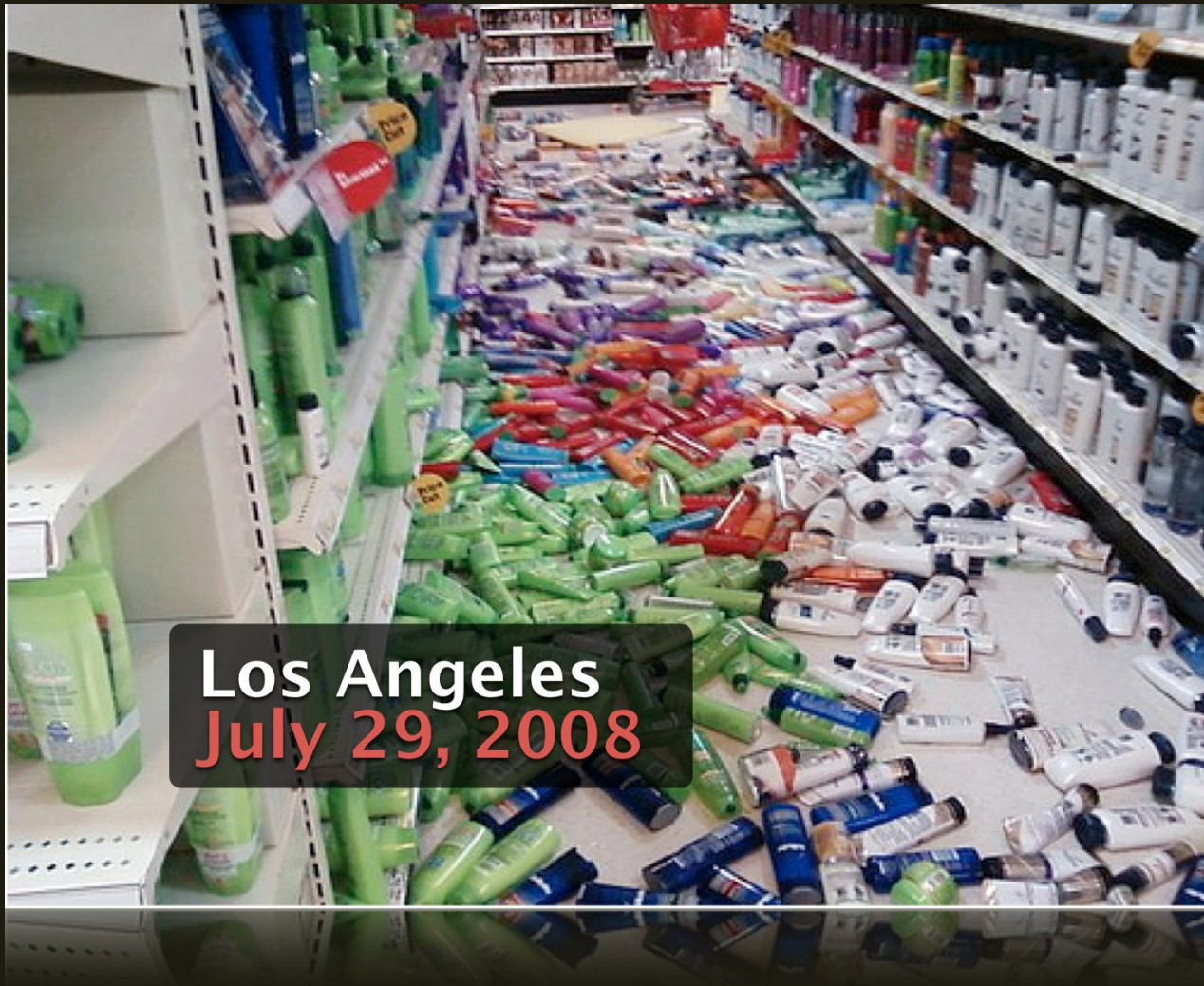


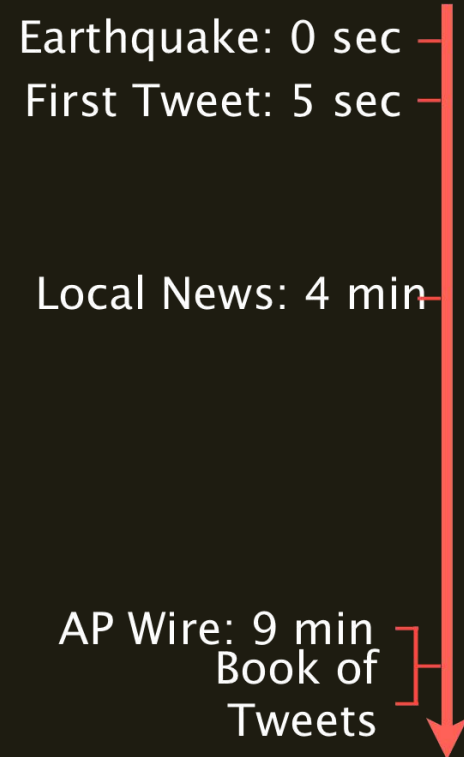
Large Scale (Machine) Learning at Twitter

Aleksander Kołcz
Twitter, Inc.

