

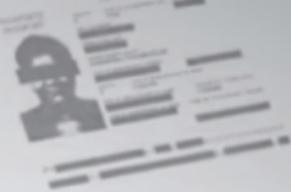






# PANAMA PAPERS

# MOSSACK FONSECA



	1999	2000	2001
Assets	Likely to be delivered	Revenue	Commitments and contingent liabilities
(D'000)	AED'000	AED'000	AED'000
444,750	1,523,153	298,950	
357,647	11,078	60,342	
804,775	8,843	57,379	
56,995	-	43,639	
	1,543,072	448,230	
2,256,504	70,418	18,362	
82,633	167,031	273,016	
107,863	4,447	65,040	
171,779	-	91,872	
45,385	1,301,126	-	
	1,543,072	448,230	
13 7,664,167			

[John Doe]

Hello. This is John Doe.  
Interested in data?

[Süddeutsche Zeitung]

[John Doe]

We're very interested.  
There are a couple of conditions. My life is in danger.  
We will only chat over encrypted files.  
No meeting, ever.  
The choice of stories is obviously up to you.

[Süddeutsche Zeitung]

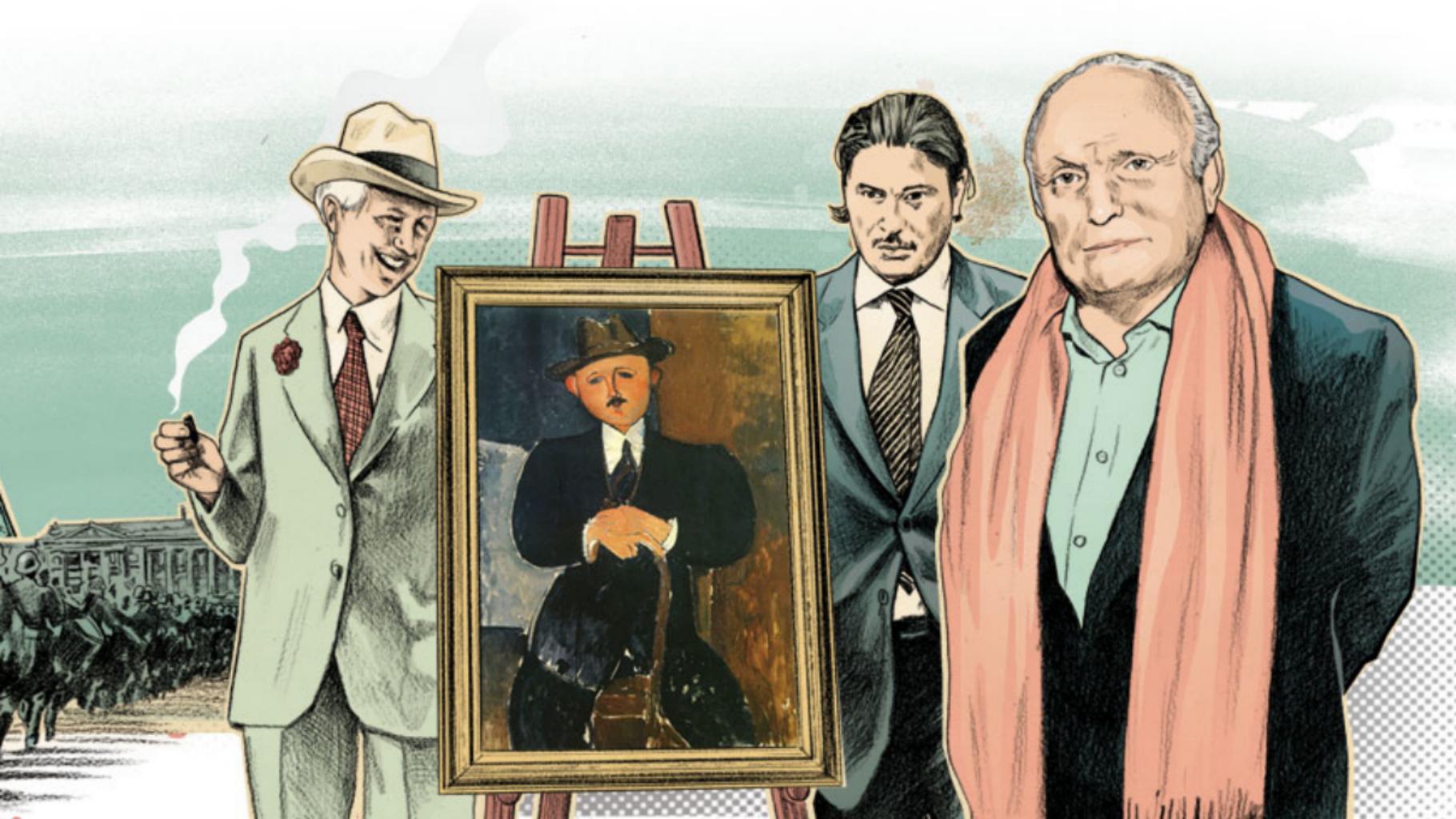
Why are you doing this?

[John Doe]

I want to make these crimes public.









# FIFA



SAY NO TO RACISM

CAPE TOWN

# The scale of the leak

Volume of data compared to previous leaks

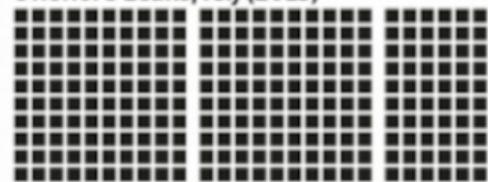
**1,7 GB**

Cablegate/Wikileaks (2010)



**260 GB**

Offshore Leaks/ICIJ (2013)



**4 GB**

Luxemburg Leaks/ICIJ (2014)



**3,3 GB**

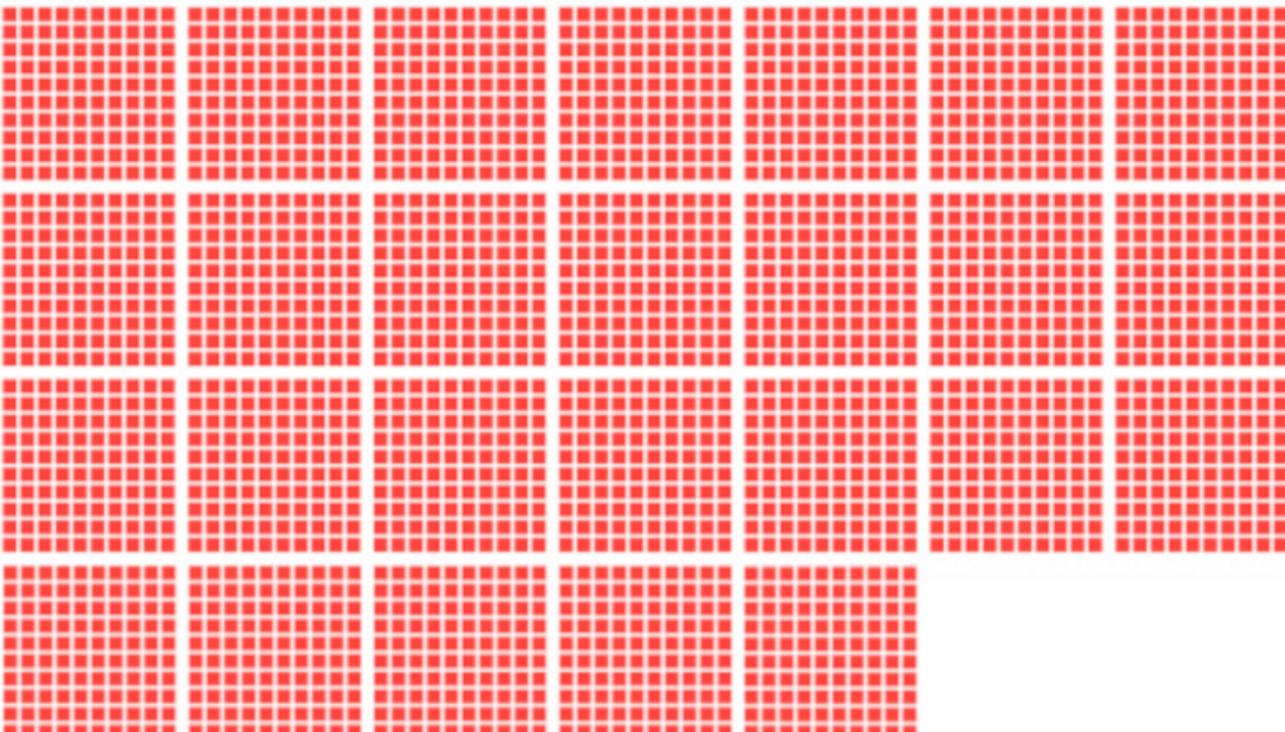
Swiss Leaks/ICIJ (2015)

