most famous novel

This town is known as "Sin City" & its downtown is "Glitter Gulch"

As of 2010, this is the only former Yugoslav republic in the EU

99 cents got me a 4-pack of Ytterlig coasters from this Swedish chain

question classification & decomposition $\longrightarrow$ knowledge back-ends

WIKIPEDIA
The Free Encyclopedia

DBpedia

freebase

YAGO

D. Ferrucci et al.: Building Watson: An Overview of the DeepQA Project. AI Magazine, Fall 2010.

www.ibm.com/innovation/us/watson/index.htm

# Application: Machine Reading

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

uncleOf

Blomkvist visits Henrik Vanger at his estate on the island of Hedeby.

same    same

The old man draws Blomkvist in by promising solid evidence against Wennerström.

same

Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is

owns

home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist

uncleOf    hires

becomes acquainted with the men of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Harriet.

enemyOf    affairWith

After discovering that Salander has hacked into his computer, he persuades her to assist

same    same

him with research. They eventually become lovers, but Blomkvist has trouble getting close

affairWith

to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings

headOf

supports herself by doing deep background investigations for Dragan Armansky, who, in turn, worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

same

**O. Etzioni, M. Banko, M.J. Cafarella: Machine Reading, AAAI ,06**

# Named Entity Disambiguation

Sergio talked to
Ennio about
Eli's role in the
Ecstasy scene.
This sequence on
the graveyard
was a highlight in
Sergio's trilogy
of western films.

**Mentions
(surface names)**

**?**

**Eli (bible)**

**Eli Wallach**

# Named Entity Disambiguation



Sergio talked to Ennio about Eli's role in the Ecstasy scene. This sequence on the graveyard was a highlight in Sergio's trilogy of western films.

**Mentions (surface names)**

**KB**

Sergio means  Sergio_Leone
Sergio means  Serge_Gainsbourg
Ennio  means  Ennio_Antonelli
Ennio  means  Ennio_Morricone
Eli  means  Eli_(bible)
Eli  means  ExtremeLightInfrastructure
Eli  means  Eli_Wallach
Ecstasy  means  Ecstasy_(drug)
Ecstasy  means  Ecstasy_of_Gold
trilogy  means  Star_Wars_Trilogy
trilogy  means  Lord_of_the_Rings
trilogy  means  Dollars_Trilogy

**Entities (meanings)**

Eli (bible)

Eli Wallach

Ecstasy (drug)

Ecstasy of Gold

Star Wars Trilogy

Lord of the Rings

Dollars Trilogy

# Mention-Entity Graph
## weighted undirected graph with two types of nodes

*Sergio talked to Ennio about Eli's role in the Ecstasy scene. This sequence on the graveyard was a highlight in Sergio's trilogy of western films.*

**Eli (bible)**

**Eli Wallach**

**Ecstasy (drug)**

**Ecstasy of Gold**

**Star Wars**

**Lord of the Rings**

**Dollars Trilogy**

**Popularity (m,e):**
- freq(e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL (context(m), context(e))

**KB+Stats**

# Mention-Entity Graph

**weighted undirected graph with two types of nodes**

*Sergio talked to Ennio about Eli's role in the Ecstasy scene. This sequence on the graveyard*

**was a highlight in Sergio's trilogy of western films.**

**bag-of-words or language model: words, bigrams, phrases**

**Eli (bible)**

**Eli Wallach**

**Ecstasy (drug)**

**Ecstasy of Gold**

**Star Wars**

**Lord of the Rings**

**Dollars Trilogy**

**Popularity (m,e):**
- freq(e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL (context(m), context(e))

**KB+Stats**

# Mention-Entity Graph
## weighted undirected graph with two types of nodes

*Sergio* talked to

*Ennio* about

*Eli*'s role in the

*Ecstasy* scene.

This sequence on

the graveyard

was a highlight in

*Sergio*'s *trilogy*

of western films.

