# Big Data Mining
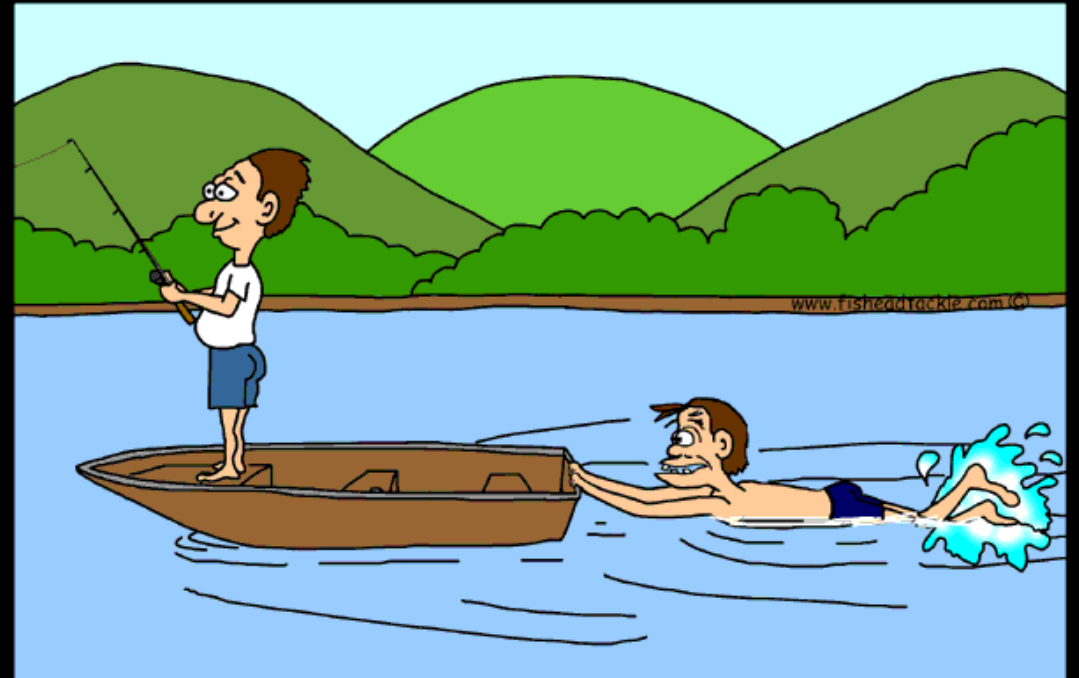


Zhifei Zhang

# Outline

- **Big data vs. Small data**

- **Software for big data**

- **An example**

# Small Data vs. Big Data

# Big Data

## Attributes
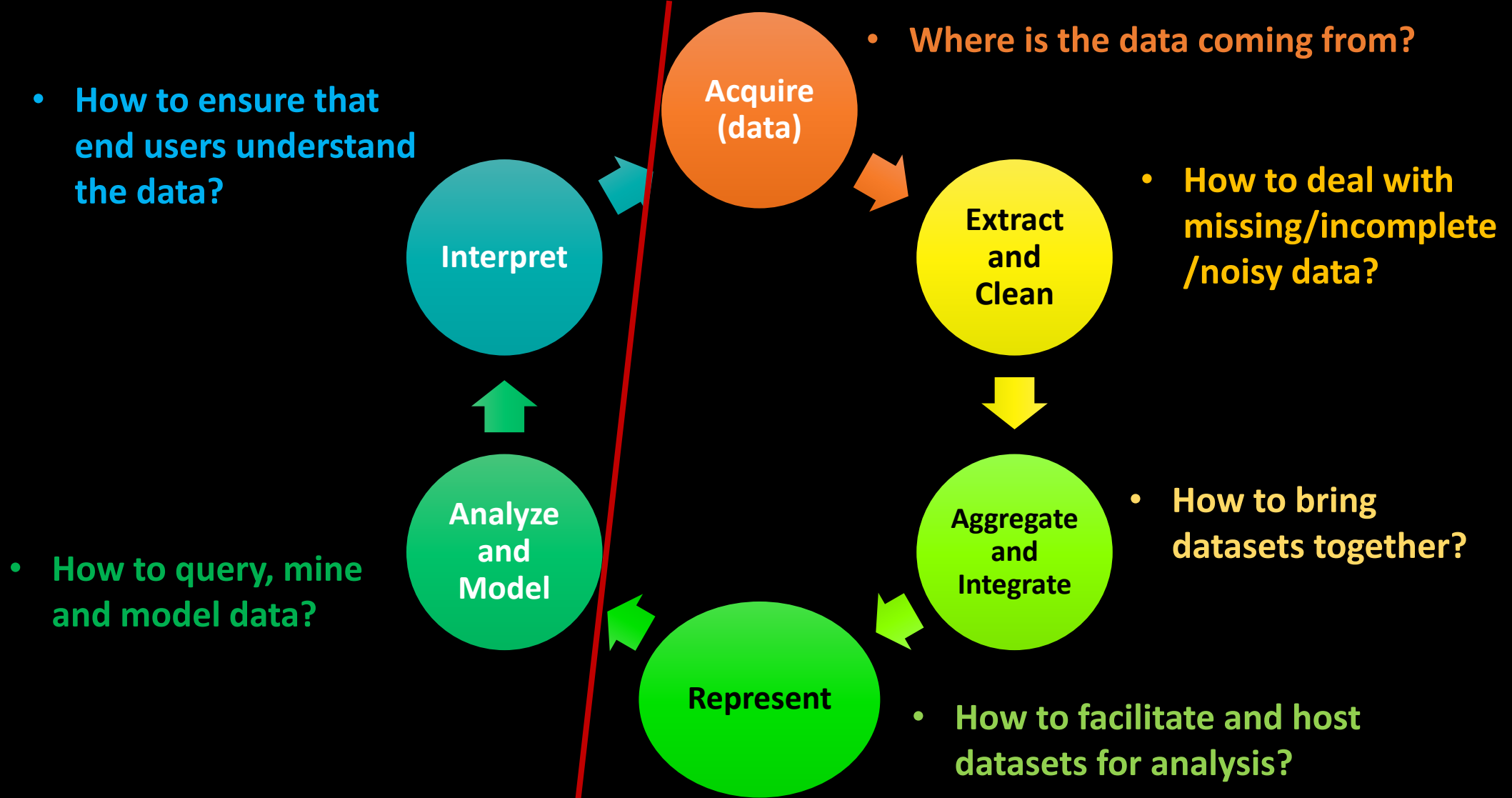
- Volume
- Variety
- Velocity
- Variability
- Veracity

## Types of Datasets

- High dimensional
- Sparse
- Graph
- Infinite/streaming
- Labeled

## Compute Models

- Map Reduce
- Streams / Online
- Single machine in-memory

# Big Data

- **Where is the data coming from?**

**Acquire (data)**

- **How to ensure that end users understand the data?**

**Interpret**

**Extract and Clean**

- **How to deal with missing/incomplete /noisy data?**

- **How to query, mine and model data?**