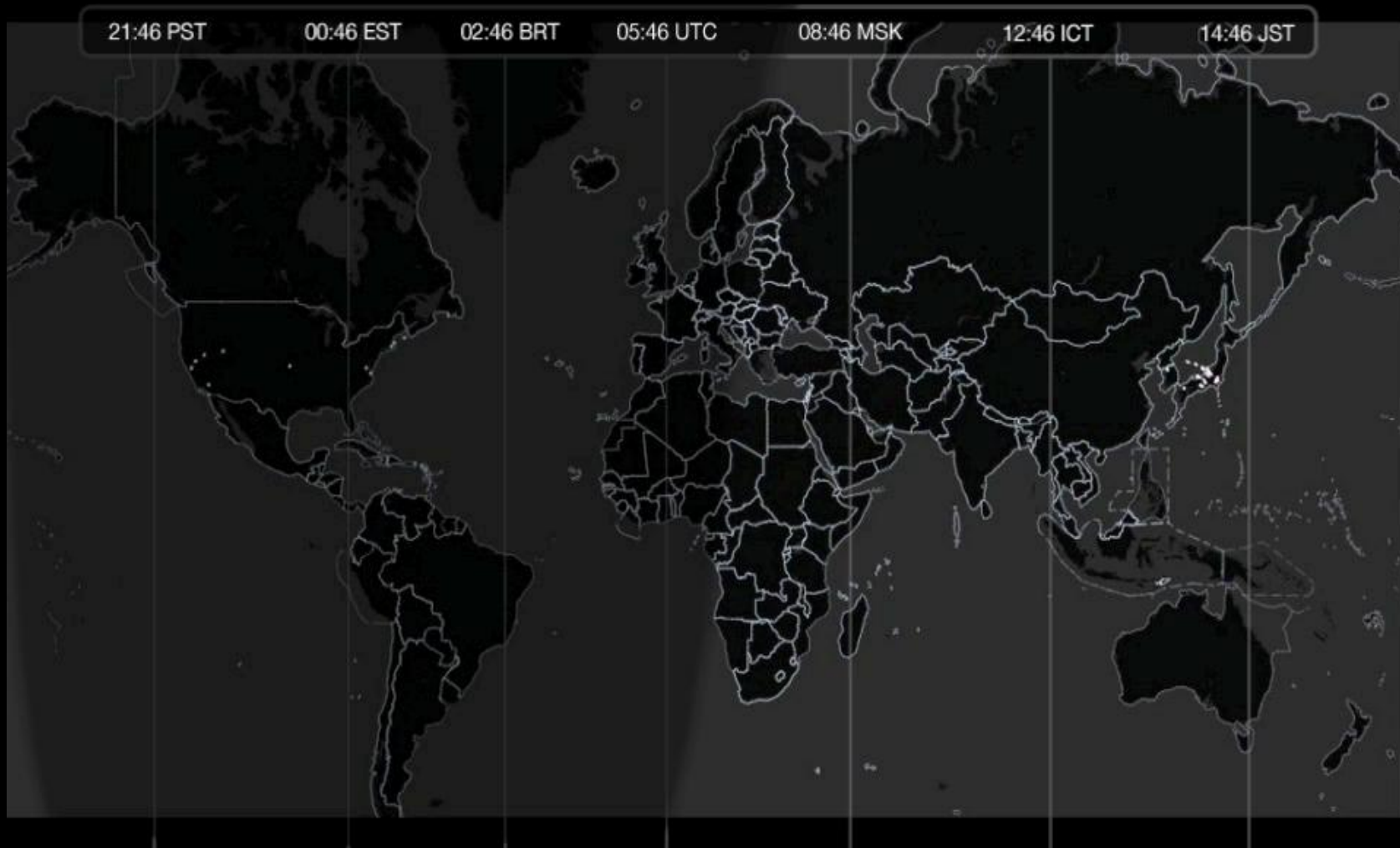




Los Angeles
July 29, 2008





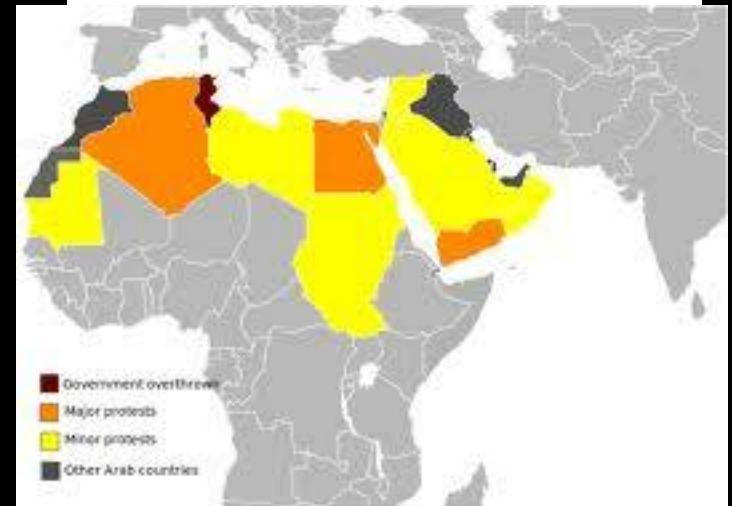


What is Twitter

- Micro-blogging service
- Open exchange of information/opinions
 - Private communication possible via DMs
- Content restricted to 140 characters
- Seems tiny but ...



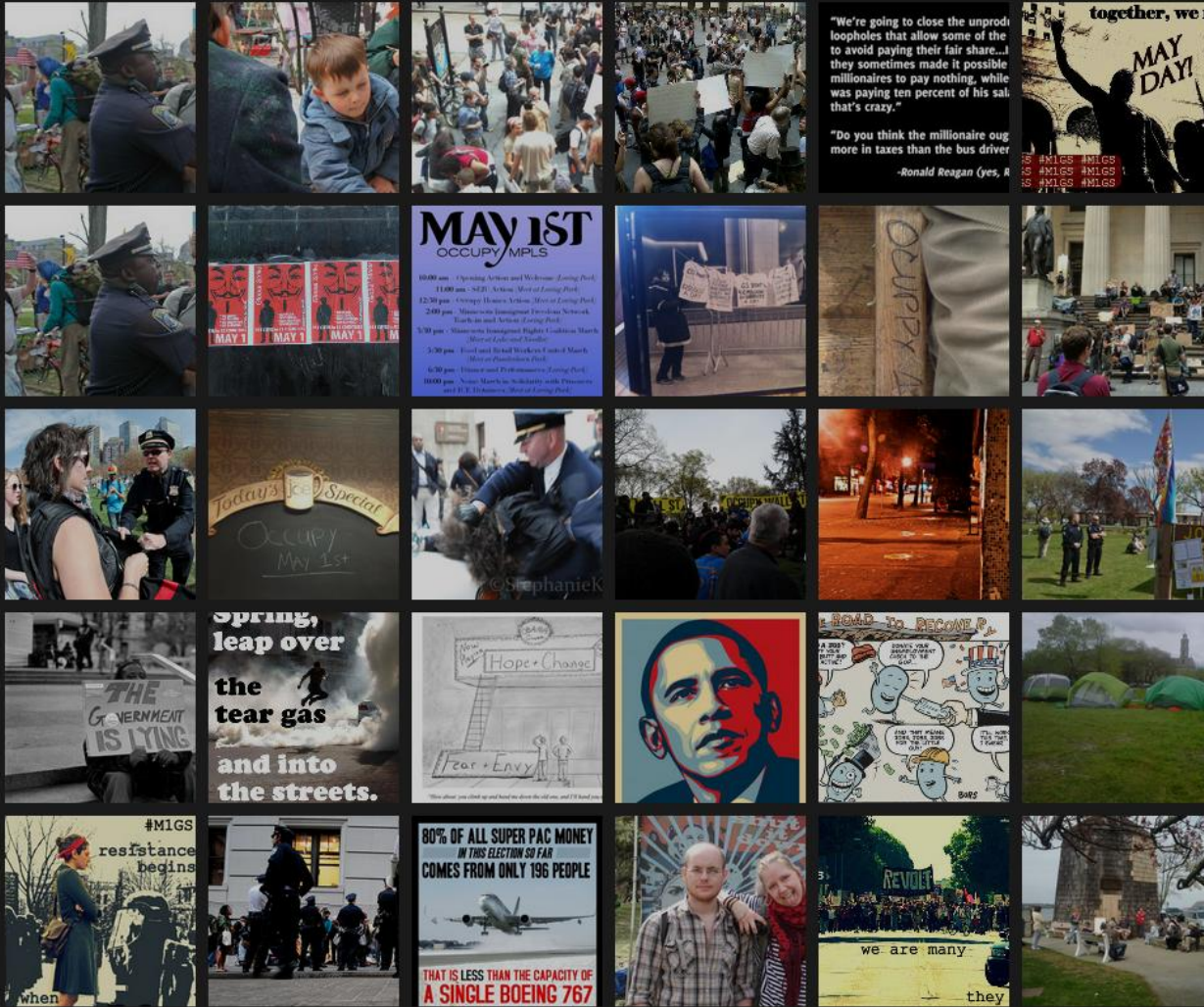
#egypt, #tunisia, #libya, #syria ...



