

Joint Non-negative Matrix Factorization for Learning Ideological Learning on Twitter

Preethi Lahoti¹, Kiran Garimella², Aristides Gionis²

¹ Max Planck Institute for Informatics, ² Aalto University

Access to Diverse Information Around World ...



Filtered and Cherry Picked Content ...



Twitter

(Retweet or Follow other users)



User – User
social graph

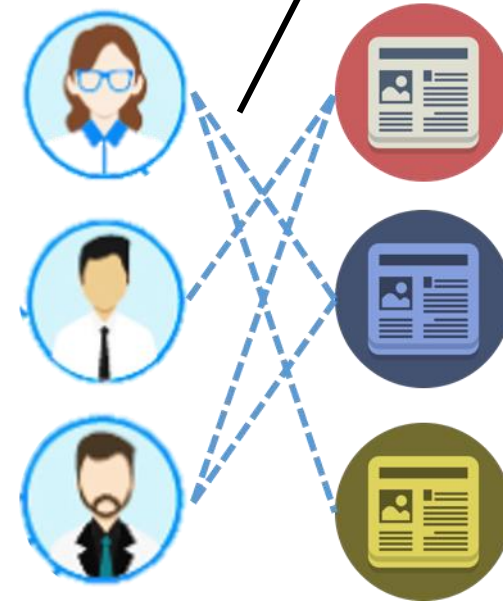
Twitter

(Retweet or Follow other users)



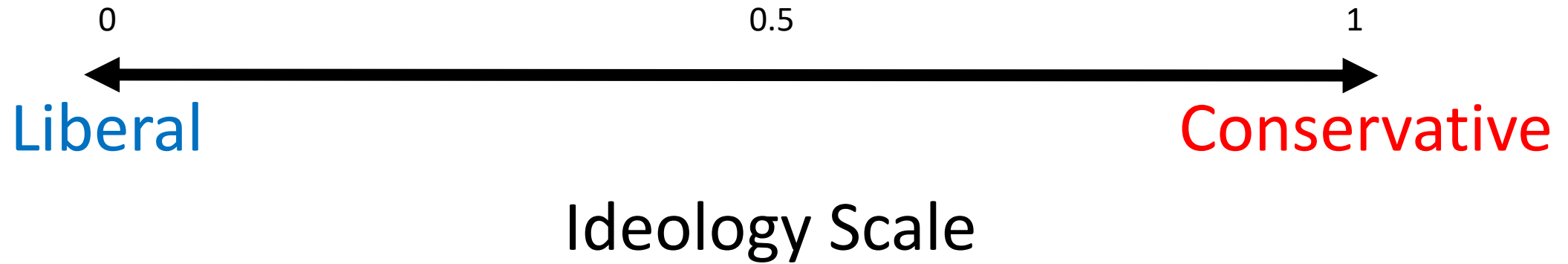
User – User
social graph

(Retweet a news article)



User - Content
content graph

User and Content Ideology



User and Content Ideology

Clinton puts Trump on defense at first debate
By [Stephen Collinson, CNN](#)
Updated 1517 GMT (2317 HKT) September 27, 2016

DONALD TRUMP WINS SECOND DEBATE; CNN SAYS IT DOESN'T MATTER



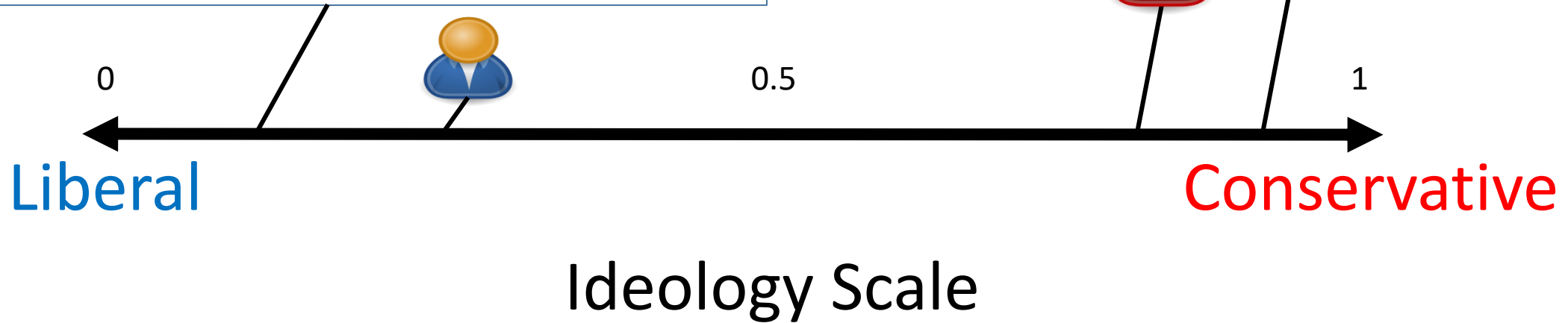
Ideology Scale

User and Content Ideology

Clinton puts Trump on defense at first debate

By Stephen Collinson, CNN
Updated 1517 GMT (2317 HKT) September 27, 2016

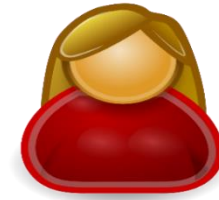
DONALD TRUMP WINS SECOND DEBATE; CNN SAYS IT DOESN'T MATTER



Filter bubble...



User A



User B

Topic: Presidential Debate

Filter bubble ...

Clinton puts Trump on defense at first debate



By Stephen Collinson, CNN

Updated 1517 GMT (2317 HKT) September 27, 2016



User A



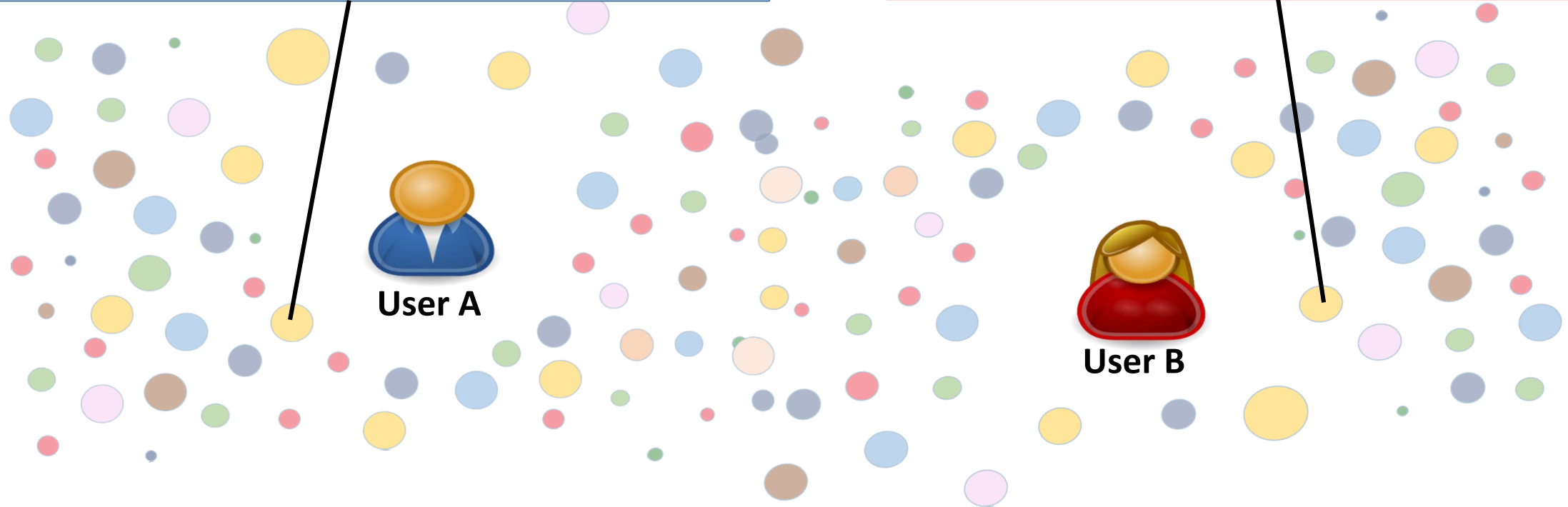
User B

Topic: Presidential Debate

Filter bubble...