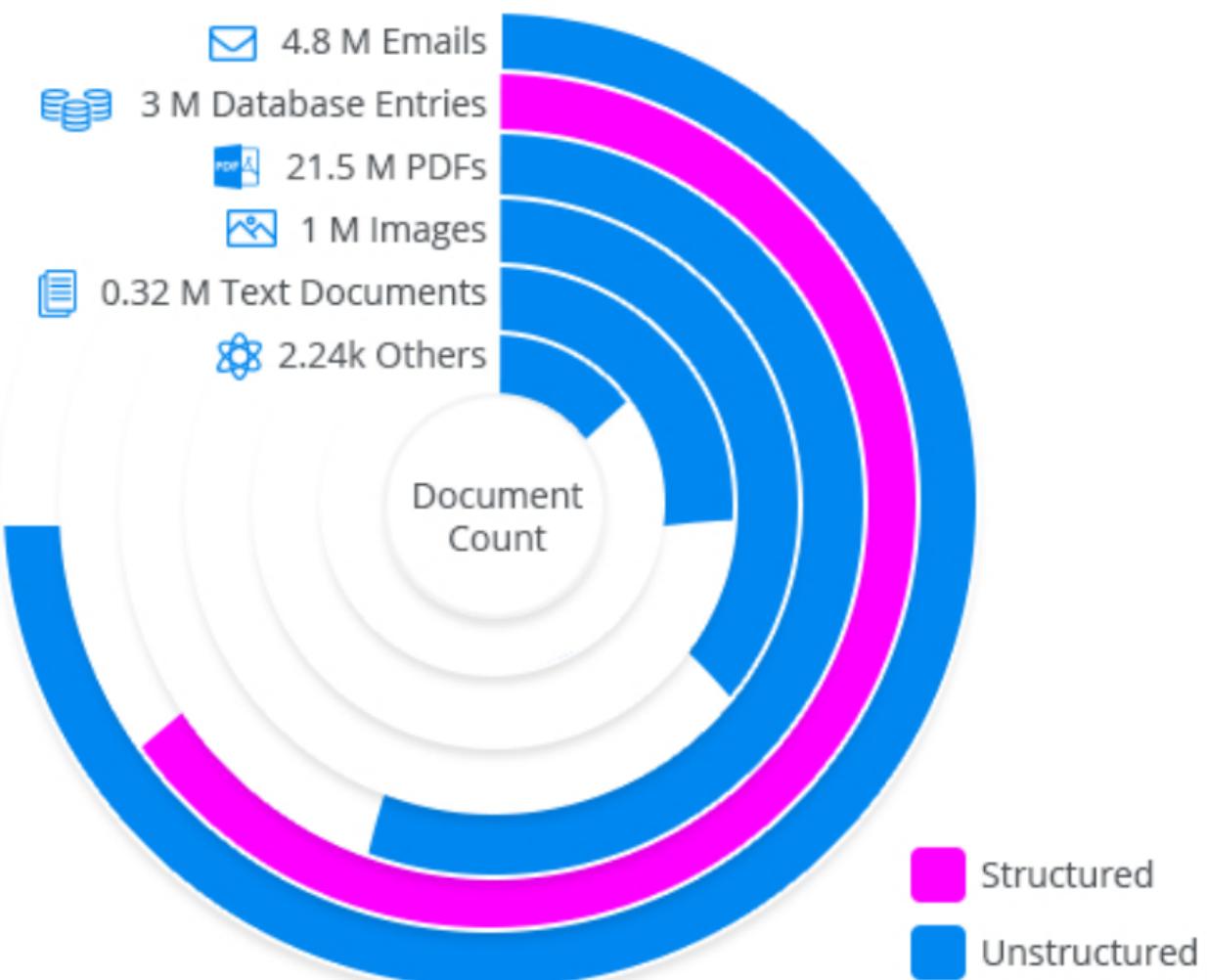


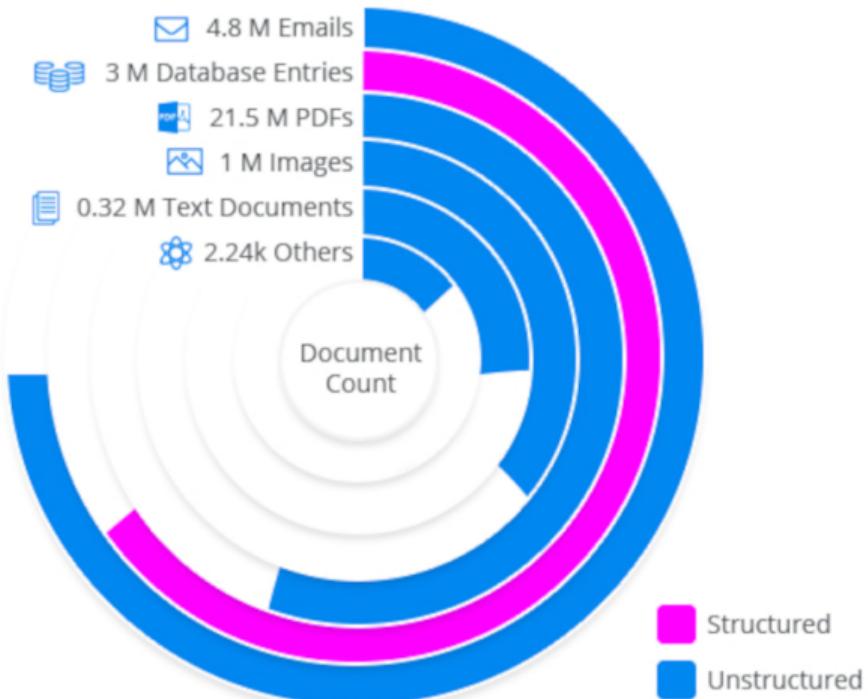
**≈ 2,6 TB**

Panama Papers/ICIJ (2016)

