

O'REILLY®

Strata CONFERENCE

Making Data Work

O'REILLY®

Strata + HADOOP CONFERENCE + WORLD™

Tools and Techniques That Make Data Work

New York, NY

October 28-30, 2013

#strataconf

Welcome to the O'Reilly Strata Online Conference

*Big Data and the
Ethics and Challenges
of Living in a
Connected Society*

September 27, 2013
9:00 AM -11:00 AM (PT)

40

million
incremental revenue
in the first year
with MapR



90

billion
real-time auctions
per day
conducted with MapR



1.7

trillion
events processed
per month
using MapR



come
visit
booth
102

Agenda

9:00 – 9:10 Intro by Alistair Croll

9:10– 9:30 What Makes Us Human? A Tale of Advertising Fraud – Claudia Perlich (Dstillery)

9:30 – 9:50 Machine Learning Applications: Recommendation Engines Using Multiple Behavior Sources – Ted Dunning (MapR)

9:50 – 10:10 Real-time Recommendations for Retail: Architecture, Algorithms, and Design - Juliet Hougland and Jonathan Natkins (WibiData)

10:10 – 10:30 Leveling Up With Hadoop: How Blizzard's Business Intelligence Supports The Worlds of Warcraft, Starcraft, and Diablo – Amanda Gerdes (Blizzard Entertainment)

10:30 – 10:50 Interactive Visualization of "Big" Data – Jeffrey Heer (Tifacta)

10:50 – 11:00 Closing by Alistair Croll

Intro

9:00 – 9:10 am PT



Alistair Croll
Solve for Interesting
Strata Program Chair

What Makes Us Human? A Tale of Advertising Fraud

9:10 – 9:30



Claudia Perlich
Chief Scientist, Dstillery



What makes us human?

A tale of advertising fraud.



Alan Turing
(1912-1954)



“A computer would deserve to be called intelligent if it could deceive a human into believing that it was human.”

Dark Side of Display Advertising

A screenshot of a web browser window. The address bar shows "advertising spider". The main content is a news article from PCWorld AUSTRALIA. The headline reads "Click fraud botnet defrauds advertisers up to \$6 million". The article discusses a botnet called Chameleon that generates bogus clicks on display advertisements, defrauding advertisers of up to \$6 million per month. It quotes Spider.io stating that the botnet has infected about 120,000 residential computers in the U.S. and perpetrates click fraud on 202 websites that collectively deliver 14 billion ad impressions. The article also notes that click fraud cheats Web advertisers by depriving them of customers and revenue.

PCWorld AUSTRALIA

ORECK Magnesium RS JUST 7.7 POUNDS

All Categories ▾ Shop & Compare Business Centre Jobs Guides

Back to Uni Survival Guide

News ▶ Click fraud botnet defrauds advertisers up to \$6 m...

Click fraud botnet defrauds advertisers up to \$6 million

The 'Chameleon' botnet generates high traffic on low-quality websites

Jeremy Kirk (IDG News Service) — 20 March, 2013 00:55

Share +1 0 Tweet 0 Be the first to comment 0

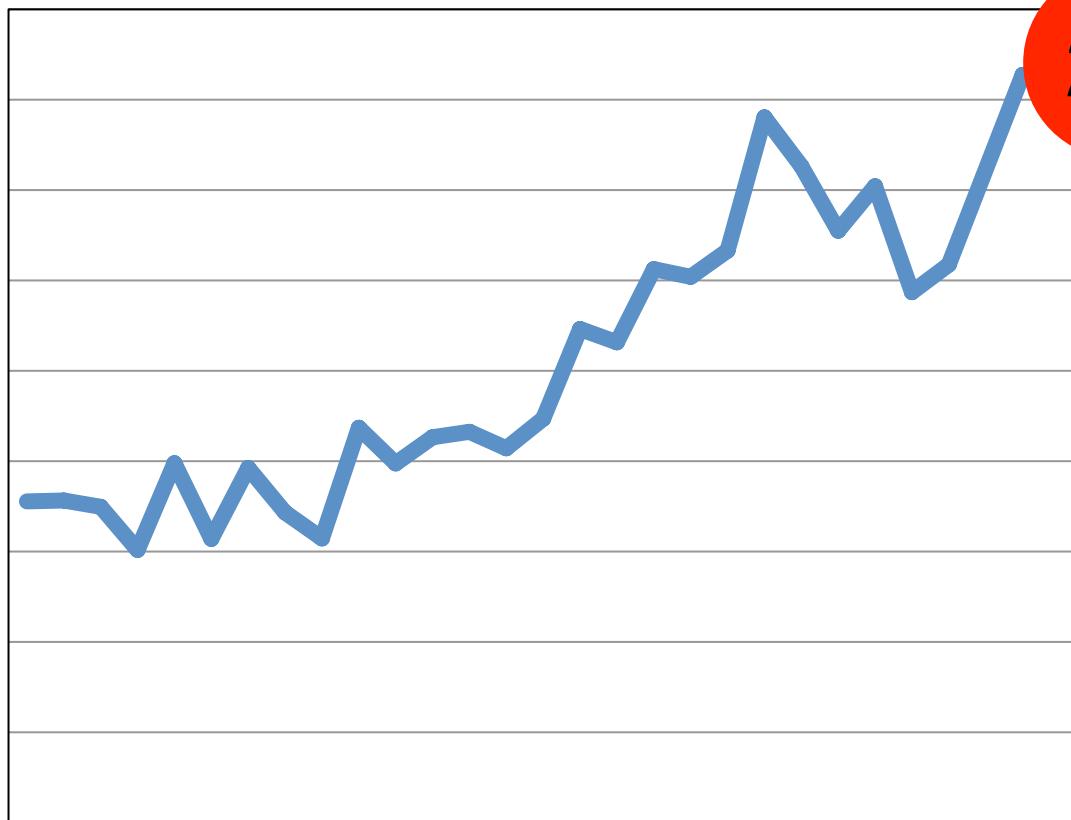
An advertising analytics company said it has discovered a botnet that generates upwards of US\$6 million per month by generating bogus clicks on display advertisements.

Spider.io, based in the U.K., wrote that the botnet code, called Chameleon, has infected about 120,000 residential computers in the U.S. and perpetrates click fraud on 202 websites that collectively deliver 14 billion ad impressions. Chameleon is responsible for 9 billion of those impressions, Spider.io said.

Click fraud cheats Web advertisers by making them pay for clicks on ads that are not legitimate, depriving them of customers and revenue. Spider.io said advertisers pay an average of \$0.69 per one thousand impressions.

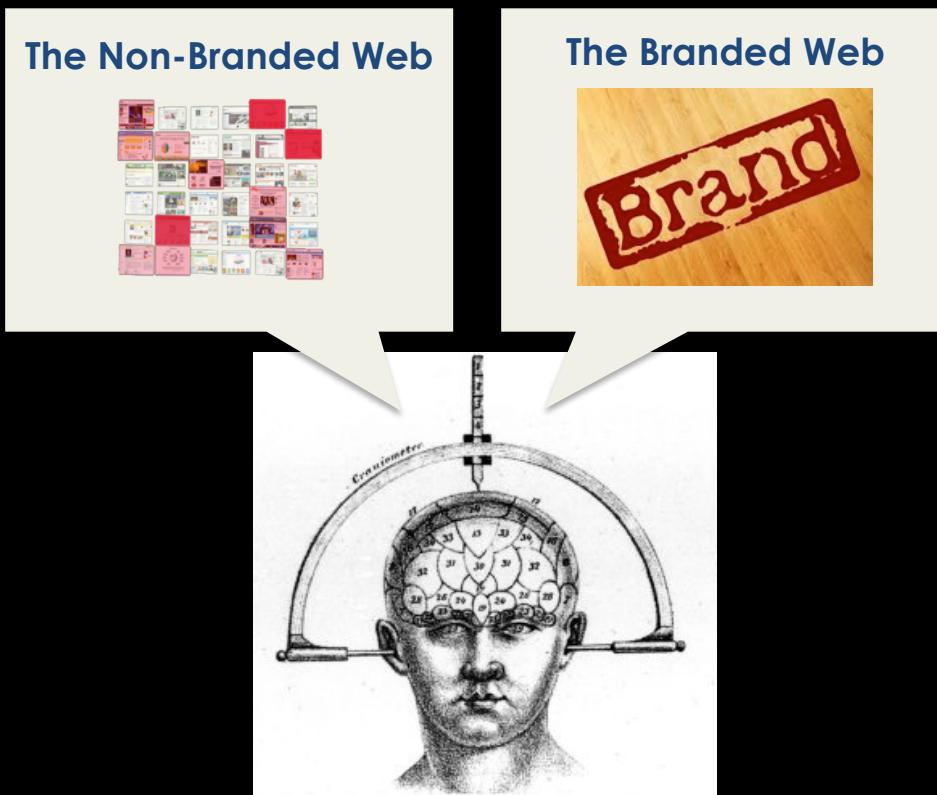
Unreasonable Performance Increase Spring 12

Performance Index



Dstillery Targeting: Agnostic Data

A consumer's online activity → gets recorded like this:



Purchases
Encoded
date1 3012L20
date 2 4199L30
...
date n 3075L50

Browsing History

Hashed URL's:

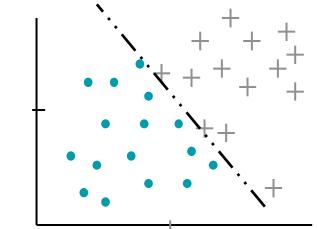
date1 abkcc
date2 kkllo
date3 88iok
date4 7uiol
...

