

Of microscopes and telescopes in the digital humanities and computational social sciences

Talk at Future of Historical Network Research

Jana Diesner, PhD

Assistant Professor

The iSchool and Department of Computer Science

University of Illinois Urbana-Champaign

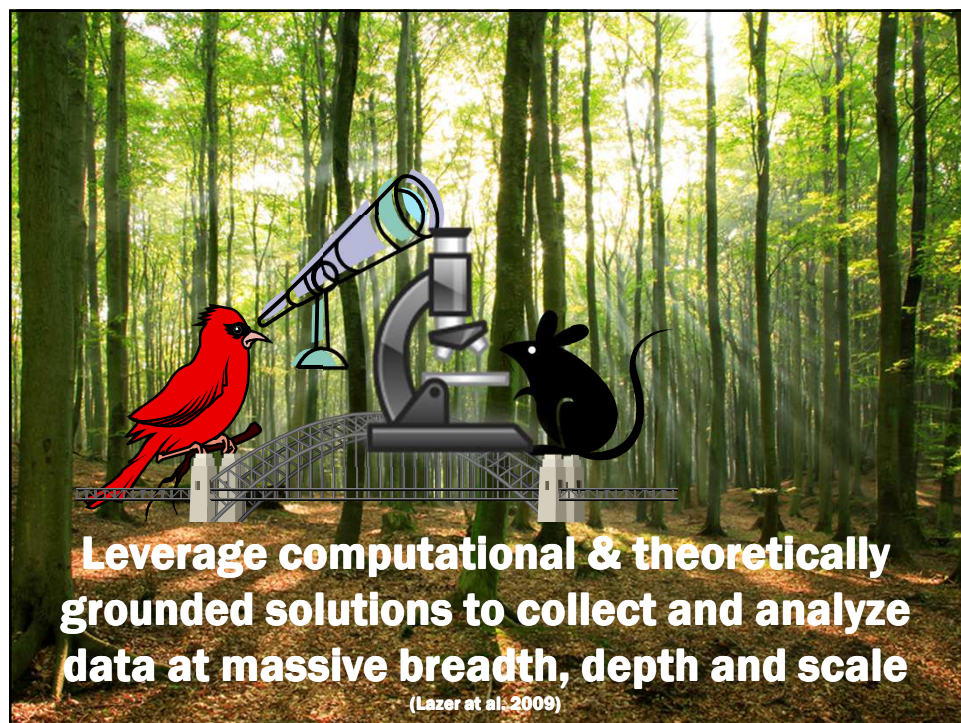
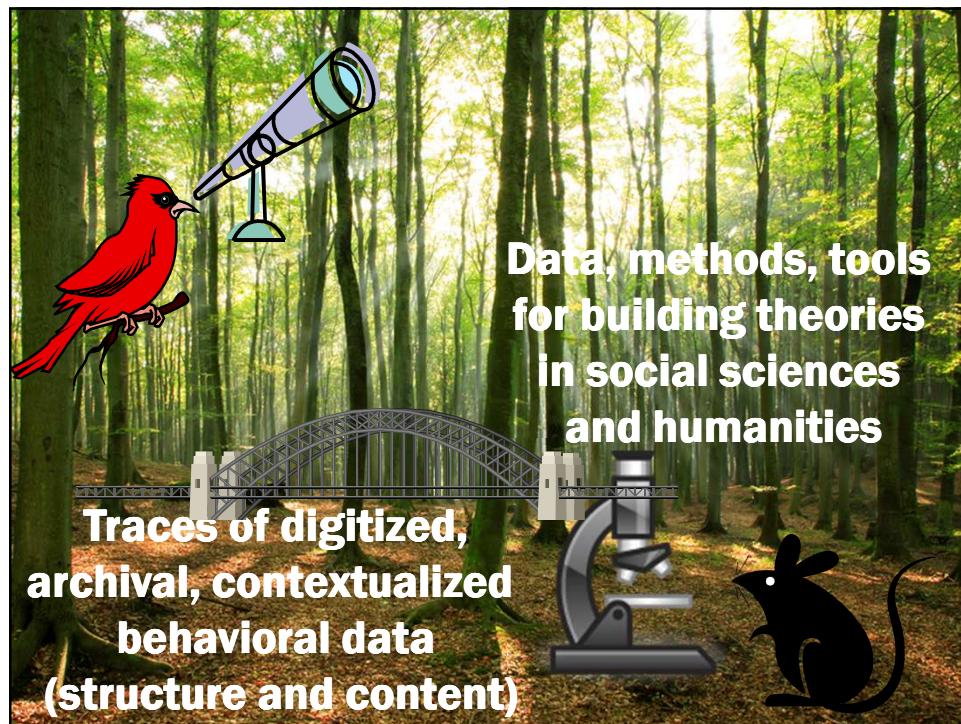


ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

GRADUATE SCHOOL OF LIBRARY AND
INFORMATION SCIENCE
The iSchool at Illinois

1



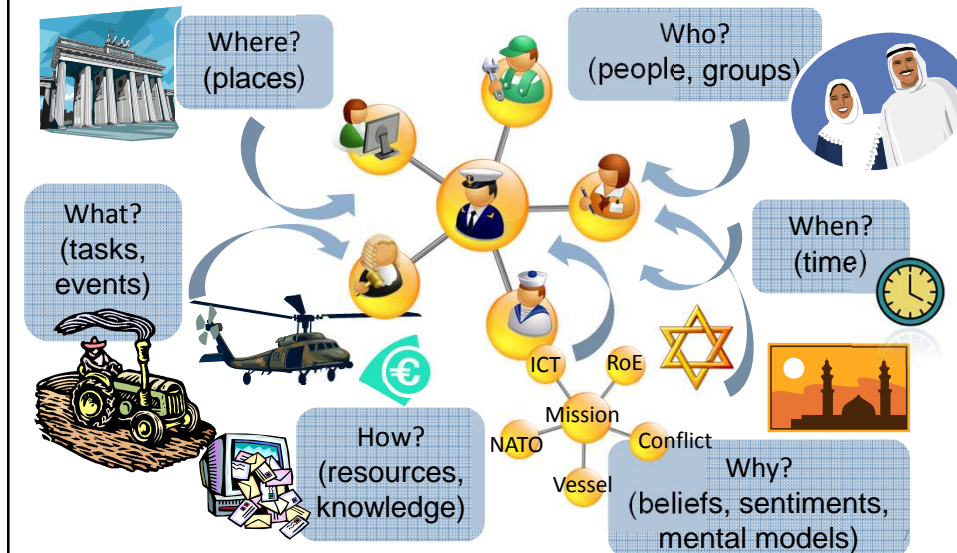


The diagram illustrates the four main components of Data Science, arranged in a circle with arrows indicating a clockwise flow between them:

- Social Science/Humanities** (Red quadrant):
 - Theories and Models
- Network Analysis** (Green quadrant):
 - Construct & analyze social- tech-semantic networks
- Machine Learning** (Purple quadrant):
 - Building Prediction Models
- Text Mining** (Blue quadrant):
 - Information & Relation Extraction

- Diesner J, Tamabyong L, Carley KM (2012) Mapping socio-cultural networks of Sudan from open-source, large-scale text data. *Computational and Mathematical Organization Theory (CMOT)*

From Words to Networks: What node classes to consider?



How to find and categorize nodes in text data? Basic recipe for probabilistic solution

- Get some labeled ground-truth data (BBN)
- Build a **classifier**/model (h) that **for every sequence of words (x)** and label per word (y) **predicts one category per word** ($y = h(x)$), incl. for new and unseen text data
- Exploit clues from text data (lexical, syntactic, statistical)
- Train and validate the model
- Get good accuracy (compare to intercoder reliability) (we made model available in end-user product AutoMap)
- Apply prediction model to text data ($\sim 80,000$ files)
- Link nodes (e.g. based on co-occurrence, proximity)
- Network data! Analysis!