Clinton puts Trump on defense at first debate
By Stephen Collinson, CNN
Updated 1517 GMT (2317 HKT) September 27, 2016

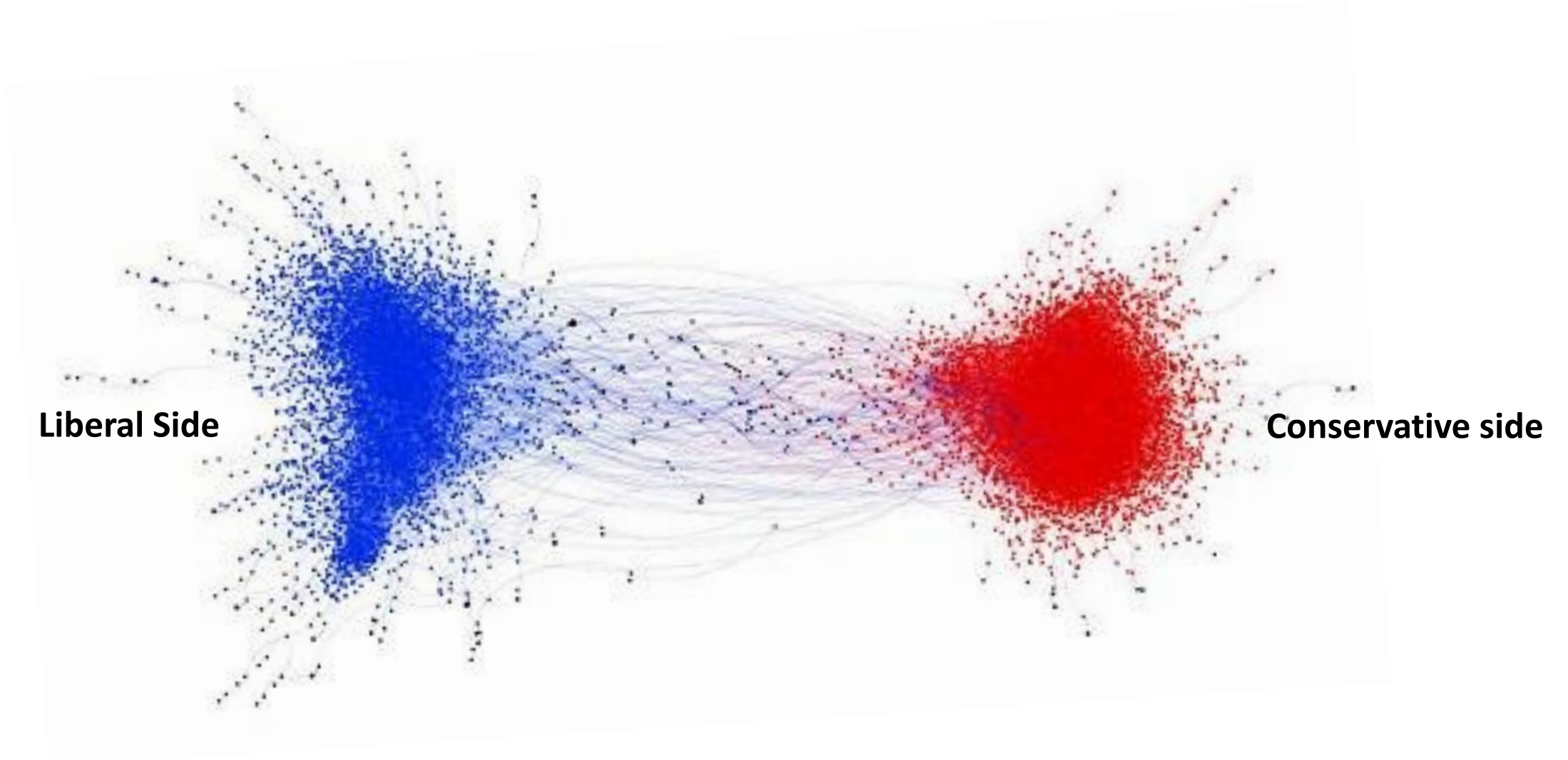
DONALD TRUMP WINS SECOND DEBATE; CNN SAYS IT DOESN'T MATTER