computers do not 'understand' who you are



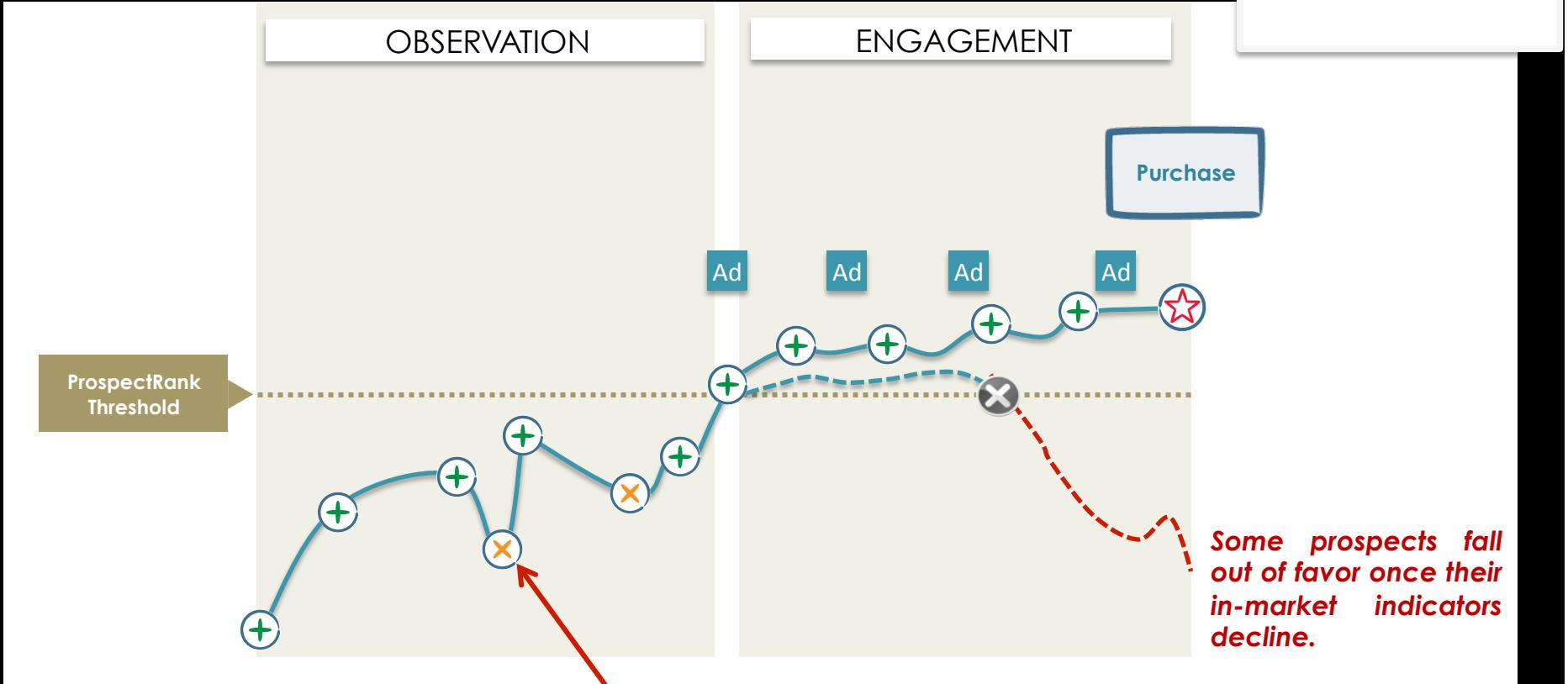
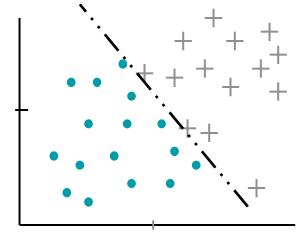
Model of browsing histories

Using Logistic Regression, we estimate statistical correlations between 10s of millions of web URLs and 1000s of actions.



$$p(\text{buy}|\text{urls}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Real-time Scoring of a User



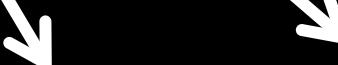
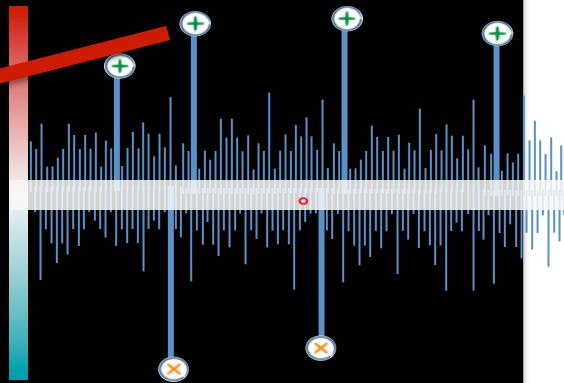
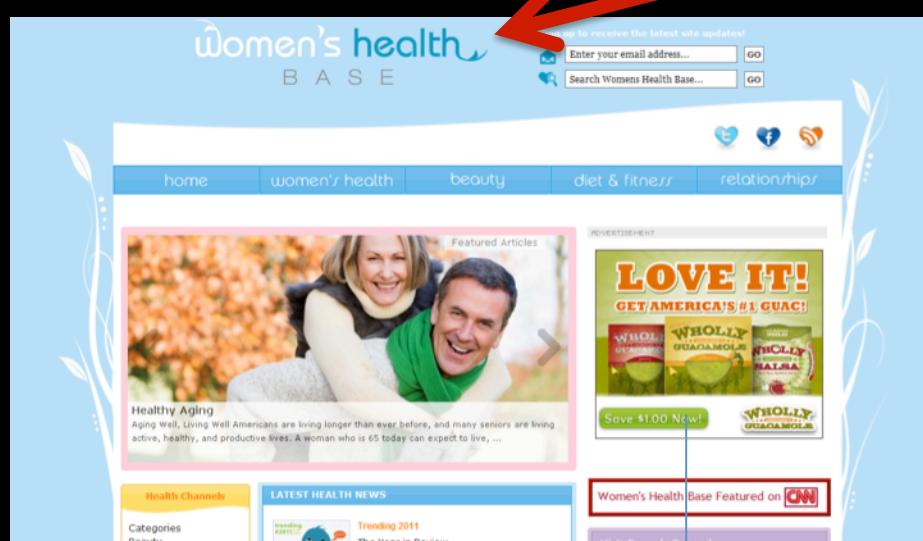
site visit with positive correlation



site visit with negative correlation

$$p(\text{buy} | \text{urls}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Oddly predictive websites?



Appearances are deceiving ...

women's health
B A S E

Sign up to receive the latest site updates!

Enter your email address...

Search Womens Health Base...

[home](#) [women's health](#) [beauty](#) [diet & fitness](#) [relationships](#)

[Featured Articles](#)

 [Healthy Aging](#)
Aging Well, Living Well Americans are living longer than ever before, and many seniors are living active, healthy, and productive lives. A woman who is 65 today can expect to live, ...

ADVERTISEMENT

LOVE IT!
GET AMERICA'S #1 GUAC!

Save \$1.00 Now!

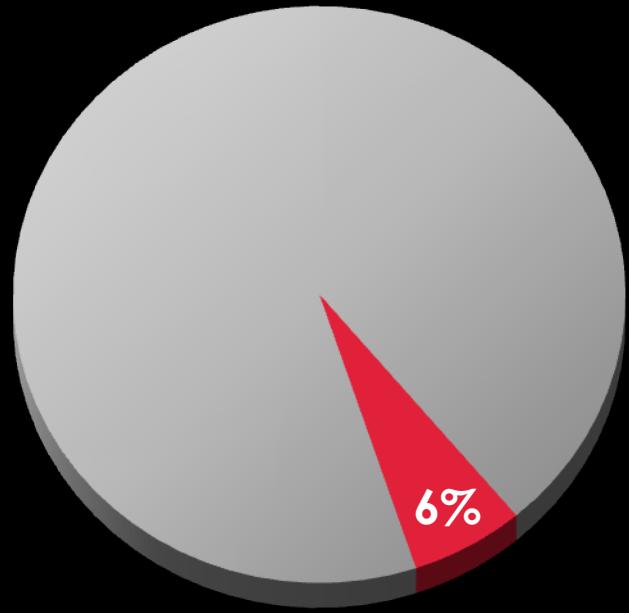
[Health Channels](#)

[Categories](#) [People](#)

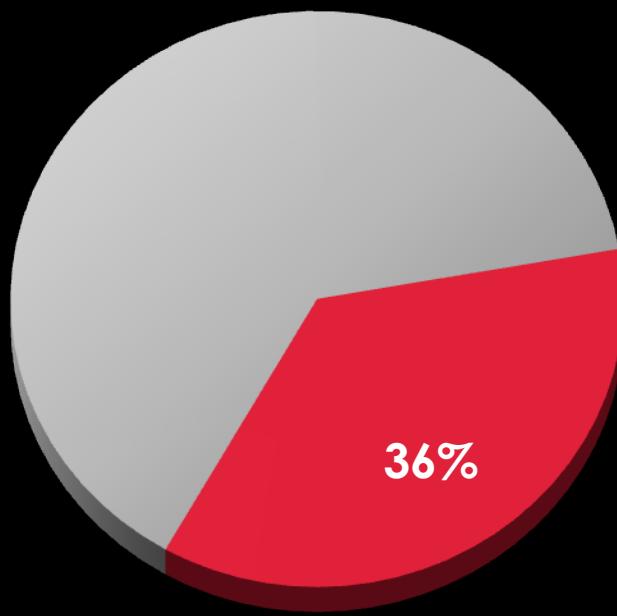
LATEST HEALTH NEWS

 Women's Health Base Featured on CNN

36% traffic is Non-Intentional



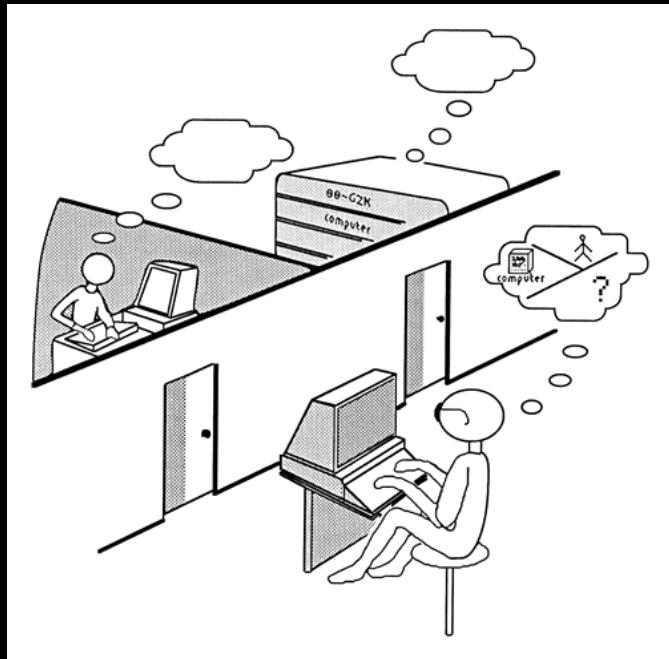
2011



2012

Telling the difference between an algorithm and a human

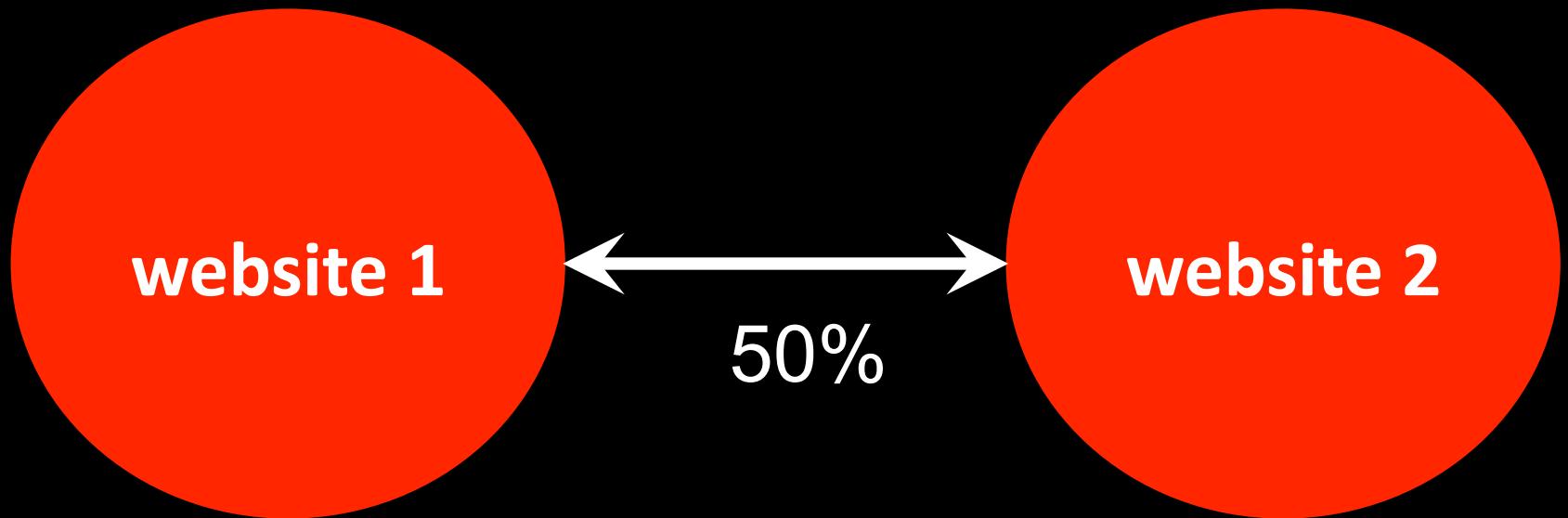
Turing test



KAPTCHA



Traffic patterns are ‘non - human’

