



 ntent™

Agenda

- WSDM's Birth
- Analysis of Ten Years
- Current & Future Challenges
 - Contextual Multilingual Semantic Search
 - Bias in the Web

WSDM's Birth and Governance (1/2)

- **May 12, 2006:** e-mail from Ziv Bar-Yossef (Google Haifa) and Junghoo Cho (UCLA)
- **Initial SC:**

Rakesh Agrawal (Microsoft)	Ricardo Baeza-Yates (Yahoo)	Krishna Bharat (Google)
Andrei Broder (Yahoo)	Soumen Chakrabarti (IIT Bombay)	Monika Henzinger (EPFL, Google)
Jon Kleinberg (Cornell)	Rajeev Motwani (Stanford)	Prabhakar Raghavan (Yahoo)
- **May 25, 2006:** meeting at WWW 2006 in Edinburgh.
Most of the SC plus Sue Dumais (Microsoft Research), Ravi Kumar & Andrew Tomkins (Yahoo), Ronny Lempel (IBM), Yoelle Maarek (Google), Mark Manasse (Microsoft), Marc Najork (Microsoft), and Torsten Suel (Polytechnic University)
- **SIGs proposed:** SIGACT, SIGIR, SIGKDD, SIGMOD, and SIGWEB
- **Plan A: September 2007,** with the advice from ACM President, Stu Feldman, to request the SIGWEB sponsorship.
- **Early 2007 conversations:** SIGWEB (Ronny) and SIGMOD (Andrew), plus SIGIR & SIGKDD (myself).
- Asked to be conference chair for the first conference with Andrei and Soumen as PC-Chairs, Ravi as treasurer, and Utkarsh Srivastava of Stanford as local organizer



www.ntent.com | @withntent | 877.861.2230

3

WSDM's Birth and Governance (2/2)

- Our plans changed during my negotiation with SIGIR (Jamie Callan), as the conference was between SIGIR and CIKM, two SIGIR sponsored conferences
- **Plan B: second week of February 2008,** with the second conference planned in **Barcelona**.
As I was the natural chair for 2009, I stepped down for 2008, being replaced by Marc Najork
- **Driving team:** Andrei, Andrew, Junghoo, Marc, Ravi, Ricardo, Ronny, Soumen & Zvi
- **April 2007:** formal support of SIGKDD and SIGMOD
However SIGIR wanted to have a steering committee with its own representative and clear governance rules
- **May 2007:** Andrei, Andrew, Marc, Ravi, Ricardo, Soumen, and Ziv
discussed the governance rules at WWW 2007 in Banff, Canada
- **Steering committee rules:** 8 members, 4 of them representing the sponsoring SIGs,
balancing industry and academia members, with periods of 4 years
- With the new rules, SIGIR and SIGWEB became sponsors too



www.ntent.com | @withntent | 877.861.2230

4

WSDM's Formal Steering Committee

June 2007:

Rakesh Agrawal (Microsoft, KDD rep.)	Ricardo Baeza-Yates (Yahoo, SIGIR rep., chair)
Ziv Bar-Yossef (Google)	Soumen Chakrabarti (IIT Bombay)
Monika Henzinger (EPFL, Google)	Jon Kleinberg (Cornell)
Rajeev Motwani (Stanford)	Prabhakar Raghavan (Yahoo)

January 2010: Marc Najork (Microsoft) replaced Prabhakar, and soon after
Hector Garcia-Molina (Stanford) joined representing SIGMOD

May 2013: SC Renewal, almost 2 years late!

Ricardo Baeza-Yates (Yahoo, SIGIR rep., chair)	Paolo Boldi (Univ. de Milano)
Andrei Broder (Google)	Brian Davison (Lehigh Univ., SIGWEB rep.)
Nick Koudas (Univ. of Toronto, SIGMOD rep.),	Bing Liu (UIC, SIGKDD rep.)
Marc Najork (Microsoft)	

January 2014: Hang Li (Huawei) joined to complete the SC and soon after I proposed
Marc Najork (Google) as new chair, who was elected by acclamation!

2017: SC should be partially renewed



www.ntent.com | @withntent | 877.861.2230

5

WSDM's History

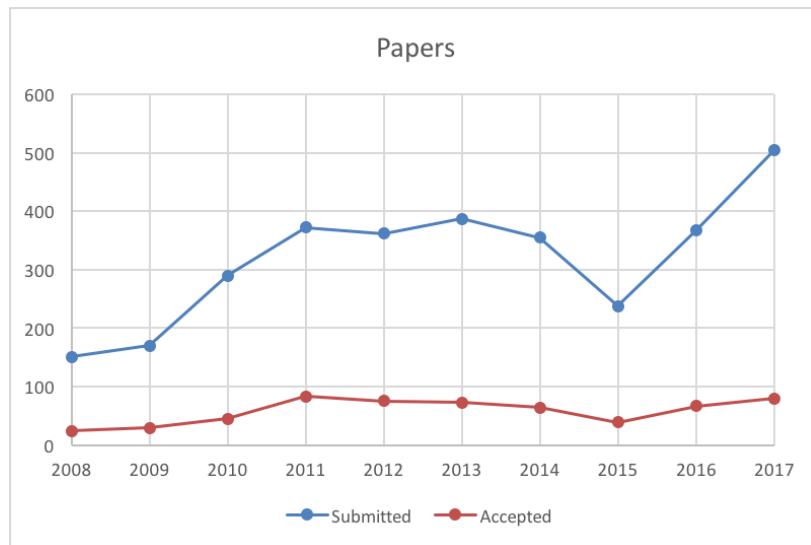
Year	Location	Conference Chair	PC Chairs
2008	Palo Alto	Marc Najork	Andrei Broder; Soumen Chakrabarti
2009	Barcelona	Ricardo Baeza-Yates	Paolo Boldi; Berthier Ribeiro-Neto
2010	New York	Brian Davison; Torsten Suel	Nick Craswell; Bing Liu
2011	Hong Kong	Irwin King	Wolfgang Nejdl; Hang Li
2012	Seattle	Eytan Adar; Jaime Teevan	Eugene Agichtein; Yoelle Maarek
2013	Rome	Stefano Leonardi; Alessandro Panconesi	Paolo Ferragina; Aristides Gionis
2014	New York	Ben Carterette; Fernando Diaz	Carlos Castillo; Donald Metzler
2015	Shanghai	Xueqi Cheng; Hang Li	Evgeniy Gabrilovich; Jie Tang
2016	San Francisco	Paul Bennet; Vanja Josifovski	Jennifer Neville; Filip Radlinski
2017	Cambridge	Milad Shokouhi; Maarten de Rijke	Andrew Tomkins; Min Zhang