Diesner J, Carley KM (2008) Conditional Random Fields for Entity Extraction and Ontological Text Coding. Journal of Computational and Mathematical Organization Theory (CMOT), 14(3), 248 – 262.

How to find and categorize nodes in text data?

- Model relationship among hidden states (y) as **Markov Random Field** (MRF) conditioned on observed data (x) (Lafferty et al. 2001)
- Compute **conditional distribution** of entity sequence y and observed sequence x as normalized product of potential functions M_i :

$$M_i(y_{i-1}, y_i | x) = \exp\left(\sum_{\alpha} \underbrace{\lambda_{\alpha}}_{\text{weight}} \underbrace{f_{\alpha}(y_{i-1}, y_i, x)}_{\text{feature}} + \sum_{\beta} \underbrace{\mu_{\beta}}_{\text{weight}} \underbrace{g_{\beta}(y_i, x)}_{\text{feature}}\right)$$

$$P_{\theta}(y | x) = \frac{\prod_{i=1}^{n+1} M_i(y_{i-1}, y_i | x)}{\prod_{i=1}^{n+1} M_i(x)_{start, stop}}$$

- **Edge and transition features** plus **node and emission features**
- f, g : boolean feature vectors with **learned** weights
- Tool: CRF project page, training data: BBN

9

How good is it?

Class model	Class	Specificity	Subtype	Example	States, Edges	Time (300)
1	x			agent	11, 121	17.5 hours
2	x	x		agent, specific	16, 256	1.25 days
3	x		x	agent, political	32, 1024	3.1 days
4	x	x	x	agent, spec., pol.	45, 2025	5 days

Variable	Accuracy	Training Time
Baseline	large	small
Syntax features (POS)	small	small
Lexical features (dict, hard match)	large	small
Iteration rate	large	large
Complexity of class model	small	large

	Boundary Model	Class model 1 (cat.)	Class model 2 (cat. + spec.)	Class model 3 (cat.+ sub.)	Class model 4 (cat+spec+sub)
Precision	93.2%	91.4%	91.9%	90.4%	90.8%
Recall	92.5%	89.7%	90.0%	88.6%	88.9%
F	92.9%	90.6%	90.9%	89.5%	89.8%
Bound. & Class (rules)		Model 1	Model 2	Model 3	Model 4
Precision	n.a.	89.7%	90.0%	88.6%	88.9%
Recall		87.7%	87.7%	86.4%	86.5%
F		88.7%	88.8%	87.5%	87.7%

Network Data! Analysis!

Activity:

Control:

Close to power:

Degree Centrality	03	04	05	06	07	08	09	10
Omar al-Bashir	3	3	2	1	1	1	1	1
Ali Osman Taha	1	2	3	4	3	3	3	3
John Garang	2	1	1	3	3	4	6	8
Salva Kiir Mayardit	8	10	4	2	2	2	2	2
Hosni Mubarak	4	7	5	6	9	8	4	6
Sadiq al-Mahdi	6	5	10	9	5	7	8	4
Hassan al-Turabi	5	6	7	10	5	8	9	5
Abdul Wahid al Nur	10	9	9	8	7	4	5	7
Yoweri Museveni	7	8	7	6	10	7	8	5
Kofi Annan	9	4	6	5	8	10	7	7
Deng Alor	10	10	10	10	6	9	8	9

Triads	03	04	05	06	07	08	09	10
Omar al-Bashir	1	1	1	1	1	1	1	1
Ali Osman Taha	2	3	3	4	4	3	2	2
John Garang	3	2	2	2	2	6	7	7
Salva Kiir Mayardit	7	10	4	3	3	2	3	3
Hosni Mubarak	7	4	5	6	6	8	4	5
Sadiq al-Mahdi	4	7	7	7	6	7	7	3
Abdul Wahid al Nur	10	9	9	7	4	5	5	7
Kofi Annan	7	5	5	5	5	7	7	7
Yoweri Museveni	6	6	8	9	9	10	6	5
Hassan al-Turabi	5	8	9	9	8	9	7	7
Deng Alor	10	10	9	9	10	4	7	7