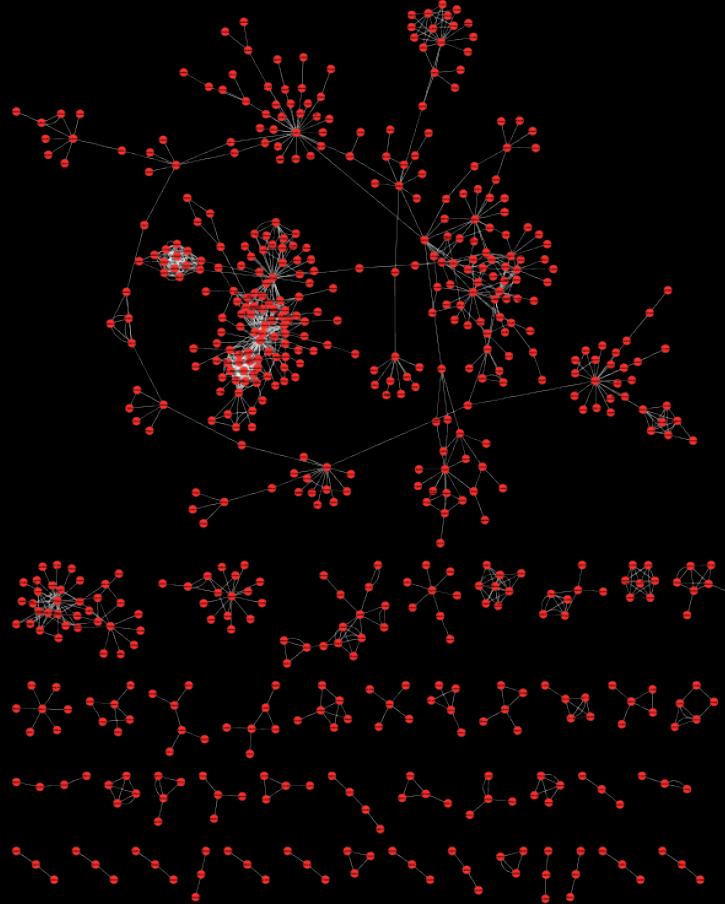


Data from Bid Requests in Ad-Exchanges

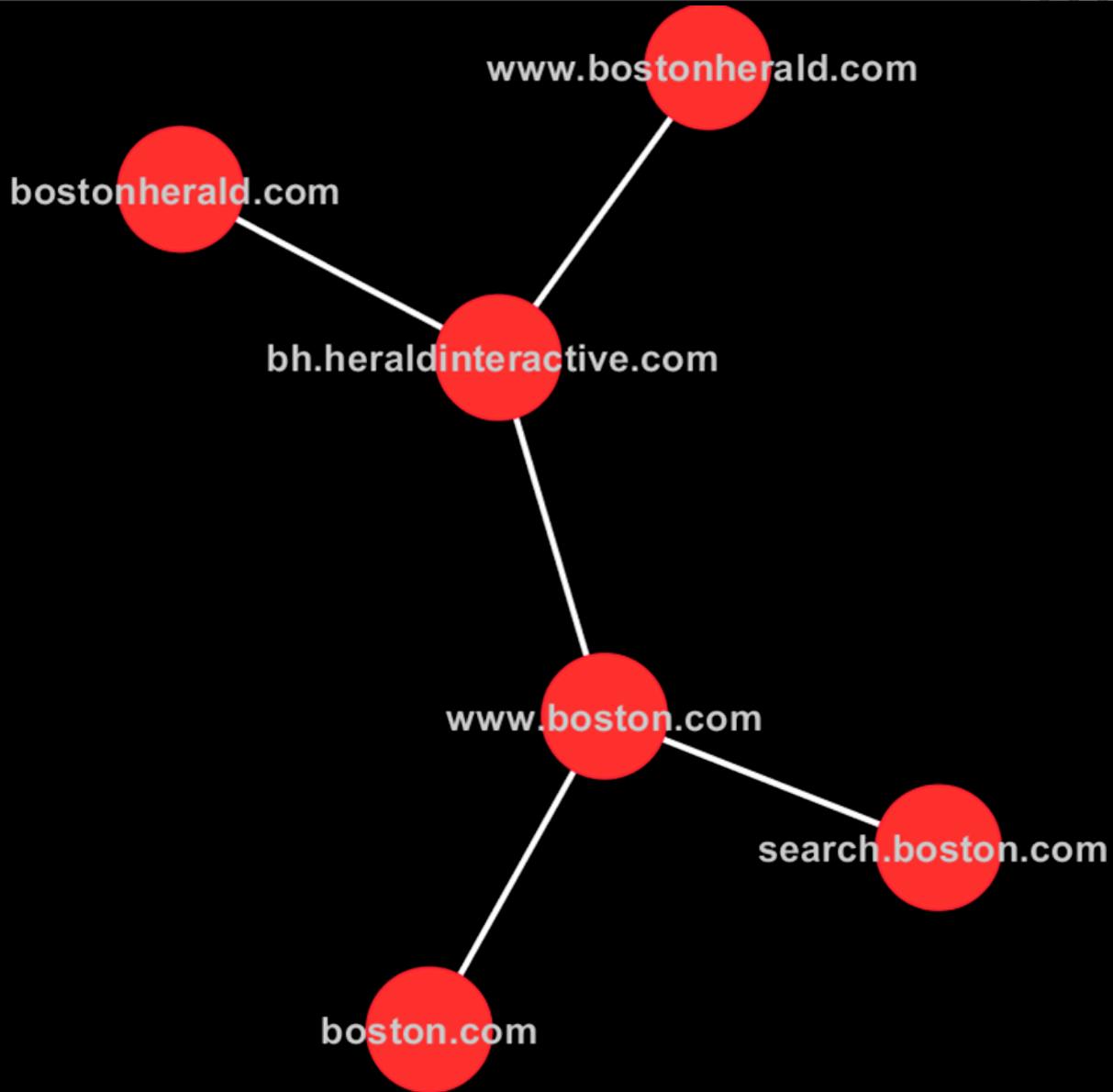
WWW 2010

Node:
hostname

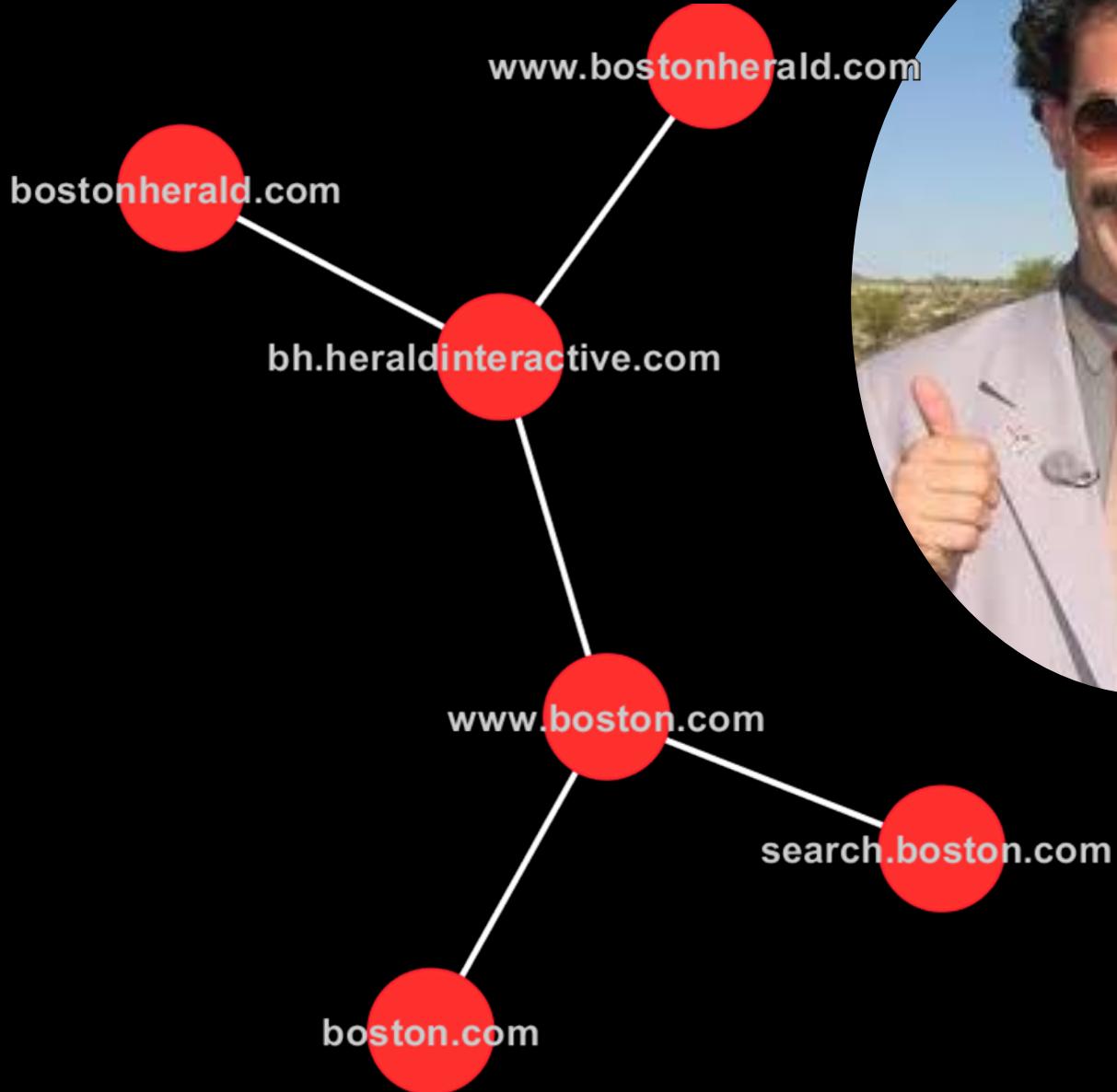
Edge:
50% co-visitation



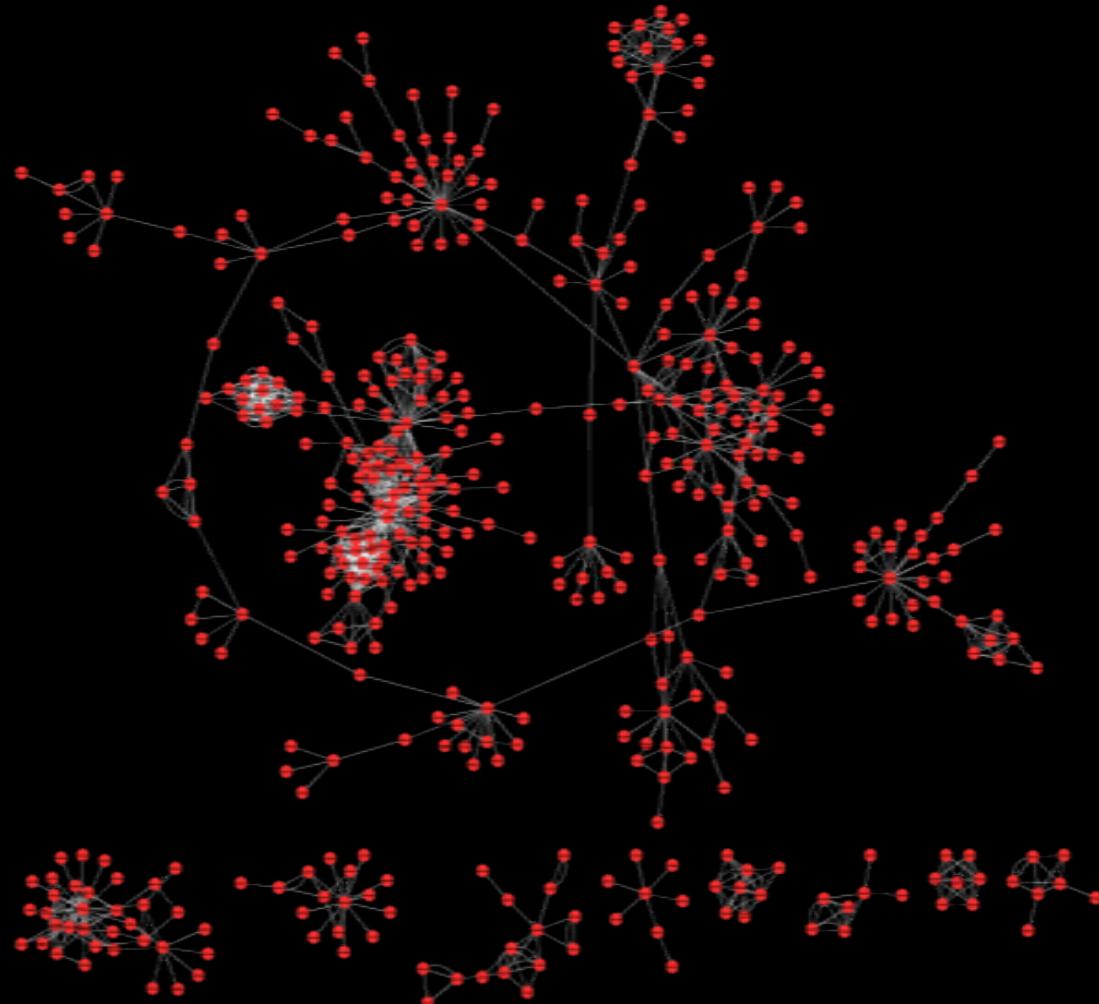
Boston Herald

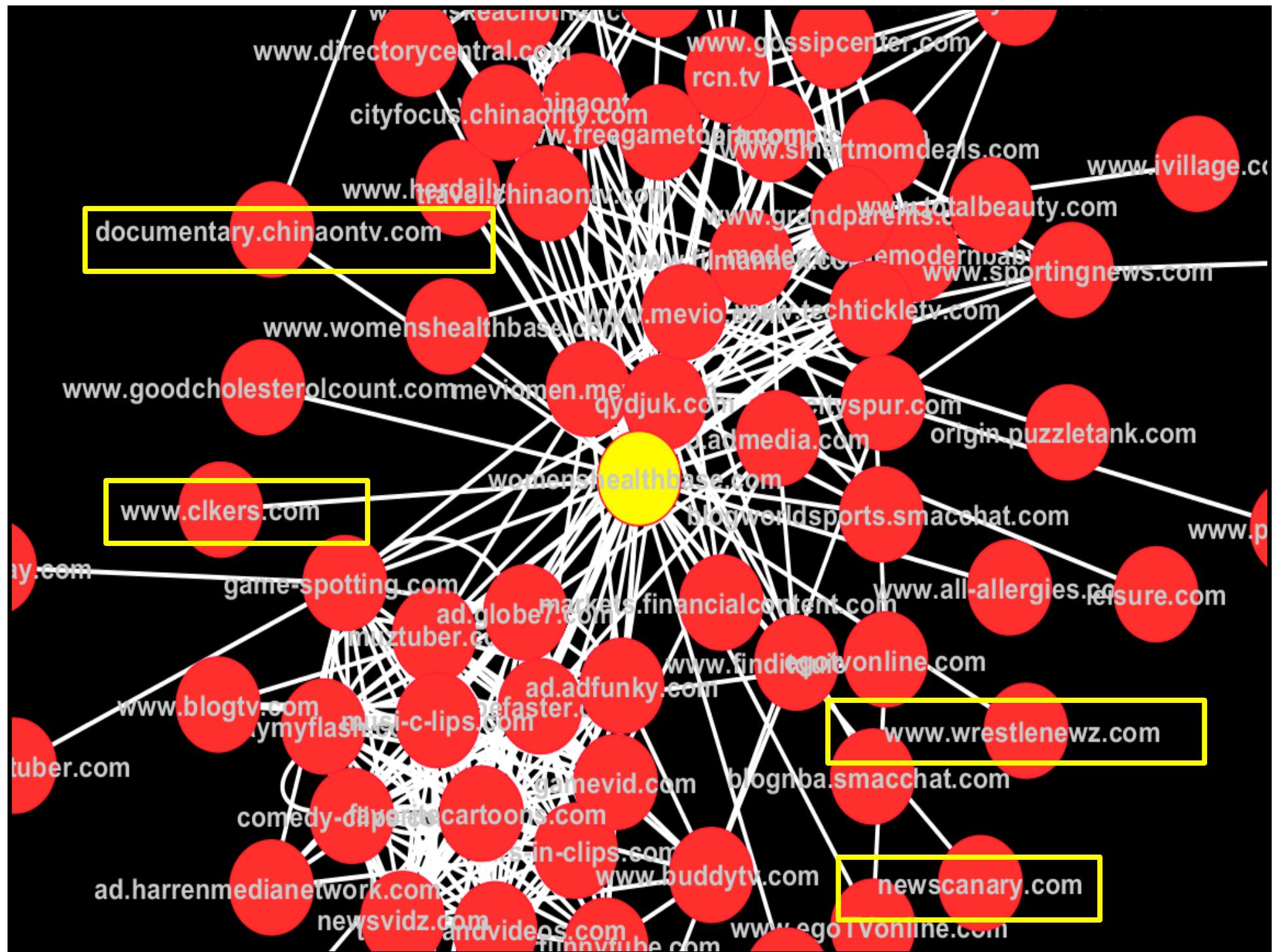


