

Question Answering Systems

Named entity recognition and disambiguation

Rishiraj Saha Roy

Max Planck Institute for Informatics, Germany

4

Lecture overview

- Logistics
- Named entity basics
- NERD System 1: TAGME
- NERD System 2: AIDA
NED

Logistics: Exam

- All 21 papers *→ clarify doubts* *→ Syllabus*
- ~15 minutes per student *10-20*
- ~3-5 topics *→ questions* *→ content covered in class*
- Question test your understanding *of QA sys.*
- Relative grading *→ Difficult*

Logistics: Assignment_s

6 ECTS

- Self-assessment, no tutorials
- Grades available every week
- Sample good reviews (try)
- Extensive guidelines subtle factors
- Assignments contribute to final grade
- Exact formula only at end
- Our decision is final, vetted by two

Logistics: Material

- Reading material
 - <https://www.mpi-inf.mpg.de/question-answering-systems/>
- Slides and recordings
 - https://drive.google.com/drive/folders/1Z0IjVSjymCD6IX_TcN4dNInuz8Kb42BL
- Assignment grades *shared folder*
 - https://docs.google.com/spreadsheets/d/e/2PACX1vRN0qyr-ooE1JGLfLoPM89ipdPRBprwRUAKLkaRXPqCmDcj0Ht9T5LfGlqEe3gLk3sHS9YIHndDjRLI/pubhtml?gid=1036917168&single=true&urp=gmail_link *→ shared spreadsheet*

Question of the day

How can we disambiguate named entities
present in questions?

You'll find this covered in

- ① ▪ Fast and Accurate Annotation of Short Texts with Wikipedia Pages *TAGME*
- Ferragina and Scaiella
 - CIKM 2010 + IEEE Software 2011 *CIKM, SIGIR, WSDM, WWW : IR+ Web*

Links → ▪ <https://arxiv.org/pdf/1006.3498.pdf>

- ② ▪ Robust Disambiguation of Named Entities in Text *AIDA*
- Hoffart et al.
 - EMNLP 2011 *ACL, EMNLP, NAACL, EACL, COLING : NLP*
 - <https://www.aclweb.org/anthology/D11-1072.pdf>

*↑
official*

Entities: Basics

- Entities and relationships
- Entities and named entities

- People
- Organizations
- Locations

- More

(ever expanding: movies, food, animals...)

- Concepts and classes

- Here: What is there in Wikidata/Wikipedia ☺

obj. facts

NE Michelle is married to Barack

NE India is the birthplace of Gandhi

NE? Pandas eat bamboo

verb-med R
noun-med R
noun E

E, P/R, C, L

consensus

think about generalizing entities:
lists → category
concepts

taxes?
marriage
marriage laws
bank
investment

Entities: Tasks

NLP + IR

- Named entity recognition : NER in test

- Named entity disambiguation / linking
NED/NEL

- Named entity typing : NET

- Main ideas: Similarity and coherence

in NERD
 with Nolan directed film won ...
 ?

Stanford corenlp.run

Michelle is married..

NER ↑

Barcelona FC is the current champion in UEFA

NER

points

to KB / Wikipedia

Film club
Football club

BFC

~M

Answer typing

Not exactly NET

Where does Messi play?

type of ans

Entities: Applications

- 1 ■ Question understanding
 - 2 ■ News readability
 - 3 ■ Information extraction
- Many many more in IR + NLP
 - Can you think of some..?

Which Nolan films ...
→ correct name ⇒ right answer.

^{Mexico FC}
Mexico defeated Sunderland ...
↓ FC

Research paper 1

Fast and Accurate Annotation of Short Texts with Wikipedia Pages


[Fast and accurate annotation of short texts with wikipedia pages](#)

P Ferragina, U Scaiella

IEEE software 29 (1), 70-75

815 * 2011
version

Try it out!



TAGME is a powerful tool that is able to identify *on-the-fly* meaningful short-phrases (called "spots") in an unstructured text and link them to a pertinent [Wikipedia page](#) in a fast and effective way. This annotation process has implications which go far beyond the enrichment of the text with explanatory links because it concerns with the *contextualization* and, in some way, the *understanding* of the text.

Try **TAGME** now!


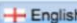

You can play with the demo interface below or check the [TAGME RESTful API](#) we are currently supporting.

Currently **TAGME** is available in English, German and in Italian and it is based on Wikipedia snapshots of April, 2016.

NEWS! TAGME is now hosted by the D4Science infrastructure. Check the [RESTful API page](#) for details.


Developed by [Paolo Ferragina](#) and Ugo Scaiella at [A³ Lab](#)
[Dipartimento di Informatica](#), [University of Pisa](#).

Input Text

 Italiano  English  Deutsche

On this day 24 years ago Maradona scored his infamous "Hand of God" goal against England in the quarter-final of the 1986

Many links



Few links

Reset

TAGME!

URL

<https://tagme.d4science.org/tagme/>

tunable knob

Thanks to Prof. Paolo Ferragina for the slides

Tagged text Topics

On this day 24 years ago Maradona scored his infamous "[Hand of God](#)" [goal](#) against [England](#) in the quarter-final of the 1986

Less links

Wikipedia

Tagged text Topics

On this day 24 years ago [Maradona](#) scored his infamous "[Hand of God](#)" [goal](#) against [England](#) in the [quarter-final](#) of the 1986

More links

API

One issue: Synonymy

challenges

classical approach:
word vectors

He is using Microsoft's browser

He is a fan of Internet Explorer

same

Another issue: Polysemy

the paparazzi photographed the star
the astronomer photographed the star

→ celebrity

→ astronomical body

Wikipedia is a rich source of instances

Title

Steve Jobs

From Wikipedia, the free encyclopedia

For the biography, see [Steve Jobs \(book\)](#).

Steven Paul "Steve" Jobs (/ˈdʒɒbz/; February 24, 1955 – October 5, 2011)^{[5][6]} was an Arab-American^[7] entrepreneur^[8] and inventor,^[9] who was the co-founder, chairman, and CEO of **Apple Inc.** Through Apple, he was widely recognized as a charismatic pioneer of the **personal computer revolution**^[10] and for his influential career in the computer and **consumer electronics** fields, transforming "one industry after another, from computers and smartphones to music and movies..."^[12] 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, when Disney acquired Pixar. Jobs was among the first to see the commercial potential of **Xerox PARC** mouse-driven graphical user interface, which led to the creation of the **Apple Lisa** and, one year later, the **Macintosh**. He also played a role in introducing the **LaserWriter**, one of the first widely available laser printers, to the market.^[13]

After 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. In 1986, he acquired the computer graphics division of **Lucasfilm**, which was spun off as **Pixar**.^[14] He was credited in *Toy Story* (1995) as an executive producer. He served as CEO and majority shareholder until Disney's purchase of Pixar in 2006.^[15] In 1996, after Apple had failed to deliver its operating system, **Copland**, **Gil Amelio** turned to NeXT Computer, and the **NeXTSTEP** platform became the foundation for the **Mac OS X**.^[16] Jobs returned to Apple as an advisor, and took control of the company as an interim CEO. Jobs brought Apple from near bankruptcy to profitability by 1998.^{[17][18][19]}

anchors



infobox

PARC (company)

From Wikipedia, the free encyclopedia
(Redirected from **PARC User Interface**)

