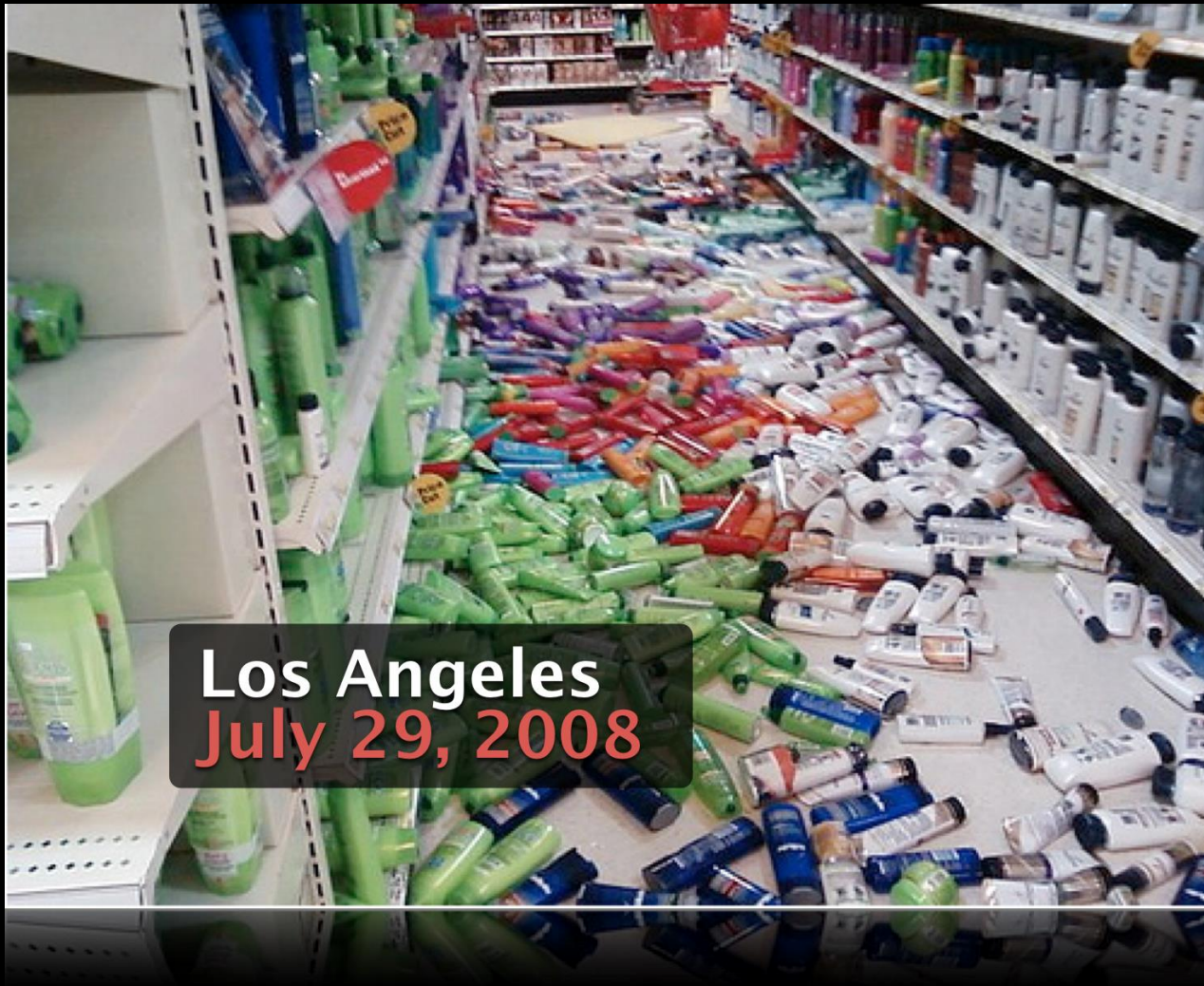


Research at Twitter

Aleksander Kołcz

Twitter, Inc.





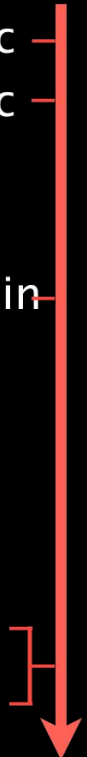
Los Angeles
July 29, 2008

Earthquake: 0 sec

First Tweet: 5 sec

Local News: 4 min

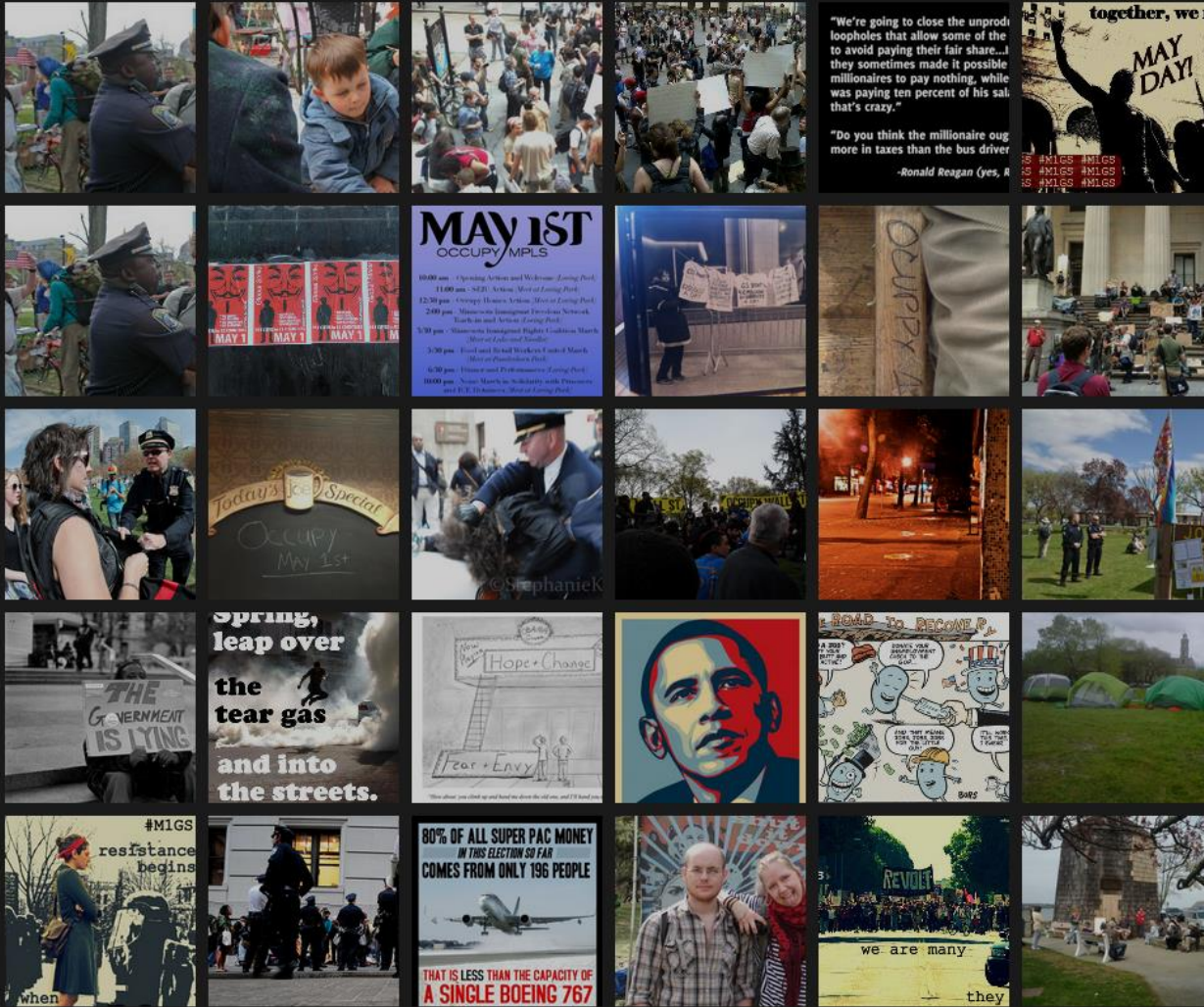
AP Wire: 9 min
Book of
Tweets



#ocws, #occupywallstreet, ...

[← Back to search results](#)

Top images for #ows



The Scale of Twitter

- > 200M active users
- Approx. 400M Tweets sent/day
- > 400M/month unique visitors to twitter.com
- Support for > 35 languages
- 70% of users outside US
- 33,388 record TPS



Internal vs External R&D

- A lot of Twitter data is open (in small quantities)
- It is easy to create small-scale specialized collections and institutions can acquire larger collections via Gnip
- There are lots of academic publications analyzing interesting properties of Twitter data
- Internal R&D is focused on improving the user experience



Overview

- Key problem areas
 - Product features
 - Infrastructure
- Key research problems/subproblems
 - Evolving landscape
- Collective experiences
 - A few lessons learned



Problems we are trying to solve

- Relevance
 - ranking in search

Results for **president obama**

Tweets Top / All / Timeline

 **Obama 2012** @Obama2012 11h
"In this country, prosperity does not trickle down. Prosperity grows from the bottom up." —**President Obama** speaking in Elyria, Ohio today

 **Barack Obama News** @ObamaNews 7m
Press Release: **President Obama** Signs Hawaii Disaster Declaration bit.ly/14dN0q

 **Barack Obama News** @ObamaNews 4h
Blog Post: **President Obama** Talks About Investing in Training American Workers bit.ly/13JvLj

 **Donna Brazile** @donnabrazile 4h
For the record, I support **President Obama's** re-election efforts. But, I am not a surrogate for the campaign or the spokesperson for the DNC.