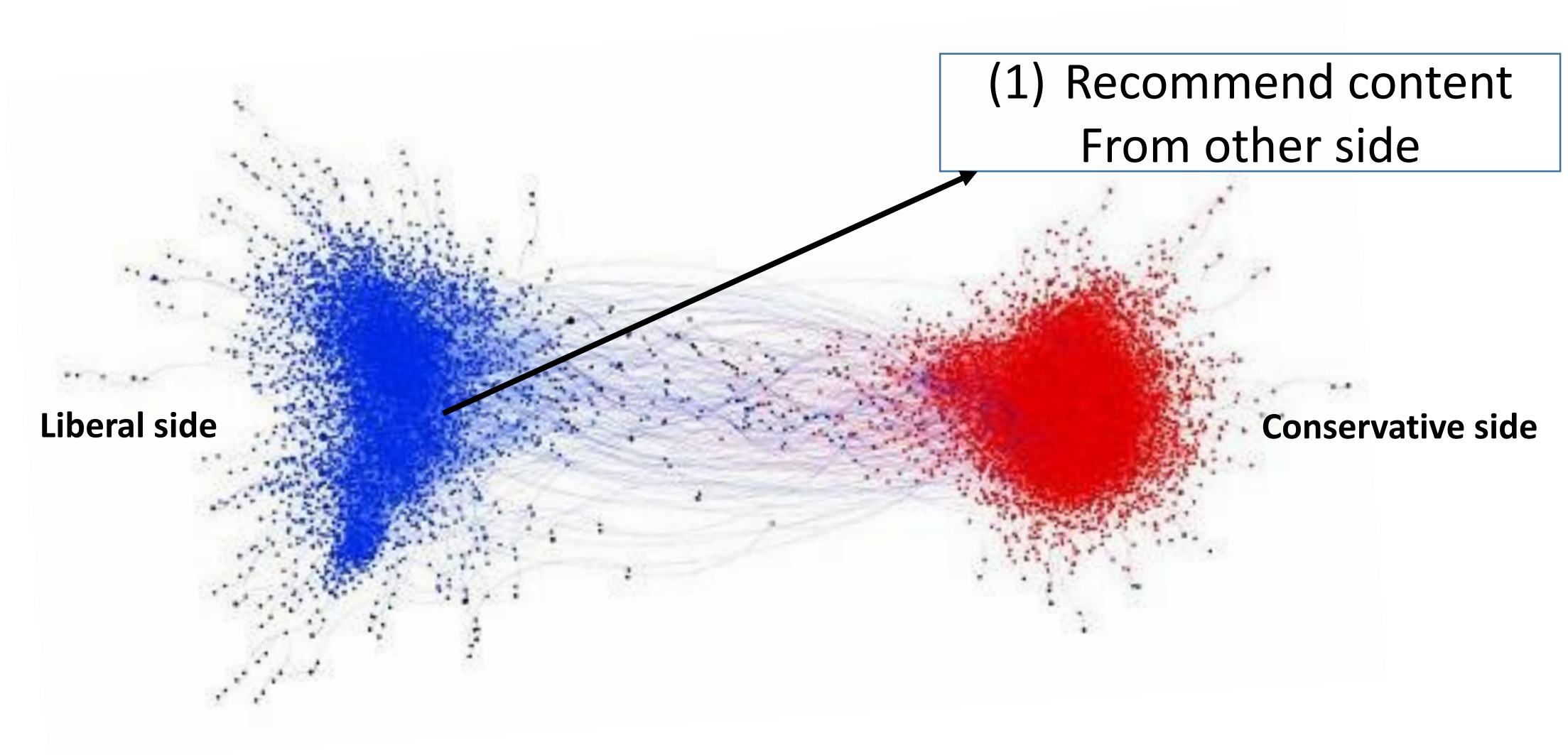
Topic: Presidential Debate

Motivation

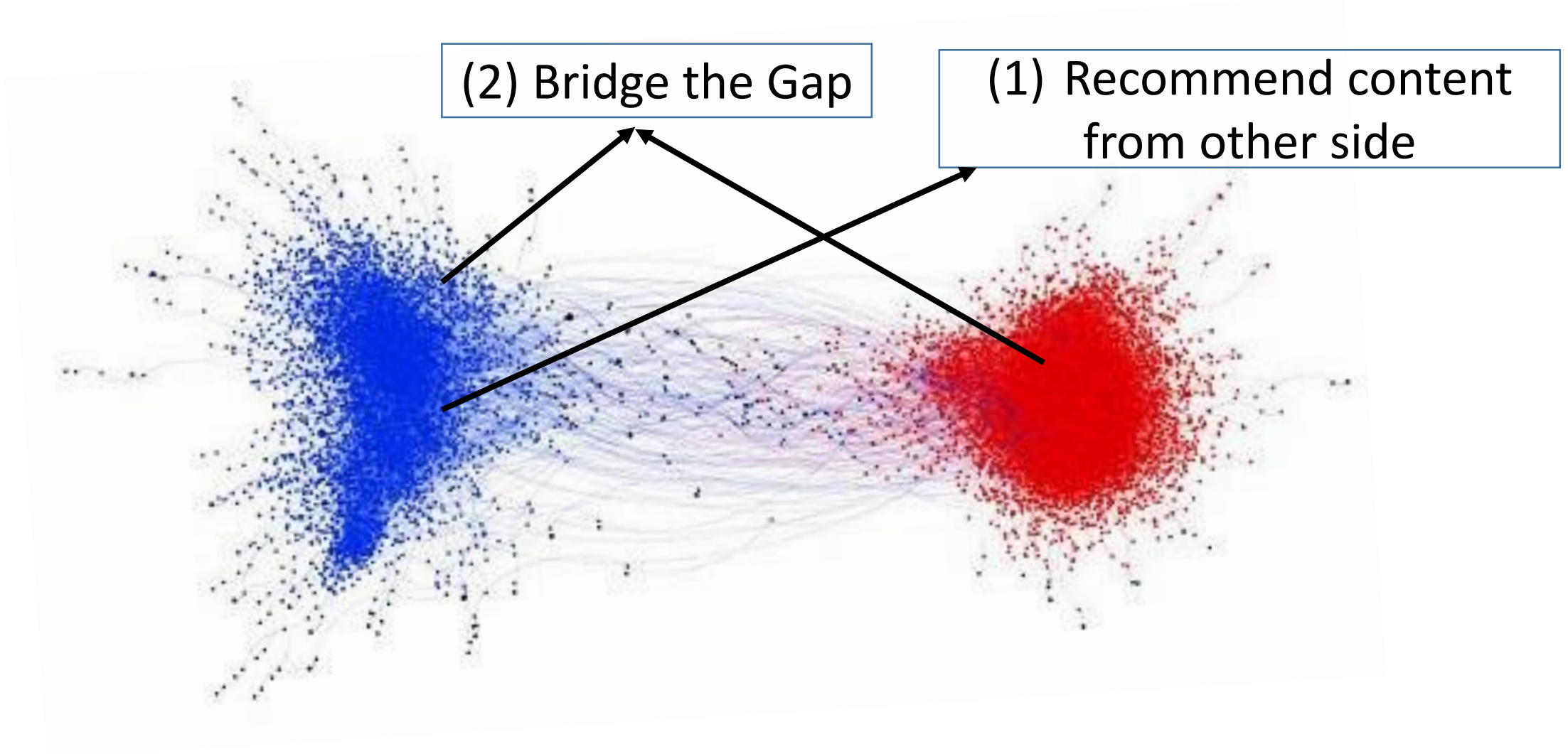
Ideological segregation, polarization, biased views



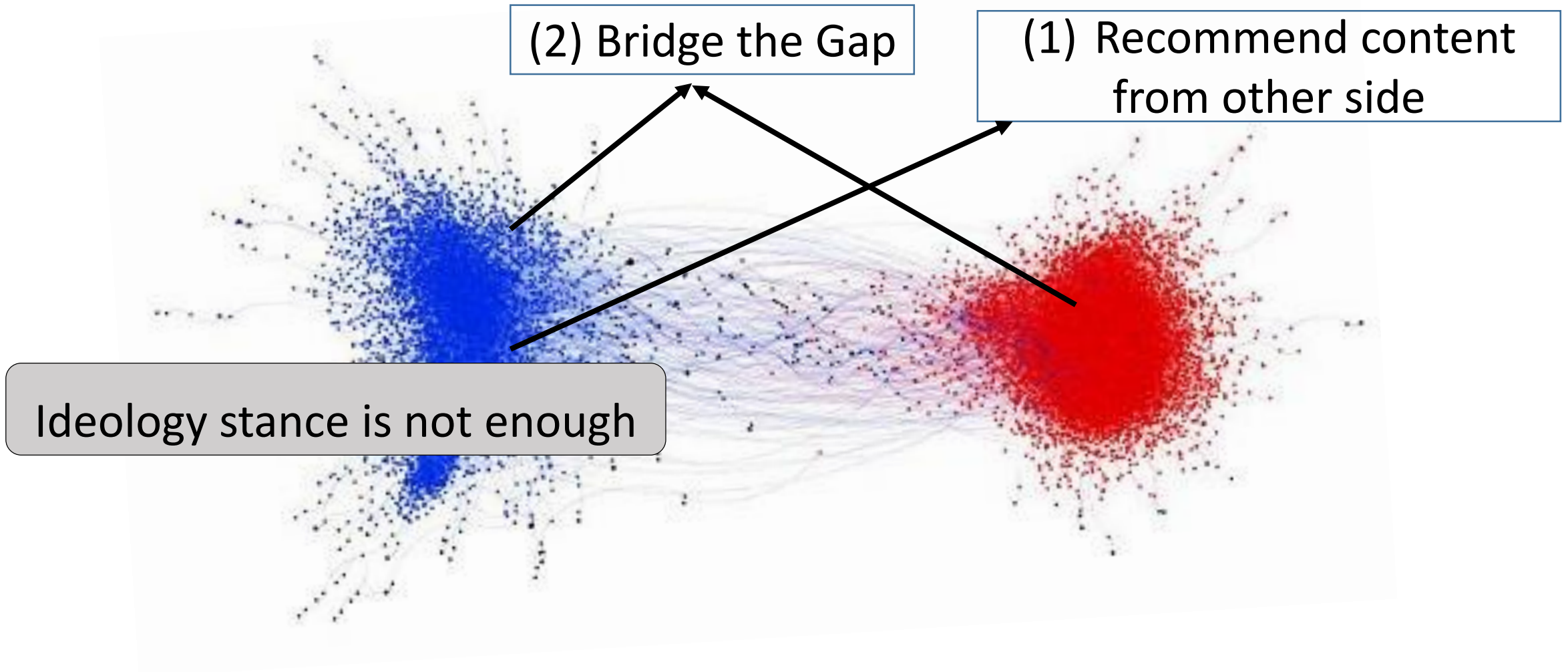
Bursting the filter bubble...