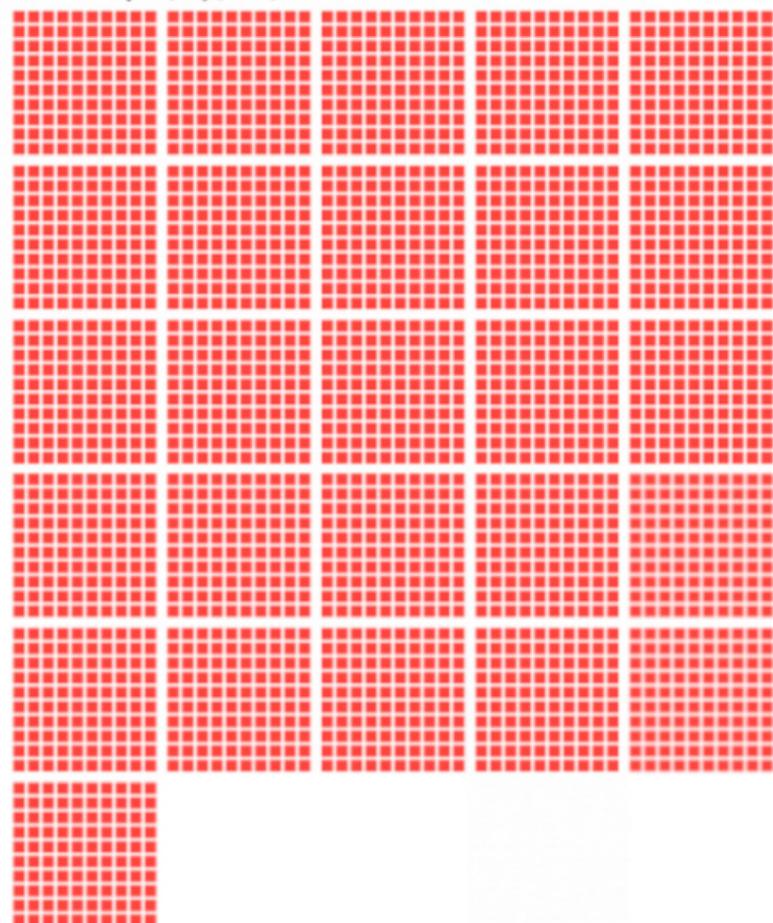
■ = 1 GB

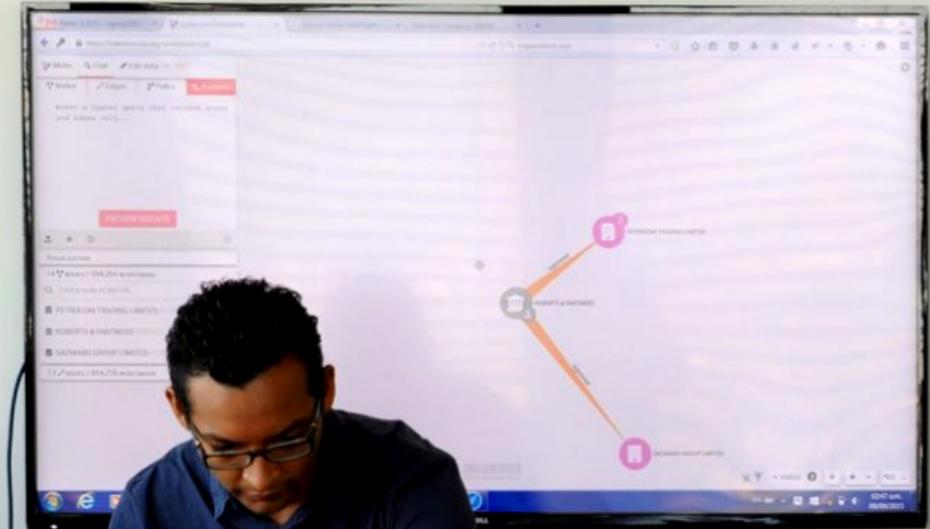


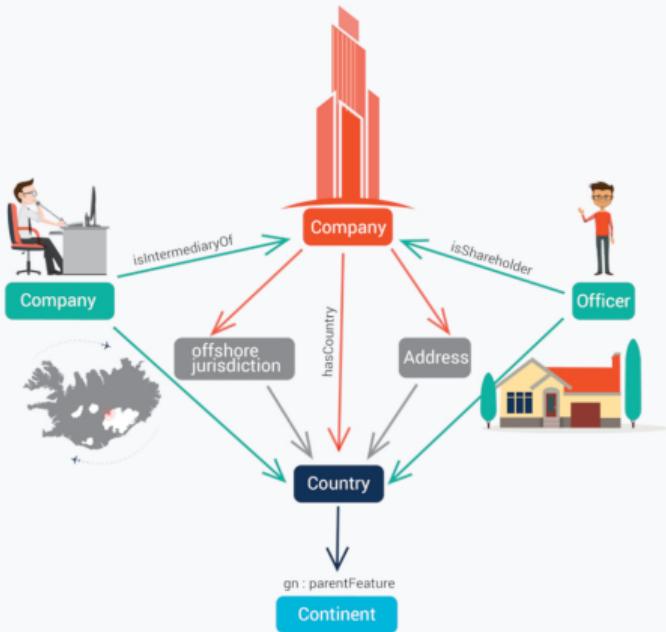




≈ 2,6 TB  
Panama Papers/ICIJ (2016)



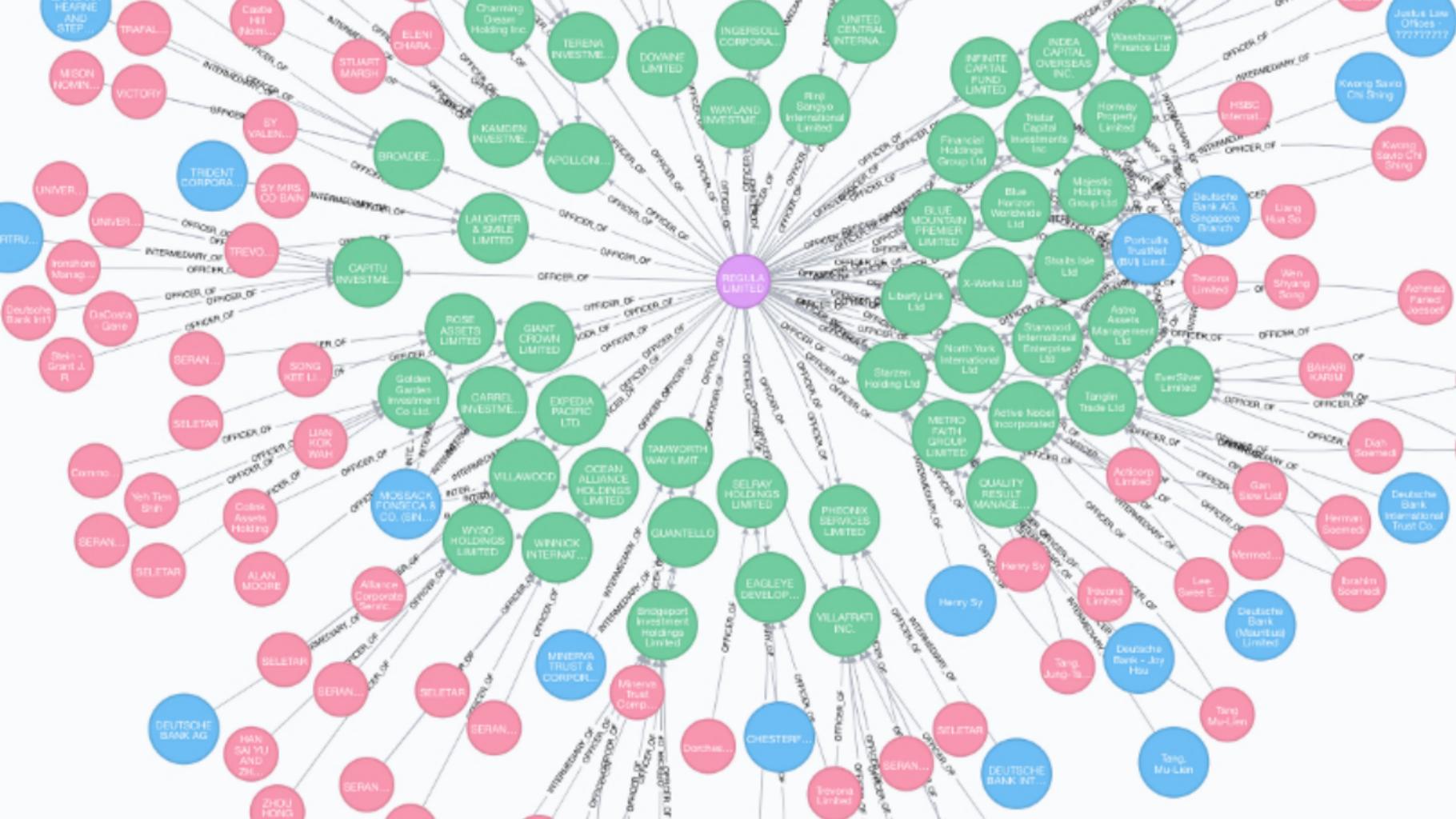






(:Address)--(:Officer)--(:Entity)--(:Intermediary)--(:Address)

```
MATCH (:Address)--(:Officer)--(e:Entity)
MATCH (e)--(:Intermediary)--(:Address)
MATCH (e)--(:Address)
WHERE o.name CONTAINS 'Emma Watson'
```



# CLUSTERING GRAPHS

Applying a Label Propagation Algorithm to  
Detect Communities in Graph Databases

