

This application is a small brother of the **Collaboration Spotting project @ CERN**

- Collspotting Project leader: Dr Jean-Marie Le Goff
- 1 fellowship, 3 PhD students
- <http://collspotting.web.cern.ch>



Leveraging insight into your data network by viewing co-occurrences while navigating across different perspectives.

Graph:

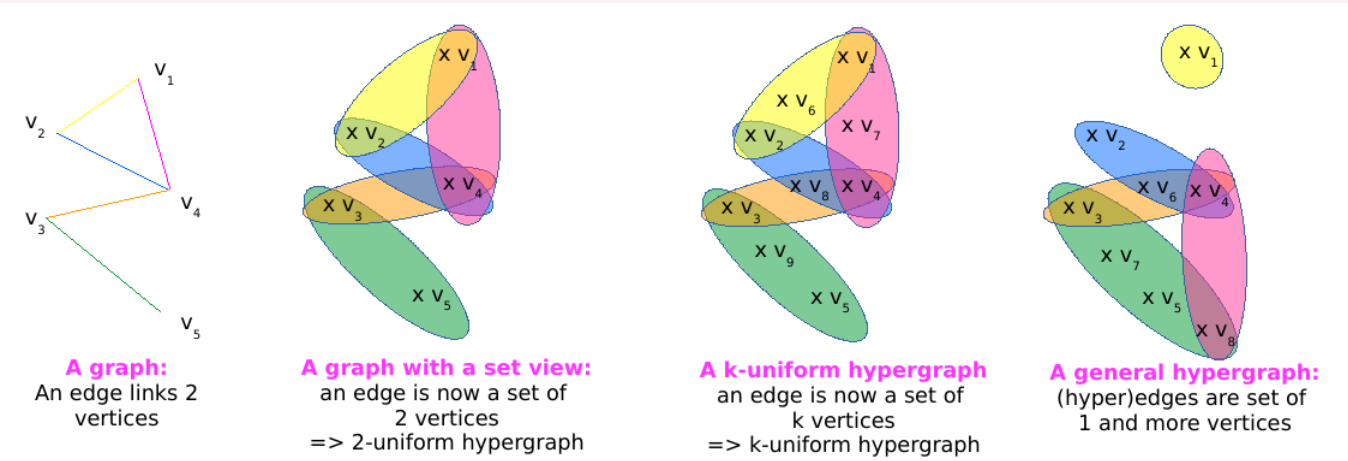
- Set of vertices and set of edges.
- An edge links two vertices : pairwise relationship.

Sets:

- Regroup elements with no repetition and no order.

Hypergraphs:

- Extend graphs.
- Allow relations between multiple vertices.
- Family of hyperedges of unempty subsets of a vertex set.

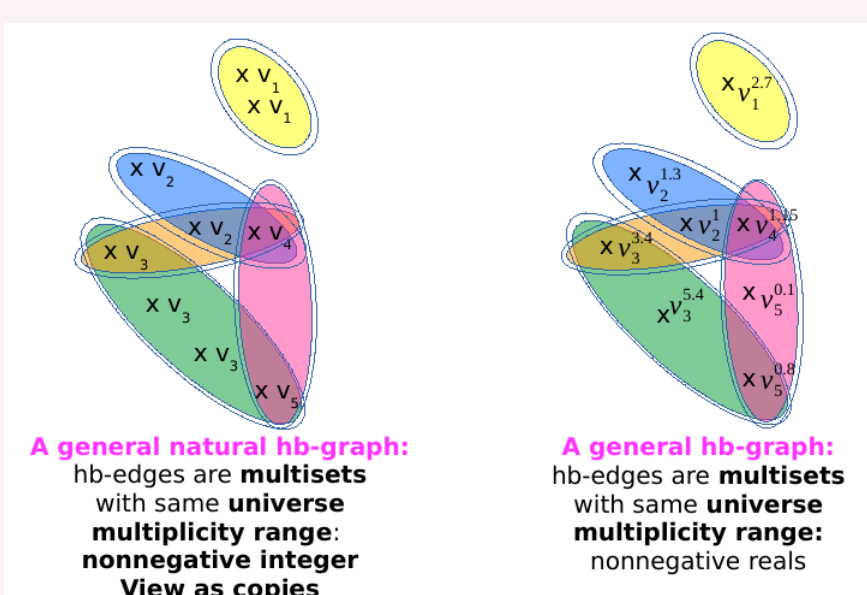


Multisets:

- Defined by a universe and a multiplicity function
- Natural multisets: multiplicity function values are integer => Collection of objects with allowed repetitions.