**Analyze and Model**

**Aggregate and Integrate**

- **How to bring datasets together?**

**Represent**

- **How to facilitate and host datasets for analysis?**

# Algorithms for Big Data

- **Classification: Bayes, clustering.**

- **Regression: polynomial, SVM.**

- **Dim. Reduction: PCA, SVD.**

# Example 1: Naïve Bayes

$$P(y|x) \propto p(x|y)P(y)$$

$$P(x|y) = \frac{C(X = x, Y = y)}{C(Y = y)}$$

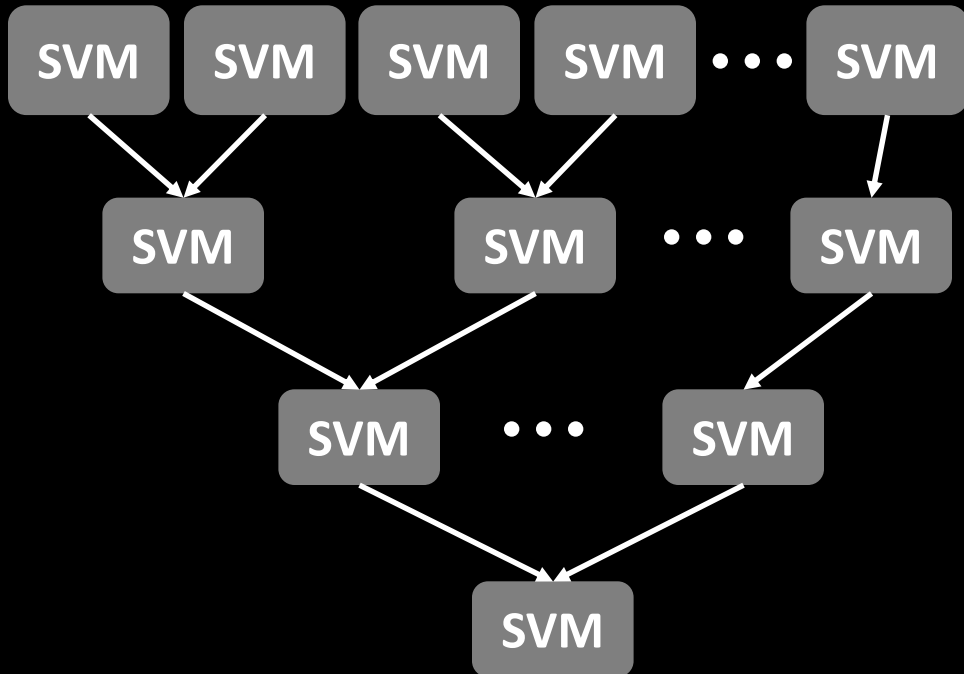$$P(y) = \frac{C(Y = y)}{C(Y = any)}$$

**Given a data point** $(x, y)$:

$$C(Y = ang) + +$$

$$C(Y = y) + +$$

$$C(Y = y, X = x) + +$$

# Example 2: SVM

$$\min_{\alpha} \left( \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) - \sum_i \alpha_i \right)$$

