

Word!
Yo.



LITERARI.LY

A LITERARY ASSISTANT
-ADAM LENART, ANKITH GUNAPAL & SANDIP PALIT

PROBLEM STATEMENT

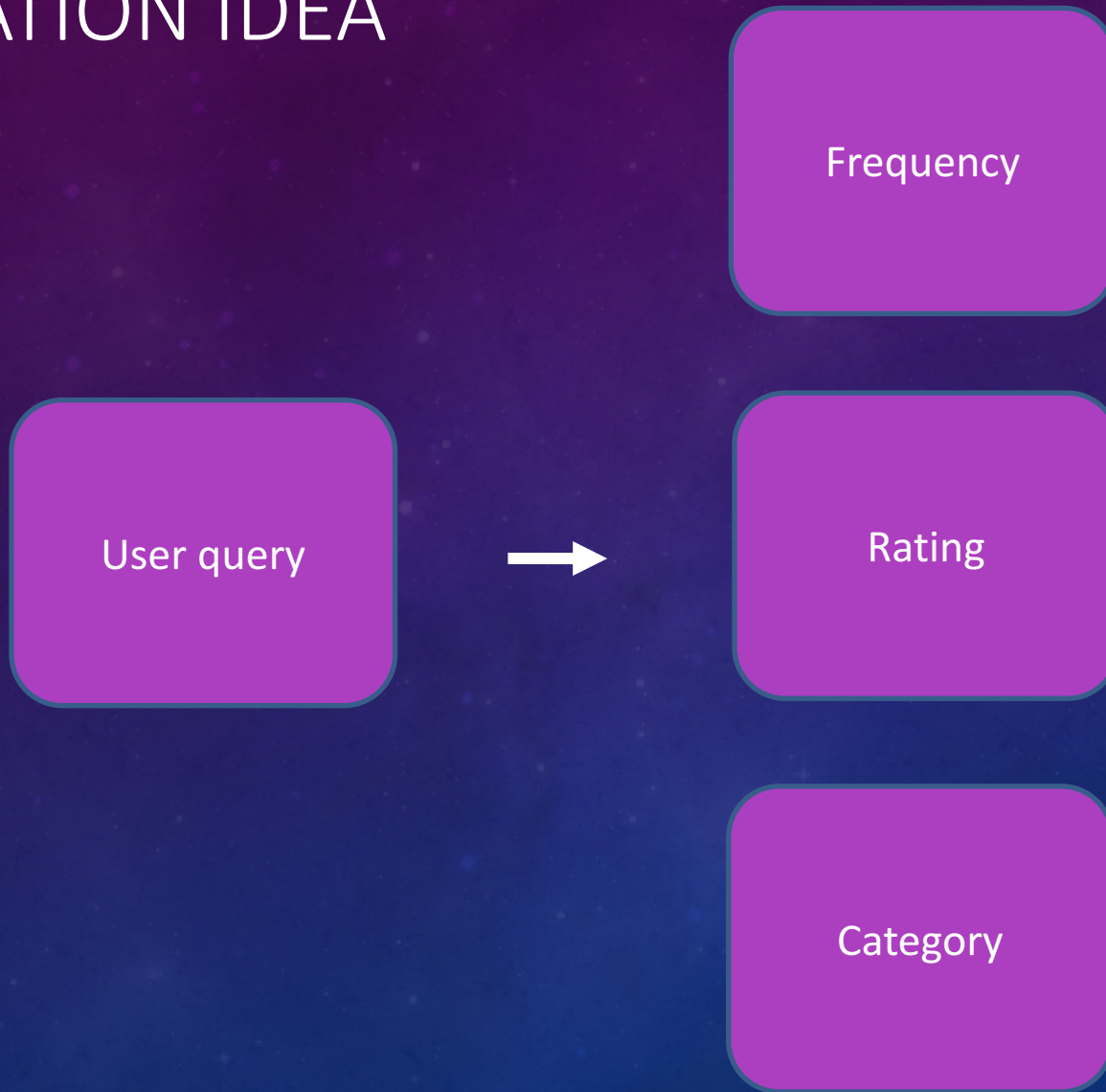
How do you write a successful piece of literature?

