

User-Driven Re-Ranking for Adapting the Variety in Search Results

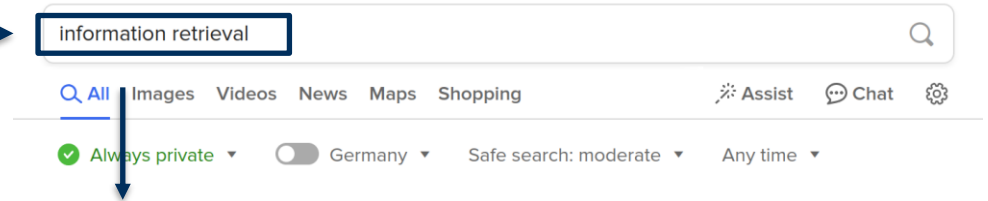
Daphne Auer, Divyasha Sunil Naik, Birgitta König-Ries

Open Search Symposium 2024, 11.10.2024

Information
Need

What is information
retrieval about?

Query



Document Ranking

[https://en.wikipedia.org › wiki › Information_retrieval](https://en.wikipedia.org/wiki/Information_retrieval)

Information retrieval - Wikipedia

Information retrieval (IR) in computing and information science is the task of identifying and retrieving information system resources that are relevant to an information need. The information need can be specified in the form of a search query. In the case of document retrieval, queries can be based on ful...

[https://www.geeksforgeeks.org › what-is-information-retrieval](https://www.geeksforgeeks.org/what-is-information-retrieval)

What is Information Retrieval? - GeeksforGeeks

Sep 19, 2023 · Information Retrieval (IR) can be defined as a software program that deals with the organization, storage, retrieval, and evaluation of information from document repositories, particularly textual information. Information Retrieval is the activity of obtaining material that can usually be...

[https://www.britannica.com › technology › information-retrieval](https://www.britannica.com/technology/information-retrieval)

Information retrieval | Definition, Methods, & Facts | Britannica

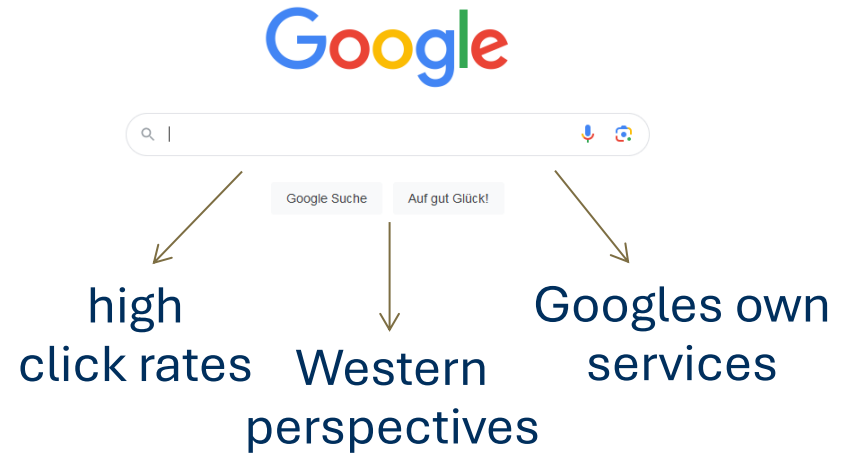
information retrieval, recovery of information, especially in a database stored in a computer. Two main approaches are matching words in the query against the database index (keyword searching) and traversing the database using hypertext or hypermedia links. Evolving information-retrieval technique...

Why should we care about **variety**?

The screenshot shows a news article on the website 'GROUND'. The article title is "Nobel Prize in medicine awarded to US duo Victor Ambros, Gary Ruvkun for discovery of microRNA". Below the title is a progress bar with three segments: "Left 35%", "Center 51%", and "R 14%".

Erdogan · Albania
Turkey's Erdogan inaugurates a
ground.news

The screenshot shows a news article on the website 'BLINDSPOT'. The article title is "Hamis chief Sinwar revives 'suicide bombings' after...". Above the title is the text "Stories disproportionately reported by the Left or the Right". Below the title is a progress bar with three segments: "10% Center 30%" and "Right 60%".



Why should we care about **variety**?

Range of perspectives
→ Diversity

Nobel Prize in medicine awarded to US duo Victor Ambros, Gary Ruvkun for discovery of microRNA

Left 35% | Center 51% | R 14%

BLINDSPOT™

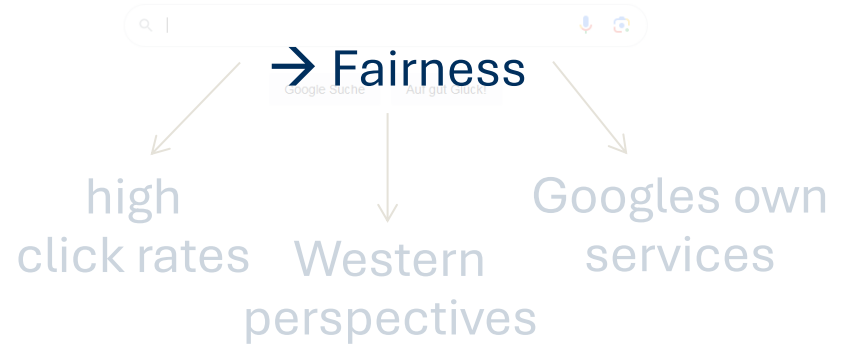
Stories disproportionately reported by the Left or the Right

Blindspot 10 Sources

Hamas chief Sinwar revives 'suicide bombings' after...

10% Center 30% | Right 60%

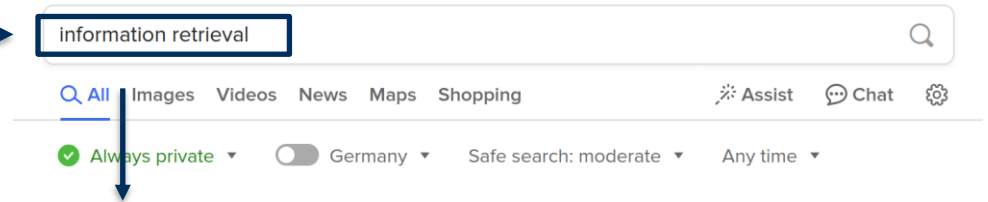
Characterization using 'protected attributes'



Information
Need

What is information
retrieval about?

Query



Document Ranking

[https://en.wikipedia.org › wiki › Information_retrieval](https://en.wikipedia.org/wiki/Information_retrieval)

Information retrieval - Wikipedia

Information retrieval (IR) in computing and information science is the task of identifying and retrieving information system resources that are relevant to an information need. The information need can be specified in the form of a search query. In the case of document retrieval, queries can be based on ful...

[https://www.geeksforgeeks.org › what-is-information-retrieval](https://www.geeksforgeeks.org/what-is-information-retrieval)

What is Information Retrieval? - GeeksforGeeks

Sep 19, 2023 · Information Retrieval (IR) can be defined as a software program that deals with the organization, storage, retrieval, and evaluation of information from document repositories, particularly textual information. Information Retrieval is the activity of obtaining material that can usually be...