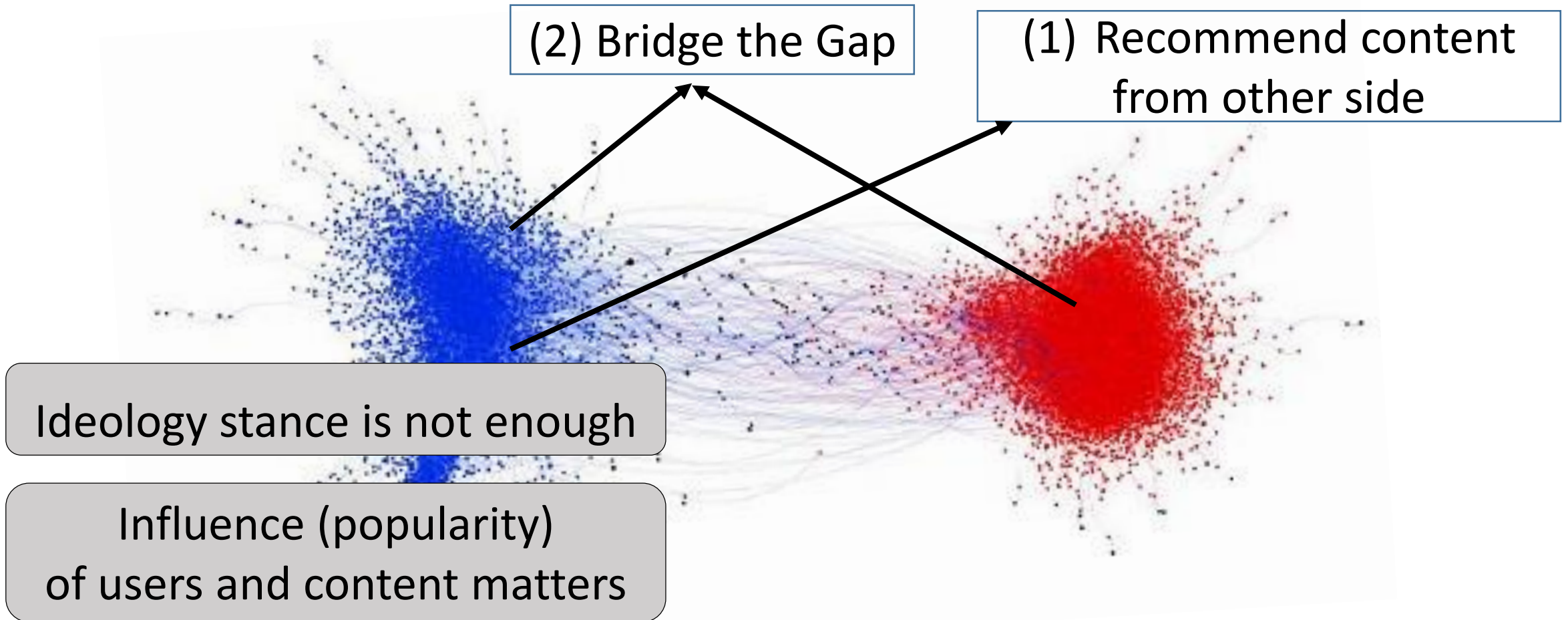
Bursting the filter bubble...



Bursting the filter bubble...



Bursting the filter bubble...



Problem Statement

- Input:



social graph
(A)



content graph
(C)

Problem Statement

- Input:



social graph
(A)



content graph
(C)

- Learn the **shared latent space** between A and C

Problem Statement

- Input:



social graph
(A)



content graph
(C)

- Learn the **shared latent space** between A and C
- Discover **ideology-popularity** latent dimensions

Problem Statement

- Input:



social graph
(A)

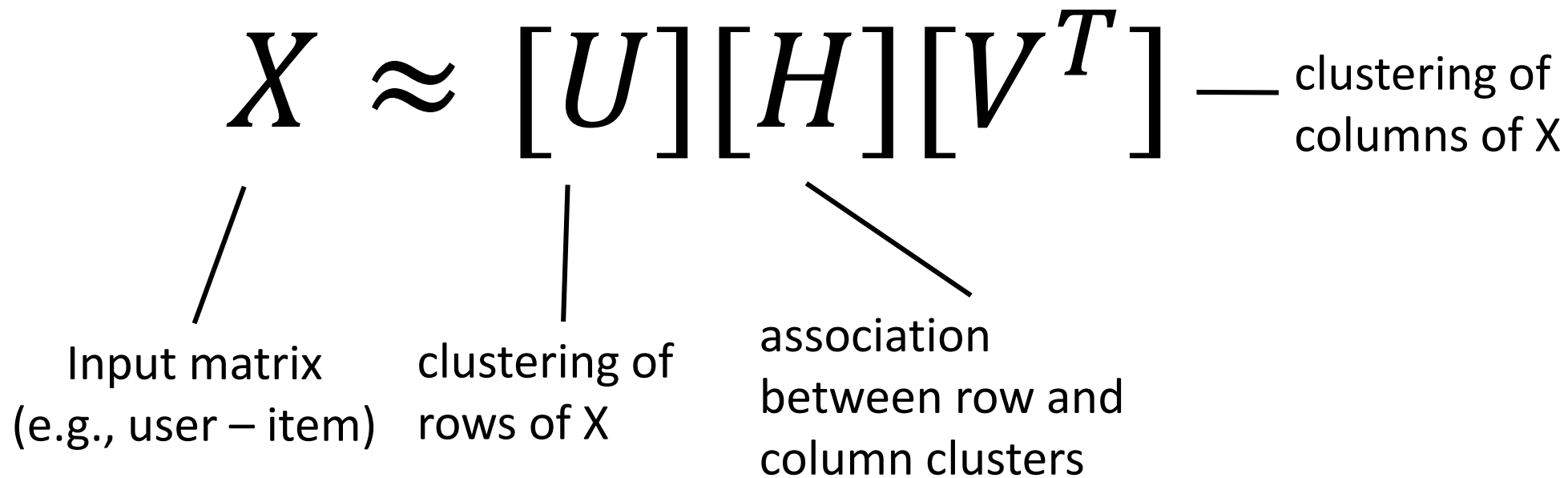


content graph
(C)

- Learn the **shared latent space** between A and C
- Discover **ideology-popularity** latent dimensions
- Estimate **ideology and popularity scores** for users and content

Proposed Methodology

Orthogonal Non-negative Matrix Factorization as **Co-Clustering Model** [Ding et al]



Combining Link and Content

$$J = \| \mathbf{A} - UH_u U^T \|_F^2 + \| \mathbf{C} - UH_s V^T \|_F^2$$

|
user-user
matrix

/
user-content
matrix

Shared Latent Space:

- A and C are related via users
- row datatype of matrix A is the same as rows of C

Combining Link and Content

$$J = \underbrace{\|A - UH_u U^T\|_F^2}_{\text{user-user matrix}} + \underbrace{\|C - UH_s V^T\|_F^2}_{\text{user-content matrix}} \quad \text{Joint Matrix factorization}$$

Shared Latent Space:

- A and C are related via users
- row datatype of matrix A is the same as rows of C

Learning Shared Latent Space

$$J = \|A - UH_u U^T\|_F^2 + \|C - UH_s V^T\|_F^2$$

Joint Matrix factorization

user-user matrix user-content matrix Shared latent factors (U and V)

Shared Latent Space:

- A and C are related via users
- row datatype of matrix A is the same as rows of C

Learning Hidden Manifolds in The Data

$$J = \|A - UH_u U^T\|_F^2 + \|C - UH_s V^T\|_F^2 + \lambda \text{Tr}(U^T L_u U) + \lambda \text{Tr}(V^T L_s V)$$

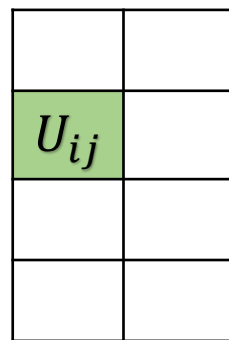
User Manifold and content Manifold
are tied together

- Users connected in social graph tend to be ideologically similar
- Ideologically similar users share similar content (and ideologically similar content is shared by similar users)

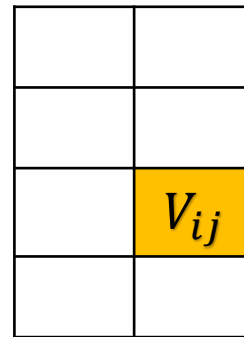
Latent factors have a probabilistic interpretation

$$J = \|A - UH_uU^T\|_F^2 + \|C - UH_sV^T\|_F^2$$

Latent factors U and V



U



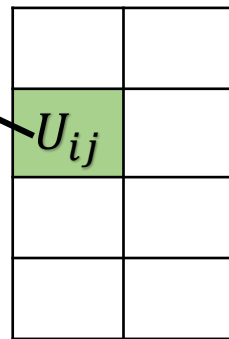
V

Latent factors have a probabilistic interpretation

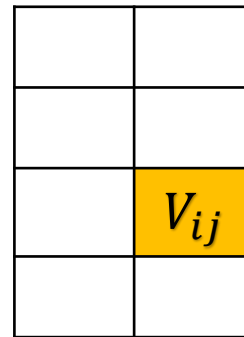
$$J = \|A - UH_uU^T\|_F^2 + \|C - UH_sV^T\|_F^2$$