Helping authors choose the best expression to reach their target audience

- Books
- Blogs
- Technical Publications



APPLICATION IDEA



DATA SOURCES

Frequency

Google N-grams

Rating, category

Google Books

+

Good Reads

- Google n-grams dataset: fixed dataset with the frequency of 1 to 5 word expressions in published books and magazines over time from about 1500 until now.
- Google Books API : text searches returning bibliographical information
- Good Reads API: reviews, ratings, tags of books

VOLUME * VARIETY * VELOCITY

Volume:

- * Google N-gram: 3.5 TB data stored in our S3 bucket (only English bigrams)
- * we used a 15 GB subset of it for bigrams starting with “AB...”

Variety:

- * JSONs from Google Books, GoodReads API calls
- * Google N-gram CSVs

Velocity:

- * Kibana dashboard refreshes search history every 5 seconds

Services used:

- * Amazon EC2
- * Amazon S3
- * Amazon RDS
- * Amazon ES Services

Tools used:

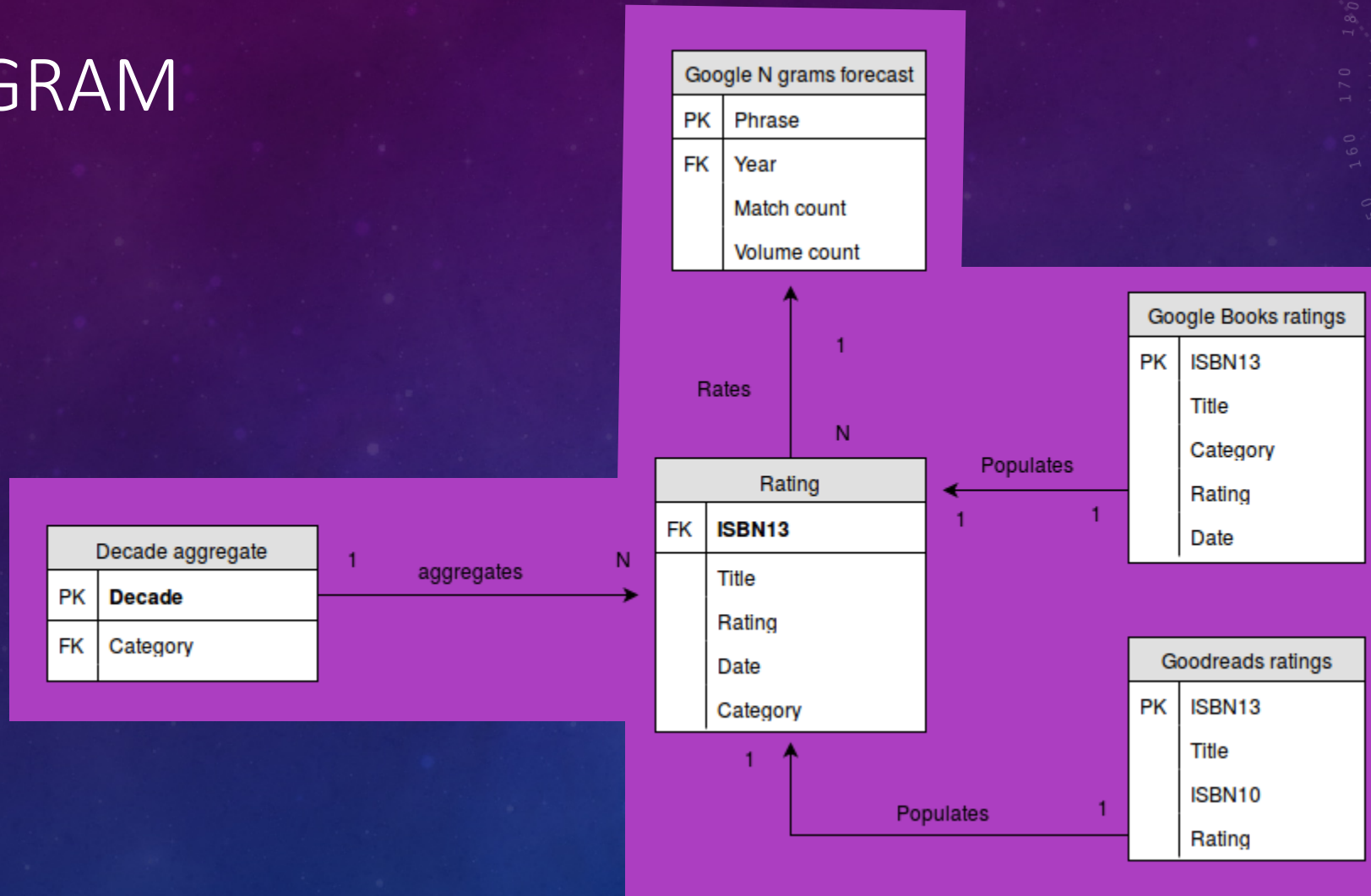
- * Apache Spark (PySpark, SparkR)
- * PostgreSQL
- * ElasticSearch, Kibana
- * Shiny, D3
- * Python, R

