

# **DJ BYG DATA**

Diego  
Trajche  
Tribhuvanesh

# Continuing...

- Goal: Is it possible to determine the genre of a song based solely on its lyrics?
- K-Means Clustering
- Scaling up
  - > 237,000 tracks
  - > 20000 artists
  - 5000 unique (stemmed) words
  - 5-20 clusters

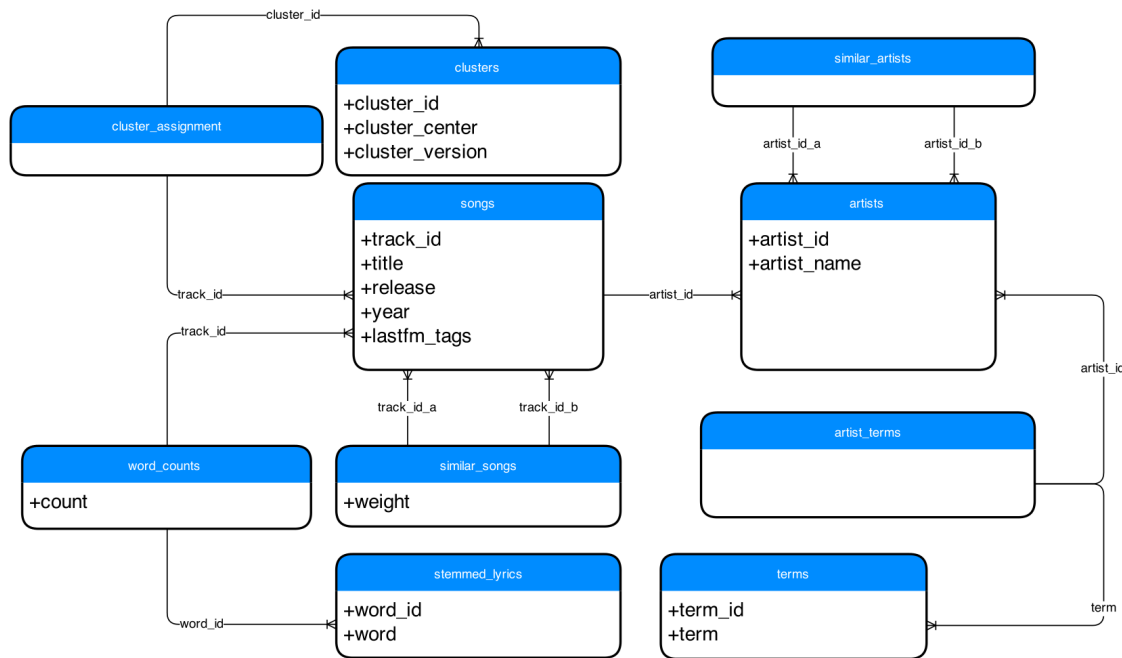
# The Big Data Challenge

- Problem 1: Preprocessing  
Process and clean 280GB of data
- Problem 2: Clustering  
280k songs to cluster, vectors of dimension 5000
- Problem 3: Visualization and Analysis  
9M pairs, 280k songs, 20.5k artists

