

IM<sub>3</sub>

INTEGRATED  
MULTISECTOR  
MULTISCALE  
MODELING

# Strategic Systematic Review and Exploration of the Research Area of MultiSector Dynamics using Natural Language Processing, Graph Machine Learning, and Large Language Models

Chris R. Vernon and Patrick M. Reed  
on behalf of the MSD CoP facilitation team

This research is supported by the U.S. Department of Energy, Office of Science, as part of research in MultiSector Dynamics, Earth and Environmental System Modeling Program



Cornell University



THE UNIVERSITY  
of NORTH CAROLINA  
at CHAPEL HILL



## MSD COP FACILITATION TEAM



**Chris Vernon**  
PNNL



**Patrick Reed**  
Cornell University



**Erwan Monier**  
UC Davis



**Antonia Hadjimichael**  
Penn State



**Rohini Gupta**  
Cornell University



**Lillian Lau**  
Cornell University



**Sequoia Alba**  
UC Davis



**Gabriela Gesualdo**  
Penn State



**Hamsa Ganapathi**  
UC Davis



## MSD COP SCIENTIFIC STEERING GROUP (SSG)



**Nathalie Voisin,**  
PNNL  
Core Member



**Klaus Keller,**  
Dartmouth  
Core Member



**Nicole Jackson,**  
Sandia  
Core Member



**Casey Burleyson,**  
PNNL  
Core Member



**Jen Morris,** MIT  
Core member



**Andy Jones,**  
UC Berkeley  
Core member



**Rebecca Saari,**  
U. of Waterloo  
Core Member



**David McCollum,** Oak  
Ridge  
WG Representative



**Wei Peng,**  
Princeton University  
WG Representative



**Julia Szinai,** LBNL  
WG representative



**Vivek Srikrishnan,** Cornell  
University  
WG representative



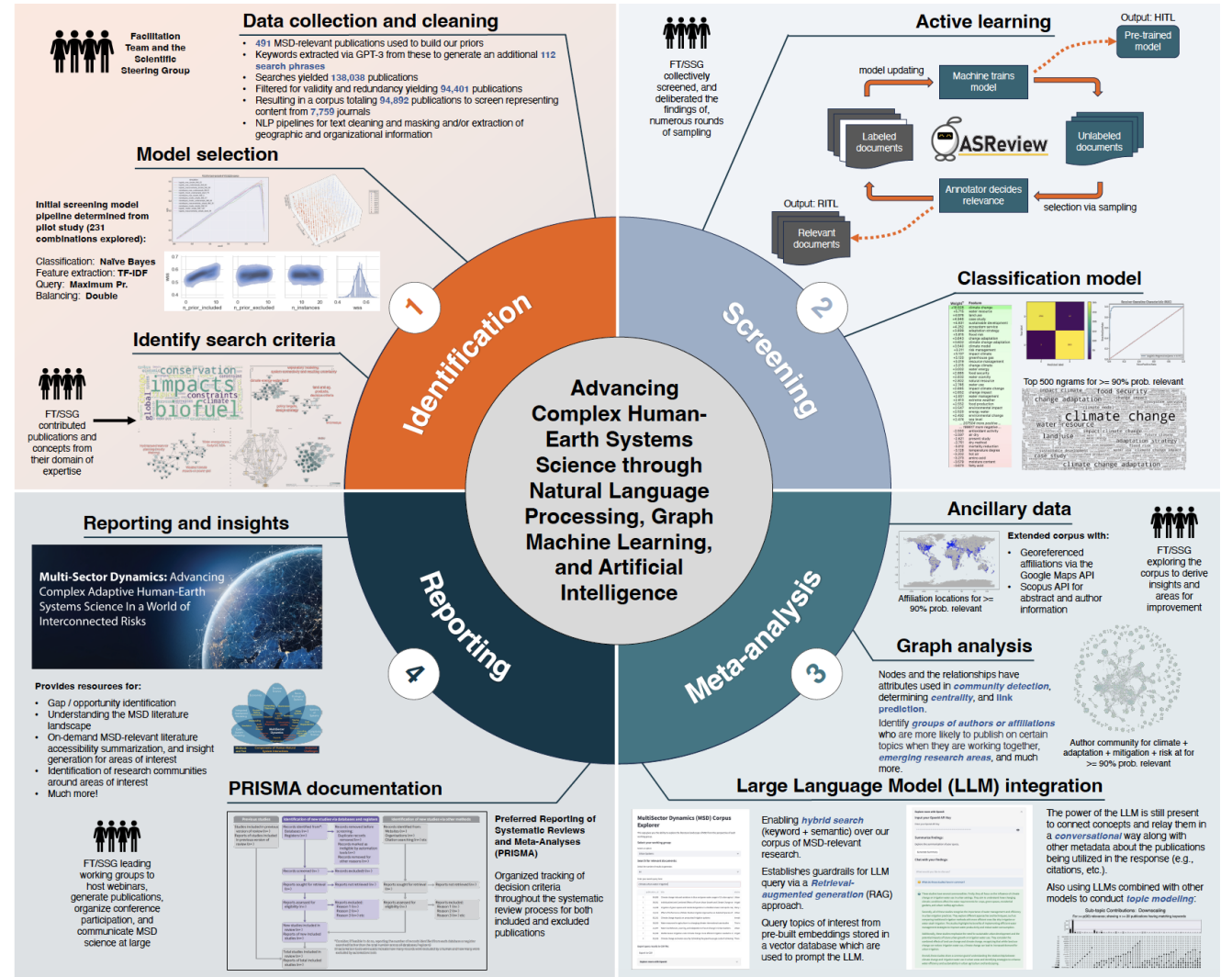
**Christa Brelsford,**  
LANL  
WG representative



**Jim Yoon,** PNNL  
WG representative

# MSD-FUTURES: AI-ENABLED SYSTEMATIC REVIEW OF MSD

- Collaborative effort with our SSG and FT members
- **Fully transferable framework** using Large Language Model and Graph, and other ML innovations to understand the MSD literature landscape at large for **105,336** publications
- Identify gaps and opportunities for collaboration
- **MSD-FUTURES: Foresights for Understanding Thematic Unity in Reviews of Emergent Science**







Facilitation Team and the Scientific Steering Group

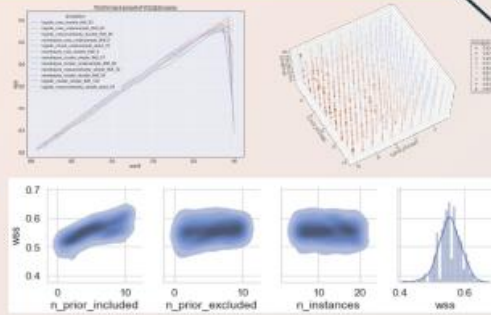
## Data collection and cleaning

- 491 MSD-relevant publications used to build our priors
- Keywords extracted via GPT-3 from these to generate an additional 112 search phrases
- Searches yielded 138,038 publications
- Filtered for validity and redundancy yielding 94,401 publications
- Resulting in a corpus totaling 94,892 publications to screen representing content from 7,759 journals
- NLP pipelines for text cleaning and masking and/or extraction of geographic and organizational information

## Model selection

Initial screening model pipeline determined from pilot study (231 combinations explored):

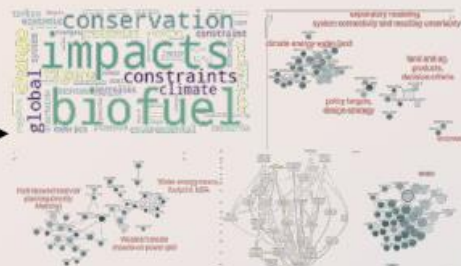
Classification: Naive Bayes  
 Feature extraction: TF-IDF  
 Query: Maximum Pr.  
 Balancing: Double