 **Barack Obama** @BarackObama 5h
President Obama met with some Ohioans who are benefitting from community college job training programs today: OFA.BO/Wmcw83





Problems we are trying to solve

- Who to follow

Who to follow


Twitter accounts suggested for you based on who you follow and more.



USGS  @USGS

Earth science knowledge is just a tweet away. Tweets do not = endorsement: <http://on.doi.gov/pgwuoY>


Followed by USDA Food Safety , The Economist and Emergency_In_SF .



Hilary Mason @hmason


chief scientist @bitly. Machine learning; I ♥ data and cheeseburgers.


Followed by Gregory Piatetsky , Ian Soboroff and Eugene Agichtein .



Adam Rugel @Adam


Trazzler, Reston, Syracuse University, Sandwich



Google Research  @googleresearch

At Google, research is performed company wide, not just in isolated labs. We produce and leverage research to build systems that are used in the real world.

Followed by Tao Tao , Kurt Smith and SIGKDD/KDD News .



Problems we are trying to solve

Content recommendation
(stories/media)

Stories

Twitter Empowers Engineers With New Patent Agreement



Twitter, in what it says is an act of good will to its engineers and designers, announced a new patent agreement that gives control back to inventors in...

bits.blogs.nytimes.com/2012/04/17/twi...



Tweeted by **Matt Cutts**

Report: Sony's Image URL Accidentally Reveals God of War IV



God of War IV Å is coming. It's obvious at this point that Sony will be unveiling God of War IV Å soon. The next entry in this franchise has had a slew of rumors and...

technobuffalo.com/gaming/platfor...



Trending Tweets about **God of War**

Striking New Photos Of Great 1906 Earthquake Emerge



On the anniversary of the Great San Francisco Earthquake of 1906, the San Francisco Municipal Transportation Agency has released a stunning new set...

sfist.com/2012/04/18/new...

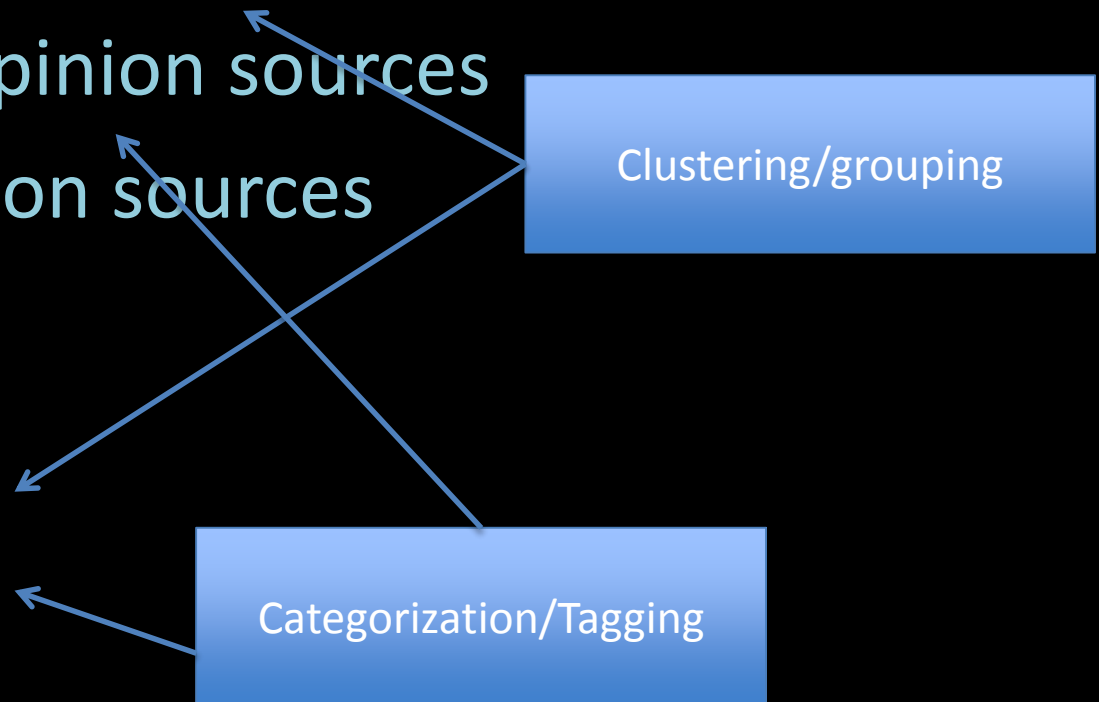


Tweeted by people who share your interests



Recommendation/Personalization

- Other users – friends
- Other users – opinion sources
- News/information sources
- Specific tweets
- Specific urls
- Lists/hashtags

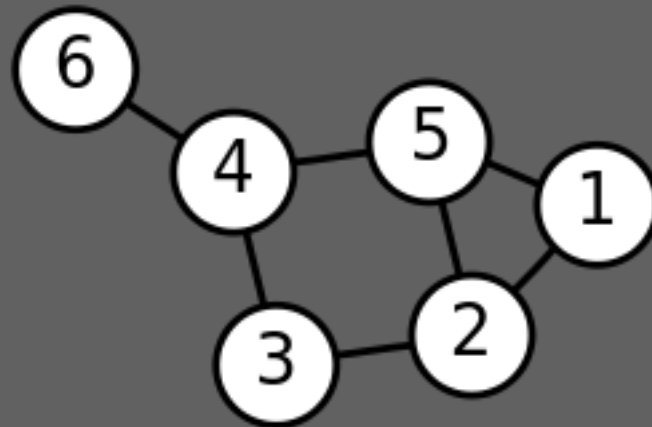


Key Research Areas

- Graph mining
- Content/user tagging/classification
- Content/user recommendation
- Search/IR
- Integration/Infrastructure (best practices)



Graph mining



Graph mining (Follow graph)

- How does it vary across locales
- How to use it to judge user importance
- How to use it to provide user recommendations
- The role of reciprocated edges (mutual follows)
- How to use it to analyze information spread (diffusion)



Follow graph

- 20B edges connecting active users
- Over 1K users having 1M+ followers
- > 25 users having 10M+ followers



wtf


(who not to follow)

Who to follow · [refresh](#) · [view all](#)



freshbooks FreshBooks · [Follow](#)



 Promoted · Followed by @zappos and others.



alanwarms Alan Warms · [Follow](#)



Followed by @fredwilson and others.



Mozzie21 Moises Henriques · [Follow](#)



can eat

Similar to @ryanhall3 · [view all](#)



RunnerSpace_com RunnerSpace.com · [Follow](#)

RunnerSpace.com has the latest in news and media...



chrislieto chris lieto · [Follow](#)

Chris Lieto is a top ranked World Class Triathlete, ...



runningtimes runningtimes · [Follow](#)

Source: Gupta et al. "WTF: The Who to Follow Service at Twitter", WWW'13

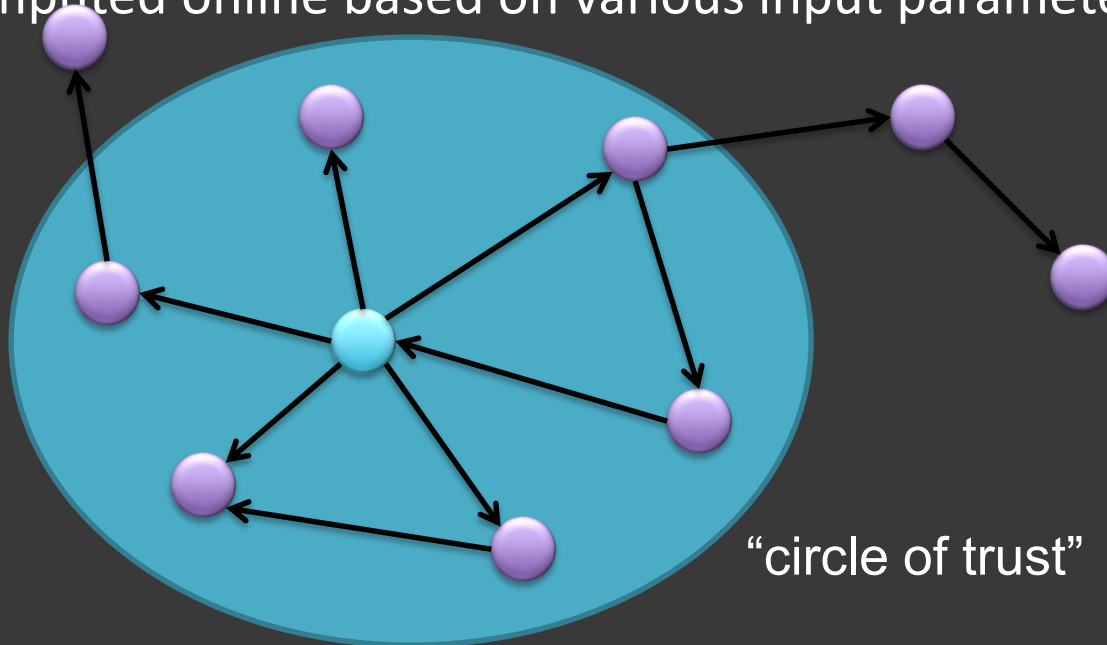


“Circle of Trust”