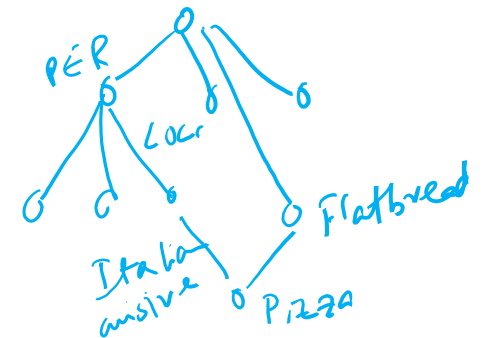
Wikipedia's categories contain classes

research prototypes

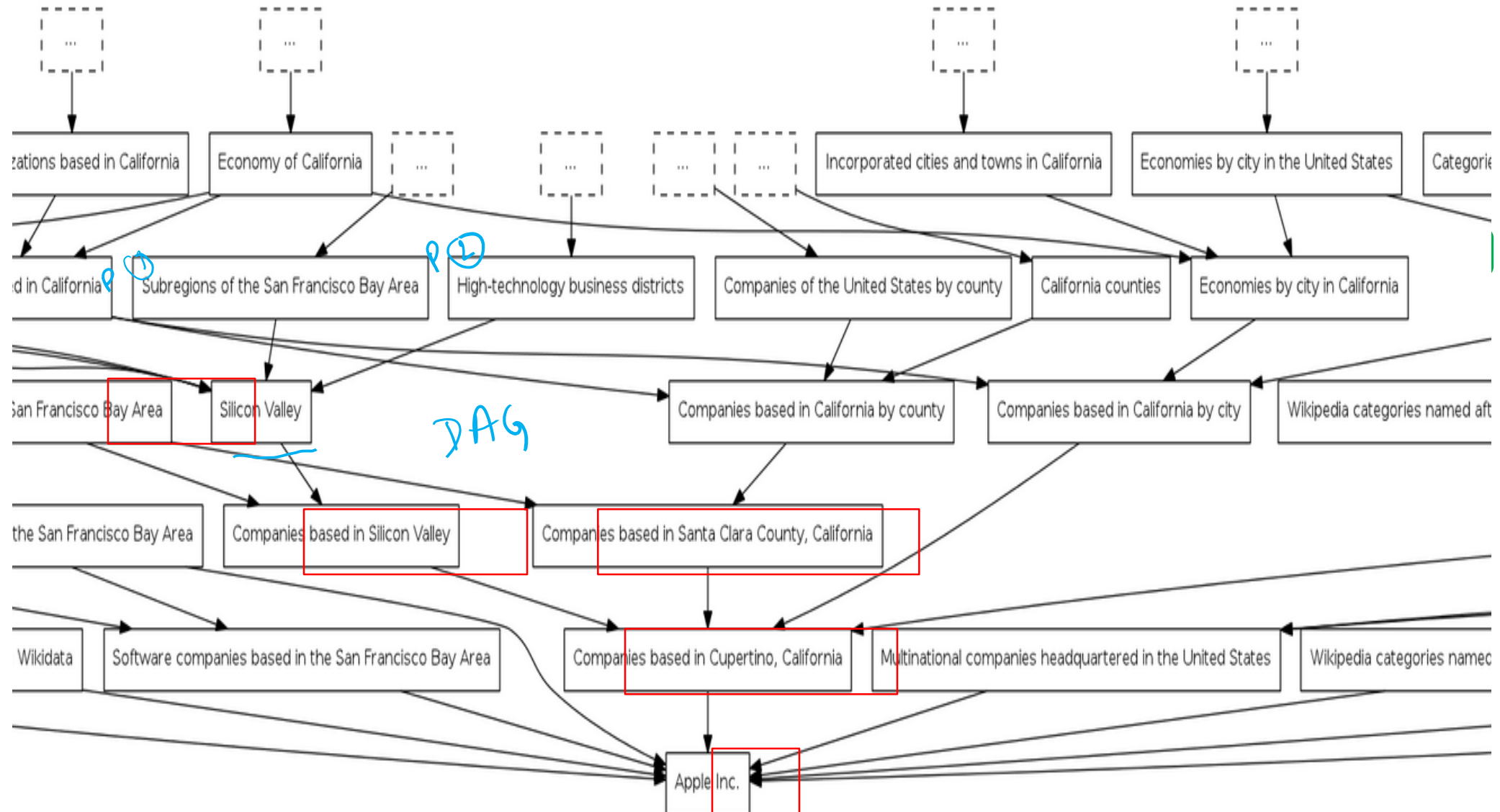
Categories: Steve Jobs | 1955 births | *not class* 2011 deaths | American adoptees | American billionaires ✓
American chief executives | American computer businesspeople | American industrial designers
American inventors | American people of German descent | American people of Swiss descent
American people of Syrian descent | American technology company founders | American Zen Buddhists *classes*
not class 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

ontology

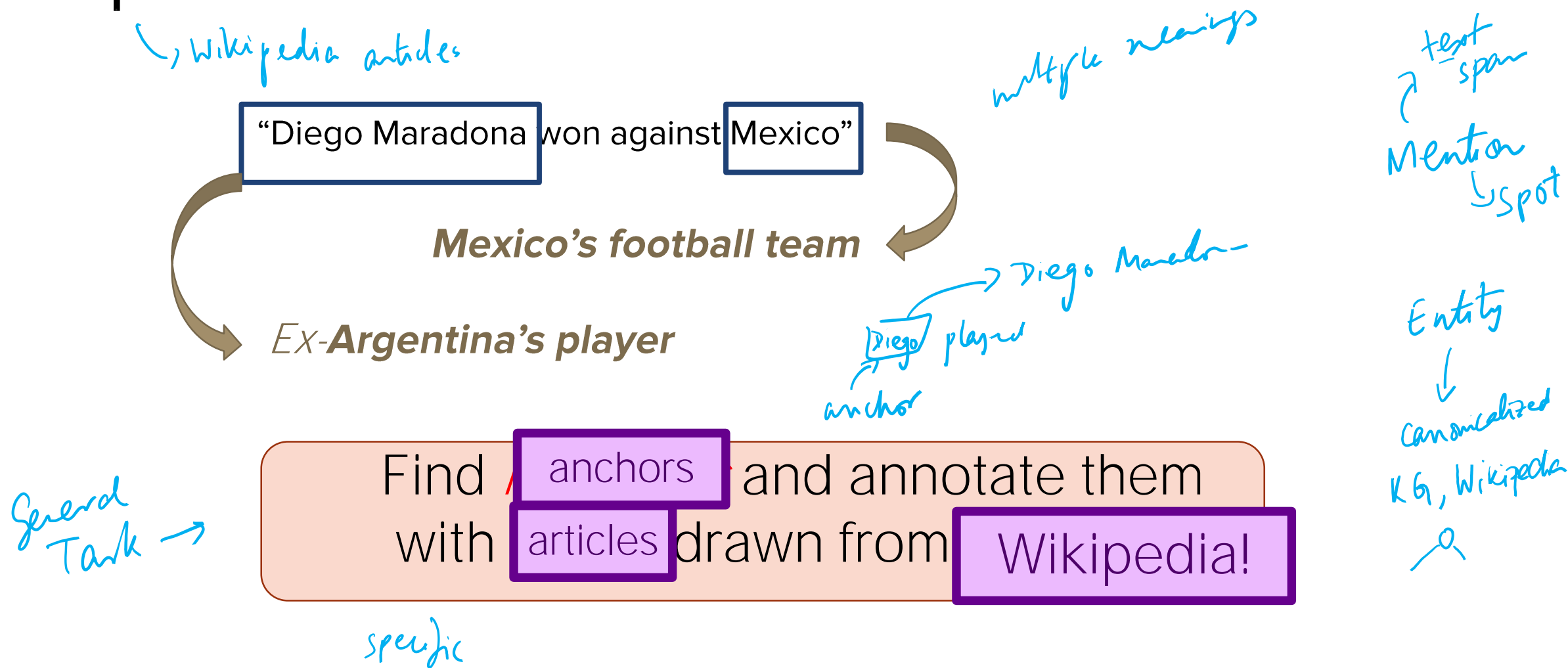
Categories typically form a taxonomic DAG



DAG of categories



Topic-based annotation



Synonymy

Internet Explorer

From Wikipedia, the free encyclopedia

Internet Explorer^[a] (formerly **Microsoft Internet Explorer**^[b] **Microsoft** and included in the **Microsoft Windows** line of opera were available as free downloads, or in-service packs, and in browser is discontinued, but still maintained.^[4]

He is using Microsoft's browser


She plays with Internet Explorer

Polysemy

 *Celebrity is a person who is famously recognized ...*

the paparazzi photographed the **star**

the astronomer photographed the **star**

 *Star is a massive, luminous ball of plasma ...*

Why is it a difficult problem?

NER + NED



NED

The literature

- ▶ TagMe (Univ. Pisa)
- ▶ DBPedia Spotlight (Univ. Berlin)
- ▶ Illinois Wikifier (Univ. Illinois)
- ▶ AIDA (Max Planck Institute for Informatics)
- ▶ CNMS (Univ. Amsterdam)
- ▶ Wikipedia Miner (Univ. Waikato)

Many commercial software: AlchemyAPI, DBpedia Spotlight, Extractiv, Lupedia, OpenCalais, Saplo, SemiTags, TextRazor, Wikimeta, Yahoo! Content Analysis, Zemanta.