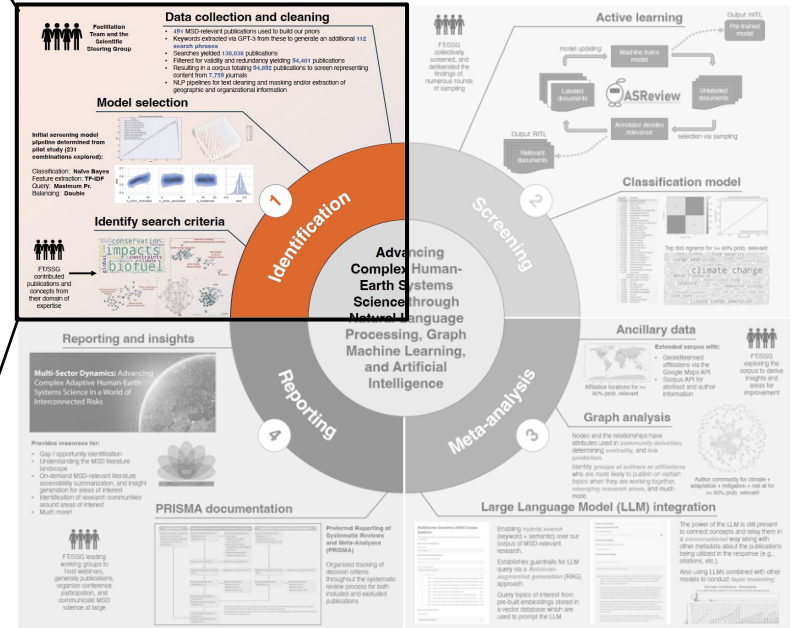
## Identify search criteria



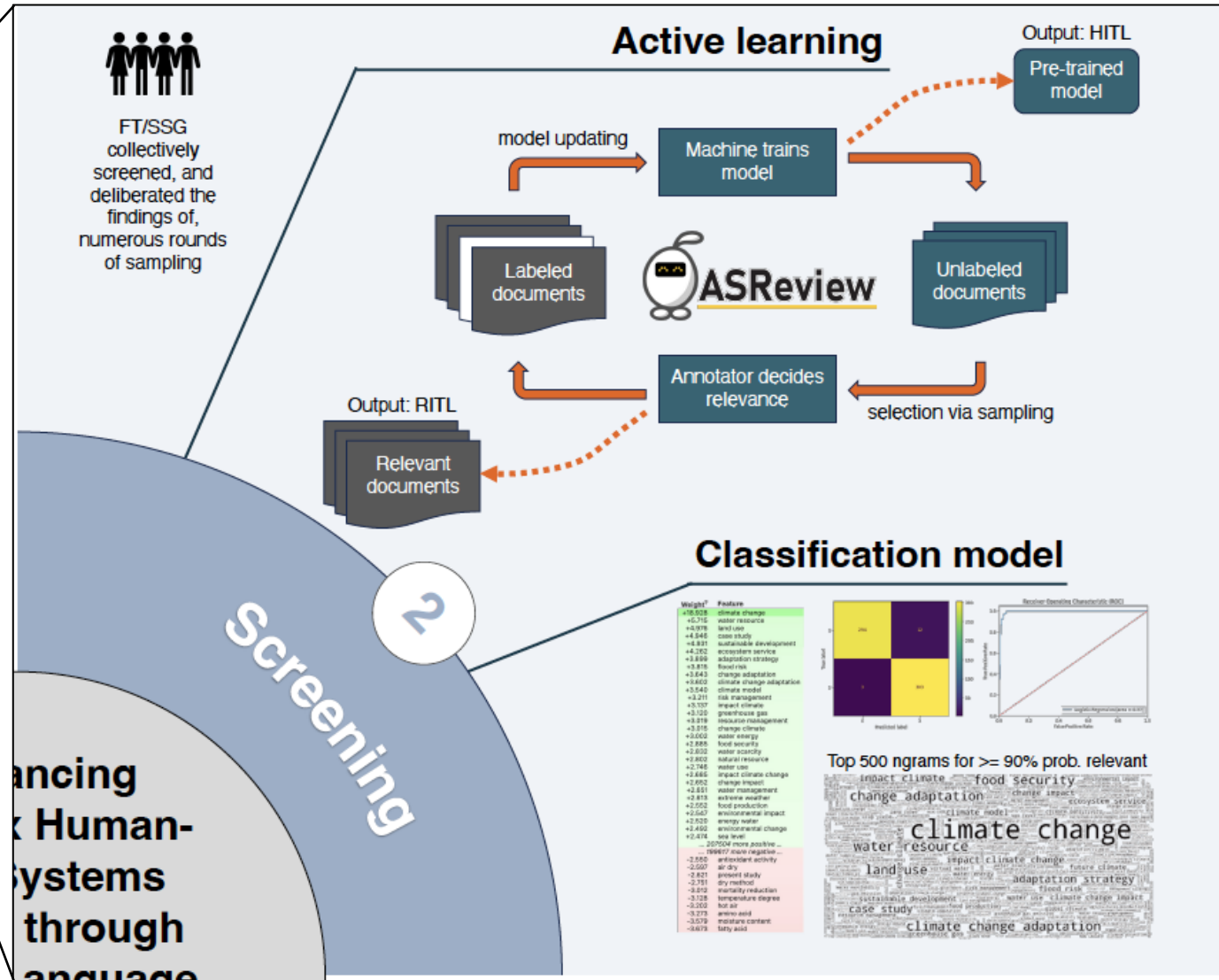
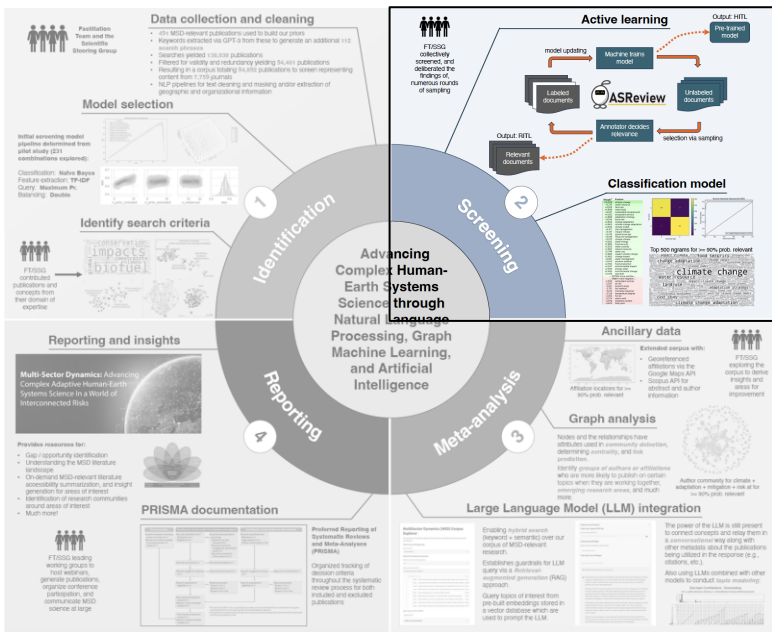
FT/SSG contributed publications and concepts from their domain of expertise



Identification

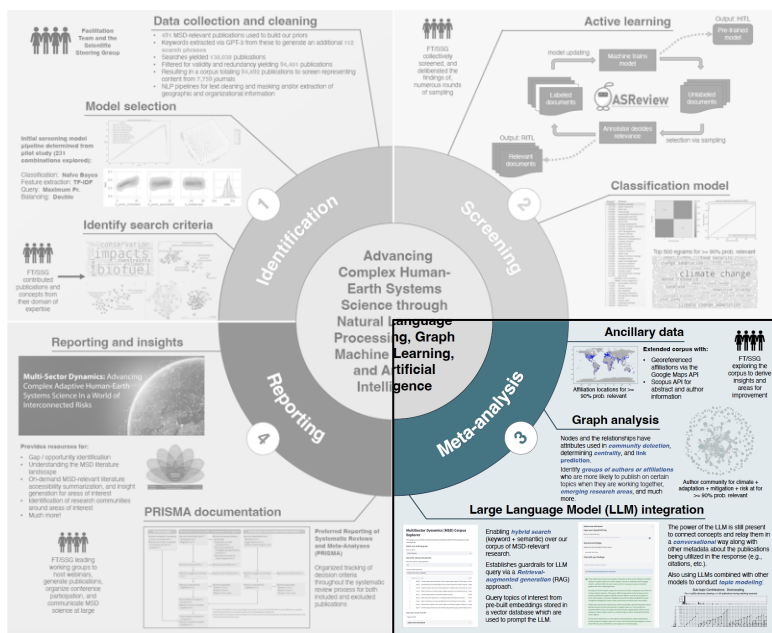


# MSD-FUTURES: AI-ENABLED SYSTEMATIC REVIEW OF MSD





# MSD-FUTURES: AI-ENABLED SYSTEMATIC REVIEW OF MSD



## Language Learning, Graph Learning, Artificial Intelligence

### Ancillary data

Affiliation locations for >= 90% prob. relevant

**Extended corpus with:**

- Georeferenced affiliations via the Google Maps API
- Scopus API for abstract and author information

FT/SSG exploring the corpus to derive insights and areas for improvement

# Meta-analysis

3

### Graph analysis

Nodes and the relationships have attributes used in *community detection*, determining *centrality*, and *link prediction*.

Identify *groups of authors or affiliations* who are more likely to publish on certain topics when they are working together, *emerging research areas*, and much more.

Author community for climate + adaptation + mitigation + risk at for >= 90% prob. relevant

### Large Language Model (LLM) integration

Enabling *hybrid search* (keyword + semantic) over our corpus of MSD-relevant research.

Establishes guardrails for LLM query via a *Retrieval-augmented generation (RAG)* approach.

Query topics of interest from pre-built embeddings stored in a vector database which are used to prompt the LLM.

MultiSector Dynamics (MSD) Corpus Explorer

System prompt with Search: Input your OpenAI API Key

Summarize findings: Input the summarization of your query

Chat with your findings: What is the distribution of...

The power of the LLM is still present to connect concepts and relay them in a *conversational* way along with other metadata about the publications being utilized in the response (e.g., citations, etc.).

Also using LLMs combined with other models to conduct *topic modeling*:

Sub-topic Contributions: Downscaling (For >= 90% relevance, showing n = 20 publications having matching keywords)

## Reporting and insights

**Multi-Sector Dynamics: Advancing Complex Adaptive Human-Earth Systems Science In a World of Interconnected Risks**

**Provides resources for:**

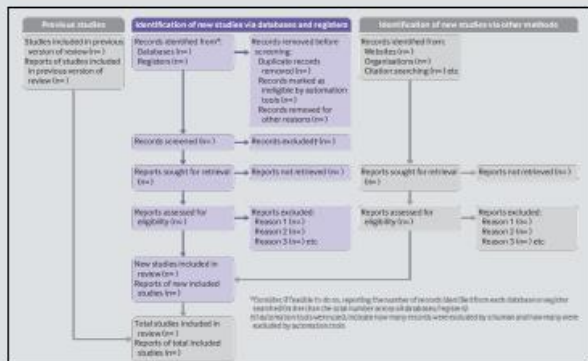
- Gap / opportunity identification
- Understanding the MSD literature landscape
- On-demand MSD-relevant literature accessibility summarization, and insight generation for areas of interest
- Identification of research communities around areas of interest
- Much more!



