



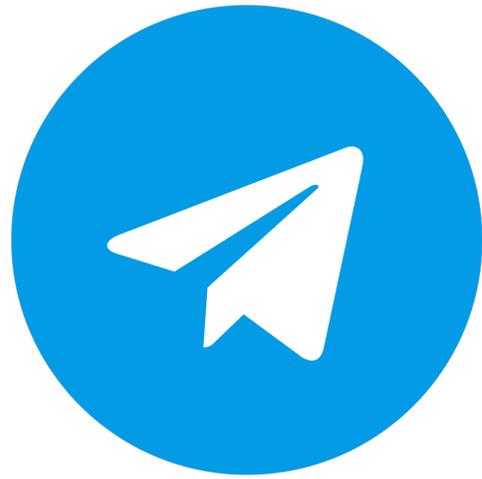
Mapping **Dark Platforms**: Modelling Information Flows on Telegram

Rupert Kiddle¹, Mónika Simon¹, Kasper Welbers², Anne Kroon¹, Damian Trilling¹

1: The University of Amsterdam (UvA)

2: The Vrije University (Amsterdam)





Telegram



Private: end-to-end encrypted messaging.



Public: *chats* (many-to-many) and *channels* (few-to-many).

A home for the **deplatformed**: permissively governed, low-moderated space.

Relatively free of **algorithmic interference**.

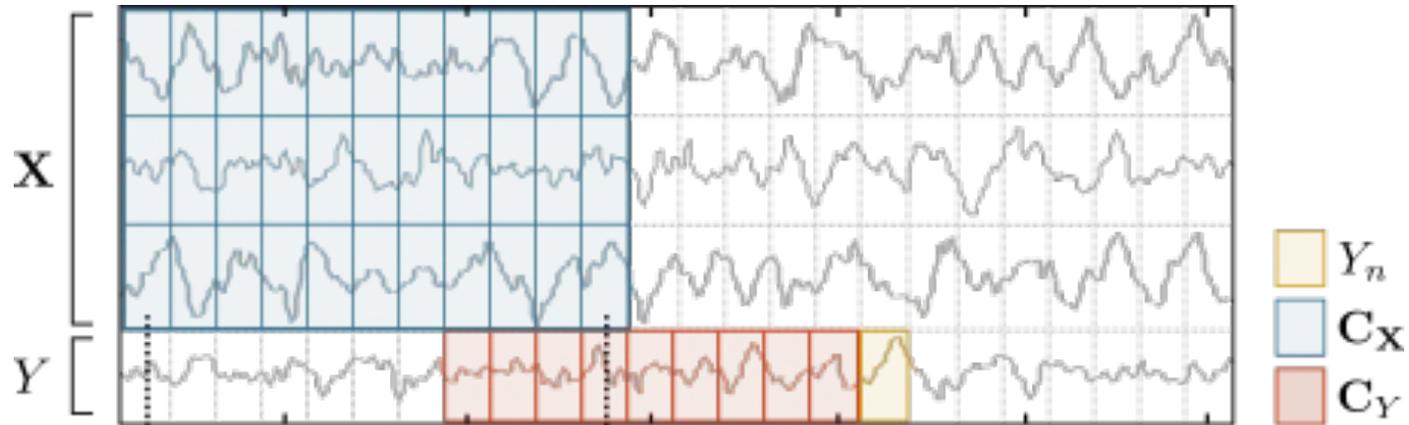


Information about network structure is very limited!

What we have: full messaging text-streams of chats/channels.

What we lack: relational data describing user, chat/channel associations.

from agents to flows -
Multivariate Transfer Entropy (mTE) -



IDTxl package for Python: P. Wollstadt, J. T. Lizier, R. Vicente, C. Finn, M. Martinez-Zarzuela, P. Mediano, L. Novelli, M. Wibral (2018). *IDTxl: The Information Dynamics Toolkit xl: a Python package for the efficient analysis of multivariate information dynamics in networks.* Journal of Open Source Software, 4(34), 1081. <https://doi.org/10.21105/joss.01081>.