www.ntent.com | @withntent | 877.861.2230

6

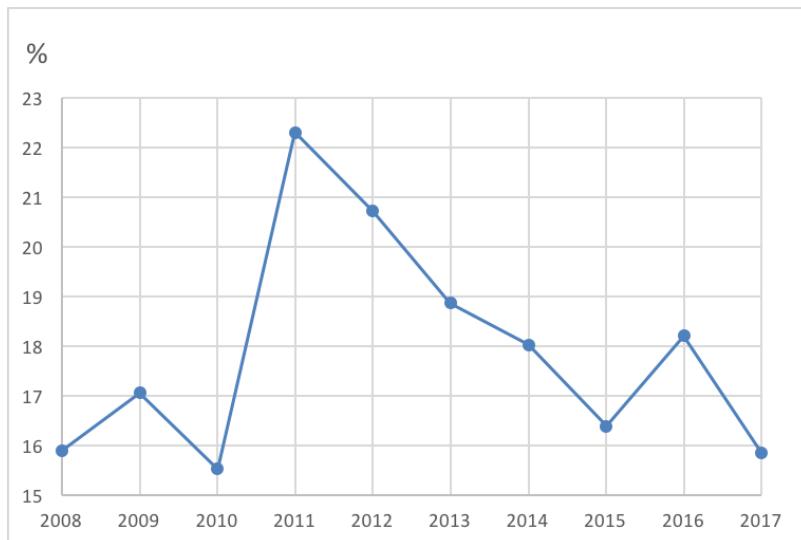
WSDM's Growth: Papers



www.ntent.com | @withntent | 877.861.2230

7

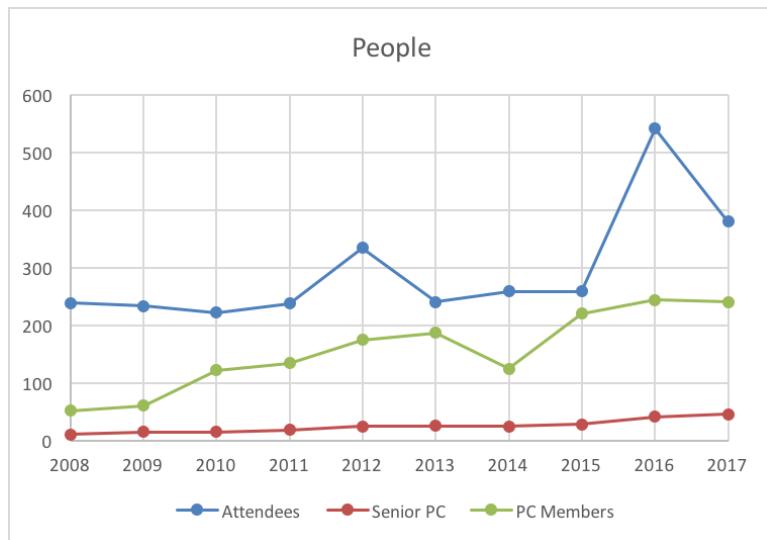
WSDM's Quality: Acceptance Rate



www.ntent.com | @withntent | 877.861.2230

8

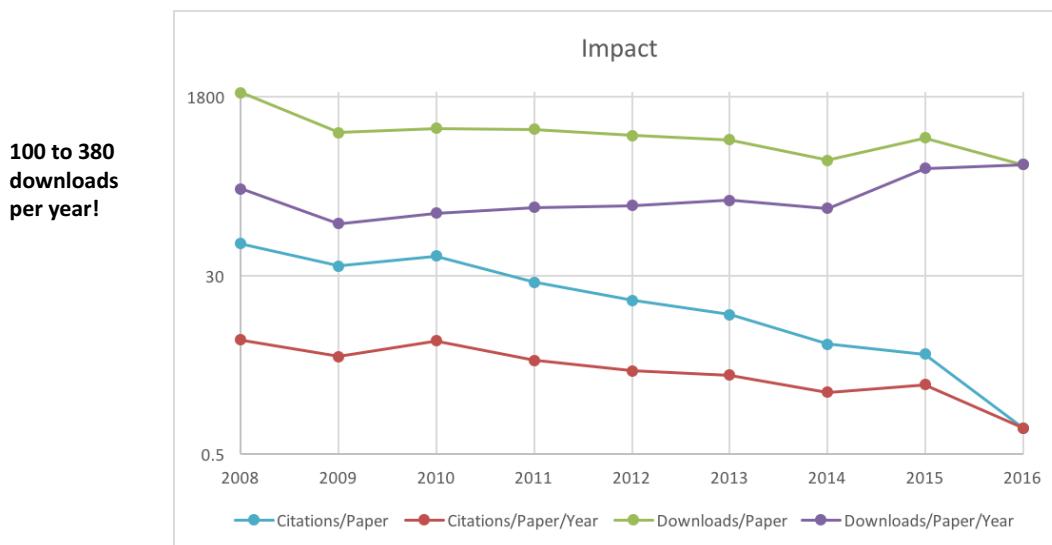
WSDM's Growth: People



www.ntent.com | @withntent | 877.861.2230

9

WSDM's Impact: ACM Citations & Downloads (Dec 2016)



www.ntent.com | @withntent | 877.861.2230

10

WSDM's Places

2008	2011	2014	2017
University Research Microsoft Yahoo IBM Laboratory Indiana Illinois Technology State	University Research Yahoo Microsoft Laboratories Institute Stanford Technology Corporation Chinese	University Research Microsoft Institute Google UIUC Technology Yahoo Tsinghua CMU	University Research Technology Microsoft Science(s) California State Yahoo UIUC UI at Chicago



www.ntent.com | @withntent | 877.861.2230

11

WSDM's Topics

2008	2011	2014	2017
Search Web Ranking Mining Classification Advertising Models Document Graph Analysis	Web Search Social Mining Query Analysis Temporal Online Learning Data	Search Web Data Social Networks Modeling Learning Systems Click Advertising Mining	Social Networks Search Learning Time IR Modeling Text Data Mining Recom. Systems Embedding



www.ntent.com | @withntent | 877.861.2230

12

Current & Future Challenges

- Contextual Multilingual Semantic Search
- Fake web content
- Bias in the Web



www.ntent.com | @withntent | 877.861.2230

13

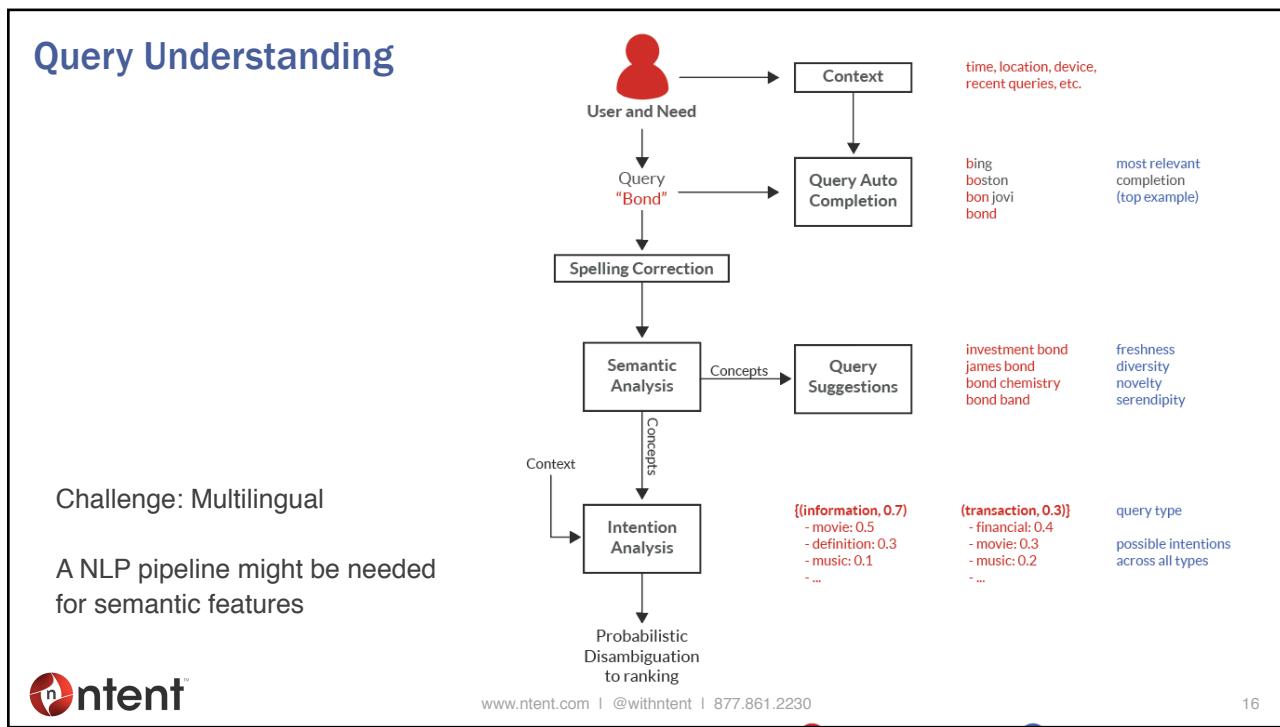
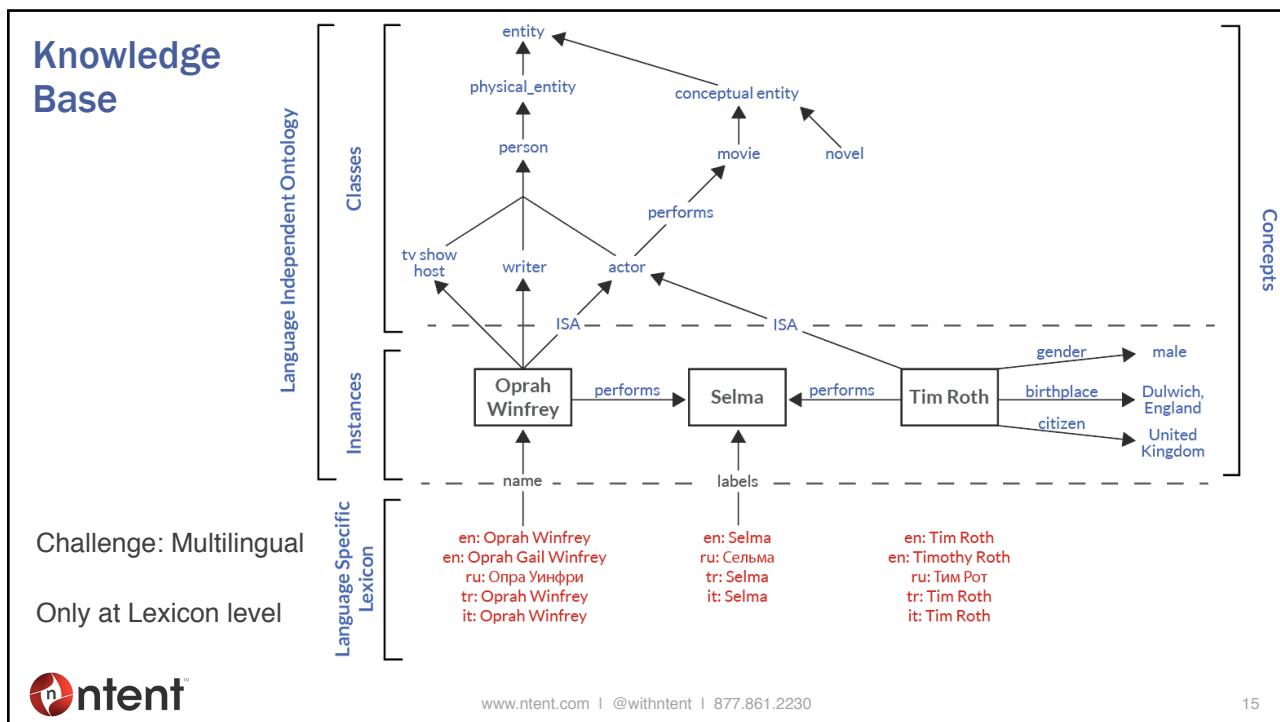
Contextual Multilingual Semantic Search