[https://www.britannica.com › technology › information-retrieval](https://www.britannica.com/technology/information-retrieval)

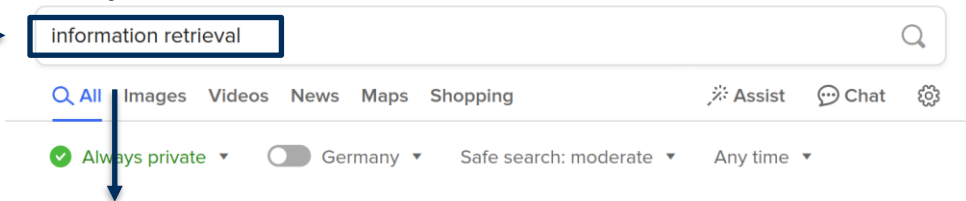
Information retrieval | Definition, Methods, & Facts | Britannica

information retrieval, recovery of information, especially in a database stored in a computer. Two main approaches are matching words in the query against the database index (keyword searching) and traversing the database using hypertext or hypermedia links. Evolving information-retrieval technique...

Information
Need

What is information
retrieval about?

Query



Variety
Setting

Diversity Degree δ



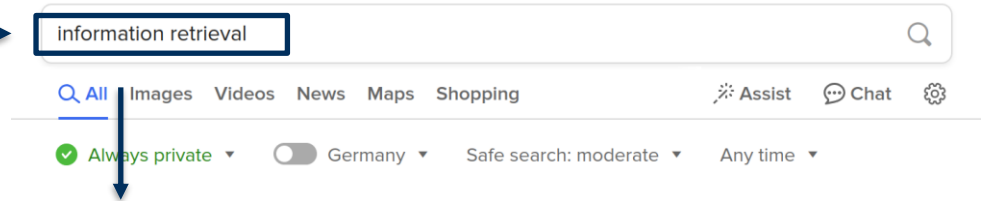
Document Ranking

A screenshot of search results for the query "information retrieval". The results are ranked and displayed in a list. The first result is from Wikipedia, titled "Information retrieval - Wikipedia". The second result is from GeeksforGeeks, titled "What is Information Retrieval? - GeeksforGeeks". The third result is from Britannica, titled "Information retrieval | Definition, Methods, & Facts | Britannica". Each result includes a brief description of the topic.

Information
Need

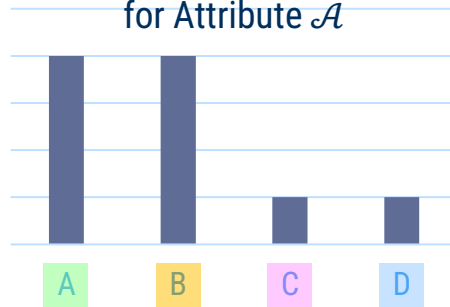
What is information
retrieval about?

Query



Variety
Setting

Target Distribution
for Attribute A



Document Ranking

A screenshot of search results for the query 'information retrieval'. The results are ranked and include:

- https://en.wikipedia.org/wiki/Information_retrieval
Information retrieval - Wikipedia
Information retrieval (IR) in computing and information science is the task of identifying and retrieving information system resources that are relevant to an information need. The information need can be specified in the form of a search query. In the case of document retrieval, queries can be based on ful...
- <https://www.geeksforgeeks.org/what-is-information-retrieval>
What is Information Retrieval? - GeeksforGeeks
Sep 19, 2023 · Information Retrieval (IR) can be defined as a software program that deals with the organization, storage, retrieval, and evaluation of information from document repositories, particularly textual information. Information Retrieval is the activity of obtaining material that can usually be...
- <https://www.britannica.com/technology/information-retrieval>
Information retrieval | Definition, Methods, & Facts | Britannica
information retrieval, recovery of information, especially in a database stored in a computer. Two main approaches are matching words in the query against the database index (keyword searching) and traversing the database using hypertext or hypermedia links. Evolving information-retrieval technique...

Contributions

1. Include demand of a greater **range of perspectives** in search results using a **topic hierarchy**
2. Include demand of a more **fair representation** in search results using a **protected attribute**
3. Evaluate the fairness and diversity in search results

Data

The **supply of clean water** is one of the biggest social challenges of the future, on a global scale and in Germany alike.

To **prevent water scarcity** and the potential for social conflict that is linked to it, **new alliances** across research, the economy, the public sector and civil society are necessary.

<https://www.thwic.uni-jena.de/en>; 19.08.2024

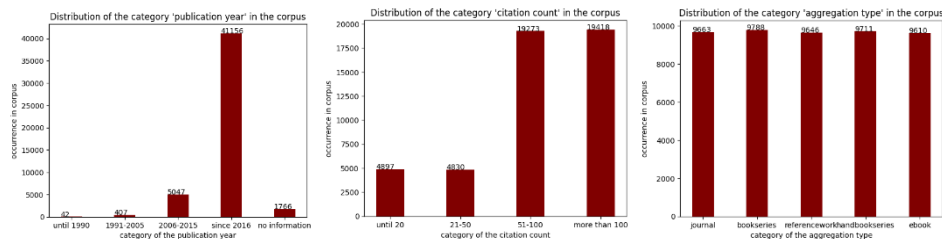
Literature Search

- crawled 48'000 publications from Elsevier/ScienceDirect
- 13 queries from chemistry
- Relevance judgments by one expert

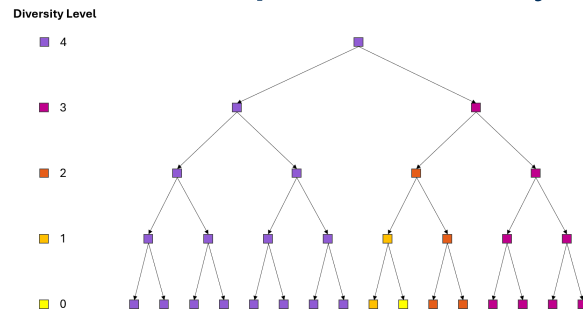
Literature Search

- crawled 48'000 publications from Elsevier/ScienceDirect
- 13 queries from chemistry
- Relevance judgments by one expert

