



Causal Influence Pathway Quantification on Social Networks

Abstract & System Flowchart

Quantifying and characterizing the spread of narrative content and sentiment from source communities to audience communities has important applications in communication and marketing as well as in national security such as countering disinformation. This poster presents a novel method to estimate the causal influence of pathways between individuals across communities on a social network. To enable data-driven influence quantification, discovery of narrative content and aspect-based sentiment classification are automated using large language models. We demonstrate the utility of this method in identifying influential accounts both as sources and bridges of the two opposing communities competing for the attention of the audience community, on the spread of Biolabs-in-Ukraine narratives, on Twitter in 2022. Validation using tweet following statistics, external corroboration, and predictions shows that causal influence reveals hidden influencers and is more accurate than existing metrics.

Influence Quantification System

Input
Social and news media data sources

Data Ingest
Targeted Queries

Targeted Collection

AI-Assisted Narrative Discovery

Narrative Network Characterization

Causal Influence Estimation

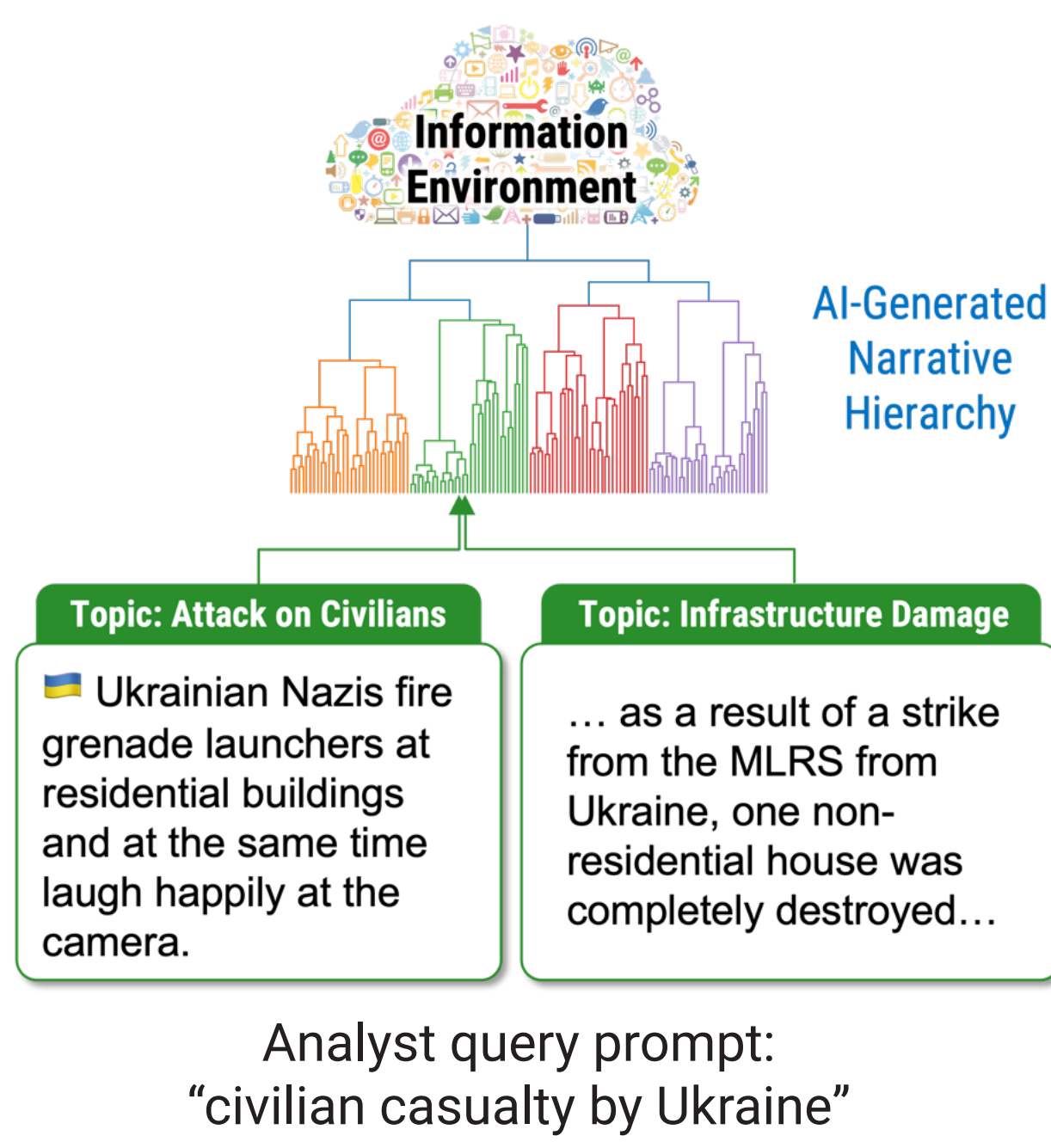
Output

- Narrative content
- Network mapping
- Causal influence
- Influence pathways
- Predictive influence

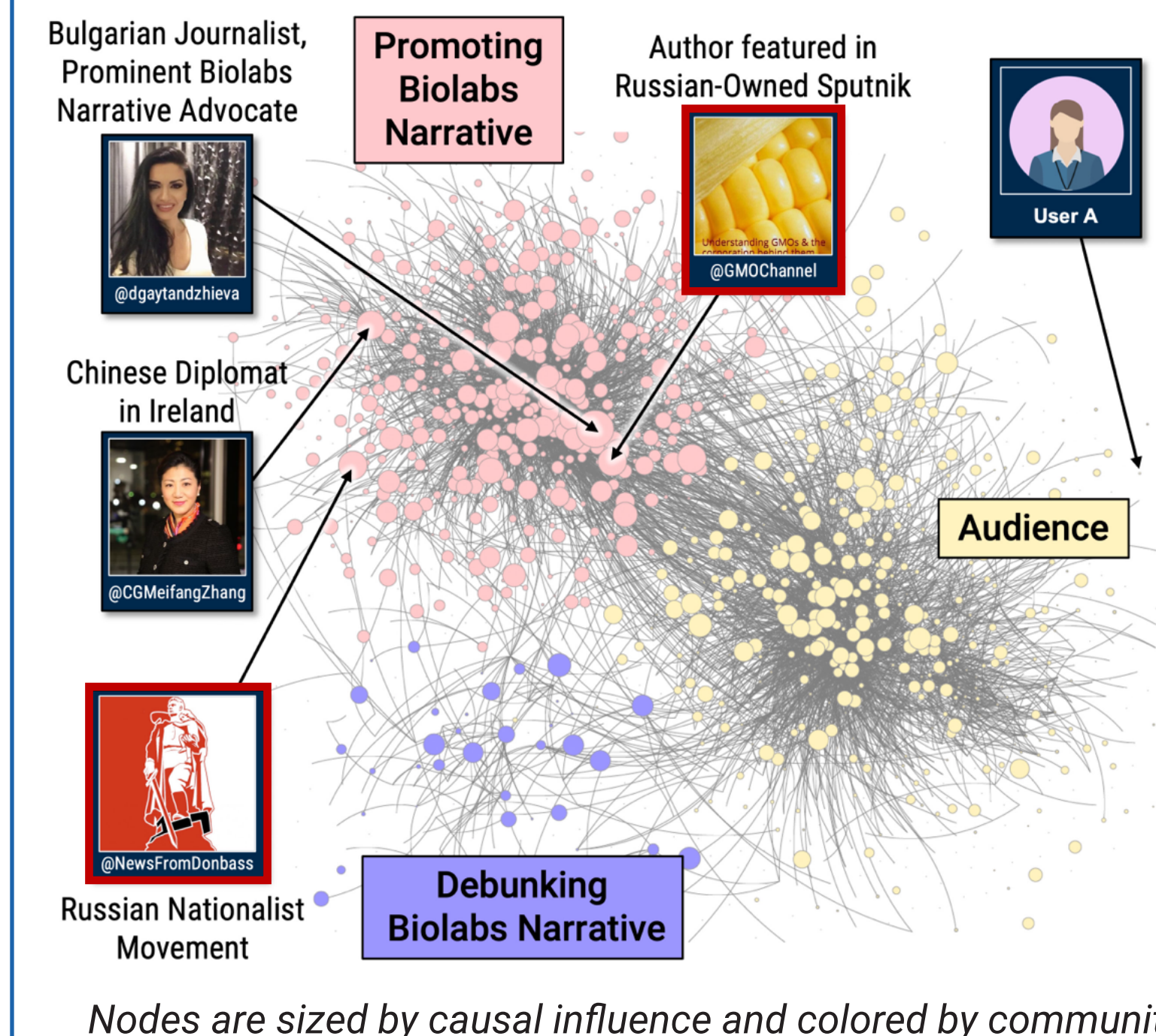
Human-AI Teaming for Narrative Discovery