The TAGME system

- Designed for **short texts**
 - news, blogs, search-results snippets, tweets, ads, etc
 - competitive on long texts too

→ unique
what are long texts ??
new problem!

- Achieves **high accuracy**
 - Massive experimental test on millions of short texts

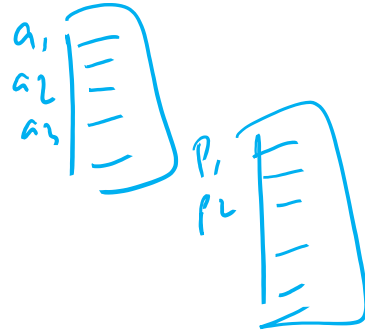
- **Fast**
 - More than 10x faster than others
 - 100% Java

Good research problem
- everything done ??
- look deeper!
- Problem + quality of results
find out → speed? / often unsaid papers
Memory?

TAGME: Distilled information

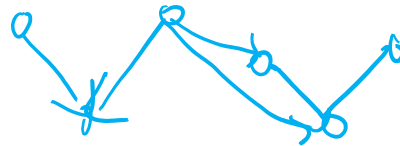
Luceme / Solver / Indri

- Anchor dictionary



- Page catalog

- In-link graph



clean notation

page anchor
 $p, a, Pg(a)$

$link(a), freq(a)$ $a = cat$

$Pr(p|a)$

$\checkmark l_p(a) = \frac{link(a)}{freq(a)}$

$p_a \in Pg(a)$

$a \mapsto p$

TAGME: Overview / Anatomy

- Anchor parsing
- Anchor disambiguation
- Anchor pruning

TagMe is NERD \rightarrow NER + NED
find mention \rightarrow parsing

w_1	w_2	w_3	w_4
-------	-------	-------	-------

nested strings

$a_1 \in a_2$
X

Trojan War

Exception: the act
vs act

overlapping mention?

ton cruise ship scene
X

$lp(a_1) > lp(a_2) \Rightarrow link(a_1) \gg link(a_2)$

Features used to link $a \rightarrow p$

prior

Commonness of a page p wrt an anchor a

$$\Pr(p | a) = \frac{\# a \text{ linked to } p}{\# a \text{ as anchor}}$$

Context of a around the mention

$$T = \dots w_1 w_2 w_3 \text{ a } w_4 w_5 w_6 \dots$$

left *right*

and the content of a *page/entity*

$$\text{Page } p = z_1 z_2 z_3 z_4 z_5 z_6 \dots$$

Link inside p

Link probability of an anchor a

$$lp(a) = \frac{\text{freq of } a \text{ as anchor}}{\text{freq of } a \text{ in the text}}$$

Graph-based features

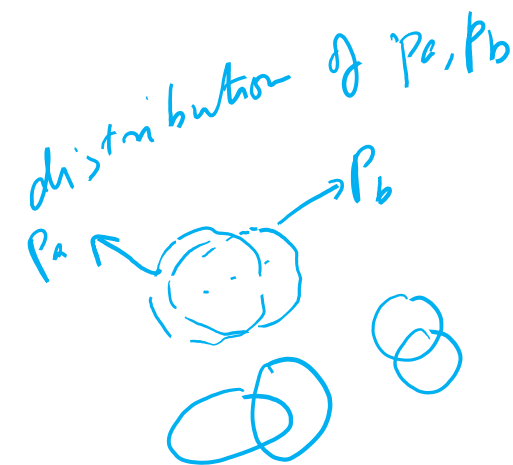
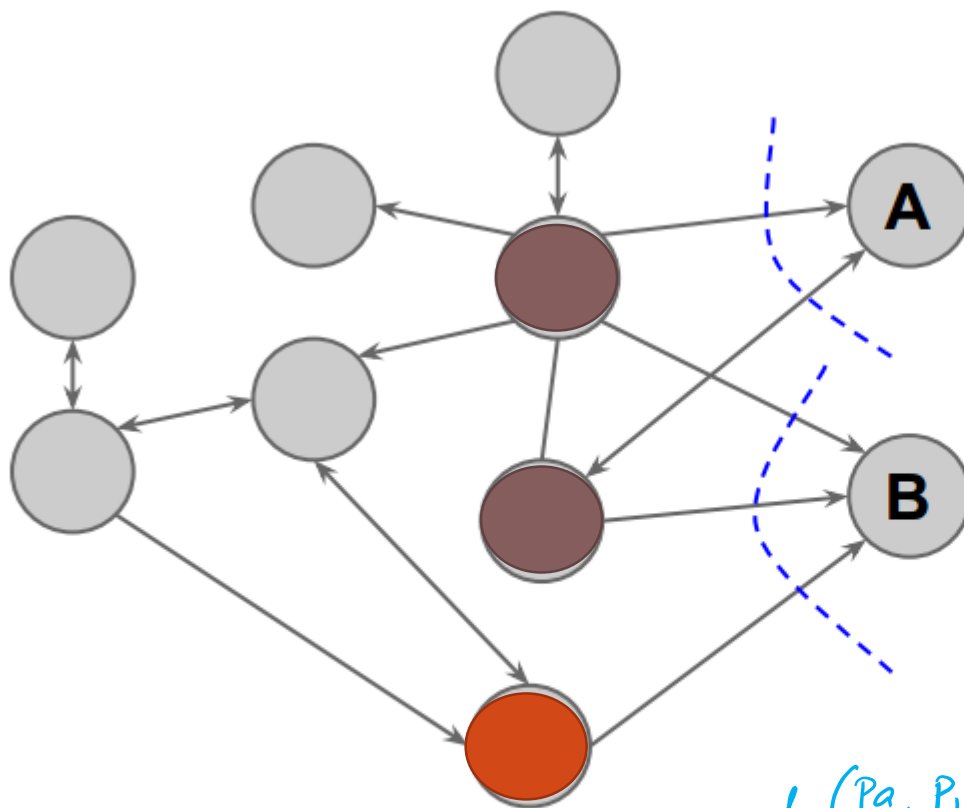
a is a mention-node