**joint mapping**

**Eli (bible)**

**Eli Wallach**

**Ecstasy (drug)**

**Ecstasy of Gold**

**Star Wars**

**Lord of the Rings**

**Dollars Trilogy**

**KB+Stats**

**Popularity (m,e):**
- freq(e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL (context(m), context(e))

# Mention-Entity Graph
**weighted undirected graph with two types of nodes**



**Popularity (m,e):**
- freq(m,e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL (context(m), context(e))

**KB+Stats**

**Coherence (e,e'):**
- dist(types)
- overlap(links)
- overlap (anchor words)

# Mention-Entity Graph
## weighted undirected graph with two types of nodes

*Sergio talked to Ennio about Eli's role in the Ecstasy scene.*

*This sequence on the graveyard was a highlight in Sergio's trilogy of western films.*

**Eli (bible)**

**Eli Wallach**

**Ecstasy (drug)**

**Ecstasy of Gold**

**Star Wars**

**Lord of the Rings**

**Dollars Trilogy**

American Jews
film actors
artists
Academy Award winners

Metallica songs
Ennio Morricone songs
artifacts
soundtrack music

spaghetti westerns
film trilogies
movies
artifacts

**Popularity (m,e):**
- freq(m,e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL (context(m), context(e))

**KB+Stats**

**Coherence (e,e'):**
- dist(types)
- overlap(links)
- overlap (anchor words)

# Mention-Entity Graph

**weighted undirected graph with two types of nodes**

*Sergio talked to Ennio about Eli's role in the Ecstasy scene.*

*This sequence on the graveyard was a highlight in Sergio's trilogy of western films.*

**Eli (bible)**

**Eli Wallach**

**Ecstasy (drug)**

**Ecstasy of Gold**

**Star Wars**

**Lord of the Rings**

**Dollars Trilogy**

http://.../wiki/Dollars_Trilogy
http://.../wiki/The_Good,_the_Bad,_
http://.../wiki/Clint_Eastwood
http://.../wiki/Honorary_Academy_A

http://.../wiki/The_Good,_the_Bad,_
http://.../wiki/Metallica
http://.../wiki/Bellagio_(casino)
http://.../wiki/Ennio_Morricone

http://.../wiki/Sergio_Leone
http://.../wiki/The_Good,_the_Bad,_
http://.../wiki/For_a_Few_Dollars_M
http://.../wiki/Ennio_Morricone

**Popularity (m,e):**
- freq(m,e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL (context(m), context(e))

**KB+Stats**

**Coherence (e,e'):**
- dist(types)
- overlap(links)
- overlap (anchor words)

# Mention-Entity Graph

## weighted undirected graph with two types of nodes

*Sergio talked to Ennio about Eli's role in the Ecstasy scene.*

*This sequence on the graveyard was a highlight in Sergio's trilogy of western films.*

**Eli (bible)**

**Eli Wallach**

The Magnificent Seven
The Good, the Bad, and the Ugly
Clint Eastwood
University of Texas at Austin

**Ecstasy (drug)**

**Ecstasy of Gold**

Metallica on Morricone tribute
Bellagio water fountain show
Yo-Yo Ma
Ennio Morricone composition

**Star Wars**

**Lord of the Rings**

**Dollars Trilogy**

For a Few Dollars More
The Good, the Bad, and the Ugly
Man with No Name trilogy
soundtrack by Ennio Morricone

**Popularity (m,e):**
- freq(m,e|m)
- length(e)
- #links(e)

**Similarity (m,e):**
- cos/Dice/KL
  (context(m),
   context(e))

**KB+Stats**

**Coherence (e,e'):**
- dist(types)
- overlap(links)
- overlap
  (anchor words)

# Graph Algorithm for **Joint** Mapping



- **Compute dense subgraph to**
    **maximize min weighted degree among entity nodes**
  **such that:**
    **each m is connected to exactly one e (or at most one e)**
- **NP-hard** $\rightarrow$ **approximation algorithms**