## Parallel-Hierarchical SVM



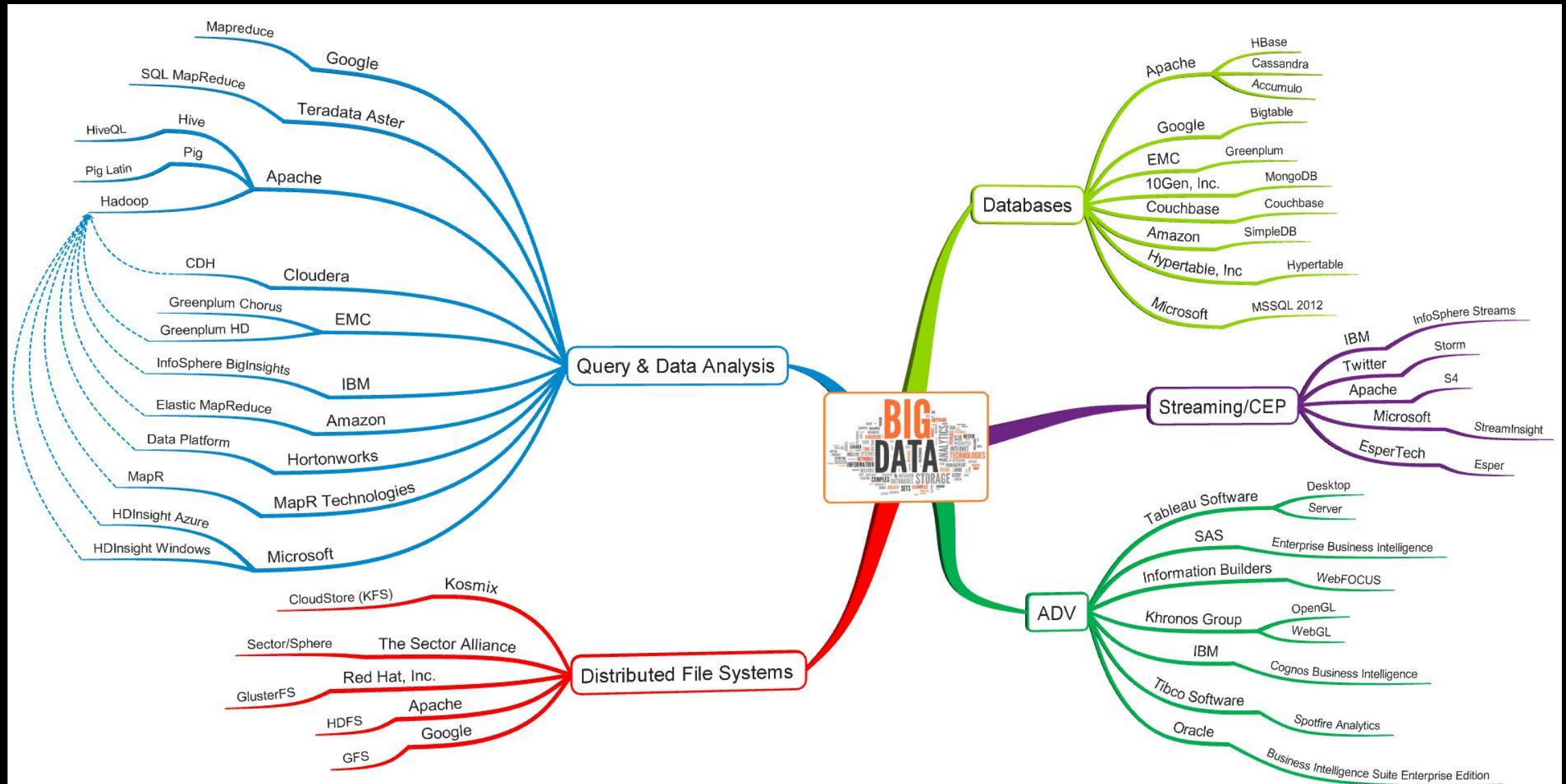## Given a partial data set:

**Implement SVM**

**Save the local model**

**Combine local models**

# Algorithms for Big Data

- **On-line version**
- **High efficiency**
- **Less iteration**
- **Approximation**