p is an entity-node

Links $a \rightarrow p$ and paths between pages



Relatedness between pages



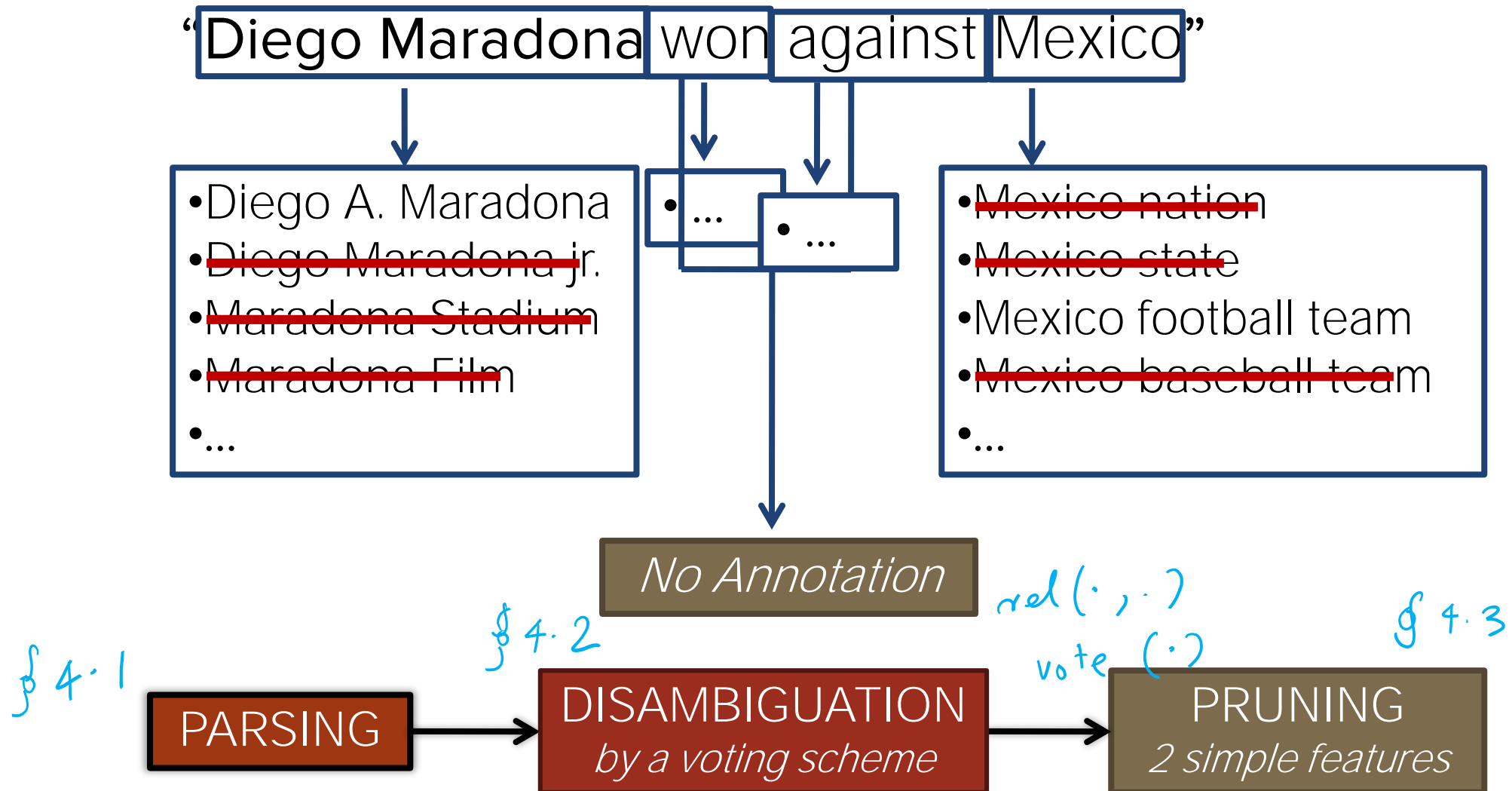
$rel(P_a, P_b)$

$$sr(a, b) = \frac{\log(\max(|A|, |B|)) - \log(|A \cap B|)}{\log(|W|) - \log(\min(|A|, |B|))}$$

$\log(10^6) \sim 6$ constant

Milne & Witten

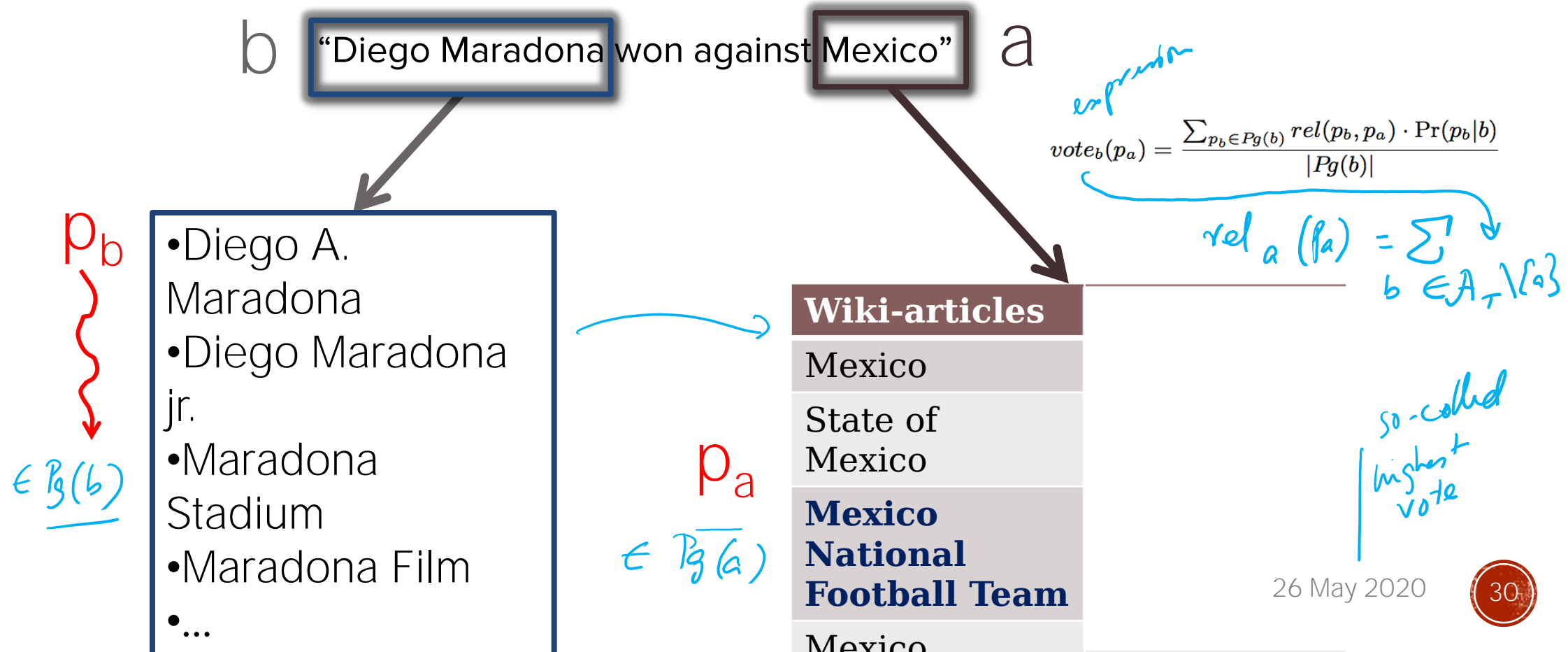
How TAGME works



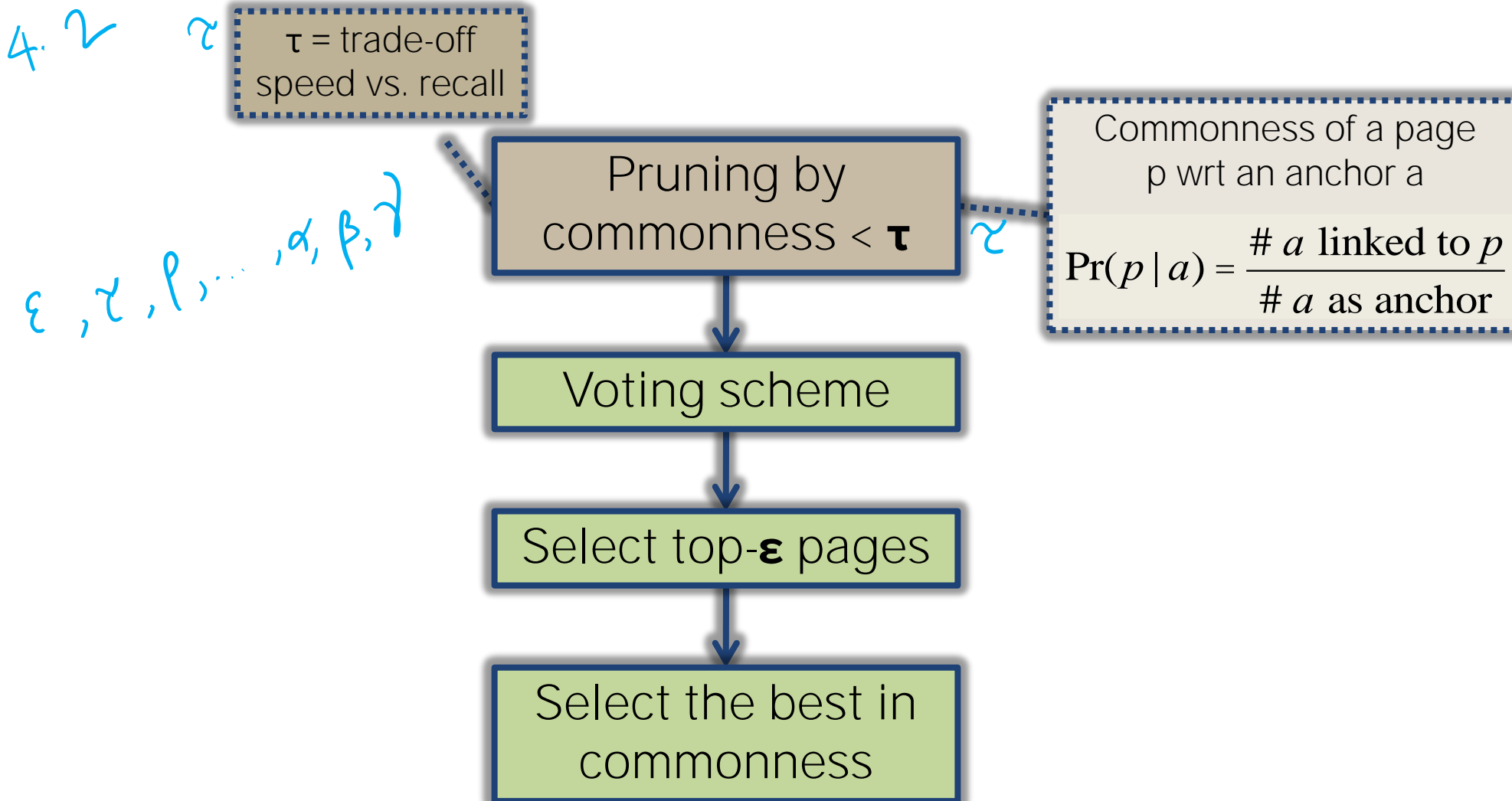
Disambiguation: The Voting Scheme

4. 2

Collective agreement among topics via voting



Disambiguation: All steps



Pruning §4.3

- Use 2 features:
 - link probability
 - coherence wrt context
v. imp. in NED
- Compute a p score via
 - 3 classifiers, avg, linear combination