#ocws, #occupywallstreet, ...

[← Back to search results](#)

Top images for #ows



The Scale of Twitter

- Twitter has more than 140 million active users
- 340 million Tweets are sent per day
- Over 400 million monthly unique visitors to twitter.com
- 50 million people log into Twitter every day



Large scale infrastructure of information delivery

- Users interact via web-ui, sms, and various apps
- Over 55% of our active users are mobile users
- Real-time redistribution of content
- Plus search and other services living on top of it
 - E.g., 2.3B search queries/day (26K/sec)

Support for user interaction

- Search
 - Relevance ranking
- User recommendation
 - WTF or Who To Follow
- Content recommendation
 - Relevant news, media, trends
- ML is an important component but really just a part of a larger whole

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 .

twitter 

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...



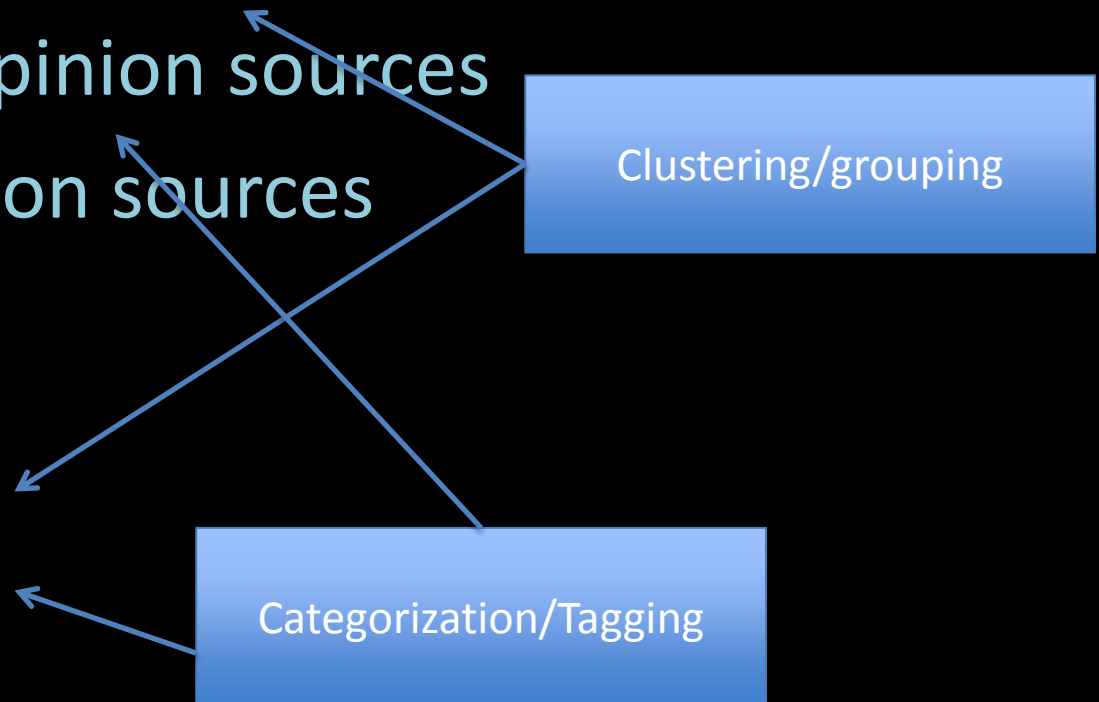
Tweeted by people who share your interests

(other) problems we are trying to solve

- Trending topics
- Language detection
- Anti-spam
- Revenue optimization
- User interest modeling
- Growth optimization

Recommendation/Personalization

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



Is this BIG data ?



Challenges

- Twitter is International
 - 70% of accounts are outside the US
 - Twitter supports more than 28 different languages
- Twitter is real time
 - The “now” factor is very important for people interacting with Twitter
 - The concept of relevance is highly time dependent
- 140 characters : “documents” are very short and vocabulary contains many non-standard acronyms
 - @XXXX_ Cool
 - Please please please



What type of machine learning?

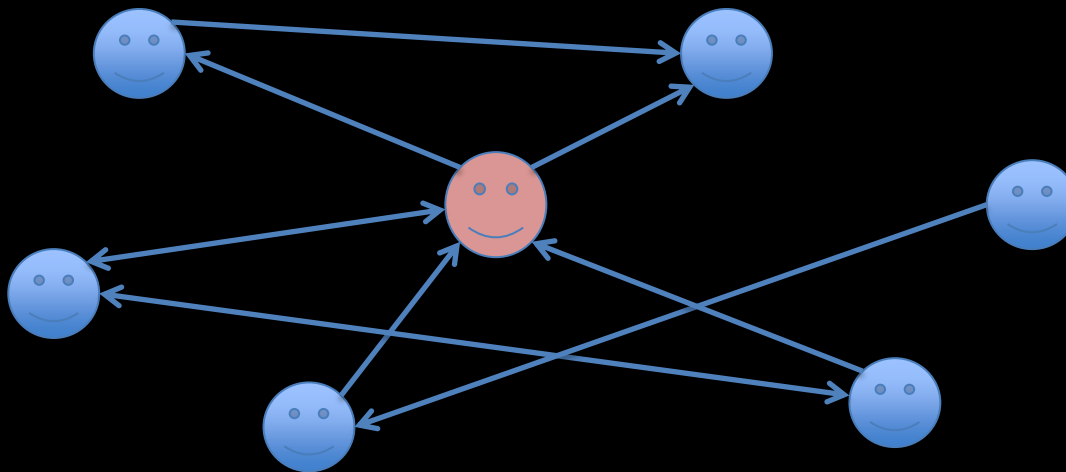
- At this point ML is an extensive and diverse field
- Algorithms differ widely in terms of implementational and operational complexity
- Decisions which techniques to use are driven by the nature of the problems and infrastructure/operational constraints

ML for social networks

- Power of simple models
 - It has been found again and again that with enough data relatively simple models (e.g., Naïve Bayes or Logistic Regression) work remarkably well
 - Ease of integration with data processing flows (both training and inference) trumps model sophistication

ML for social networks

- Power of graph aggregation
 - Even an imperfect signal (e.g., classifier) can be useful when its effect on a graph node is taken in the context of the node's neighborhood



ML for social networks

- Challenges of scalability and adaptive processing
 - Learning and inference need to handle data streams and combinations thereof
 - Models need to be able to quickly adapt to data stream change

Are there wheels not to reinvent?