# Software for Big Data

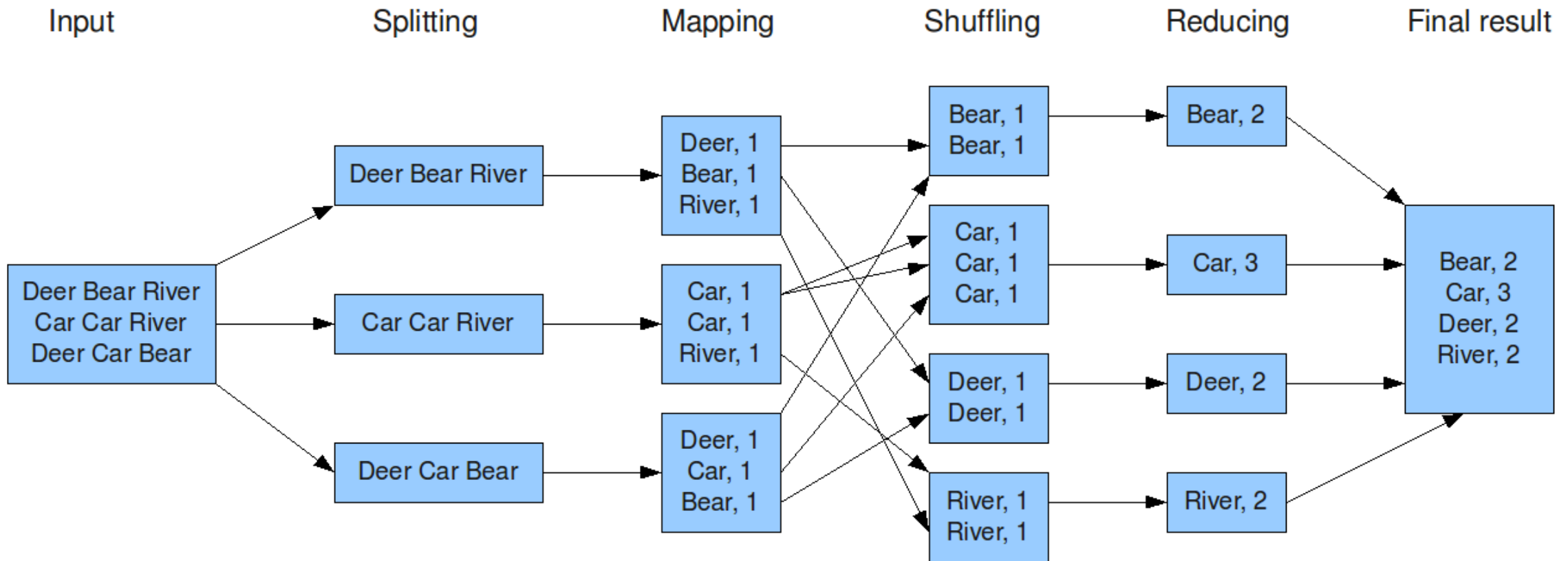# Software for Big Data

- **Common Software**
- **Machine Learning Lib**
- **Faster Framework**

# Map & Reduce



The overall MapReduce word count process