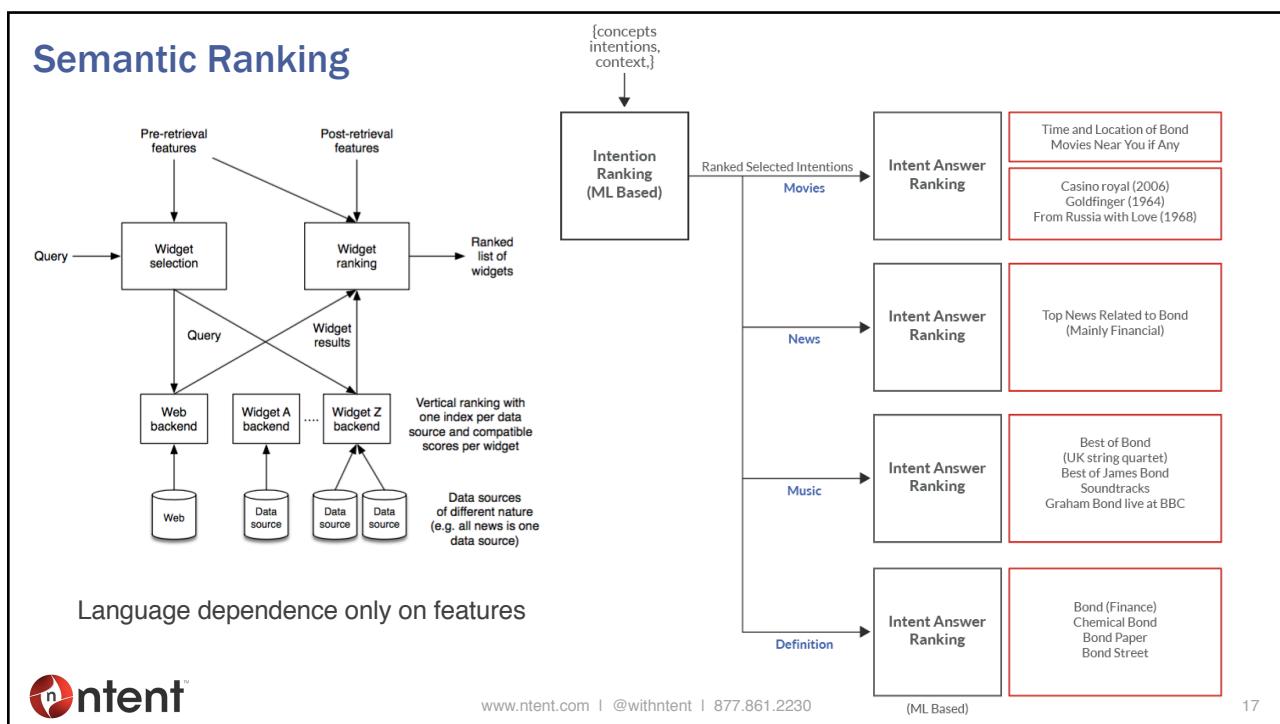
- Multilingual Knowledge Base
- Contextual Query Understanding
- Semantic Ranking



www.ntent.com | @withntent | 877.861.2230

14





Fake Content & Bias

- British Prime Minister Benjamin Disraeli:
 - "There are three kinds of **lies**: **lies**, damned **lies**, and **statistics**.

UTC professor says "Everyone has bias"

BY HANNAH LAWRENCE | FRIDAY, JULY 8TH 2016

We all have biases and preconceptions about certain subjects or groups of people according to one Chattanooga researcher.

Buzzfeed News

TOP POST
173,877 VIEWS

'16

Here Are 50 Of The Biggest Fake News Hits On Facebook From 2016

One fake news entrepreneur says we should expect even more Trump hoaxes in 2017

posted on Dec. 30, 2016, at 2:12 p.m.

Craig Silverman
BuzzFeed News Media Editor

Bias: significant deviation from a prior (unknown) distribution

www.ntent.com | @withntent | 877.861.2230

19

(Observational) Human Data has Bias

Goal: Bias Awareness

- o Gender
- o Racial
- o Sexual
- o Religious
- o Social
- o Linguistic
- o Geographic
- o Political
- o Educational
- o Economic
- o Technological

- from Noise or Spam
- Validity (e.g. temporal)
- Completeness
- Gathering process
-

Attempt of an unbiased (personal) view on bias in the Web

Many people extrapolate results of a sample to the whole population (e.g., social media analysis)

In addition there is bias when measuring bias as well as bias towards measuring it!