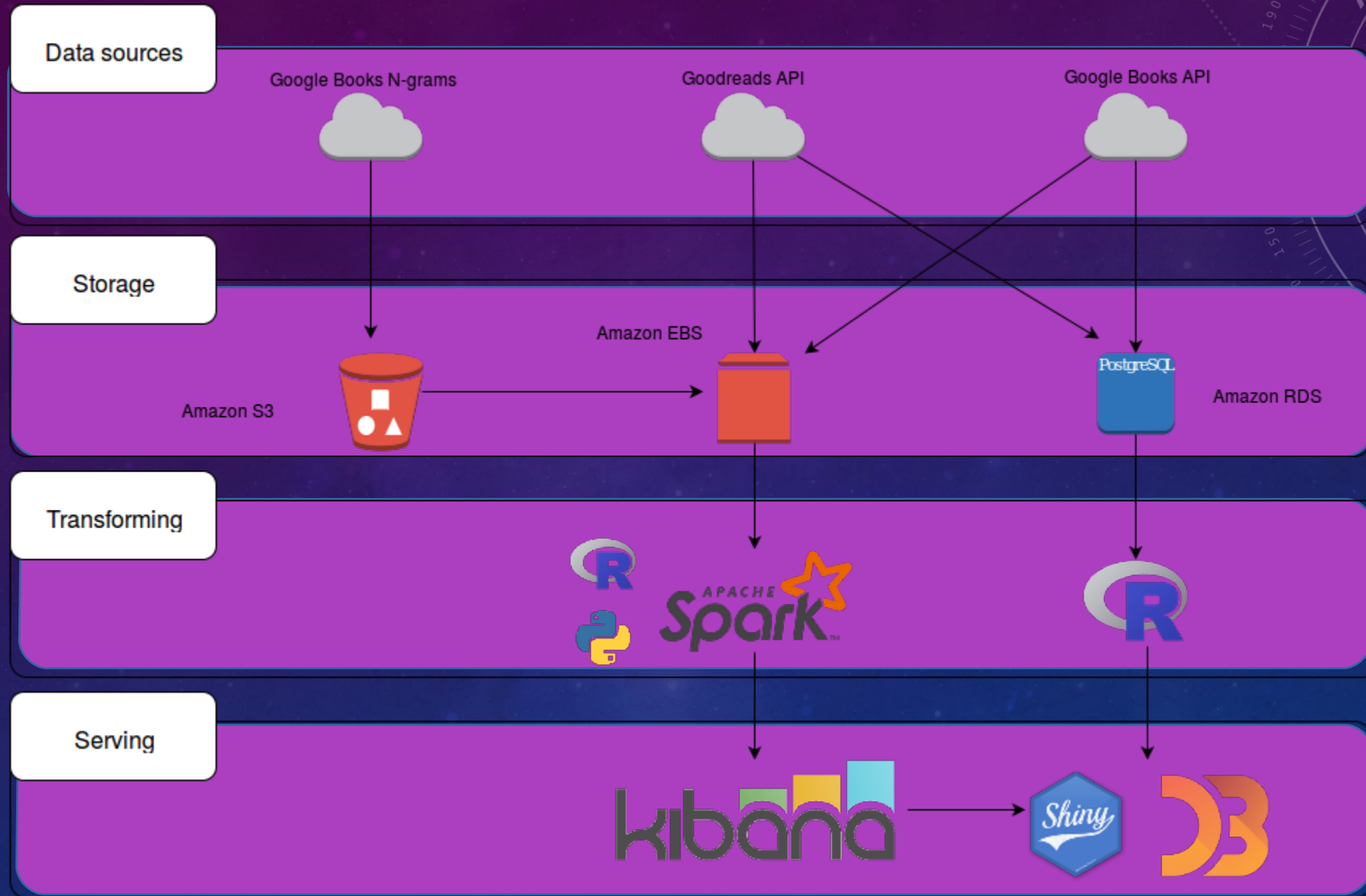
DATA CHALLENGES



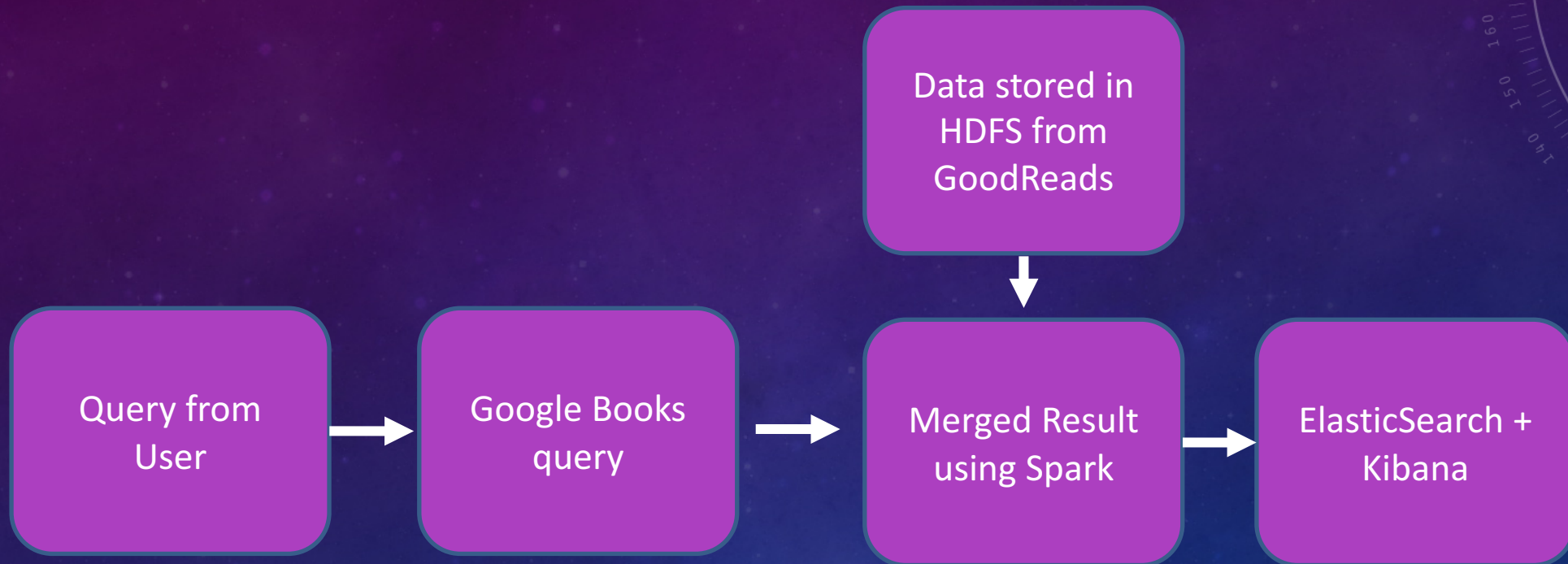
- Thousands of compressed CSV files in Google N-gram database, with some files taking 2-3 GB of space
- Good Reads has a restriction of 1 API call per second
- Filtering data from Good Reads such that we get Ratings for books written in English with a valid ISBN
- Filtering data from Google Books such that we get Ratings, ISBN and valid category for books .
- Amazon S3 maintains the Google N-gram database in an encoded format, requires building Hive/Spark with LZOCCodec.