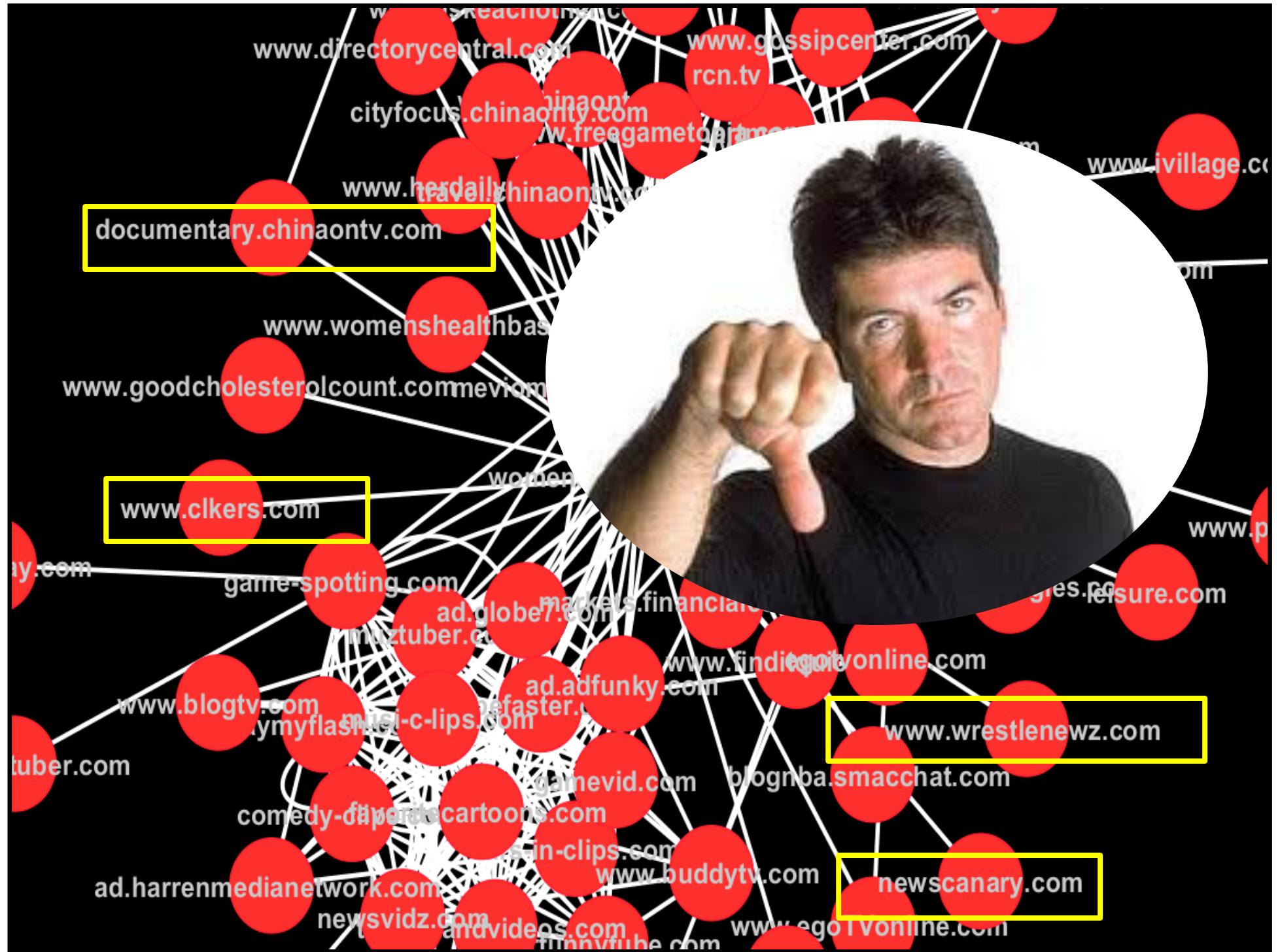
Boston Herald



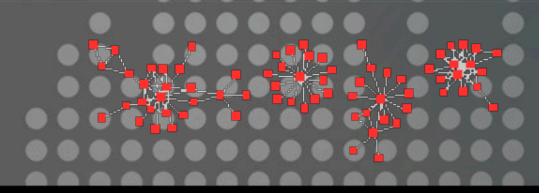
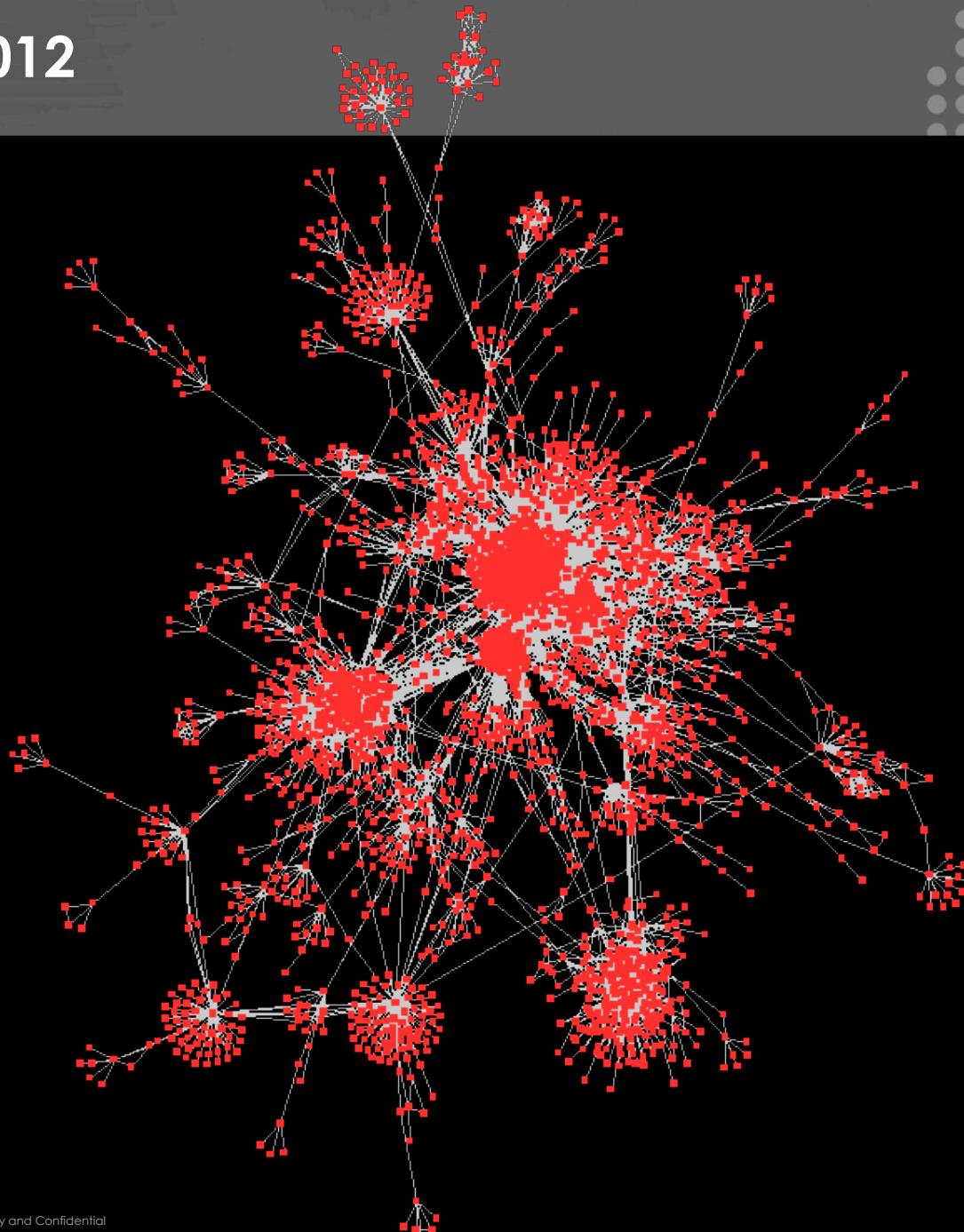
womenshealthbase?





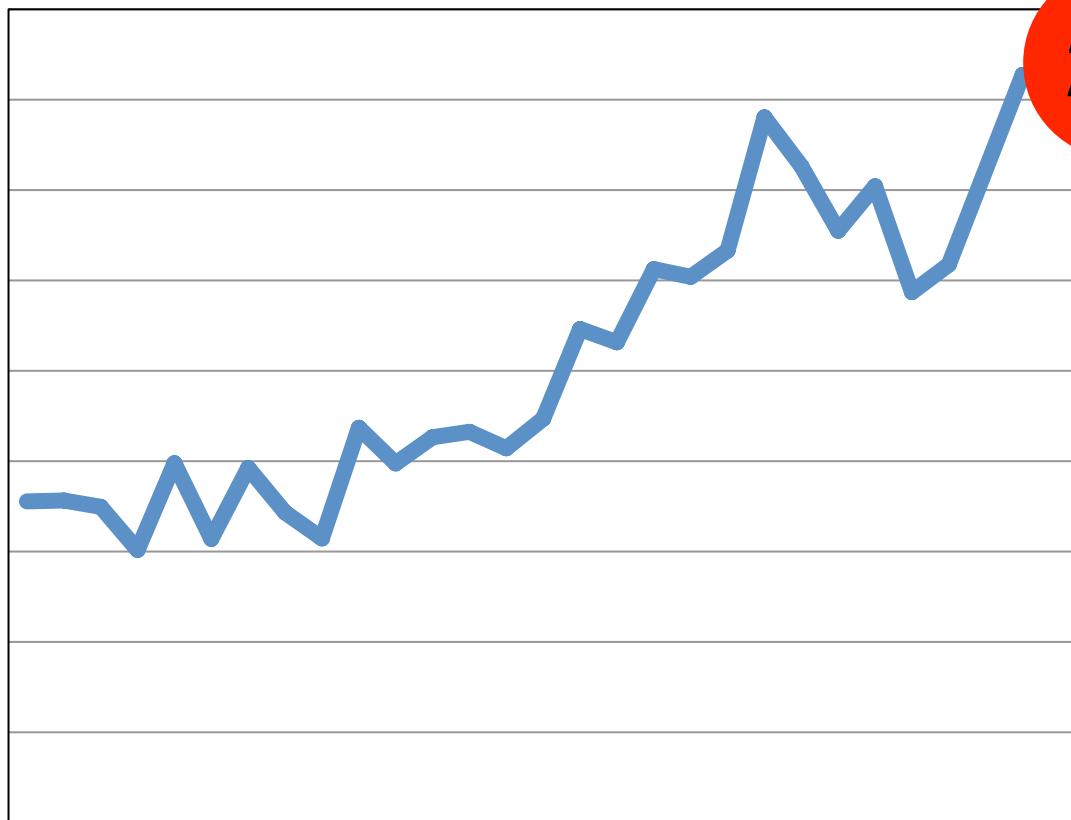


WWW 2012



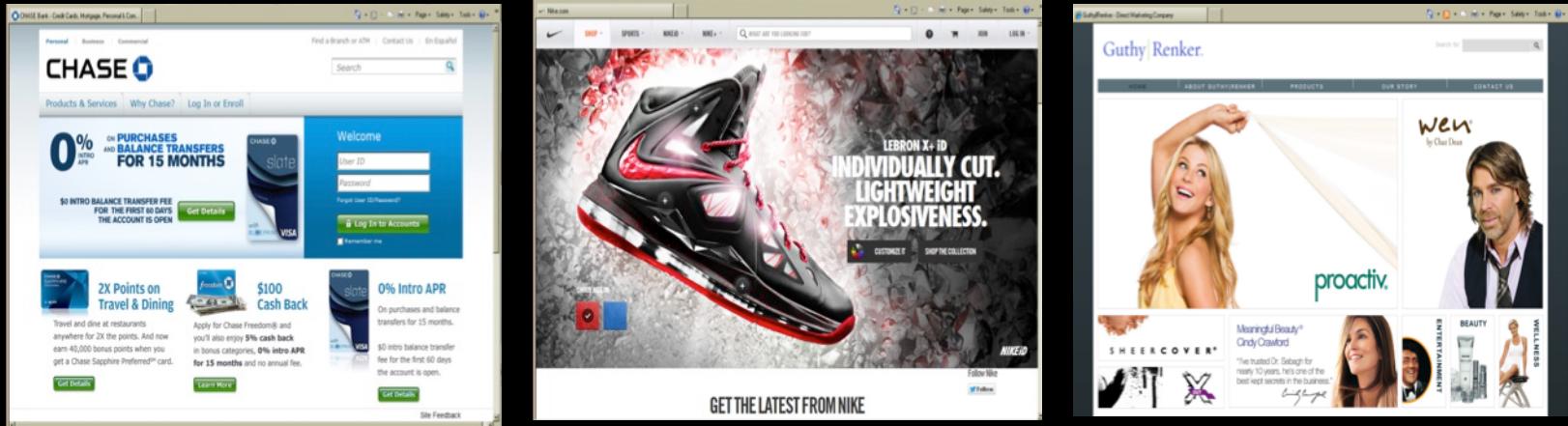
Unreasonable Performance Increase Spring 12