Name Surname

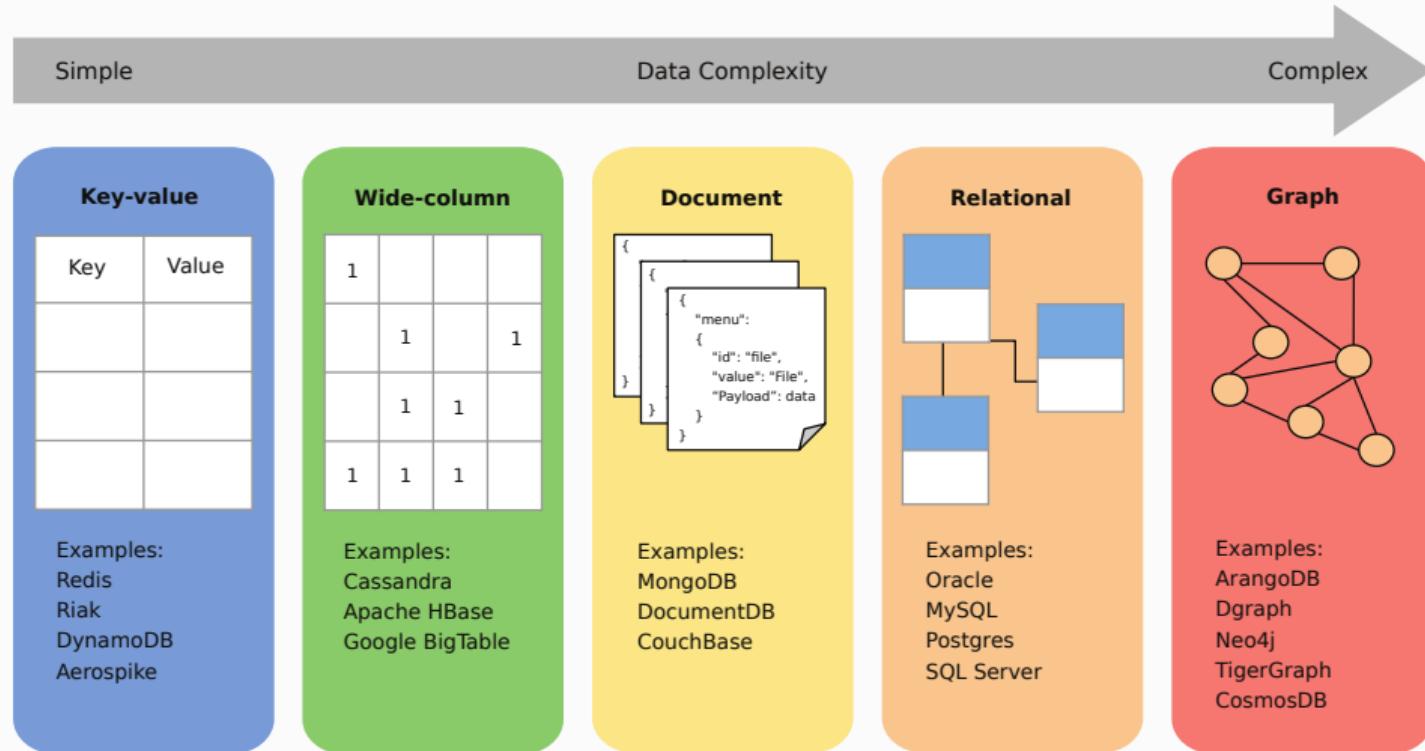


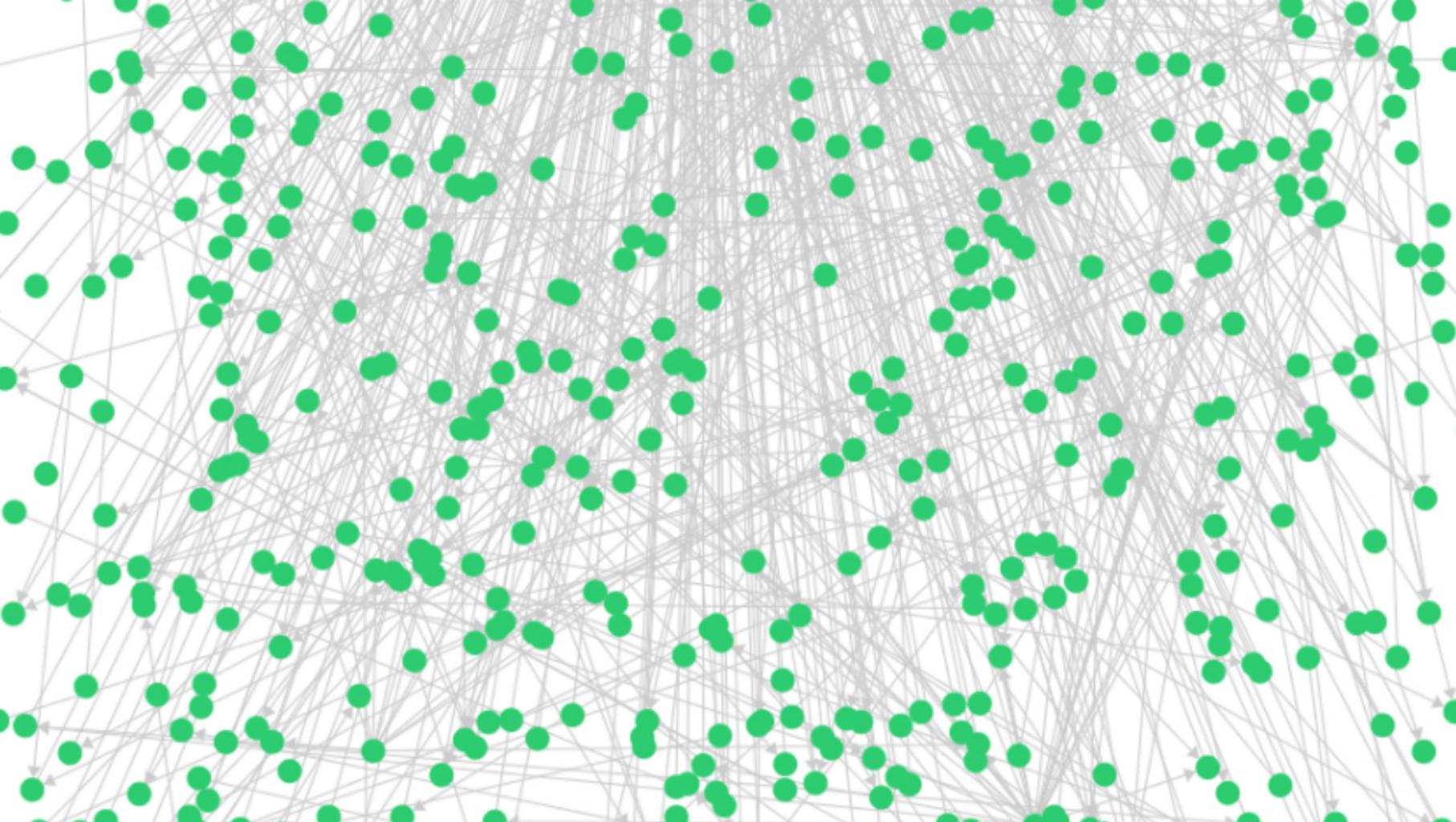
## The work done

---

1. **Literature review** of graph theory, graph databases and clustering algorithms.
2. **dblp.org Dataset download, conversion & import** in ArangoDB Graph DBMS.
3. **Data transformations** to obtain vertices, edges and the complete graph.
4. **Community Detection Algorithm application** on the graph for clustering.
5. **Web Application development** to display the results of the clustering.

# Graph databases







# The dataset: dblp.org

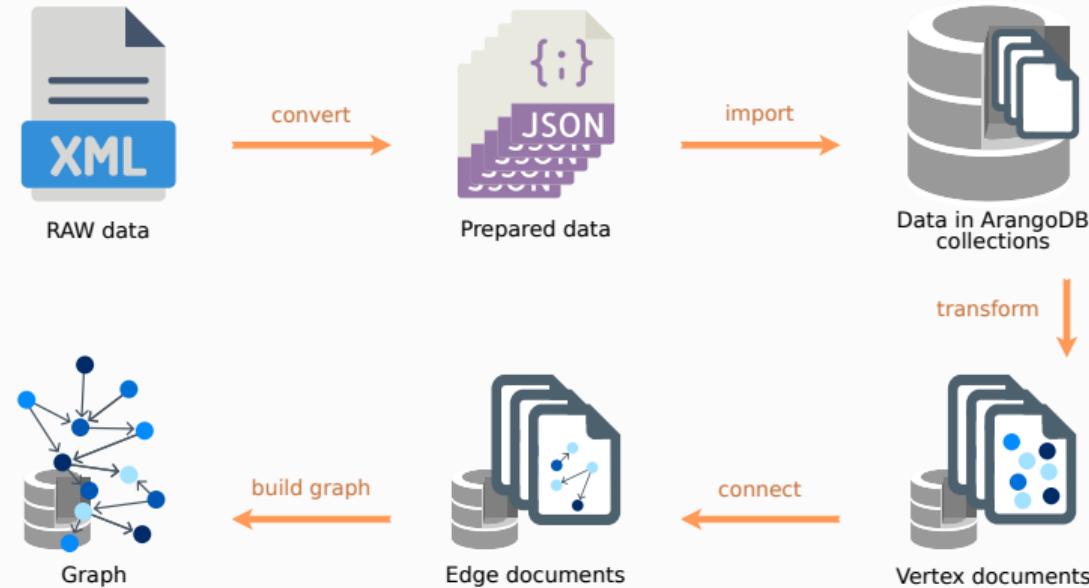
The screenshot shows a web browser window with the title bar "Index of /xml/release". The address bar contains the URL "https://dblp.org/xml/release/". The main content area displays a table titled "Index of /xml/release" with columns: Name, Last modified, Size, and Description. The table lists several XML files and their metadata. The first row is a "Parent Directory". Subsequent rows show files like "dblp-2021-09-01.xml.gz" and "dblp-2021-09-01.xml.gz.md5" with various modification dates and sizes.

Name	Last modified	Size	Description
<a href="#">Parent Directory</a>	-	-	
<a href="#">dblp-2021-09-01.xml.gz.md5</a>	2021-09-01 23:44	57	
<a href="#">dblp-2021-09-01.xml.gz</a>	2021-09-01 23:44	623M	
<a href="#">dblp-2021-08-01.xml.gz.md5</a>	2021-08-02 00:44	57	
<a href="#">dblp-2021-08-01.xml.gz</a>	2021-08-02 00:44	617M	
<a href="#">dblp-2021-07-01.xml.gz.md5</a>	2021-07-02 00:04	57	
<a href="#">dblp-2021-07-01.xml.gz</a>	2021-07-02 00:04	611M	
<a href="#">dblp-2021-06-01.xml.gz.md5</a>	2021-06-02 00:31	57	
<a href="#">dblp-2021-06-01.xml.gz</a>	2021-06-02 00:31	606M	
<a href="#">dblp-2021-05-03.xml.gz.md5</a>	2021-05-03 23:33	57	

The dataset:

- Compressed archive of 623 MB.
- Once extracted: Single XML file of 3.2 GB.
- 8.5 million XML entries on publications, authors, journals, institutions, citations etc.

