# Strata+Hadoop WORLD

PRESENTED BY
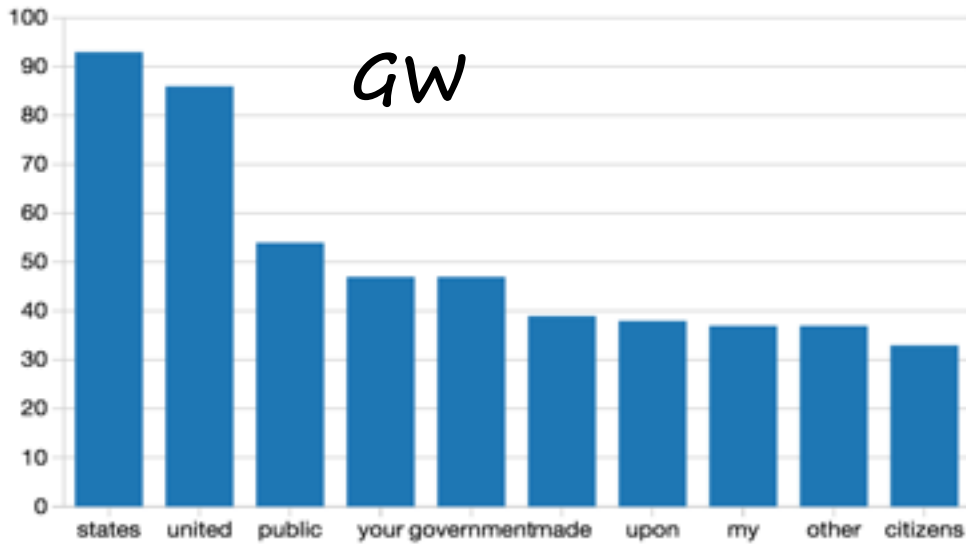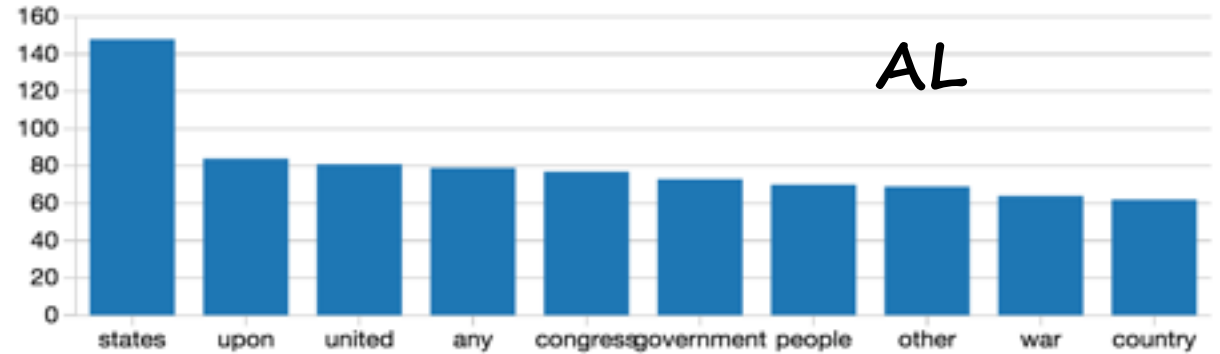
O'REILLY® cloudera®