# Tracking Drinking Behavior from Twitter Data
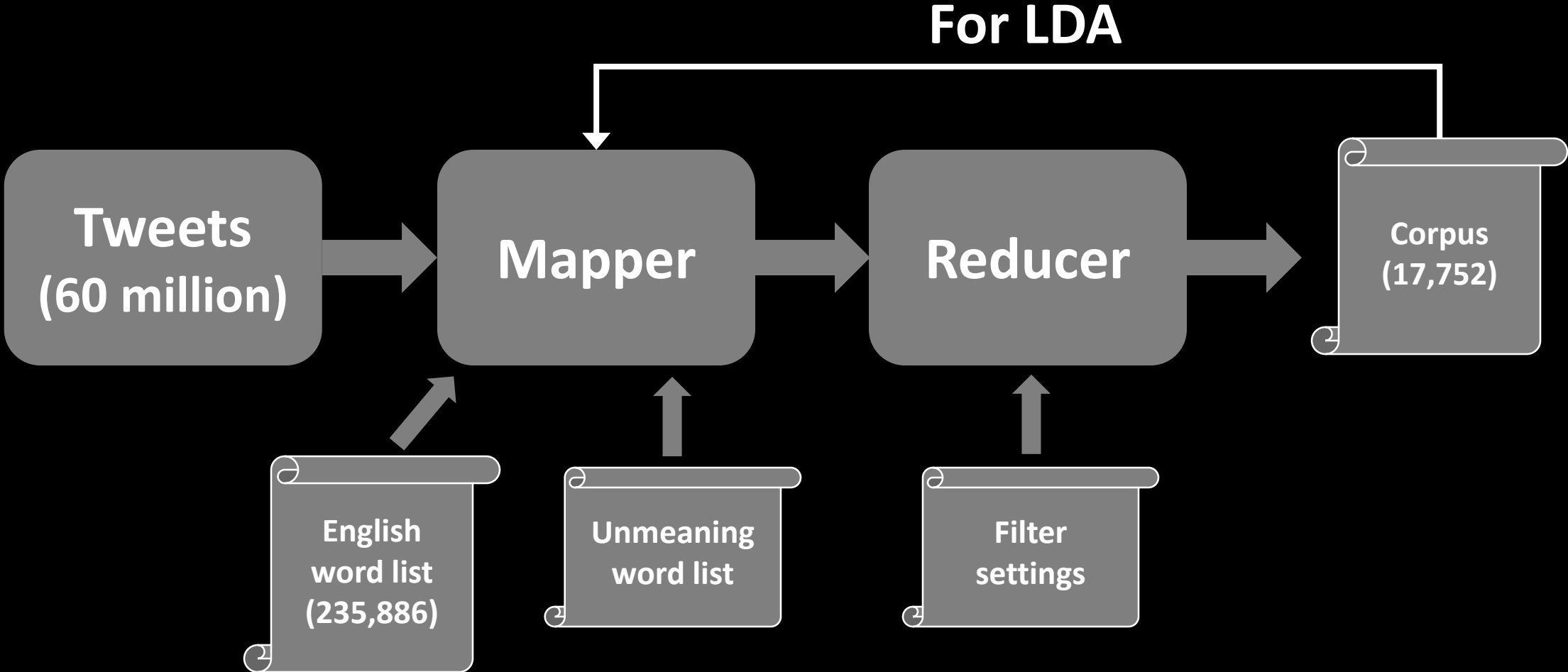
# Twitter Data

3/27/2014
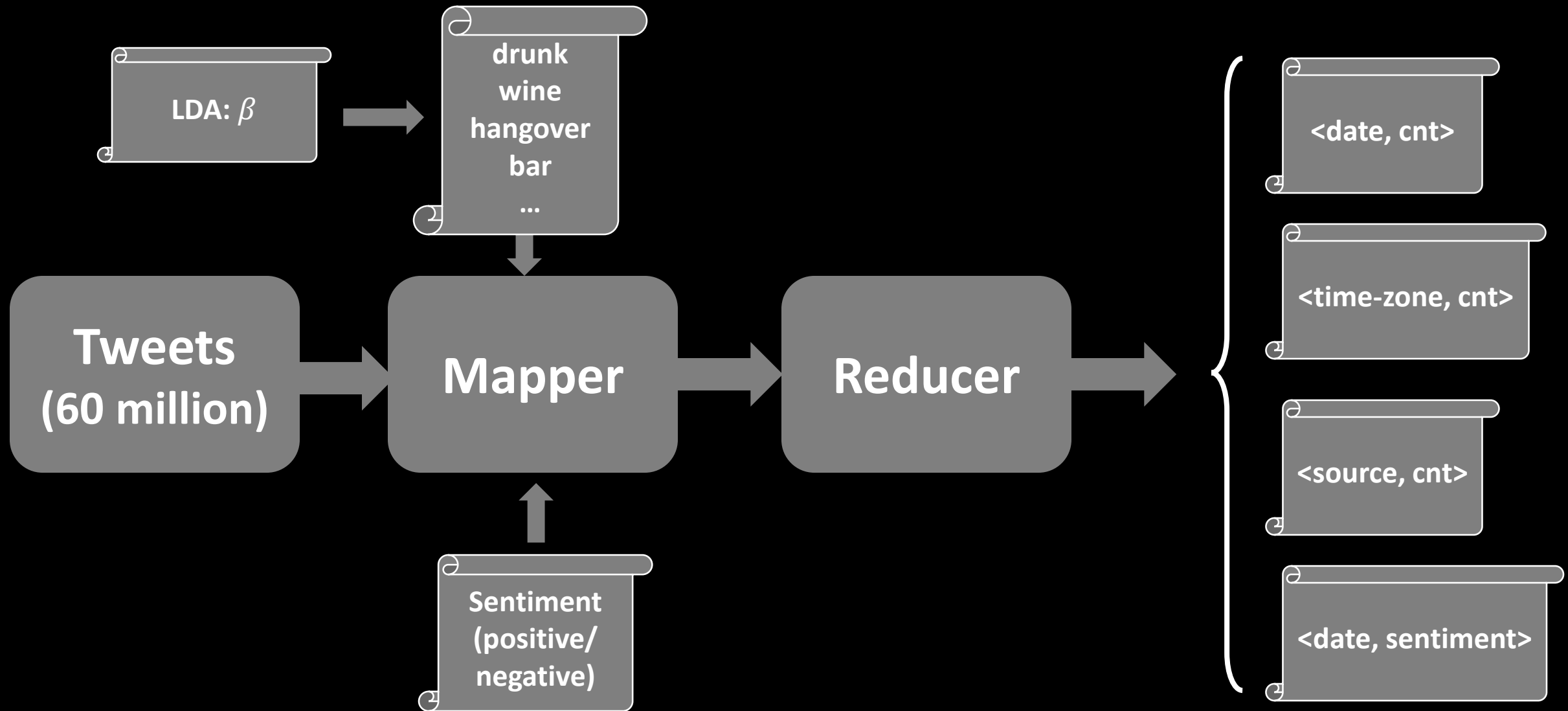
60, 000, 000
Tweets

5/1/2014

Drinking
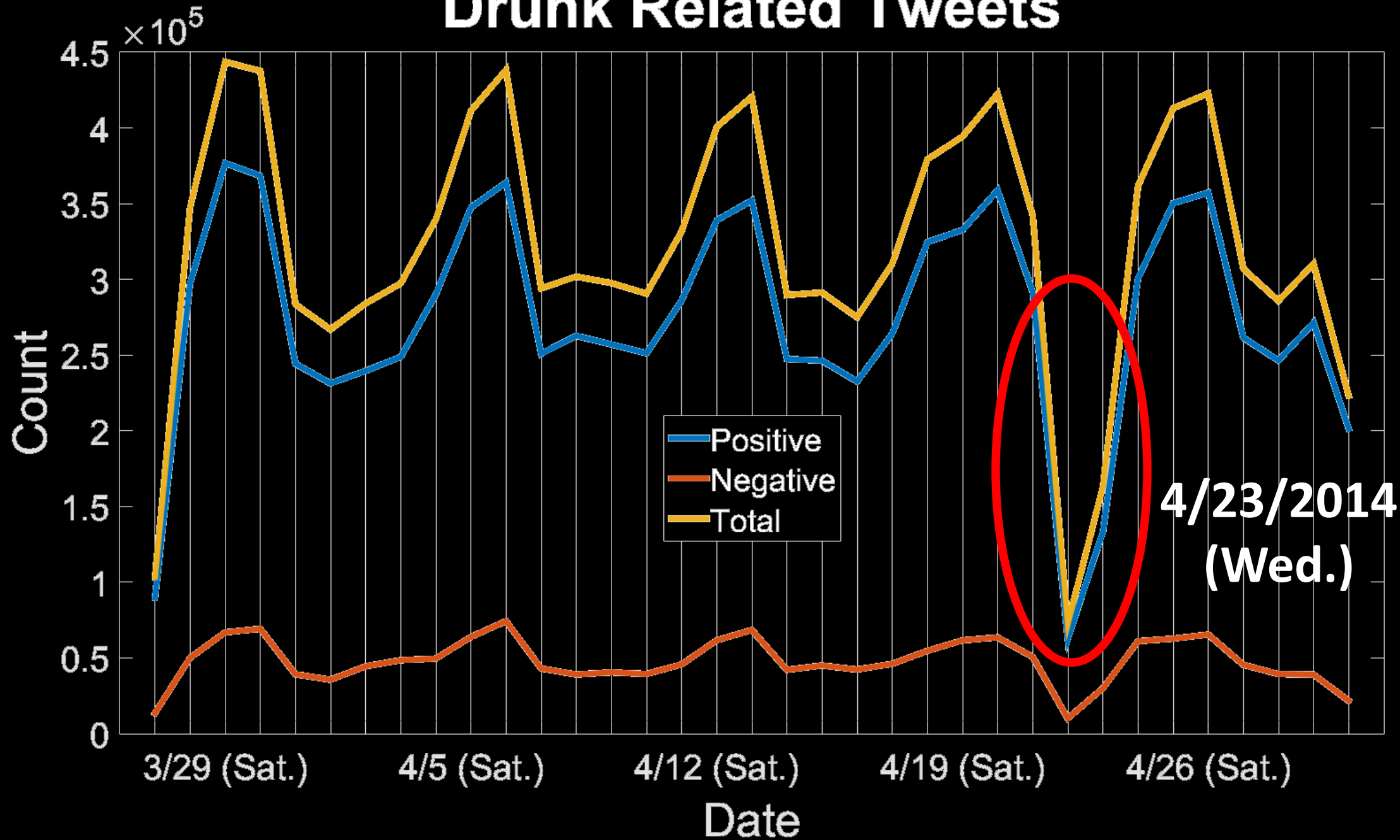
11,255,207
Tweets

# LDA + Hadoop → Drinking Words

For LDA

**Tweets (60 million)** → **Mapper** → **Reducer** → Corpus (17,752)

English word list (235,886)

Unmeaning word list

Filter settings

# Hadoop → Drinking Tweets

LDA: $\beta$

drunk
wine
hangover
bar
...