training step

if $p < p_{NA}$ then prune!

threshold

Link probability of an anchor a

$$lp(a) = \frac{\text{freq of } a \text{ as anchor}}{\text{freq of } a \text{ in the text}}$$

Avg. relatedness btw the assigned concept to the others

$$coh(a \rightarrow p_a) = \frac{1}{|S| - 1} \sum_{p_b \in S \setminus \{p_a\}} rel(p_a, p_b)$$

balance precision vs. recall

Research paper 2

Robust Disambiguation of Named Entities in Text

Robust disambiguation of named entities in text

J Hoffart, MA Yosef, I Bordino, H Fürstenau, M Pinkal, M Spaniol, ...

Proceedings of the Conference on Empirical Methods in Natural Language ...

833

2011

+

Aida: An online tool for accurate disambiguation of named entities in text and tables

MA Yosef, J Hoffart, I Bordino, M Spaniol, G Weikum

Proceedings of the VLDB Endowment 4 (12), 1450-1453

158

2011

Disambiguating to KG entities

- TAGME is fast and effective
- Works well for short texts
- Does not go all the way!
- Wikipedia more general, but we need KG-linking!
- Lookup KG entities using Wikipedia links?
- Harness KG properties! Enter AIDA.

medi- → wikipedia → KG (w.kidolc)

search KG
Christopher Nolan (Wikipedia)
↓
role in KG
Chr.

✓ Try it out!

AIDA → Ambiverse startup

AIDA ↓
AIC Onl Disa NE
→ 6-1

AMBIVERSE
Text to Knowledge

AmbiverseNLU Demo

Jack founded Alibaba in Hangzhou with investments from SoftBank and Goldman.

Settings *time*

Confidence threshold: 0.075






Language: Auto

Enter any text in English, German, Spanish, or Chinese.

ANALYZE

☒ Visual ☐ JSON *concept*

Jack founded Alibaba in Hangzhou with investments from SoftBank and Goldman.

Person	Organization	Location	Organization	Organization
 Jack Ma Chinese businessman	 Alibaba Group Hangzhou-based group of Internet-based companies	 Hangzhou capital of Zhejiang Province, China	 SoftBank Group Japanese company	 Goldman Sachs American multinational investment banking firm

✓ <https://ambiversenlu.mpi-inf.mpg.de/>

Disambiguating names to entities

AIDA
Framework
Mention & ambiguity
Entity candidates
Popularity prior
Context sim. & m.
Coherence
Overall obj. fn.

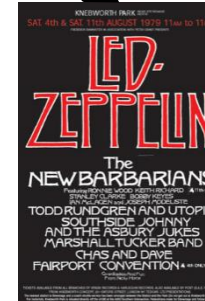
mentions



NERD is
non-trivial

When Page played Kashmir at Knebworth, his Les Paul was uniquely tuned.