- Temptation is always there
- Data processing infrastructure?
- Core ML algorithm libraries ?
- Data pre/post processing ?
- Visualization ?
- Glue code ?



Analytics Ecosystem

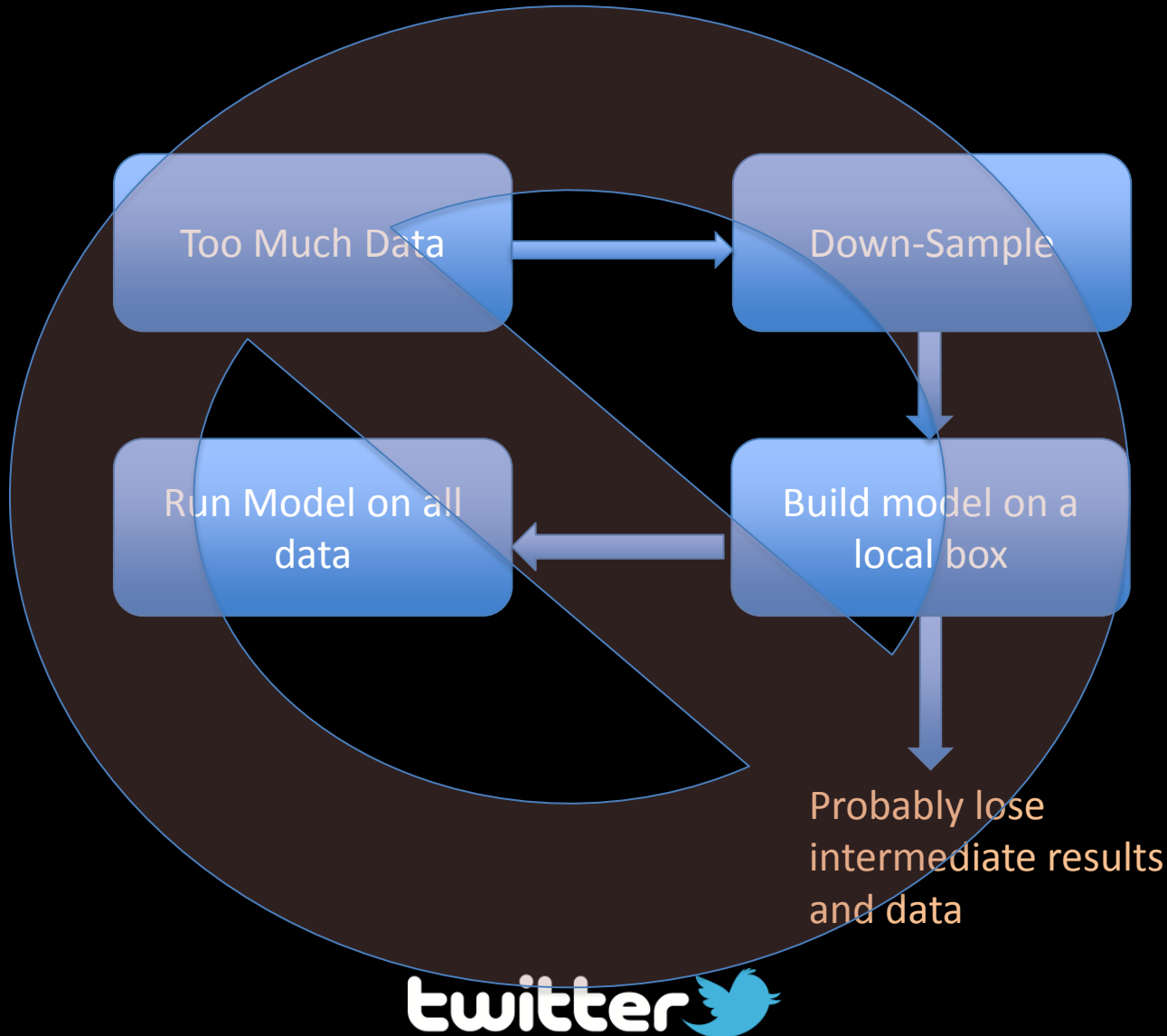
- Apache Mesos



Maximizing the use of Hadoop

- We cannot afford too many diverse computing environments
- Most of analytics job are run using Hadoop cluster
 - Hence, that's where the data live
 - It is natural to structure ML computation so that it takes advantage of the cluster and is performed close to the data

AVOID: “janky” analysis of messy data



Leveraging off-line tools

- While the data is big, a lot of useful feature components can be learned with smaller datasets
- The ML tool ecosystem provides a wide range of options for optimization and tuning various models
- Once tuned, the models can be applied to big data in a distributed fashion
 - often not as final models but as feature extractors
- The key is not to rely primarily on ad-hockery in production pipelines