# Text Analytics

*In which we analyze the mood of the nation & then the 2016 US Presidential primary debates*
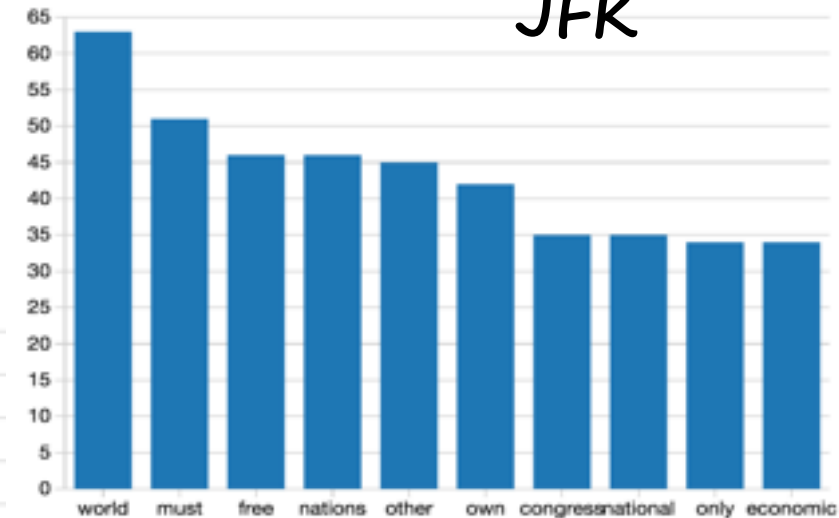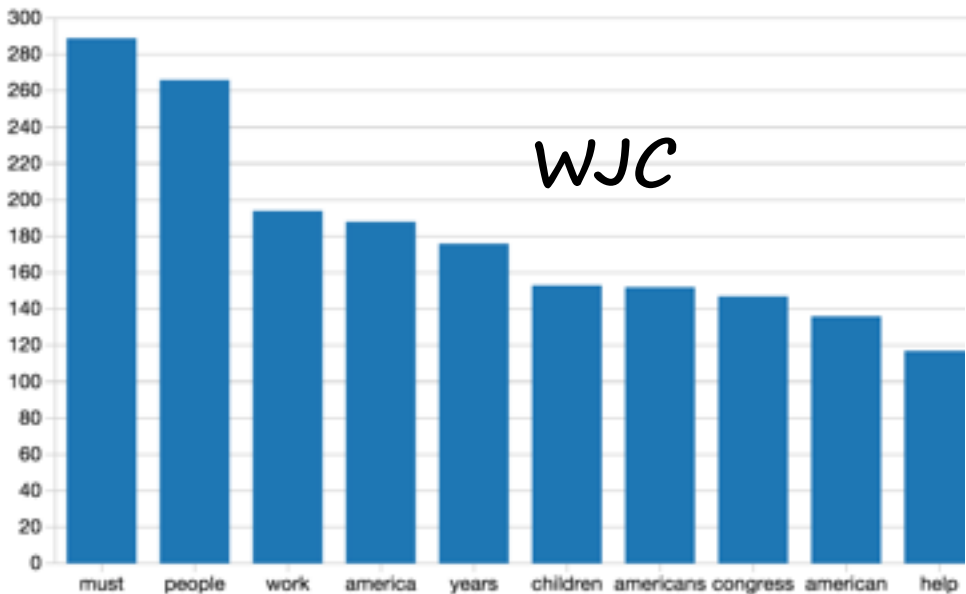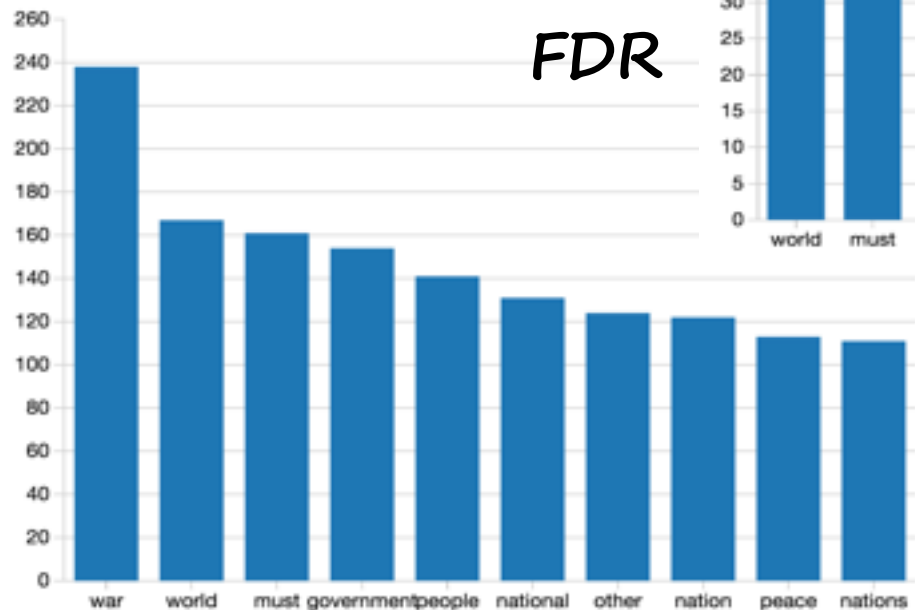
*(Don't worry, we will keep it civilized)*