HyperBag-graphs (Hb-graphs):

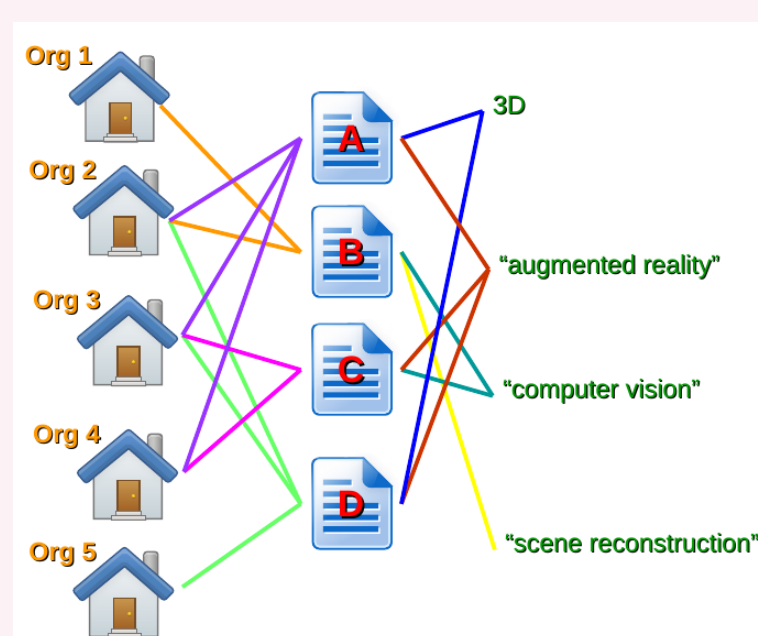
- First presented in Ouvrard et al. 2018
- Family of multisets - called **hb-edges** - of same universe called the **vertex set**
- **Natural hb-graphs**: use natural multisets => Allow duplication of elements
- **Hypergraphs** are **particular hb-graphs** with {0;1} for multiplicity function ranges.



In a Scientific Publication Database:

- Metadata store information on the structure of the DataBase
- Metadata have types, that can be used either as dimension or as reference
- Data instances attached to one type of metadata can be regrouped by using a reference => we talk about co-occurrences
- **Co-occurrences** are **n-adic relationships**.
- Co-occurrences are **multisets**, often reduced to **sets**.

On an example:

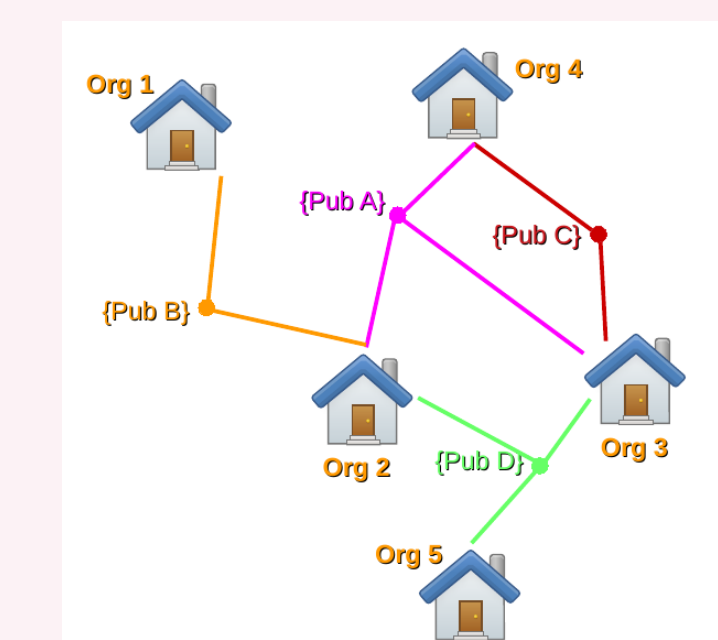


We can retrieve in a simplified way a **hypergraph**: we use as **reference** the publication and the **co-occurrences** of organisations seen as subsets of all organisations:

Reference: Publication, Facet: Organisation

Pub A	Org 2, Org 3, Org 4
Pub B	Org 1, Org 2
Pub C	Org 3, Org 4
Pub D	Org 2, Org 3, Org 5

that we can visualize with an **extra-node representation**:

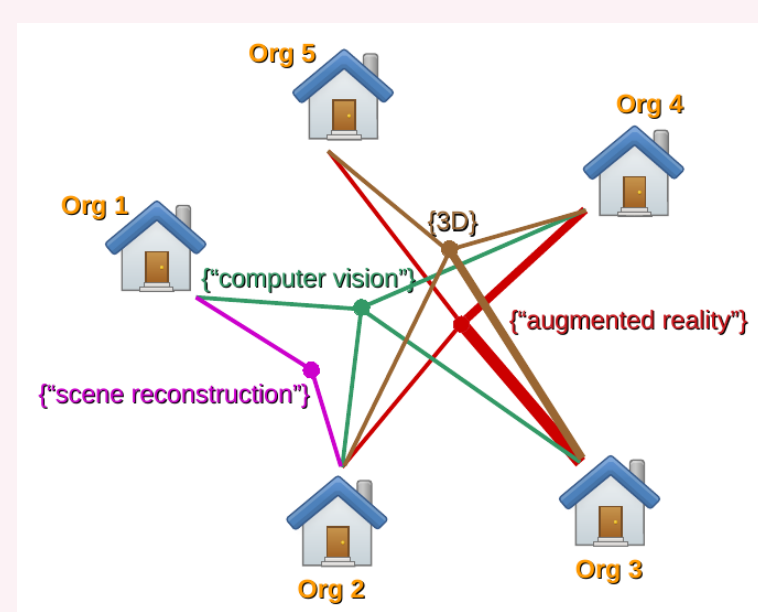


Choosing as reference keywords, we retrieve **family of multisets** of organisations, called a **hb-graph**.