Attributes and Categories

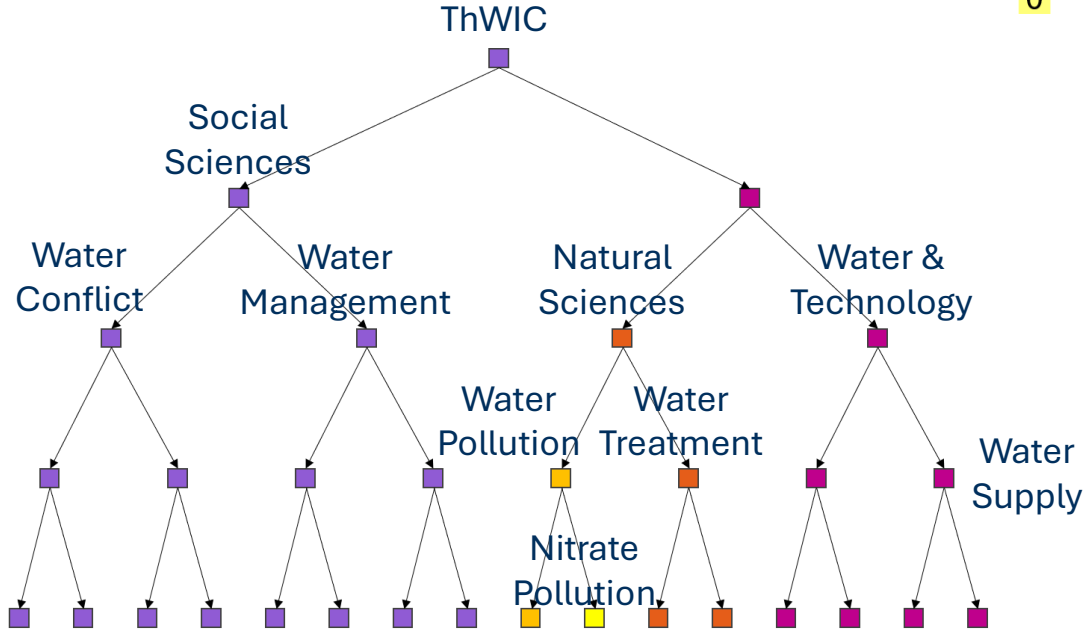
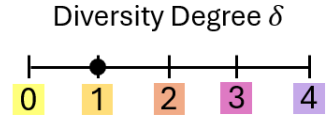


Topic Hierarchy



ThWIC Hierarchy

Diversity Level



ThWIC Hierarchy

Diversity Level

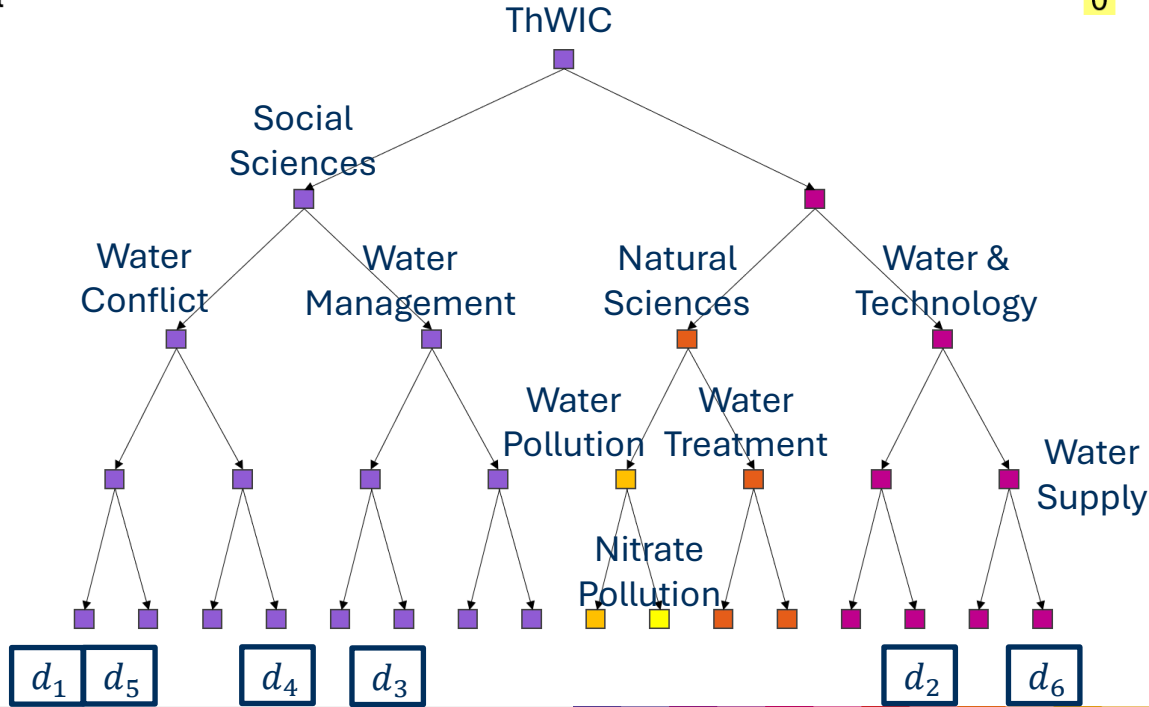
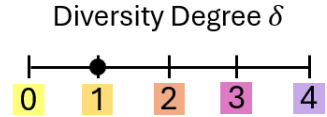
4

3

2

1

0



Diversity-Based Ranking

Diversity

The aim of diversity in search results is to present various aspects of a query.