Ordered set of important neighbors for a user

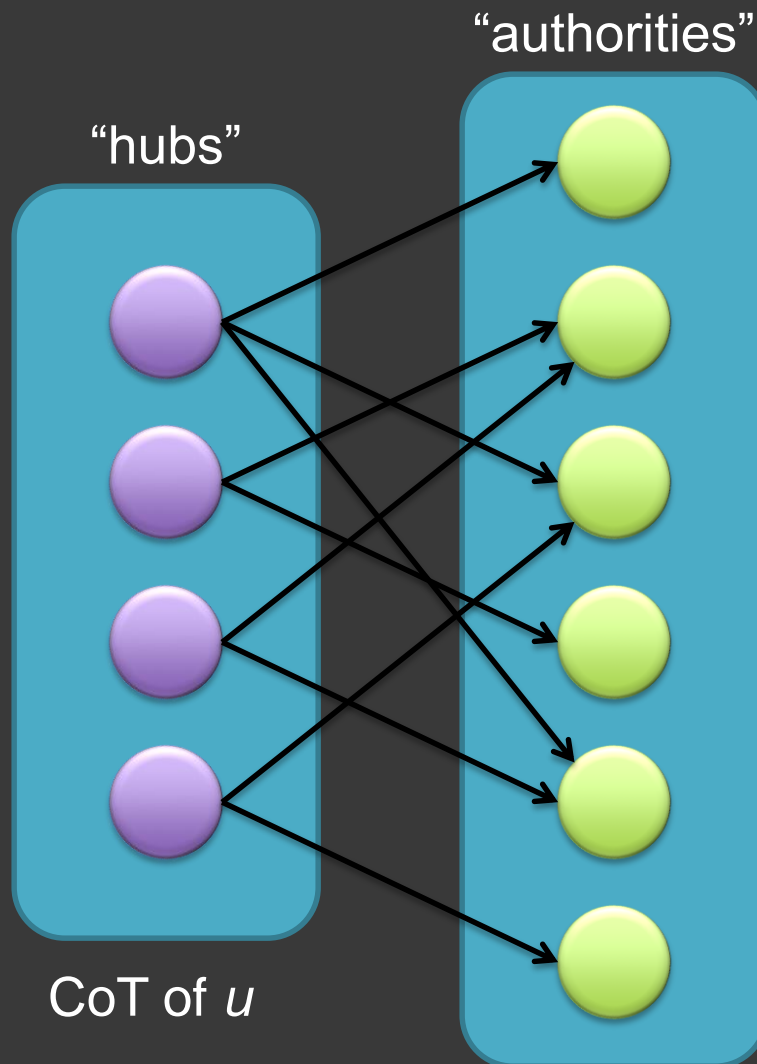
Result of egocentric random walk

Computed online based on various input parameters



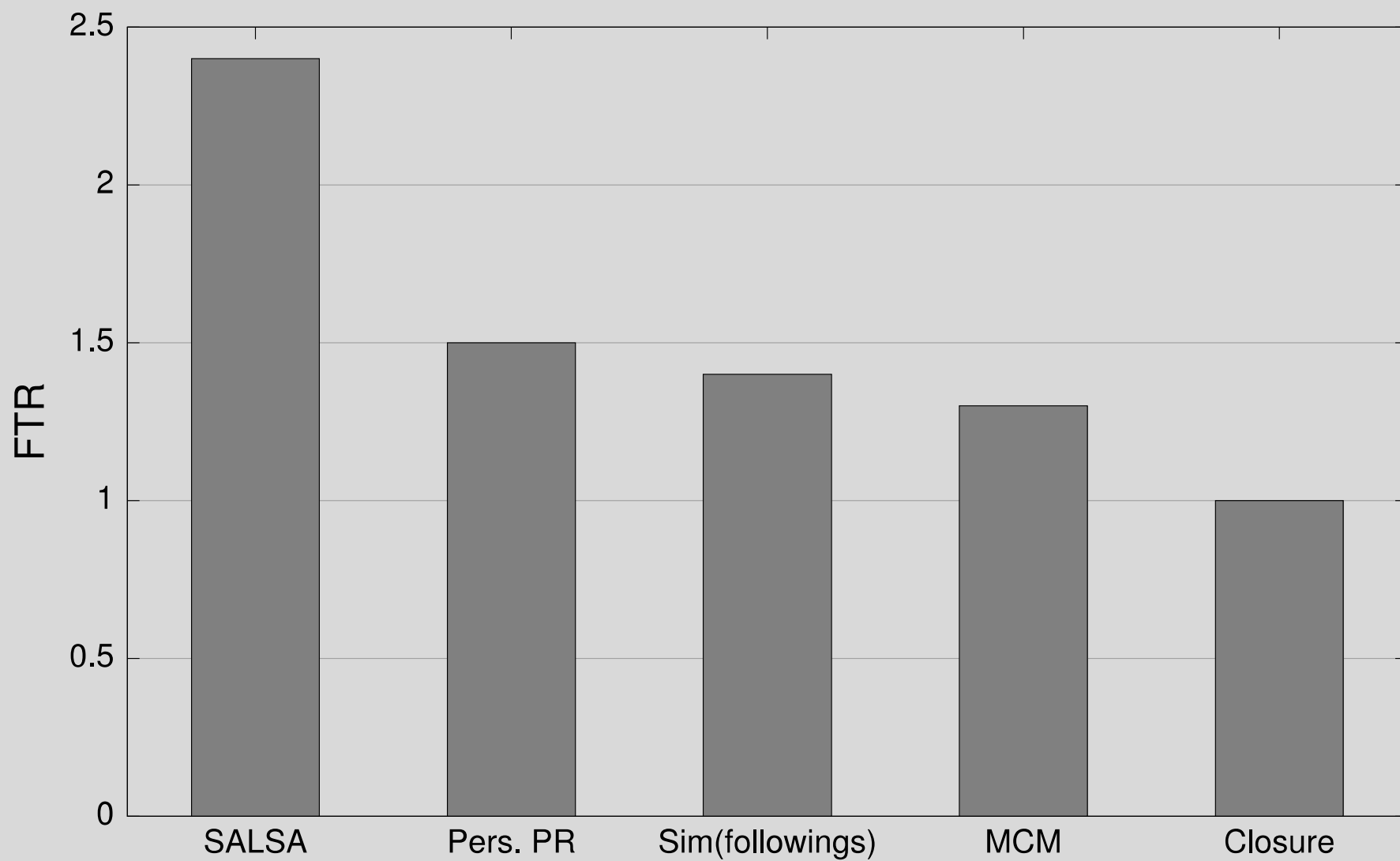
One of the features used in search

SALSA for Recommendations



hubs scores:
similarity scores to u

authority scores:
recommendation scores for u



Source: Gupta et al. "WTF: The Who to Follow Service at Twitter", WWW'13

Graph mining (Engagement)

- Connecting user actions to content (users, tweets, urls, images, video)
- Implicit relevance judgment
- Optimizing the performance of search and recommender systems



Graph mining (Heterogeneous)

- Recommendation of users/content may depend not just on user connectivity according to follow behavior but also on their actions and perceptions by others
- What sources of information to combine (clicks, tweets, favs, lists)
- How best to integrate this information (weighting factors, combination method)



User/content categorization



User Interest Modeling

- Users are interested in topics embodied by users they follow
- Users interact (read, retweet, fav) with tweets that seem interesting
- Users click/read urls corresponding to interesting content
- Users search for items that are of interest