Performance Index



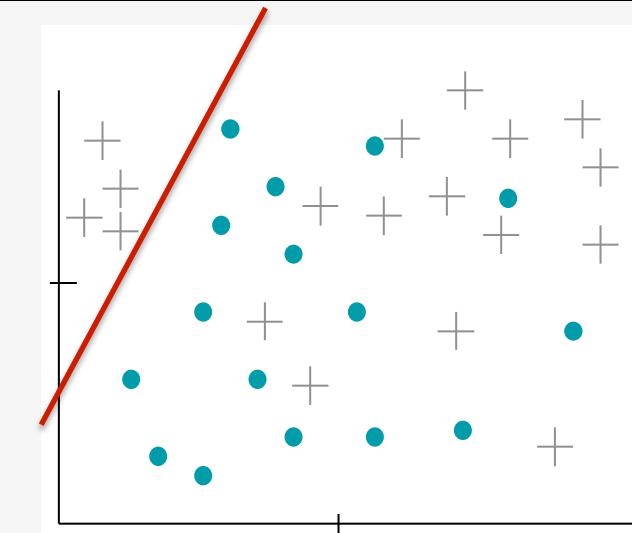
2 weeks

Now it is coming also to brands



- ‘Cookie Stuffing’ increases the value of the ad for all the retargeters
- Messing up your Web analytics ...
- Messes up my models because a botnet is easier to predict than a human