Tweets
(60 million)

Mapper

Sentiment
(positive/
negative)

Reducer

<date, cnt>

<time-zone, cnt>

<source, cnt>

<date, sentiment>

# What Happened around Apr. 23rd, 2014 ?




Unrest in Egypt


Car bomb in Kenya, killing 4 people


Sinking of the South Korean ferry killed 304 passengers


Unrest in Ukraine


6.6M earthquake in Canada

# The Truth is ...



| drwxr-xr-x | zzhang61 | hdusers | 0 B | 0 | 0 B | 2014-04-22 |
|---|---|---|---|---|---|---|
| drwxr-xr-x | zzhang61 | hdusers | 0 B | 0 | 0 B | 2014-04-24 |

## There is no data on 4/23

**Source of Drunk Related Tweets**

iPhone: 50%

Others: 31%

Windows: < 1%

Android: 18%

# Drinkers Love iPhone More?

2014 Worldwide Smartphone Market Share (IDC)

- 2.70%
- 1.10%
- 13.80%
- 82.30%

Legend: Android, iOS, Windows Phone, Other OS

Source of Drunk Related Tweets

- iPhone: 50%
- Others: 31%
- Windows: < 1%
- Android: 18%

# Drinkers Love iPhone More?

248 Results of "wine" in Google play

2985 Results of "wine" in App Store

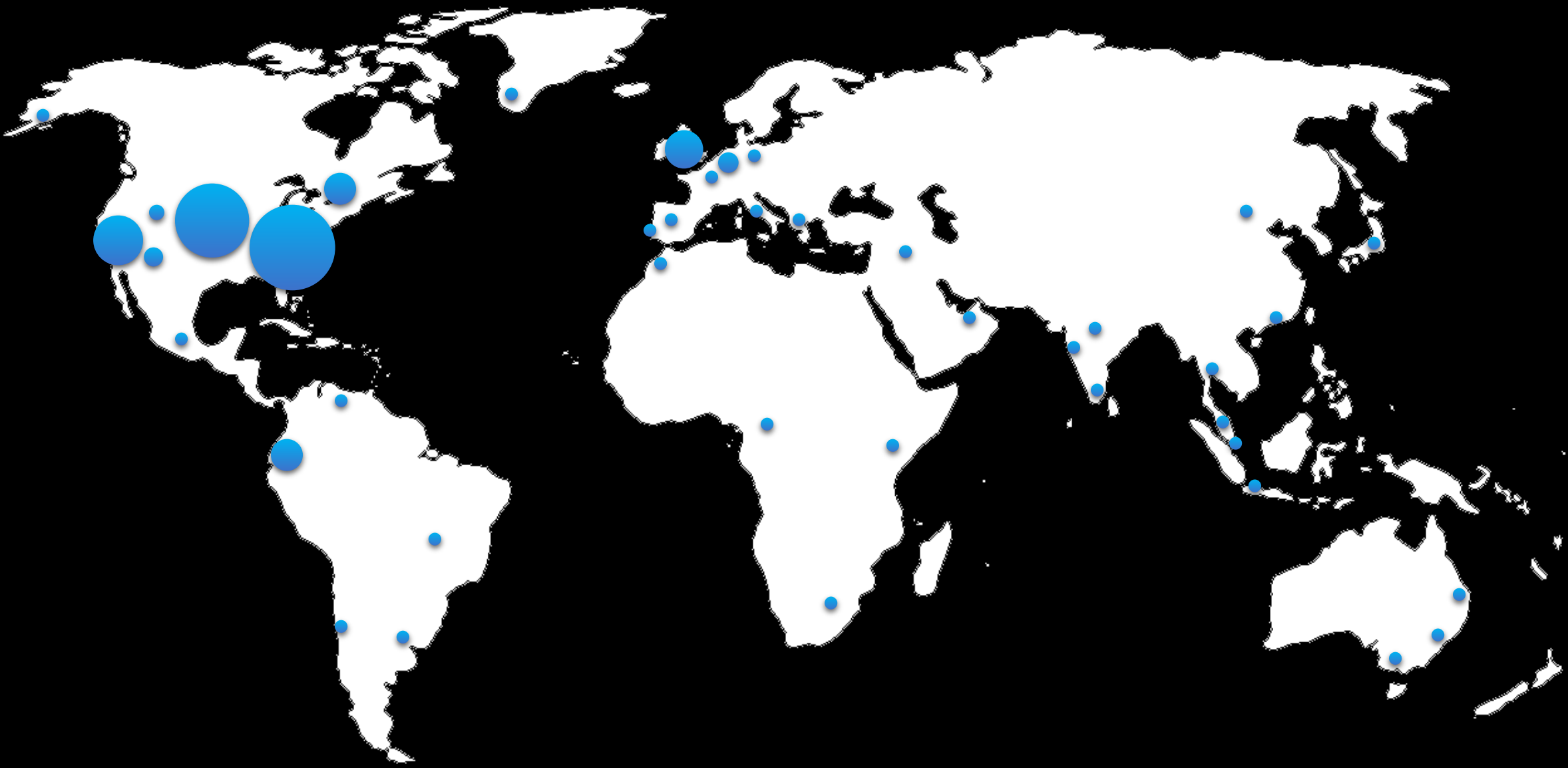# Location of Drinking Related Tweets

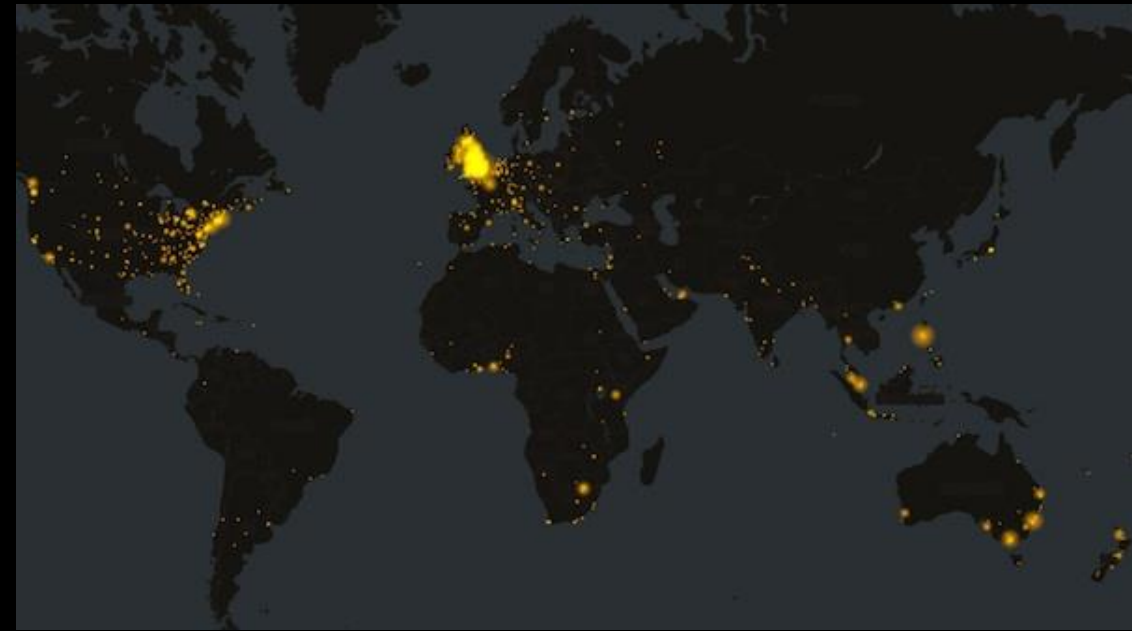## 46 zones have over 10,000 drunk related tweets

| | | | | | |