ntent™

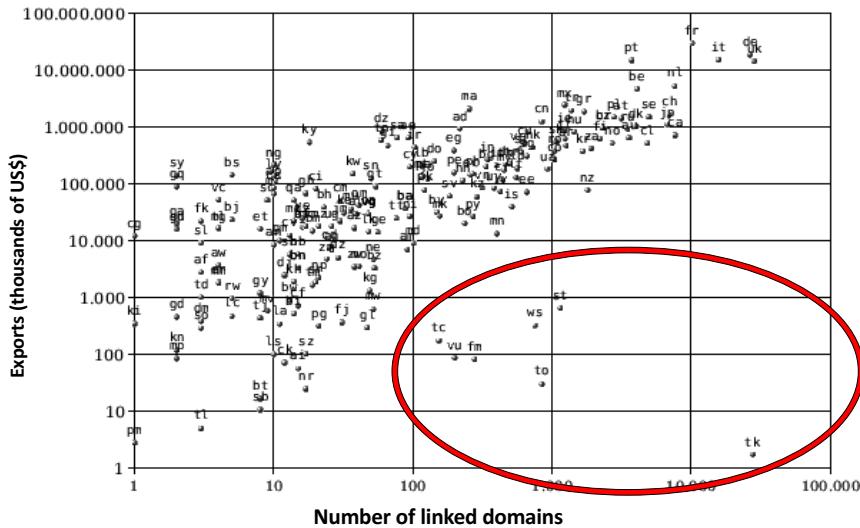
Bias in the Web

Web

Data bias

ntent™

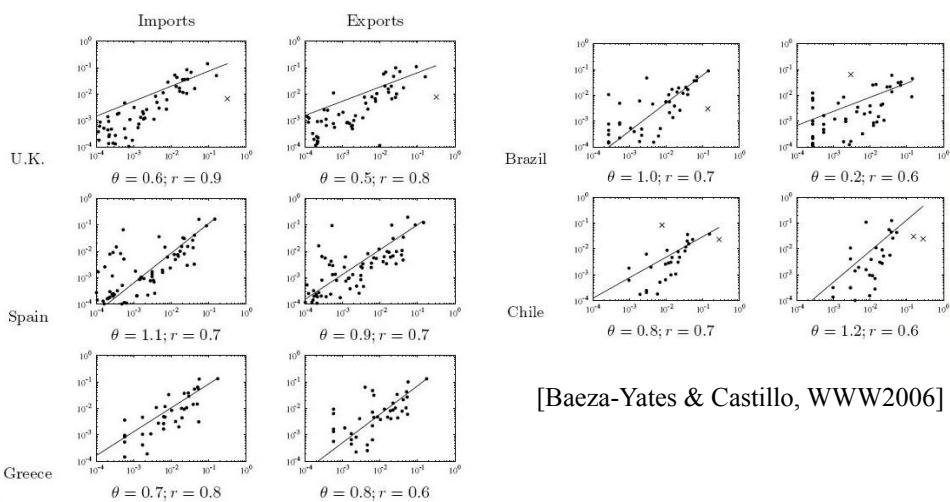
Economic Bias in Links



[Baeza-Yates, Castillo & López. Characteristics of the Web of Spain. Cybermetrics, 2005]

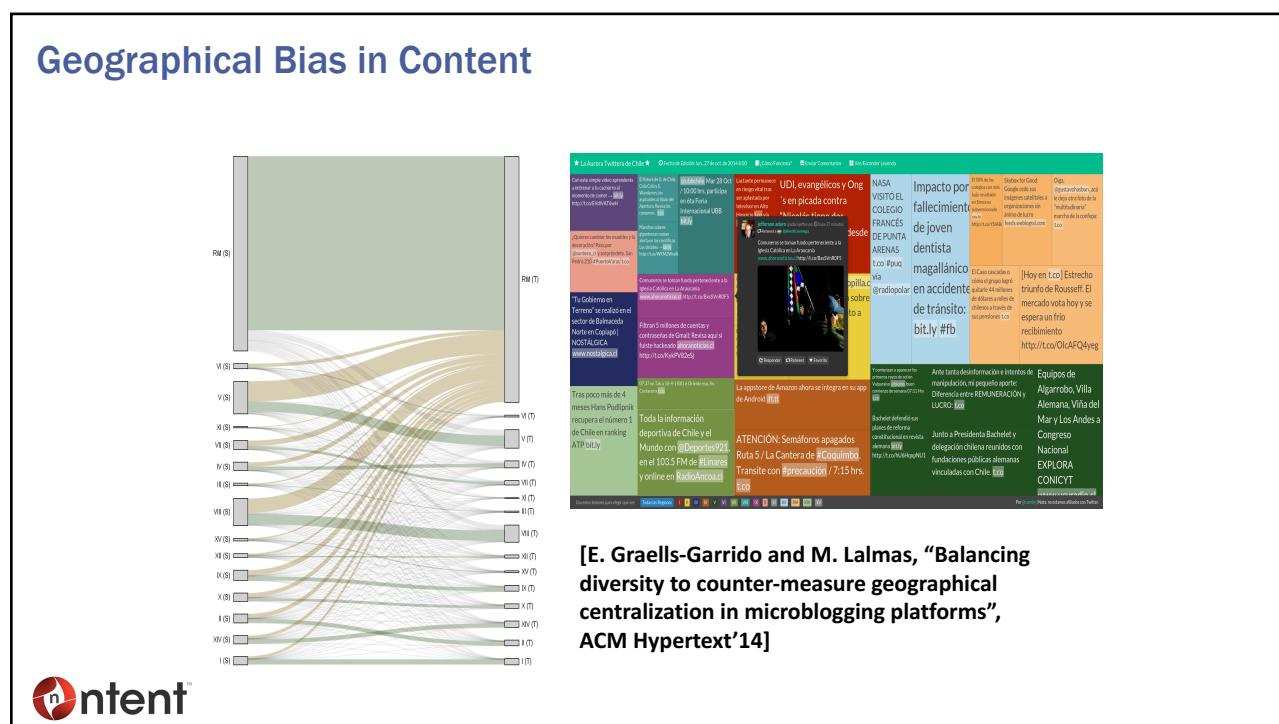
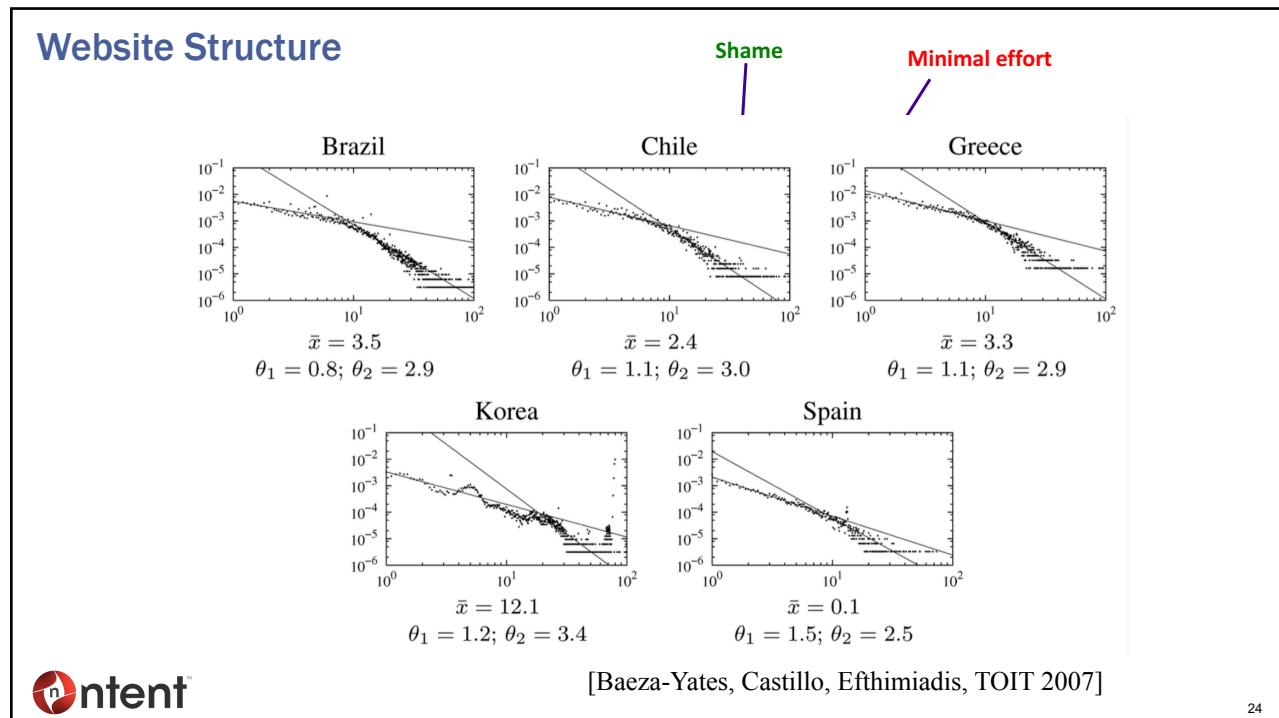
22

Economic Bias in Links



ontent[™]

23



Gender Bias in Content

- Word embedding's in w2vNEWS

Gender stereotype *she-he* analogies.

sewing-carpentry	register-nurse-physician	housewife-shopkeeper
nurse-surgeon	interior designer-architect	softball-baseball
blond-burly	feminism-conservatism	cosmetics-pharmaceuticals
giggle-chuckle	vocalist-guitarist	petite-lanky
sassy-snappy	diva-superstar	charming-affable
volleyball-football	cupcakes-pizzas	hairdresser-barber

Gender appropriate *she-he* analogies.

queen-king	sister-brother	mother-father
waitress-waiter	ovarian cancer-prostate cancer	convent-monastery

Most journalists are men?