- Large language model embedding puts texts with similar meaning in nearby semantic space
- Clustering in the semantic space reveals the information environment structure
- This semantic structure enables rapid discovery of narrative content without knowing specific keywords
- One hour of human-AI interaction identified 58k of Biolabs narrative tweets out of a 8.5M tweet corpus with 81% precision at 86% recall

AI enables content discovery by semantic meaning, instead of keyword search



Causal Influence Quantification Reveals Hidden Influencers¹



Screen name	T	RT	F	Earliest time	Pagerank Centrality	Causal Influence
@dgyatandzhieva	7.3	34.6	49k	01/25/2022	5.6	1,088
@CGMeifangZhang	7.0	16.7	38k	03/08/2022	7.6	822
@User A	8.9	0	800	03/09/2022	0.1	73
@NewsFromDonbass	3.1	0.4	3k	03/06/2022	0.2	846
@GMOChannel	4.7	1.0	5k	03/07/2022	0.2	822

Tweets (T), Retweets (RT), Followers (F)

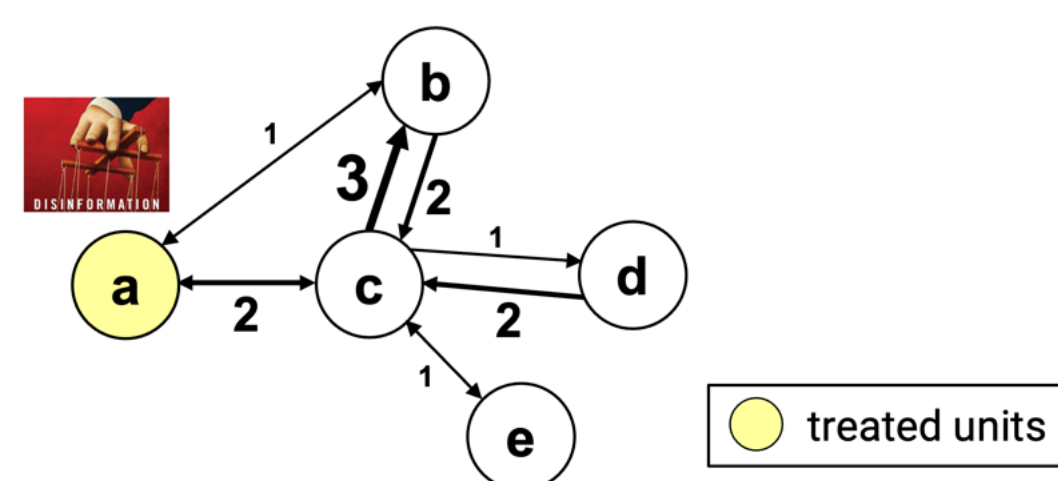
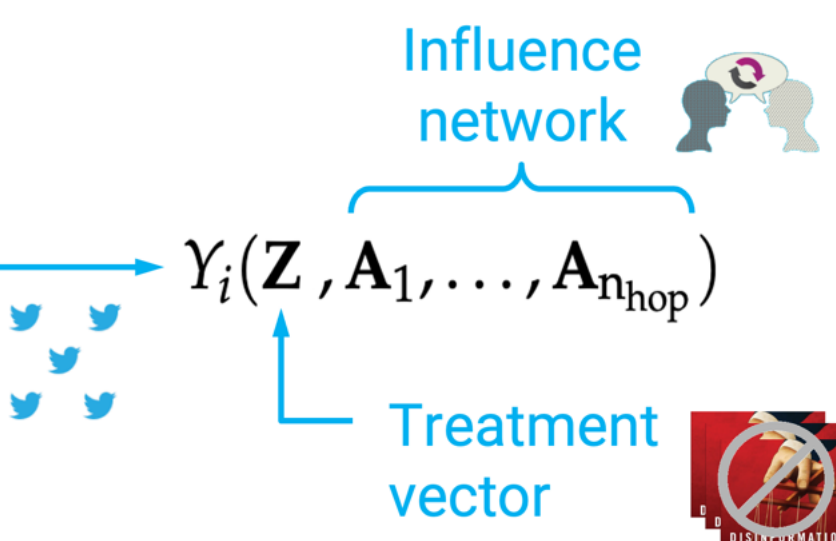
- Causal influence score identifies key influencers that do not stand out in traditional activity and network metrics
- Causal influence models tweet propagation on network and disentangles causal effects from social confounders