Fraud pollutes my models



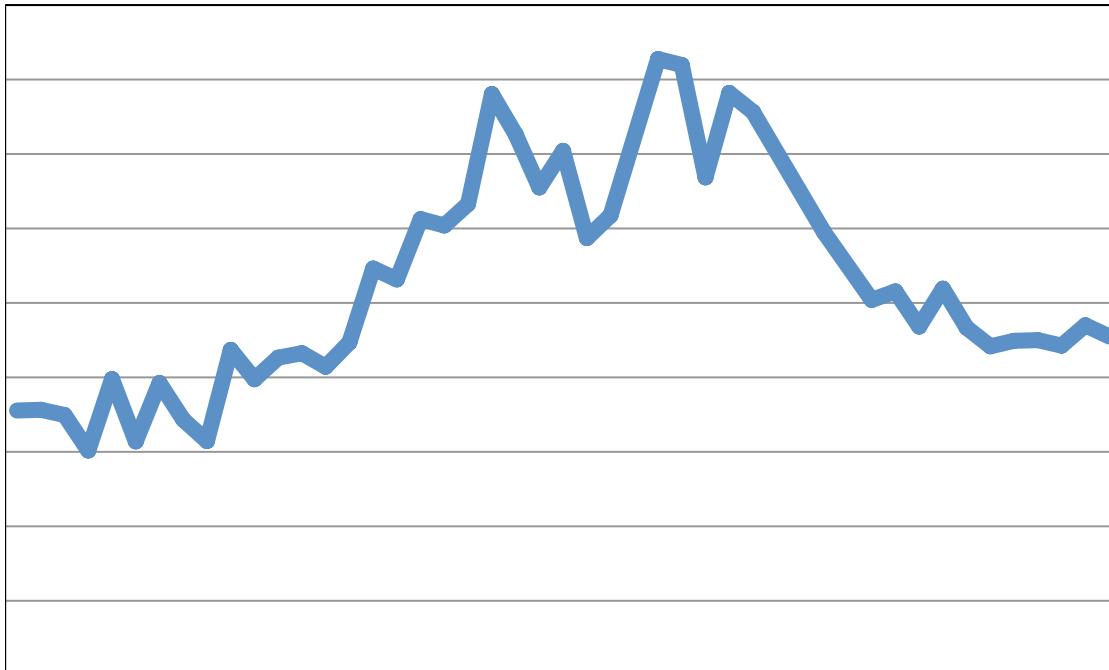
- Don't show ads on those sites
- Don't show ads to a high jacked browser
- Need to remove the visits to the fraud sites
- Need to remove the fraudulent site visits

When we see a browser on caught up in fraudulent activity:
send him to the penalty box where we ignore all his actions



Using the penalty box: all back to normal

Performance Index



3 more weeks in spring 2012

Fun Facts: Trying to look human?



Thanks! To the Data science Team

More technical details:

<http://m6d.com/who-we-are/datascience/>

Machine Learning Applications: Recommendation Engines Using Multiple Behavior Sources

9:30 – 9:50



Ted Dunning
Chief Application Architect, MapR



Multi-input Recommendations

Topic For Today

- What is recommendation?
- What makes it different?
- What is multi-model recommendation?
- How can I build it using common household items?

Oh ... Also This

- Detailed break-down of a live machine learning system running with Mahout on MapR
- With code examples

I may have to
summarize

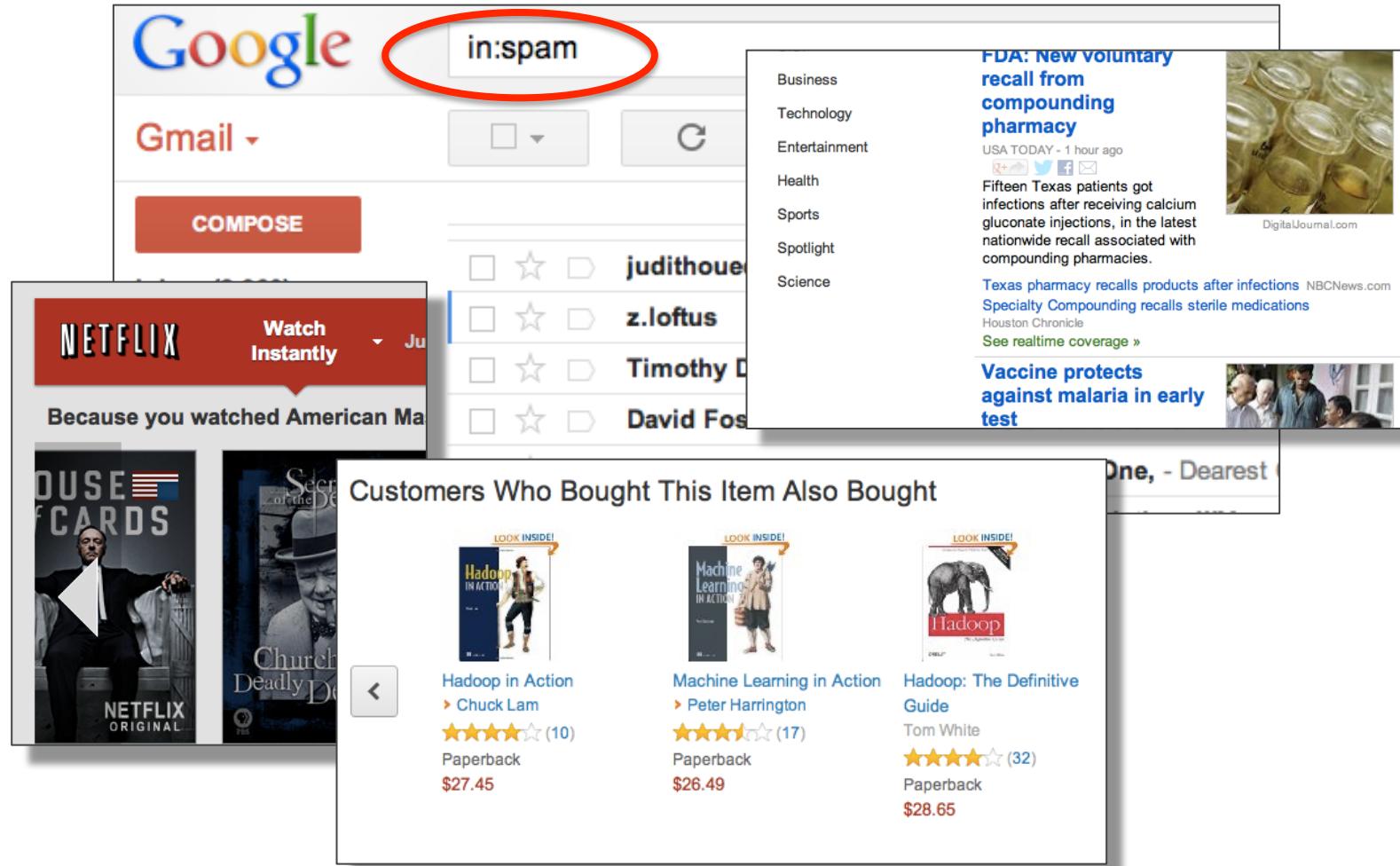
I may have to
summarize
just a bit

Part 1: 5 minutes of background

Part 2: 5 minutes: I want a pony

Part 1: 5 minutes of background

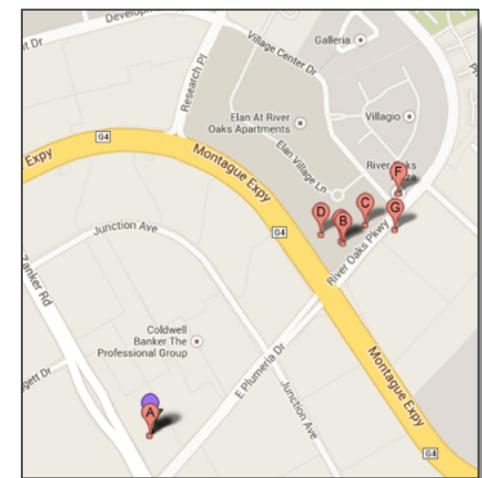
What Does Machine Learning Look Like?



Recommendations as Machine Learning

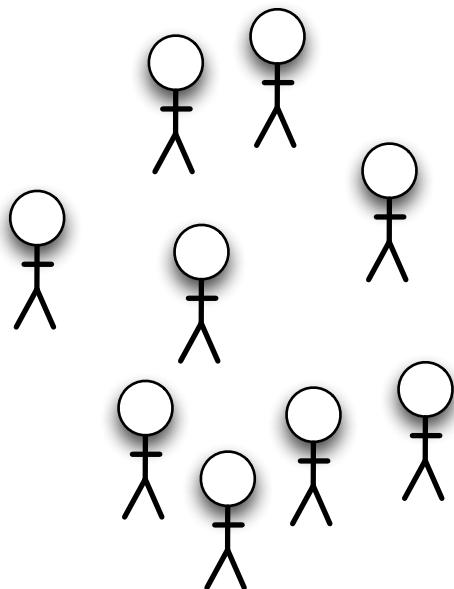
- **Recommendation:**

- Involves observation of interactions between people taking action (users) and items for input data to the recommender model
- Goal is to suggest additional appropriate or desirable interactions
- **Applications include:** movie, music or map-based restaurant choices; suggesting sale items for e-stores or via cash-register receipts



Part 2: How recommenders work (I still want a pony)

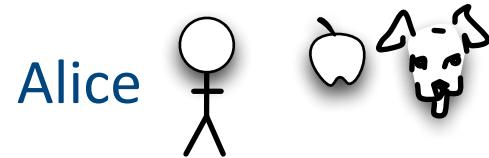
Recommendations



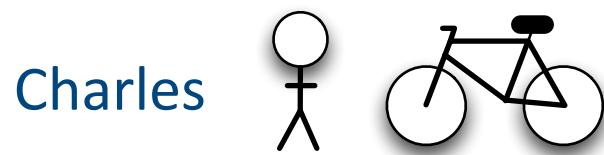
Recap:

Behavior of a crowd helps us understand what individuals will do

Recommendations

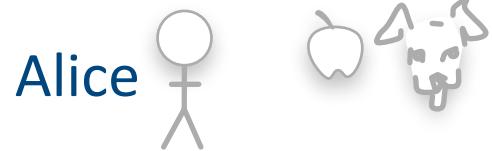


Alice got an apple and a puppy

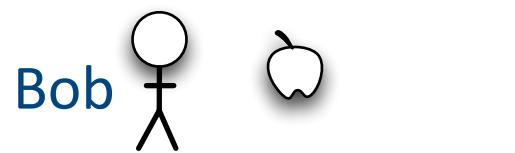


Charles got a bicycle

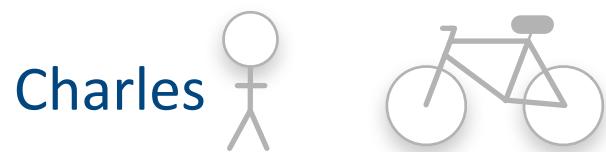
Recommendations



Alice got an apple and a puppy

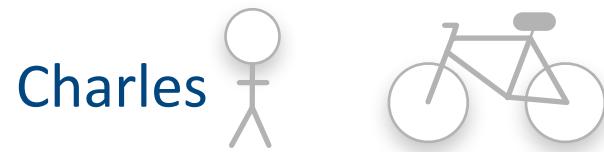
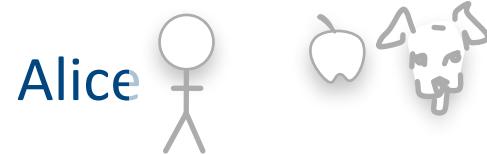


Bob got an apple

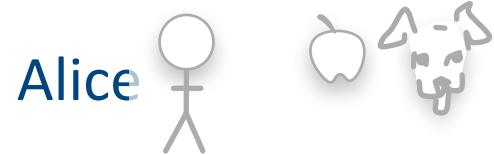


Charles got a bicycle

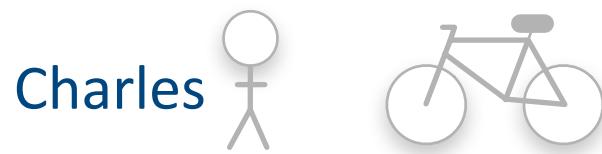
Recommendations



Recommendations

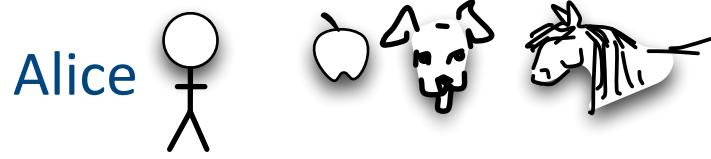


A puppy, of course!

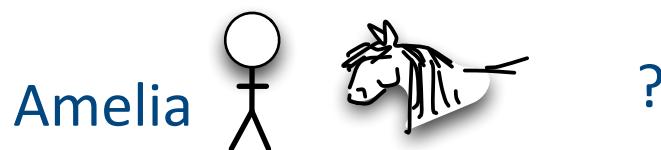
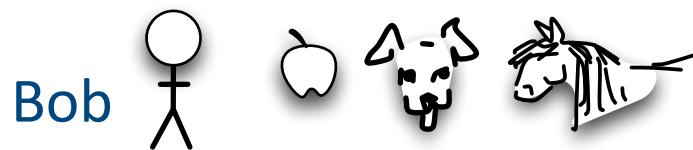


You get the idea of how
recommenders work...
(By the way, like me, Bob
also wants a pony)

Recommendations



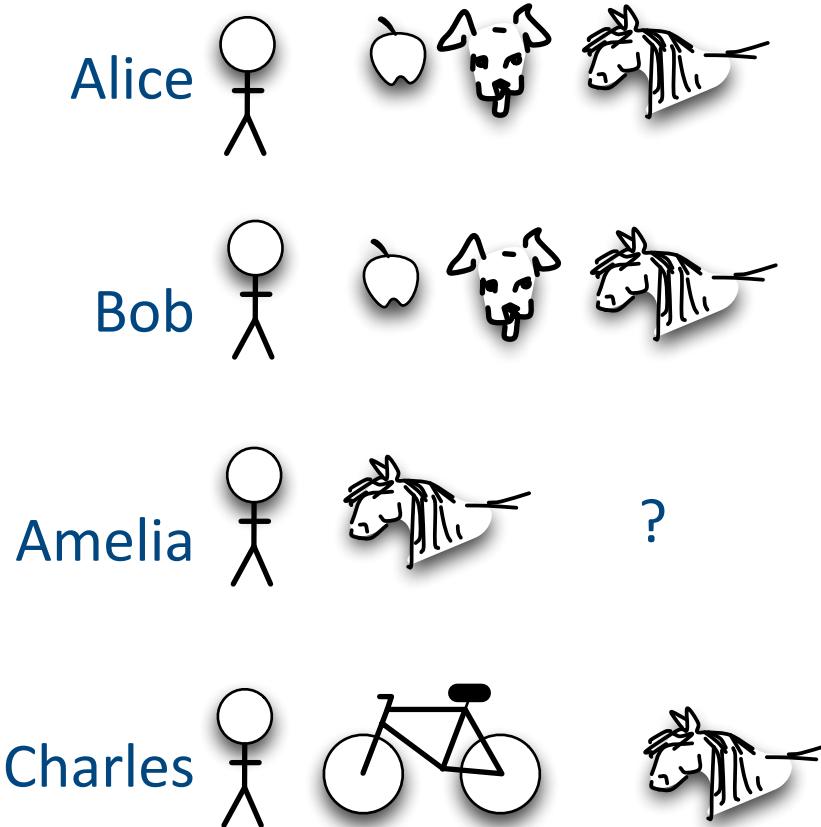
What if everybody gets a pony?



What else would you recommend for Amelia?



Recommendations



If everybody gets a pony, it's not a very good indicator of what to else predict...

Problems with Raw Co-occurrence

- **Very popular items co-occur with everything (or why it's not very helpful to know that everybody wants a pony...)**
 - Examples: Welcome document; Elevator music
- **Very widespread occurrence is not interesting as a way to generate indicators**
 - Unless you want to offer an item that is constantly desired, such as razor blades (or ponies)
- **What we want is *anomalous* co-occurrence**
 - This is the source of interesting indicators of preference on which to base recommendation

Get Useful Indicators from Behaviors

1. Use log files to build **history matrix** of users x items
 - Remember: this history of interactions will be sparse compared to all potential combinations
2. Transform to a **co-occurrence matrix** of items x items
3. Look for *useful* co-occurrence by looking for *anomalous* co-occurrences to make an **indicator matrix**
 - **Log Likelihood Ratio (LLR)** can be helpful to judge which co-occurrences can with confidence be used as indicators of preference
 - RowSimilarityJob in Apache Mahout uses LLR

Log Files

Alice	
Charles	
Charles	
Alice	
Alice	
Bob	
Bob	

History Matrix: Users by Items

				
Alice	✓	✓	✓	
Bob	✓		✓	
Charles			✓	✓

Co-occurrence Matrix: Items by Items

How do you tell which co-occurrences are useful?.

The diagram illustrates a 4x4 co-occurrence matrix for four items: Apple, Dog, Horse, and Bike. The items are arranged in a grid, with the Apple at the top-left, Dog above it, Horse to its right, and Bike at the bottom-right. Each item is accompanied by a small icon. The matrix itself is a 4x4 grid of cells, each containing a numerical value representing the co-occurrence count between the corresponding row and column items. The values are: Apple (1, 2, 0), Dog (1, 1, 0), Horse (2, 1, 1), and Bike (0, 0, 1).

	Apple	Dog	Horse	Bike
Apple	1	2	0	
Dog	1	1	0	
Horse	2	1	1	
Bike	0	0	1	

Co-occurrence Binary Matrix

	apple	not	apple
not	1		
dog	1	1	1

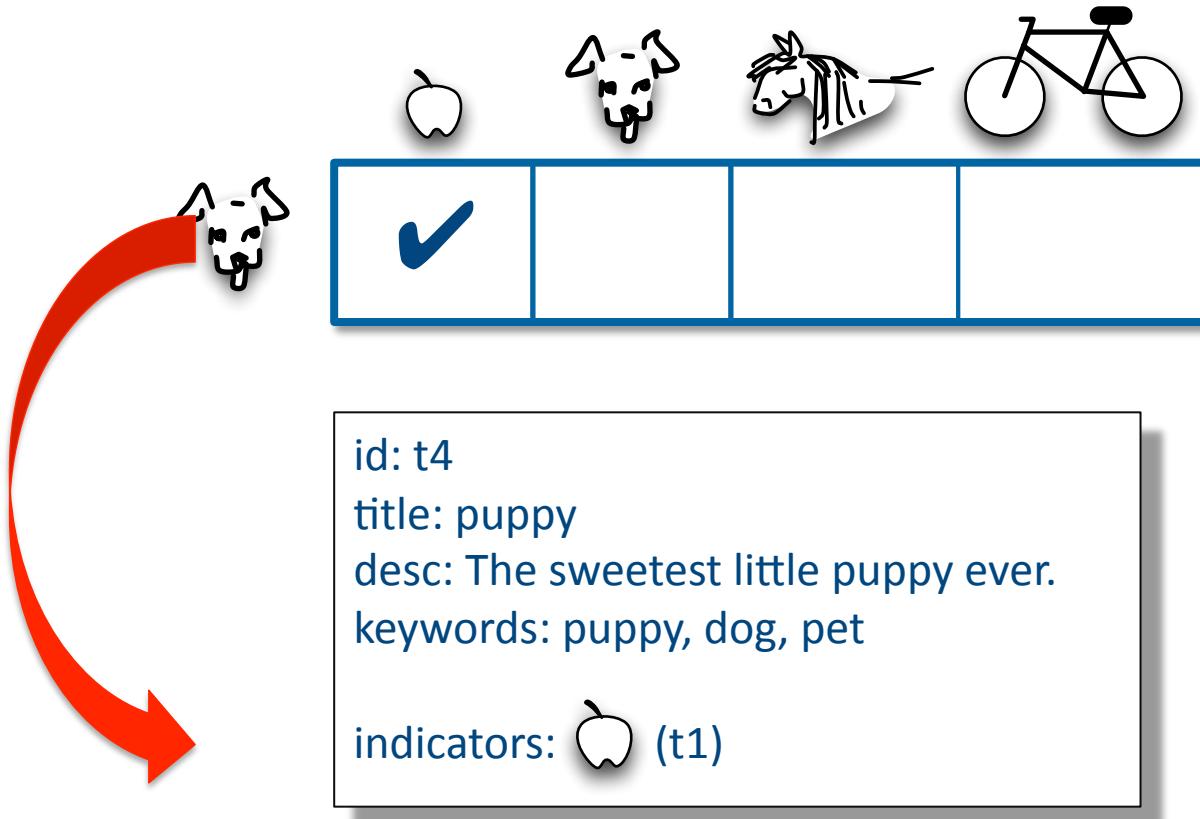
Indicator Matrix: Anomalous Co-Occurrence

Result: The marked row will be added to the indicator field in the item document...