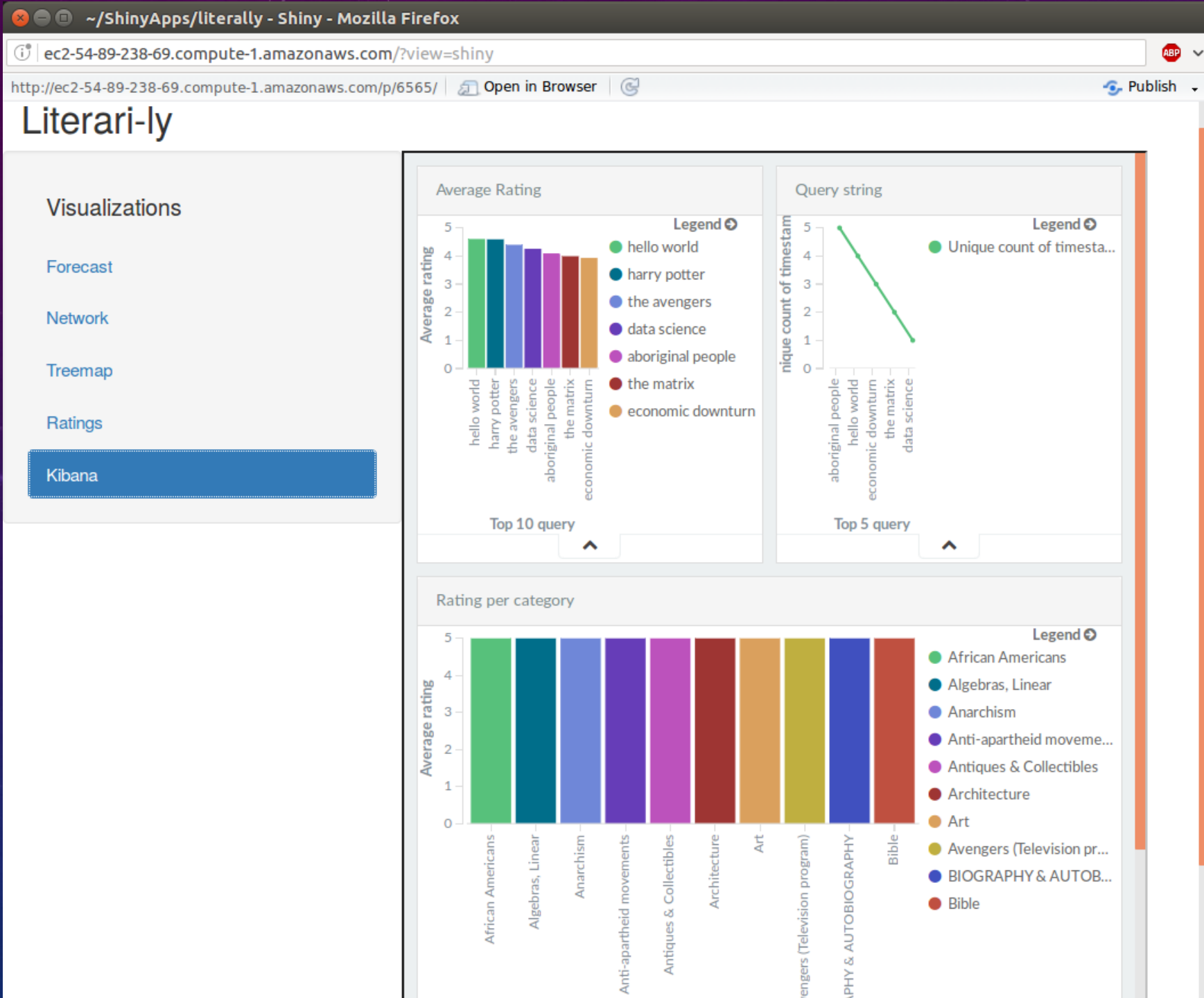
ER DIAGRAM



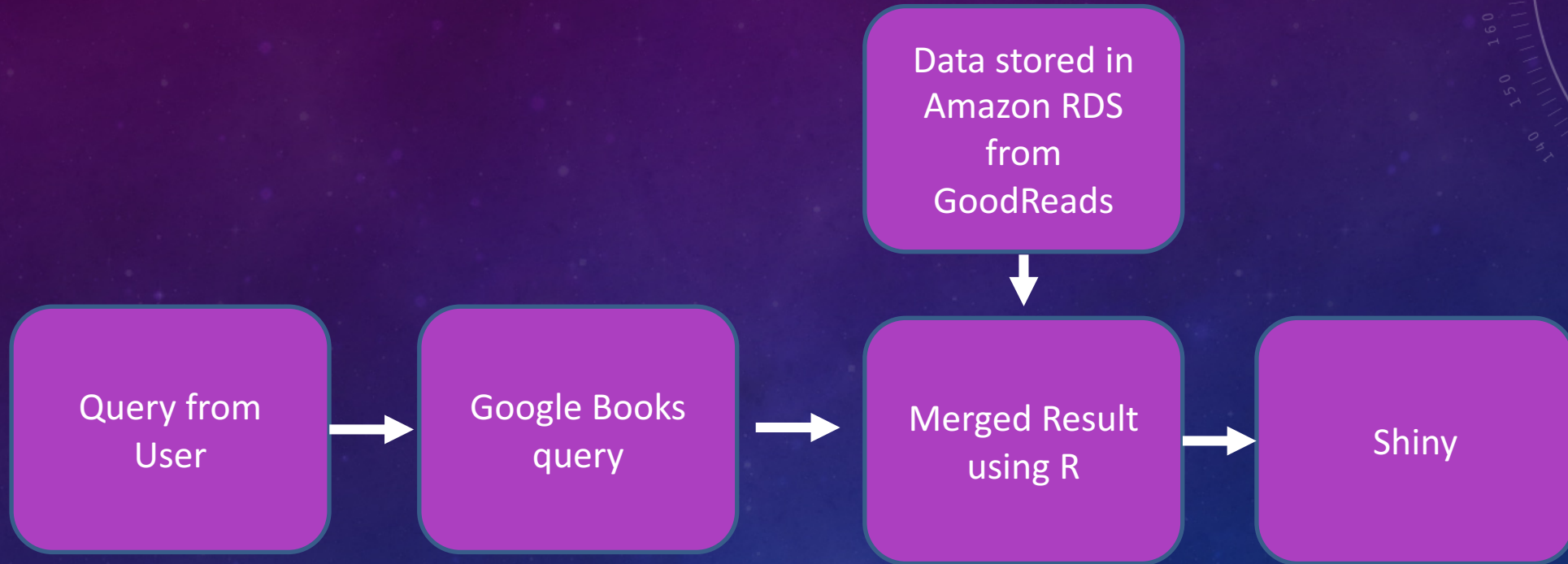


ARCHITECTURE * KIBANA DASHBOARD





ARCHITECTURE * RATING AND CATEGORIES



Literari-ly

Visualizations

Forecast

Network

Treemap

Ratings

Kibana

Links between categories and decades of publication