strataconf.com

#StrataHadoop

Krishna Sankar
https://www.linkedin.com/in/ksankar
March 29, 2016

Mood Of The Nation
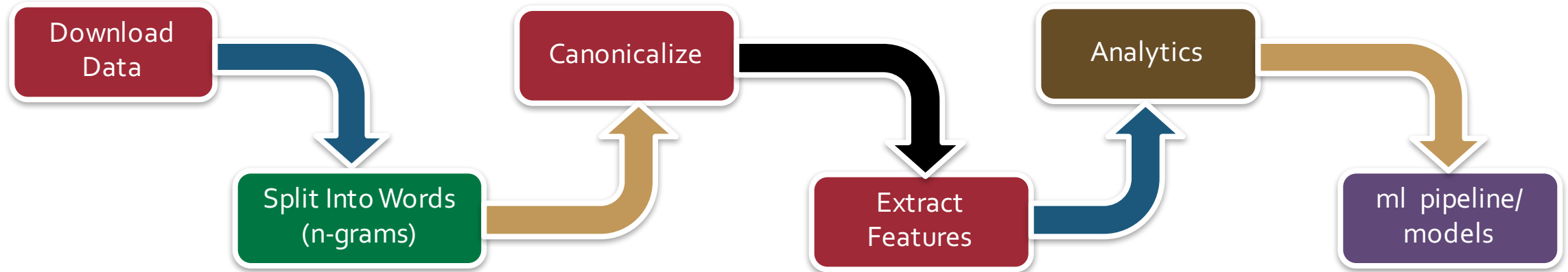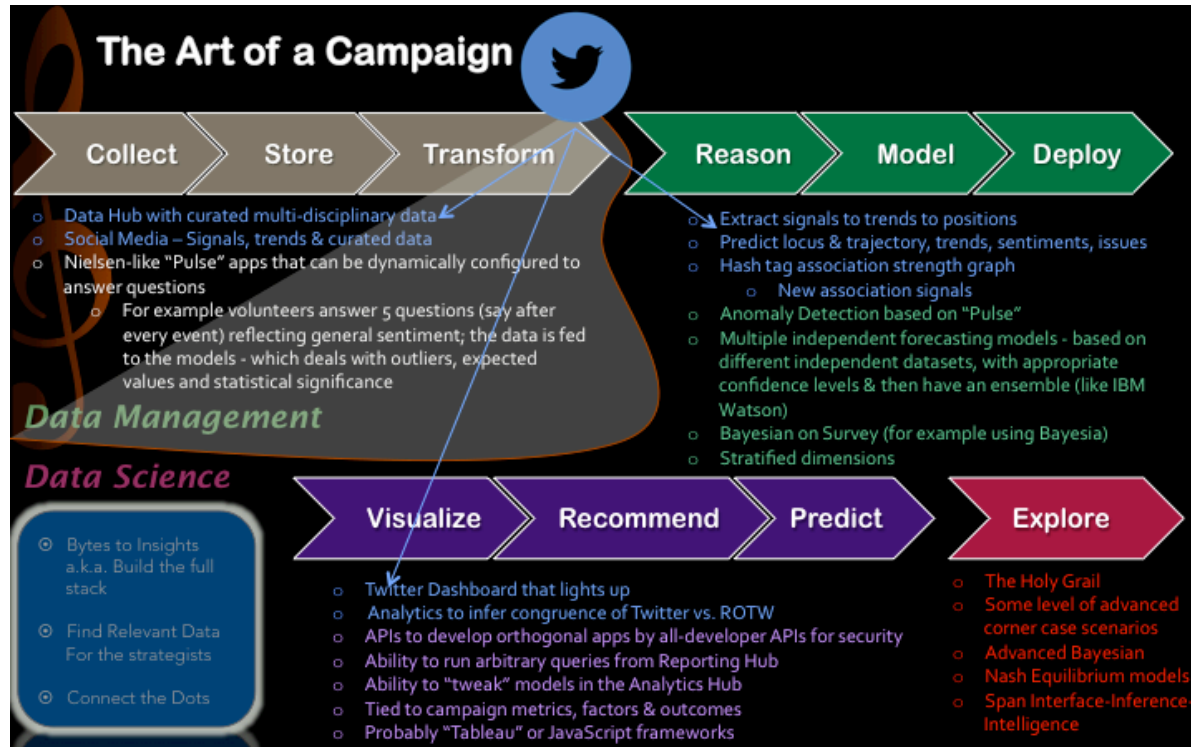
#StrataHadoop