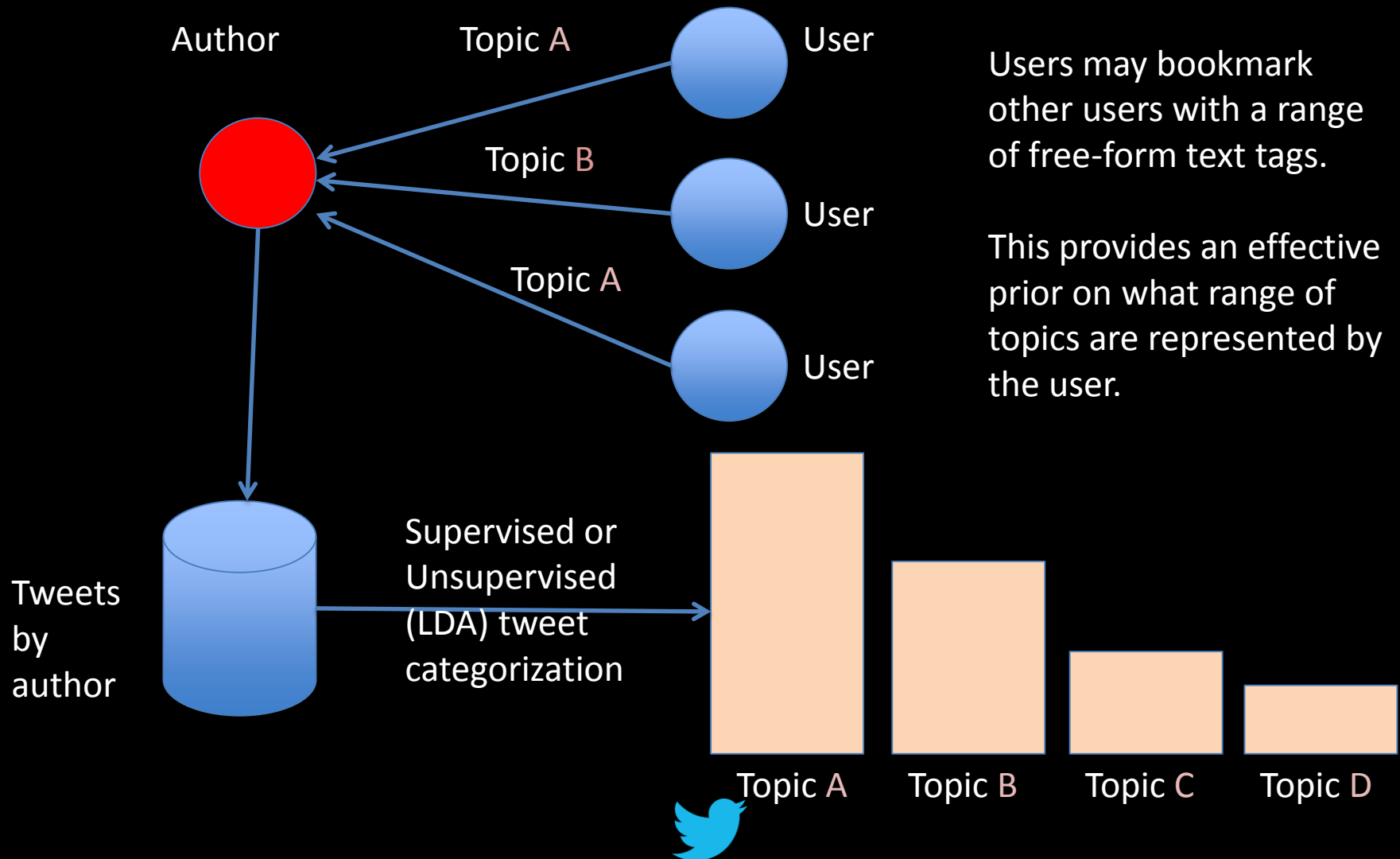


Challenges

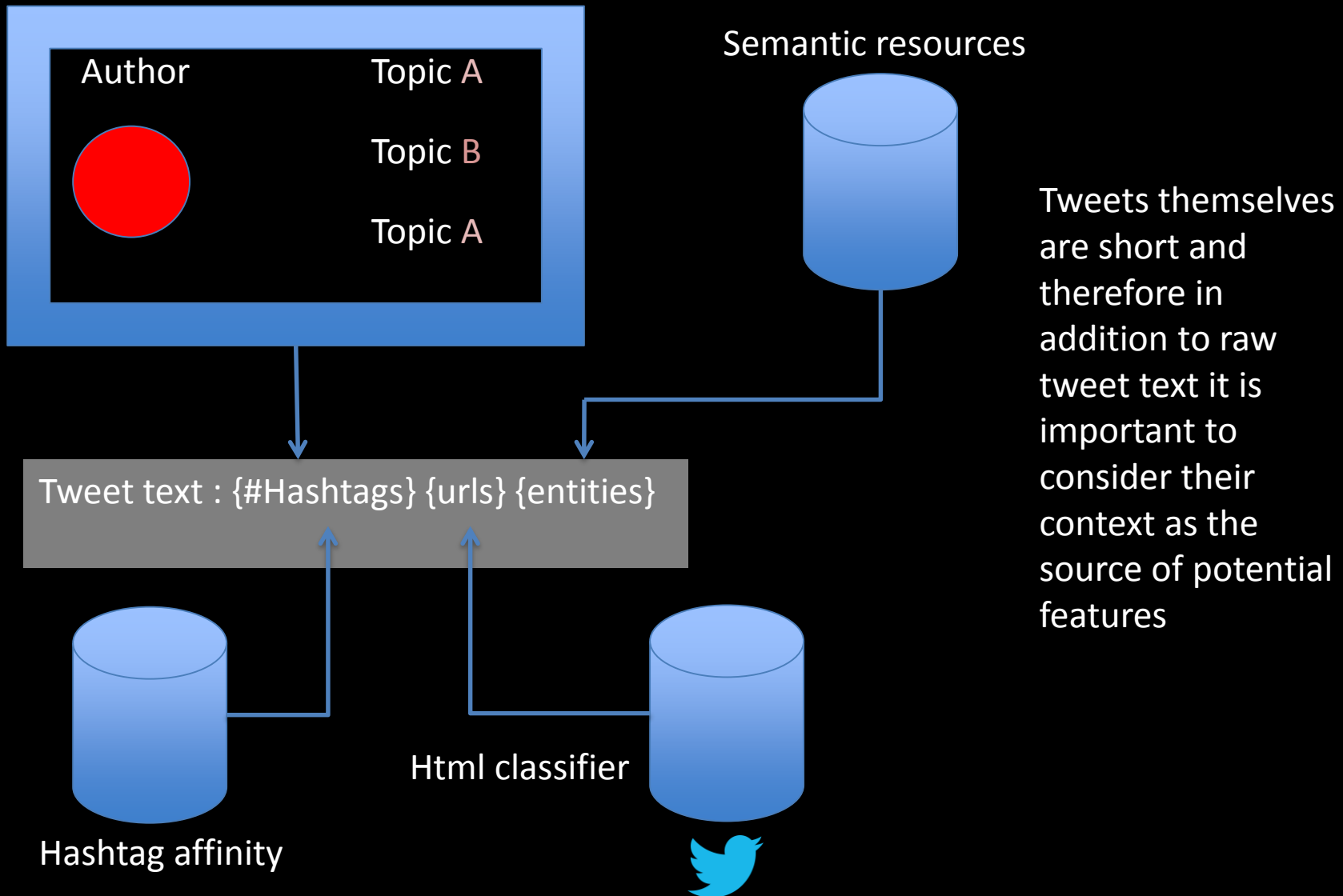
- A user may be followed for a number of reasons, e.g., family/friend relationships
 - Can we infer the reason?
- A user may tweet on a range of topics, and often “chatter” about their personal life
 - It is hard to assign labels to tweets by using the author
- Tweets are short and noisy
 - Can we infer the topic just looking at a tweet without context?
- User interests are time sensitive
 - Vary with trends, fads, current events, etc



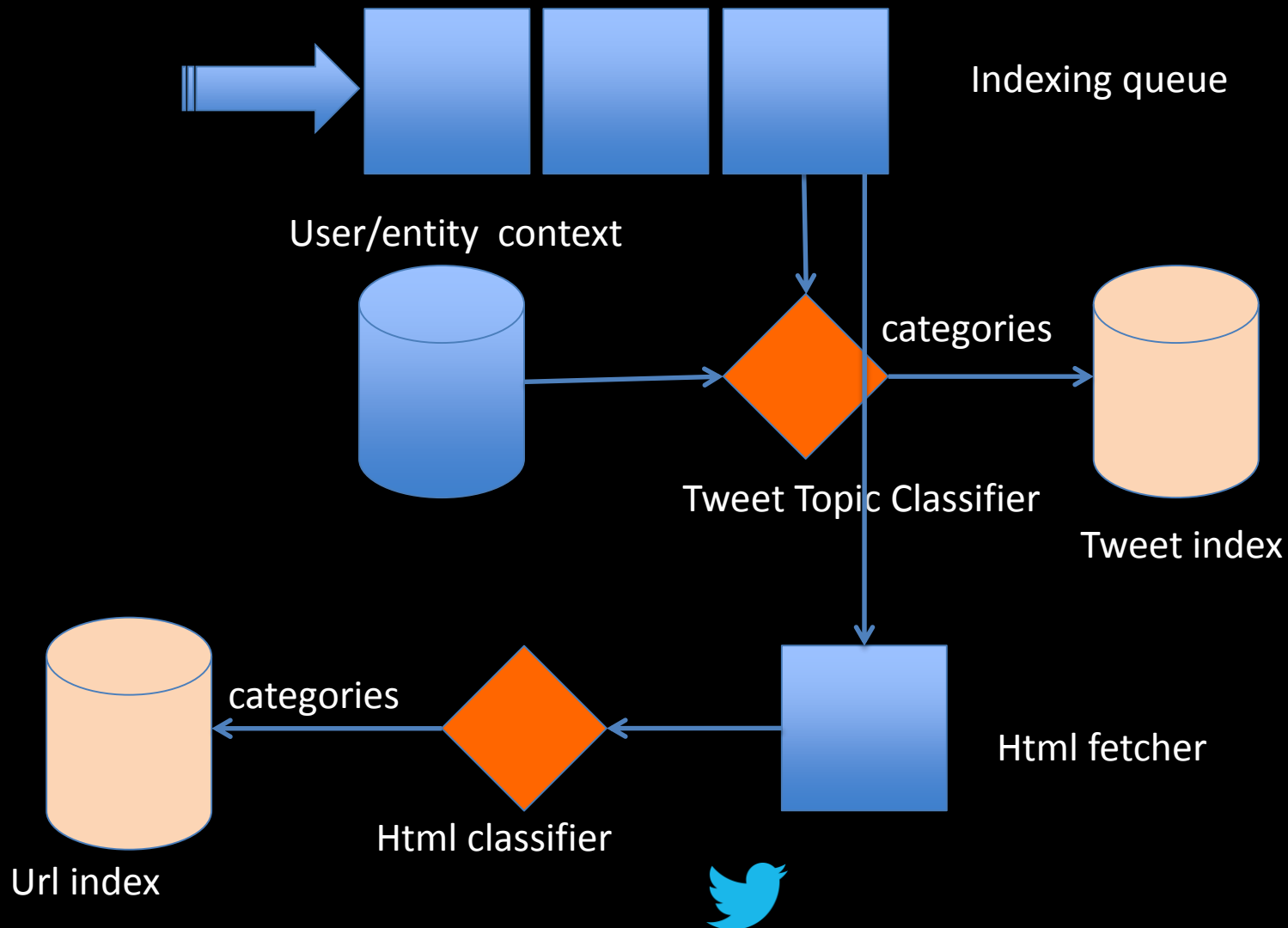
User topic attribution: producer side



Tweet topic attribution



Integration of tweet/url models



Multi-label classification at scale - 1

- Tweets are tagged with millions of rapidly evolving hashtags
- Users are tagged with millions of unique list names
- These are just two examples of folksonomy-type tag spaces that can be used to derive models to classify new content
- Each item can naturally belong to multiple categories, but it is unlikely that an item is labeled with all categories it can belong to
- How to curate labeled data to facilitate training and evaluation
- How to deal with label correlation