By default, the 5 highest frequency categories are plotted but the data will likely have many more.

Search expression

aboriginal people

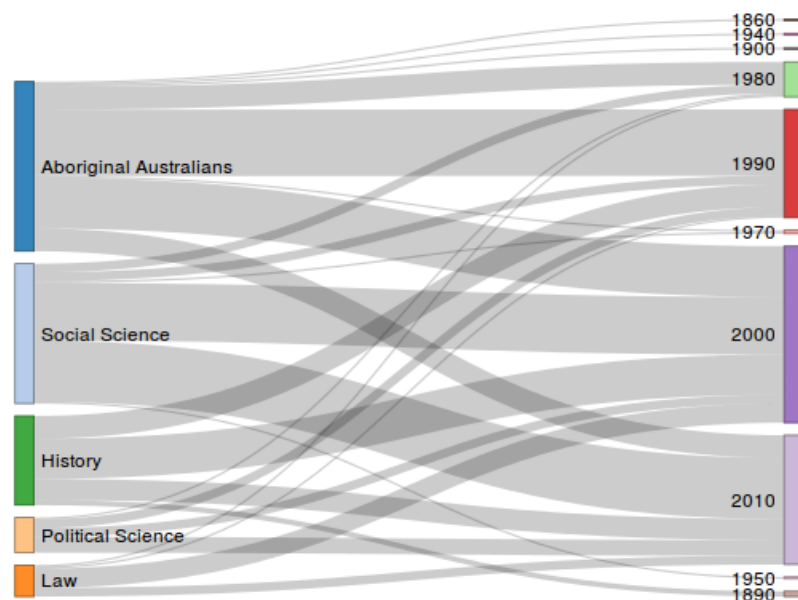
Search

Font

A horizontal number line with tick marks every 1 unit, labeled from 1 to 20. A blue slider is positioned at the number 15. The number 15 is also displayed in a blue box above the slider.