## PRISMA documentation



FT/SSG leading working groups to host webinars, generate publications, organize conference participation, and communicate MSD science at large

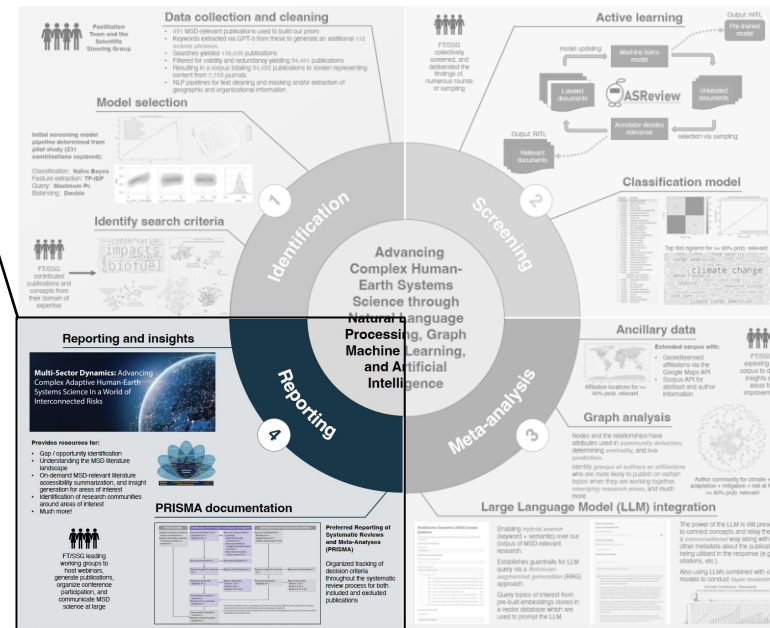


### Preferred Reporting of Systematic Reviews and Meta-Analyses (PRISMA)

Organized tracking of decision criteria throughout the systematic review process for both included and excluded publications

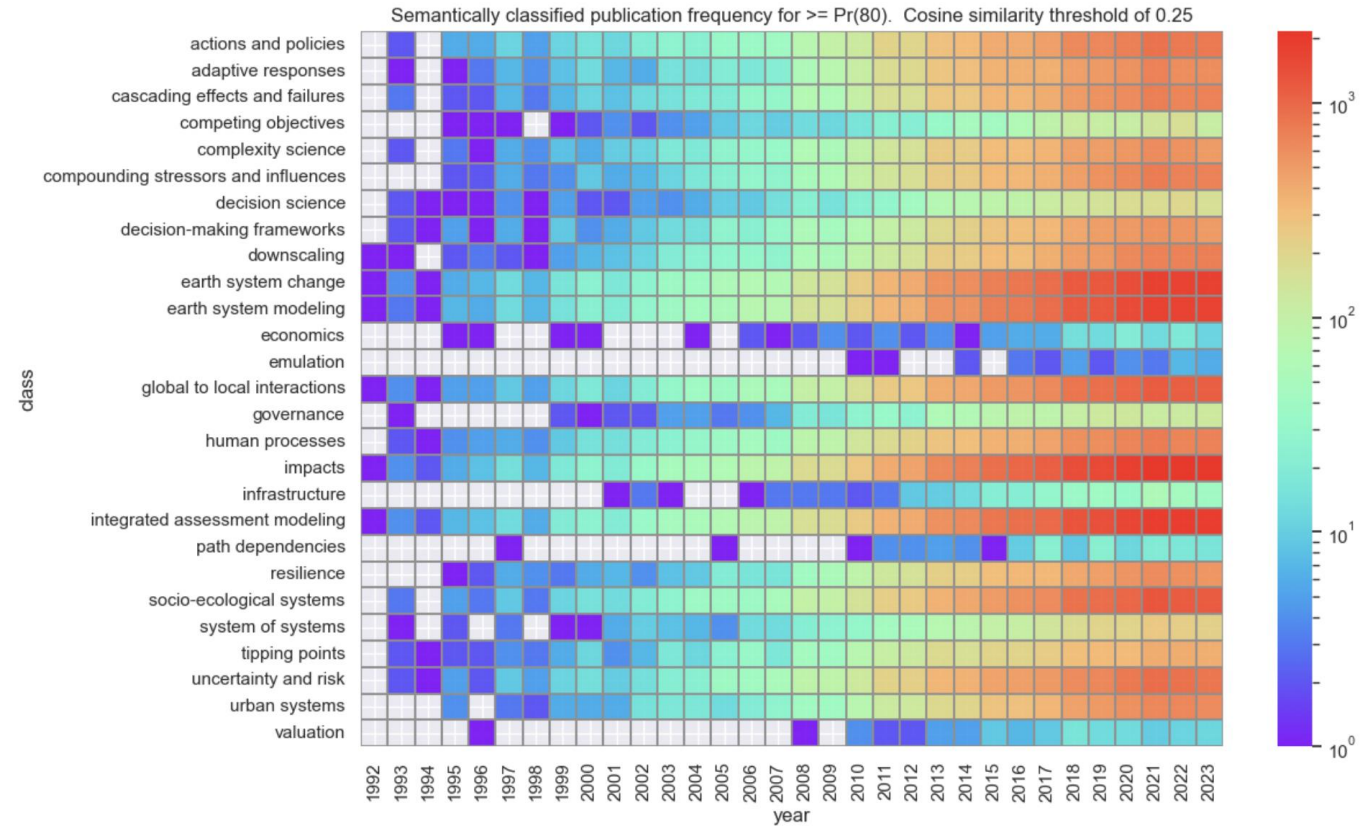
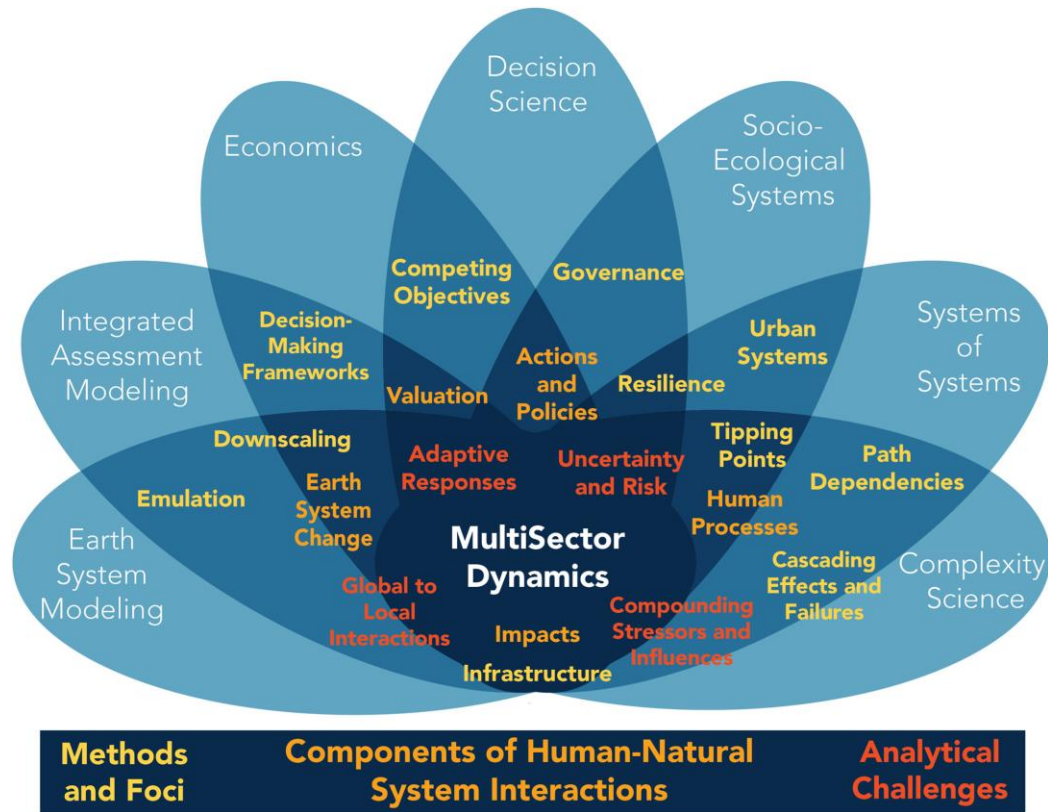
# Reporting

Natural Language Processing Machine Learning and Artificial Intelligence





## FROM QUALITATIVE TO QUANTITATIVE



Reed et al., 2022, Earth's Future  
<https://doi.org/10.1029/2021EF002621>

Shows raw count of publications per year per feather for  $P(X \geq 80)$  under a higher similarity threshold; number of documents = 19031

# LARGE LANGUAGE MODEL (LLM) INTEGRATION

## MultiSector Dynamics (MSD) Corpus Explorer

This app gives you the ability to explore the literature landscape of MSD from the perspective of each working group.

### Select your working group:

Select an option:

Urban Systems

### Search for relevant documents:

Select the number of results to generate:

10

Enter your search query here:

climate urban water irrigation

	publication_id	title	abstra
0	84,806	Climate change-induced variations in blue and green water usage in US urban agricul	Urban
1	39,561	Individualized and Combined Effects of Future Urban Growth and Climate Change on Irrigati	Irrigati
2	12,038	Irrigation of green spaces and residential gardens in a Mediterranean metropolis: Ga	Many c
3	13,412	Effect of Performance of Water Stashes Irrigation Approaches on Selected Species of	Urban
4	83,161	Climate change impacts on pressurised irrigation systems	Invest
5	36,620	Water resources for agriculture in a changing climate: international case studies	This in
6	21,497	Water Use Behavior, Learning, and Adaptation to Future Change in Urban Gardens	Urban
7	52,549	Mediterranean irrigation under climate change: more efficient irrigation needed to cc	Irrigati
8	92,616	Climate change and water security: Estimating the greenhouse gas costs of achieving	There

Export query results to CSV file:

Export to CSV

Explore more with OpenAI

Explore more with OpenAI

Input your OpenAI API Key

Enter your OpenAI API Key

.....



Summarize findings:

Explore the summarization of your query.

Generate Summary

Chat with your findings:

What would you like to discuss?

🗣️ What do these studies have in common?

🗣️ These studies have several commonalities. Firstly, they all focus on the influence of climate change on irrigation water use in urban settings. They aim to understand how changing climatic conditions affect the water requirements for crops, green spaces, residential gardens, and urban rooftop agriculture.

Secondly, all of these studies recognize the importance of water management and efficiency in urban irrigation practices. They explore different approaches and techniques, such as comparing traditional irrigation methods with more efficient ones like drip irrigation or water-stash irrigation. The studies highlight the benefits of implementing efficient water management strategies to improve water productivity and reduce water consumption.