|---|---|---|---|---|---|
| 'Eastern/Time/(US/&/Canada)' | 1349324 | 'Casablanca' | 187121 | 'Dublin' | 38960 |
| 'Central/Time/(US/&/Canada)' | 1141305 | 'Alaska' | 129952 | 'Singapore' | 38276 |
| 'Pacific/Time/(US/&/Canada)' | 783803 | 'Athens' | 102358 | 'Santiago' | 37012 |
| 'London' | 583376 | 'Brasilia' | 92126 | 'Sydney' | 31784 |
| 'Atlantic/Time/(Canada)' | 497300 | 'Beijing' | 74283 | 'Madrid' | 31394 |
| 'Quito' | 490727 | 'Greenland' | 52823 | 'Kuala/Lumpur' | 29112 |
| 'Amsterdam' | 325605 | 'Chennai' | 51933 | 'Paris' | 27678 |
| 'Arizona' | 291043 | 'Bangkok' | 49231 | 'Pretoria' | 26745 |
| 'Mountain/Time/(US/&/Canada)' | 234895 | 'Edinburgh' | 44594 | 'Mexico/City' | 26023 |
| 'Hawaii' | 198108 | 'Buenos/Aires' | 43854 | 'Caracas' | 25682 |

**Zones with over 10,000 Drinking Related Tweets**

# American Drink More?



Twitter user distribution

Drinking Tweet distribution

# Pros & Cons

**Pros:**
- **Through big data mining, we may find something hard to be recognized in daily life.**
- **Identify effect of certain event on the public.**

**Cons:**
- **How to interpret the result? (Misleading)**
- **Only reflect pubic behavior but not for individuals.**