**[J. Hoffart et al.: EMNLP'11,  M. Yosef et al.: VLDB'11,  M. Yahya et al.: WWW'12]**

# AIDA: Accurate Online Disambiguation

# AIDA: Accurate Online Disambiguation

# AIDA: Accurate Online Disambiguation



**://www.mpi-inf.mpg.de/yago-naga/aida/**

# AIDA: Accurate Online Disambiguation

# AIDA: Accurate Online Disambiguation

# AIDA: Accurate Online Disambiguation



**http://www.mpi-inf.mpg.de/yago-naga/aida/**

# AIDA: Accurate Online Disambiguation



http://www.mpi-inf.mpg.de/yago-naga/aida/

# General Word Sense Disambiguation

**{songwriter, composer}**

**{cover, perform}**

**{cover, report, treat}**

**{cover, help out}**

Which

song writers

covered

ballads

written by

the Stones ?

# Outline

✓ **Knowledge for Machines**

✓ **Construction of Knowledge Bases**

✓ **KB Population from Text & Web Pages**

✓ **Knowledge for Intelligent Applications**

★ **Opportunities & Challenges**

★ **Conclusions**

# Discovering "Unknown" Knowledge

**so far KB has relations with type signatures**
**<entity1, relation, entity2>**

**< CarlaBruni  marriedTo  NicolasSarkozy>**  ∈ **Person × R × Person**
**< NataliePortman  wonAward  AcademyAward >**  ∈ **Person × R × Prize**

**Open and Dynamic Knowledge Harvesting:**
**would like to discover new entities and new relation types**
**<name1, phrase, name2>**

*Madame Bruni in her happy marriage with Nicolas....*
*The first lady is in passionate love with the French president ....*
*Natalie was honored by the Oscar ...*
*Bonham Carter was disappointed that her nomination for the Oscar ...*

# Temporal Knowledge

for **all people** in Wikipedia (300 000) gather **all spouses**, incl. divorced & widowed, and corresponding **time periods**!

>95% accuracy, >95% coverage, in one night

1) **recall: gather temporal scopes for base facts**
2) **precision: reason on mutual consistency**



28 January 1955 (age 53)
Paris, France

Nicolas Paul Stéphane Sarközy

| Political party | RR (?–2002) |
| | UMP (2002–) |
| Spouse | Marie-Dominique Culioli (div.) |
| | Cécilia Ciganer-Albéniz (div.) |
| | Carla Bruni |
| Children | Pierre (by Culioli) |
| | Jean (by Culioli) |
| | Louis (by Ciganer-Albéniz) |
| Residence | Élysée Palace |
| Alma mater | University of Paris X: Nanterre |
| Occupation | Lawyer |
| Religion | Roman Catholic |

1. Catherine of Aragon — *Divorced*
2. Anne Boleyn — *Beheaded*
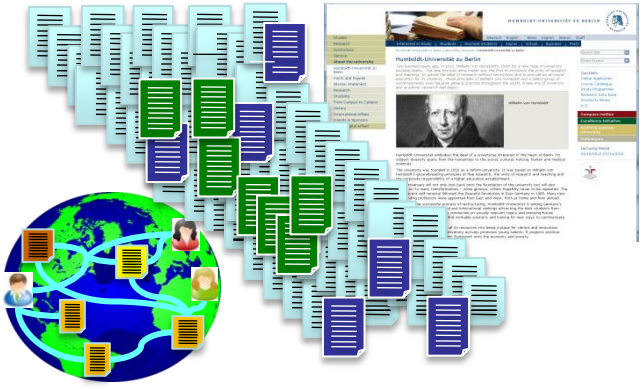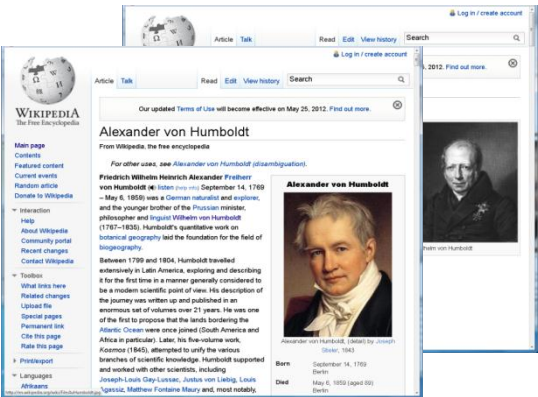3. Jane Seymour — *Died*