Latent factors U and V

degree to which user i
belongs to ideology j



U



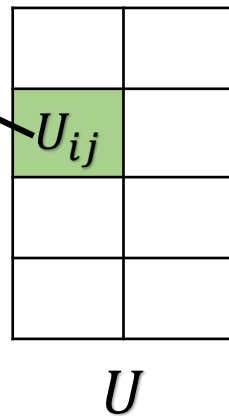
V

Latent factors have a probabilistic interpretation

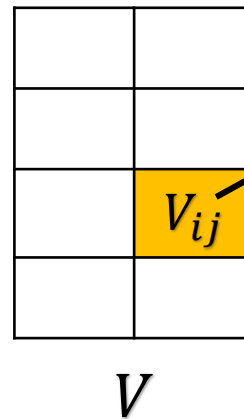
$$J = \|A - UH_uU^T\|_F^2 + \|C - UH_sV^T\|_F^2$$

Latent factors U and V

degree to which user i
belongs to ideology j



degree to which content i
belongs to ideology j



Latent factors have a probabilistic interpretation

Latent factors U and V

degree to which user i
belongs to ideology j

U_{ij}	

U

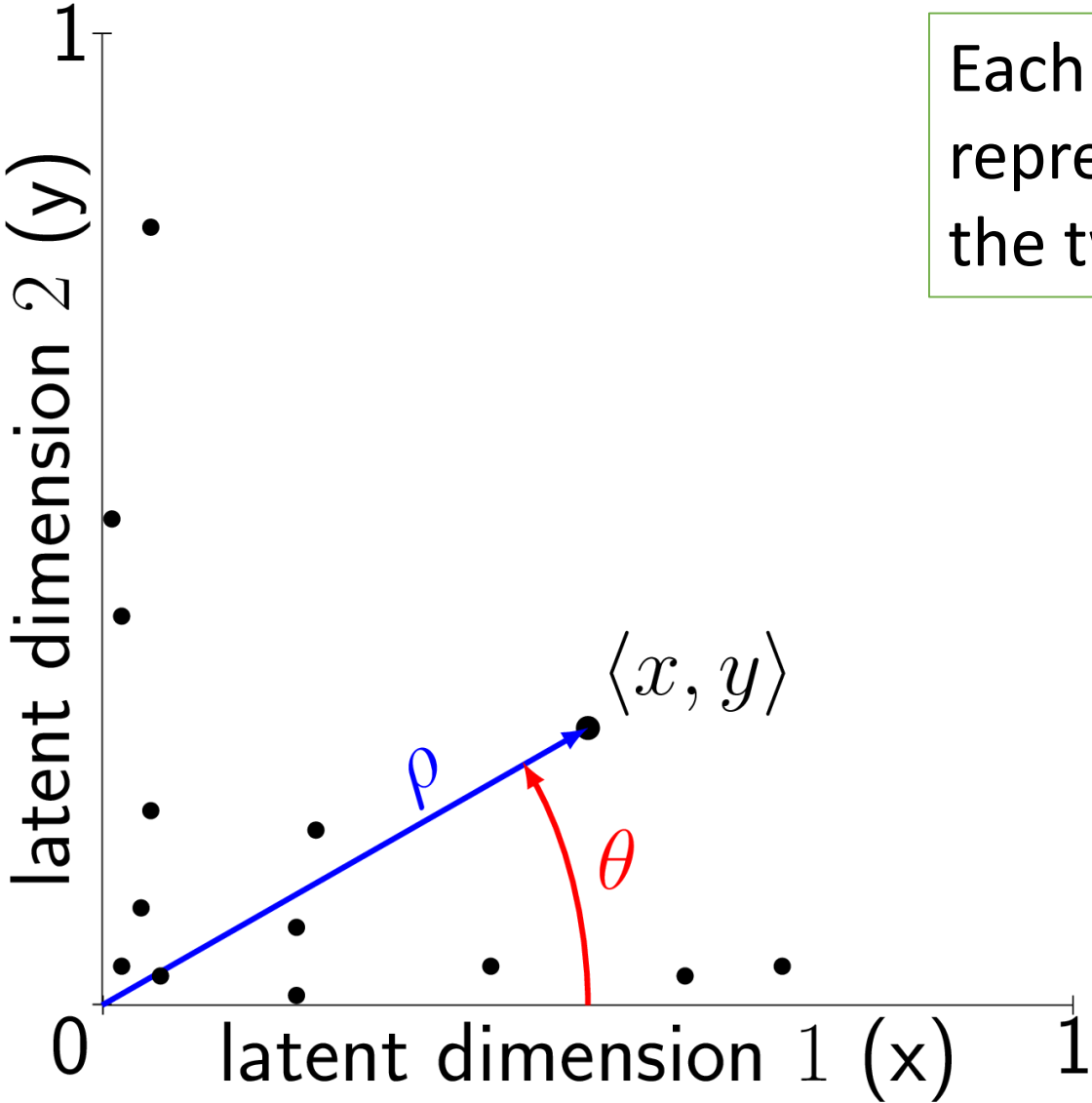
degree to which content i
belongs to ideology j

	V_{ij}

V

Each row of U and V can be represented as a two dimensional vector $\langle x, y \rangle$

Estimating Ideological Leaning



Each row of U and V can be represented as a vector $\langle x, y \rangle$ in the two dimensional space

Evaluation

Dataset

Twitter Streaming API (2011 - 2016)

- 7000 users
- 19 million tweets

Three controversial topics

- Gun control
- Abortion
- Obamacare

Ground Truth

- 500 news media channels
- Bayesian point estimate using large annotated data set [Barbera et al]

Baselines

Link (user-user)

- Graph partitioning (Retweet/Follows)
- NMF [Lee & Seung]

Content (user-content)

- ONMTF [Ding et al]
- DMCC [Gu et al]

Combined (link+ content)

- IFD / IFD-NGR [Proposed method]
- BIAS WATCH [Lu et al]
- KULSHRESTHA [Kulshrestha et al]

Baselines

Link
(user-user)

- Graph partitioning (Retweet/Follows)
- NMF [Lee & Seung]

No ideology scores

Content
(user-content)

- ONMTF [Ding et al]
- DMCC [Gu et al]

Combined
(link+ content)

- IFD / IFD-NGR [Proposed method]
- BIAS WATCH [Lu et al]
- KULSHRESTHA [Kulshrestha et al]

Baselines

Link
(user-user)

- Graph partitioning (Retweet/Follows)
- NMF [Lee & Seung]

No ideology scores

Matrix Factorization
based approaches

Content
(user-content)

- ONMTF [Ding et al]
- DMCC [Gu et al]

Combined
(link+ content)

- IFD / IFD-NGR [Proposed method]
- BIAS WATCH [Lu et al]
- KULSHRESTHA [Kulshrestha et al]

Baselines

Link
(user-user)

- Graph partitioning (Retweet/Follows)
- NMF [Lee & Seung]

No ideology scores

Matrix Factorization
based approaches

Content
(user-content)