	Apple	Dog	Horse	Bike
Apple				
Dog		✓		
Horse				
Bike				

Indicator Matrix

That one row from indicator matrix becomes the indicator field in the Solr document used to deploy the recommendation engine.



Note: data for the indicator field is added directly to meta-data for a document in Solr index. You don't need to create a separate index for the indicators.

Internals of the Recommender Engine

The screenshot shows the LucidWorks dashboard for artists. At the top, there is a navigation bar with links for admin, Sign out, Users, and Settings. Below the navigation bar, the title "Dashboard : Artists" is displayed, along with tabs for Status, Indexing, Querying, Access Control, Search (which is selected), and Advanced.

In the search bar, the query "indicator_artists:2122 OR indicator_artists:303" is entered, and the "Search" button is clicked. Below the search bar, there are links for Create Alert and View Alerts.

The search results show 152 items, with results 1 - 10 listed:

Artist ID	MBID	Name	Area	Gender	Indicator IDs
1710	592a3b6d-c42b-4567-99c9-ecf63bd66499	Chuck Berry	United States	Male	386685,875994,637954,3418,1344,315154,694400,789739,1460,630993,716411,735157,694612,12108,12100,910856,774838,24107,4111,858195,29,745742,3331,162
541902	983d4f8f-473e-4091-8394-415c105c4656	Charlie Winston	United Kingdom	None	,997727,109601,815,830794,59588,900,2591,311520,696268,6238,56964,809858,42884,821586,22532,805582,112605,2122,185581,242,303,719703,2129,125800,825

For each result, there is a link to "Find similar - Explain". On the right side of the search results, there are sections titled "Narrow your search" with dropdown menus for "Area" and "Gender".

Internals of the Recommender Engine

The screenshot shows the LucidWorks dashboard for 'Artists'. The top navigation bar includes links for 'Status', 'Indexing', 'Querying', 'Access Control', 'Search' (which is highlighted in orange), and 'Advanced'. The main content area displays search results for the query "indicator_artists:2122 OR indicator_artists:303". The results list includes two entries:

ID	mbid	name	area	gender	indicator_artists
1710	592a3b6d-c42b-4567-99c9-ecf63bd66499	Chuck Berry	United States	Male	386685,875994,637954,3418
541902	983d4f8f-473e-4091-8394-4	Charlie Winston	United Kingdom	None	,997727,109601,815,830794,59588,900,2591,311520,696268,6238,56964,809858,42884,821586,22532,805582,112605,2122,185581,242,303,719703,2129,125800,825

A red circle highlights the 'indicator_artists' field value for the first entry (1710). Below the results, there are 'Find similar - Explain' links for both entries.

A Quick Simplification

- Users who do h (a vector of things a user has done)

$$\mathbf{A}\mathbf{h}$$

\mathbf{A} translates things into users

- Also do r

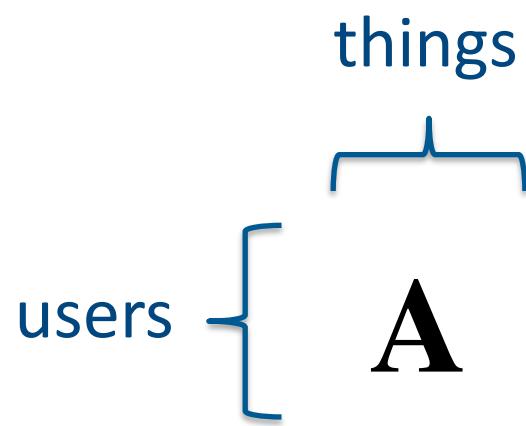
$$\mathbf{A}^T (\mathbf{A}\mathbf{h})$$

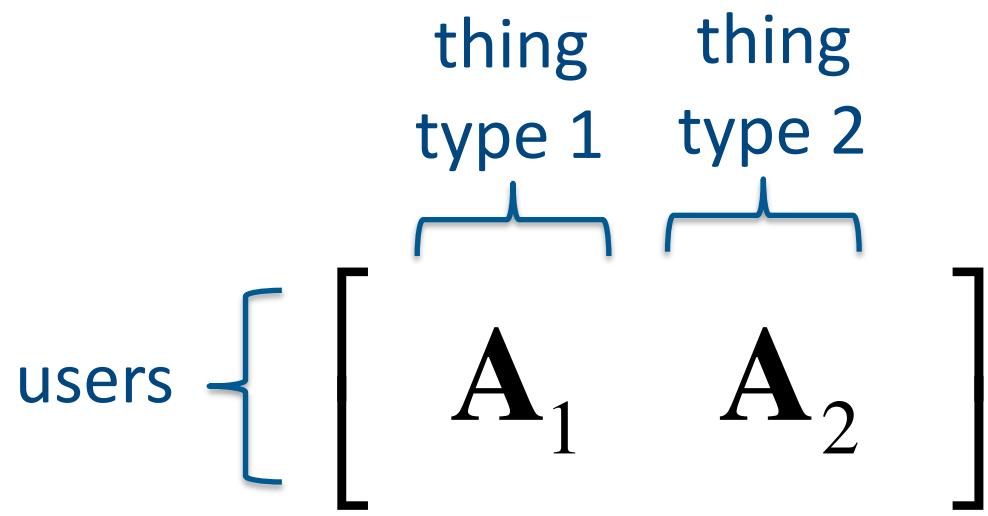
User-centric recommendations
(transpose translates back to things)

$$(\mathbf{A}^T \mathbf{A})\mathbf{h}$$

Item-centric recommendations
(change the order of operations)

Nice. But we
can do better?





Now again, without
the scary math

For example

- Users enter queries (A)
 - (actor = user, item=query)
- Users view videos (B)
 - (actor = user, item=video)
- $A^T A$ gives query recommendation
 - “did you mean to ask for”
- $B^T B$ gives video recommendation
 - “you might like these videos”

The punch-line

- B^TA recommends videos in response to a query
 - (isn't that a search engine?)
 - (not quite, it doesn't look at content or meta-data)

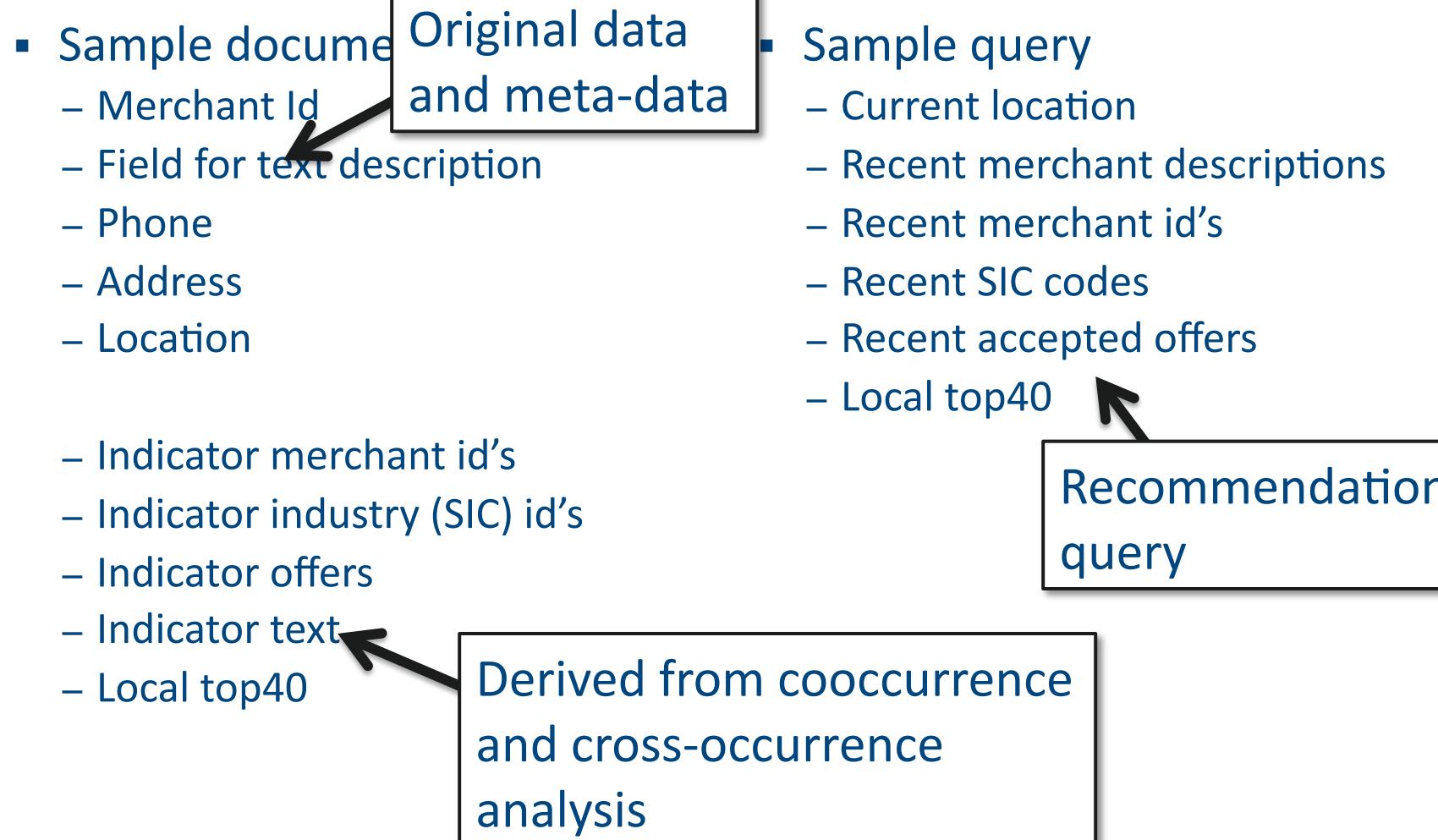
Real-life example

- Query: “Paco de Lucia”
- Conventional meta-data search results:
 - “hombres del paco” times 400
 - not much else
- Recommendation based search:
 - Flamenco guitar and dancers
 - Spanish and classical guitar
 - Van Halen doing a classical/flamenco riff

Real-life example

	<p><u>CONCIERTO CIUDAD DE LAS IDEAS PARTE FINAL</u> Music 58 views</p>
	<p><u>Siudy / Buleria</u> Music 722 views</p>
	<p><u>Vicente Amigo 2ª parte Ciudad de las Ideas</u> Music 124 views</p>
	<p><u>Van Halen's Eruption</u> Music 4400 views</p>
	<p><u>Freestyle Flamenco</u> Music 653 views</p>

Search-based Recommendations



Me, Us

- Ted Dunning, Chief Application Architect, MapR
Committer PMC member, Mahout, Zookeeper, Drill
Bought the beer at the first HUG
- MapR
Distributes more open source components for Hadoop
Adds major technology for performance, HA, industry standard API's
- Info
Hash tag - #mapr
See also - @ApacheMahout @ApacheDrill
@ted_dunning and @mapR