# Data conversion, import and transformations



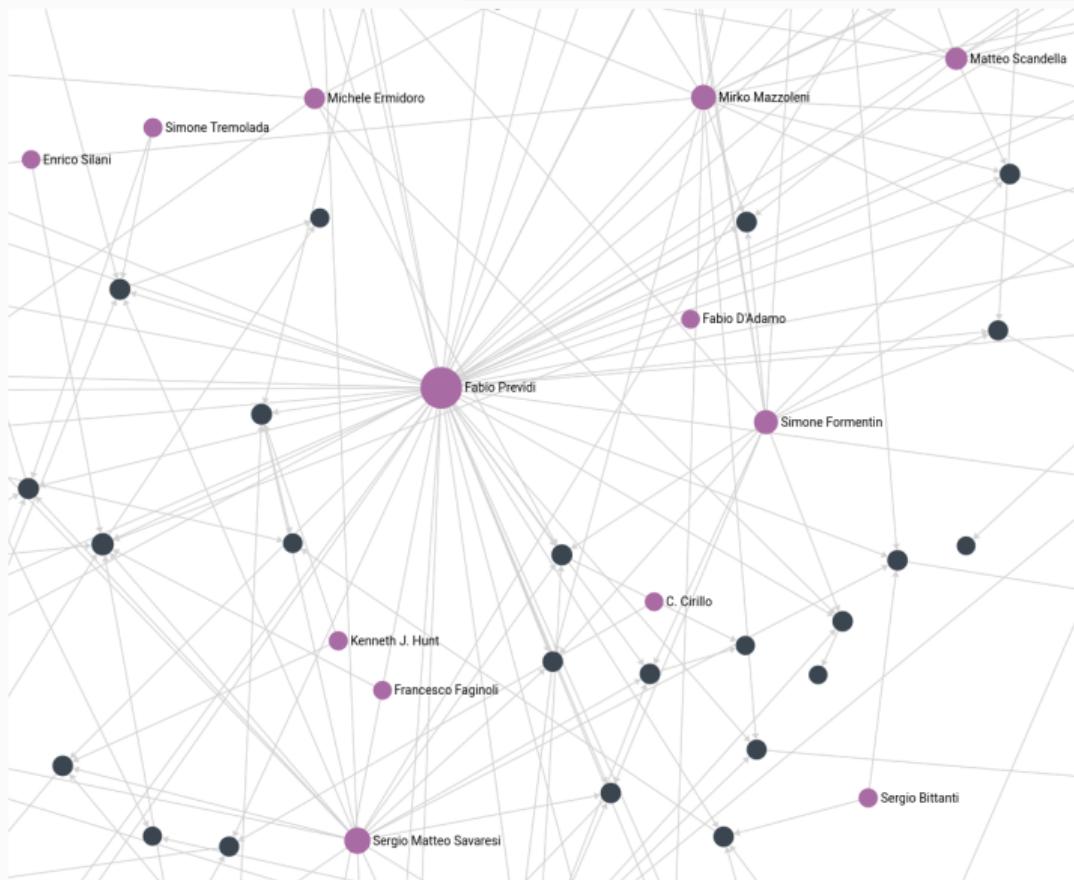
The steps:

1. Conversion from XML to line JSON;
2. Import in ArangoDB collection;

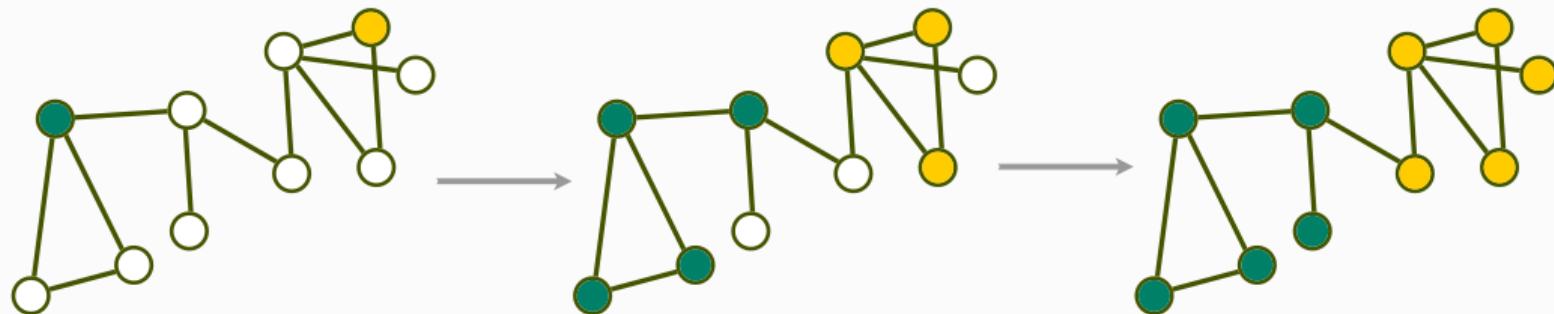
3. Transform imported data to vertex documents;
4. Link vertices by edges;
5. Build the graph.

# The graph

A subgraph of the obtained final graph:



## Pregel Label Propagation Community Detection Algorithm



The algorithm:

1. Initially every vertex is labeled with a unique label.
2. The labels are propagated from vertex to vertex iteratively.

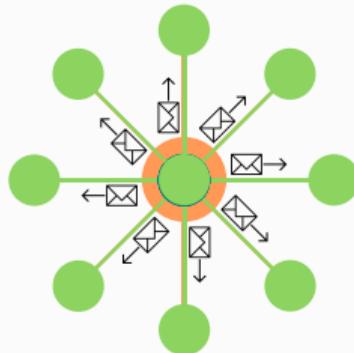
3. At the end of an iteration, each vertex updates its label to the one with maximum weight of neighbor vertices.
4. Convergence is reached when each vertex is labeled as the majority of its neighbors.

# A generic Pregel algorithm

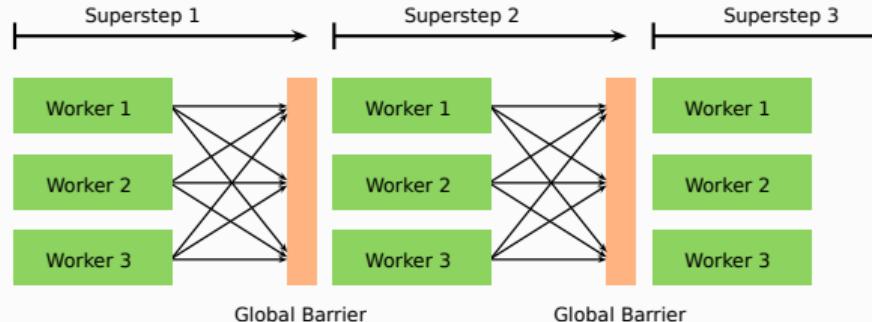
Pregel's supersteps:

**Diffusion:** information is propagated from vertex to neighbors.

**Fusion:** information is aggregated from neighbors to a set of entities



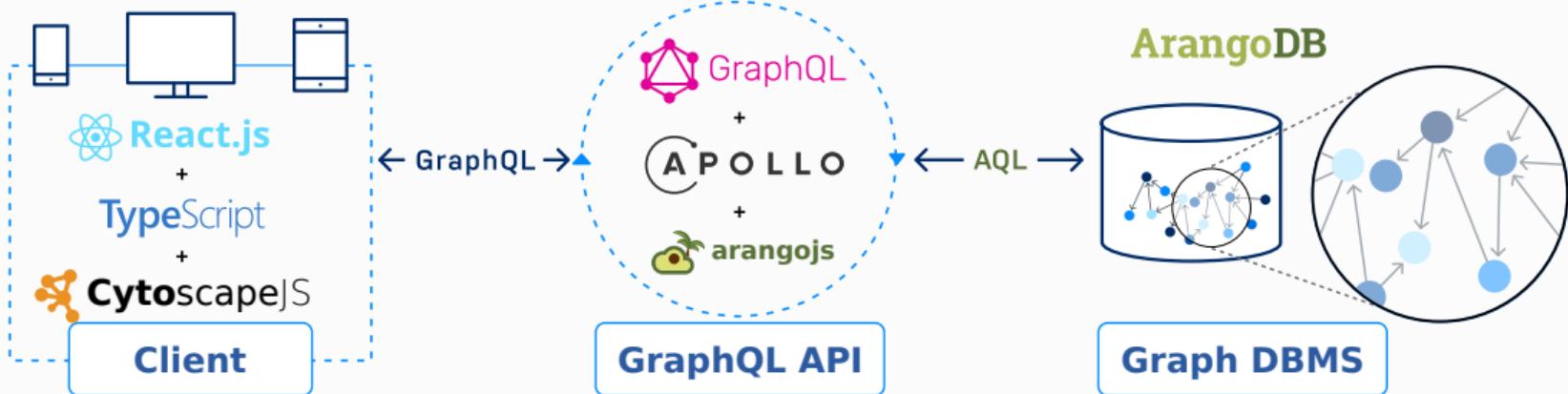
```
while not converged do
    /* Superstep */
    communicate:
        send(msg: own_value, to: neighbors)
    compute:
        own_value ← max(value_from_all_messages ∪ own_value)
end while
```



## Clustering results

Vertex type	Number of vertices	Number of detected communities
author	2786113	177592
editor	43644	9837
institution	56918	25415
journal	1905	1896
publication	5662747	141939
publisher	2292	1437
school	2098	1677
series	1742	934
all types	8557459	187451

# Web Application's Architecture



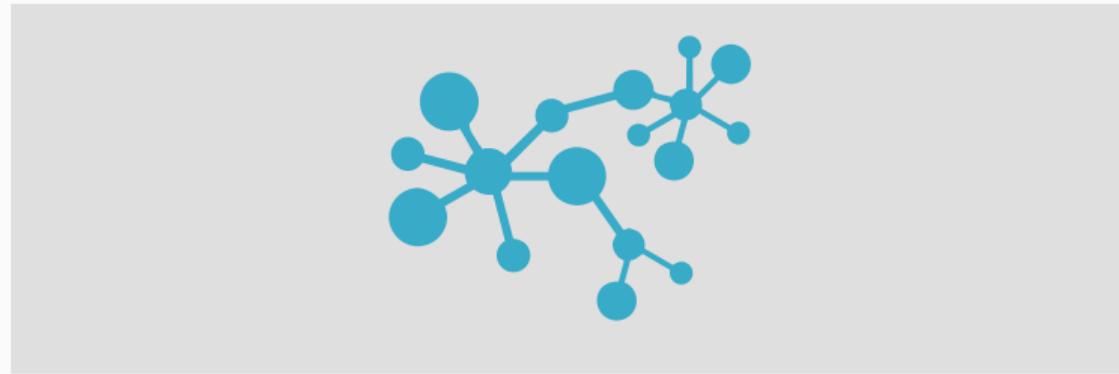
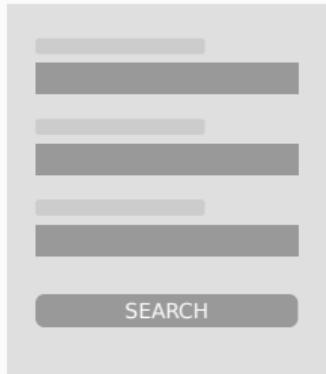
# GraphQL API



The API, initially REST then converted to a GraphQL API, is made of two resolvers:

1. A resolver function to handle the search form autocomplete suggestions.
2. A resolver function to provide collaboration graph's data.