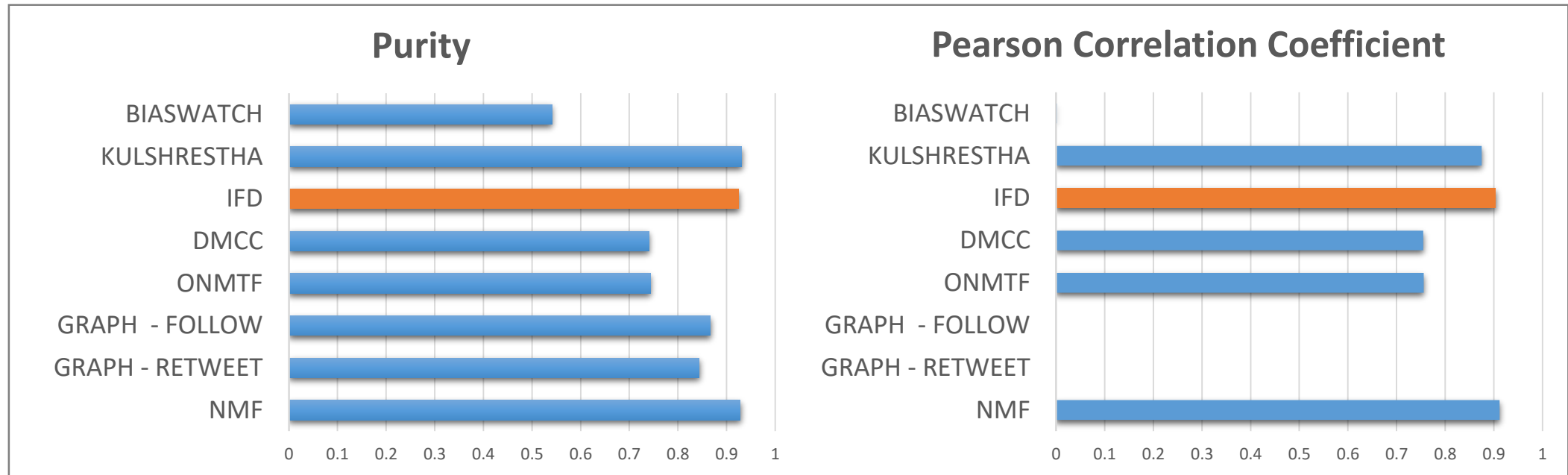
- ONMTF [Ding et al]
- DMCC [Gu et al]

Combined
(link+ content)

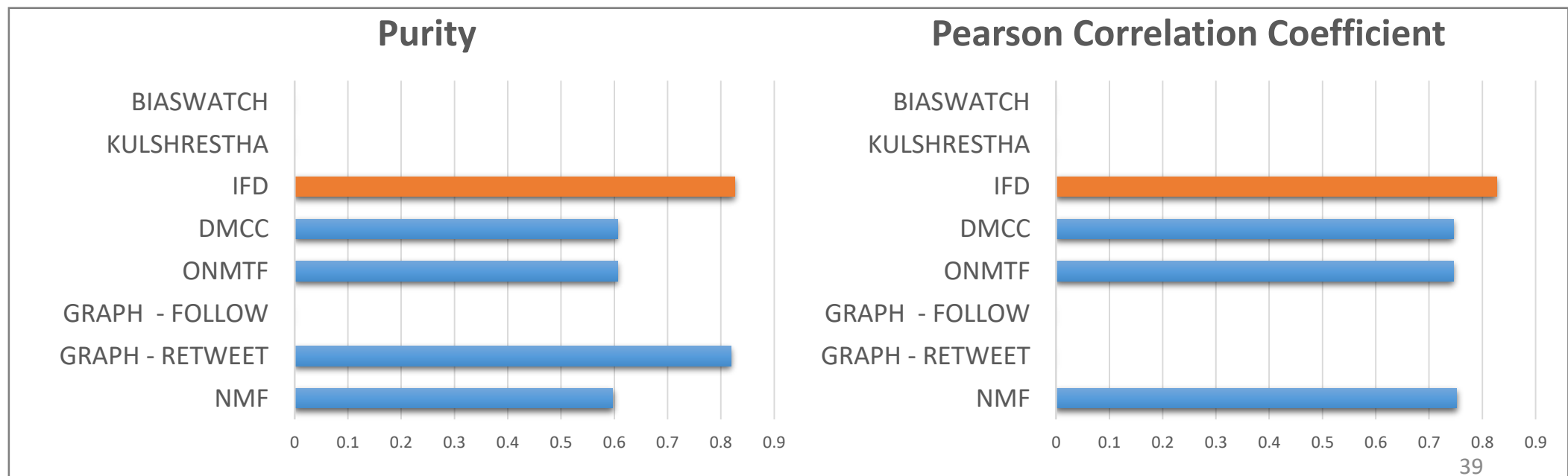
- IFD / IFD-NGR [Proposed method]
- BIAS WATCH [Lu et al]
- KULSHRESTHA [Kulshrestha et al]

No ideology scores
for media channels

Twitter Users (Ideology)

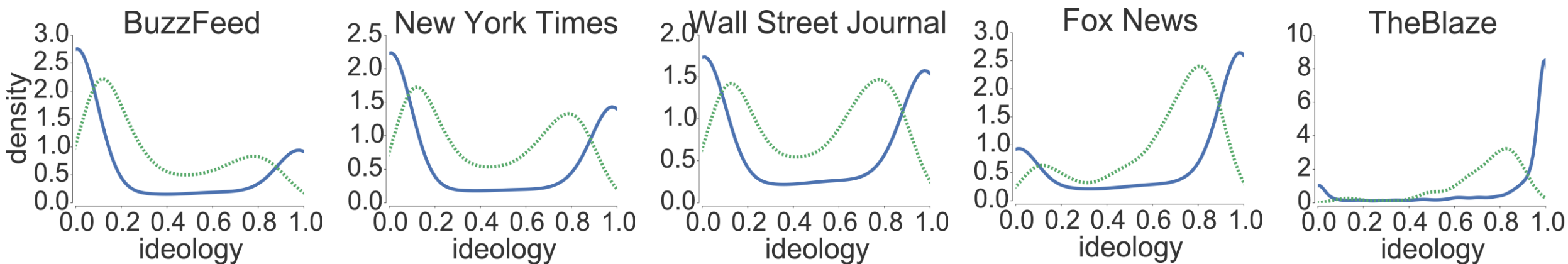


News Media (Ideology)



Estimated Ideology scores of high quality across the ideology spectrum (including center)

