# Data Recommender System: Improving the Discovery of Environmental Datasets through Text Analytics and Usage Mining

Anusuriya Devaraju, Uwe Schindler, Robert Huber, Michael Diepenbroek
{adevaraju, uschindler, rhuber, mdiepenbroek}@marum.de
PANGAEA Data Publisher for Earth & Environmental Science,
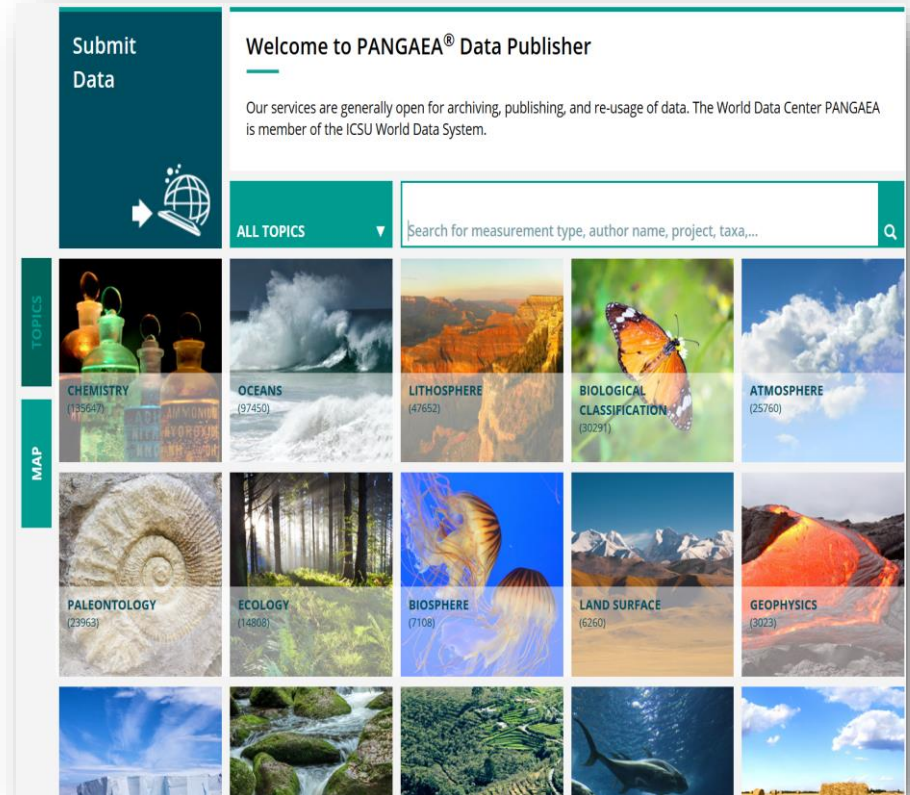MARUM - Center for Marine Environmental Sciences, University of Bremen, Germany.

# Presentation Outline

- PANGAEA Digital Data Repository

- Data Recommender System
    a. Metadata-based data recommendation
    b. Users interaction-based data recommendation

- Online Evaluation & Results

- Conclusions

# PANGAEA Digital Data Repository

- Founded in 1993

- Jointly managed by AWI and MARUM.

- Datasets from researchers, projects, research centres and infrastructures (national & international) published with DOIs

- Data types, e.g., time series, spatial, images, audio, video.

**> 386000** **scientific datasets (30.08.2019)**



PANGAEA Data Portal (https://pangaea.de/)

# Data Search in PANGAEA

- PANGAEA offers tools/APIs for meta(data) access and discovery.



Search functions in PANGAEA

Dataset and its related research objects

Devaraju, A. et al., Data Recommender System: Improving the Discovery of Environmental Datasets through Text Analytics and Usage Mining, Geospatial Sensing Conference 2019.

# Data Search in PANGAEA

# Search Meets Discovery

- "Search is often struggling to deliver meaningful results, unless you're very explicit and goal oriented…" (Bibblio, Search vs Discovery, 2015)

- How about users ..
  - may not know what/how to search
  - who are not aware of the range of available datasets
  - if presented with many datasets may not be able to choose the datasets-of-interest.



PANGAEA.

| ALL TOPICS ▼ | Water temperature and salinity |

61900 datasets found on search for »Water temperature and ...«

Keyword search brings too many almost identical datasets; diversity is missing!

1. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1068896.
2. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1068897.
3. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1068898.
4. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1068899.
5. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176596.
6. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176597.
7. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176598.
8. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176599.
9. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176696.
10. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176697.
11. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176698.
12. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1176699.
13. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1474196.
14. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1474197.
15. **Clarke, RA (2006):** Water temperature and salinity from profiling float 1474198.

## How can users discover relevant and 'novel' datasets on the portal?

# Recommender System

A recommender system is an information filtering system that provide users with personalized contents and services.

# PANGAEA Data Recommender



Data Users

Dataset

View

**1** Metadata-based Recommendations
Datasets with similar metadata

**2** Interaction-based Recommendations
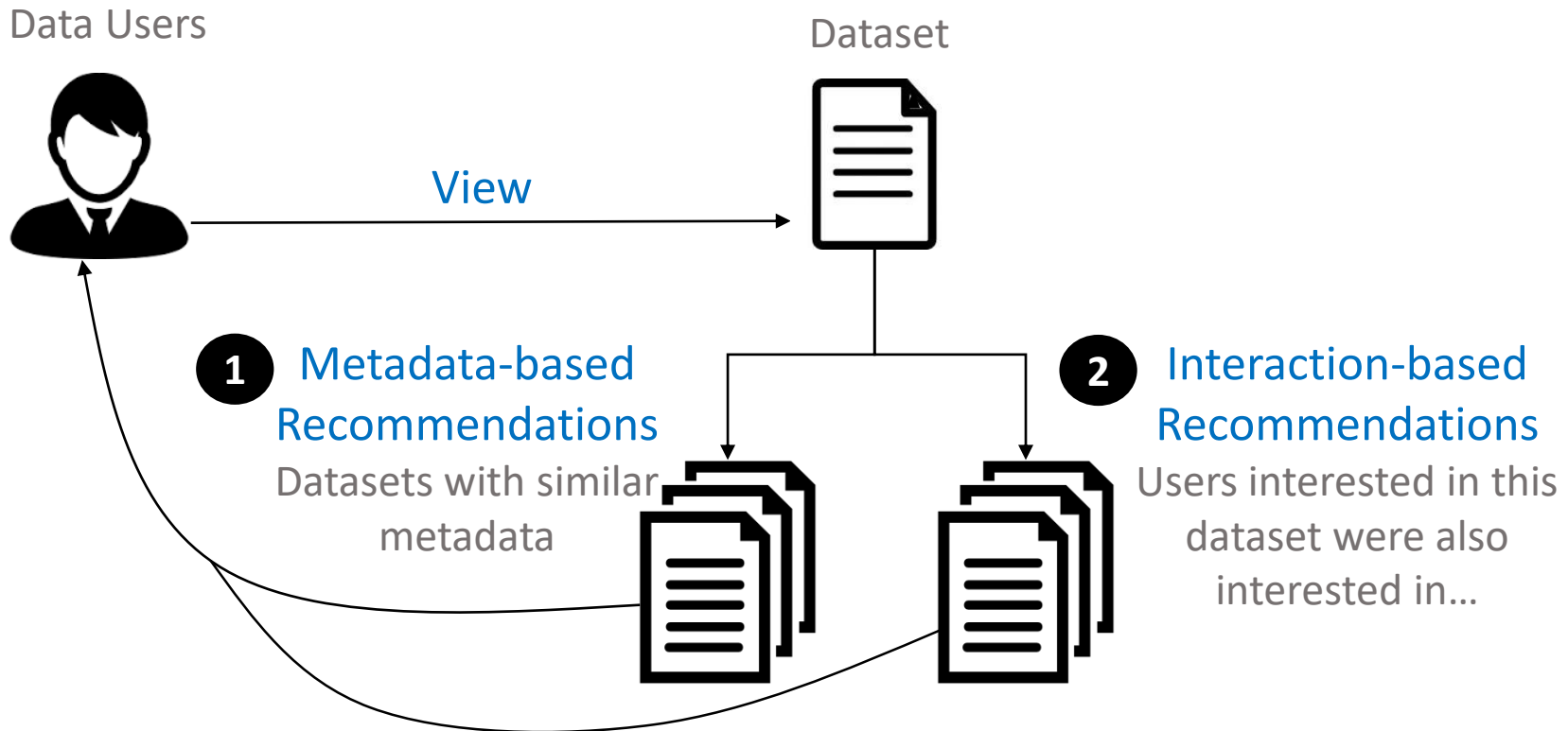Users interested in this dataset were also interested in…