Number of most frequent categories to keep

5



Literari-ly

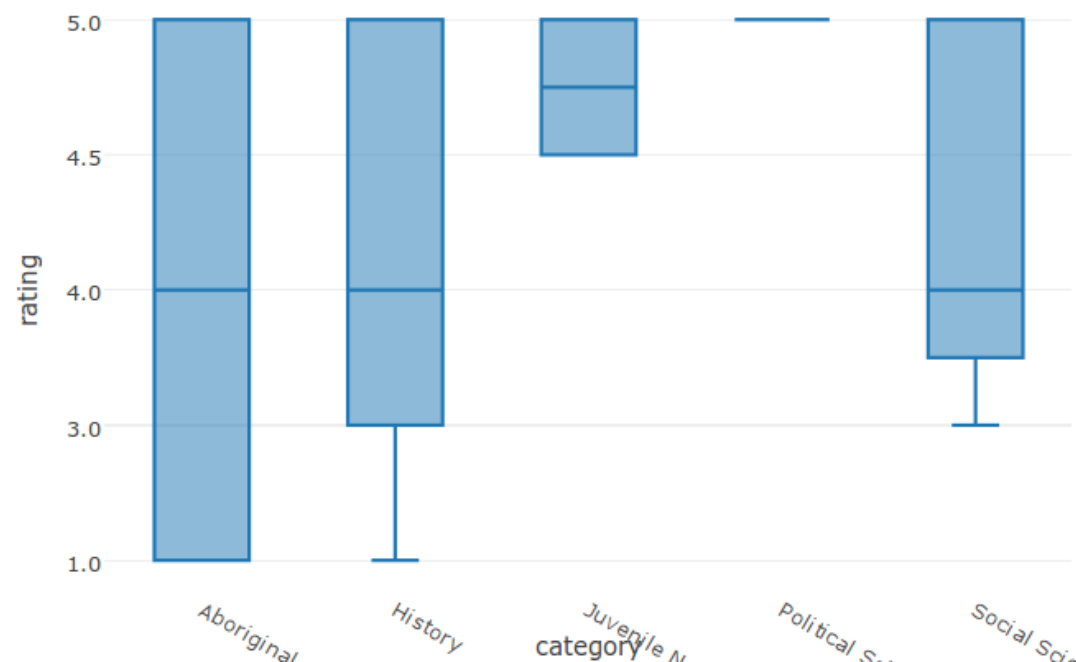
Visualizations

[Forecast](#)[Network](#)[Treemap](#)[Ratings](#)[Kibana](#)