Multi-label classification at scale - 2

- Challenges of running tens of thousand of models (parallel evaluation)
- Automatic detection of synonymous topics (related hashtags, lists) even if surface similarity of tags is low
- Dealing with concept hierarchy (auto-generation)
- The role of external semantic resources (Wikipedia, Freebase, etc)



Chatter detection

- A lot of tweets correspond to personal updates
- These are certainly important but perhaps to a rather small audience
- In recommendation it is important to identify items of general interests (e.g., oriented about well defined topics)



[illegible]

Topical vs chatter

Topical Topics

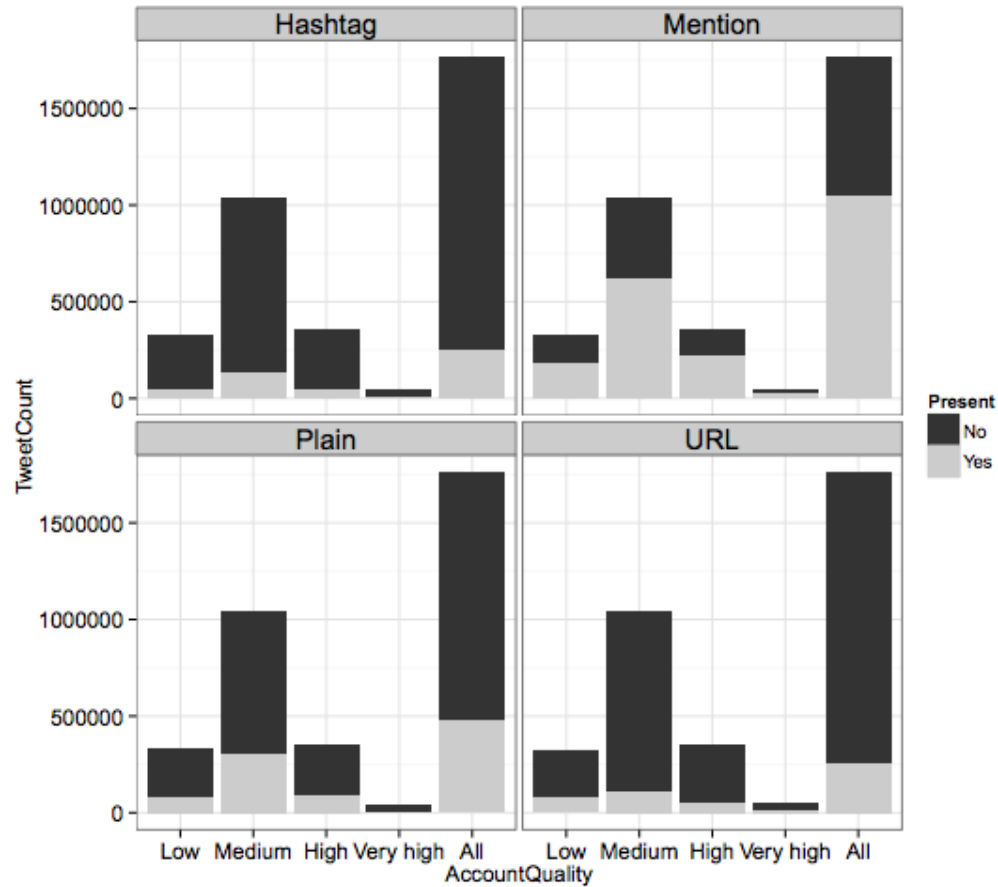
- * potter, harry, tumblr, fandom, weirdly, gifs, gif, hogwarts, fictional, rewatch, hp, rowlingpuns, obvs, deathly
- * rep., ballot, officials, mayor, gov, votes, investigation, ballots, mayor's, committee, capitol, district, courthouse, senate
- * donation, charity, donated, awareness, donations, donate, autism, raise, help, funds, fundraiser, fundraising, donating, support
- * investors, financial, banks, debt, markets, goldman, finance, banking, stocks, economic, earnings, ipo, bank, equity, investment

Chatter Topics

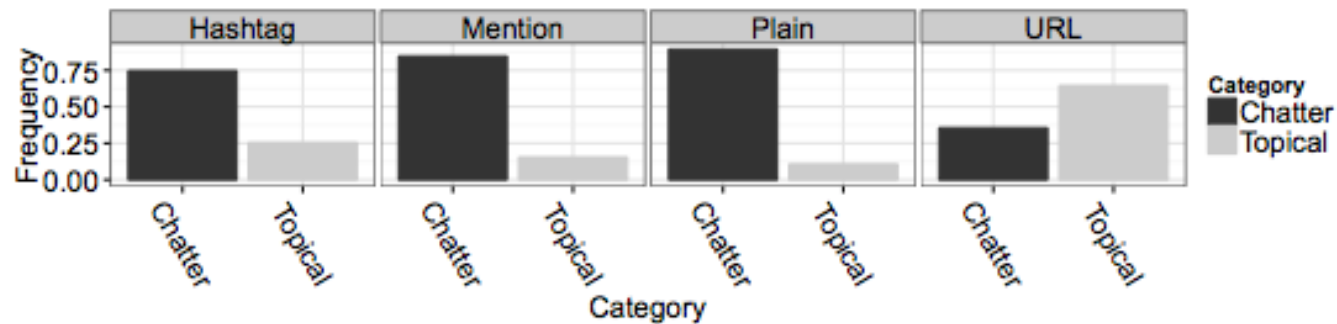
- * foh, lmaooo, lmaooooo, lml, lmaoo, deadass, henny, lmaooooo, niggas, djzeeti, smfh, nah, nigga, tho
- * coworker, washer, dangit, dryer, yeah, i've, beeping, oh, 6ish, i'll, 10ish, nope, 4am, kinda, probably
- * thanks, enjoyed, congrats, big, next, week, soon, coming, everyone, well, incredible, wow, meeting, achievement, weekend
- * someone, because, hate, sometimes, anymore, she, person, tell, her, aren't, saying, sleep, enough, without, ask, real, money



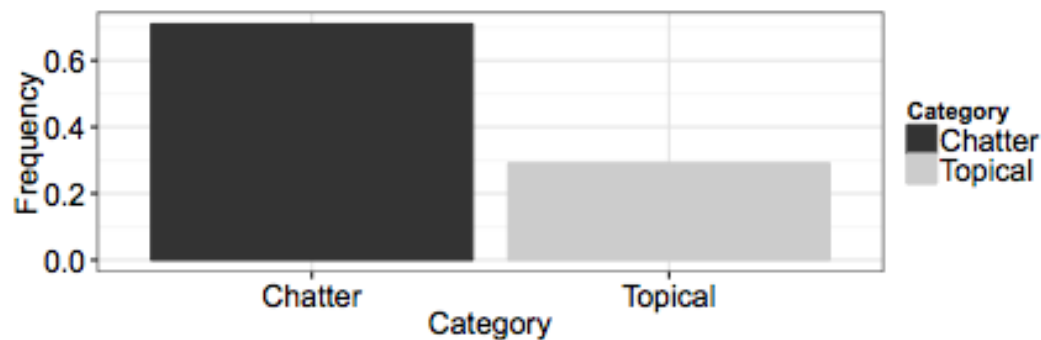
Chatter detection



Chatter detection



(b) Chatter incidence across tweets with different characteristics.



(c) Overall sample chatter incidence.



International support

- Over 70% of users reside outside of US
- Localization is important (over 35 languages)
- Content modeling in a variety of languages
- Topic distribution not the same for different languages
- Quantity (skew) of data across different languages
- Many users are multi-lingual



Bootstrapping topic models

- English tends to be the first choice in developing supervised or unsupervised topic models
- Acquisition of labeled data is expensive and time consuming
- Can an English corpus augmented with automatic translation technology be build models in other languages
- If additional labeled data in the target languages is needed (e.g., for evaluation), how to minimize the extra amount



Spammy interactions



Adversarial Classification (spam)

- Avoidance of being gamed by spammers
- Research into how spammers operate
- Types of spam
- Account creation vs purchasing



Adversarial games

- Spammers try to evade detection defenses
- The tug-of-war leads to the shift of tactics over time
- Tricks cover randomization, hiding of payload and trying to appear to be “normal”



Tweet level



dragonmxd @dragonmxd

19 Sep

dImhwloocdnnginsshvfoxf @kaylaide masvideo.info/PVkKnG

Expand Reply Retweet Favorite



H... @H...k6

#syria OMG!! This helped me laugh =)) HAHAAHAHAHAHAHA
bit.ly/RlxnFT

View summary Reply Retweet Favorite



Florida Macomb @FloridaMacomb

16 Sep

Hi there @T... ar Shirley had said you might be sincerely
interested in @... motion pofile

Expand



Florida Macomb @FloridaMacomb

16 Sep

Whats up @C... h Zachery had said you will be interested in
@... motion pofile

Expand Reply Retweet Favorite



Tweet-history level

swing 081

Expand

Whoa what a awesome day it is!

Expand

@Z... you've been chosen to be in the next D... video
click @Be... and follow the directions

[View conversation](#)

Whoa its such a awesome day;)

[Expand](#) [Reply](#) [Retweet](#) [Favorite](#)

22 Sep

once barnett1

[Expand](#) [Reply](#) [Retweet](#) [Favorite](#)

22 Sep

@The... We are filming near you for a brand new movie
that stars Justin Bleber ! We want you to be in it >>> @v...ing

[View conversation](#)

22 Sep

Tess swanson @Tess87ldn

in order that Catton123

[Expand](#)

18 Sep

18 Sep

18 Sep

23h



Norman... @norman...