Dahl, Kristina A; Oppo, Delia W (2006): (Table 3) Mg/Ca ratios of Globigerinoides ruber from Arabian Sea sediments. PANGAEA, https://doi.org/10.1594/PANGAEA.834987,

In supplement to: Dahl, KA; Oppo, DW (2006): Sea surface temperature pattern reconstructions in the Arabian Sea. Paleoceanography, 21(1), PA1014, https://doi.org/10.1029/2005PA001162
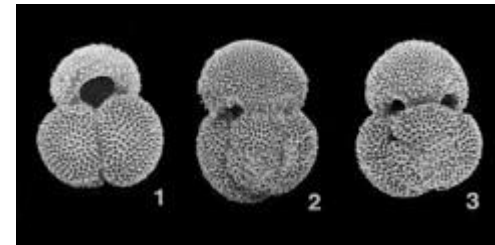
Image : World Register of Marine Species

More diverse recommendations.

## Metadata-based Recommendation

**Datasets with similar metadata**

Weldeab, S; Schneider, RR; Kölling, M et al. (2005): Mg/Ca ratios of Globigerinoides ruber of sediment core GeoB4905-4. https://doi.org/10.1594/PANGAEA.738242

Cléroux, C; Debret, M; Cortijo, E et al. (2012): Mg/Ca and Sr/Ca ratios on Globigerinoides ruber (white) in sediment core MD99-2203, Cape Hatteras. https://doi.org/10.1594/PANGAEA.776433

Tian, J; Pak, DK; Wang, P et al. (2006): (Appendix 2) Mg/Ca ratios of Globigerinoides ruber from ODP Site 184-1143. https://doi.org/10.1594/PANGAEA.707839

## Interaction-based Recommendation

**Users interested in this dataset were also interested in**

Sirocko, F; Garbe-Schönberg, C-D; Devey, CW (2000): Composition of sediments from the Arabian Sea. https://doi.org/10.1594/PANGAEA.728741

Schulz, H (1995): Planktic foraminiferal assemblage for the 10kyr time slice from different sediment cores. https://doi.org/10.1594/PANGAEA.51969

Munz, P; Siccha, M; Lückge, A et al. (2015): Distribution of planktic foraminifera in surface sediments in the northeastern Arabian Sea. https://doi.org/10.1594/PANGAEA.853966

# Metadata-based Data Recommendation

- Leverages ElasticSearch More Like This (MLT) with boosting.

- MLT returns datasets that are similar to a provided data based on metadata elements, e.g., title, abstract, related publication, authors, topics, projects, devices, campaign, location, time.
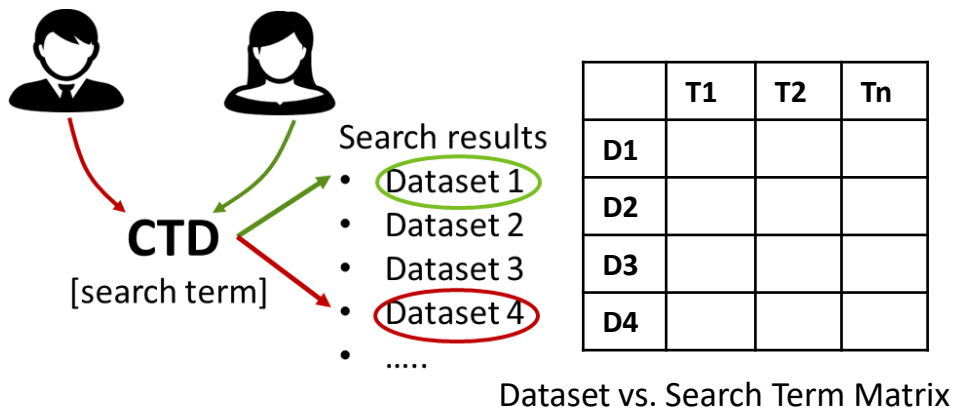


Devaraju, A. et al., Data Recommender System: Improving the Discovery of Environmental Datasets through Text Analytics and Usage Mining, Geospatial Sensing Conference 2019.