¹ Smith et al. Automatic Detection of Influential Actors in Disinformation Networks, Proceedings of the National Academy of Sciences (PNAS) (2021).

Causal Influence Pathway Quantification Using Network Potential Outcomes

Network Potential Outcomes:²

Network potential outcome of unit i (e.g. number of tweets on IO narrative)



Causal Influence of Account j over community m :

$$\zeta_j^m = \frac{1}{|C_m|} \sum_{i \in C_m} \zeta_{ji}$$

$$\zeta_{ji} \doteq Y_i(Z = z_{j+}, A_1, \dots, A_{n_{hop}}) - Y_i(Z = z_{j-}, A_1, \dots, A_{n_{hop}})$$

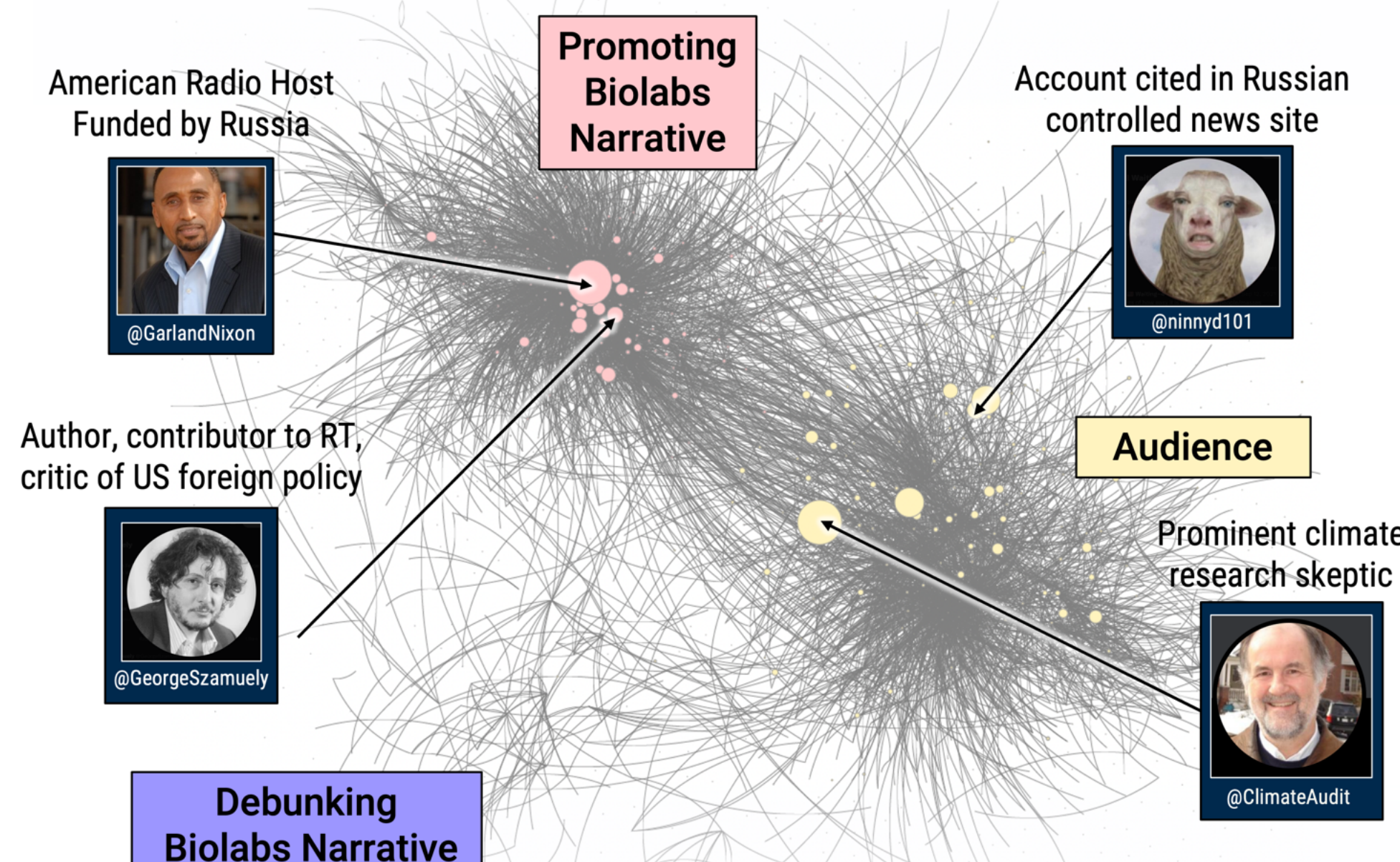
Example Path-Specific Causal Estimand – Bridge Influence Estimand:

$$\zeta_{ji}^{\ell bm} = \sum_{j \in C_\ell} \frac{1}{|C_m|} \sum_{i \in C_m} \zeta_{ji}^{\ell bm}$$

Influence network manipulation counterfactual: Outcome of account i if paths between community ℓ and community m via bridge account b did not exist

$$\zeta_{ji}^{\ell bm} \doteq Y_i(Z = z_{j+}, A_1, \dots, A_{n_{hop}}) - Y_i(Z = z_{j+}, A_1^{\bar{c}_{\ell \rightarrow b}}, \dots, A_k^{\bar{c}_{\ell \rightarrow b} \rightarrow b \rightarrow C_m}, \dots, A_{n_{hop}}^{\bar{b} \rightarrow C_m})$$

Top Bridge Accounts from Biolabs-Promoting to Audience Community



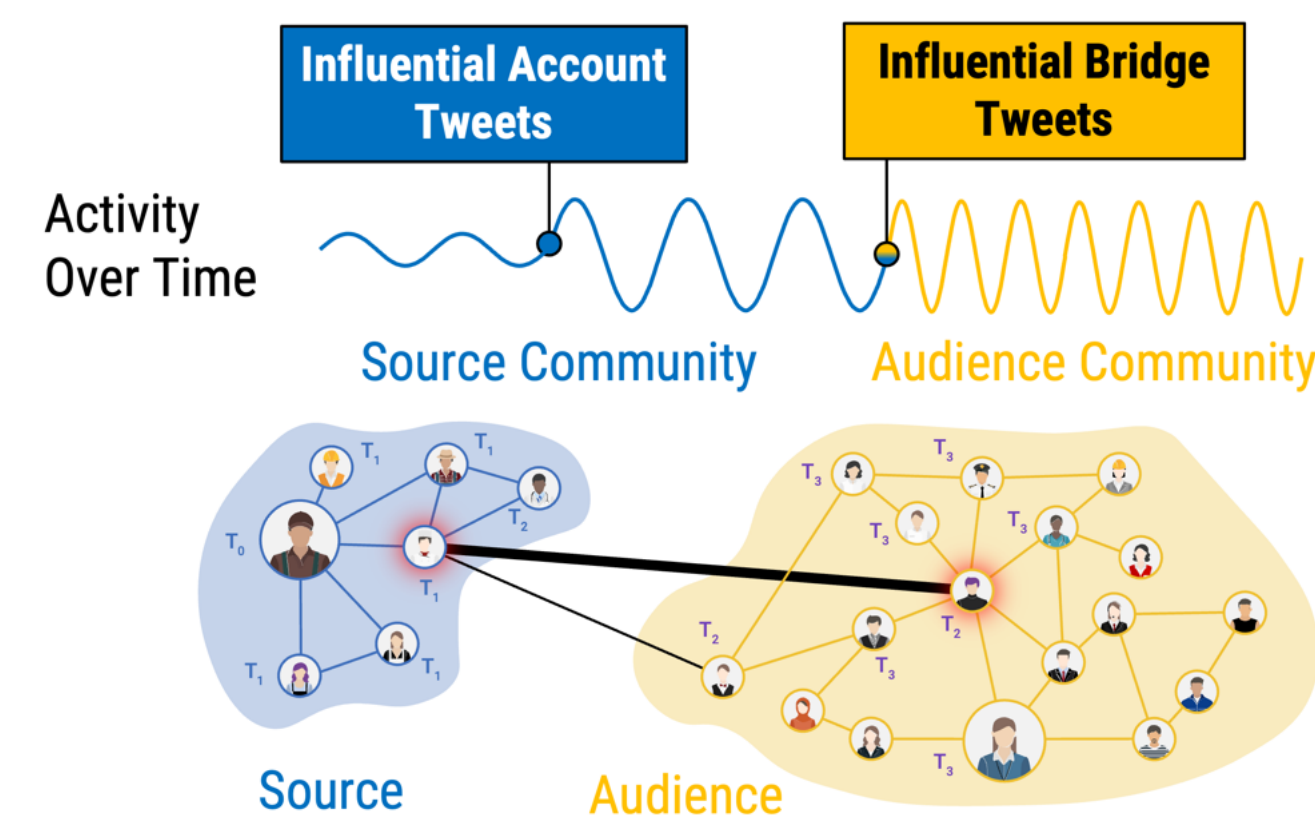
Nodes are sized by causal influence and colored by community

² Kao, Causal inference under network interference: A framework for experiments on social networks. Ph.D. Thesis, Harvard University (2017).