This year's report has just been released and there are some
interesting trends that personal branders can draw from the report

[Expand](#)



Norman... @norman...

I will give you 5000 twitter followers in less than 24 hours

bit.ly/OQssfM @catton123 @Fam...ing

[Expand](#)

23h



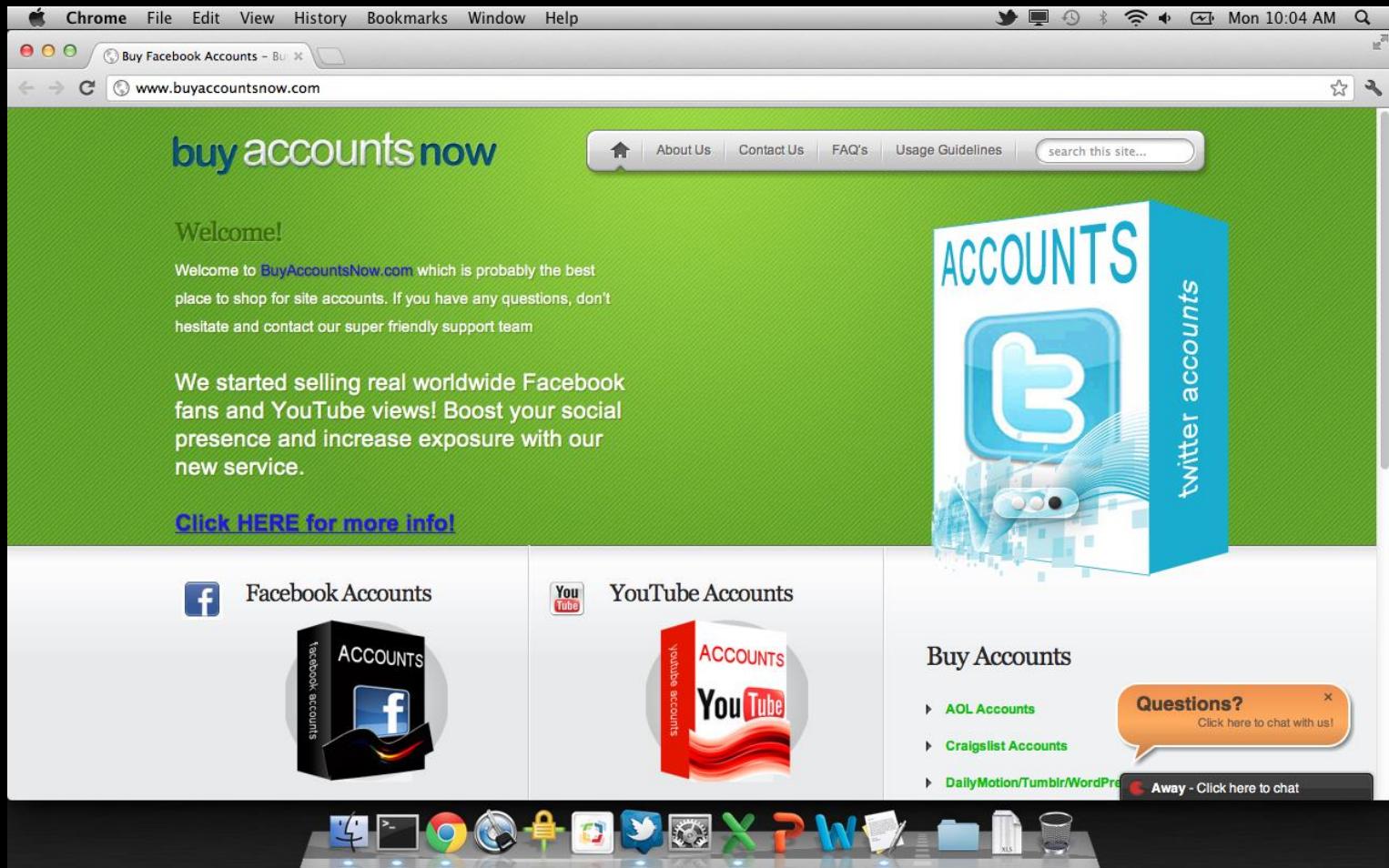
Norman... @norman...

If you are a Coca cola lover or a Coca cola drinker then you must see
this post in this post i have shown

[Collapse](#) [Reply](#) [Retweet](#) [Favorite](#)

23h

Account sales



Pricing Per Thousand Accounts



\$15-60



\$5



\$100



Legitimacy signals

- Email confirmation rate: 80%
- Unique IP signup rate: 87%
- Accounts are “pre-aged”
- Markets exist for both blank accounts and accounts with a social network
 - Prices vary



Big Data vs Fast Data

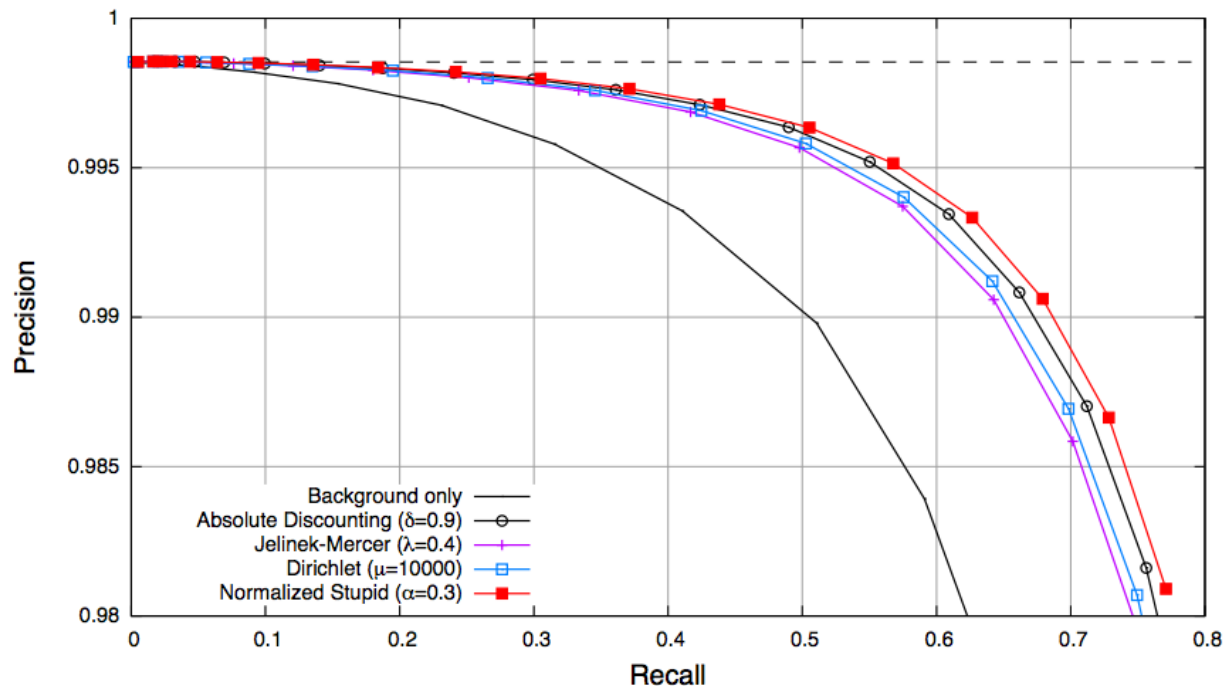


Hashtag topic tracking

Extrinsic Evaluation: Results

Unigram LM

Topic 1: #nfl



Normalized stupid backoff is at least as good as other smoothing techniques

“Real-time” human computation

- Spiking queries often correspond to new events/topics.
 - Cold start problem
- Semantics of the intent may be quite different from the “usual” semantics of the query.
 - “Clint Eastwood” during the 2012 presidential debate
- The lifetime of the spike is short-lived.
 - A lot can be learned from user actions clicks, but by that time the query may have stopped trending



“Real-time” human computation (crowdsourcing)

- Solutions – ask humans!
- Crowdsourcing systems allow one to have question answered with short turnaround time
 - Low cost allows one to ask many questions.
- Humans-in-the-loop can be automated
 - Use of MTurk and internal APIs
- Quality control
 - Custom pool of judges