Economic
perspective

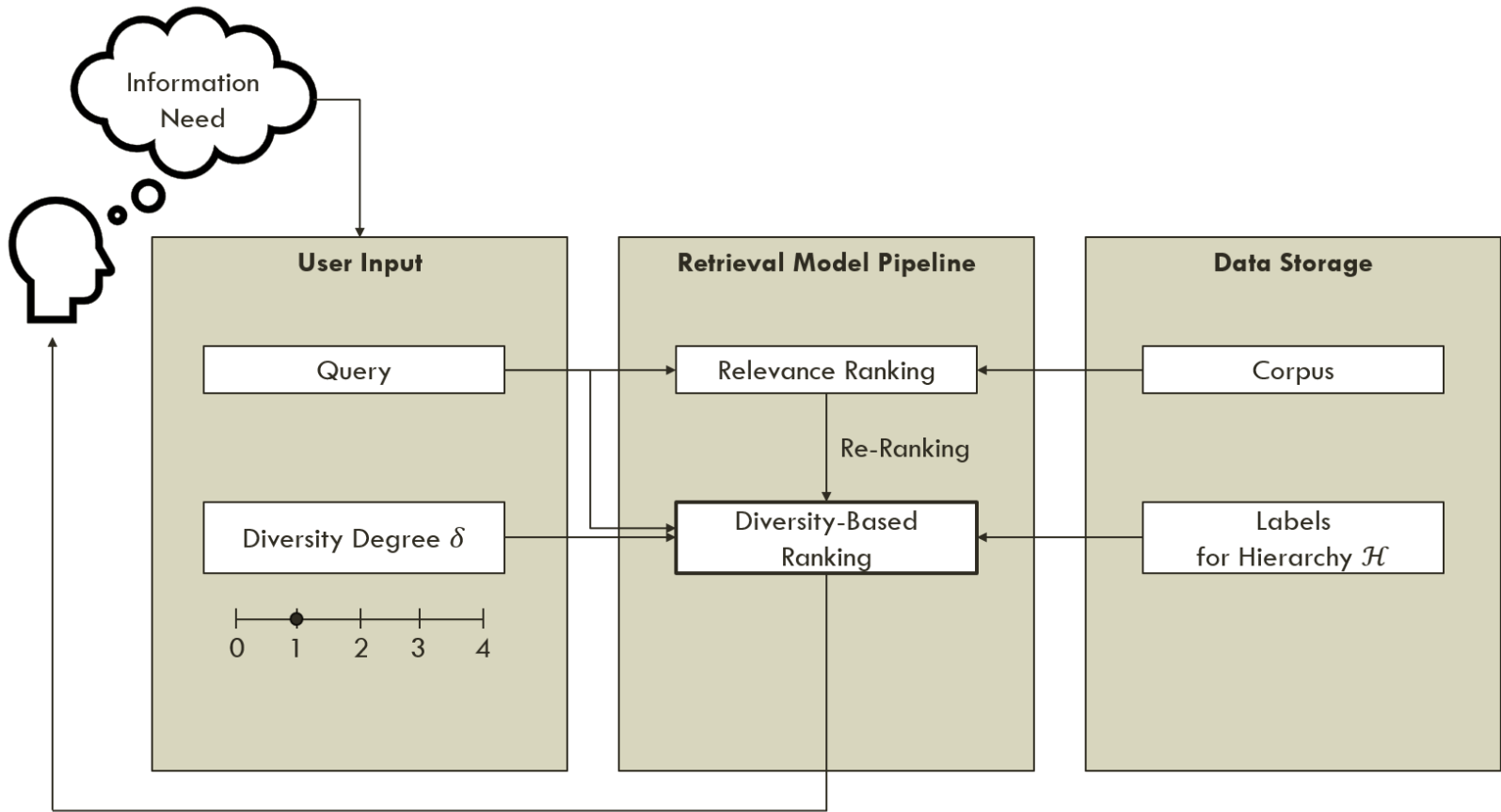
Chemical perspective

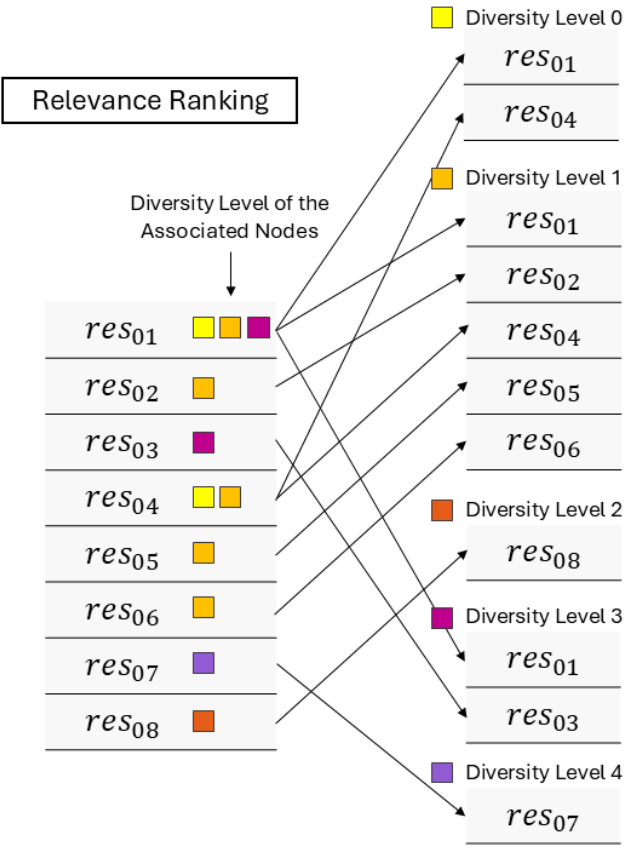


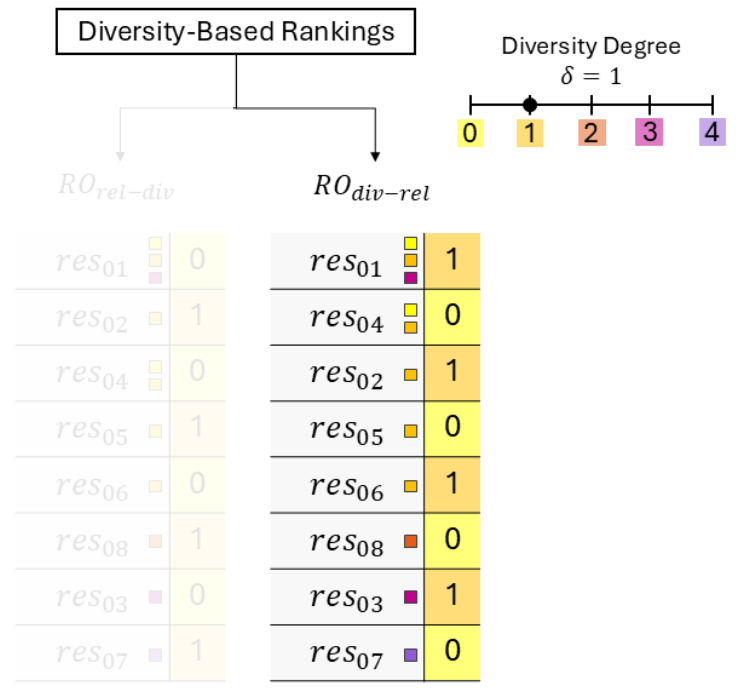
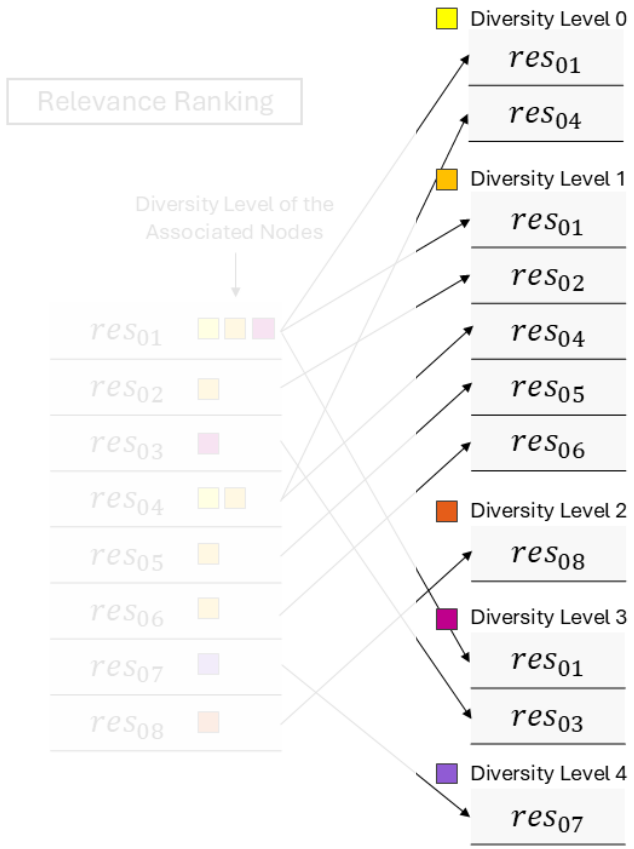
Social perspective

Geological
perspective

H Wu et al.: Result Diversification in Search and Recommendation: A Survey. arXiv, 2024







nDCHierarchical Precision

Ranking

1.	$sim(q, d_5) = 0.50$
2.	$sim(q, d_2) = 1.00$
3.	$sim(q, d_6) = 0.75$
4.	$sim(q, d_1) = 0.75$
5.	$sim(q, d_3) = 0.00$
6.	$sim(q, d_4) = 0.25$

$$\begin{aligned} DCHP'@k & \\ &= 0.50 * 1.00 \\ &+ 1.00 * 0.63 \\ &+ 0.75 * 0.50 \end{aligned}$$

weight

1.00
0.63
0.50
0.43
0.39
0.36

Most Similar Ranking

1.	$sim(q, d_2) = 1.00$
2.	$sim(q, d_6) = 0.75$
3.	$sim(q, d_1) = 0.75$
4.	$sim(q, d_5) = 0.50$
5.	$sim(q, d_4) = 0.25$
6.	$sim(q, d_3) = 0.00$

$$\begin{aligned} DCHP^*@k & \\ &= 1.00 * 1.00 \\ &+ 0.75 * 0.63 \\ &+ 0.75 * 0.50 \end{aligned}$$

$$\longrightarrow nDCHP@k = \frac{DCHP'@k}{DCHP^*@k} \longleftarrow$$

nDCHierarchical Precision

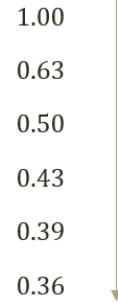
Ranking

1.	$sim(q, d_5) = 0.50$
2.	$sim(q, d_2) = 1.00$
3.	$sim(q, d_6) = 0.75$
4.	$sim(q, d_1) = 0.75$
5.	$sim(q, d_3) = 0.00$
6.	$sim(q, d_4) = 0.25$