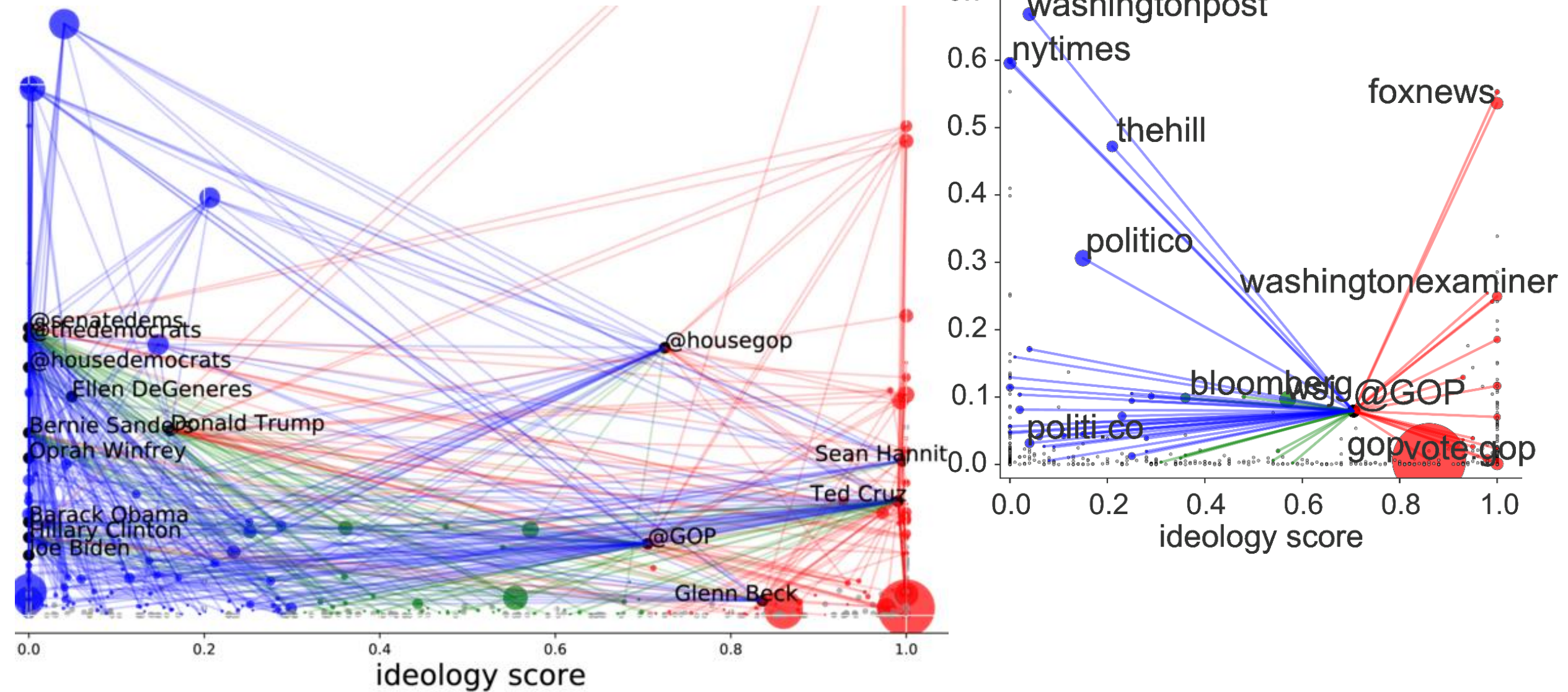
— computed score
⋯ ground truth



← Liberal

Conservative →

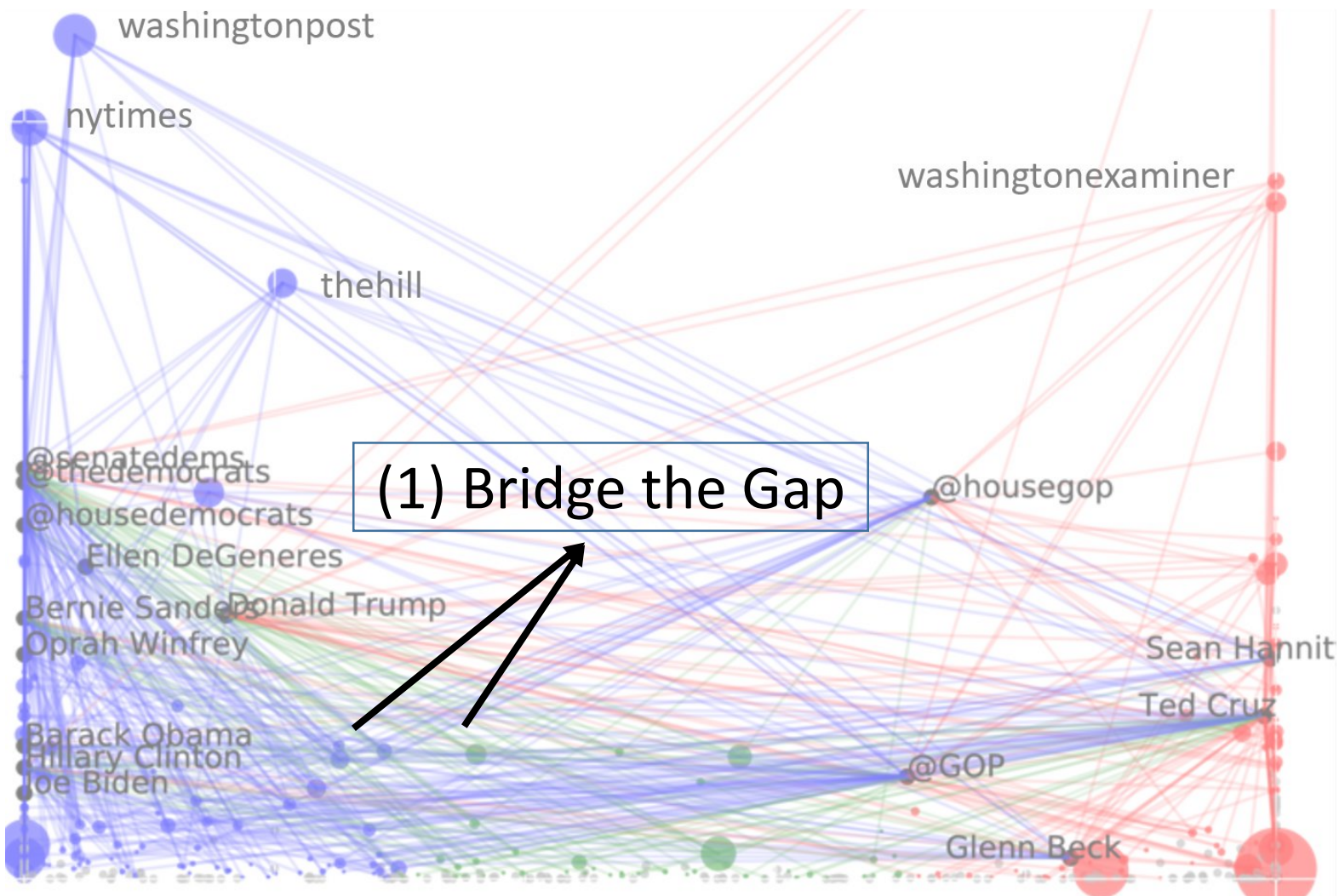
Ideological Latent Space



More Interactive Visualizations at <http://bit.ly/FilterBubbleDemo>

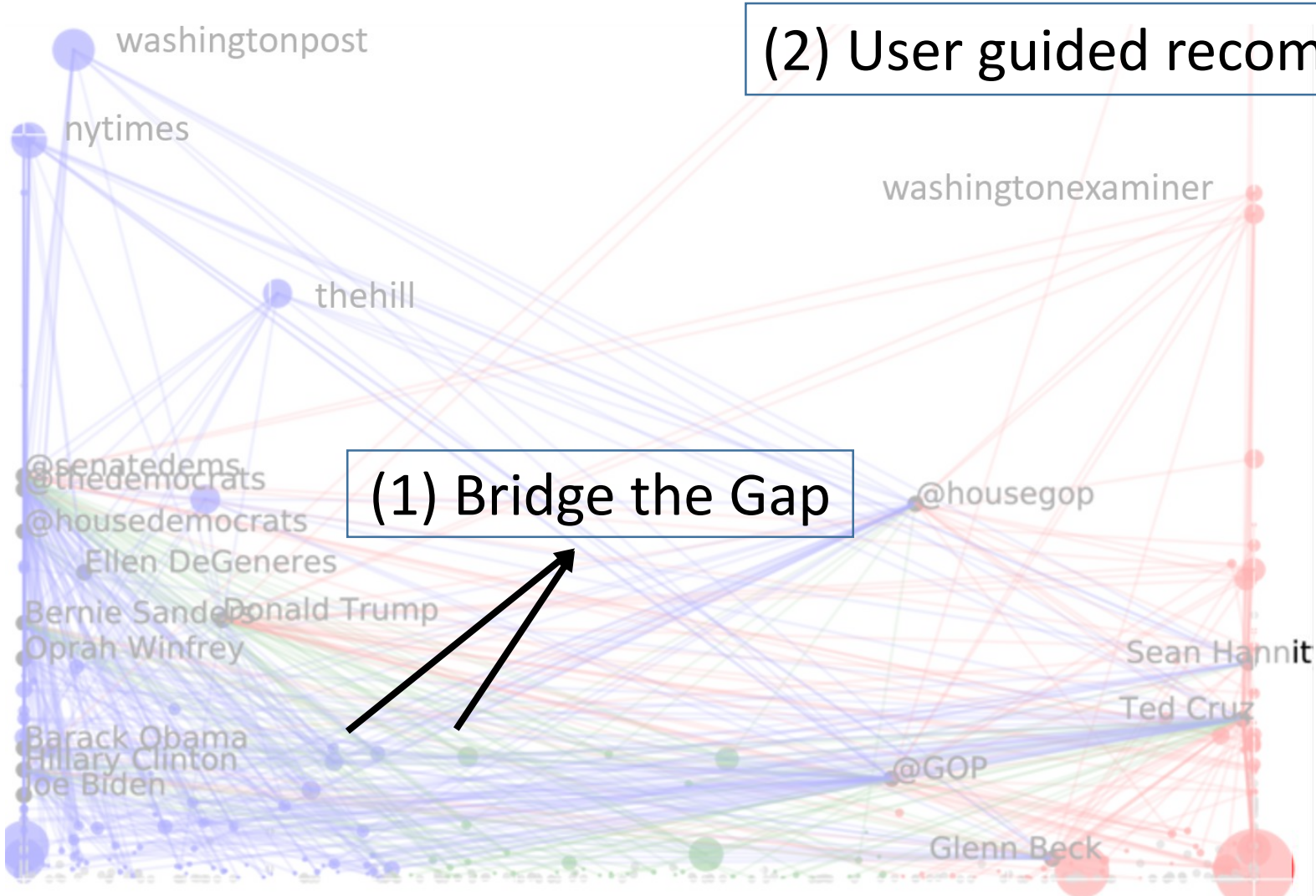
Motivation (Revisiting)