Plumbing

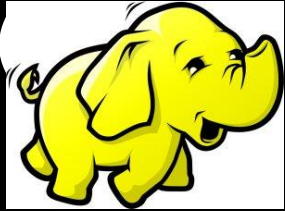


Scalable analytics/infrastructure

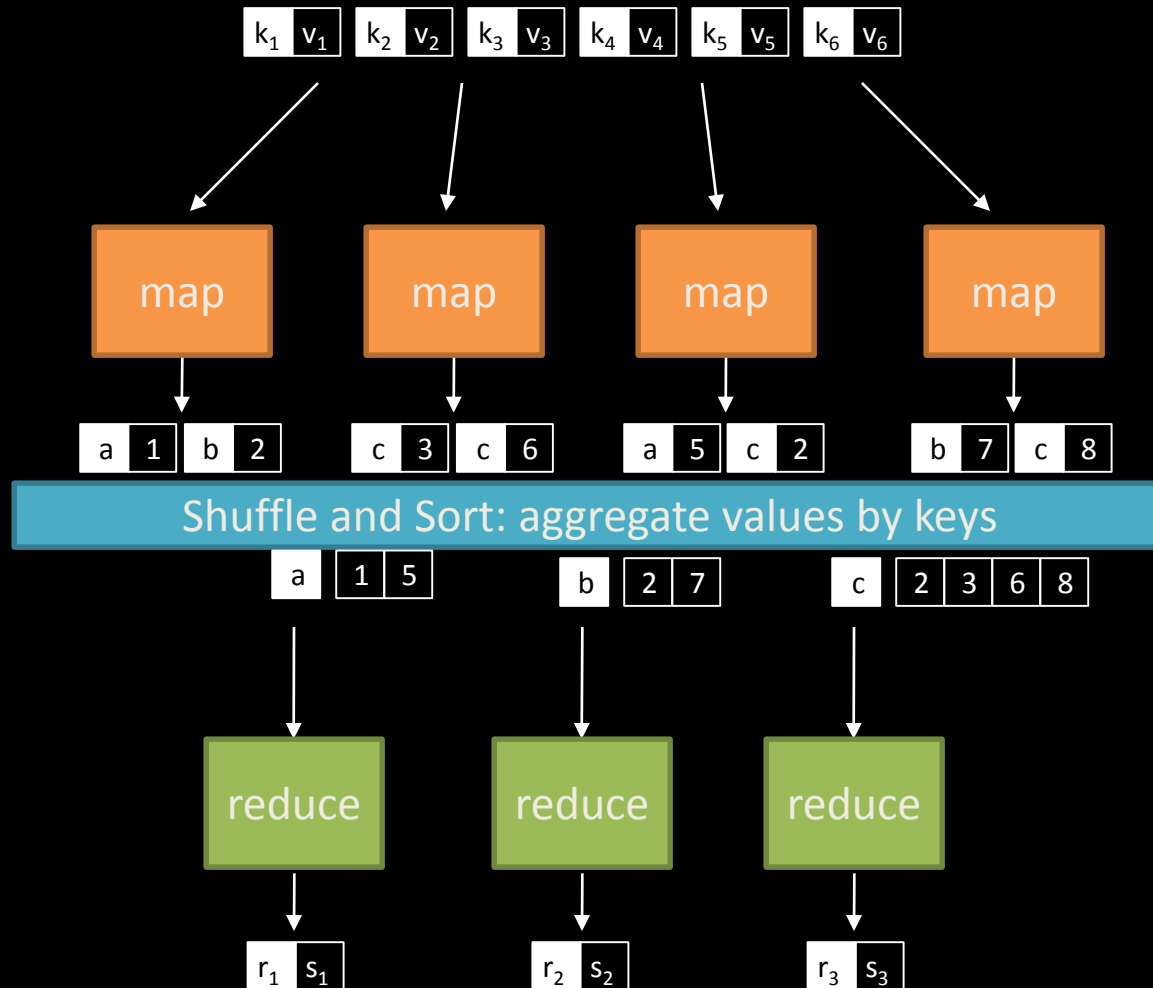
- Logging/archiving and computing at the same time
- Ability to recover from and catch up from disruptions (crashes, outages, etc)
- Planning for the needs of an ever-growing number of analytics jobs



MapReduce/Hadoop



Pig



```

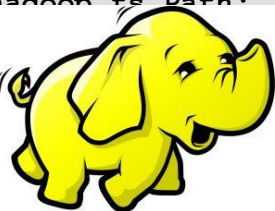
t java.io.IOException;
t java.util.ArrayList;
t java.util.Iterator;
t java.util.List;

```

```

t org.apache.hadoop.fs.Path;
t org.apache.hadoop.io.Writable;
t org.apache.hadoop.io.Text;
t org.apache.hadoop.io.LongWritable;
t org.apache.hadoop.io.IntWritable;
t org.apache.hadoop.io.FloatWritable;
t org.apache.hadoop.io.DoubleWritable;
t org.apache.hadoop.io.BytesWritable;
t org.apache.hadoop.io.NullWritable;
t org.apache.hadoop.mapred.JobConf;
t org.apache.hadoop.mapred.KeyValueTextInputFormat;
t org.apache.hadoop.mapred.Mapper;
t org.apache.hadoop.mapred.MapReduceBase;
t org.apache.hadoop.mapred.OutputCollector;
t org.apache.hadoop.mapred.Reporter;

```



```

t org.apache.hadoop.mapred.Mapper;
t org.apache.hadoop.mapred.MapReduceBase;
t org.apache.hadoop.mapred.OutputCollector;
t org.apache.hadoop.mapred.Reporter;

```

```

c class MRExam
public static c
implements

public void
    Out
    Rep

```

```

visits      = load '/data/visits' as (user, url, time);
gVisits     = group visits by url;
visitCounts = foreach gVisits generate url, count(urlVisits);
urlInfo     = load '/data/urlInfo' as (url, category, pRank);
visitCounts = join visitCounts by url, urlInfo by url;
gCategories = group visitCounts by category;
topUrls     = foreach gCategories generate
              top(visitCounts,10);

store topUrls into '/data/topUrls';

```

```

// Pull the key out
String line = val.toString();
int firstComma = line.indexOf(',');
String key = line.substring(0, firstComma);
String value = line.substring(firstComma + 1);
Text outKey = new Text(key);
// Prepend an index to the value so we know which file
// it came from.
Text outVal = new Text("1" + value);
// Collect (outKey, outVal)

```



```

if (value.charAt(0) == 'C') {
    first.add(value.substring(1));
    else second.add(value);
    reporter.setStatus("C");
}

```

```

// Do the cross product
for (String s1 : first) {
    for (String s2 : second) {
        String outval = key + s1 + s2;
        oc.collect(null, outval);
        reporter.setStatus(outval);
    }
}
}
}

```

```

public static class LoadJoined extends MapReduceBase
implements Mapper<Text, Text, Text, Text> {

    public void map(
        Text k,
        Text val,
        OutputCollector<Text, Text> oc,
        Reporter reporter) throws IOException {
        // Find the url
        String line = val.toString();
        int firstComma = line.indexOf(',');
        int secondComma = line.indexOf(',', firstComma + 1);
        String key = line.substring(0, firstComma);
        // drop the rest of the line
        // just use a 1 for the index
        Text outKey = new Text(key);
        oc.collect(outKey, new Text("1" + val));
    }
}

```



```

public static class ReduceUrls extends MapReduceBase
implements Mapper<Text, LongWritable, Text, LongWritable> {

    public void reduce(
        Text key,
        Iterator<LongWritable> values,
        OutputCollector<Text, LongWritable> oc,
        Reporter reporter) throws IOException {
    }
}

```

```

public void reduce(
    Text key,
    Iterator<LongWritable> values,
    OutputCollector<Text, LongWritable> oc,
    Reporter reporter) throws IOException {
    // Find the url
    String line = key.toString();
    int firstComma = line.indexOf(',');
    int secondComma = line.indexOf(',', firstComma + 1);
    String key = line.substring(0, firstComma);
    // drop the rest of the line
    // just use a 1 for the index
    Text outKey = new Text(key);
    oc.collect(outKey, new Text("1" + val));
}
}

```

Scalable machine learning

- Large scale machine learning on Hadoop
 - Big data streaming with stochastic gradient descent
 - Parallel training with data randomization (linear and tree-based ensembles)
- Trade-offs between learning with true MR jobs vs using MR for pre-processing and doing the learning part on a cluster of large boxes with hdfs access
 - There a limitation of RAM size of a reducer (e.g., 3-6GB)
 - Hadoop clusters tend to be composed from smaller boxes
 - Sometimes it is necessary to iterate over a dataset of significant size (e.g., 100-200G B)



What kind of models?

- Most of our data is text based
- Models should be able to consume large amount of training data
 - and possibly a large number of categories
- Emphasis on simplicity and speed



Model training UDF internals

- A single node in the Hadoop cluster does not have extensive memory resources
- The learner cannot cache too much data
- Natural fit for (**not iterating over all data**):
 - Stochastic gradient descent (SGD), possibly with mini-batching
 - Effective for streaming the whole dataset through a single learner



Supervised classification in a nutshell

Given $D = \left\{ \left(\mathbf{x}_i, y_i \right) \right\}_i^n$ label
(sparse) feature vector

Induce $f : X \rightarrow Y$ s.t. loss is minimized

empirical loss = $\frac{1}{n} \sum_{i=0}^n \ell(f(\mathbf{x}_i), y_i)$ loss function

Consider functions of a parametric form:

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(\mathbf{x}_i; \theta), y_i)$$

model parameters

Key insight: machine learning as an optimization problem!
(closed form solutions generally not possible)



Gradient Descent

$$w^{(t+1)} = w^{(t)} + \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla l(f(x_i; \theta^{(t)}), y_i)$$

“batch” learning: update model after considering all training instances

Stochastic Gradient Descent (SGD)

$$w^{(t+1)} = w^{(t)} + \gamma^{(t)} \nabla l(f(x; \theta^{(t)}), y)$$

“online” learning: update model after considering *each* (randomly-selected) training instance

In practice... just as good!

Solves the iteration problem!

What about the single reducer problem?