# Usage-based Data Recommendation

- Utilizes 3 **user interactions** (extracted from the server logs) such as *search interaction, joint download,* and *total download.*

## Search interaction

Datasets examined after launching similar searches are likely to be similar.



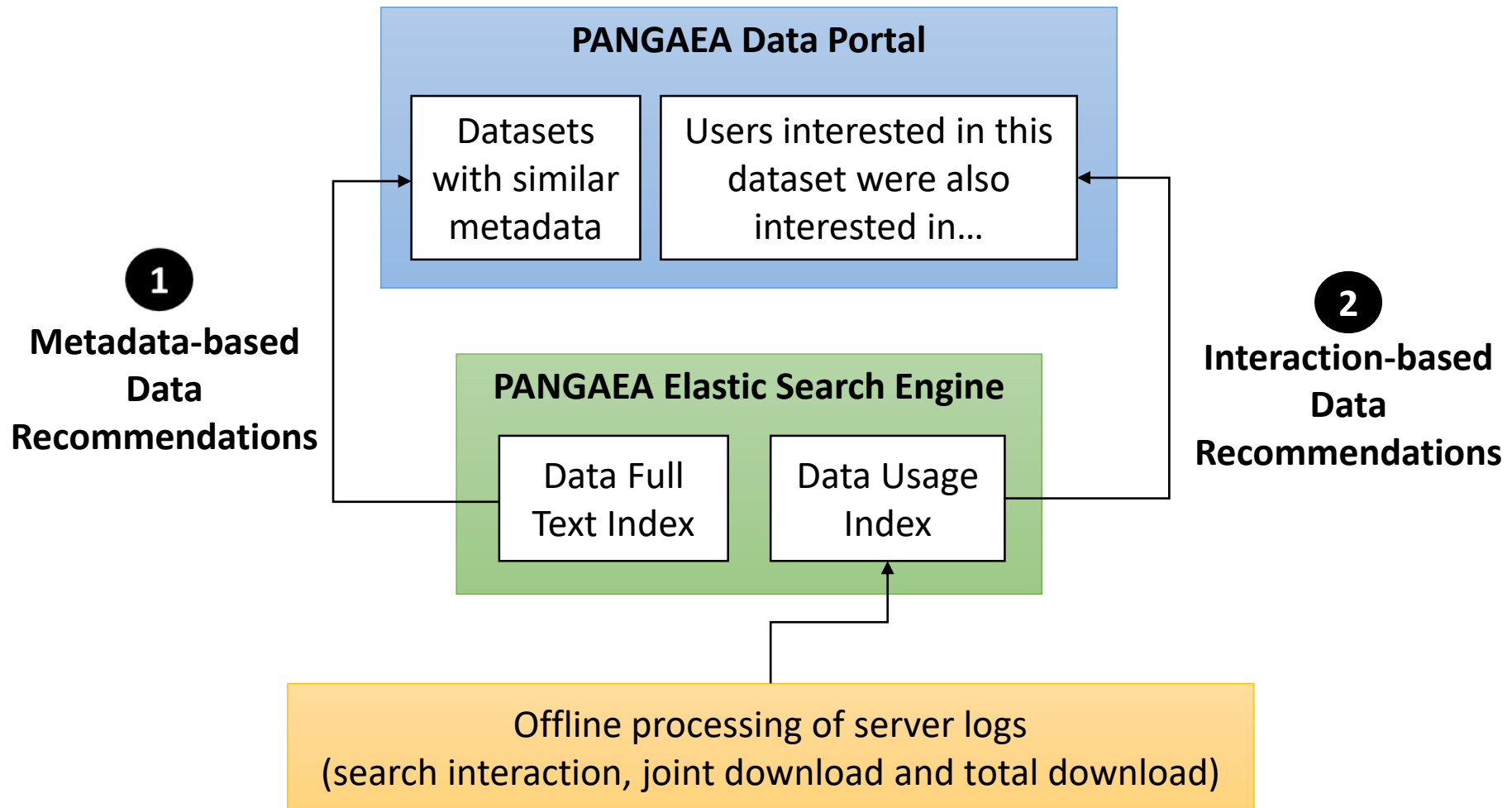Dataset vs. Search Term Matrix

## Joint download

Jointly downloaded datasets are likely to be similar.