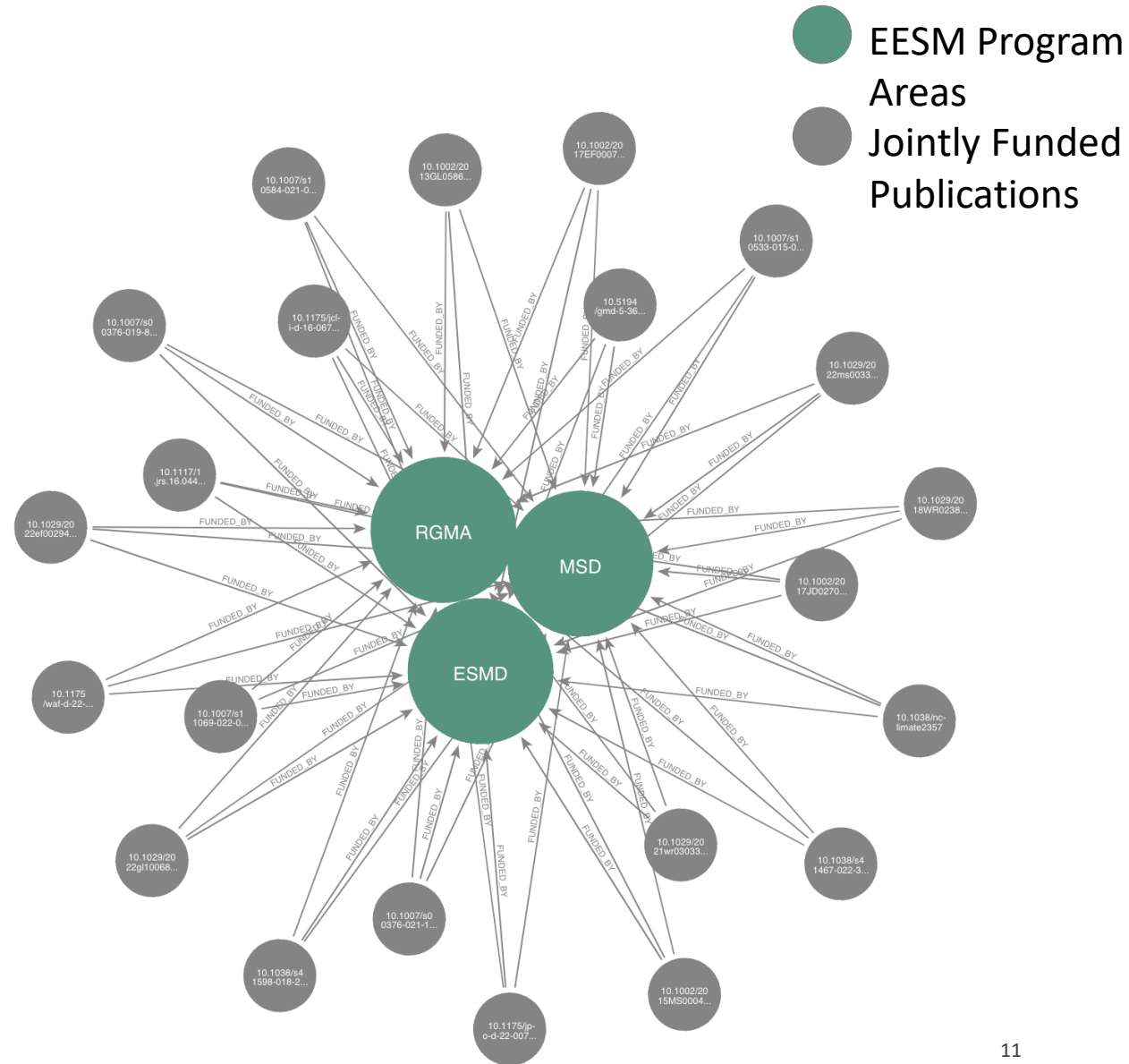
Additionally, these studies emphasize the need for sustainable urban development and the potential impacts of future urban growth on irrigation water use. They consider the combined effects of land use change and climate change, recognizing that while land use change can reduce irrigation water use, climate change can lead to increased demand for urban irrigation.

Overall, these studies share a common goal of understanding the relationship between climate change and irrigation water use in urban areas and identifying strategies to enhance water efficiency and sustainability in urban agriculture and landscaping.



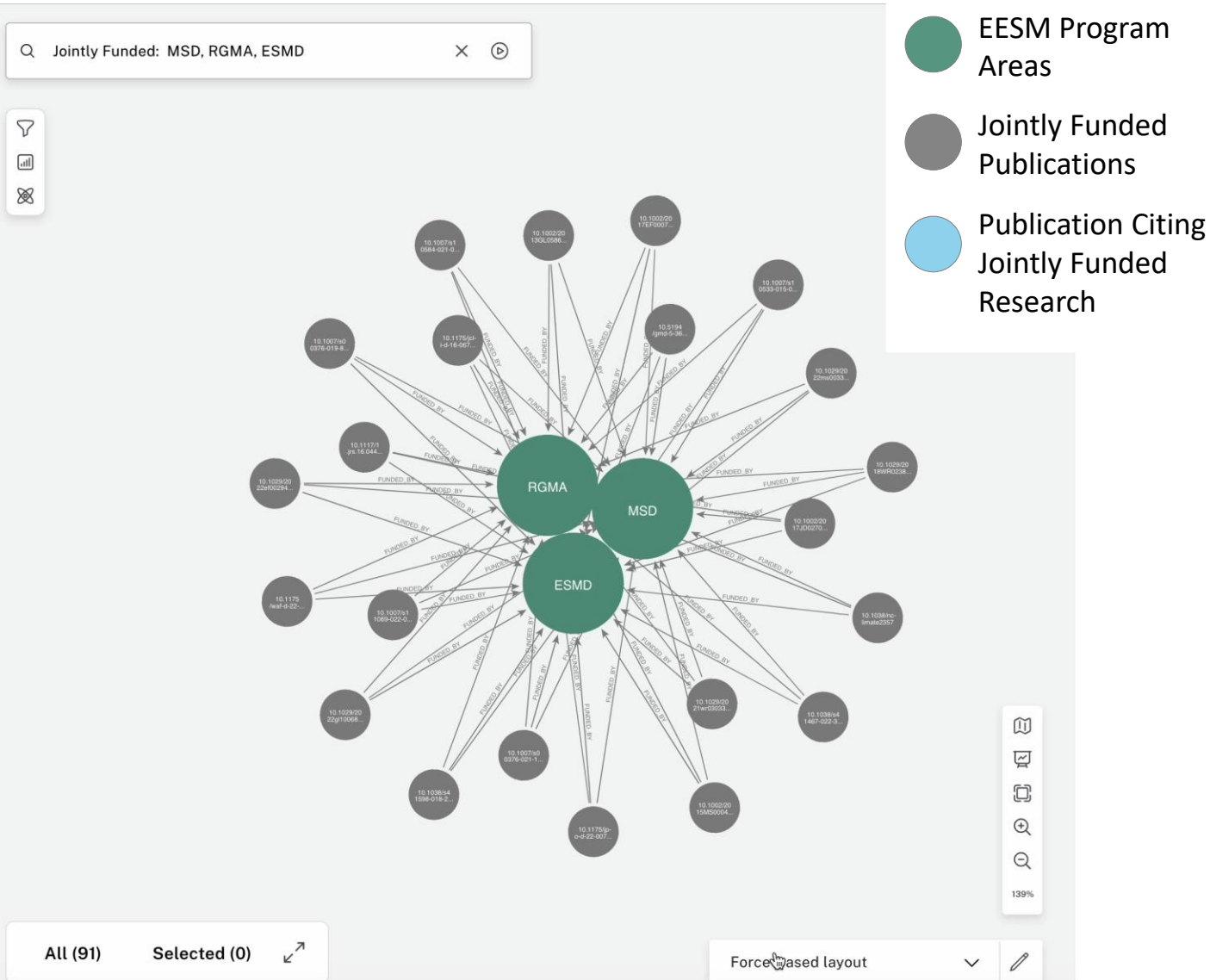
# ILLUSTRATIVE EXAMPLE: EXTENDING METHODS TO EESM AT LARGE

- Scraped the EESM publications website for all journal articles listed as funded under MSD, RGMA, or ESMD
- Produced 2,243 journal articles
- Of which, 22 were co-funded by all three program areas
- Let's talk about these for a few minutes...