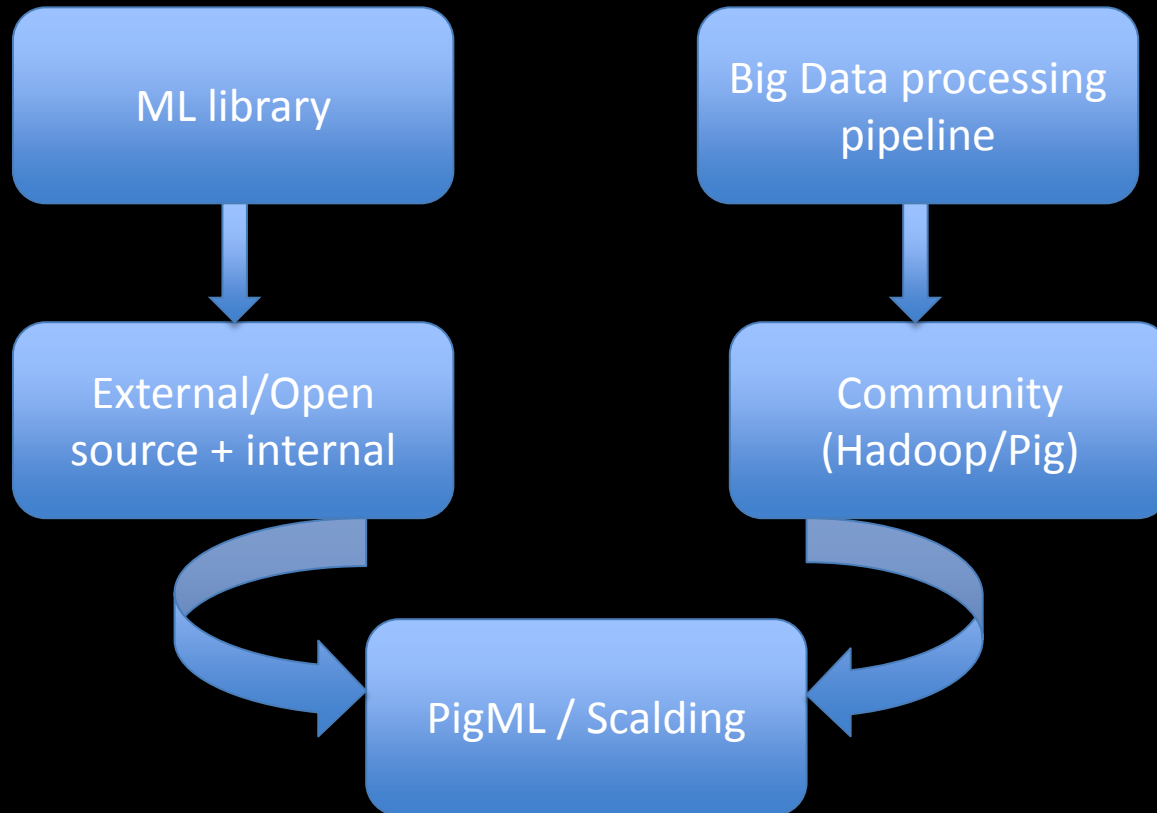
Large scale learning frameworks

- Participation in the open source community
 - There are a number of initiatives for using ML over Hadoop (e.g., Mahout)
- Upsides
 - Larger support and developer network
- Downsides
 - Not always convenient to integrate with internal analytics and data processing flows

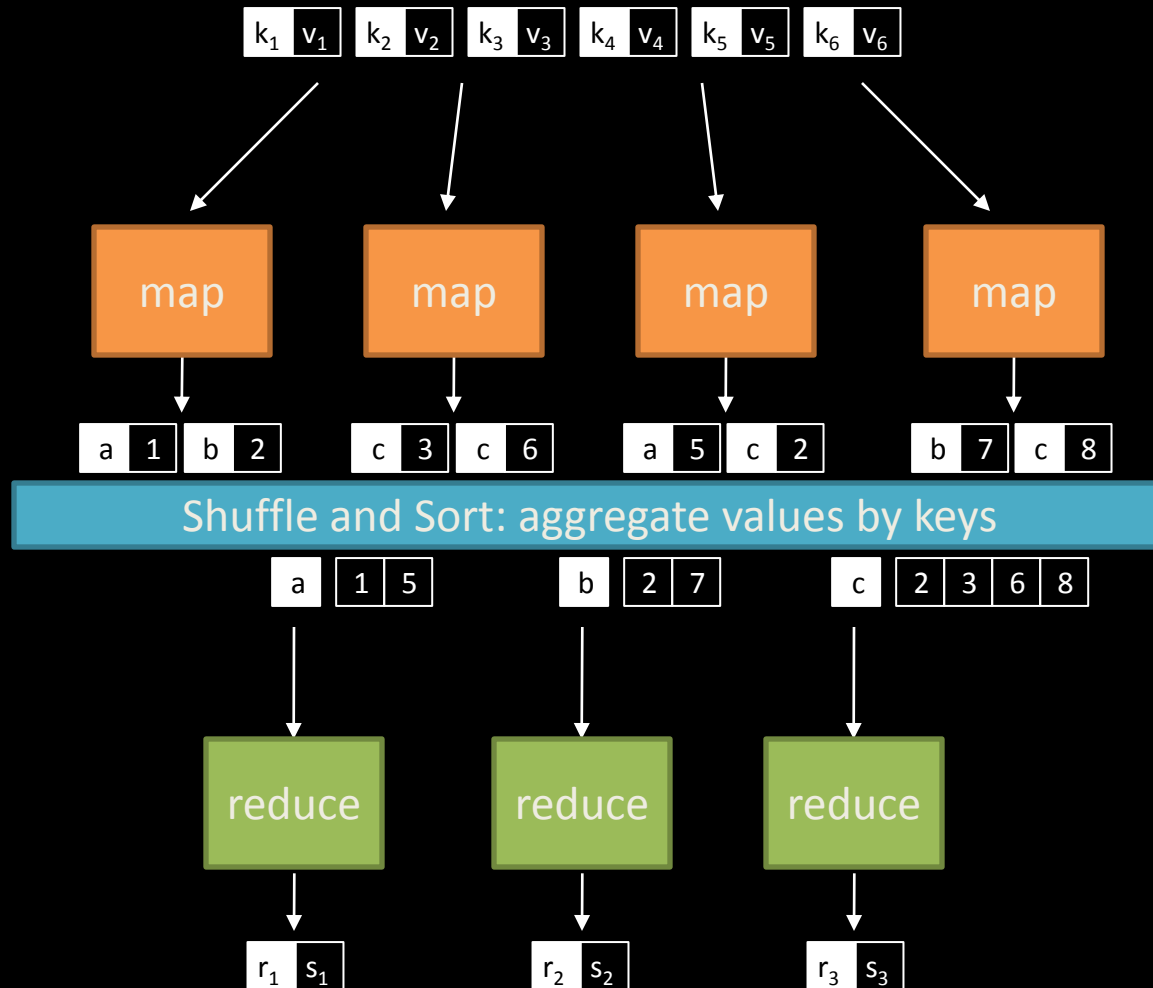
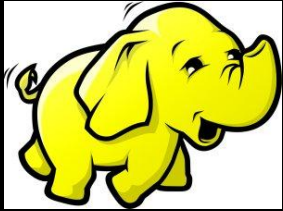
Our extensions

- PigML
 - A library of UDFs wrapping the ML functionality
- Scalding/PyCascading
 - Cascading with Scala/Jython
- ML Java lib
 - Used in both cases to capture the low level ML functionality

Build/reuse/integrate



MapReduce



```
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';
```

so we know which file

twitter

```
// Do the cross product a
for (String s1 : first) {
    for (String s2 : seco
        String outval = k
        oc.collect(null,
        reporter.setStatu
    }
}
```

```
public void reduce(
    Text key,
    Iterator<LongWritable> values,
    OutputCollector<Writable> output,
    Reporter reporter) throws IOException {
```

PigML

- Embedding ML learning/eval functionality directly in Pig
- Flows naturally with other data processing operations
- Minimal learning curve for users



Training a model in Pig

```
training = load 'trn_data' using  
  piggybank.ml.Storage() as (target: double,  
  features: map[]);
```

```
store training into 'model-LR' using  
  piggybank.ml.train.online.LRClassifierBuilder('with  
  Pegasos withLambda:0.1');
```