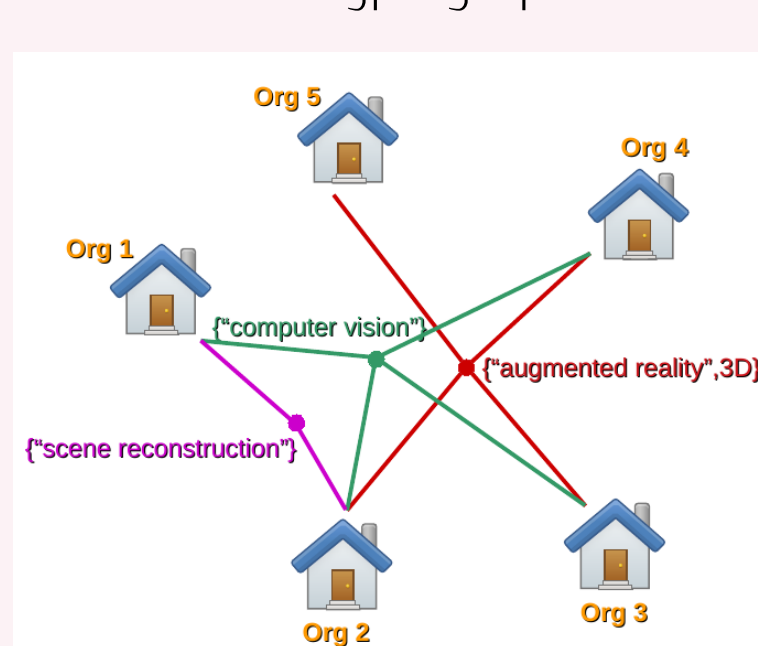
Reference: Keywords, Facet: Organisation

scene reconstruction	{Org 1 ¹ , Org 2 ¹ }
computer vision	{Org 1 ¹ , Org 2 ¹ , Org 3 ¹ , Org 4 ¹ }
augmented reality	{Org 2 ¹ , Org 3 ³ , Org 4 ² , Org 5 ¹ }
3D	{Org 2 ¹ , Org 3 ² , Org 4 ¹ , Org 5 ¹ }

Represented by a **bundled extra-node multipartite graph representation**:



Which can be reduced as a **hypergraph**:



Searching on a Scientific Publication Database

With traditional verbatim browser:

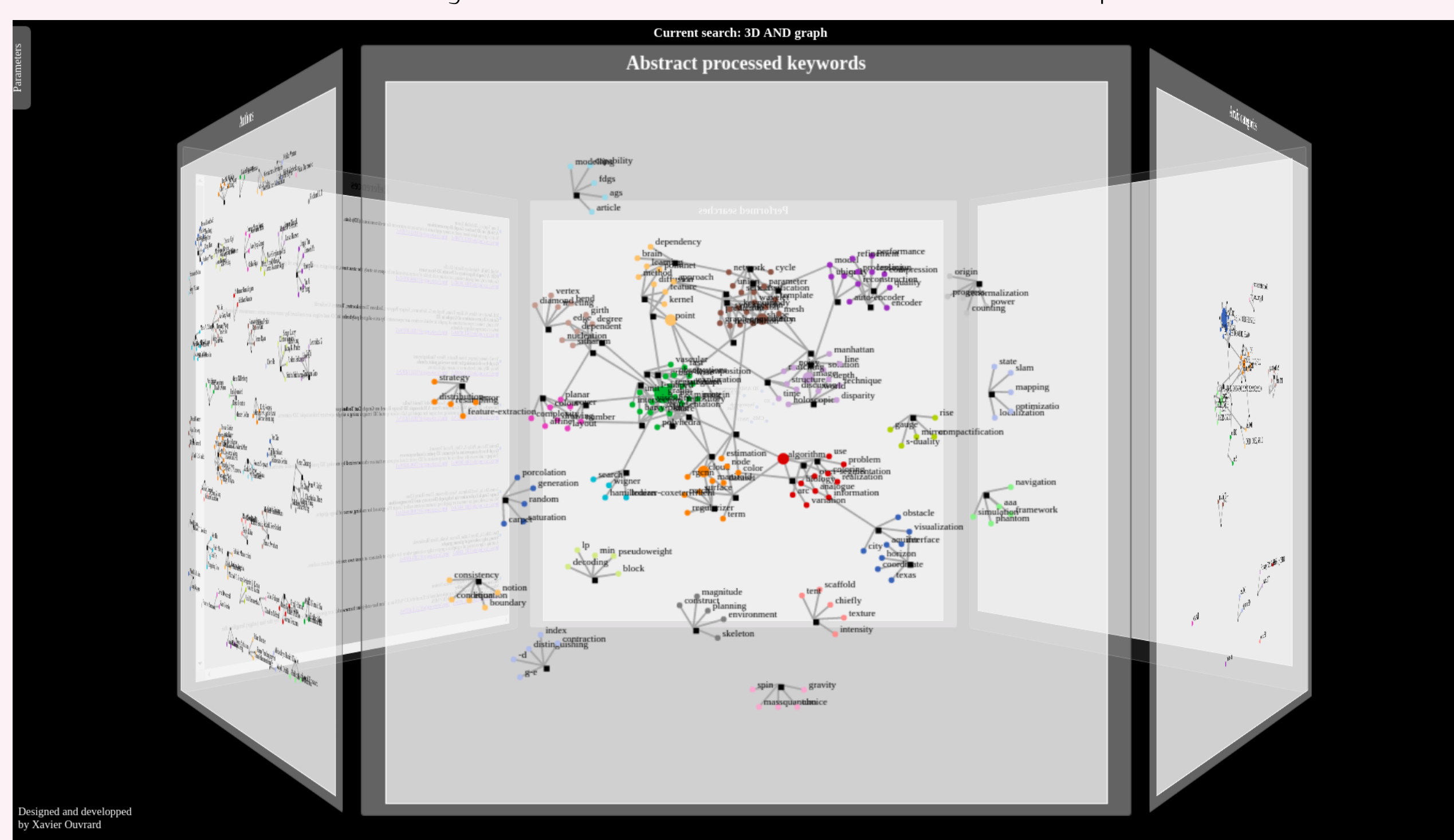
- Give **linear information**
- To **refine** information: perform a **new search**
- Making **complex query** can be **hazardous** for most people
- Accessing the **different facets** of the information space require to perform **different searches**

But in fact:

- A space of information is **multi-faceted**
- **Much more information** is available or can be extracted
- Use of natural language processing allow to extract keywords
- Hb-graphs highlight how the data instances are linked and allow additional information to be displayed