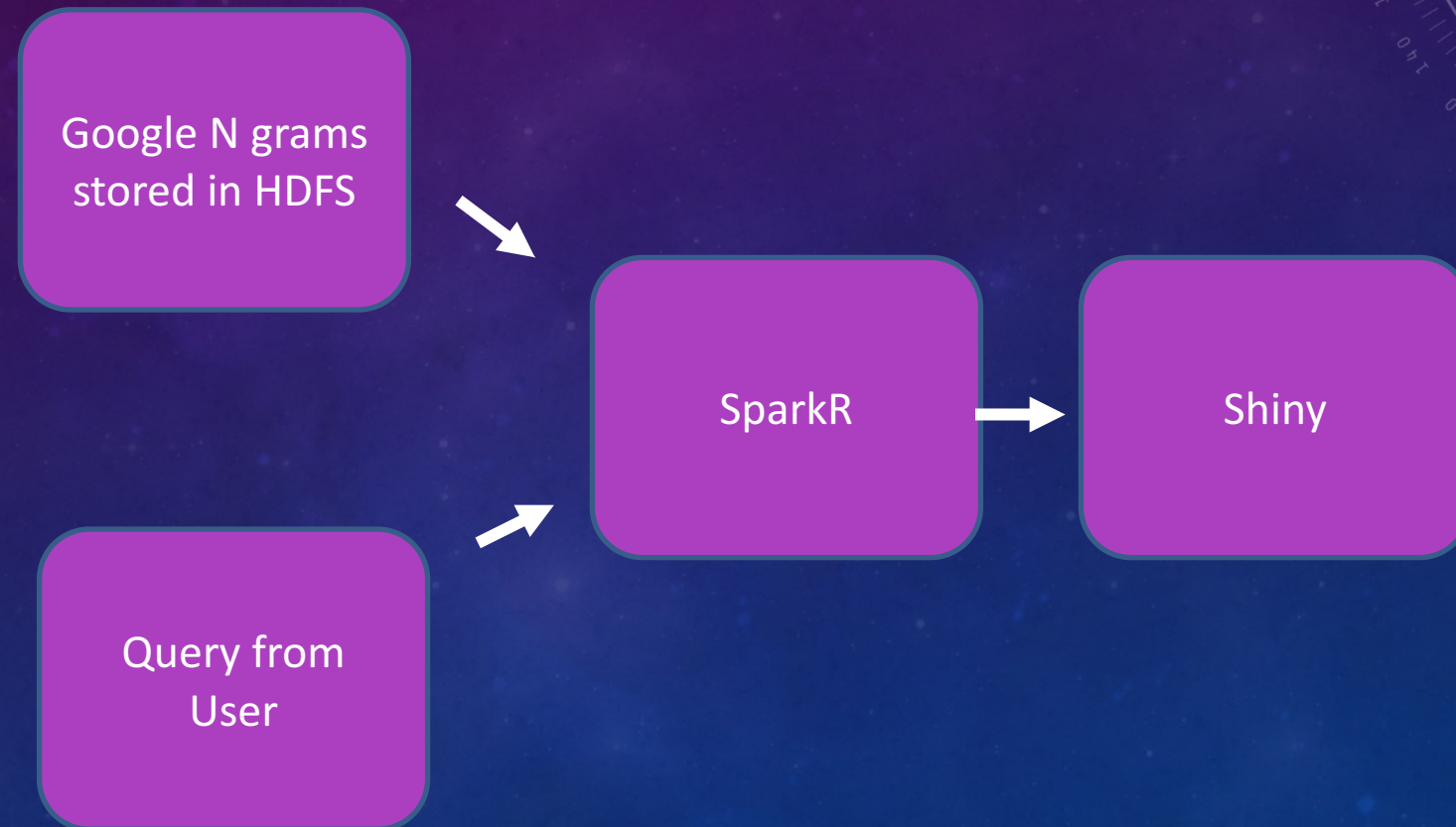
Ratings of books by categories

Search expression

Number of most frequent categories to keep



ARCHITECTURE * FORECASTING



Literari-ly

Visualizations

Forecast

Network

Treemap

Ratings

Kibana

Forecast popularity of an expression

Search expression

aboriginal people

Search

Role of first word

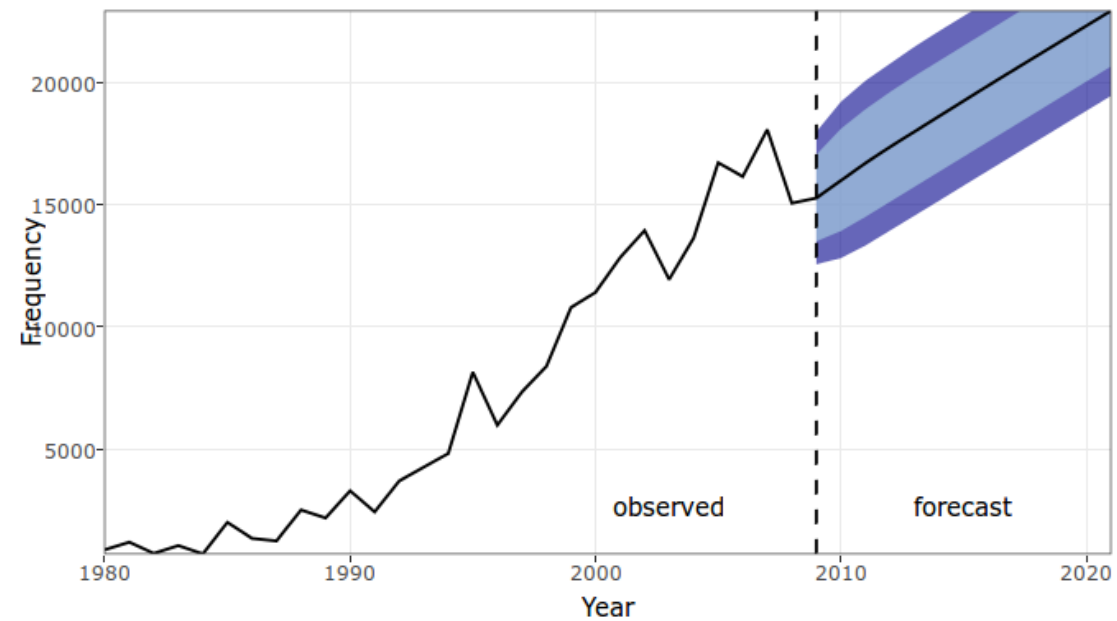
Adjective

Adjective

Adverb

Noun

Verb



FUTURE WORK

- The expressions could be queried in the Twitter API and return the distribution of user interests who tweeted the queried expressions.
- Example implementation of finding out user interests at scale:

<http://www.mpi-sws.org/~mzafar/papers/recsys14-userinterests.pdf>

<http://twitter-app.mpi-sws.org/who-likes-what/>