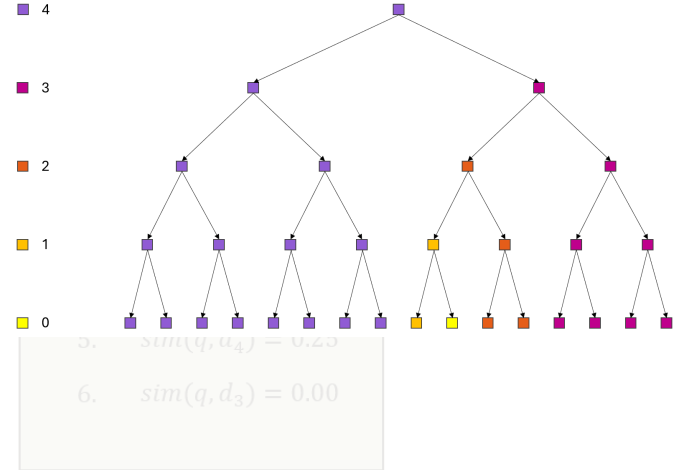
$$\begin{aligned}
 DCHP'@k & \\
 &= 0.50 * 1.00 \\
 &+ 1.00 * 0.63 \\
 &+ 0.75 * 0.50
 \end{aligned}$$

$$\longrightarrow nDCHP@k = \frac{DCHP'@k}{DCHP^*@k} \longleftarrow$$

weight



Diversity Level



$$\begin{aligned}
 DCHP^*@k & \\
 &= 1.00 * 1.00 \\
 &+ 0.75 * 0.63 \\
 &+ 0.75 * 0.50
 \end{aligned}$$

nDCHierarchical Precision

Ranking

1.	$sim(q, d_5) = 0.50$
2.	$sim(q, d_2) = 1.00$
3.	$sim(q, d_6) = 0.75$
4.	$sim(q, d_1) = 0.75$
5.	$sim(q, d_3) = 0.00$
6.	$sim(q, d_4) = 0.25$

$$\begin{aligned} DCHP'@k & \\ &= 0.50 * 1.00 \\ &+ 1.00 * 0.63 \\ &+ 0.75 * 0.50 \end{aligned}$$

weight

1.00
0.63
0.50
0.43
0.39
0.36

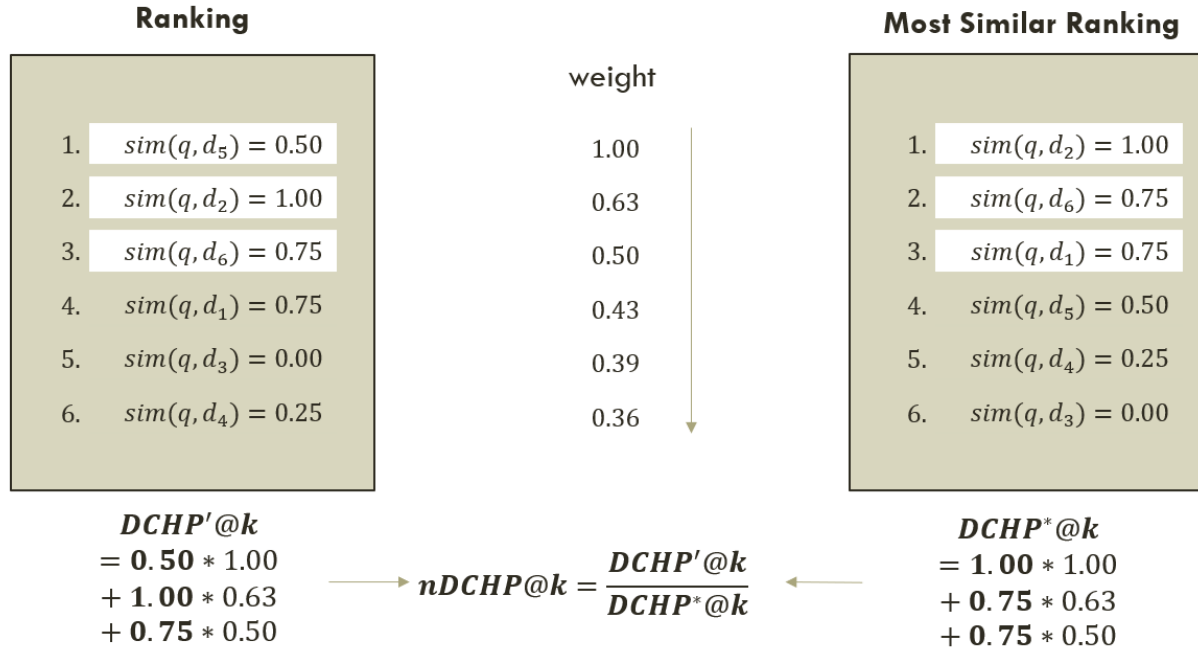
Most Similar Ranking

1.	$sim(q, d_2) = 1.00$
2.	$sim(q, d_6) = 0.75$
3.	$sim(q, d_1) = 0.75$
4.	$sim(q, d_5) = 0.50$
5.	$sim(q, d_4) = 0.25$
6.	$sim(q, d_3) = 0.00$

$$\begin{aligned} DCHP^*@k & \\ &= 1.00 * 1.00 \\ &+ 0.75 * 0.63 \\ &+ 0.75 * 0.50 \end{aligned}$$

$$\longrightarrow nDCHP@k = \frac{DCHP'@k}{DCHP^*@k} \longleftarrow$$

nDCHierarchical Precision



Results

δ		$nDCHP \downarrow$	
0	rel-ranking	0.51987	± 0.14979
	O_{spec}	0.55740	± 0.15201
	O_{div}	0.55740	± 0.15201
1	O_{spec}	0.55662	± 0.14983
	O_{div}	0.55690	± 0.15007
2	O_{spec}	0.55212	± 0.14432
	O_{div}	0.55119	± 0.14443
3	O_{spec}	0.53921	± 0.13909
	O_{div}	0.53439	± 0.13946
4	O_{spec}	0.53277	± 0.14848
	O_{div}	0.52557	± 0.15012