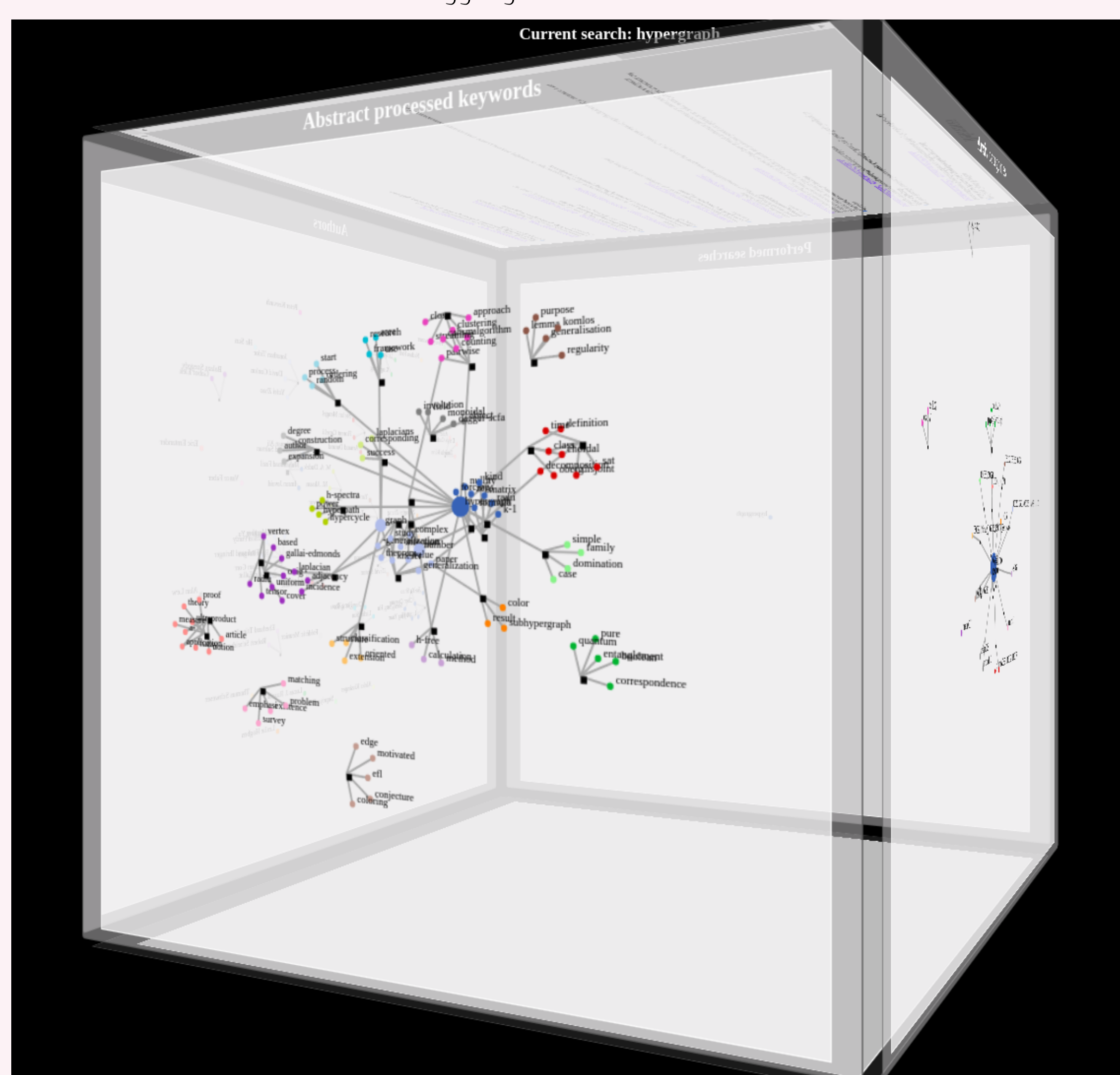


Switching between the different facets of the information space



Carousel view of the information space

Viewing to the DataEdron



The Hb-graph DataEdron: switching between the different facets of the information space

Additional information can be displayed:

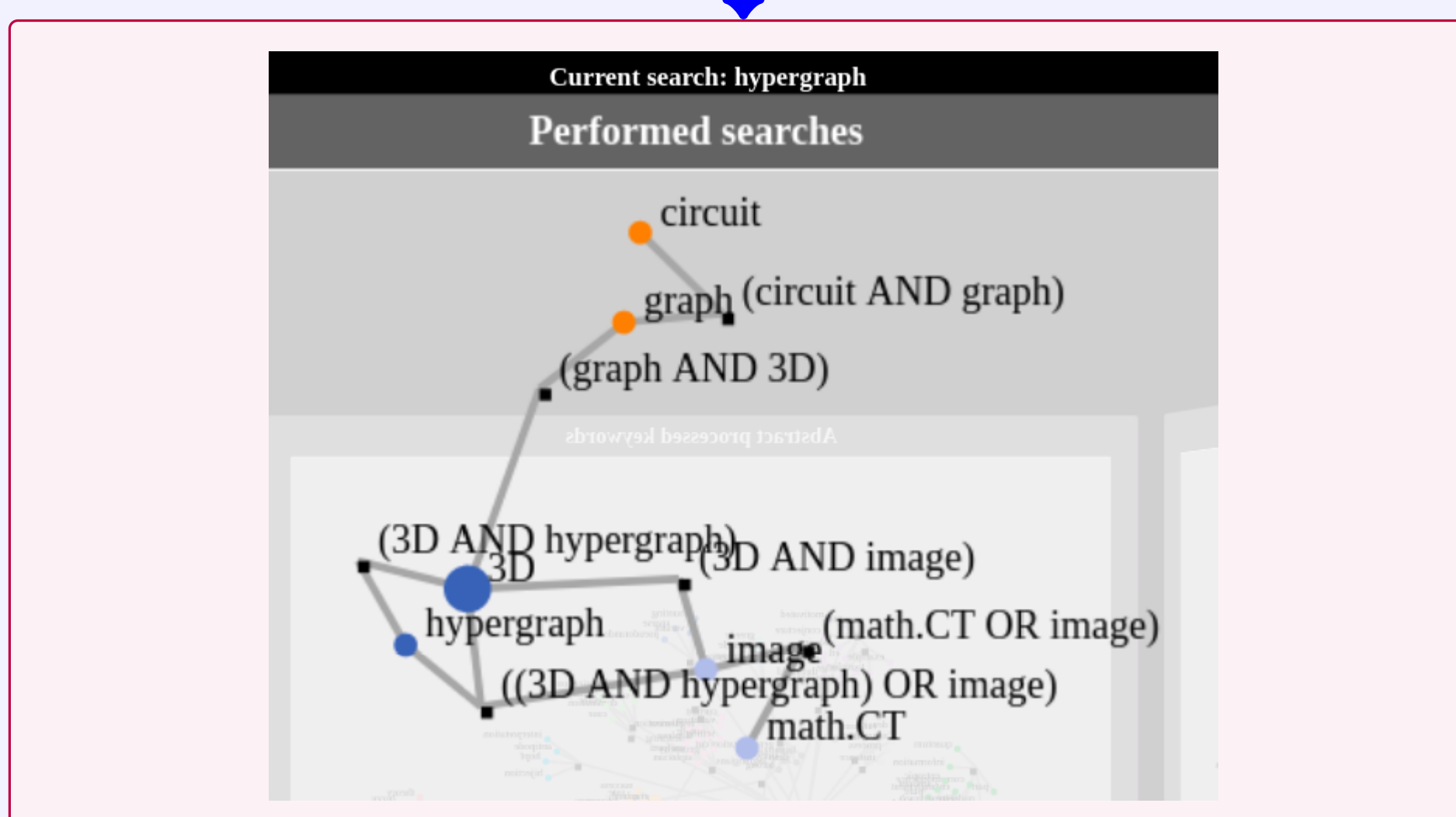
- Dblp / LinkedIn profile
- Publication abstract / full article
- Wikipedia information (for keywords)
- DuckDuckGo deambiguation and abstract...