**consistency constraints** are potentially helpful:

- **functional dependencies:** *husband, time → wife*
- **inclusion dependencies:** *marriedPerson ⊆ adultPerson*
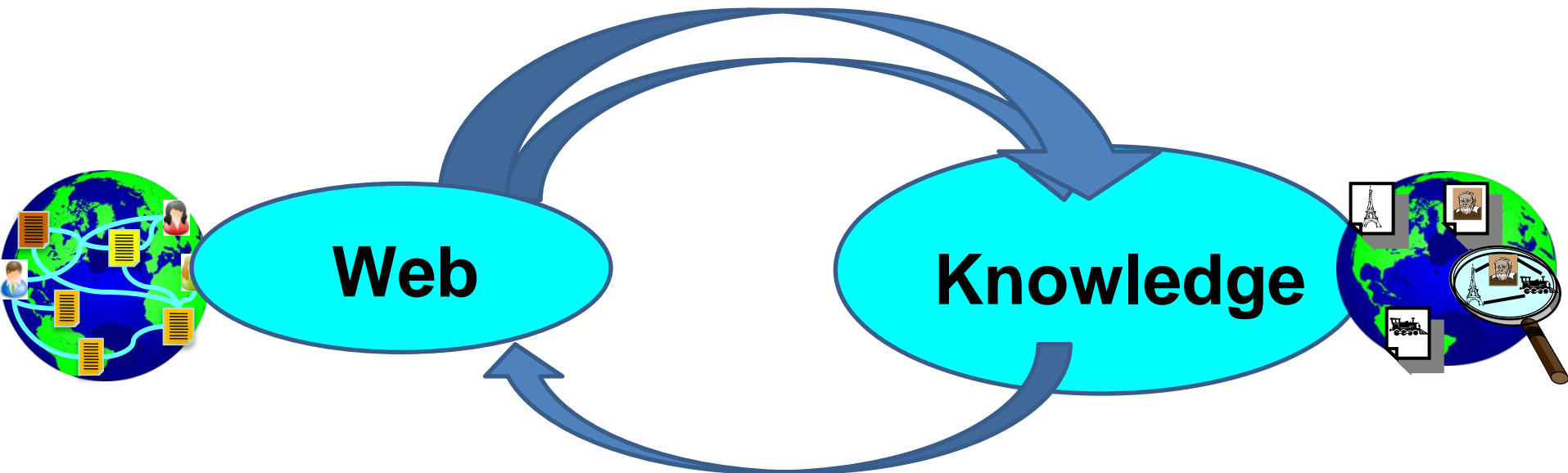- **age/time/gender restrictions:** *birthdate + Δ < marriage < divorce*

# Outline

✓ **Knowledge for Machines**

✓ **Construction of Knowledge Bases**

✓ **KB Population from Text & Web Pages**

✓ **Knowledge for Intelligent Applications**

✓ **Opportunities & Challenges**

★ **Conclusions**

# Summary

# Take-Home Message



- **Machine Knowledge** from Web Sources:
  finally real and big!
- Enabling Asset for **Intelligent Applications**:
  digital humanities, big-data analytics,
  deep machine reading, QA & HCI, etc.
- Rich Research **Opportunities & Challenges**:
  scale & robustness, temporal, multimodal,
  open & real-time knowledge discovery, etc.

# Merci Beaucoup!