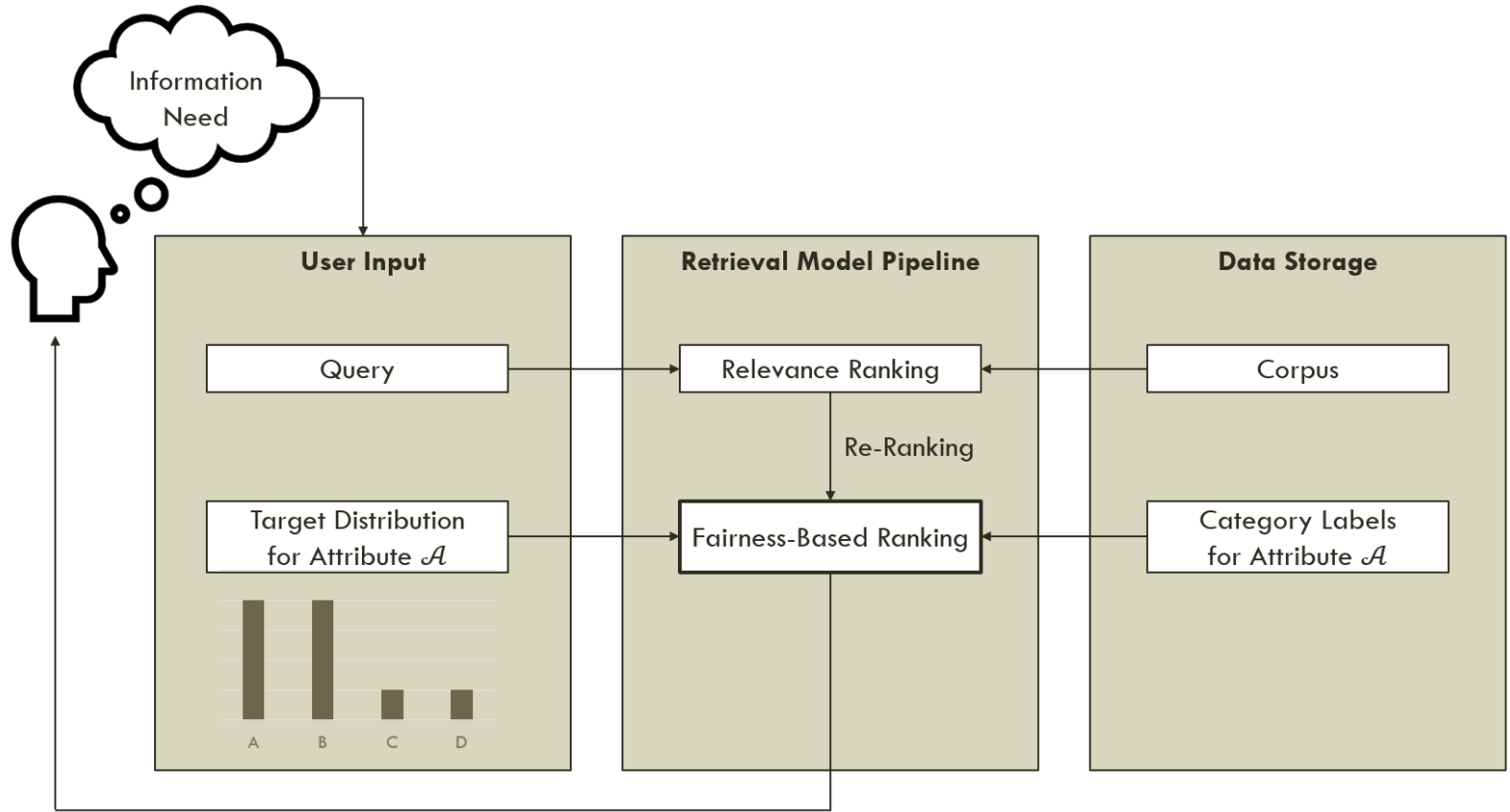
Fairness-Based Ranking

Fairness

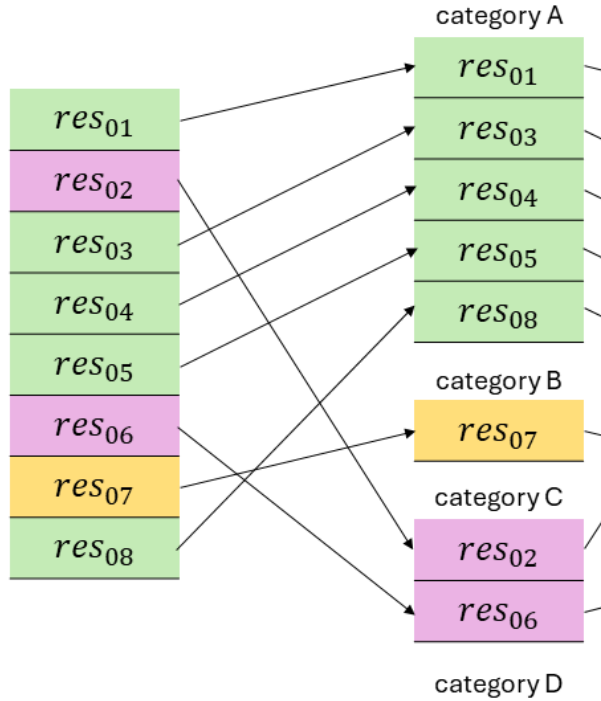
The aim of investigating algorithmic fairness is to identify, measure, and mitigate aspects that make a system unfair in a particular way.

However, “there is not one particular definition of what constitutes fairness.”

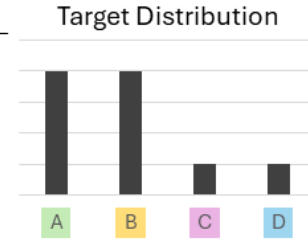
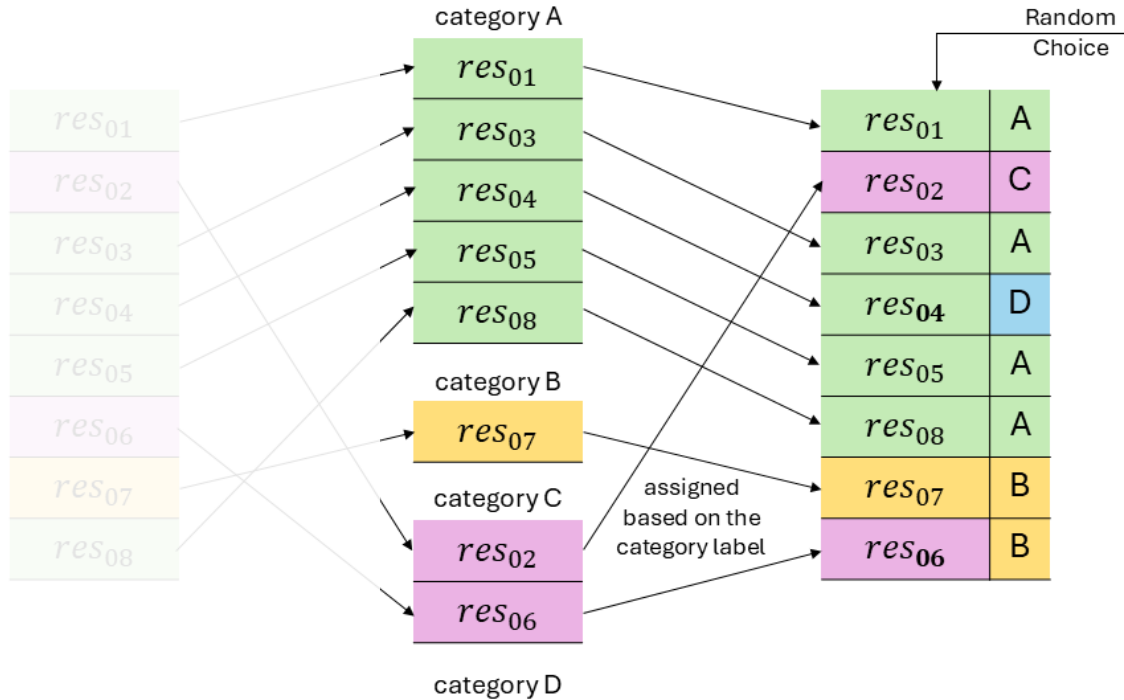
MD Ekstrand et al.: Fairness in Information Access Systems. In: Foundations and Trends in Information Retrieval 2022