Full interactivity of faces:

- **Highlight extra-nodes** through vertices on the same reference
- **Highlight vertices** involved in the highlighted references

How can we perform search?

- **Traditional text field search**
- Then queries can be **built visually** using vertices of the hb-graphs
- **AND, OR, NOT** possibilities
- The **graph of search** can be explored and reused
- **Possibility of merging** different searches on a single graph



Recomputing facets

Which database can we search on?

- Arxiv
- ... but can be applied to any databases.

How do we proceed?

- All the **queries** are made **online**
- Everything is **processed online**, including keywords if they are not provided
- **No intermediate storage**

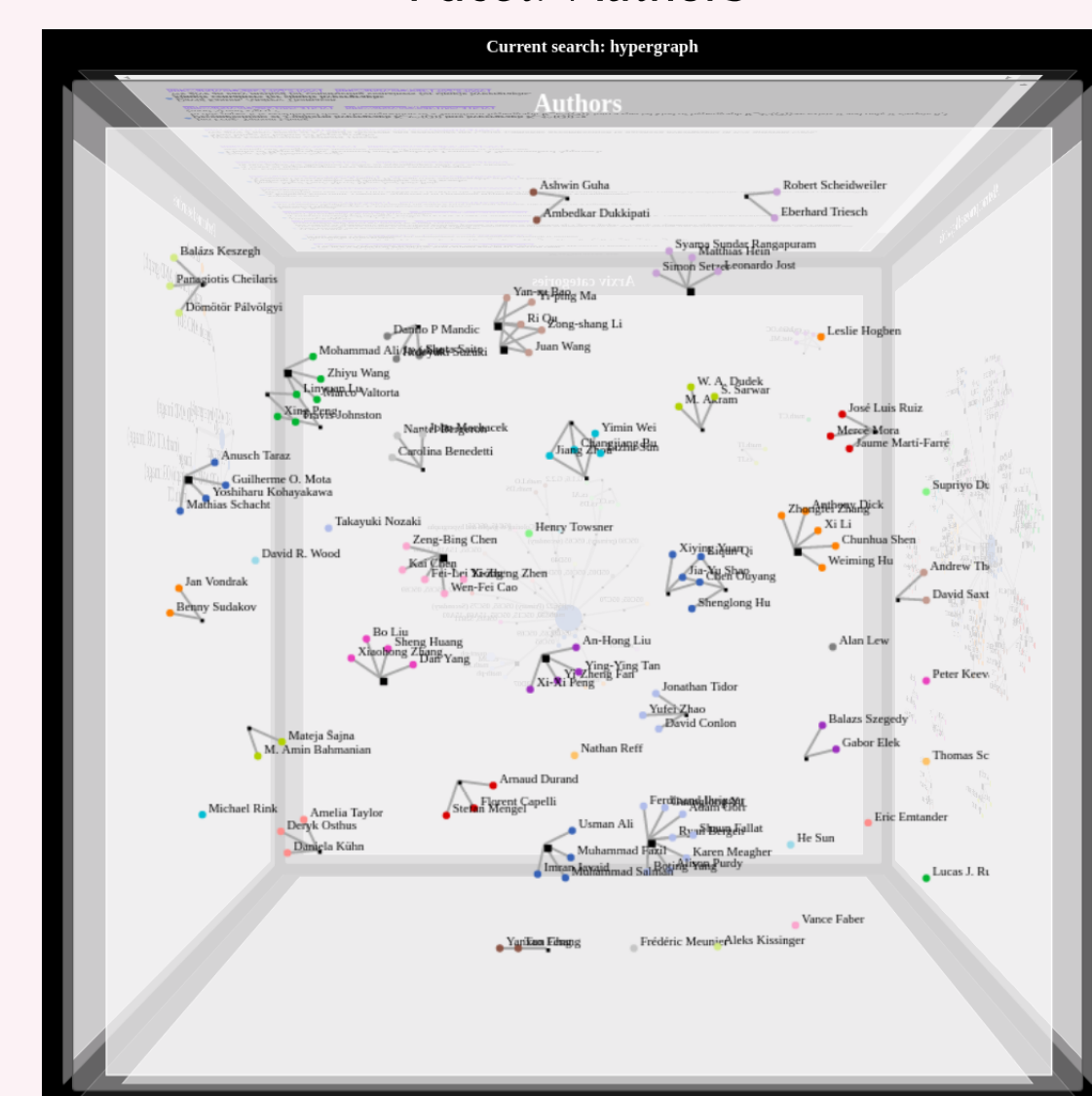
Which machine learning is used?