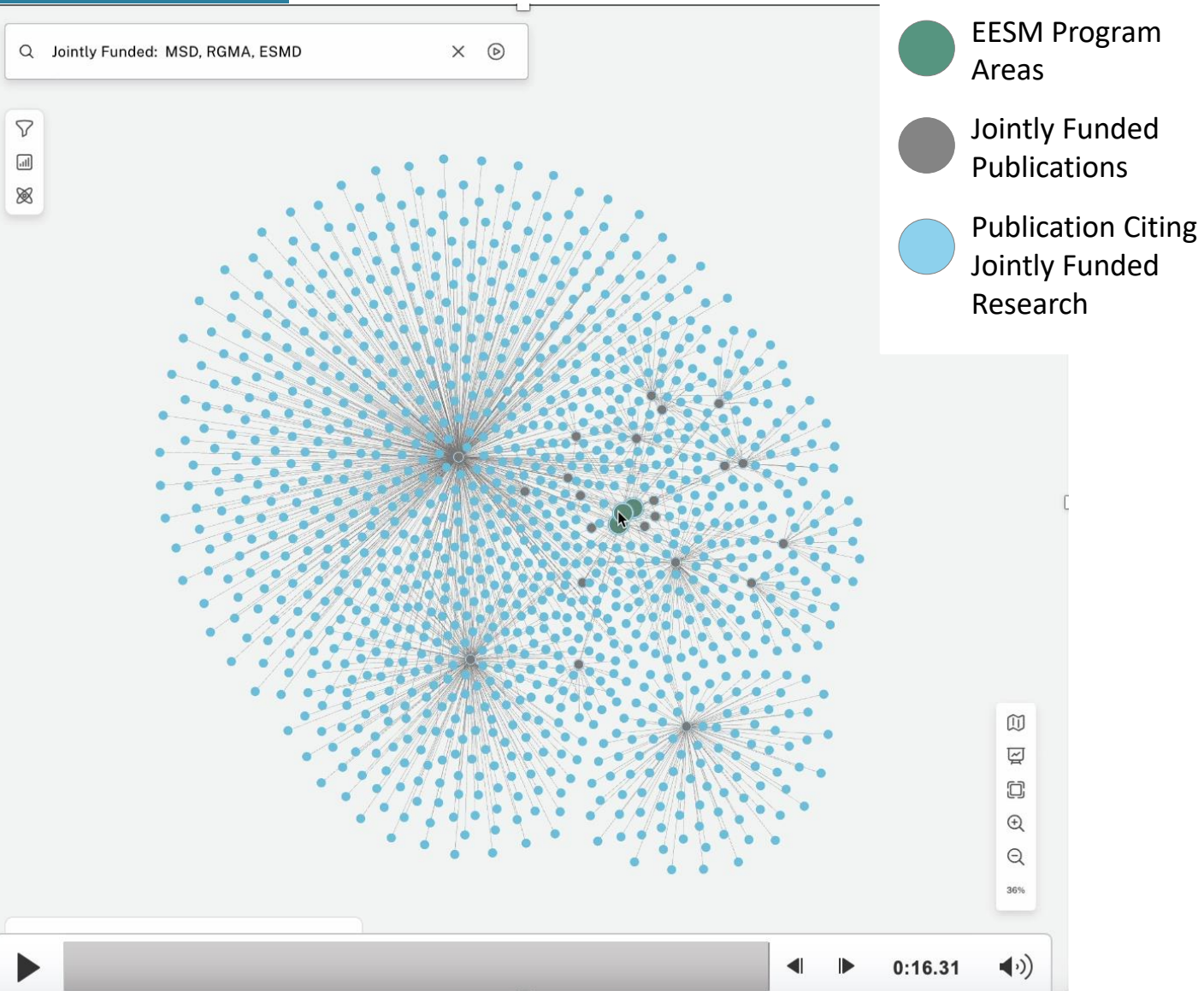


# A SNAPSHOT OF IMPACT WHEN EESM JOINTLY FUNDED

- The 22 jointly funded publications by all three program areas were cited 1,120 times



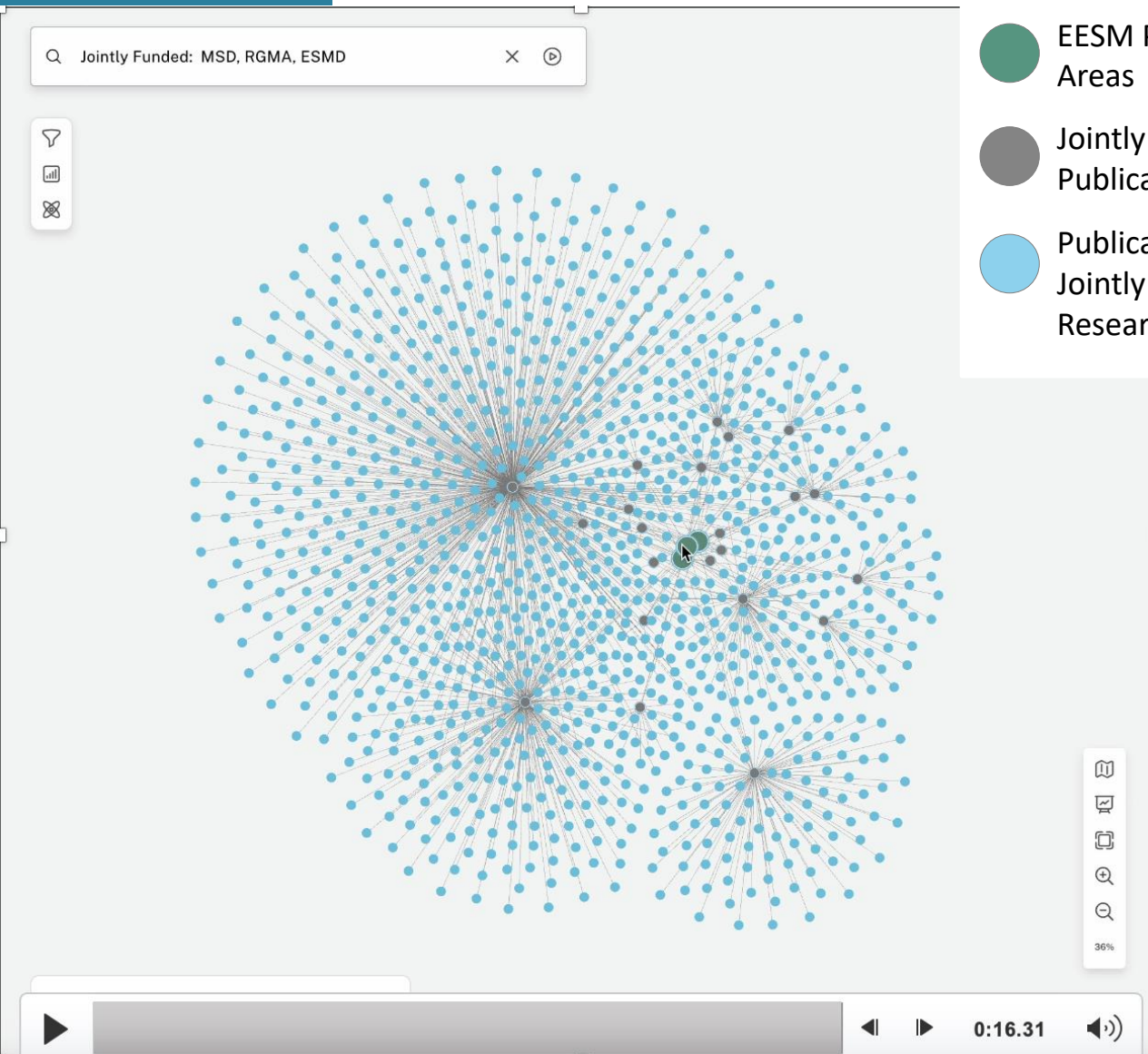
# A SNAPSHOT OF IMPACT WHEN EESM JOINTLY FUNDED



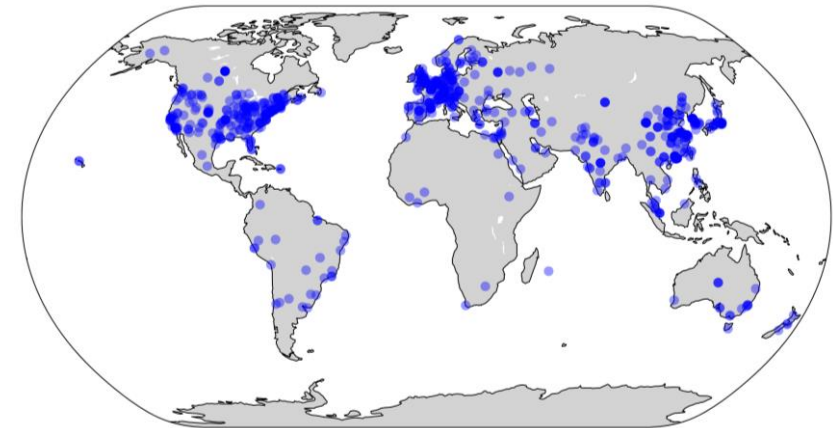
- The 22 jointly funded publications by all three program areas were cited 1,120 times
- Only 62 citations were from other EESM funded research – leaving an external impact of 1,058 citations



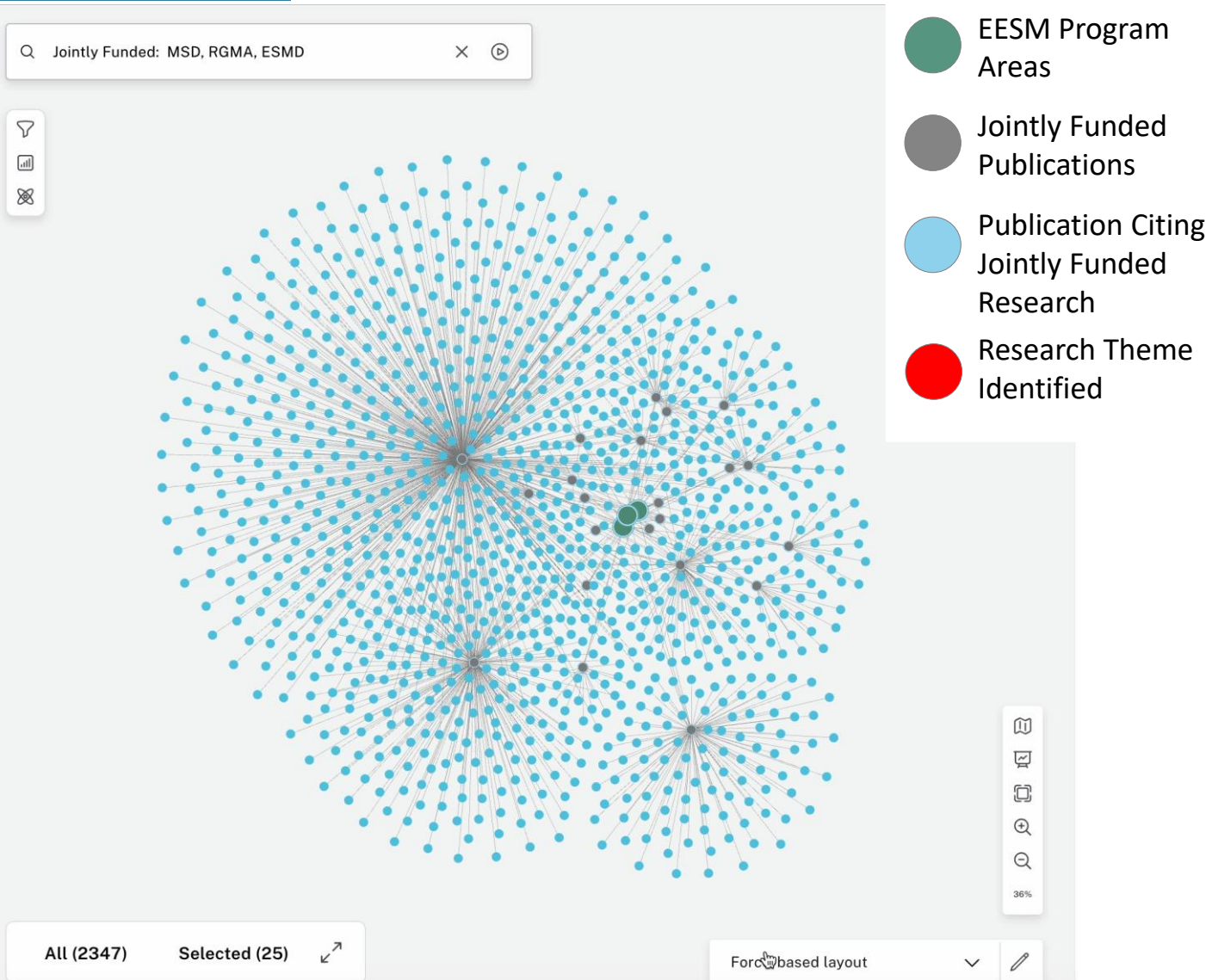
# A SNAPSHOT OF IMPACT WHEN EESM JOINTLY FUNDED