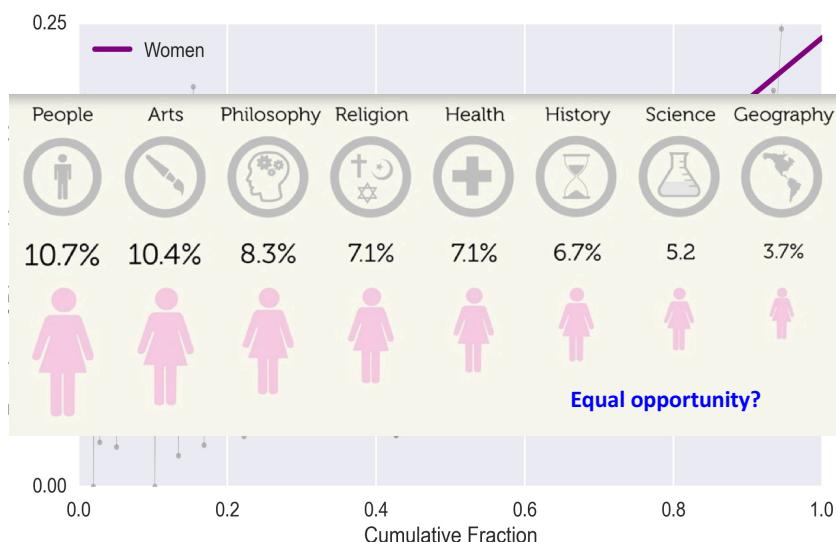
[Bolukbasi et al, ArXiv 2016]

Yes, about 60 to 70% at work
although at college is the inverse



Gender Bias in Content

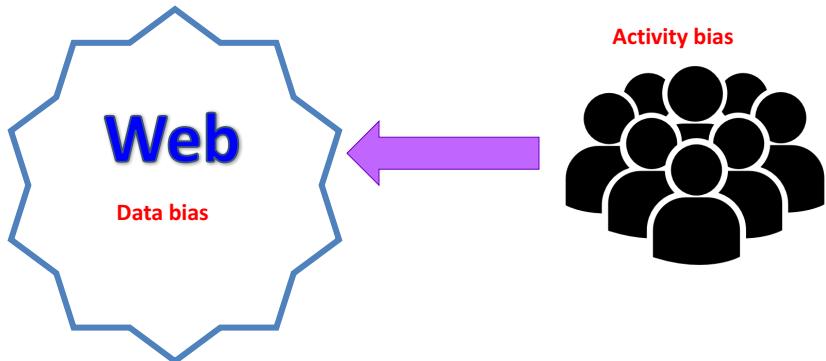
Systemic bias?



[E. Graells-Garrido et al., "First Women, Second Sex: Gender Bias in Wikipedia", ACM Hypertext'15]

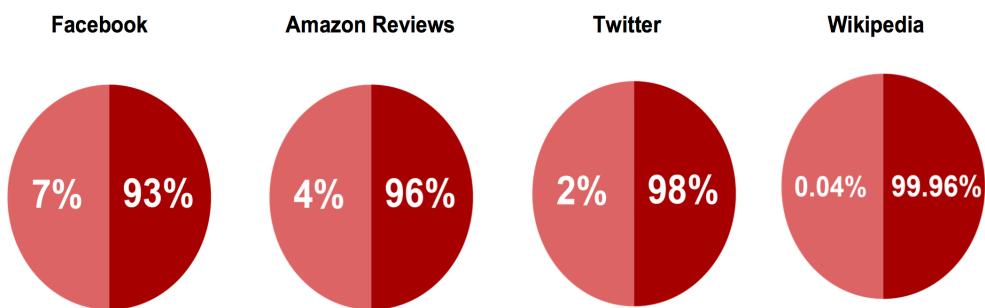


Bias in the Web



Activity Bias

Which percentage of users produce 50% of the content?



[Baeza-Yates & Saez-Trumper, ACM Hypertext 2015]



the guardian

sport football opinion culture business lifestyle fashion environment tech travel ≡ all sections

Amazon sues 1,000 'fake reviewers'

October 2015

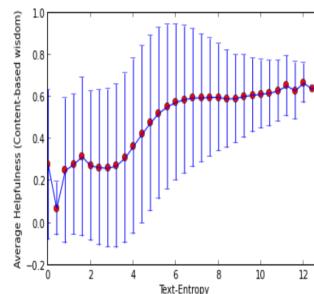
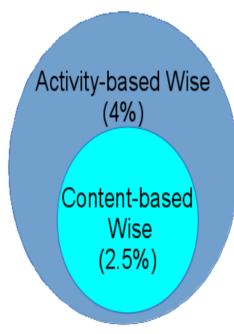
Online retailer files lawsuit in US against people whose names it says it does not know, claiming they offer reviews for sale

Amazon Continues Their Crusade Against Fake Reviews

By Tyler Lee on 04/26/2016 05:07 PDT

Quality of Content?

- Adding content implies adding wisdom?
- We used Amazon's reviews helpfulness and computed the text entropy
- Content-based-wise users
- How many of those users are being paid?



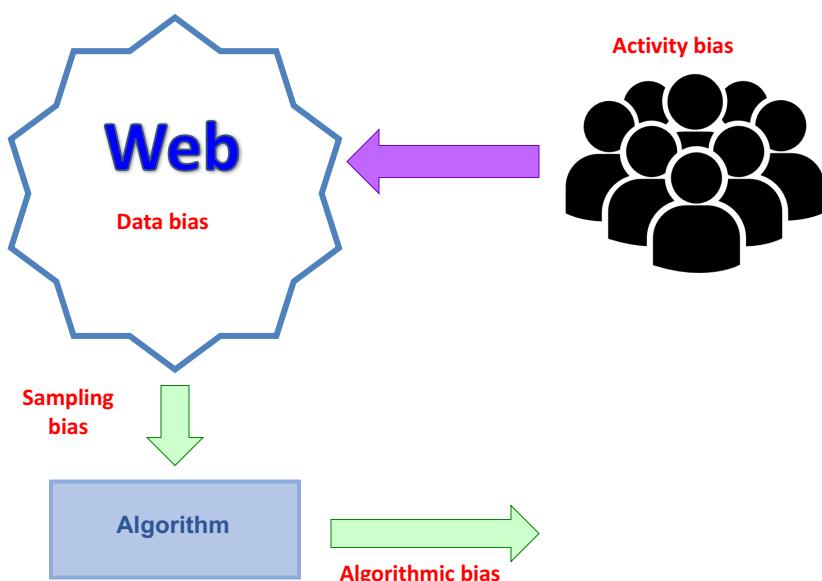
[Baeza-Yates & Saez-Trumper, ACM Hypertext 2015]

Content that is never seen: Digital Desert

- 1.1% of the Twitter content is never seen.*
- 31% of articles added/edited in May 2014 in wikipedia, were not visited in June.



[Baeza-Yates & Saez-Trumper, ACM Hypertext 2015]



Sample Size?

- If we want to estimate the frequency of queries that appear with probability at least p with a certain relative error ϵ we can use the standard binomial error formula $\sqrt{(1-p)/np}$ which works well for p near $1/2$ **but not for p near 0**
- Better is the Agresti-Coull technique (also called *take 2*) which gives:

$$n \geq Z_{1-\alpha/2}^2 \left(\frac{p'(1-p')}{\epsilon^2} - 1 \right)$$

where Z is the inverse of the standard normal distribution, $1 - \alpha$ is the confidence interval and $p' = p + Z^2/2$

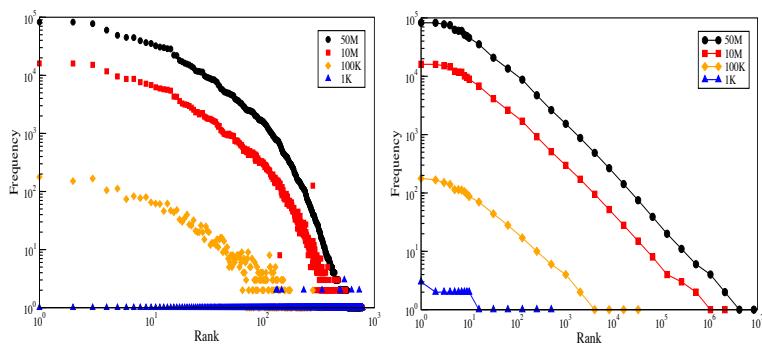
- If $p = 0.1$, $1 - \alpha$ is 90% and ϵ is 10%, we get $n = 2342$.
The standard formula gives $n = 900$!