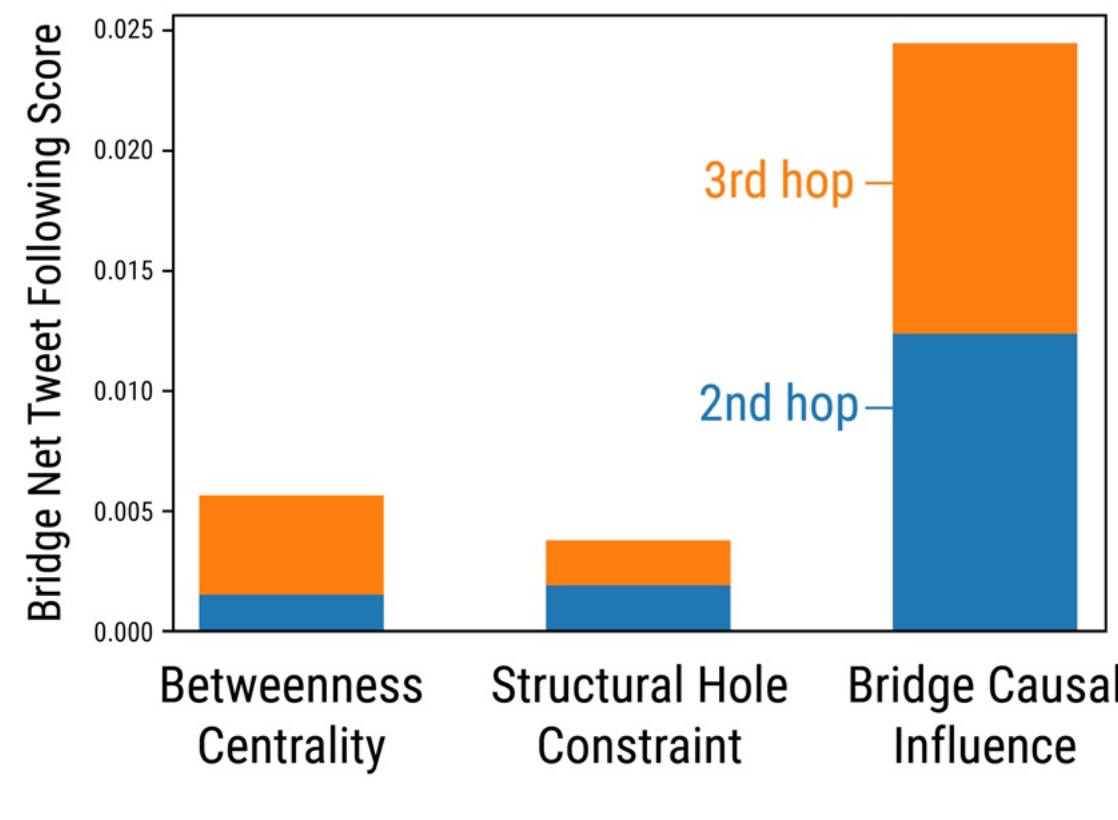
Validation with Net Tweet Following

Bridge Net Tweet Following Metric

How prominent is a bridge in the paths of net tweet following?