- The 22 jointly funded publications by all three program areas were cited 1,120 times
- Only 62 citations were from other EESM funded research – leaving an external impact of 1,058 citations
- Citations included 801 unique author affiliations that were globally distributed



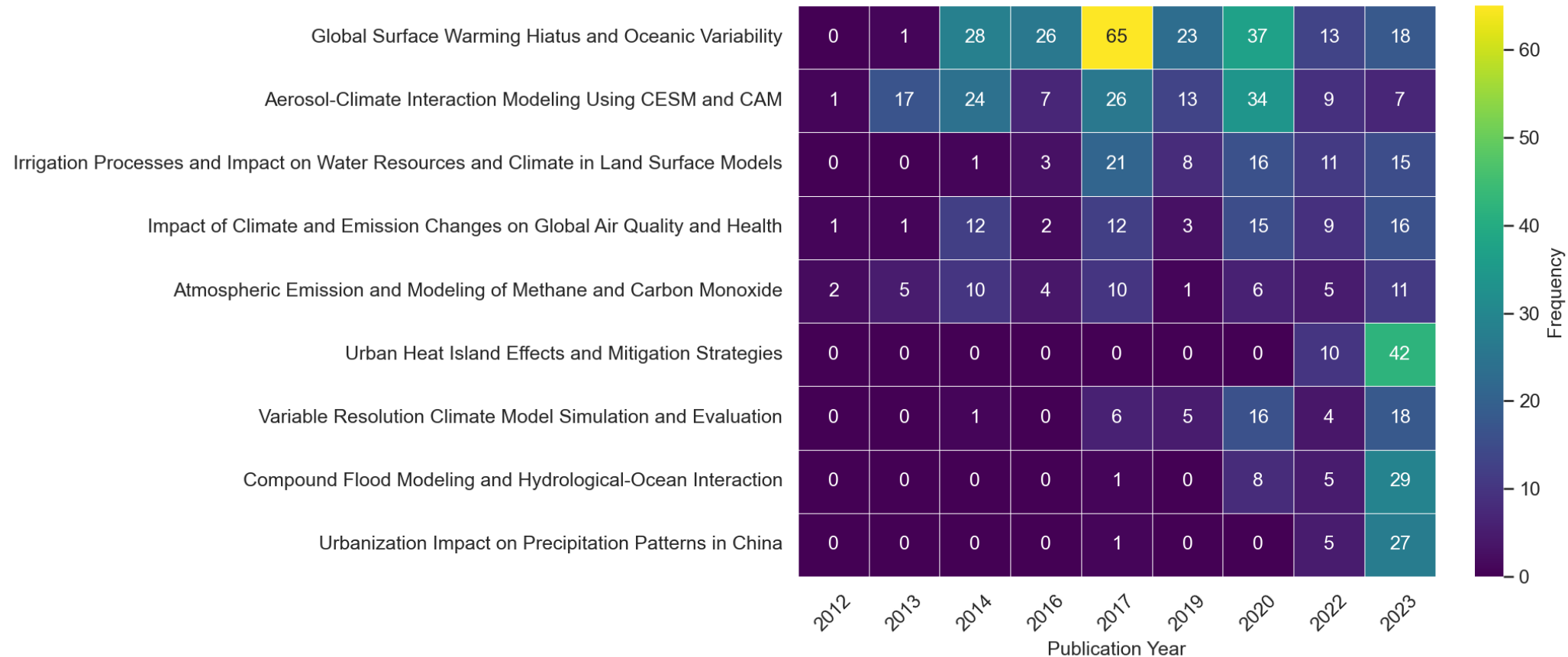
# RESEARCH THEMES FOR THOSE CITING EESM JOINTLY FUNDED PUBLICATIONS



- When we are jointly funded, which research communities are we impacting?
- Topic modeling (unsupervised) using the same semantically rich embedding model that ChatGPT uses to extract knowledge
- Seeks out common themes in research
- Generated 18 research themes, each having at least 10 publications to identify as a theme
- Let's look at the top 9 themes...

## RESEARCH THEMES FOR THOSE CITING EESM JOINTLY FUNDED PUBLICATIONS

We can look at the number of publications cited in each research theme over time to look at emergent and diminished themes



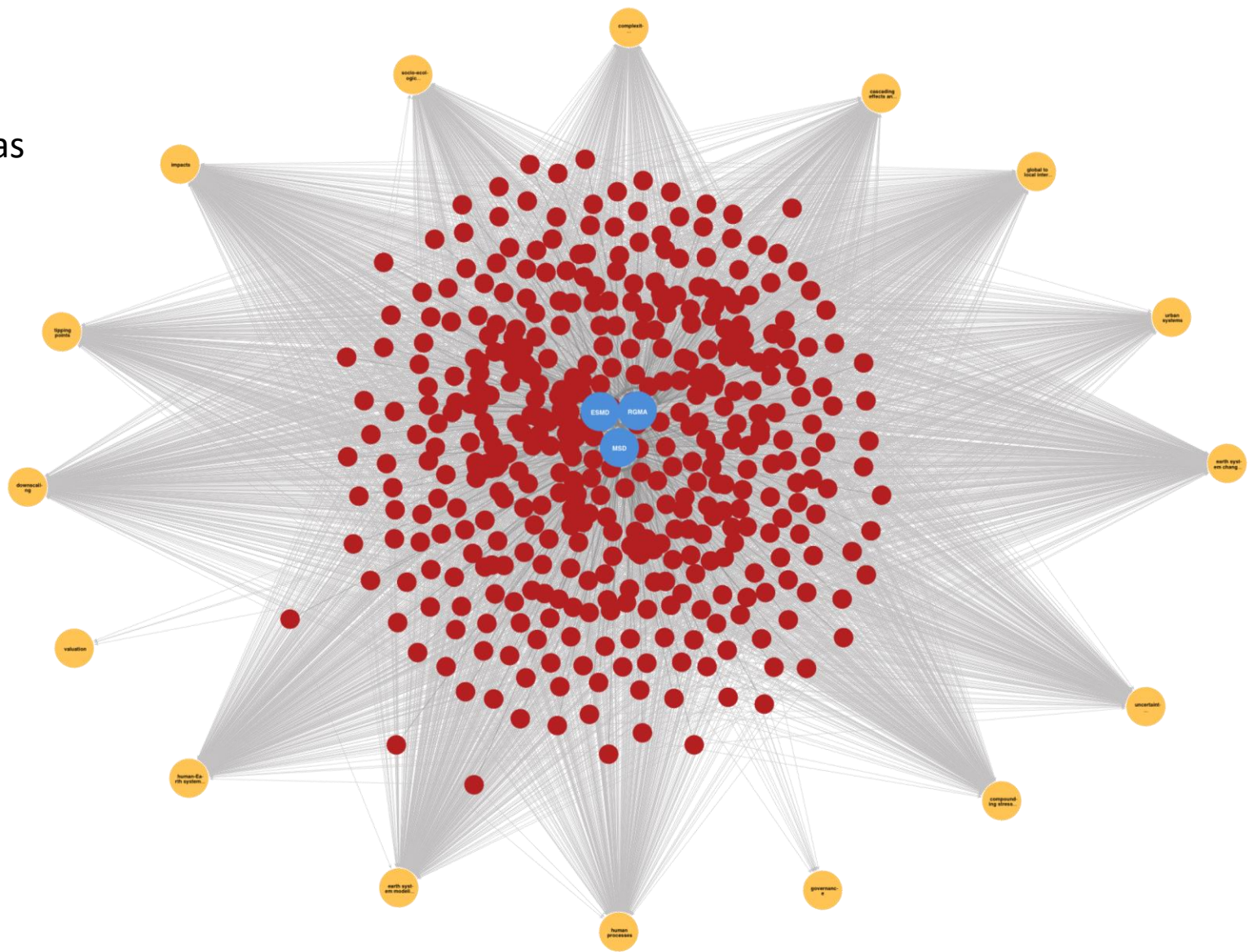


# COMMUNITY DETECTION TO IDENTIFY OPPORTUNITIES

- Explore the communities of researchers, institutions, program areas, etc. that are highly related to each other and various research areas
- Helps **identify groups of researchers who could potentially conduct research together** very easily – these groups may be completely unaware of each other
- Used to predict emerging research topics and potential growth patterns within MSD at large