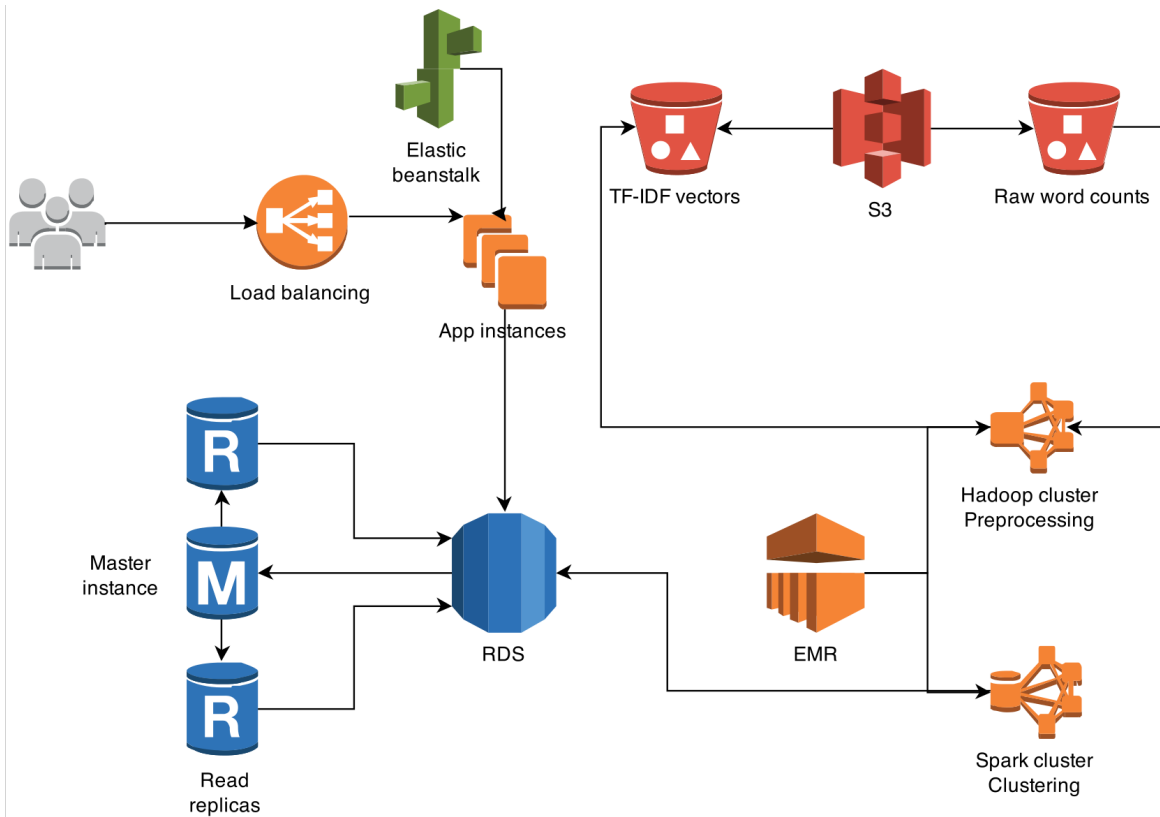
# The Big Data Challenge

- Problem 1: Preprocessing  
Solution: S3 and EMR Hadoop MapReduce
- Problem 2: Clustering  
Solution: Spark on EMR
- Problem 3: Visualization and Analysis  
Solution: Relational database, Replicated, Scalable web cluster

# Data Model - Relational



# Architecture



# Quality Metrics

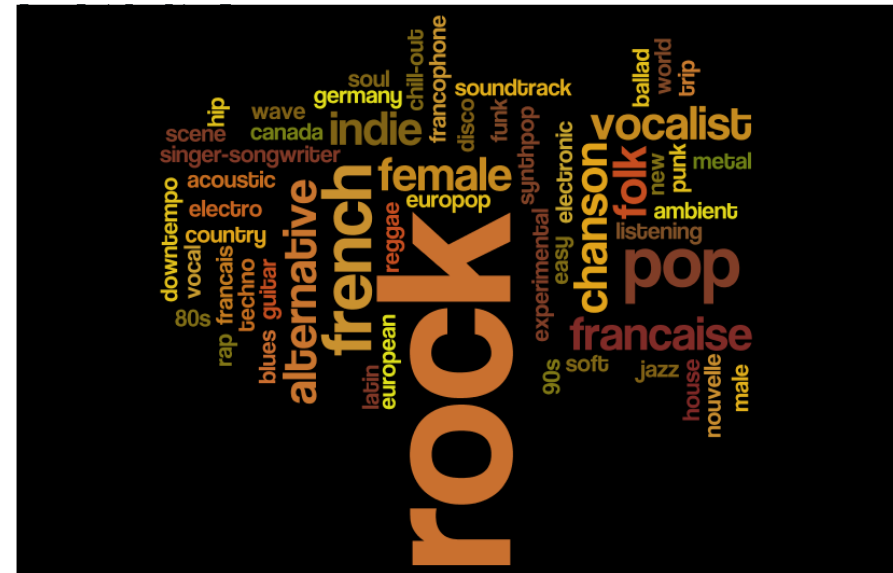
- Clustering parameters
  - Number of clusters, i.e.  $k$
  - Number of top words removed when calculating the clusters
- Using the similarity information from the MSD
  - Similar artists - Clusters should group similar artists together
  - Similar songs - Clusters should preserve similar songs graph
- Using the fact that artists are expected to be contained in a single cluster
- Using the distribution of songs among the clusters