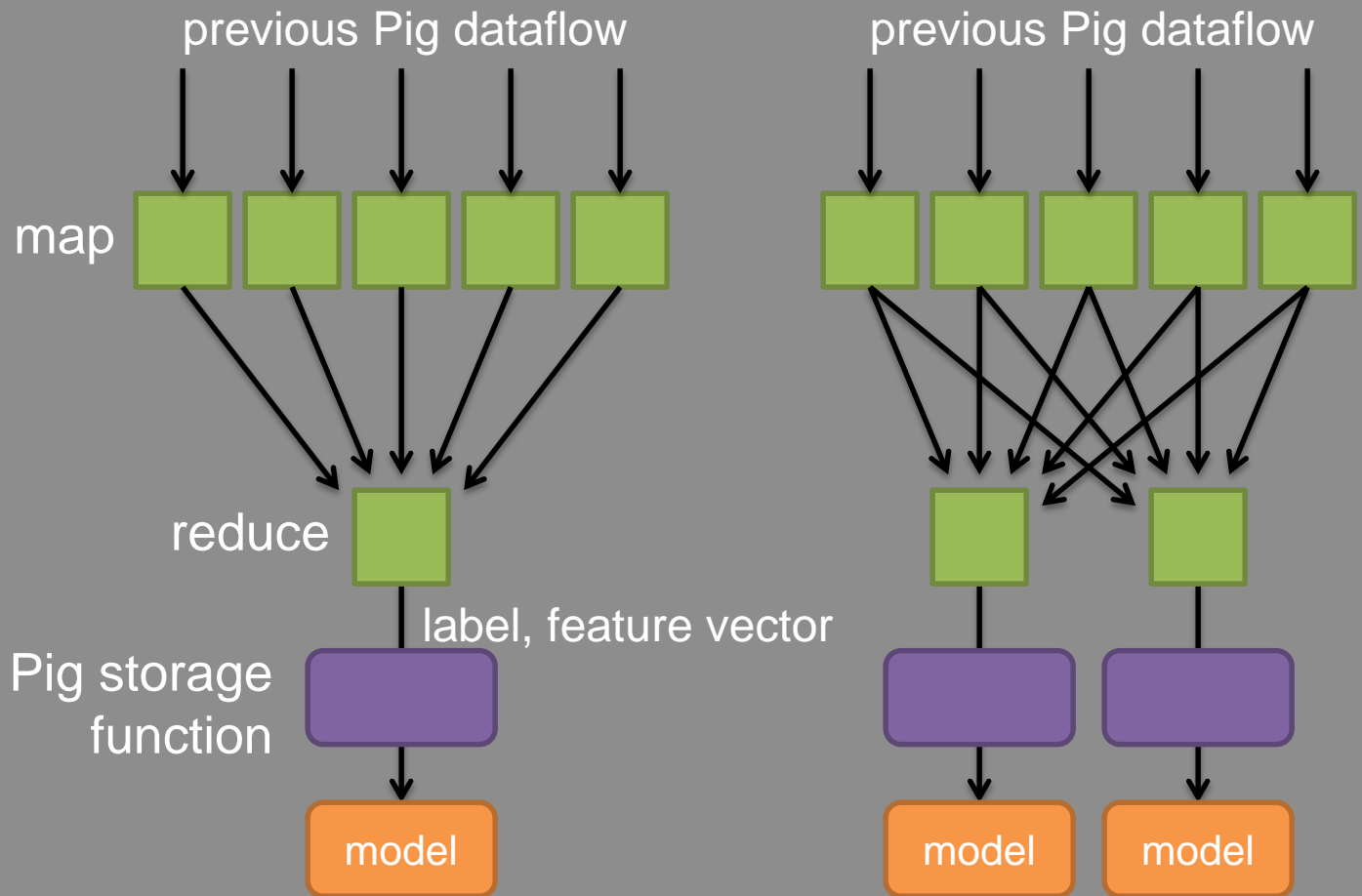


Another way - Ensembles

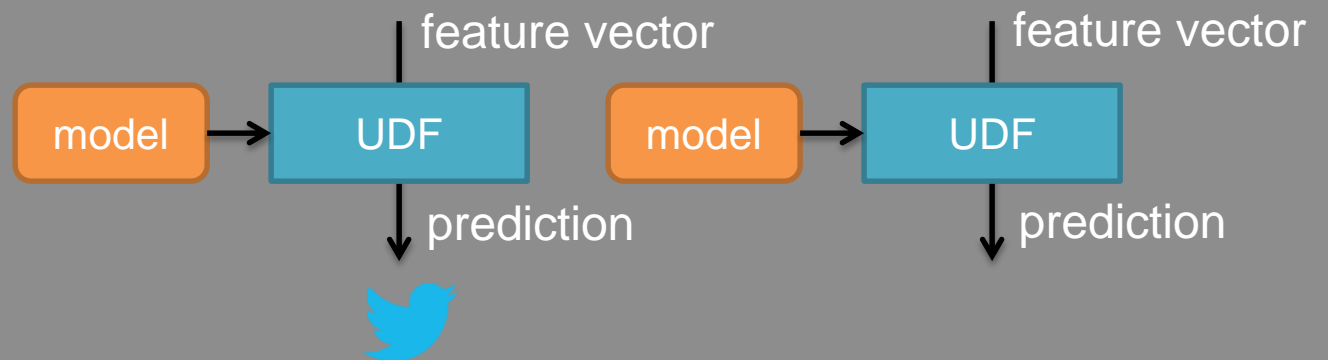
- Classifier committees are one of the best performing types of learners
- Some of these algorithms are sequential (not very MR friendly)
 - Boosting
- But others rely mostly on randomization
 - Each learner is trained over a different split (features and/or instances) of the data



Classifier Training



Making Predictions

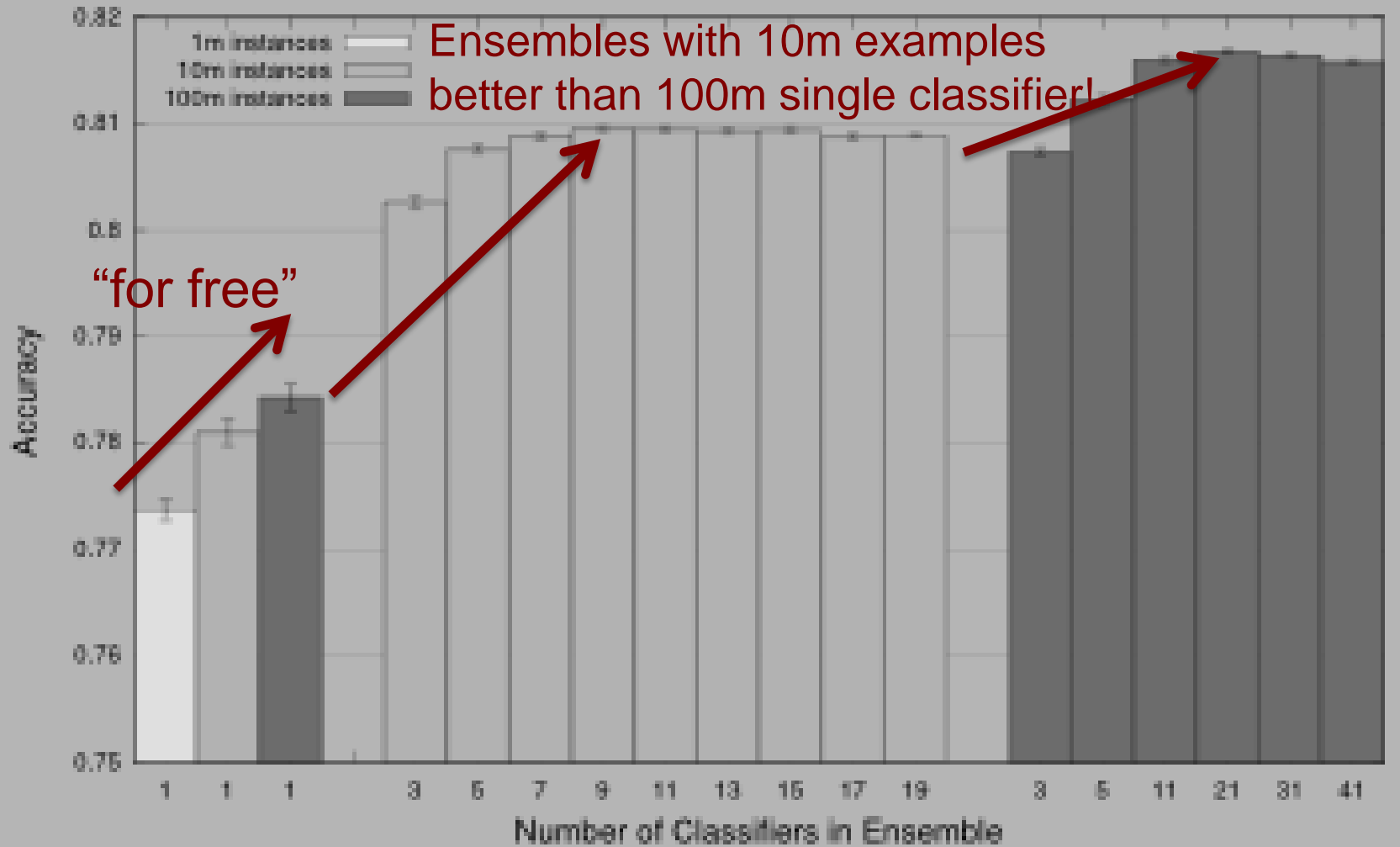


Example: tweet sentiment detection

- Training/Test data: tweets with *emoticons*
 - 😊 😞
- Emoticons provide surrogate labels for training
 - Emoticons are removed from the data
- Logistic Regression trained over character 4grams
- Single classifier vs. use of ensembles



Diminishing returns...



single classifier

10m ensembles

100m ensembles



Stream-based machine learning

- Twitter is largely about real-time
- Global and personal relevance models need to respond quickly to the changes of content distribution as well as user-action history.
- Technologies to handle real-time streams
 - Storm (open sourced)
 - Several application specific custom solutions
 - Write-through caching
- Learning
 - Updating models vs. updating features
 - Stochastic gradient descent (e.g., with Logistic Regression)



Big boxes have a role

- Cassowary
 - Twitter's WTF service served by redundant single box systems 2010-2012
- Big-box supervised learning
 - 200GB 12-core boxes at the upper end of “commodity”



THANK YOU