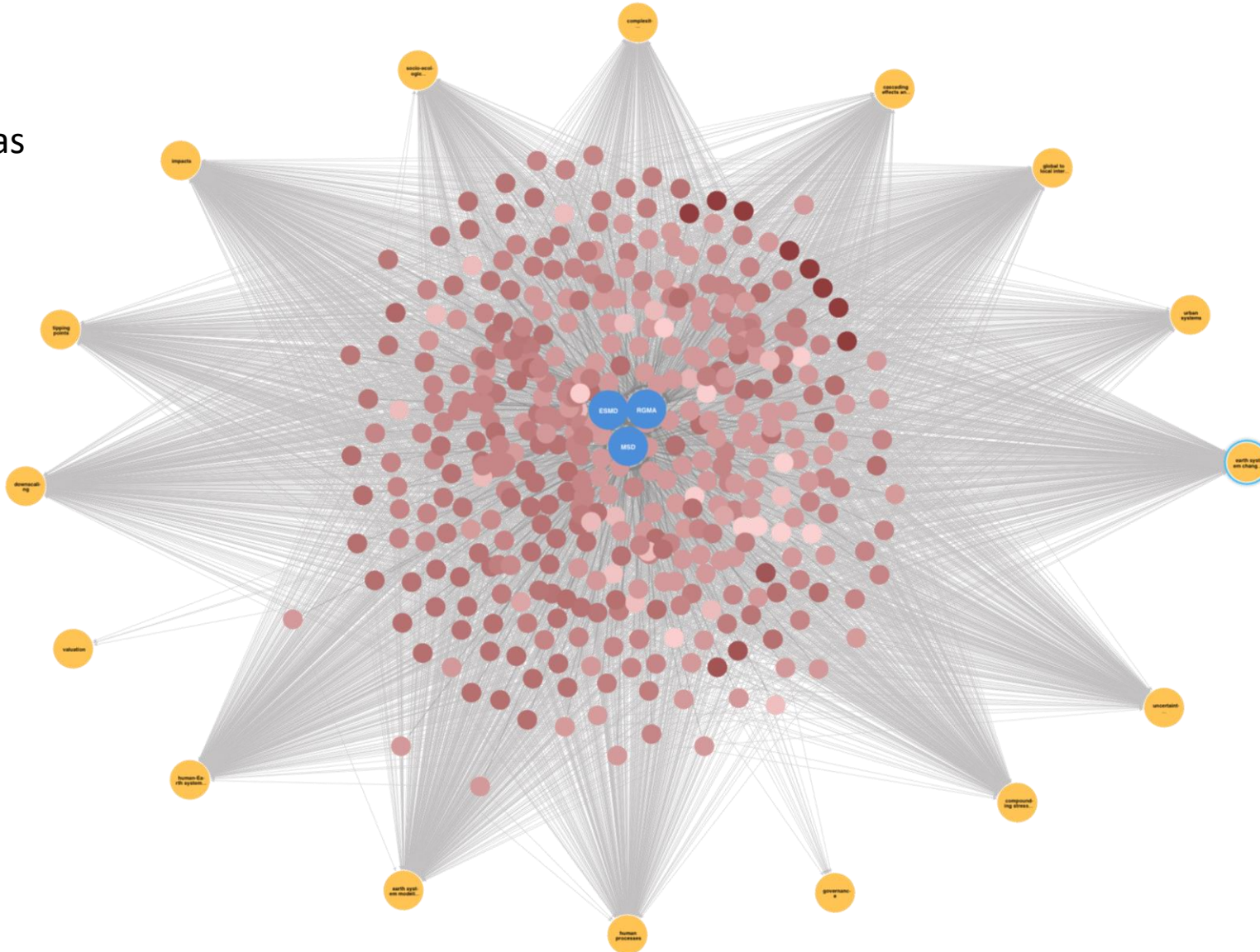


- EESM Program Areas
- Research Areas
- Authors



Preliminary Results – Snapshot of EESM funded research –  $P(X \geq 80)$  to MSD

- EESM Program Areas
- Research Areas
- Authors

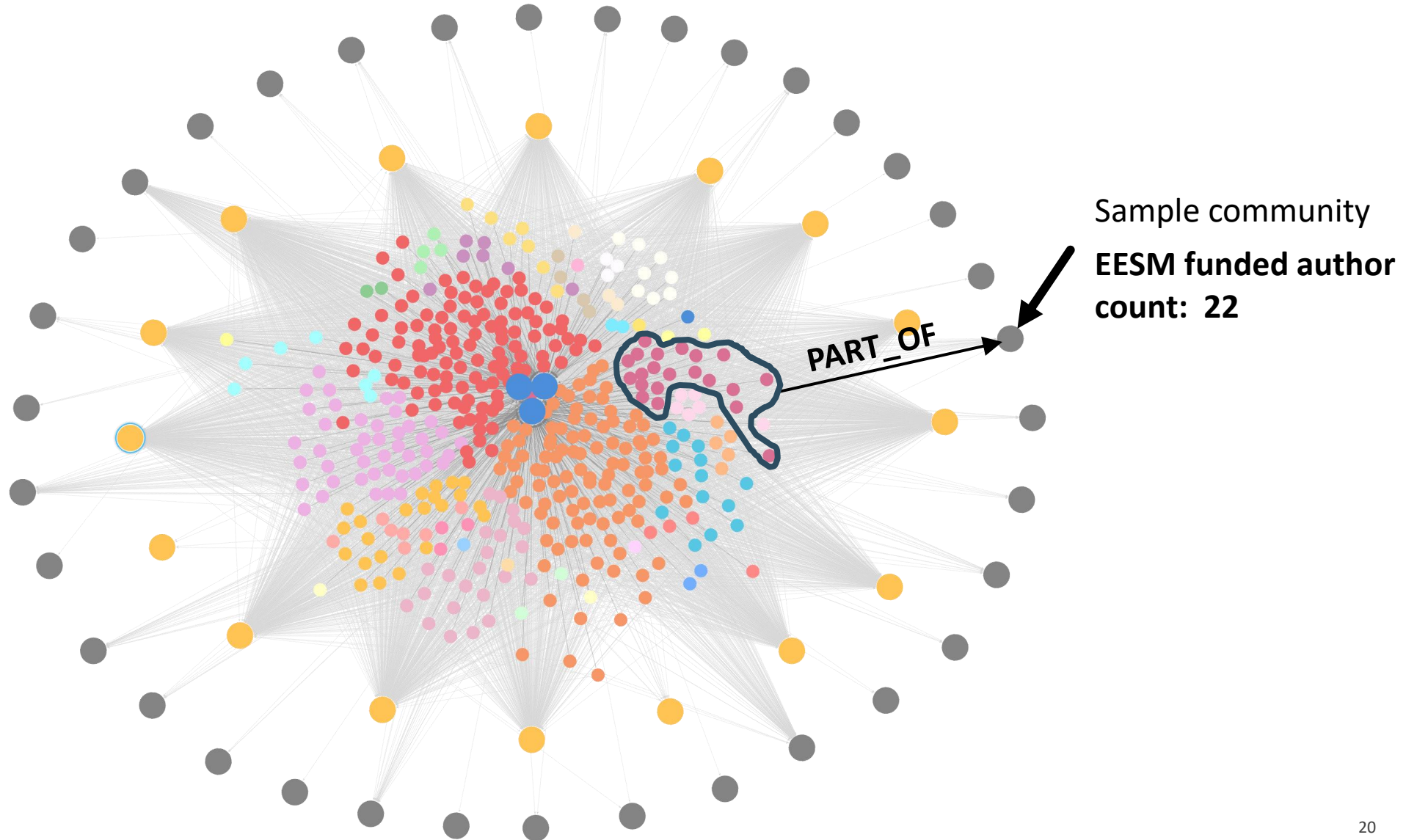


**Moving from single authors to communities of authors by existing collaboration and like research**

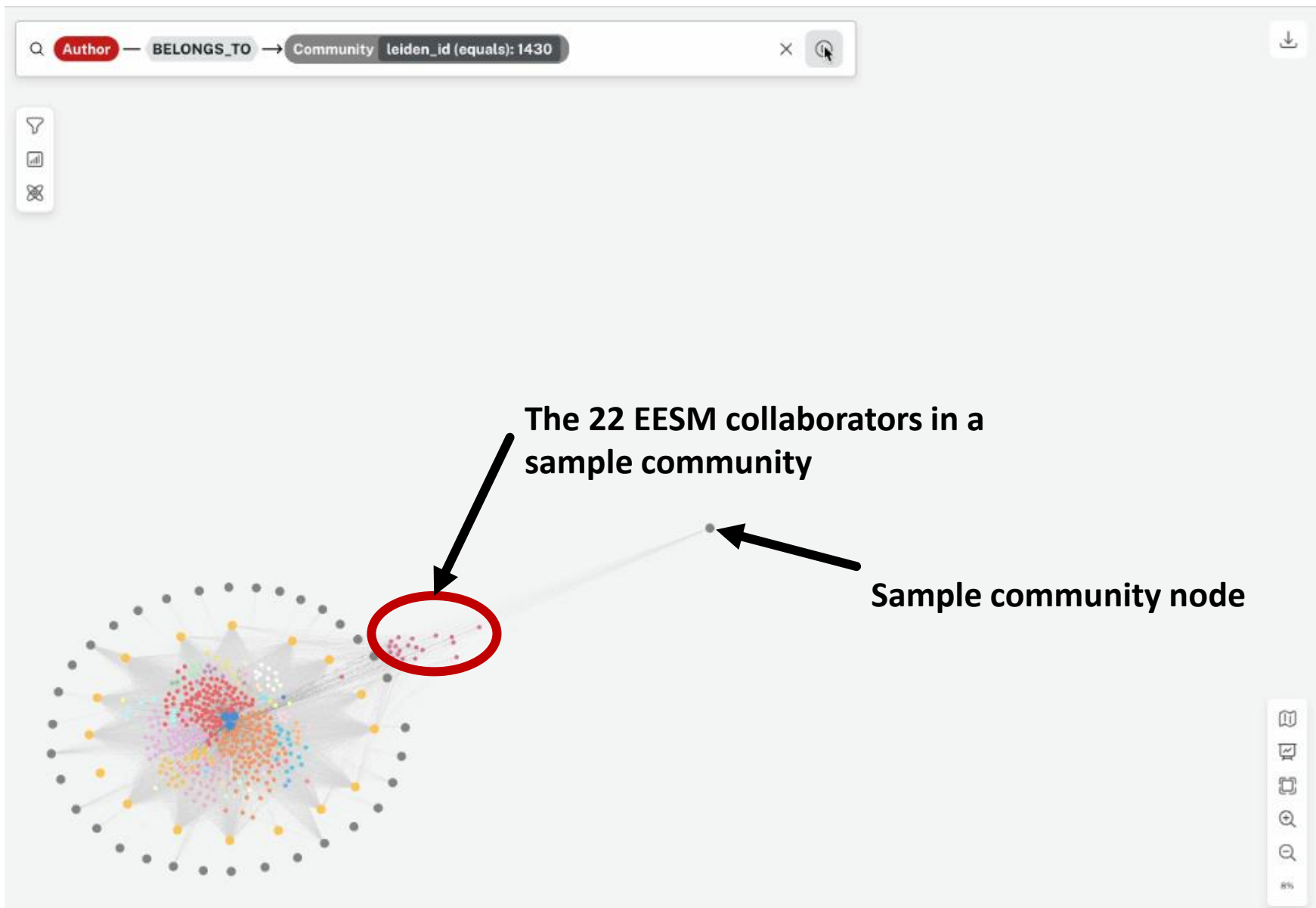


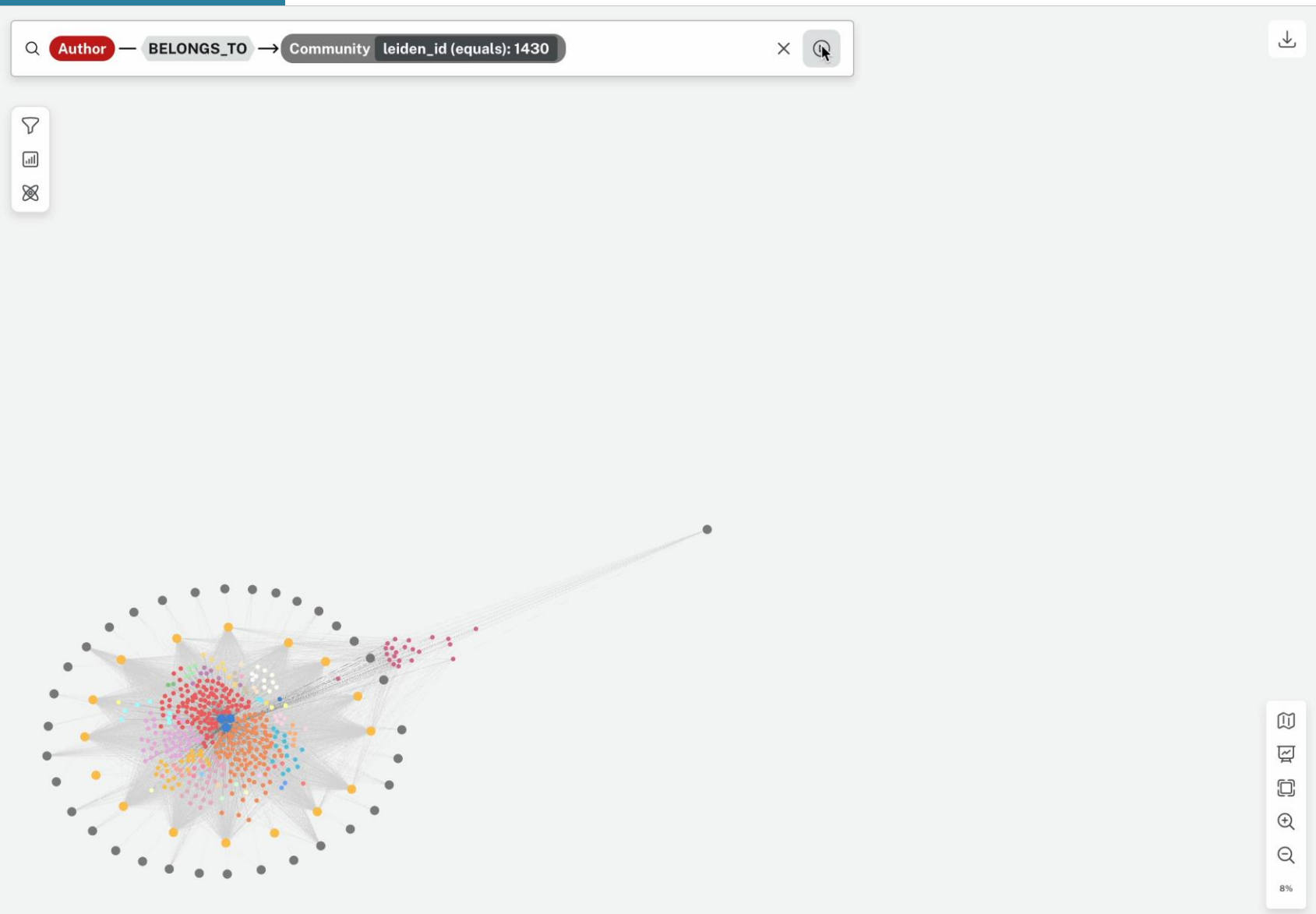
## COMMUNITY ANALYSIS

- EESM Program Areas
- Research Areas
- Community Nodes
- Authors in Communities
- Communities



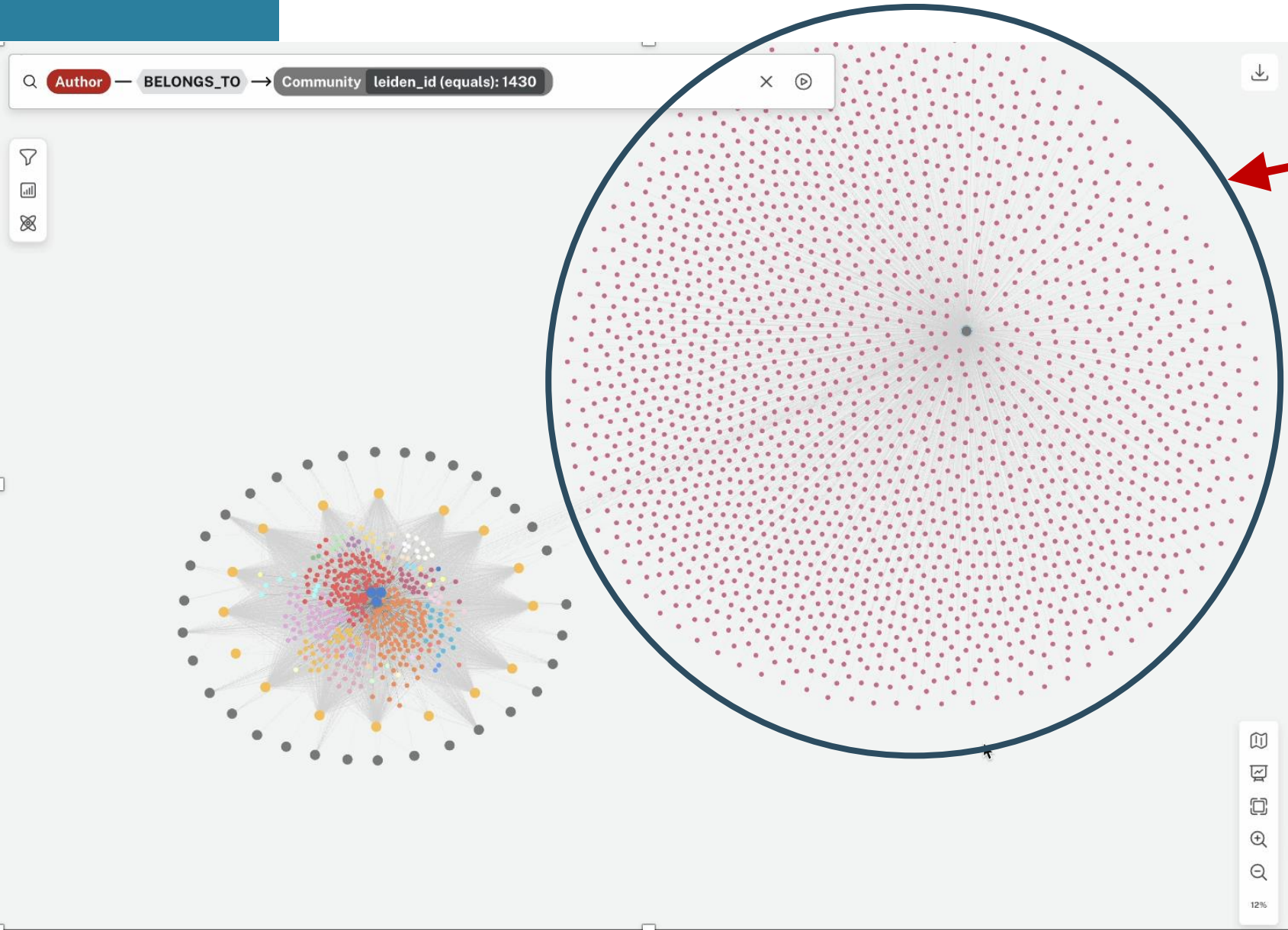
## COMMUNITY ANALYSIS





- **Quickly identify untapped potential** in our existing author communities! Can also do this for institutions, geographic areas (varying scale), topics within communities, and much more through time!
- Explore aligned / parallel research activities within, and outside of, BER to avoid redundancy and promote **informed collaboration**
- To be published in new Earth's Future Special Issue as an MSD CoP collaborative contribution





- **Quickly identify untapped potential** in our existing author communities! Can also do this for institutions, geographic areas (varying scale), topics within communities, and much more through time!
- Explore aligned / parallel research activities within, and outside of, BER to avoid redundancy and promote **informed, strategic collaboration**
- To be published in new Earth's Future Special Issue as an MSD CoP collaborative contribution



**THANKS!**