- President **North**: Known performer
- President **South**: Now established
- Legacy of **religious leaders**
- Presence of **neighboring presidents**



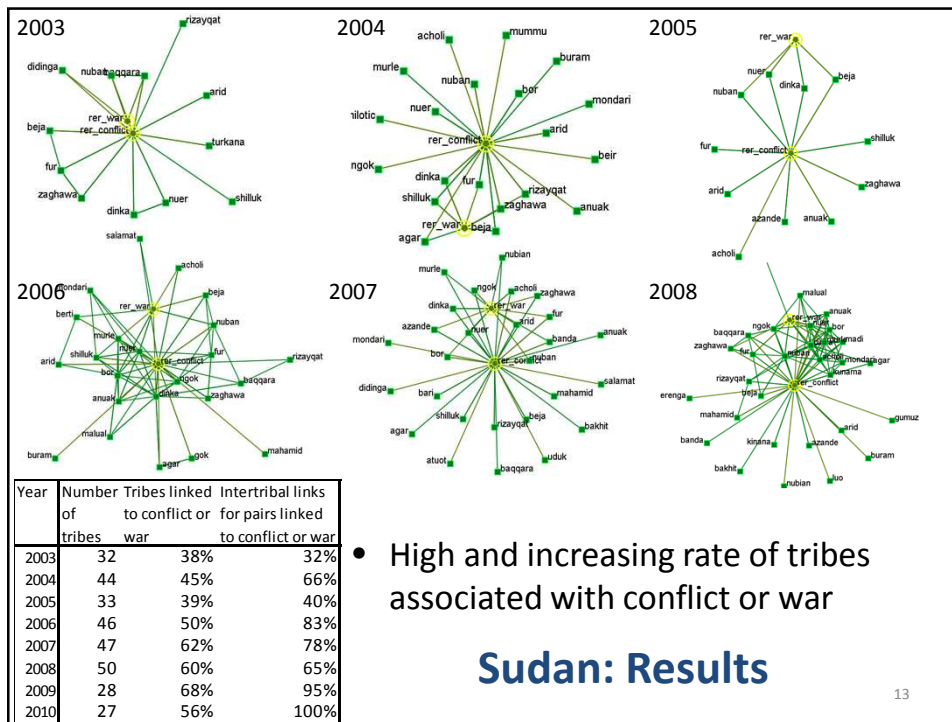
2003	2004	2005	2007	2010
Darfur conflict	Continuous civil war (since 1993)	Comprehensive Peace Agreement Garang 1 st VP, followed by Kiir Autonomous South Sudan	SPLA withdraws from government	Votum in South Sudan about Separation

Sudan: Results (Groups)