Strata+Hadoop WORLD
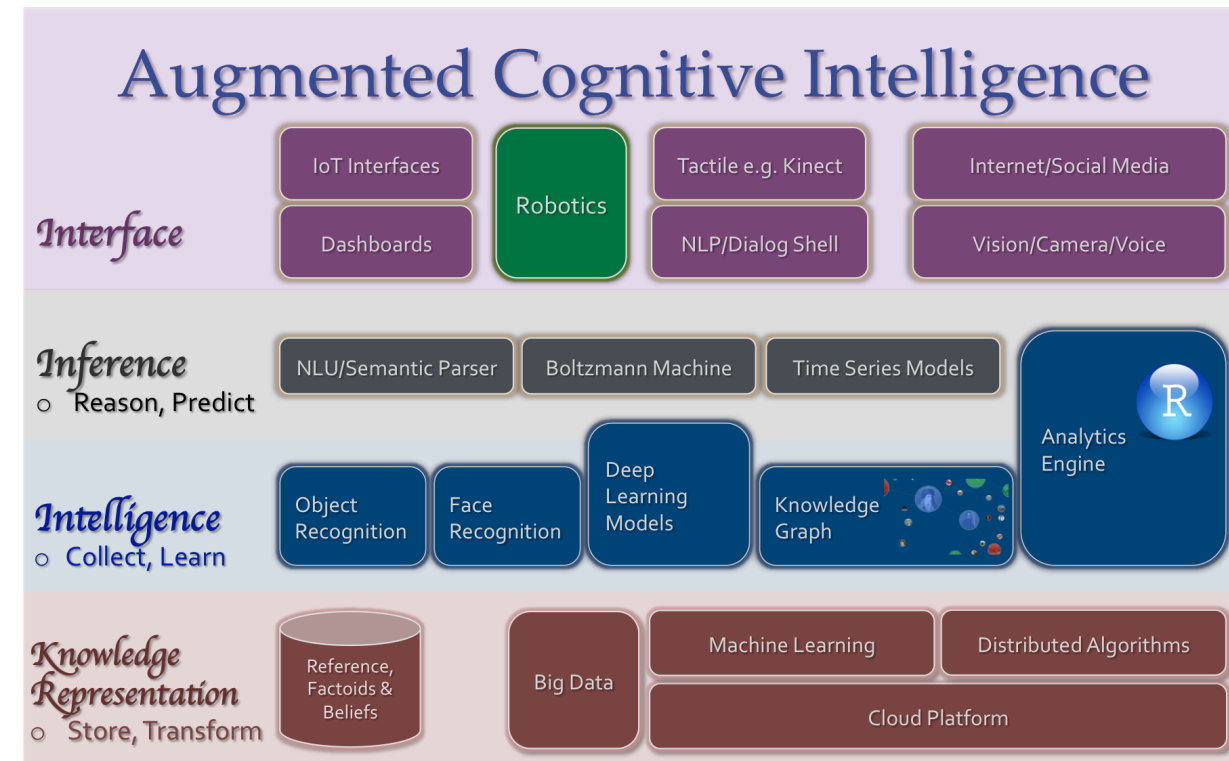
# Text Analytics Pipeline



| ○ http://stateoftheunion.o netwothree.net/texts/in dex.html | ○ Spark RDD Mechanisms | ○ Case transformation, special characters, space elimination,… | ○ Common word elimination | ○ Exploration | ○ Logistic Regression<br>○ Bayesian Models<br>○ LDA |
|---|---|---|---|---|---|
| ○ http://www.presidency.u csb.edu/debates.php | ○ Dataframes (especially for ml pipelines | ○ Stemming, Lemma | ○ TF-IDF | ○ Knowledge Representation (Knowledge Graph) | ○ Deep Learning<br>○ (Topics, Classification,…) |
| | | ○ Domain Scoping | | *word2vec* | |