# Results

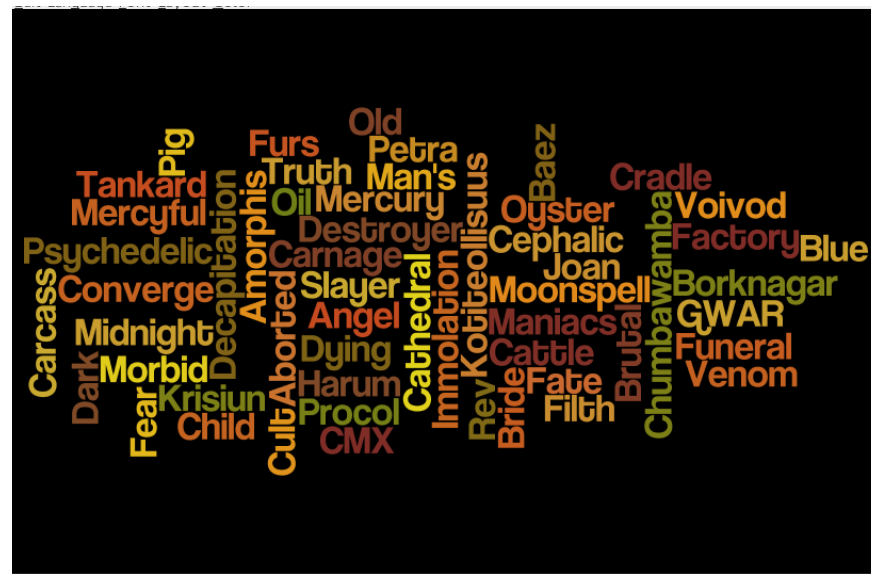




JulienDaho  
 Solaar  
 Paradis  
 Fabian  
 Brue  
 Vanessa  
 Attaque  
 Clerc  
 Hysteria  
 Les  
 Barback  
 St-Pier  
 Patrick  
 De  
 Benabar  
 Miossec  
 Ketanou  
 Obispo  
 Etienne  
 Lara  
 Goldman  
 Jean-Jacques  
 Pascal  
 Edith  
 Ogres  
 Rue  
 Alain  
 Mass  
 MC  
 Higelin  
 Breil  
 Suchon  
 Jacques  
 Natasha  
 Louise  
 Piaf  
 Bielay



A word cloud of music genres and related terms. The most prominent words are "metal" and "rock" in large, bold, orange and green fonts respectively. Other visible words include "progressive", "indie", "folk", "alternative", "classic", "pop", "vocalist", "guitar", "ambient", "states", "wave", "thrash", "80s", "black", "blues", "vocal", "united", "industrial", "experimental", "ballad", "synthpop", "funk", "male", "house", "heavy", "soul", "electro", "instrumental", "death", "jazz", "hip", "techno", "grunge", "new", "hardcore", "soundtrack", "soft", "acoustic", "american", "psychedelic", "punk", "downtempo", "electronica", "emo", "country", "doom", "female", "electronic", "hard", and "germany".



# Demo