Images taken from Wikipedia under CC BY-SA 3.0

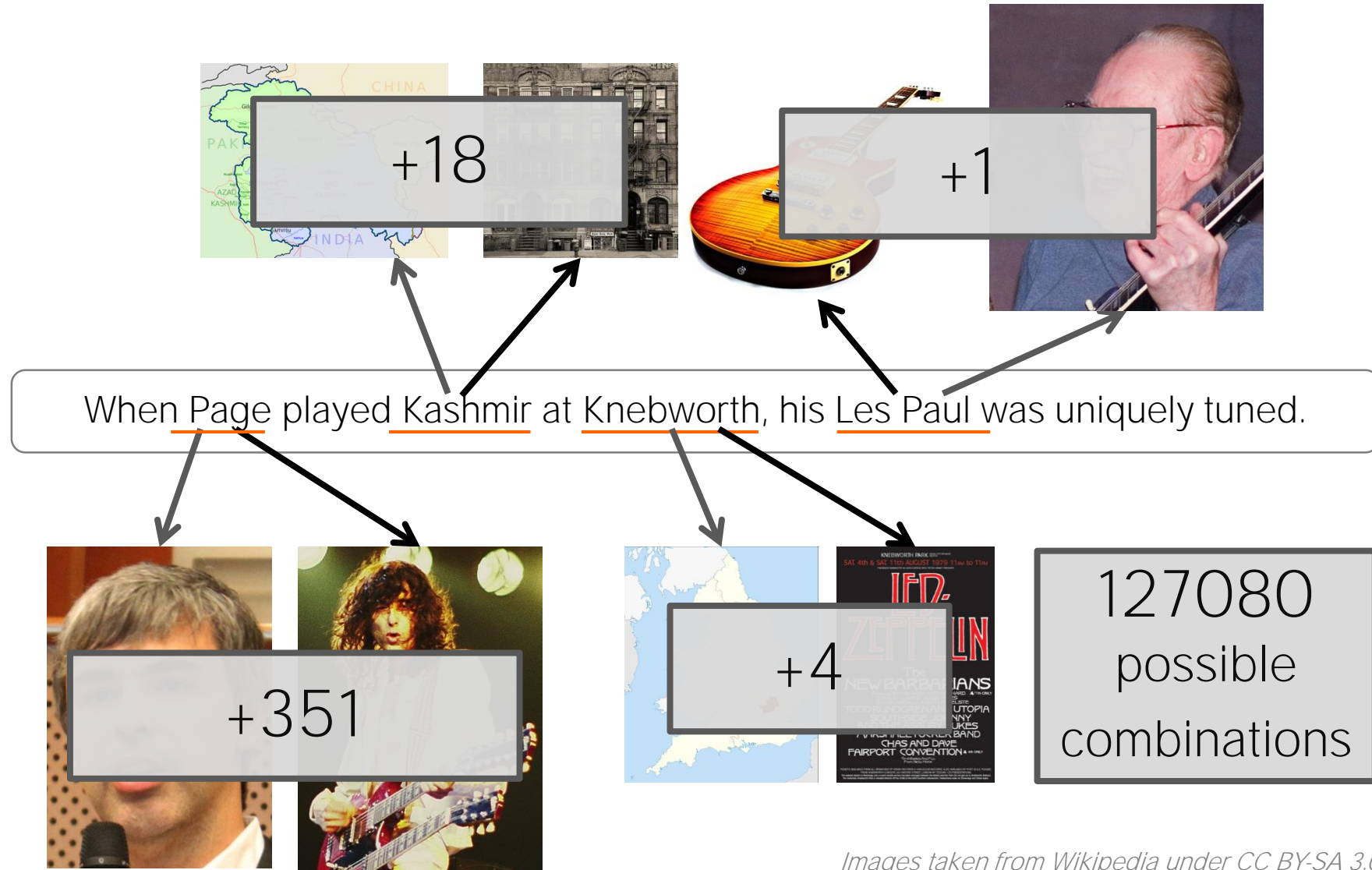


AIDA uses (Stanford) NER

NERD

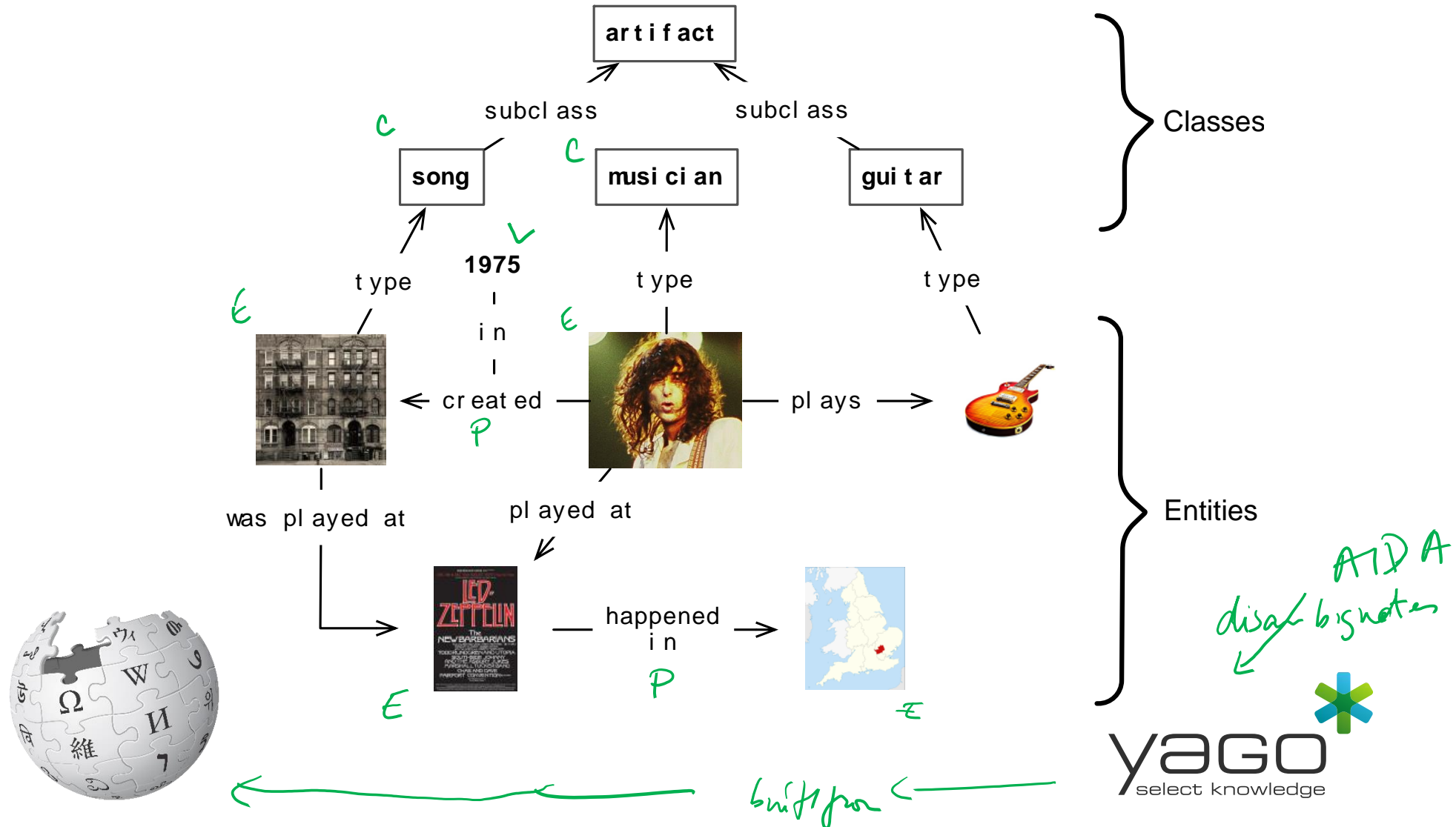
Thanks to Dr. Johannes Hoffart for the slides

Disambiguating names to entities



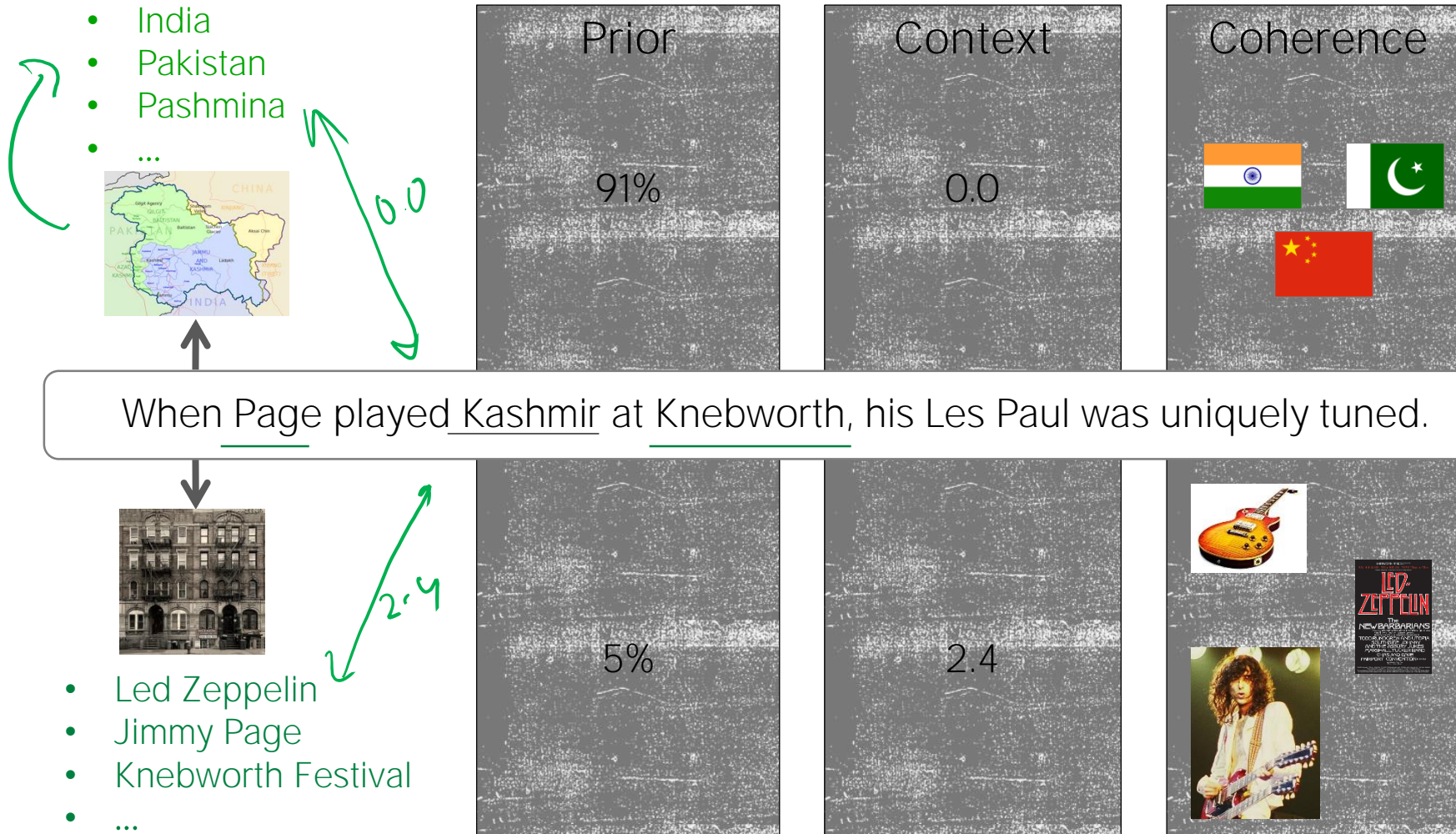
Images taken from Wikipedia under CC BY-SA 3.0

Entities in knowledge bases



AIDA features for disambiguation

Images taken from Wikipedia under CC BY-SA 3.0



How often did “Kashmir” link to this entity in Wikipedia? Are the disambiguated entities related?

