# News Recommendation System

## Introduction

In today's society a news user is often overwhelmed by the amount of available articles, spread over thousands of different topics. Reading the front page of a news site, will most likely result in articles that we do not find relevant and explicitly searching for a topic is often an effort that we do not want to make. News recommendation services therefore aim at providing articles that the user wants to read, without the need to search for them.

In this work we explores the effectiveness of integrating a user's reading history in the recommendation engine, in the form of articles that the user has liked or visited. Furthermore, we consider different sources of meta-data that can better describe the content of an article for the purpose of producing recommendations.

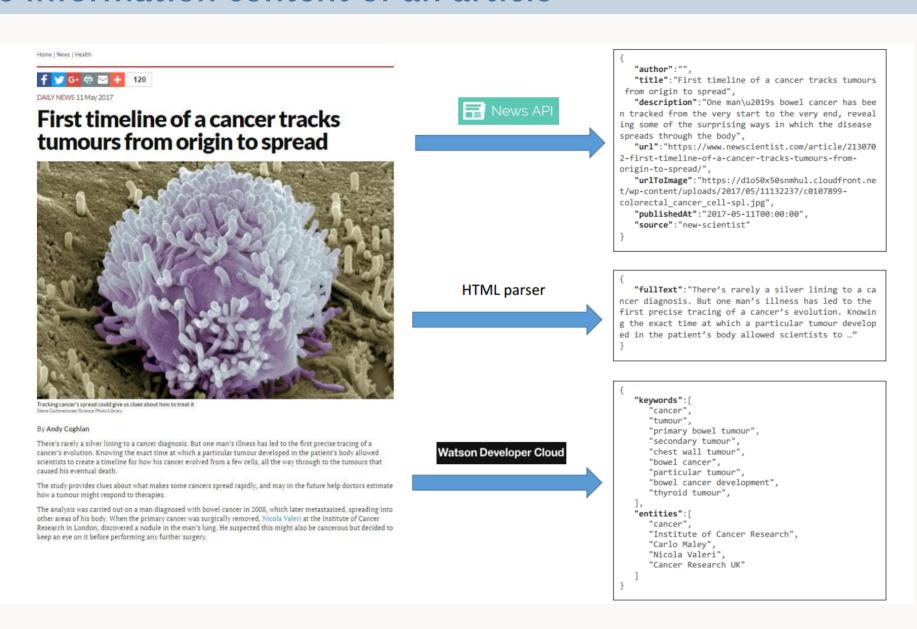## The information content of an article



Figure: Example of metadata extraction from a news

► To be able to identify the content of a piece of news we do not rely uniquely on its title of text, but we leverage several sources of meta-data.

► Simple meta-data can be the author, published date and publishing source

► More advanced (and useful) characteristics can be found extracting entities and keywords through natural language processing techniques

► Enriching the documents in storage with these information is a key step in our processes
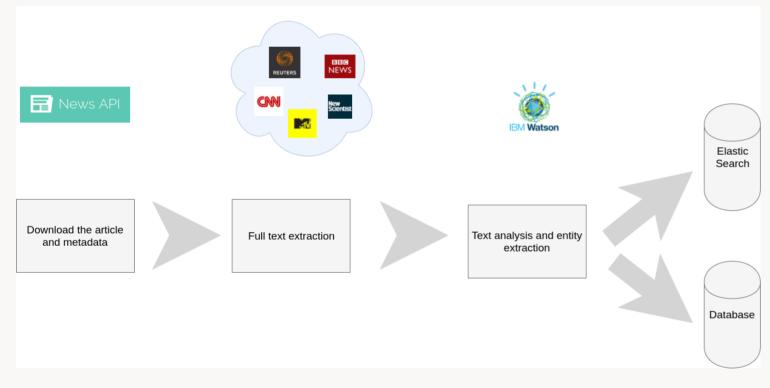
## Data ingestion pipeline



Figure: Overview of the pipeline

► First we retrieve news from the *News API* service, The API acts as an aggregator for different news websites and additionally returns meta-data such as author, publishing date, source and a short description

► After that, the full text of every article is scraped from the source using a web crawler

► Then we use the natural language processing services from *IBM Watson Developer Cloud* to extract keywords and entities from the full text

► Finally, we persist the documents in a database and we index them in an Elasticsearch instance

## Evaluation

In our evaluation we track the Discounted cumulative gain (DCG) scores assigned by the users to our recommendations as a function of the information we hold on their preferences.