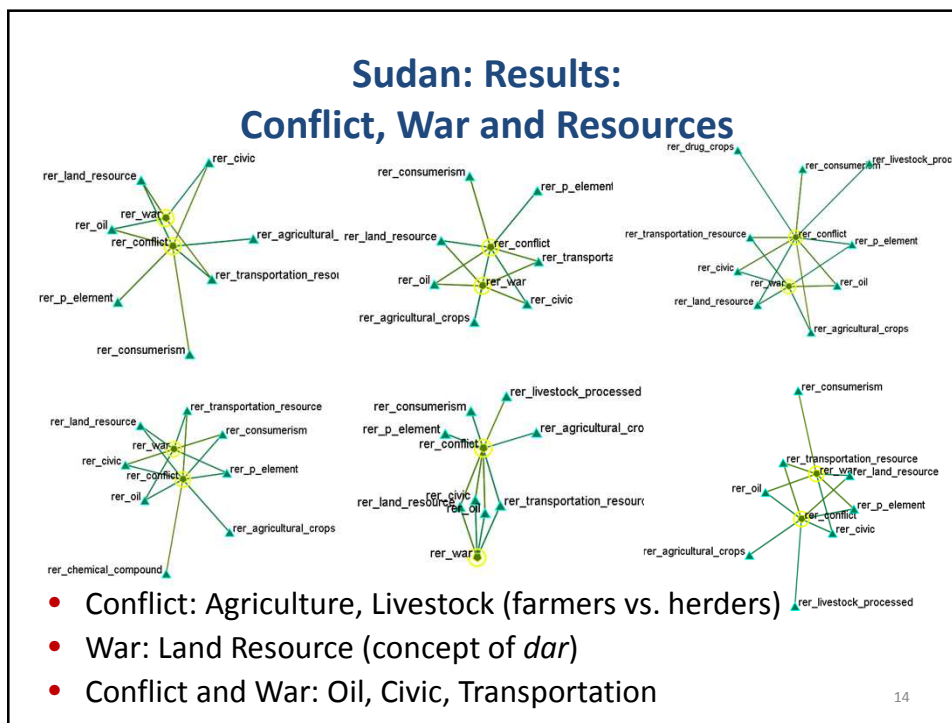
Degree Centrality	03	04	05	06	07	08	09	10
United Nations	4	2	1	1	1	1	1	5
Rebel Groups	1	1	2	3	4	3	2	3
Military	2	3	3	2	2	2	4	2
SPLA	6	5	4	3	4	3	1	1
Security Council	5	5	4	5	5	5	5	6
Sudan government	3	4	6	6	8	8	9	7
Nat. Congress Party	6	9	9	8	6	7	10	4
African Union	8	7	8	7	7	8	7	8
Inter. Criminal Court	8	7	7	7	9	6	6	9
Dinka	9	10	9	10	9	10	8	8
Churches	7	8	7	10	10	10	10	10

Triads	03	04	05	06	07	08	09	10
Military	1	1	1	1	2	1	6	1
United Nations	4	3	2	2	1	4	1	2
Rebel Groups	2	2	4	4	4	2	4	5
SPLA	5	3	3	3	3	2	4	5
Sudan government	3	4	5	7	5	7	4	6
Nat. Congress Party	5	9	10	8	6	6	10	3
African Union	8	6	6	6	7	10	7	3
Security Council	7	7	7	5	8	9	8	8
Inter. Criminal Court	8	8	9	10	5	3	7	7
Churches	6	8	9	10	9	8	10	10
Dinka	9	10	9	10	9	10	8	8

- Strong presence of armed forces
- Strong influence of external groups
- Within top 10 Sudanese groups:
 - Dinka, Nuer (ethnic groups/ tribes)



13



14

Example: Assessment of Impact of Social Justice Documentaries

- **Storytelling:** documentaries create memories and imagination, stimulate sharing (Rose 2012)
- **Impact:** documentaries aim to cause **change** in people's knowledge and/or behavior (Barrett & Leddy 2008)
- **Status quo:** funders, practitioners, scientists agree on **strong need for rigorous impact assessment**, but amount and depth of **prior reports and work limited**
- **Our questions:** How can we know if a documentary has impact? At what point in the life cycle of a film can we answer this question?

15

Our Solution

- **Solution:** develop, implement, apply, evaluate computational solution for gaining **broad** (quantitative) and **deep** (qualitative) understanding of impact (change)
- **Assumption:** films produced, screened, watched as part of large and **continuously changing ecosystems of stakeholders and themes**
- **ConText:** track, map, analyze structure, functioning and dynamics of **web of stakeholders and themes** associated with main issue addressed in a movie at scale
 - **Look no hands – empirical data analytics: network analysis** (map and analyze social networks) paired with **natural language processing/ text mining** (semantic networks, salient concepts and topics, valence of content)
 - **This is no fishing expedition – we have a theory:** developed CoMTI (content, medium, target, impact), a comprehensive measurement framework

Our Solution: Research and Development Process

Theory	<ul style="list-style-type: none"> Comprehensive lit review -> develop framework of relevant dimensions/indicators of media impact
Operationalization	<ul style="list-style-type: none"> Translate indicators into metrics and indices
Methods & Algorithms	<ul style="list-style-type: none"> Map indices to methods and algorithms suitable for analyzing large-scale, empirical data
Data Collection	<ul style="list-style-type: none"> Empirical: news coverage, social media, focus groups,
Technology	<ul style="list-style-type: none"> ConText: Comprehensive tech review -> build tool from open source code and from scratch
Analysis & Interpretation	<ul style="list-style-type: none"> Apply technology to various data sources on various movies
Evaluation	<ul style="list-style-type: none"> Assess accuracy and performance of methodology and technology

Our Methodology in a Nutshell

Understand problem space: (Where) is impact possible?