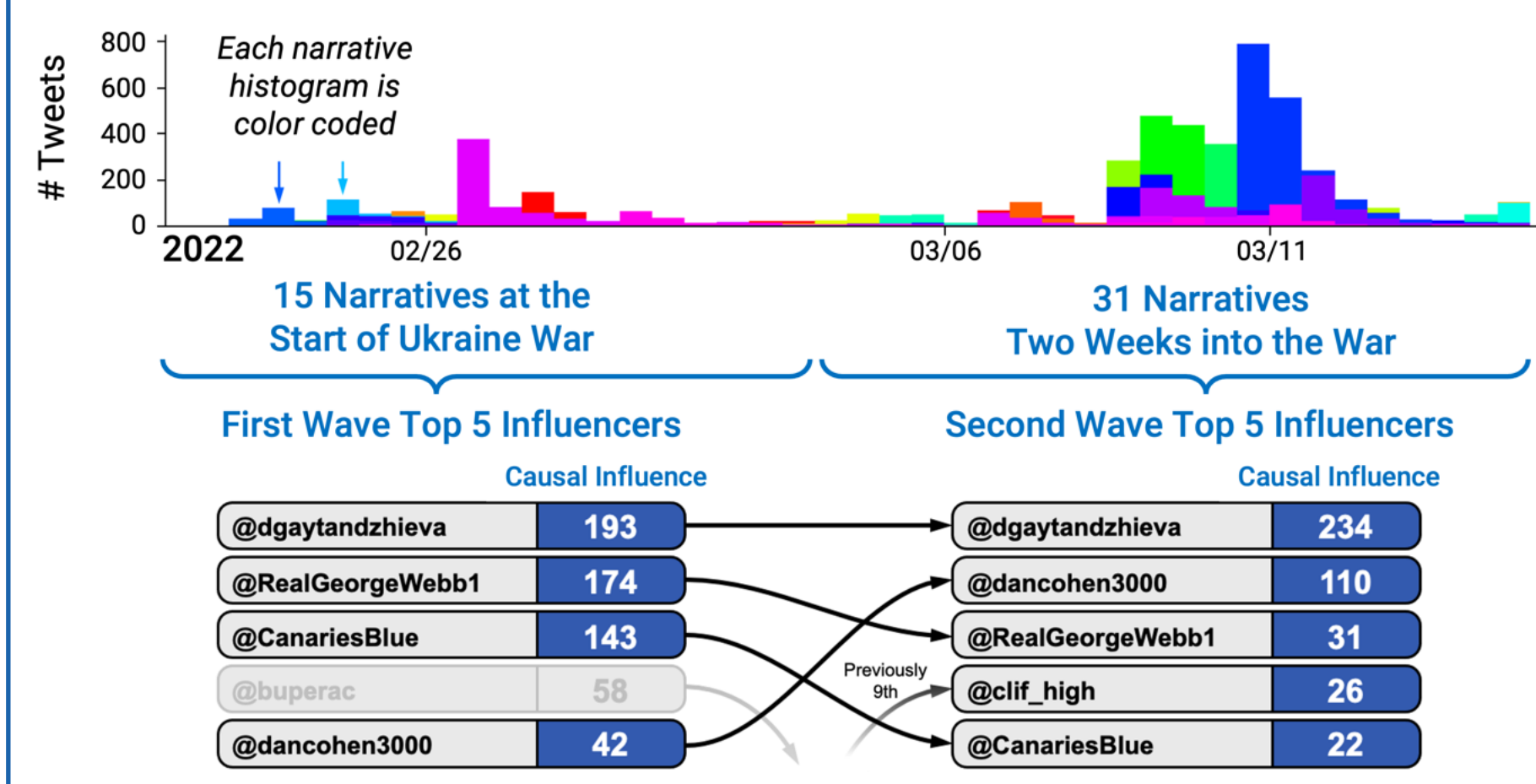
Bridge Net Tweet Following of Top 50 Accounts From Biolab Debunking to Audience Community