Entity keyphrases

✓ Kashmir (song)

From Wikipedia, the free encyclopedia

"**Kashmir**" is a song by the [English rock band Led Zeppelin](#) from their sixth album [Physical Graffiti](#), released in 1975. It was written by [Jimmy Page](#) and [Robert Plant](#) (with contributions from [John Bonham](#))

Link Anchor Texts

References

16. [^ "The 100 Greatest Rock Songs of All Time - July 2000" !\[\]\(79de0df6c6ddd2d4eb74f1cc5f48ec50_img.jpg\). VH1. Retrieved 2009-02-10.](#)

Citation Titles

- [The 100 Greatest Songs of All Time - November 2003" !\[\]\(e492b5d52ab457a7a3c2826c4091dfee_img.jpg\). Retrieved 2009-02-10.](#)

Category Names

Categories: [1975 songs](#) | [Led Zeppelin songs](#) | [Songs written by Jimmy Page](#) | [Songs written by Robert Plant](#) | [Songs by John Bonham](#) | [English-language songs](#)

Knebworth Festival

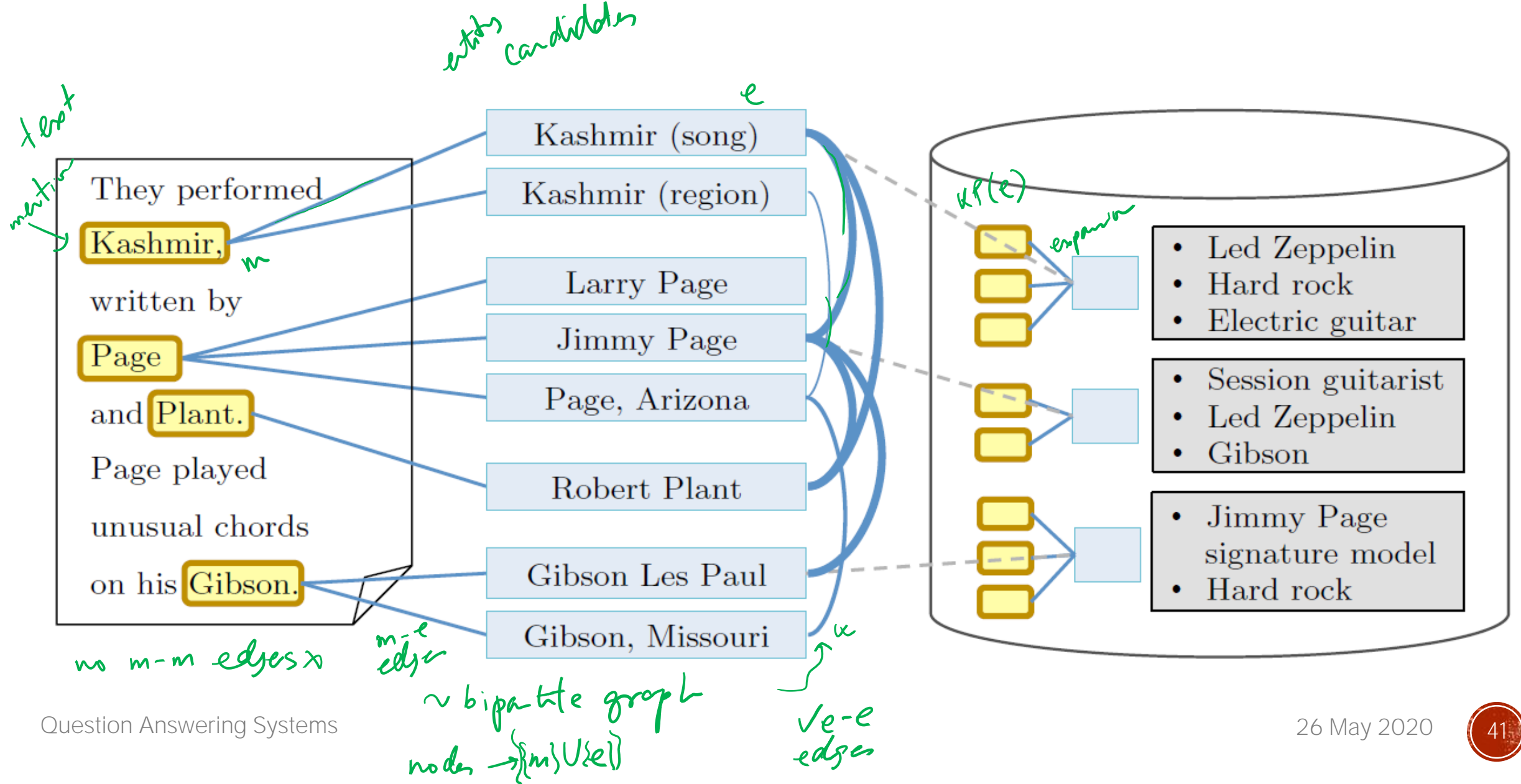
Led Zeppelin Remasters

John Paul Jones

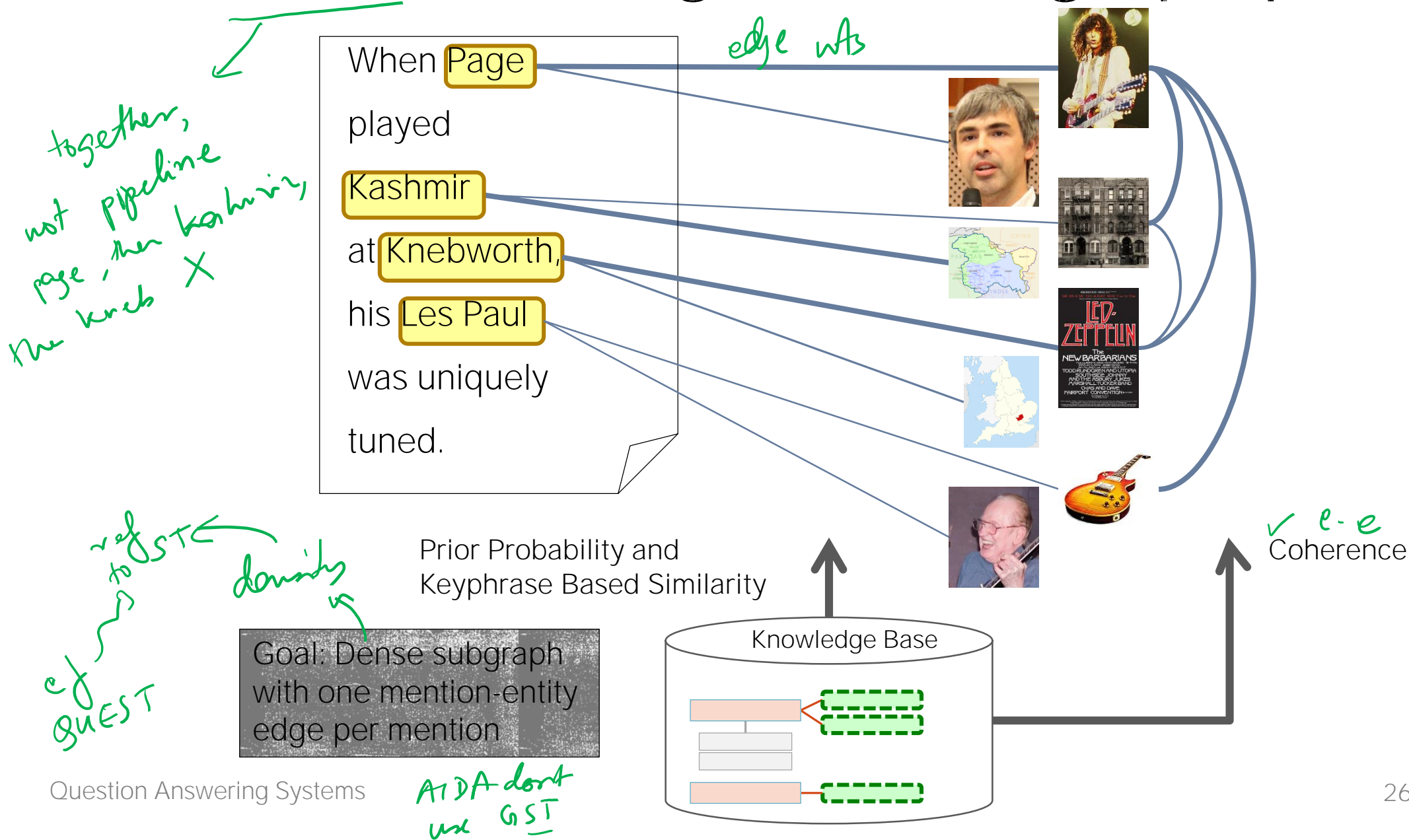
Mellotron

Titles of Linking Articles

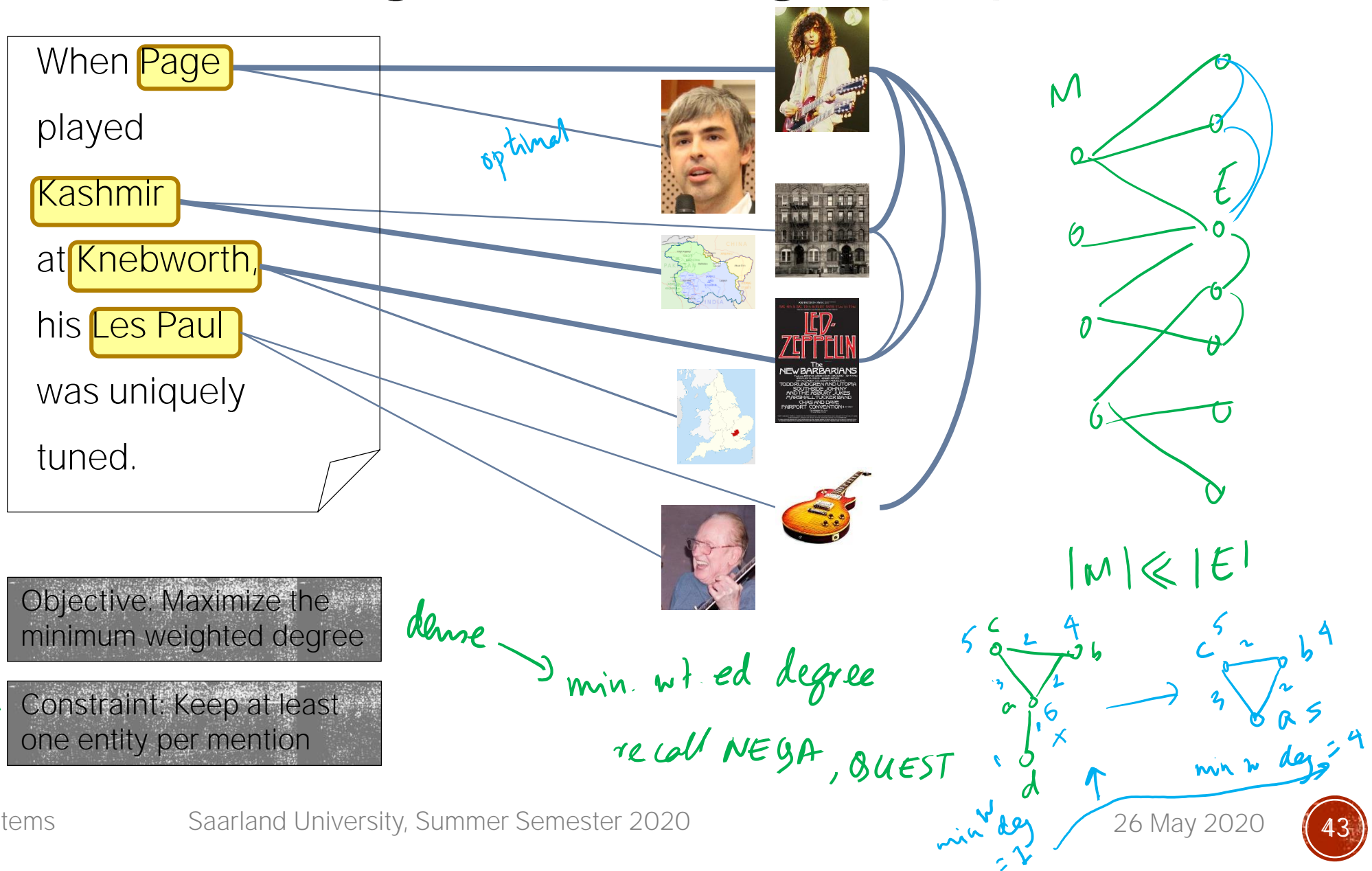
Mention-entity graph example



AIDA: Joint disambiguation as graph problem



AIDA: Joint disambiguation as graph problem



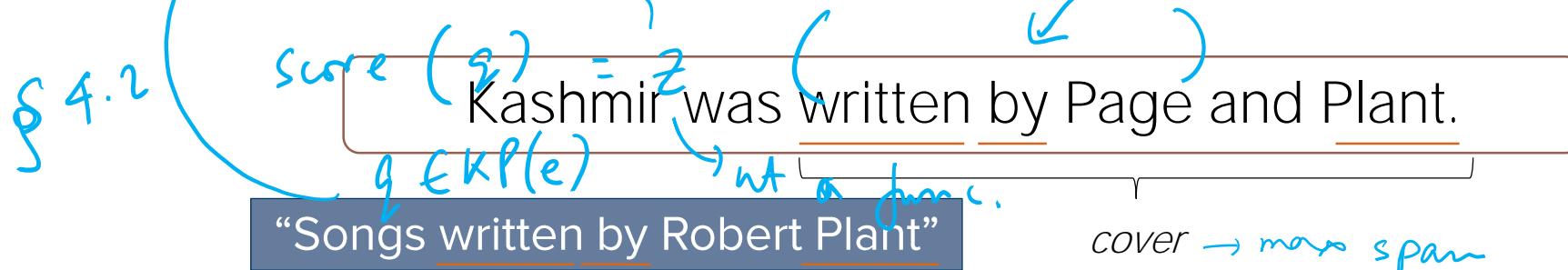
Keyphrase-based similarity

- Keyphrases (kp) commonly occur only partially

$$score(kp) = \frac{\overset{+1}{\# \text{ matching words}}}{\text{length of cover}(kp)} \left(\frac{\overset{+2}{\sum_{w \in \text{cover}} \text{weight}(w)}}{\sum_{w \in kp} \text{weight}(w)} \right)^2$$

Account for partial matches

Weight of contained tokens w



- To score an entity, all keyphrase scores are summed

$$sim\ score(m, e) = \sum_{q \in KP(e)} score(q)$$

Keyword weighting

$$score(kp) = \frac{\# \text{ matching words}}{\text{length of cover}(kp)} \left(\frac{\sum_{w \in cover} weight(w)}{\sum_{w \in kp} weight(w)} \right)^2$$

wt. (w)

wt (w)

- Global IDF of a keyphrase token w in Wikipedia *~ IRDM*
- Mutual Information of a token w and an associated entity