- Map public discourse and key players (different types of influence and power) related to main theme(s) of film prior to release (baseline model)
- Collect, analyze, fuse data from news, social media, focus groups
- Map opportunity space for linking up to people and themes, and strategic resource allocation (mobilize (social) capital)
- Identify unpromising topics early (no or uncontroversial discourse)

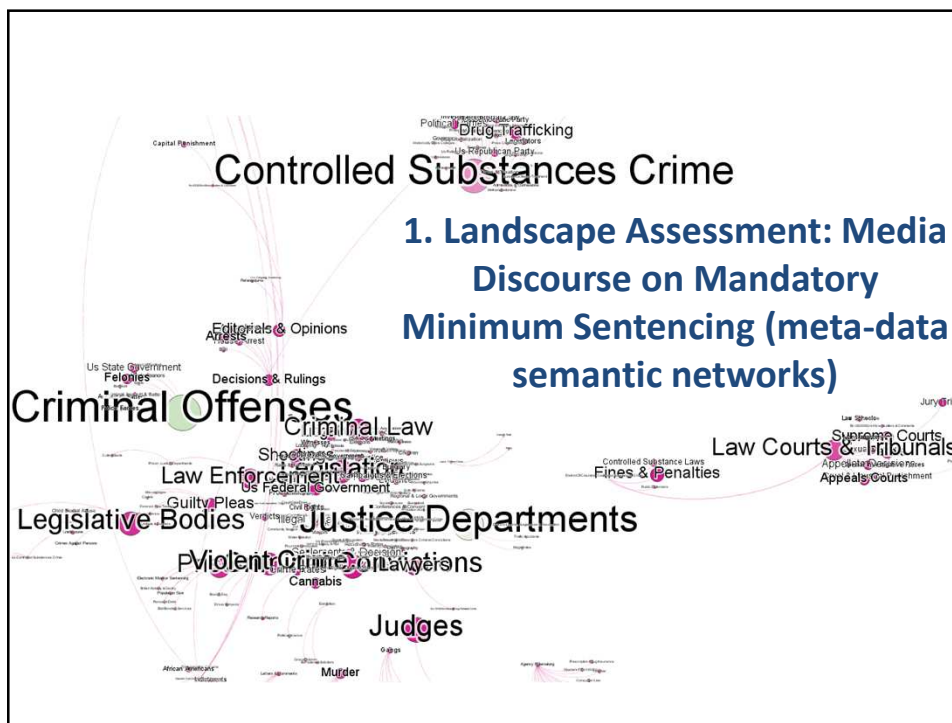
Understand film's message: Impact w.r.t. what?

- Analyze transcripts (ground truth model)

Understand film's impact: Has the needle moved?

- Reassess map of discourse and key players: New links, new themes that connect movie to baseline? Baseline connected or closer to ground truth?

ConText in Action: The House I Live In (Eugene Jarecki 2012)

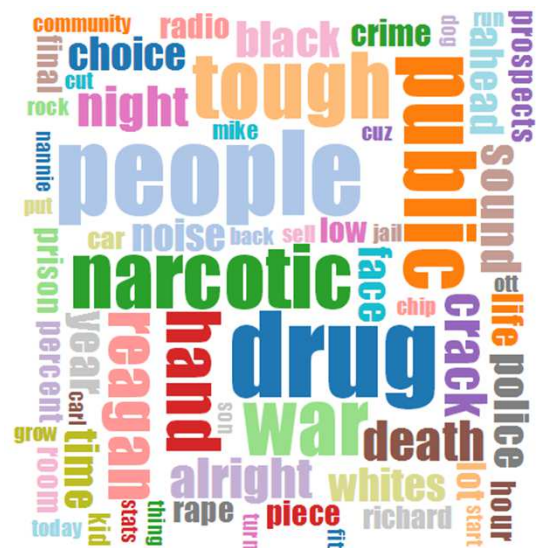


1. Landscape Assessment: Media Discourse on Mandatory Minimum Sentencing (content topic modeling)



framed as
social issue
(people in
center of
public
discourse)