Bridge causal influence identifies accounts with stronger evidence of source message amplification in audience community than baseline methods

Results on Predicting Causal Influence and Outcomes

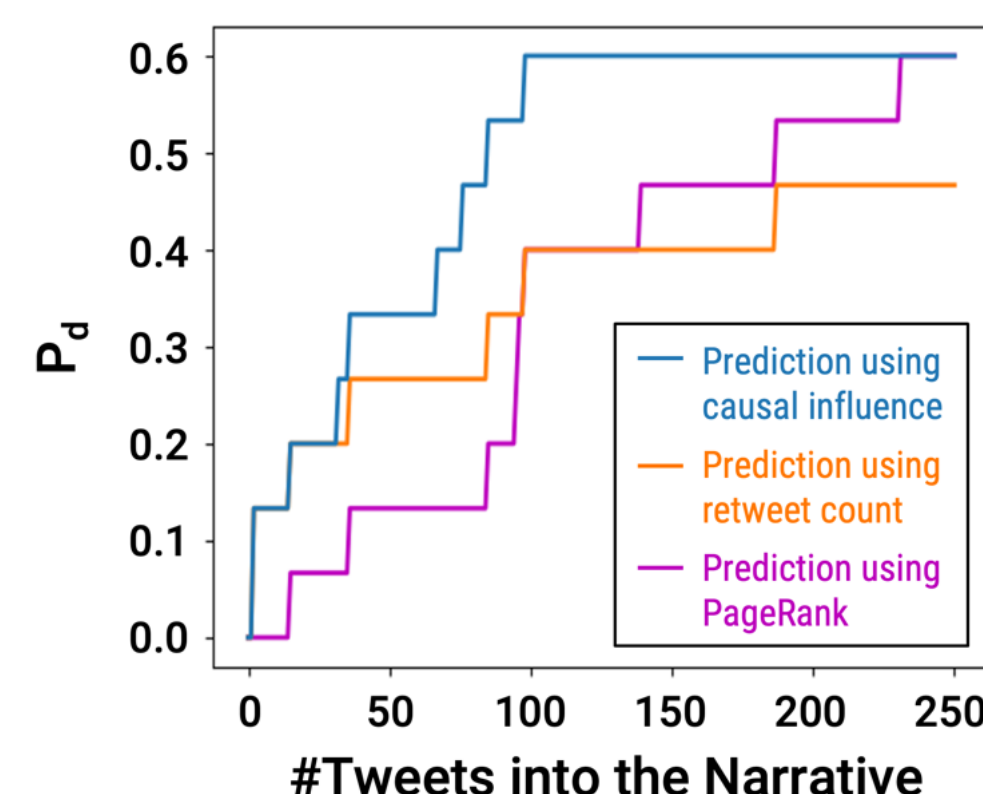
First and Second Waves of Biolabs Narratives



Past causal influence is predictive of future influence

Early Detection of Rising Narratives

Probability of detection (P_d) at constant zero false alarm rate



Participation of past causal influencers is predictive of future rising narratives