- **Natural language processing:**
 - It allows the extraction of the keywords from abstract
 - Stop words are removed
 - Lemmatization is made: only nouns are kept
 - Singularisation is made
 - Keywords are ranked using tf-idf
 - Tf: term frequency of the term in the document
 - Idf: inverse document frequency of the apparition of the term in the set of documents
- **Clustering**
 - The aim is to regroup vertices by communities
 - Vertices that are more connected than in a random graph are gathered
 - Use the modularity of Newman
 - The Louvain community detection algorithm is fast and efficient
- **Layout embedding**
 - We use a force directed algorithm
 - It attracts vertices that are connected and repels the ones that are disconnected
 - Works well on small graphs

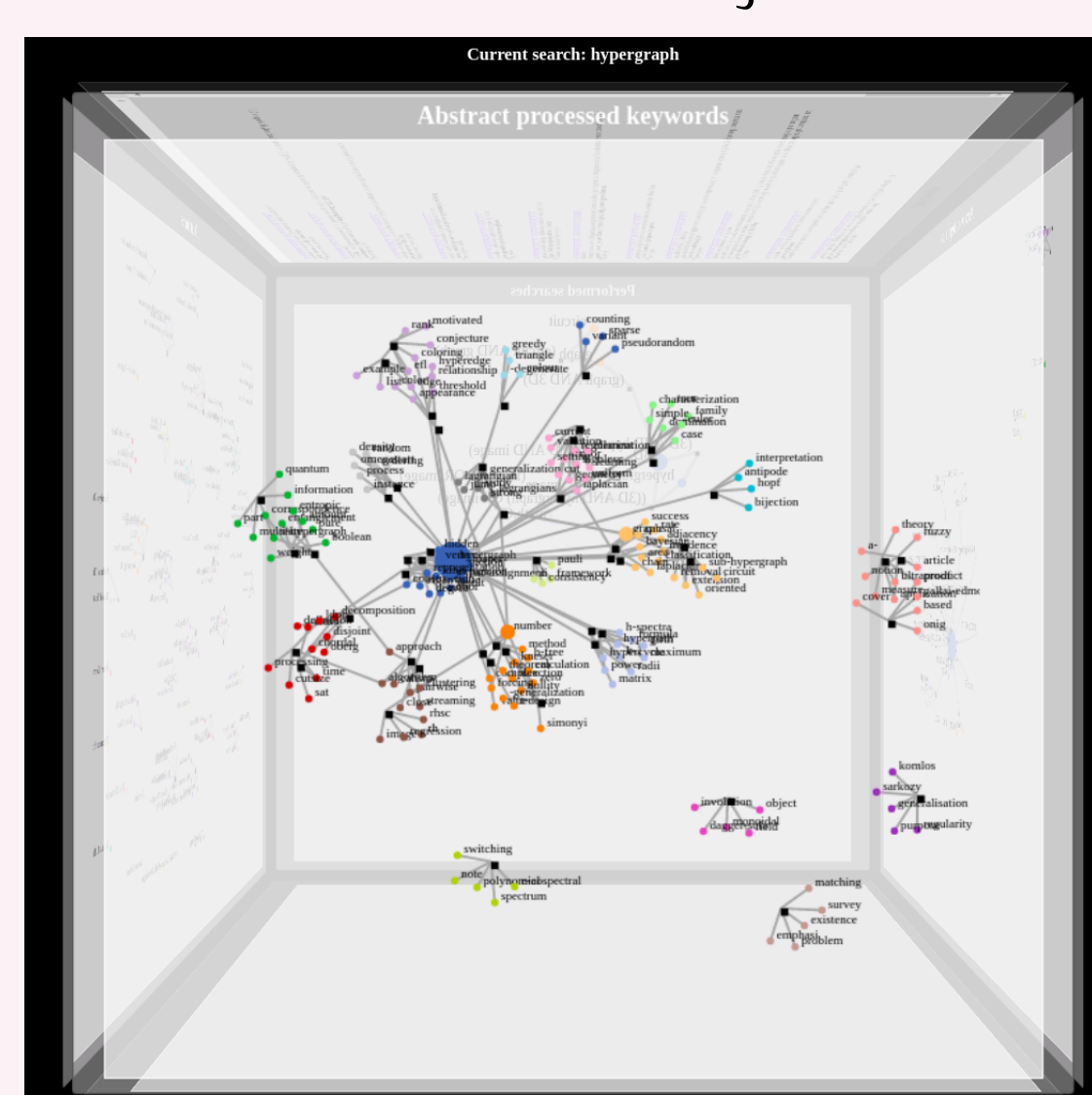
Search on the keyword: hypergraph

- **Three facets** of the information space
- **Linked by the common references**