[Baeza-Yates, SIGIR 2015, Industry track]

Sampling Techniques

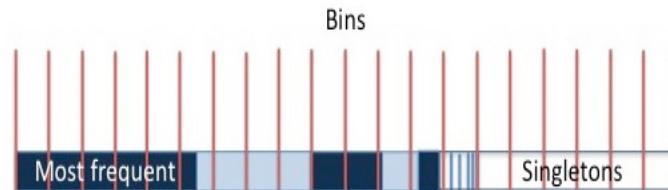
- Standard technique: $p_q \approx \hat{p}_q(\mathcal{S}) = \frac{f_q(\mathcal{S})}{\sum_{q' \in \mathcal{S}} f_{q'}(\mathcal{S})}$
- A good sample should cover well all the query distribution but this does not work with very biased distributions.



[Zaragoza et al, CIKM 2010]

Incremental Stratified Sampling

- Main goal: make good samples consistent across time
- Simple idea based in stratified sampling: bins + random start point

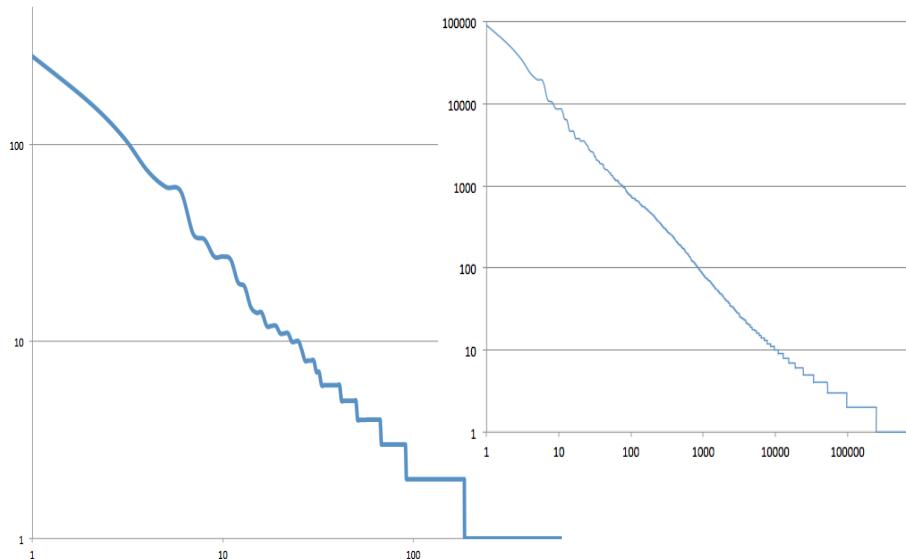


- Bin size can be found by binary search starting with a good approximation if a query frequency model is used ($b < V/n$)
- This perfectly mimics the head of the distribution, but not the tail
- Change the bins in the tail to get the right distribution



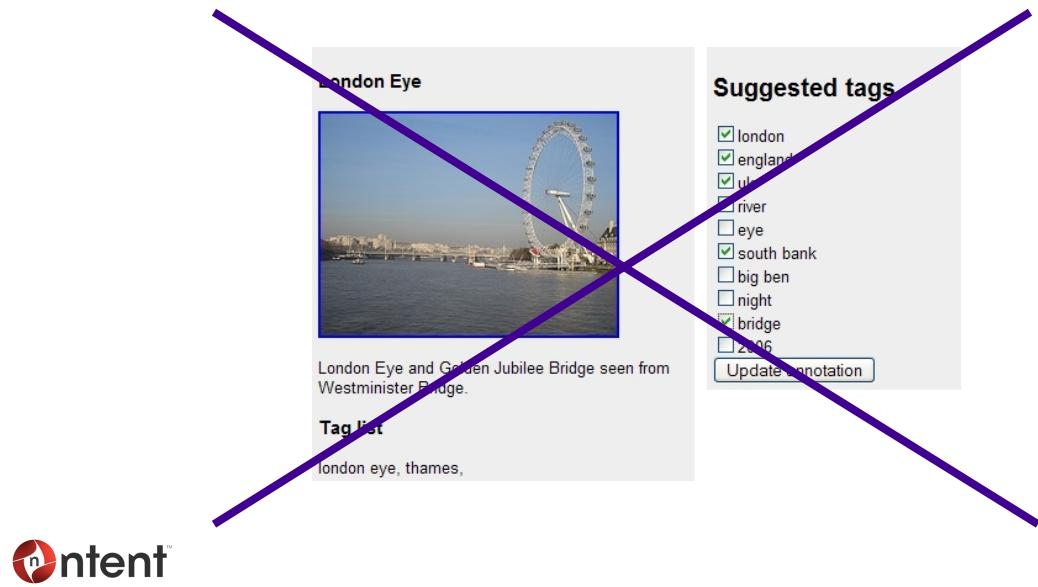
[Baeza-Yates, SIGIR 2015, Industry track] 37

Stratified Sampling Example

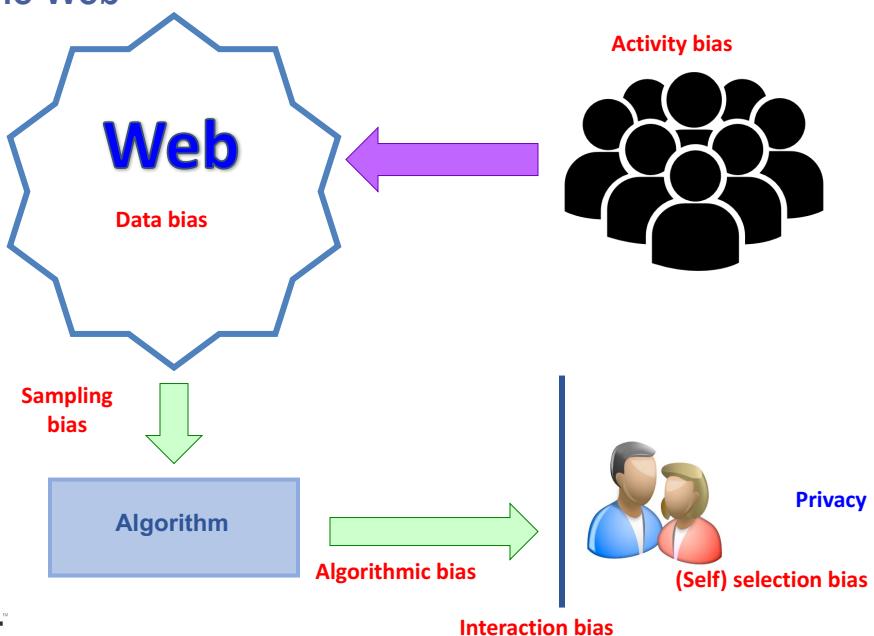


39

Extreme Algorithmic Bias



Bias in the Web



Bias in the Interaction

Related Searches: tennis racket, tennis shoes.