Dataset vs. User Matrix

Devaraju, A. et al., Data Recommender System: Improving the Discovery of Environmental Datasets through Text Analytics and Usage Mining, Geospatial Sensing Conference 2019.
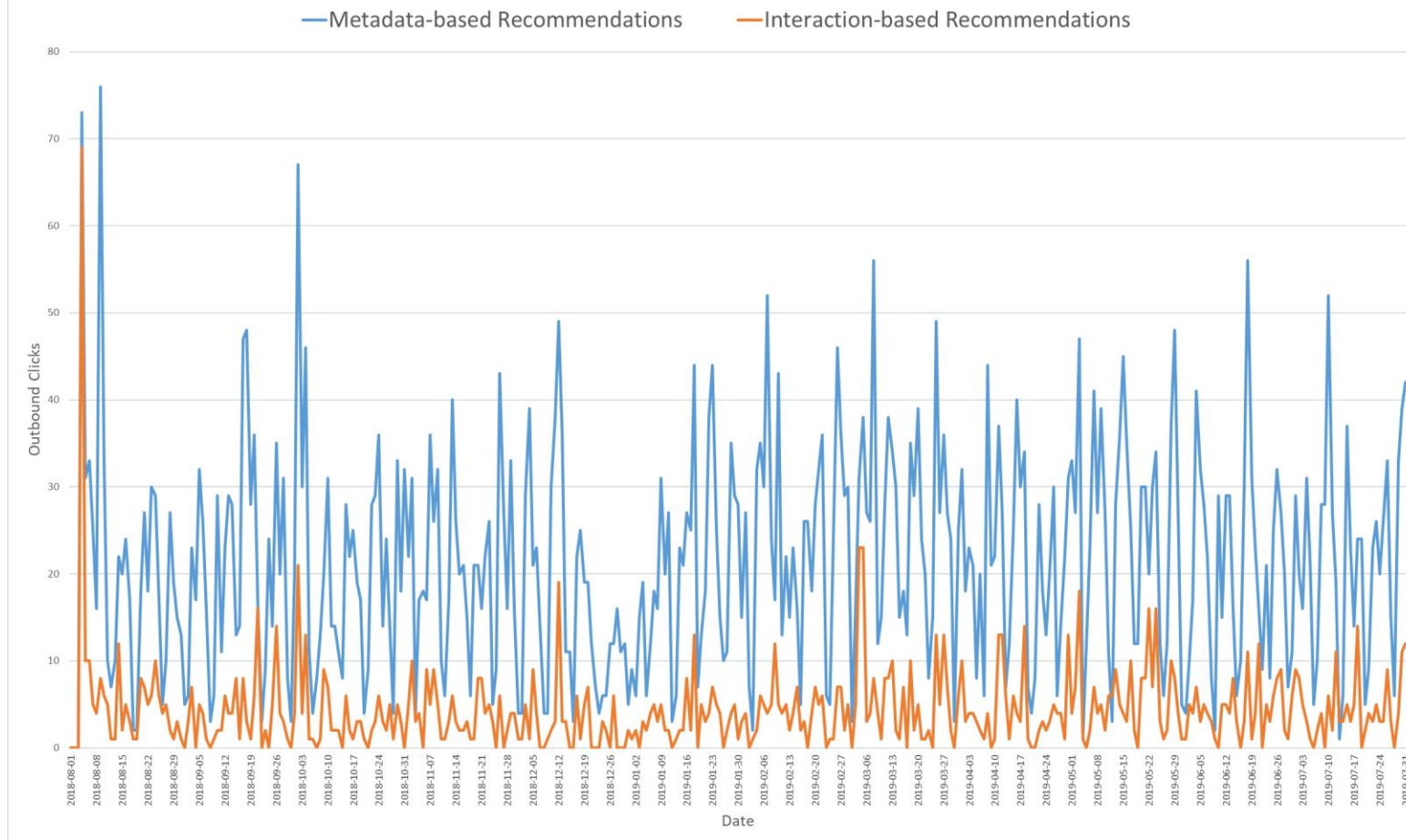
# System Implementation
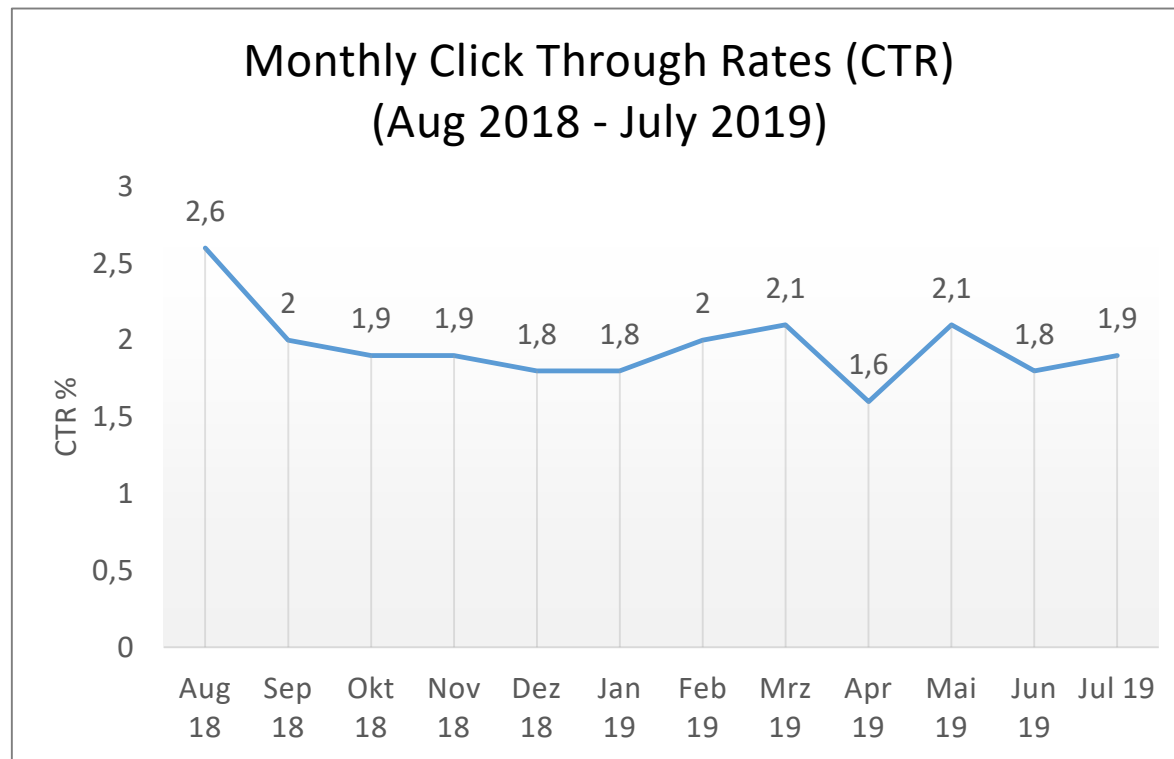
# Online Evaluation