mTE applied - **‘Network Toxicity’**

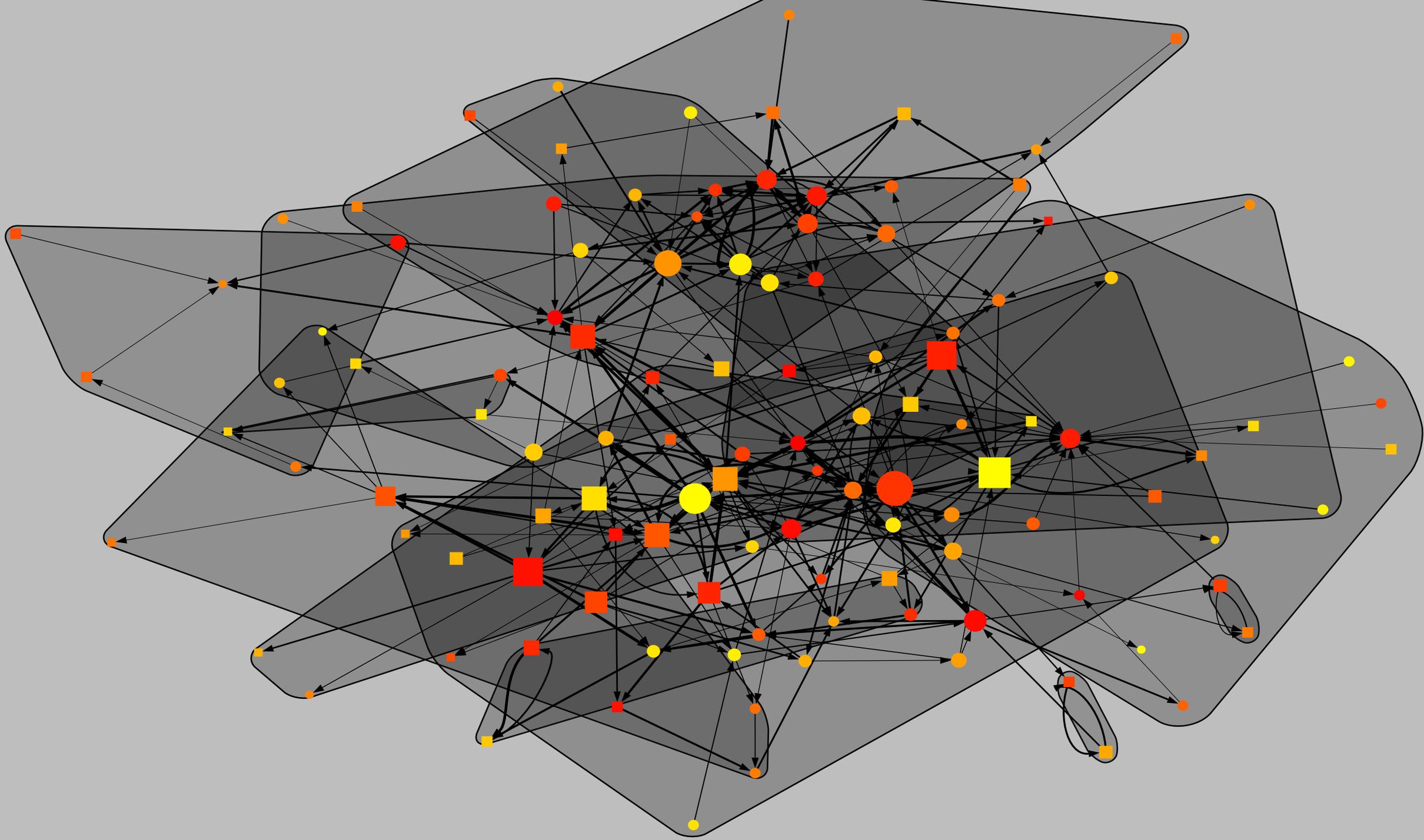
- **Aim:** to infer an *effective network* that describes the flows of toxicity between chats/channels on Telegram.
- **Toxicity:** *hateful, hurtful or marginalizing language that makes a person want to leave a conversation.*