Shop by Category

Tennis Equipment

Tennis Games

Kids' Sports

Clothing, Shoes & Jewelry

Tennis - Books

Position bias

Ranking bias

Wilson Sporting Goods Championship Extra Duty Tennis Balls (1-Can)
Jun 14, 2012
by Wilson

\$2.79 \$6.99 Add-on Item
Add to a qualifying order to get it by **Tomorrow, May 6**

More Buying Choices
\$0.99 new (18 offers)
\$7.99 used (2 offers)
See newer version

Presentation bias

Tennis Elbow Brace with Gel Comp...
\$24.50 ✓Prime
★★★★★ 7

DIMANKA Professional Table Tennis...
\$34.99
★★★★★ 9

Gamma Quick Kids 78 Ball (12 Pac...
\$19.99 ✓Prime
★★★★★ 44

Social bias

Wilson 75 Tennis Ball Pick Up Hopper
by Wilson

\$19.96 ✓Prime
Get it by **Tomorrow, May 6**

More Buying Choices
\$18.88 new (11 offers)
\$35.00 used (1 offer)

Interaction bias

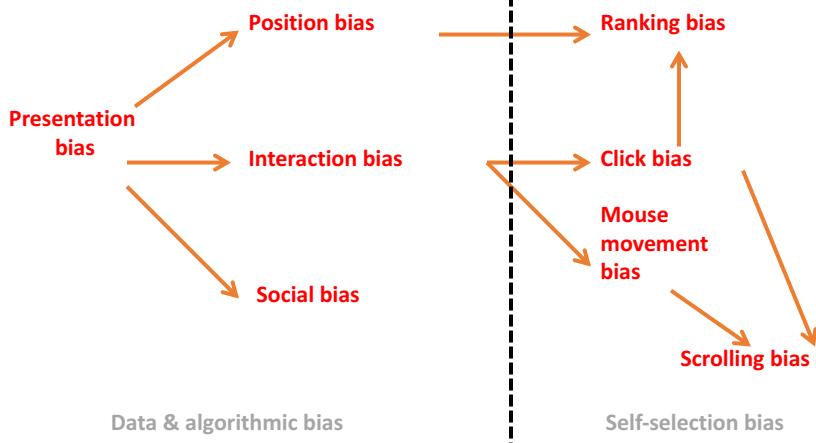
Product Features
Holds 75 tennis balls with a special no spill lid (Tennis Balls NOT included)

Sports & Outdoors: See all 60,449 items

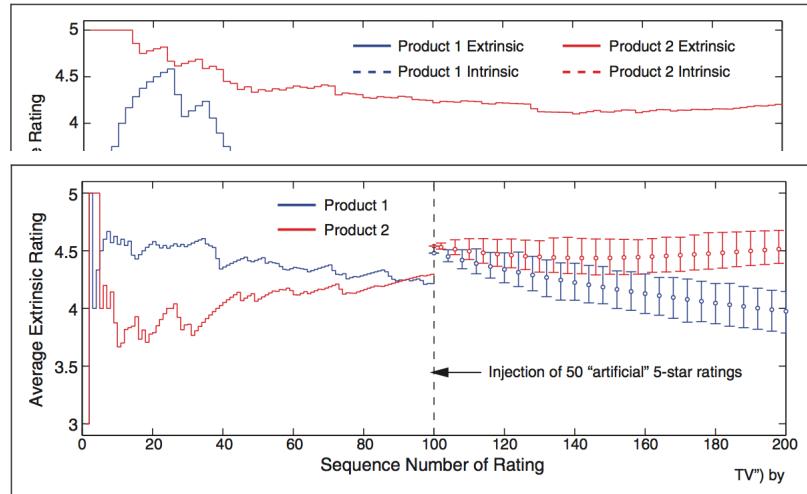
Best Seller



Dependencies: A Cascade of Biases!



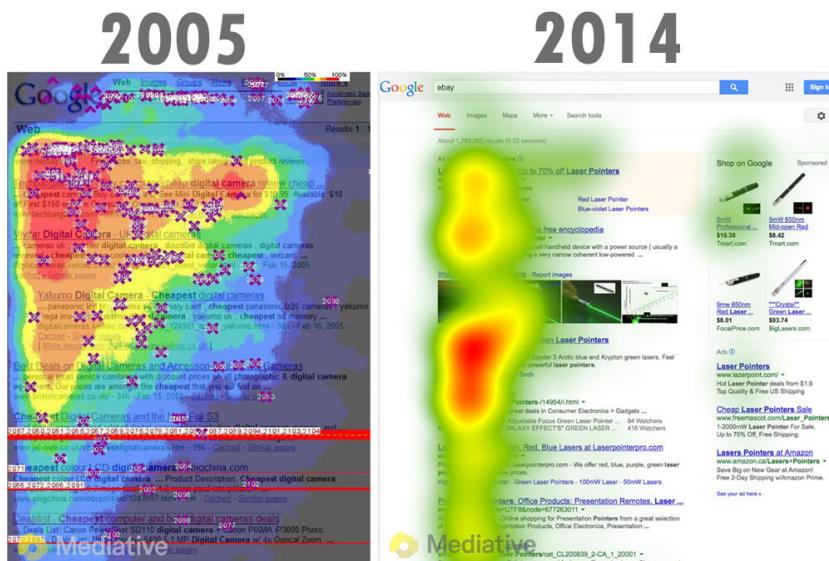
Social Bias



[WHY AMAZON'S RATINGS MIGHT MISLEAD YOU; The Story of Herding Effects
Ting Wang and Dashun Wang, Big Data, 2014]



Ranking Bias in Web Search



[Mediative Study, 2014]