We now have access to ideological position of all the users and content



(2) User guided recommendations

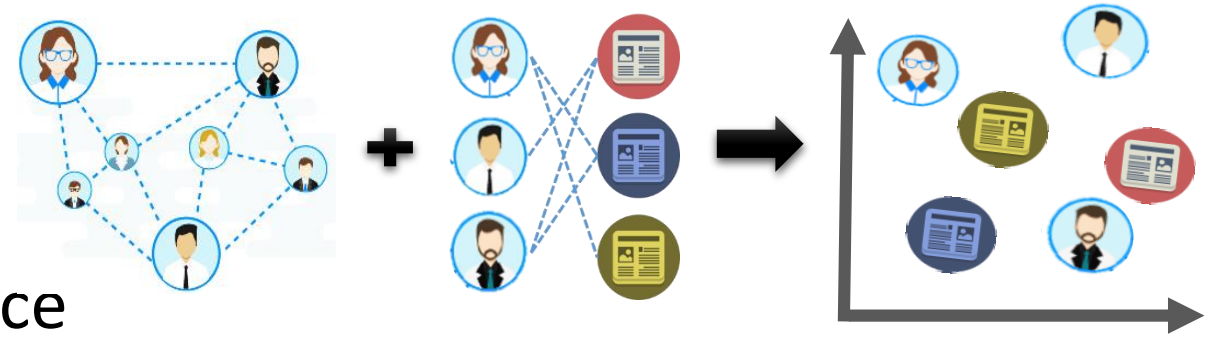
(1) Bridge the Gap



Summary

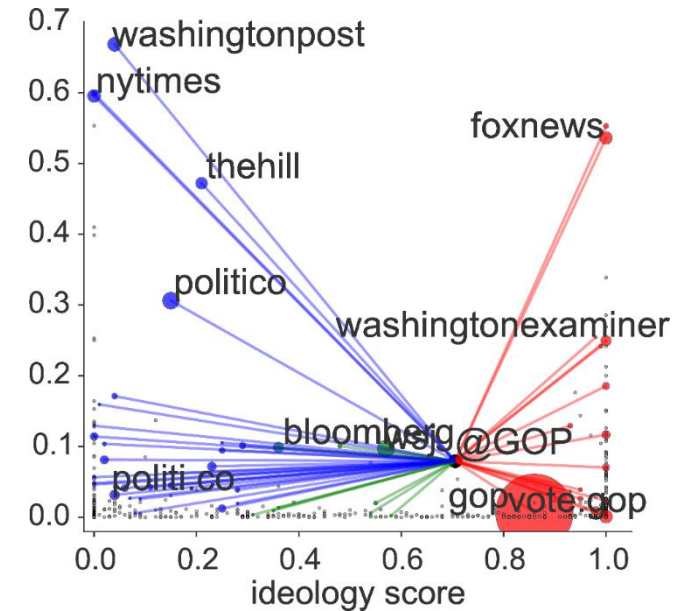
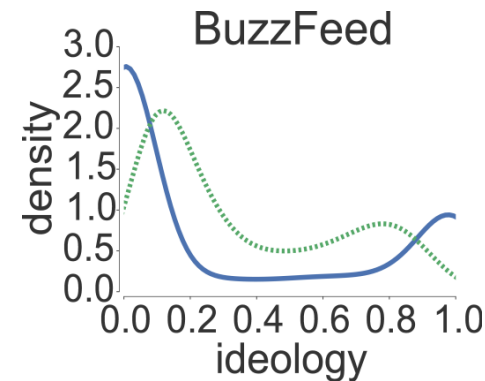
IFD framework

- Combines Link and Content Graph
- Jointly Matrix Factorization
- Shared Multidimensional Latent Space



Compared to Baselines IFD Estimates

- Ideology and Popularity Scores
- Twitter Users and Media Channels
- High Quality



THANK YOU

Learning Hidden Manifolds in The Data

Manifold assumption

If two data points x_i, x_j are close to each other in the input space then their projections in the new basis u_i, u_j are also close.

$$J = \|A - UH_u U^T\|_F^2 + \|C - UH_s V^T\|_F^2 + \lambda \text{Tr}(U^T L_u U) + \lambda \text{Tr}(V^T L_s V)$$

Graph regularization constraints [Cai et al]

where

- L_u, L_s are graph laplacians of the row and column affinity matrices of X
- $\text{Tr}(\cdot)$ is trace of the matrix

Optimization Problem

IFD (Ideology Factor Decomposition)

$$J = \|A - UH_uU^T\|_F^2 + \|C - UH_sV^T\|_F^2 + \lambda \text{Tr}(U^T L_u U) + \lambda \text{Tr}(V^T L_s V)$$

Constraints

- Bi-orthogonality ($U^T U = I ; V^T V = I$)
- Non-negativity ($U_+ ; H_{u_+} ; H_{s_+} ; V_+$)

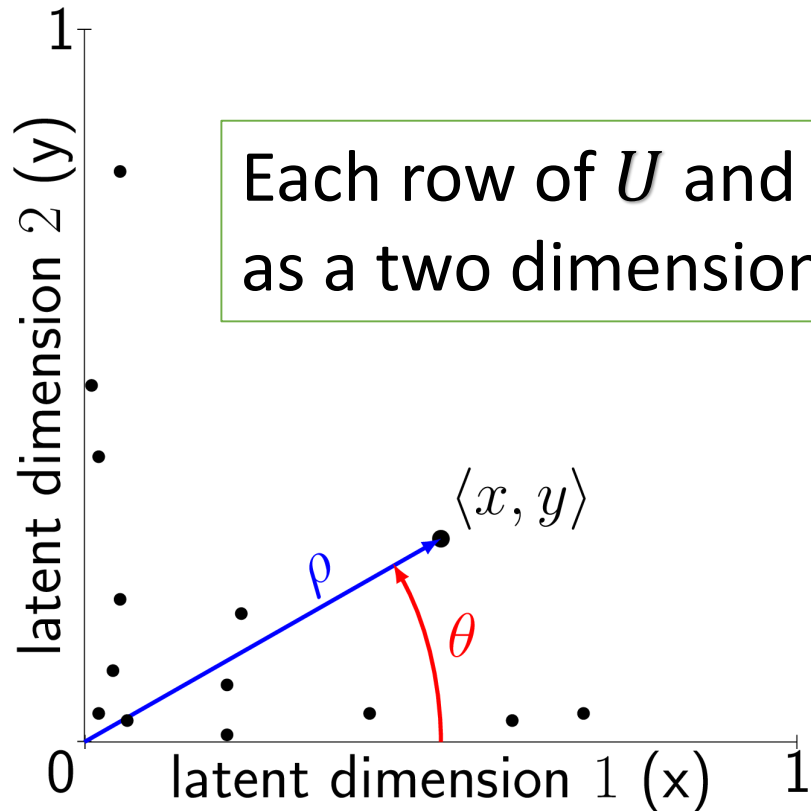
Solution

- Derive multiplicative update rules for U, V, H_u and H_s
- Iterative update algorithm
- Locally optimal solution

Estimating Ideological Leaning

The latent factors U and V have a probabilistic interpretation:

- U_{ij} : degree to which user i belongs to ideology j
- V_{ij} : degree to which content i belongs to ideology j



Each row of U and V can be represented as a two dimensional vector $\langle x, y \rangle$

Ideology:
$$i(x, y) = \frac{\theta}{\pi/2} = \frac{\arctan\left(\frac{y}{x}\right)}{\pi/2}$$

Popularity:
$$\rho(x, y) = \sqrt{x^2 + y^2}$$

Estimating Ideological Leaning