## Academic Graph Connections



Made with ❤️

# Frontend UI Search Form

The image displays two side-by-side screenshots of a search form interface, illustrating its components and functionality.

**Left Screenshot:**

- Input string:** The text "Antonio Ferramosca" is entered into the search input field.
- Autocomplete suggestions:** A dropdown menu shows suggestions for "Author: Antonio Ferramosca".
- Node name or title:** The text "Antonio Ferramosca" is highlighted in green.
- Min depth:** A dropdown menu showing the value "1".
- Min depth:** A dropdown menu showing the value "2".
- Search button:** A large, rounded rectangular button labeled "SEARCH".

**Right Screenshot:**

- Input string inserted:** The text "Antonio Ferramosca" is entered into the search input field.
- Increment and decrement buttons for min and max depth numbers:** Two sets of up/down arrow buttons are shown, one for the minimum depth (value 1) and one for the maximum depth (value 2).

# Results display (Searching for "Angelo Gargantini")

## Academic Graph Connections

Node name or title

Angelo Gargantini



Minimum depth

1

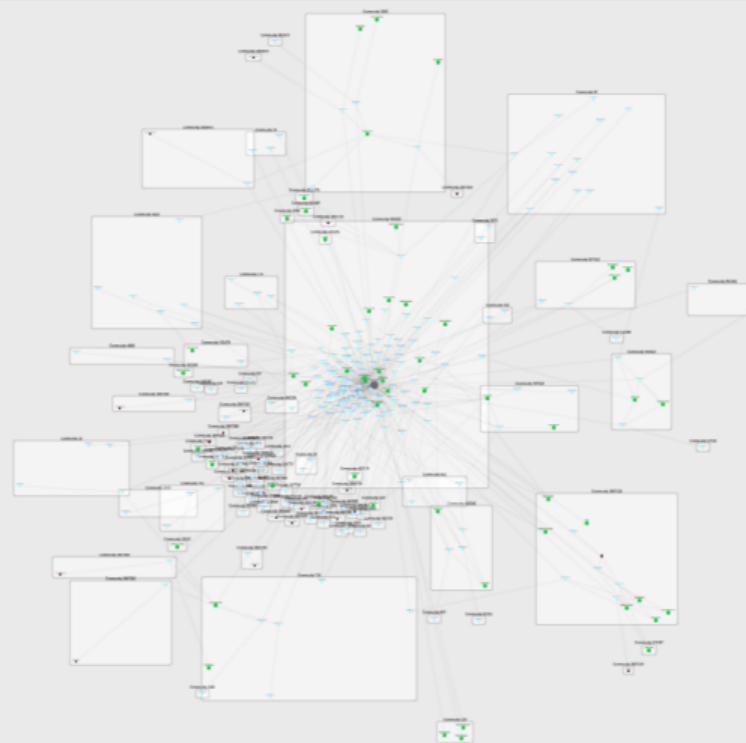


Maximum depth

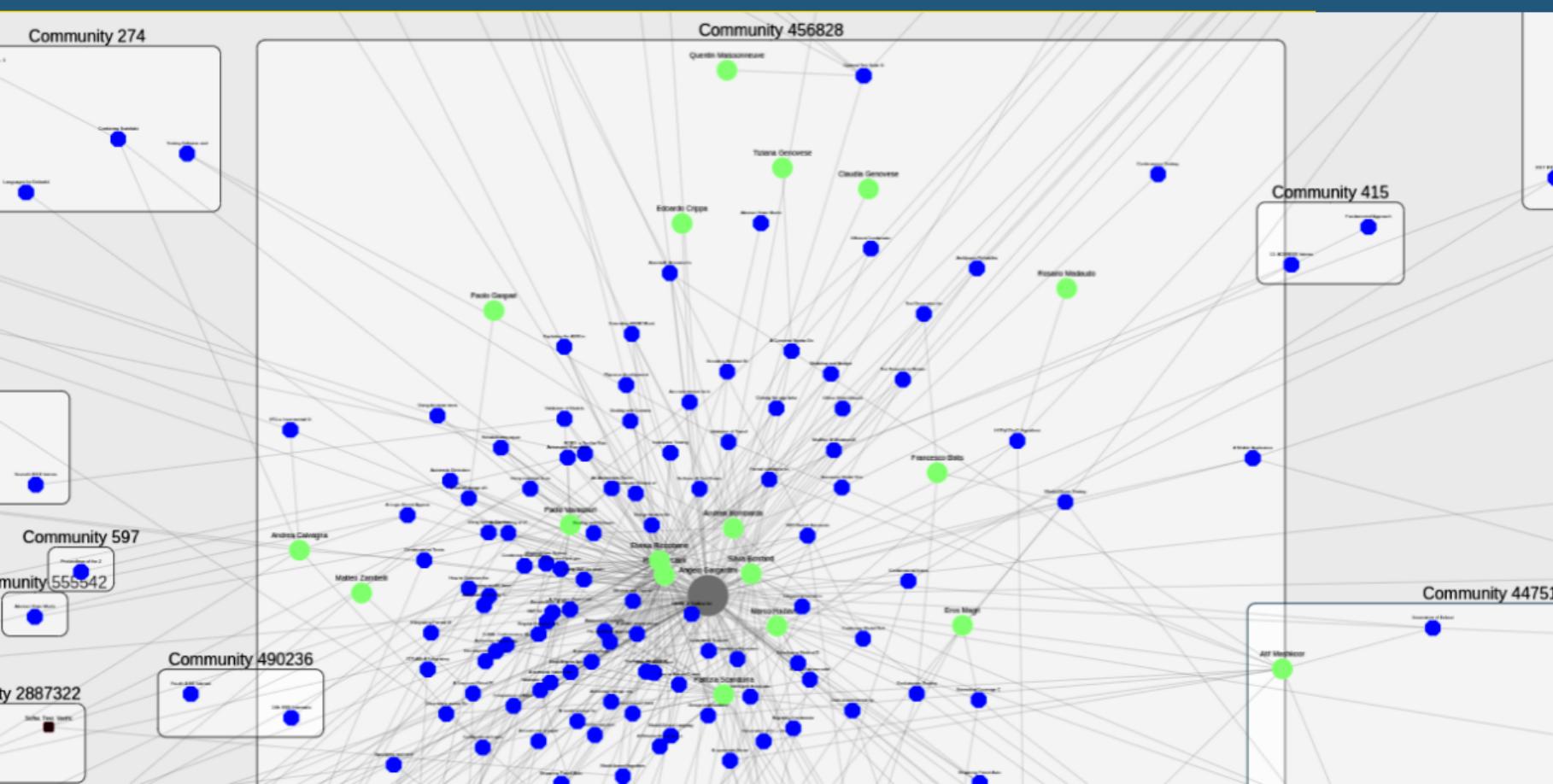
2

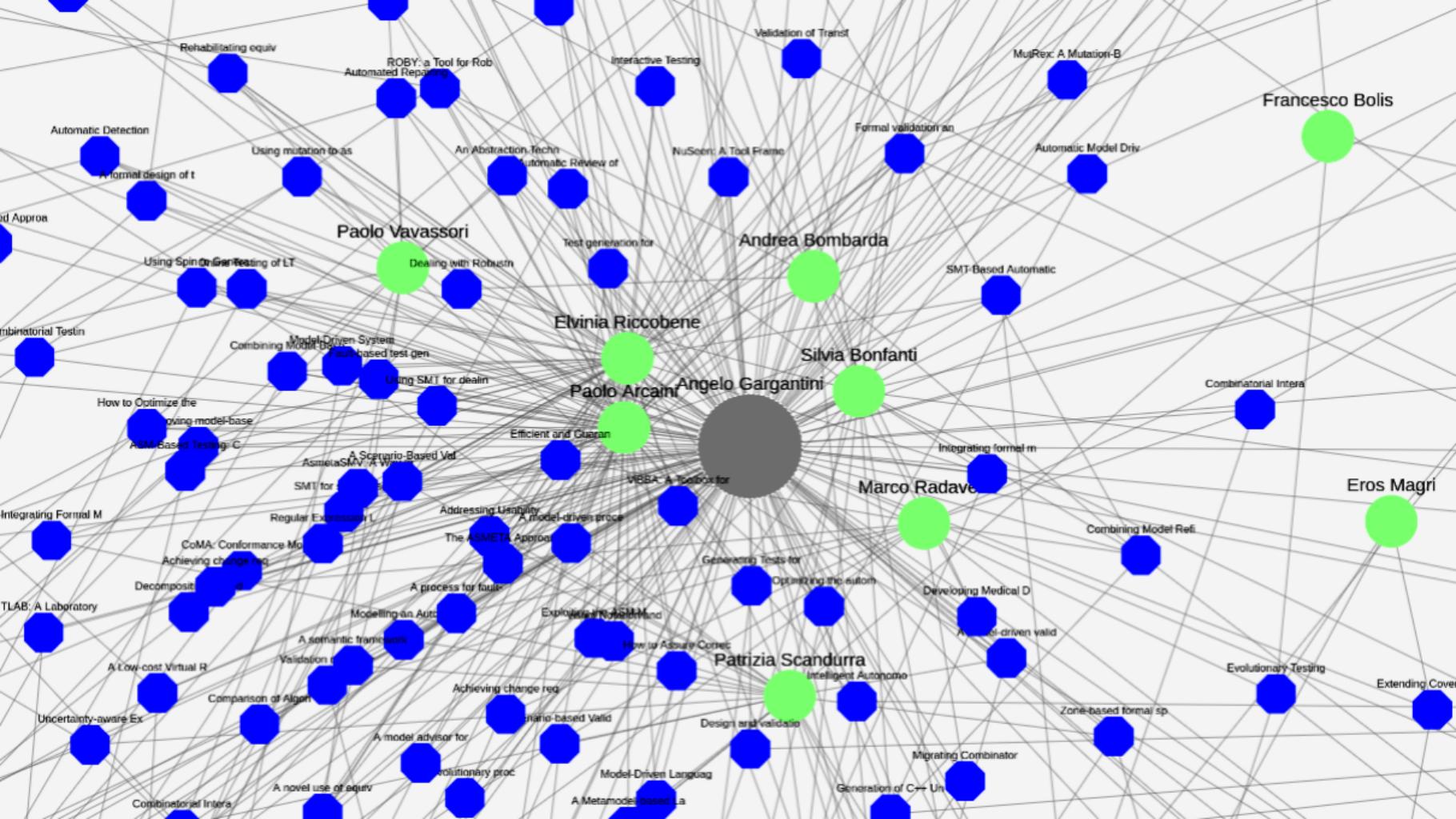


SEARCH

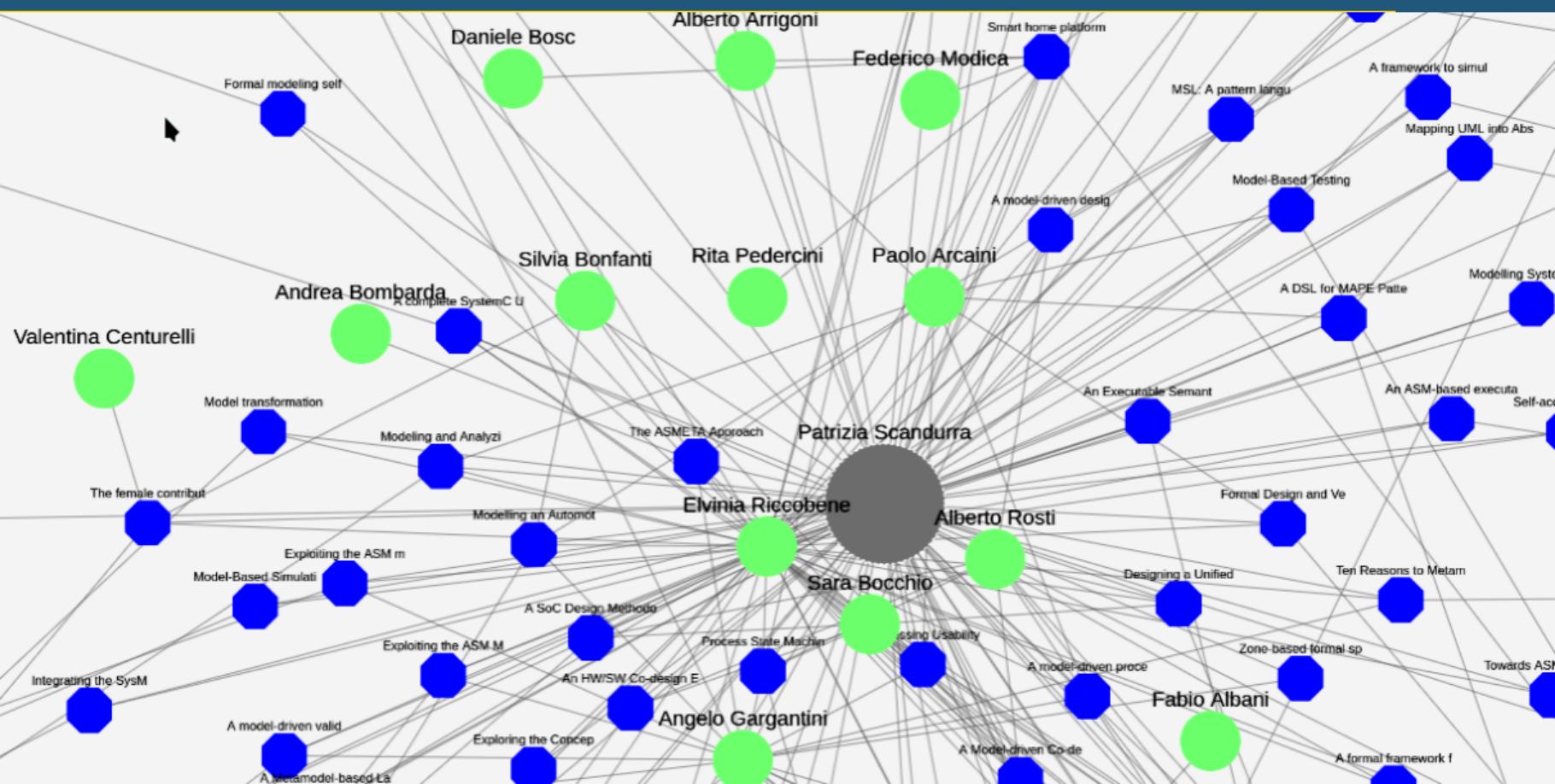