▪ **Understand the tutorial in the context of a broader reference architecture (e.g. next 2 slides), Select components as needed by the app**
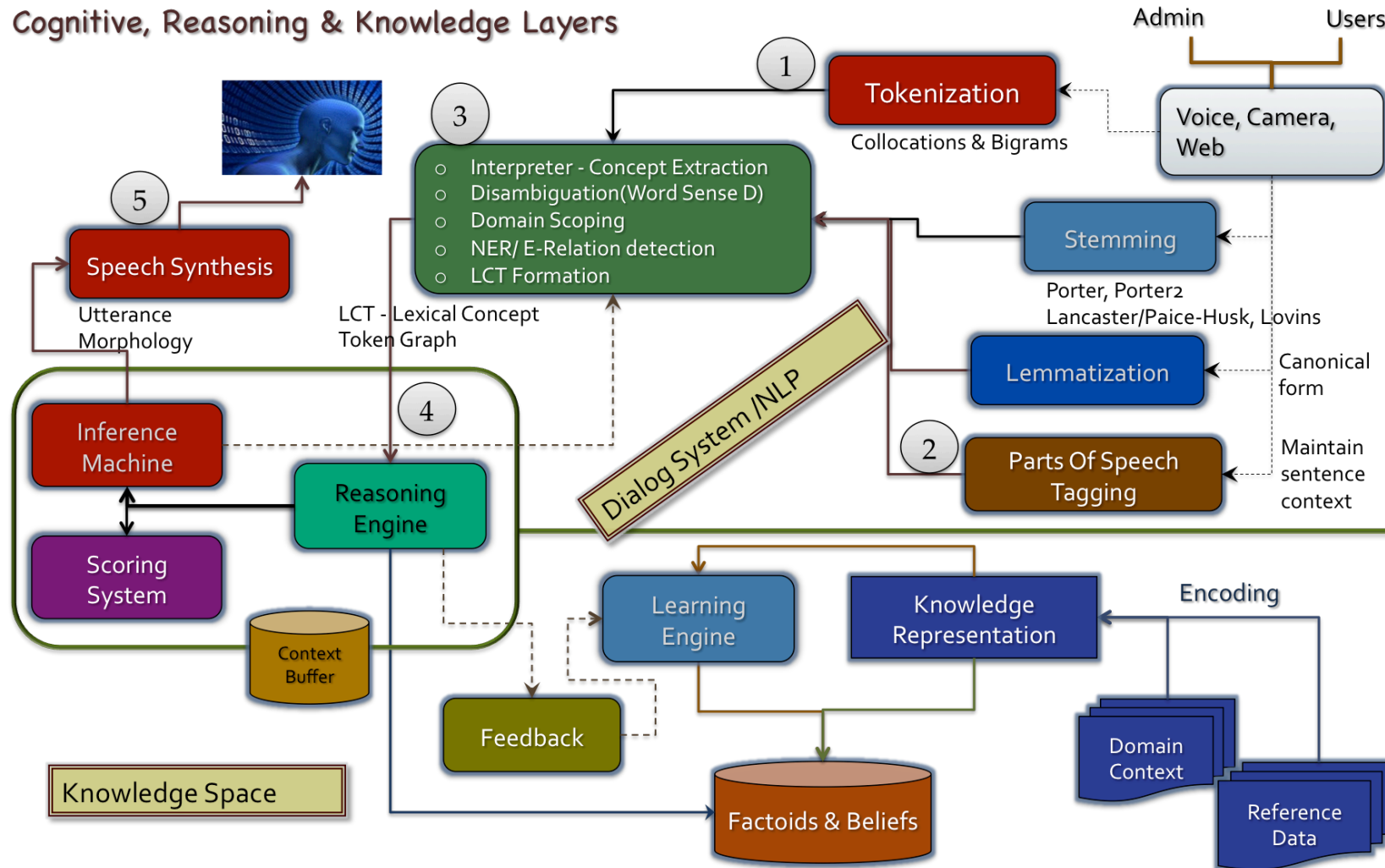
# Reference Architectures for Text Analytics



https://doubleclix.wordpress.com/2015/07/05/twitter-2-0-curated-signals-applied-intelligence-stratified-inference/

# Reference Architectures for Text Analytics

# Text Analytics

- A Data Scientist can use RDD primitives to do interesting work with text

  - Map-reduce in a couple of lines !

    - But it is not exactly the same as Hadoop Mapreduce (see the excellent blog by Sean Owen[1])

  - Set differences using substractByKey

  - Ability to sort a map by values (or any arbitrary function, for that matter)

  - Dataframe operations

- TF-IDF as Feature Extraction http://spark.apache.org/docs/latest/mllib-feature-extraction.html#tf-idf

- LDA for topic extraction/ml pipeline in next section

- *Good & Bad — mllib features span a continuum of libraries from rdds to dataframes to datasets to extraction to ml models to ml pipelines and beyond ...*

  http://blog.cloudera.com/blog/2014/09/how-to-translate-from-mapreduce-to-apache-spark/

Strata+Hadoop
WORLD

# Code Walkthru

- 03-Text Analytics Notebook