- Average clicks/day : 26
- % contribution to total events
  - Metadata : 83%
  - Interaction : 17%

## Outbound Clicks from Data Recommendations (01.08.2018-31.07.2019)

— Metadata-based Recommendations    — Interaction-based Recommendations

Devaraju, A. et al., Data Recommender System: Improving the Discovery of Environmental Datasets through Text Analytics and Usage Mining, Geospatial Sensing Conference 2019.

# Online Evaluation

- Click Through Rate (CTR) = Clicks ÷ Impressions



**Monthly Click Through Rates (CTR)**
**(Aug 2018 - July 2019)**

| Month | CTR % |
|-------|-------|
| Aug 18 | 2,6 |
| Sep 18 | 2 |
| Okt 18 | 1,9 |
| Nov 18 | 1,9 |
| Dez 18 | 1,8 |
| Jan 19 | 1,8 |
| Feb 19 | 2 |
| Mrz 19 | 2,1 |
| Apr 19 | 1,6 |
| Mai 19 | 2,1 |
| Jun 19 | 1,8 |
| Jul 19 | 1,9 |

Devaraju, A. et al., Data Recommender System: Improving the Discovery of Environmental Datasets through Text Analytics and Usage Mining, Geospatial Sensing Conference 2019.

# Conclusions

- Developed a recommender to improve data discovery, which presents users with two kinds of recommendations.

- Building a data recommender on top of the ElasticSearch enhances the scalability and maintainability of the recommender system.

Ongoing/planned work:

- More features – ontological concepts of parameters.
- Improve the presentation of recommendations.
- Extend online evaluation (conversion rate).