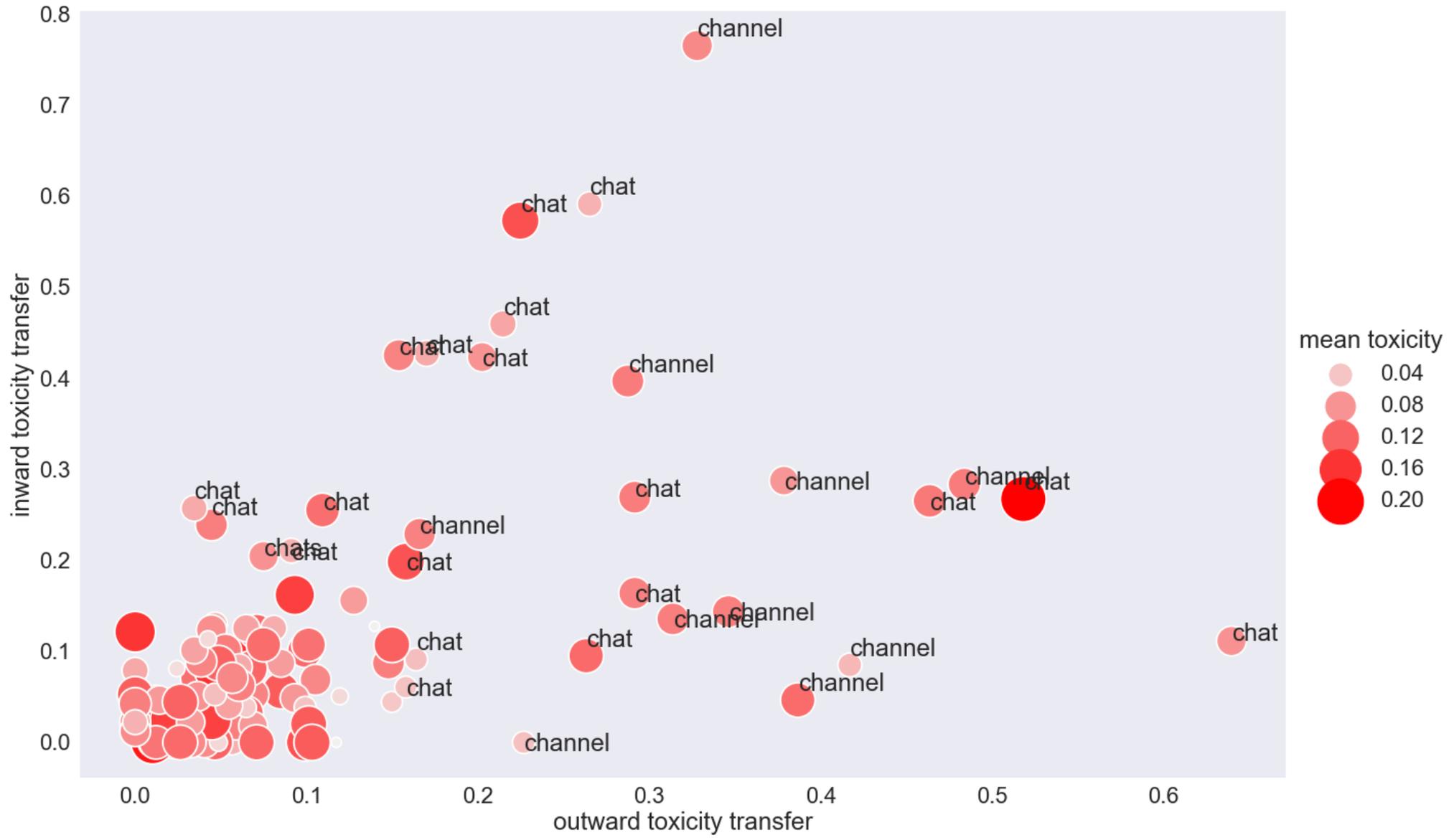


- > A growing threat to individuals & platforms.
- > Telegram users are particularly at-risk.
- > Existing approaches would struggle here.

Method -

1. **Data:** 6m snapshot (1M messages) of Dutch 'Telegramsphere' ---- *Jan '21 to Jun '21*.
2. **Classification:** messages \rightarrow *toxicity scores* with the *Perspective API (Google)*.
3. **Resampling:** *5-minute windows*; average of scores; 24-hour leeway.
4. **Estimation:** linear estimator using **IDTxI**, *30-minute maximum lag*.
5. **Aggregation:** *n for each i,j / $\max(n)$ for any i,j* = 'toxicity transfer'.
6. **Pruning:** disparity filter: 3% of edges retained ($p=0.05$).





Concluding Remarks -

- A '*model-free*' approach to understand the **dynamic consequences of toxicity**.
 - + non-linear estimators (more complex relationships)
 - + different timescales (beyond 30 minutes)
 - + validation methods ('unpacking')
- **mTE framework** for **information flows & influence** more broadly.
 - + different representations (e.g. word embeddings)
 - + beyond Telegram..

Thank you very much!



@rptkiddel



@newsflows_erc