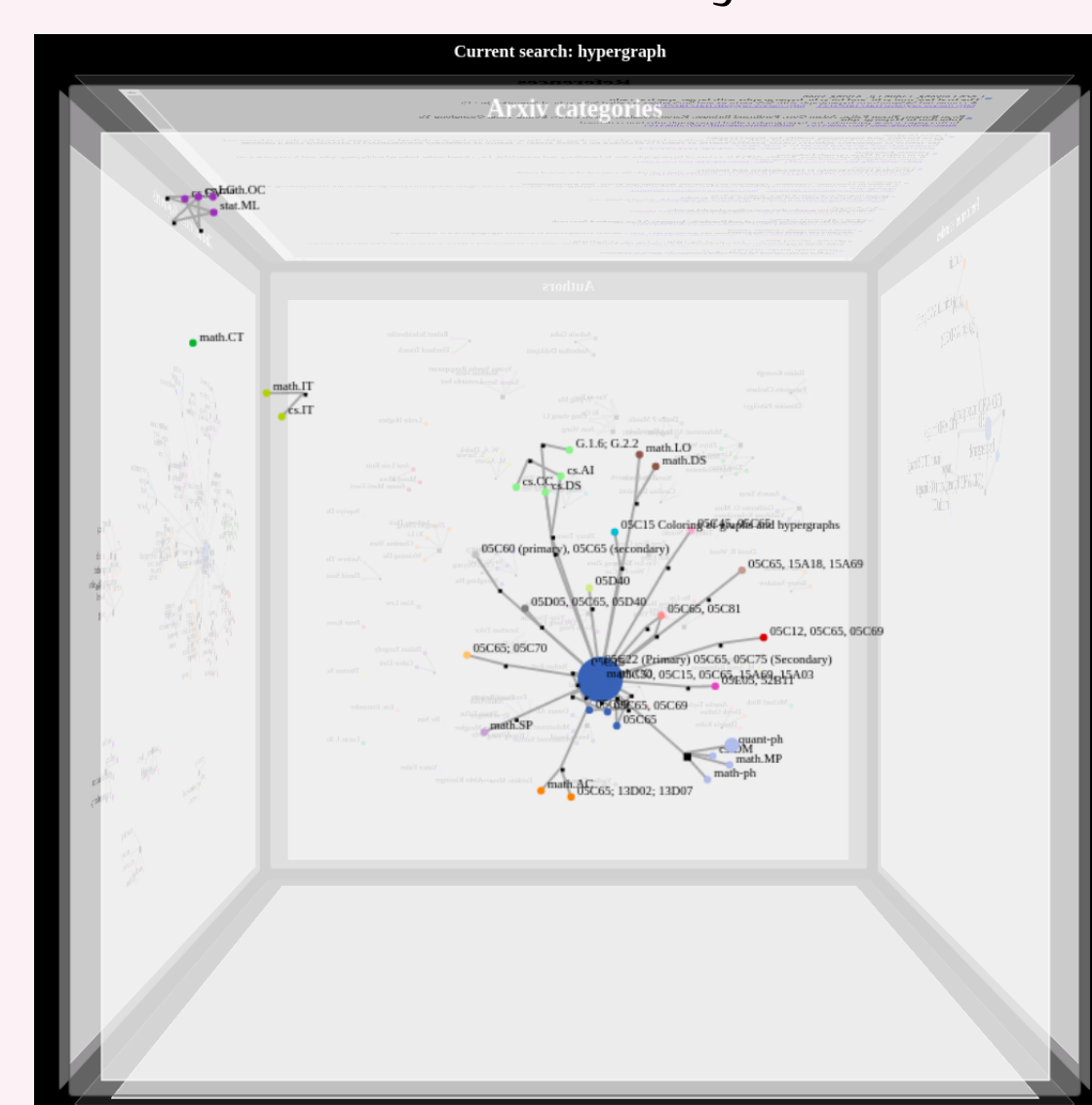
Facet: Authors



Facet: Processed keywords



Facet: Arxiv Categories



View the different facets

What are the research challenges solved?

- Modeling requires a strong framework
 - Particularly the switching of references is challenging
- Scaling up representations (for some applications):
 - Requires to fasten the computation
 - Only appropriated for some applications
- Finding important part of the representations:
 - A **diffusion process** has been proposed that allows to retrieve information on vertices and hb-edges.
 - **Weighted aggregation ranking**

What the future work is?

- Find an efficient recommendation system based on the browsing experience of the user
- Data linkage of multiple queries provenance
- Having more insights require fast extraction of information from the documents itself

Do you want more information?

- Find full details on Arxiv:1905.11695 and more on: <http://www.infos-informatique.net>
- Contact info: xavier.ouvrard@cern.ch