Relevance Ranking



Relevance Ranking



Distribution Shift

	until 1990	1991-2005	2006-2015	since 2016	unknown
document corpus	0.001	0.008	0.103	0.851	0.037
rel-ranking@500	0.000	0.006	0.094	0.874	0.026
rel-ranking@100	0.00	0.00	0.09	0.86	0.05
rel-ranking@20	0.00	0.00	0.10	0.85	0.05
rel-ranking@10	0.00	0.00	0.00	0.90	0.10
cat-pub@100	0.00	0.03	0.30	0.64	0.03
cat-pub@20	0.00	0.15	0.35	0.45	0.05
cat-pub@10	0.00	0.30	0.30	0.40	0.00
target distribution	0.10	0.20	0.30	0.40	0.00

Table 6.5: Distribution of the categories of attribute *publication year* when re-ranking the relevance-ranking of query 4. The target distribution aims to mitigate the strong bias towards category *since 2016*.

Conclusion

User-Driven Re-Ranking for Adapting the Variety in Search Results

1. Diversity-based re-ranking
with a topic hierarchy
2. Fairness-based re-ranking
with a protected attribute
3. Extension of an
evaluation metric
for hierarchies in retrieval

User-Driven Re-Ranking for Adapting the Variety in Search Results

1. Diversity-based re-ranking with a topic hierarchy
2. Fairness-based re-ranking with a protected attribute
3. Extension of an evaluation metric for hierarchies in retrieval

Fairness-Based Ranking

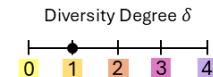


https://en.wikipedia.org/wiki/Information_retrieval
Information retrieval - Wikipedia
Information retrieval (IR) in computing and information science is the task of identifying and retrieving information system resources that are relevant to an information need. The information need can be specified in the form of a search query. In the case of document retrieval, queries can be based on full...

<https://www.geeksforgeeks.org/what-is-information-retrieval>
What is Information Retrieval? - GeeksforGeeks
Sep 19, 2023 - Information Retrieval (IR) can be defined as a software program that deals with the organization, storage, retrieval, and evaluation of information from document repositories, particularly textual information. Information Retrieval is the activity of obtaining material that can usually be...

<https://www.britannica.com/technology/information-retrieval>
Information retrieval | Definition, Methods, & Facts | Britannica
Information retrieval, recovery of information, especially in a database stored in a computer. Two main approaches are matching words in the query against the database index (keyword searching) and traversing the database using hypertext or hypermedia links. Evolving information-retrieval technique...

Diversity-Based Ranking



https://en.wikipedia.org/wiki/Information_retrieval
Information retrieval - Wikipedia
Information retrieval (IR) in computing and information science is the task of identifying and retrieving information system resources that are relevant to an information need. The information need can be specified in the form of a search query. In the case of document retrieval, queries can be based on full...

<https://www.geeksforgeeks.org/what-is-information-retrieval>
What is Information Retrieval? - GeeksforGeeks
Sep 19, 2023 - Information Retrieval (IR) can be defined as a software program that deals with the organization, storage, retrieval, and evaluation of information from document repositories, particularly textual information. Information Retrieval is the activity of obtaining material that can usually be...

<https://www.britannica.com/technology/information-retrieval>
Information retrieval | Definition, Methods, & Facts | Britannica
Information retrieval, recovery of information, especially in a database stored in a computer. Two main approaches are matching words in the query against the database index (keyword searching) and traversing the database using hypertext or hypermedia links. Evolving information-retrieval technique...

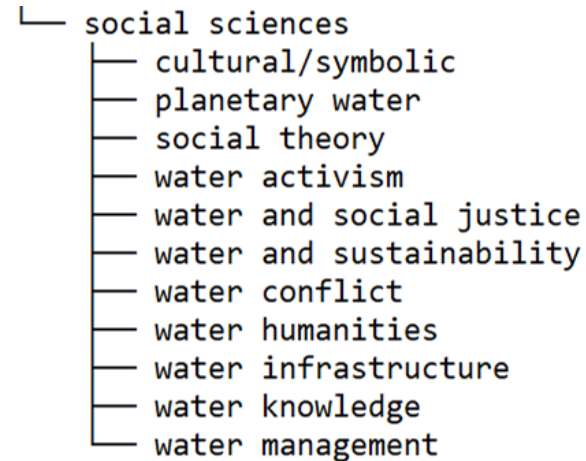
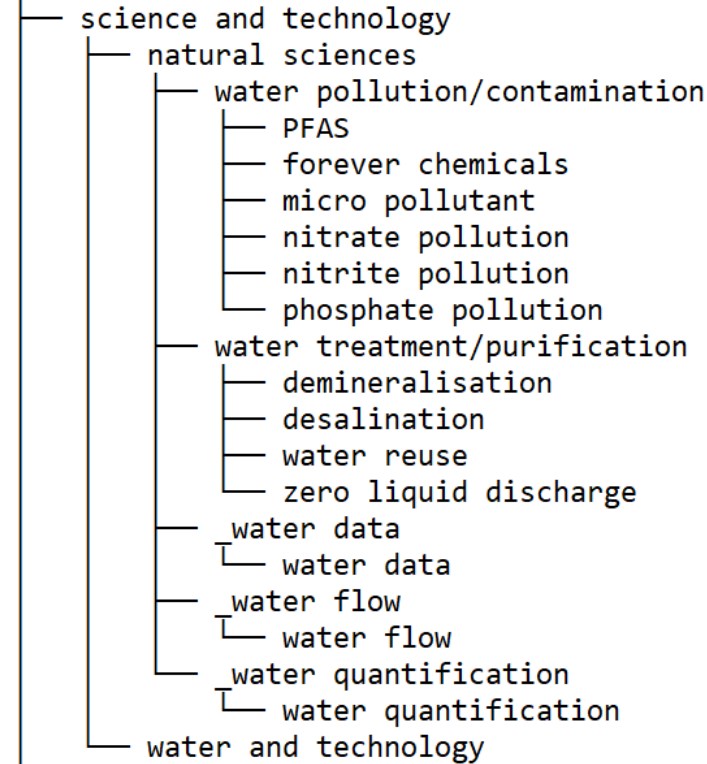
The Flute Parable

Three children want to own a flute, but only one is allowed to have it.