Training a model in Pig

```
training = load 'trn_data' using  
    piggybank.ml.Storage() as (target: double,  
    features: map[]);  
  
store training into 'model-LR' using  
    piggybank.ml.train.online.LRClassifierBuilder('with  
    Pegasos withLambda:0.1');
```


Training a model in Pig

It's just a store function!

```
training = load 'trn_data' using  
    piggybank.ml.Storage() as (target: double,  
    features: map[]);
```

```
store training into 'model-LR' using  
    piggybank.ml.train.online.LRClassifierBuilder('with  
    Pegasos withLambda:0.1');
```

Applying a model in Pig

```
DEFINE Classify
    piggybank.ml.classify.ClassifyFeaturesWithLRClassifier('model-LR');

data = load 'test_data' using piggybank.ml.Storage() as (target: double,
    features: map[]);

data = foreach data generate target, Classify(features) as prediction;

results = foreach data generate (target == prediction.label ? 1 : 0) as
    matching;

dump results;
```

Applying a model in Pig

```
DEFINE Classify
    piggybank.ml.classify.ClassifyFeaturesWithLRClassifier('model-LR');

data = load 'test_data' using piggybank.ml.Storage() as (target: double,
    features: map[]);

data = foreach data generate target, Classify(features) as prediction;

results = foreach data generate (target == prediction.label ? 1 : 0) as
    matching;

dump results;
```

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:
 - Stochastic gradient descent (SGD), possibly with mini-batching
 - Effective for streaming the whole dataset through a single learner

Supervised classification in a nutshell

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

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

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

Consider functions of a parametric form:

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

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!

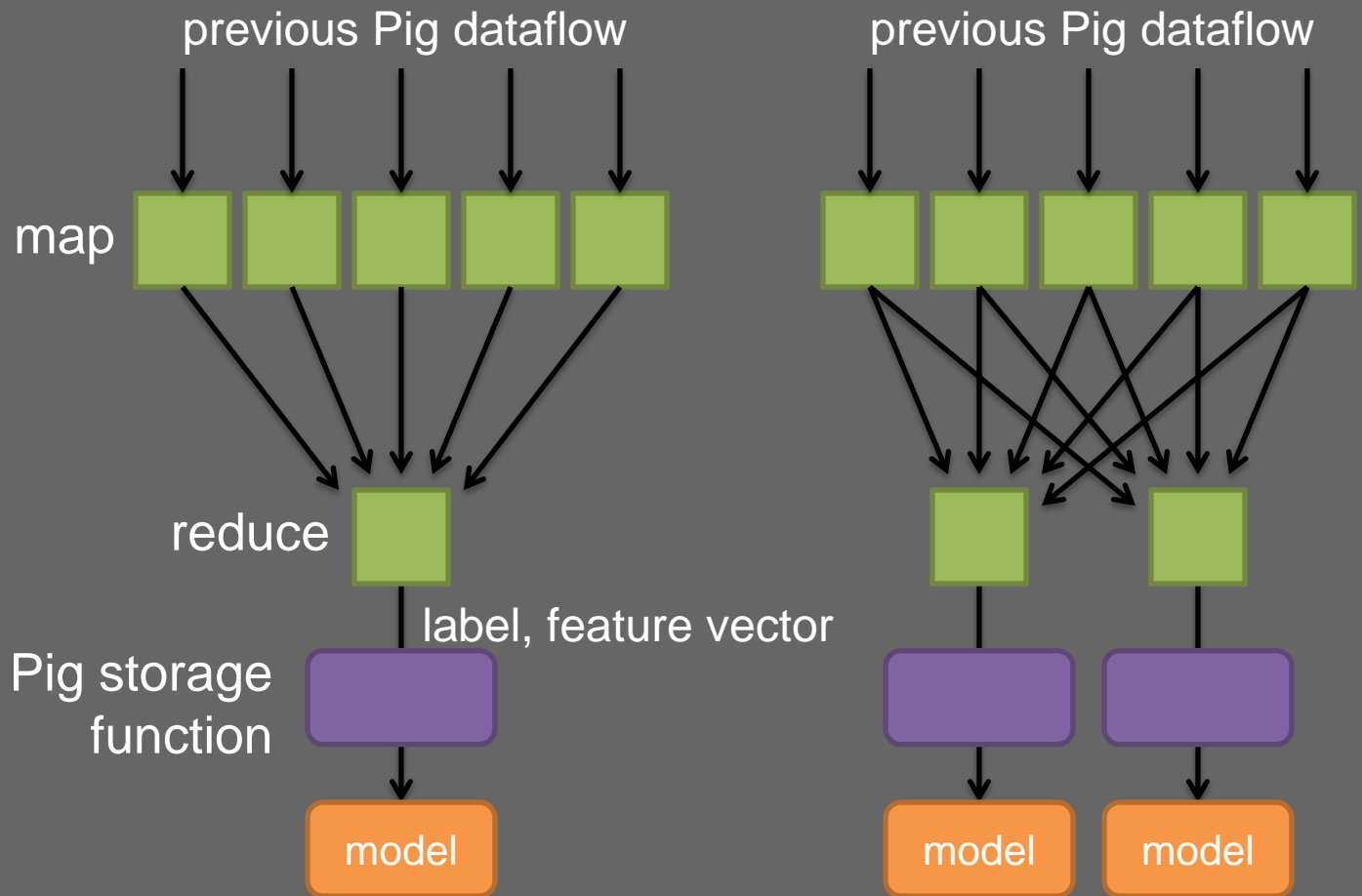
What about the single reducer problem?



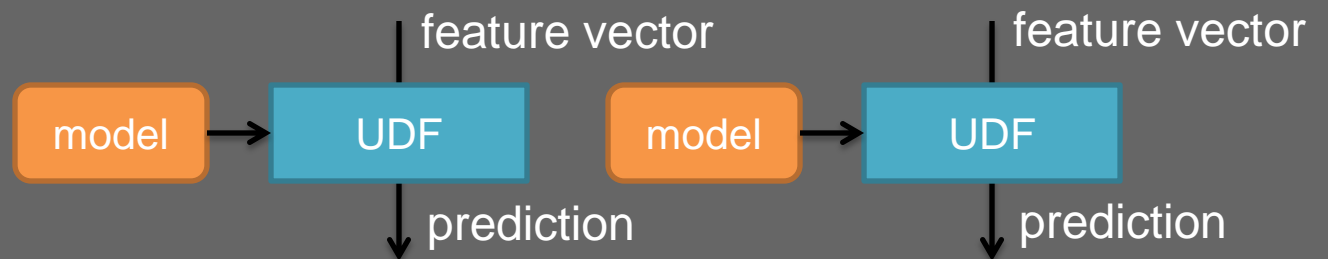
Ensembles

- 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



Ensembles: continued

- Ensembles of linear classifiers
 - E.g., each trained using SGD over different subset of features
- Ensembles of decision trees (random forest)
 - Each tree is seeing only a subset of the dataset and node split variables are randomized at each split