1. Create $N_Q$ queries to be used in the news search engine
2. Every participant sequentially execute these queries and evaluates the recommendations:
   2.1 Run a query in the news article search engine
   2.2 Consider the list of articles retrieved and choose the first 2 documents that he or she finds interesting starting from the top
   2.3 Look at the $K = 5$ articles that appear in the recommendation list and assign them a relevance score based on the aforementioned table
   2.4 Run another randomly selected query
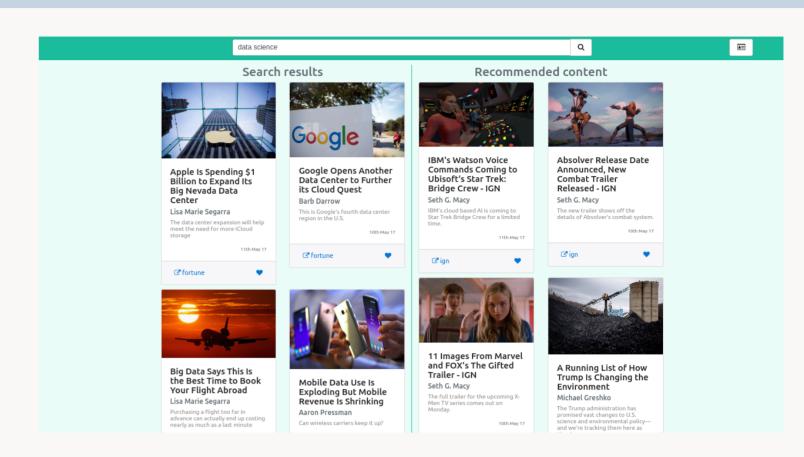
## User interface



Figure: Front-end application.

► The user can search for queries and see the results (left side), like in a normal search engine

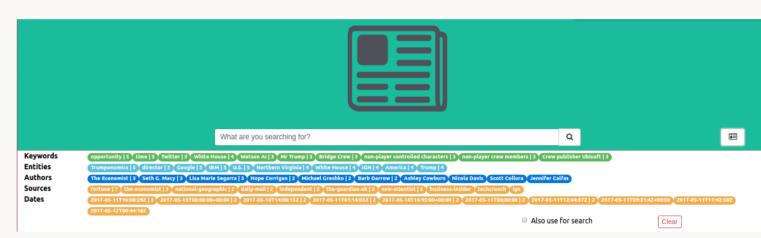► The user is given recommendation (right side), these are continuously updated.



Figure: User's interest view.

► The user profile, or interest view, is based on the articles that the user has liked or visited from the search engine results

► It is possible to clear the user's profile

## Results

► Discounted cumulative gain for a *top-K* recommendation system

$$DCG = \sum_{k=1}^{K} \frac{rel(k)}{\log_2(k+1)}$$

► The results of our user study:
   • 11 people in the test group / 11 people in the control group
   • 5 articles to rank in the recommendation list
   • 10 search queries in random order, allowing to select 2 results from each query
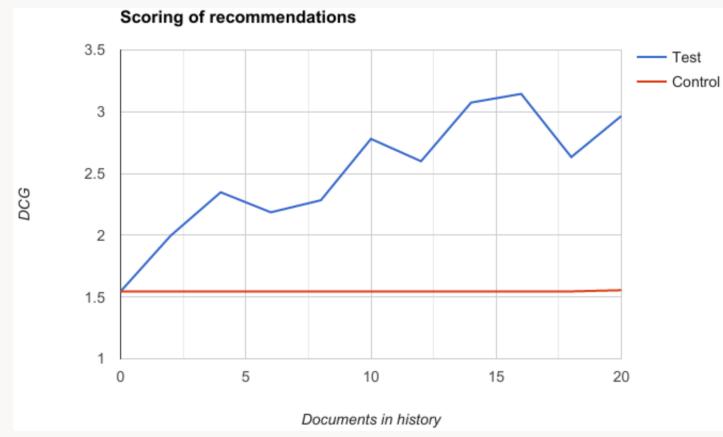


Figure: Evolution of DCG score for the test and control groups

► User profiling and customization are the base for every recommendation system

► We have validated the usage of meta-data as a key ingredient of producing recommendations

► The relative weight of every meta-data factor can be tweaked to vastly improve the results

**Federico Baldassarre, Sandra Bäckström, Sebastian Bujwid, Zimin Chen**

**KTH Royal Institute of Technology**