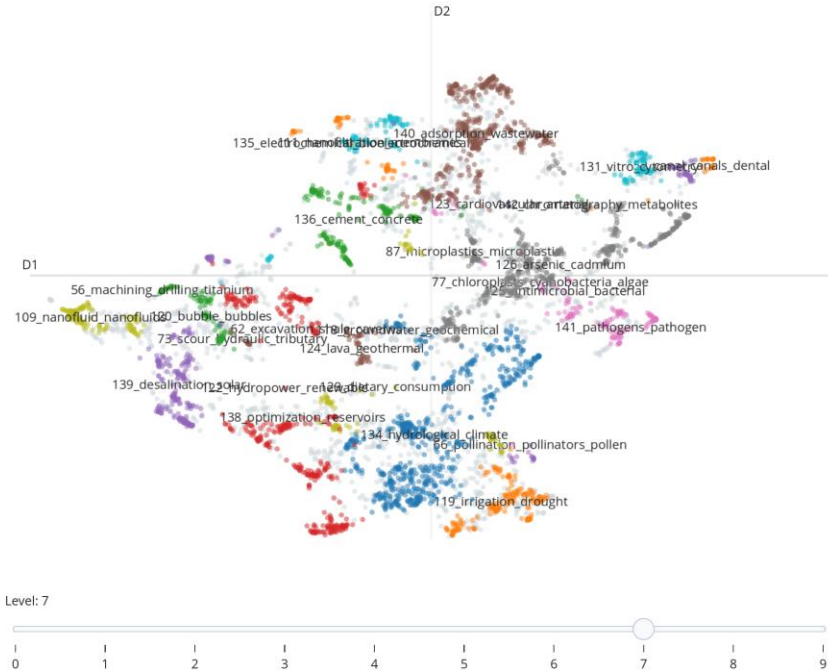
Would you give it:

- A) to the child that produced the flute,
- B) to the child that has no other toys, or
- C) to the child that can play the flute magically?

Topics in ThWIC

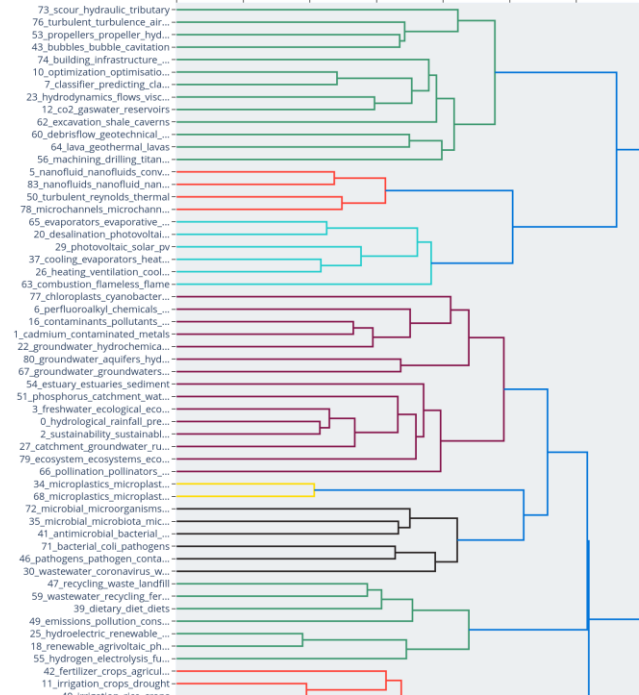


Hierarchical Documents and Topics



- 52_canal_canals_dental
- 56_machining_drilling_titanium
- 62_excavation_shale_caverns
- 66_pollination_pollinators_pollen
- 73_scour_hydraulic_tributary
- 77_chloroplasts_cyanobacteria_algae
- 87_microplastics_microplastic_wastewater
- 109_nanofluid_nanofluids_convection_nano
- 111_nanofiltration_membranes_ultrafiltra
- 118_groundwater_geochemical_methane_aqui
- 119_irrigation_drought_crops_agricultura
- 120_bubble_bubbles_hydrodynamic_flows_tu
- 122_hydropower_renewable_energy_electric
- 123_cardiovascular_arterial_vascular_dia
- 124_lava_geothermal_lavas_eruptions_geot
- 125_antimicrobial_bacterial_microbial_ba
- 126_arsenic_cadmium_toxicity_metals_poll
- 129_dietary_consumption_sustainability_s
- 131_vitro_cytometry_cells_microfluidic_c
- 134_hydrological_climate_rainfall_drough
- 135_electrochemical_bioelectrochemical_m
- 136_cement_concrete_wastewater_waste_slu
- 138_optimization_reservoirs_flow_hydraul
- 139_desalination_solar_evaporator_evapor
- 140_adsorption_wastewater_adsorbents_ads
- 141_pathogens_pathogen_wastewater_bacter
- 142_chromatography_metabolites_chromatog

Hierarchical Clustering



	journal	bookseries	referencework	handbookseries	ebook
document corpus	0.199	0.200	0.200	0.203	0.198
rel-ranking@500	0.190	0.206	0.174	0.250	0.180
rel-ranking@100	0.16	0.23	0.14	0.23	0.24
rel-ranking@20	0.30	0.25	0.20	0.10	0.15
rel-ranking@10	0.50	0.30	0.00	0.10	0.10
cat-agg@100	0.62	0.07	0.09	0.00	0.22
cat-agg@20	0.55	0.05	0.15	0.00	0.25
cat-agg@10	0.60	0.00	0.10	0.00	0.30
target distribution	0.60	0.10	0.10	0.00	0.20

Table 6.9: Distribution of the categories of attribute *aggregation type* when re-ranking the relevance-ranking of query 1. The target distribution aims to introduce a strong bias towards category *journal*.

	until 20	21-50	51-100	more than 100
document corpus	0.103	0.097	0.401	0.399
rel-ranking@500	0.098	0.088	0.408	0.406
rel-ranking@100	0.09	0.12	0.39	0.40
rel-ranking@20	0.05	0.25	0.30	0.40
rel-ranking@10	0.10	0.20	0.30	0.40
cat-cit@100	0.21	0.23	0.30	0.26
cat-cit@20	0.20	0.25	0.30	0.25
cat-cit@10	0.20	0.30	0.30	0.20
target distribution	0.25	0.25	0.25	0.25

Table 6.7: Distribution of the categories of attribute *citation count* when re-ranking the relevance-ranking of query 12. The target distribution aims at a uniform distribution of all categories.

	$AWRF_{uniform} \uparrow$	
rel-ranking	0.63905	± 0.01181
cat-pub	0.65660	± 0.01216
cat-cit	0.64513	± 0.01180
cat-agg	0.63736	± 0.01128

Chemistry data

Pipeline	nDCG	Success		Precision	
	@10	@10	@5	@10	@5
BM25→C					
cutoff@100	.810	1.00	1.00	.586	.729
cutoff@500	.763	1.00	1.00	.529	.714
DFIC→B					
cutoff@100	.341	.714	.571	.357	.400
cutoff@500	.443	.786	.786	.393	.471

Data from natural and social sciences

Pipeline	nDCG	Success		Precision	
	@10	@10	@5	@10	@5
BM25→C					
cutoff@500	.338	.462	.385	.215	.277