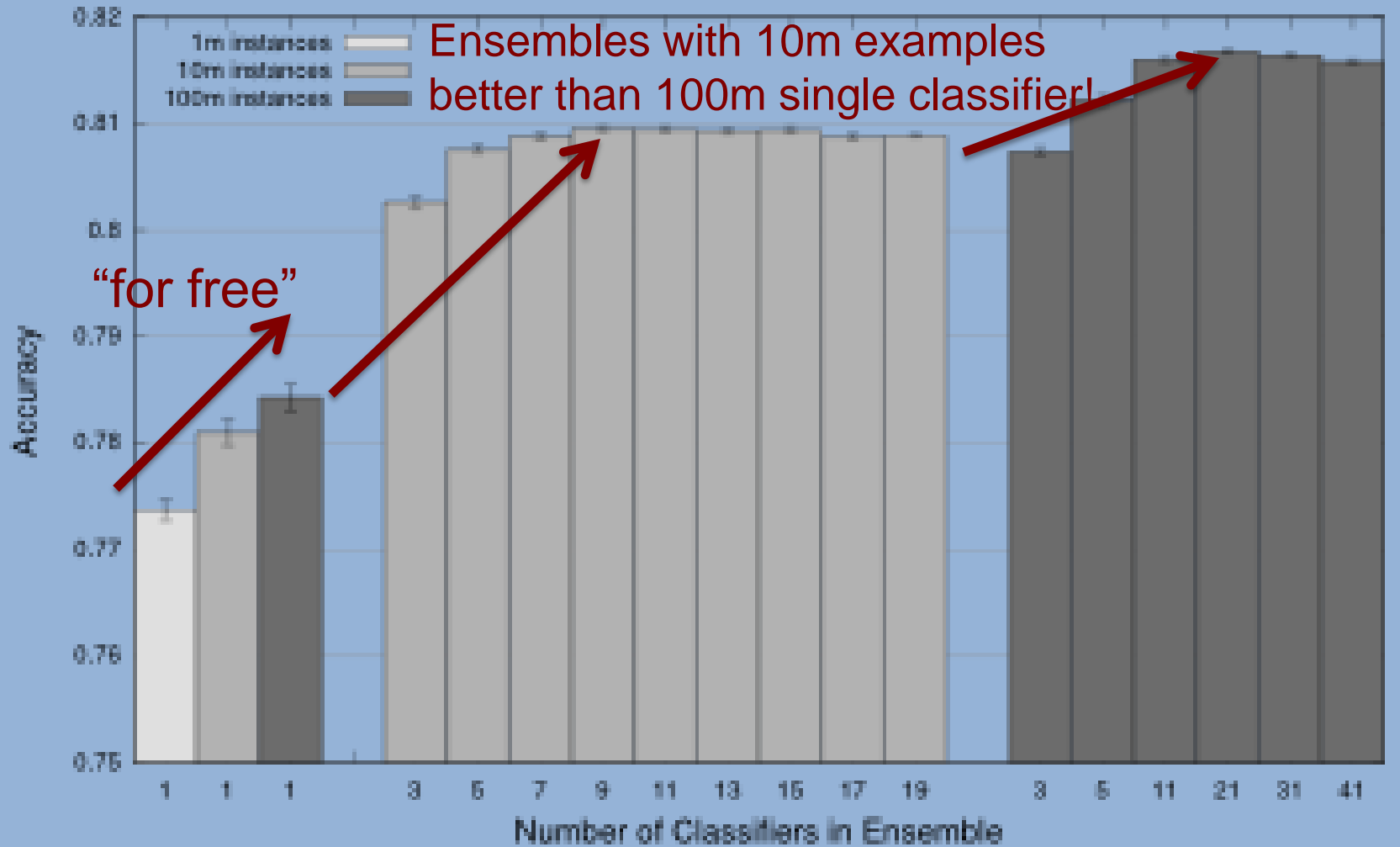
Further advantages of parallelism

- Although stream based learners look at all of the data, a number of them can be executed in parallel
 - Effective for tuning hyper-parameters
- Generative models such as Naïve Bayes are naturally well suited to distributed learning
 - Just counting

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...



