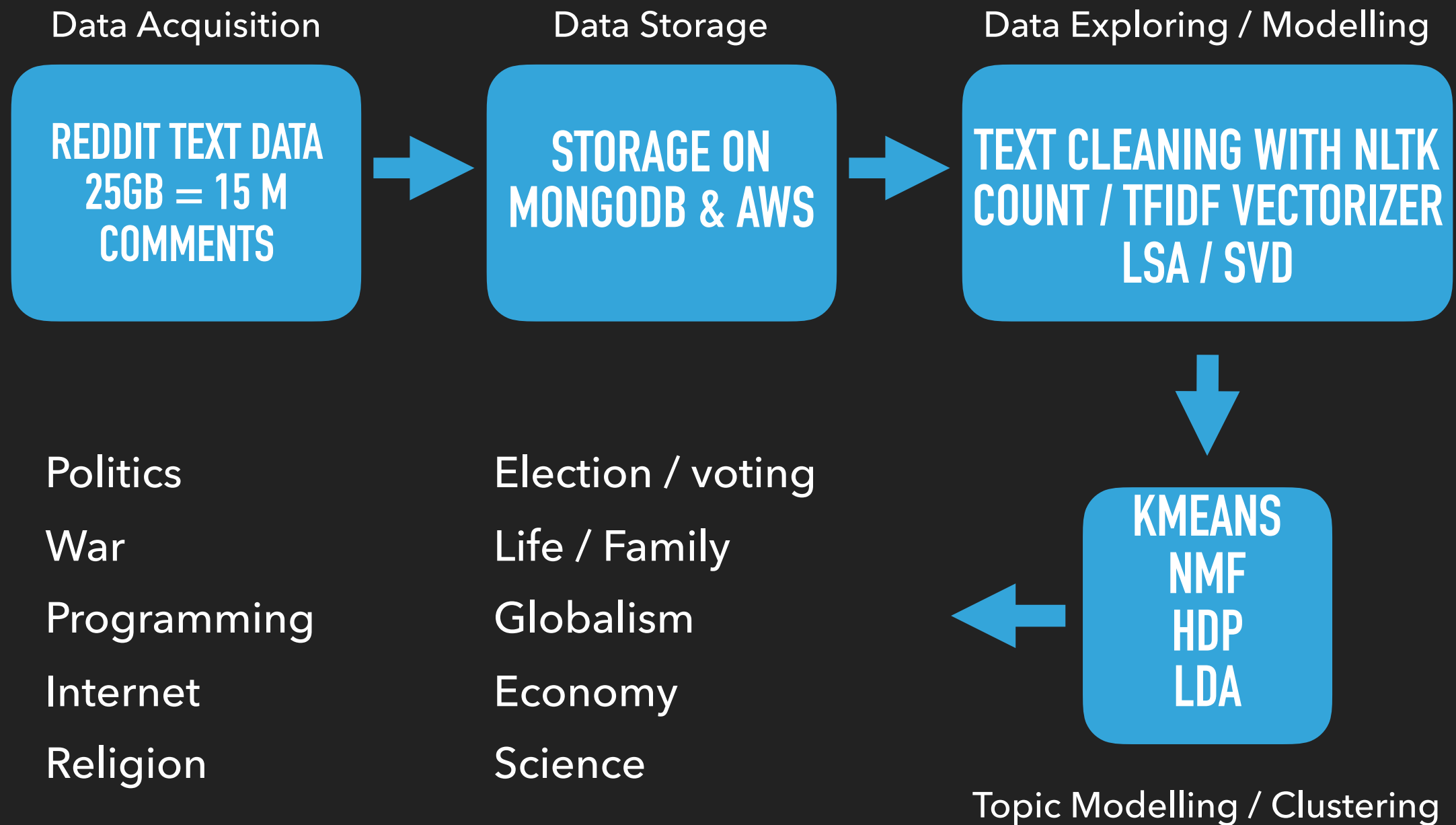Florian Philippi

06/02/2017

# REDDIT USER ANALYTICS

# NATURAL LANGUAGE PROCESSING

# OBJECTIVES

▸ Investigate reddit comments data on a user basis

▸ 300GB of text data available online

▸ Create a tool to summarize user activities using NLP

   ▸ Topics of interest

   ▸ Connection to users with similar interests

   ▸ Suggestions for new connections and topics
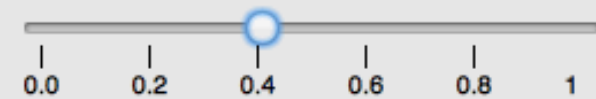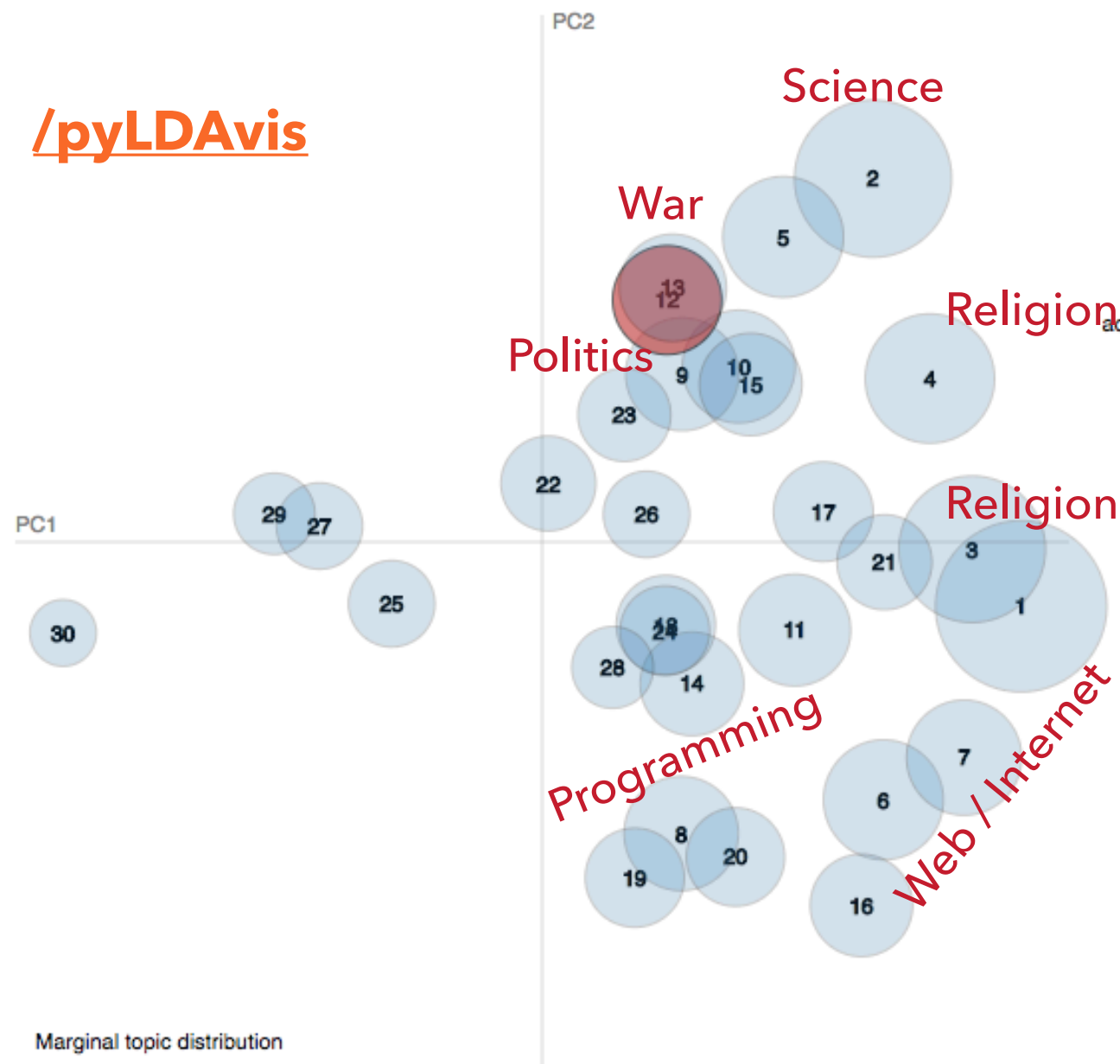
# PROCESS

Data Acquisition

**REDDIT TEXT DATA 25GB = 15 M COMMENTS**

Data Storage

**STORAGE ON MONGODB & AWS**

Data Exploring / Modelling

**TEXT CLEANING WITH NLTK COUNT / TFIDF VECTORIZER LSA / SVD**

Politics

War

Programming

Internet

Religion

Election / voting

Life / Family

Globalism

Economy

Science

**KMEANS NMF HDP LDA**

Topic Modelling / Clustering

# VISUALIZATIONS – D3

# VISUALIZATIONS – FLASK APP

# FUTURE WORKS

‣ Use more meta data as features

‣ Use all of the available data (submissions etc.)

‣ Word2vec

# SOURCES

https://files.pushshift.io/reddit/