# Results display (Angelo Gargantini)



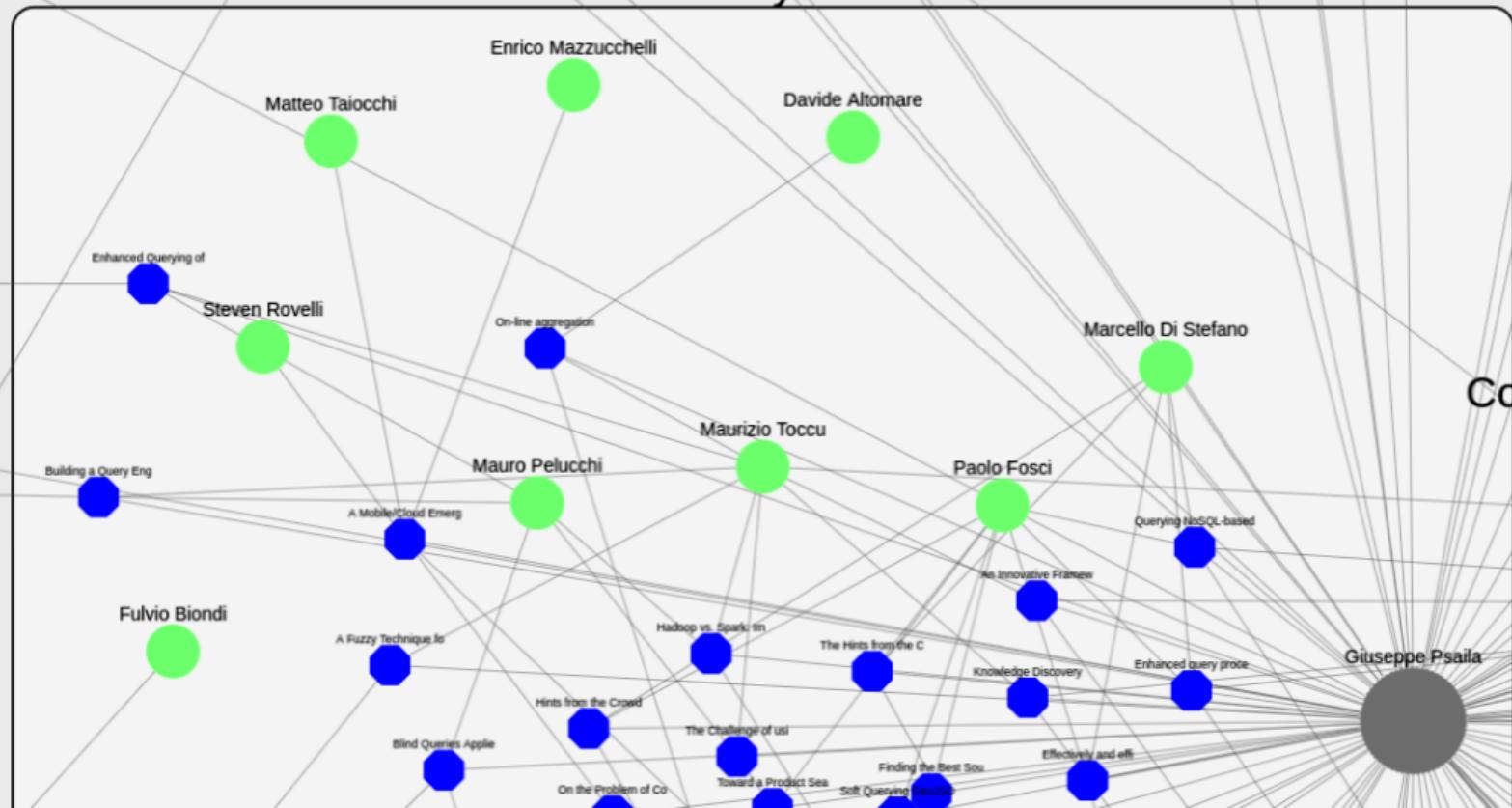


# Results display (Patrizia Scandurra)



# Results display (Giuseppe Psaila)

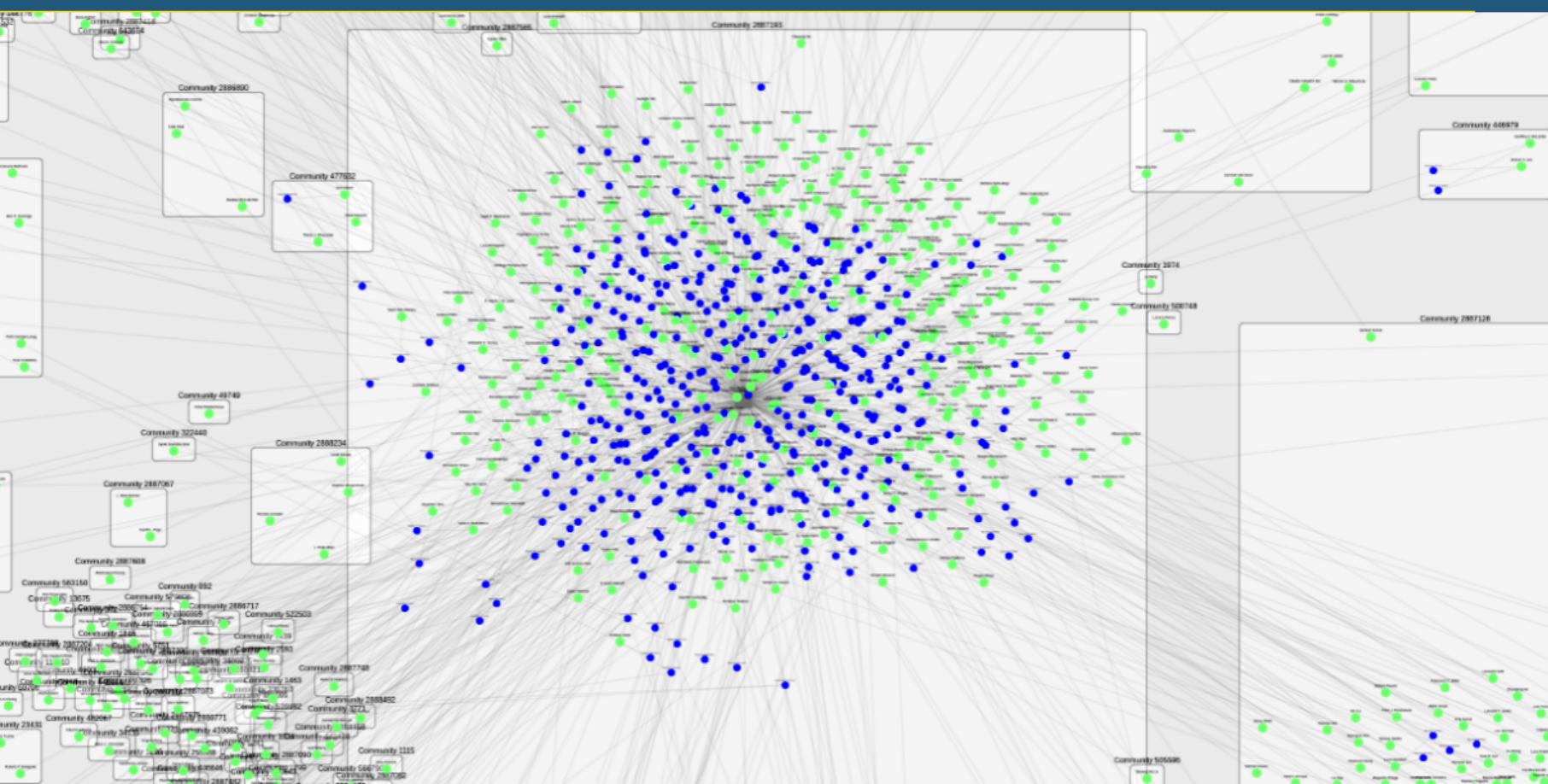
## Community 470922



Community

Designing Template

## **Results display (Statistical Methods & Applications Journal)**



# Results display (ETH Zurich)

## Academic Graph Connections

Node name or title

ETH Zurich



Minimum depth

1

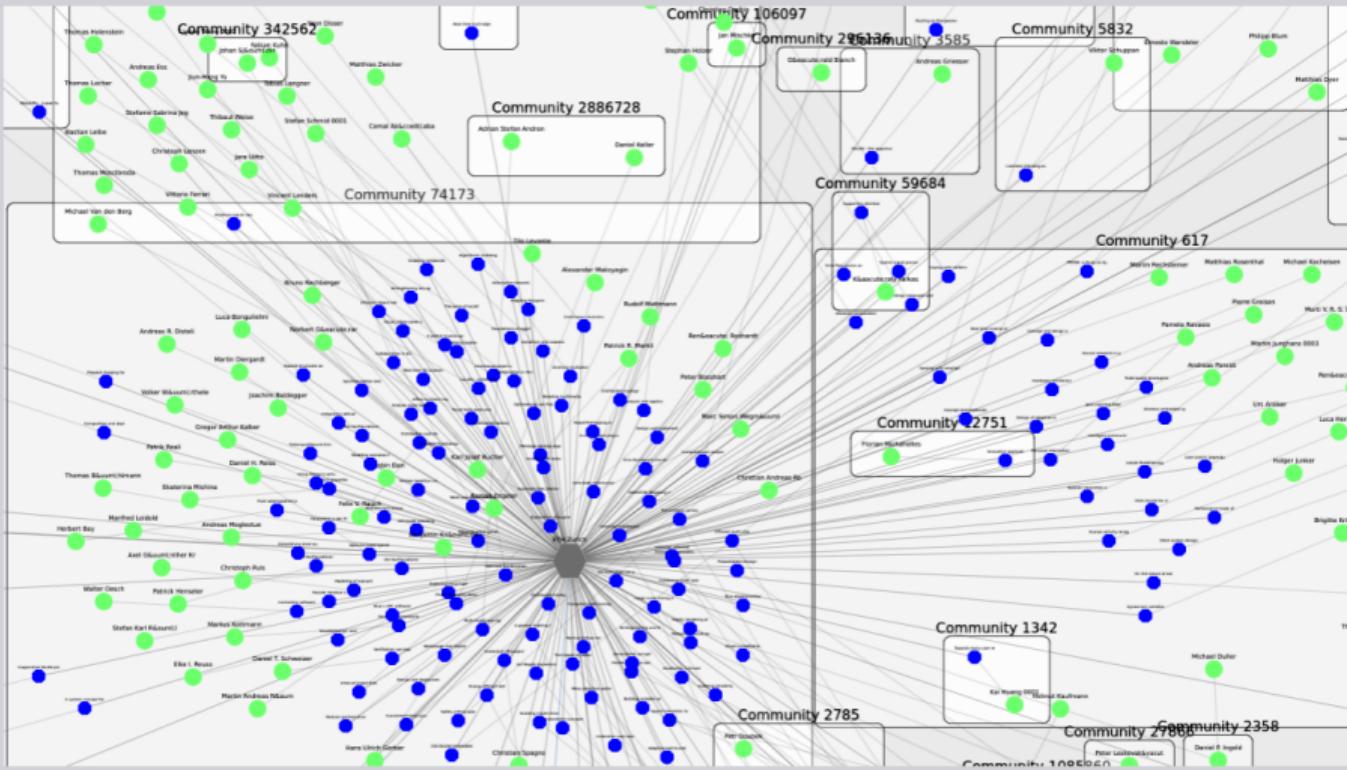


Maximum depth

2



SEARCH



Thank you!

Questions?

Name  
Surname

**CLUSTERING  
GRAPHS**

Applying a Label Propagation  
Algorithm to Detect Communities  
in Graph Databases