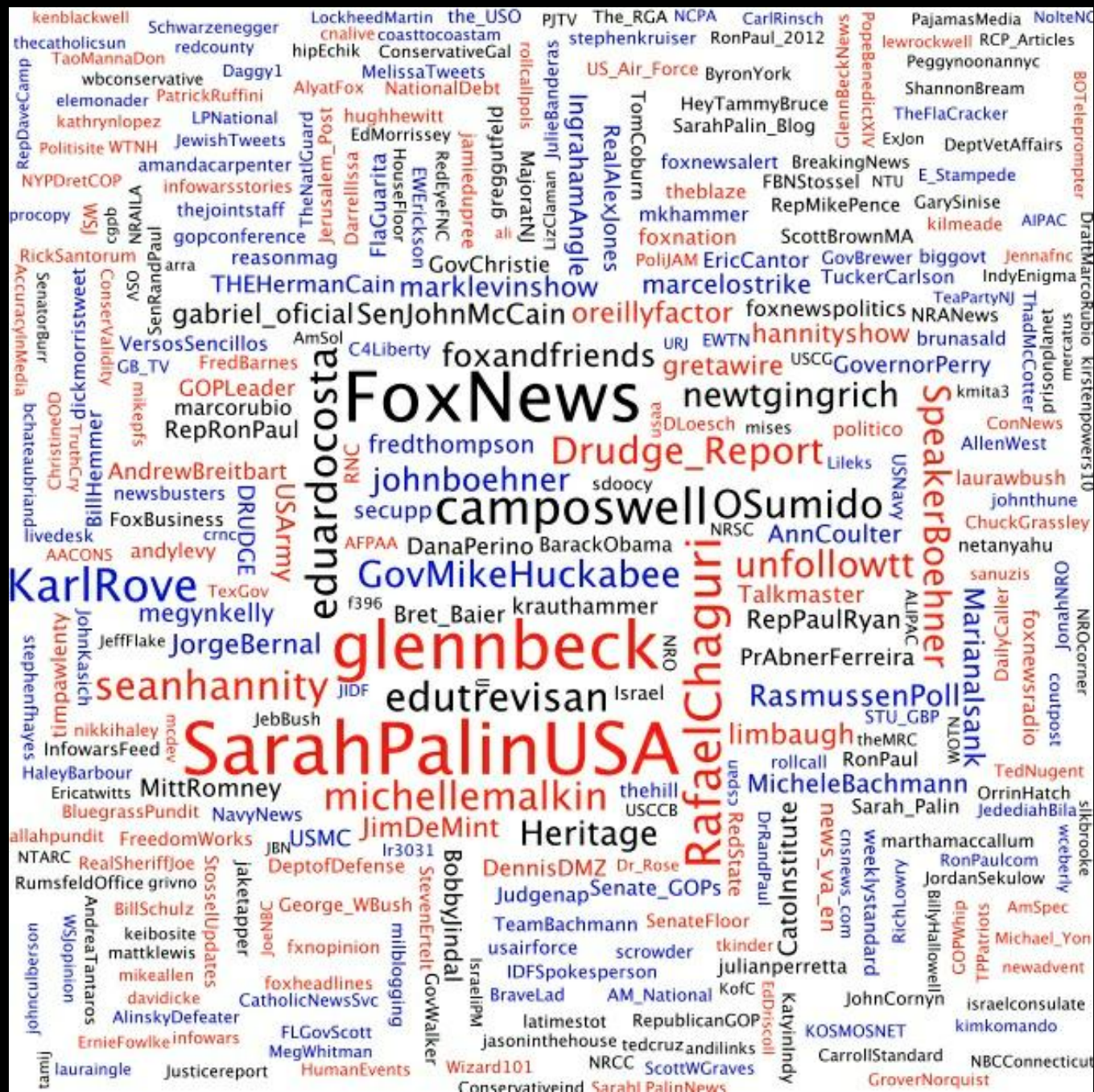
single classifier

10m ensembles

100m ensembles

Iterative algorithms

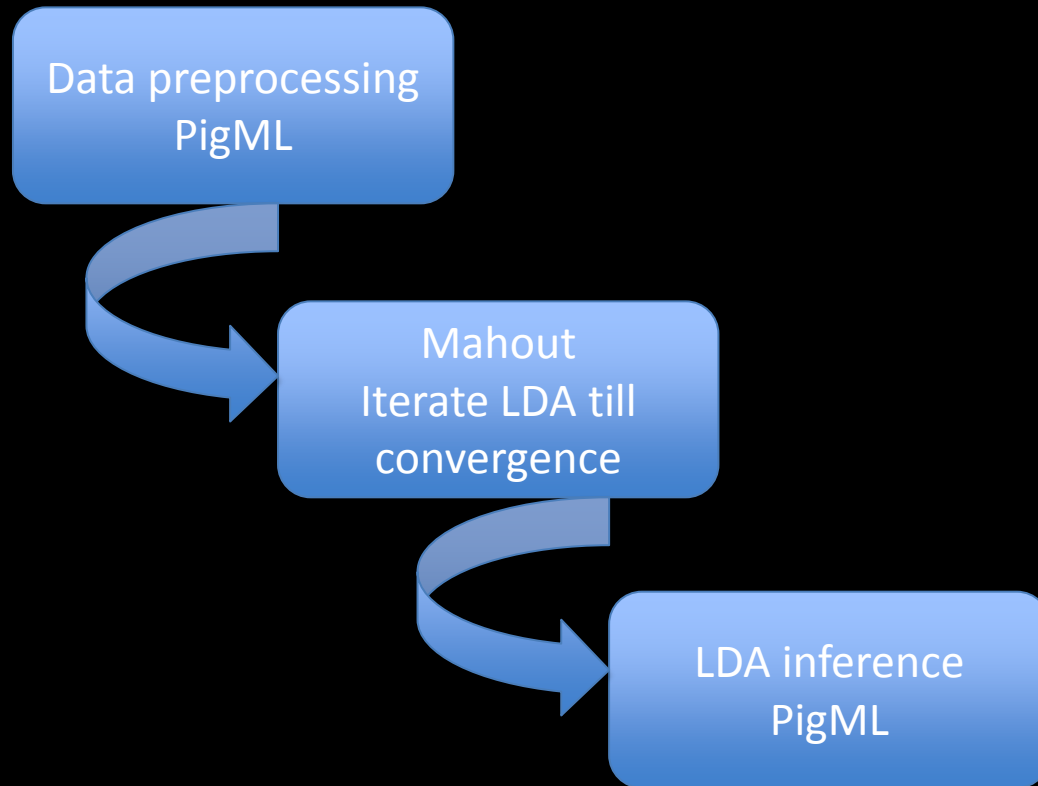
- What if one **REALLY** wants to iterate over big data till convergence?
 - Unrolling of small loops in Pig
 - Implementing the algorithm in Cascading with Scala or Python
 - Using Python Control Flow in Pig
 - Offload to custom MR job (e.g., Mahout)



Example: modeling topic distribution

- Latent Dirichlet Allocation (LDA)
 - Bayesian soft-clustering technique
 - Each document can naturally belong to multiple clusters with degree-of-membership
 - Popular in topics modeling for text data
- Why should we care?
 - Topics capture user interests
 - Knowledge of interests makes recommendation easier

Mahout/PigML integration



LDA applications

- User interest modeling facilitates user matching as well as predicting engagement.
- It does not directly solve the more general problem of finding users who might be interested in a tweet, story, web-page, etc
- LDA offers a probabilistic model of user interests (e.g., based on what they tweet about).

Anchoring LDA

- Runs of LDA can lead to significantly different results from one run to the next
- This may hamper the usefulness/interpretability of the results
- Anchoring clusters based on the known tag labels regularizes the clustering procedure
 - Also known as Labeled LDA (LLDA)

ML/Data Mining we contribute to

- Cassowary
 - Large graph mining in Java (single box)
- Pig
- Scalding
 - cascading with Scala
- PyCascading
 - Cascading with Python
- Mahout

ML outside of Twitter

- What does the academic community consider important?
 - Sentiment detection
 - Event tracking (e.g., epidemics)
 - Politics
 - Locality detection
 - Spam detection

Publications mentioning ...

- Google Scholar (late 2011)

Query	Result count
Twitter spam	9,800
Web spam	52,300
Email spam	41,600
IM spam	30,000
SMS spam	10,500

Quick search

- Number of titles on Amazon (early 2011)

Query	Result count
Twitter marketing	544
Twitter profits	100
Twitter money	236
Twitter rich	36
Twitter seo	21

Spam/spammer modeling

- Why spam Twitter (or other social media)?
- The types of spam seen
 - @replies/@mentions
 - Trend/search spam
 - Follow spam
- Long term vs. short term
 - Spammers favor “pump and dump”
 - Swift reaction to attacks is key

Example: normal interactions



Example: spammy interactions



ML @ Twitter

- Kumman Chellapila
- Jimmy Lin
- Yue Lu
- Jake Mannix
- Gilad Mishne
- Ram Ravichandran
- Miguel Rios
- Andy Schlaikjer
- Yifan Shi
- +others



THANK YOU