 - How often does the token occur in the keyphrase set of an entity?

word association measure

https://en.wikipedia.org/wiki/Mutual_information

*MI
PMI
nPMI*

Disambiguation by joint inference

Input

- Mentions
 - context of mention $cxt(m)$
 - entity candidates $e + cxt(e)$

Features

Prior	$prior(m, e)$
Similarity	$sim(cxt(m), cxt(e))$
Coherence	$coh(e_1, e_2)$

Goal

$$\alpha \cdot \sum_{i=1}^k prior(m_i, e_{j_i}) + \beta \cdot \sum_{i=1}^k sim(cxt(m_i), cxt(e_{j_i})) + \gamma \cdot coh(e_{j_1}, e_{j_2}, \dots, e_{j_k})$$

= max!

notation
 $m_i \leftrightarrow e_{j_i}$ arg max (score)
 over what?

not used in TAG ME

obj. fn.

Greedy graph algorithm

- Input: weighted graph of mentions and entities
- Output: result graph with maximum density
- Objective: maximize the minimum weighted degree
- Constraint: keep at least one entity per mention

key idea

1. Prune entities that are too distant from all mentions
2. While an entity can be removed, remove the one with the lowest weighted degree
 - Keep graph with best minimum weighted degree

Final steps

- Find subgraph maximizing total edge weight
 - If graph is small enough, enumerate all potential mention-entity mappings
 - Otherwise do local search, randomly switching mention-entity mappings for a fixed number of times

Robustness issues

§5.3

Prior may be misleading

Given the prior probability for all entity candidates, only use prior when very good indicator for one single entity (>90%)

Coherence can get hooked to a wrong subgraph

Given prior probability and similarity distribution for all entity candidates. If they are reasonably similar, fix entity for mention before running the graph algorithm.

Dataset is available at <http://www.mpi-inf.mpg.de/yago-naga/aida>

Conclusions

- Understanding entities is vital to QA
- NERD is a vital (first) cog in the QA wheel
- Often used off-the-shelf
- Many, many innovative techniques
- Mainly based on priors, ^{m-e}similarity and coherence
- Applications span beyond QA!

*Thank
you*