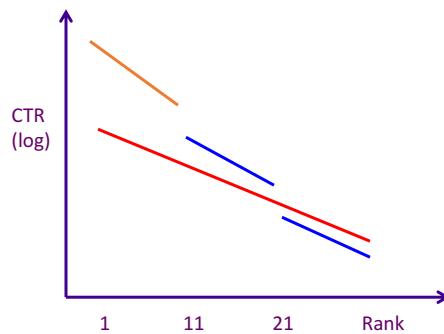
Click Bias in Web Search

- Ranking & **next page** bias



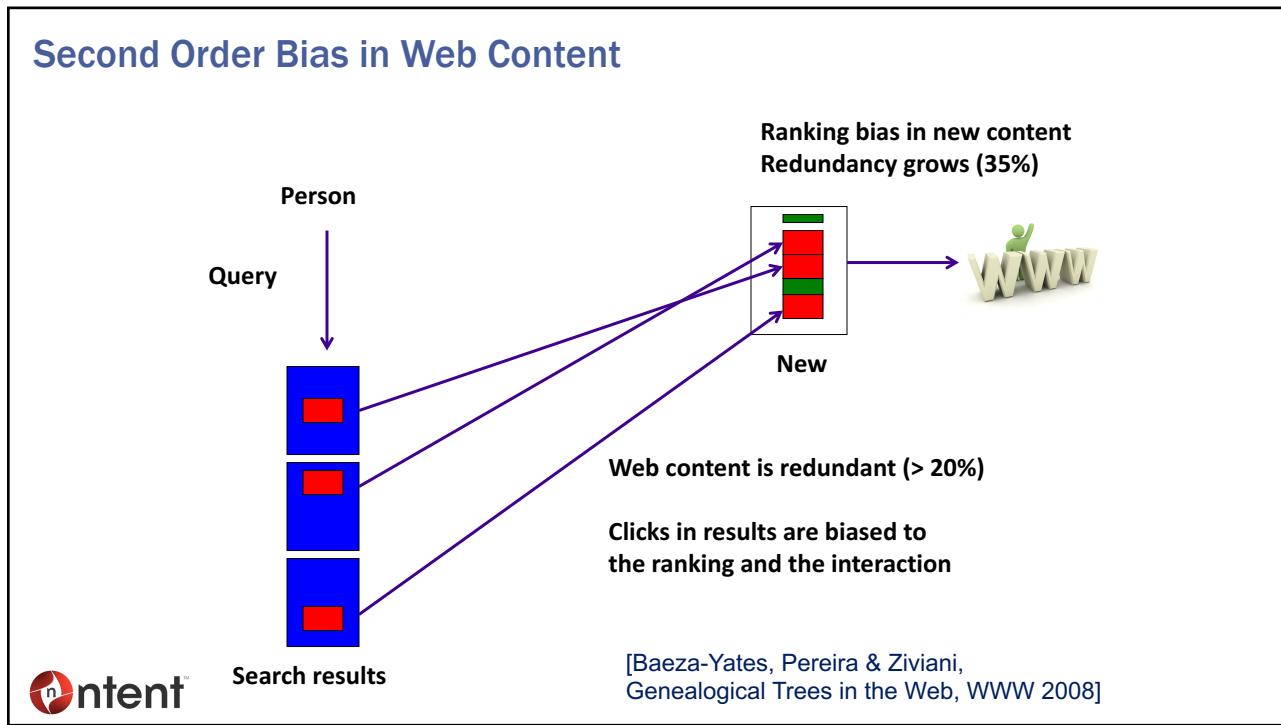
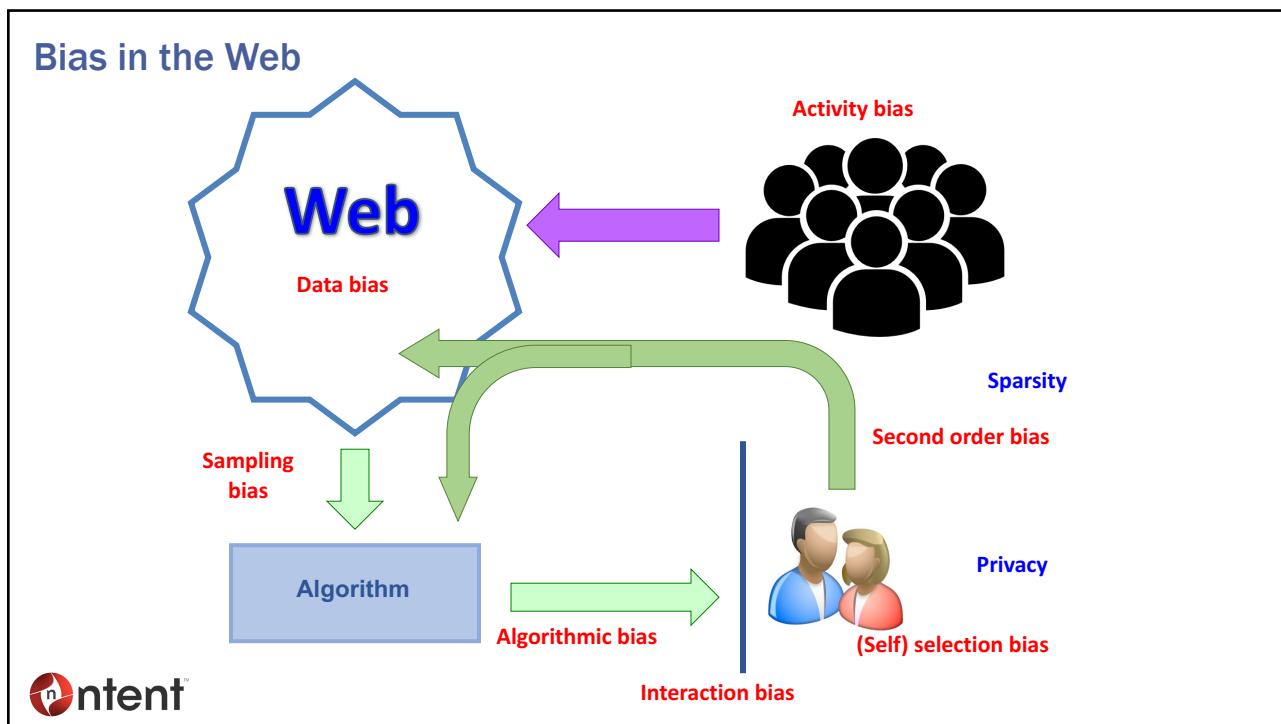
Unbiasing Search Clicks

Clicks as implicit positive user feedback



[Dupret & Piwowarski, SIGIR 2008]
[Chapelle & Zhang, WWW 2009]





Avoid Second Order Bias due to Personalization

The Filter “Bubble”, Eli Pariser (2011)

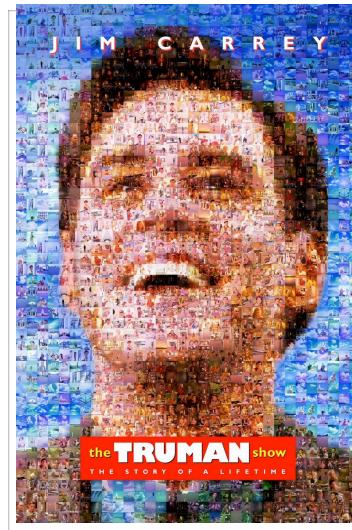
- The effect of self selection bias
- Avoid the poor get poorer syndrome
- Avoid the echo chamber
- Empower the tail

Partial solutions:

- Diversity
- Novelty
- Serendipity
- Show me the dark side

Cold start problem solution: Explore & Exploit

How much exploration is needed for presentation bias?



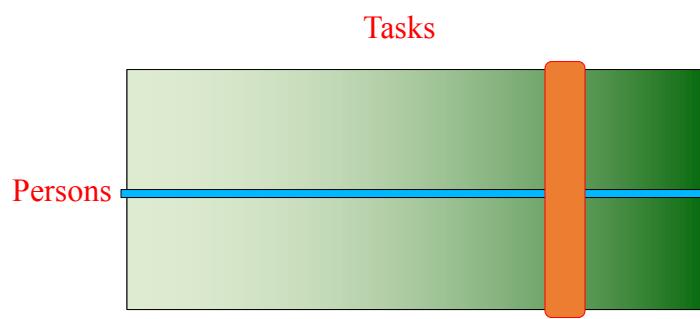
Aggregating in the Tail

- Exploit the context (and deep learning!)

91% accuracy to predict the next app you will use
[Baeza-Yates et al, WSDM 2015]

- Personalization vs. **Contextualization**

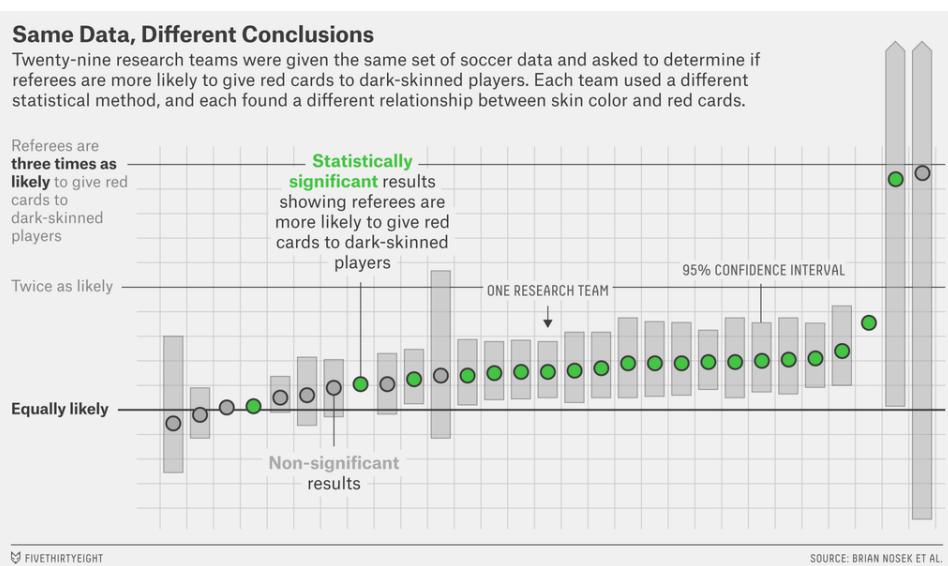
Recall that user interaction is another long tail



It's Hard to Get the Truth from Data (Professional Bias)

Same Data, Different Conclusions

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.



SOURCE: BRIAN NOSEK ET AL.



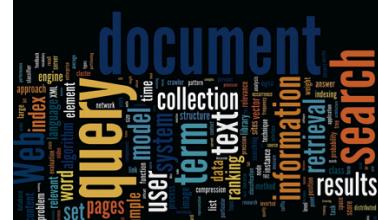
→ 61 analysts, 29 teams: 20 yes and 9 no (Univ. of Virginia, COS)

→ We need to focus on small data, not big data

Questions?

ASIST 2012 Book of the Year Award (Biased Ad)

Modern
Information Retrieval
the concepts and technology behind search
Second edition



Ricardo Baeza-Yates
Berthier Ribeiro-Neto

Contact: rbaeza@acm.org

www.baeza.cl

@polarbearby

Biased Questions?

More bias: We are hiring in Barcelona & San Diego!