2. What's in the Movie? (topic modeling of transcript)



large common denominator:
people (social issue),
difference:
media: prisons and violence,
movie: drugs and related politics

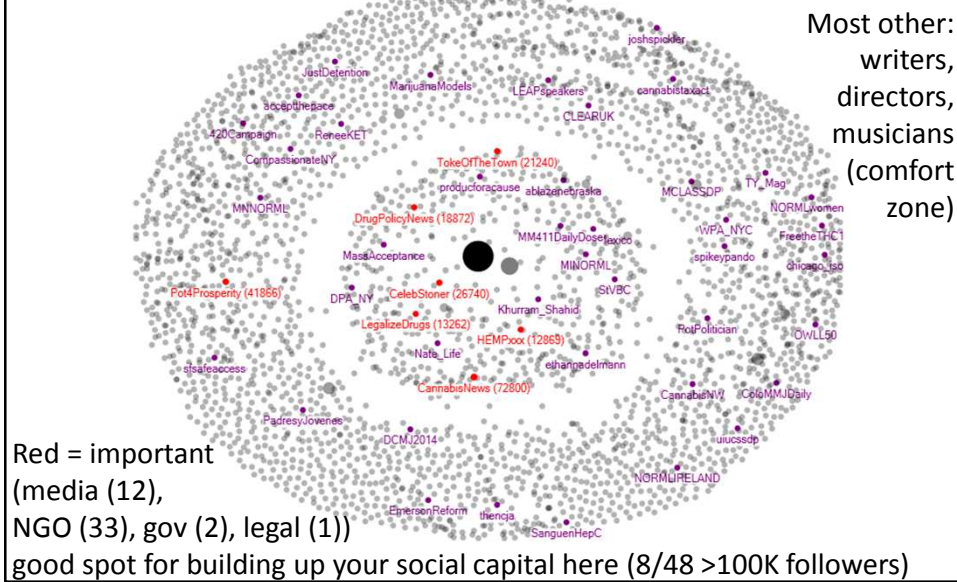
[illegible]

Followers (@DrugWarMovie)
(3314, visible if >200 followers)

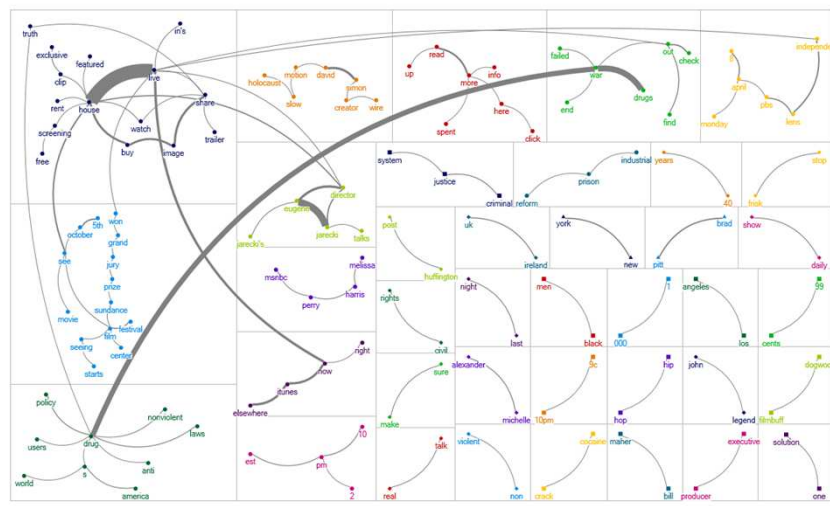
People who the movie follows (2245)

Intersection (510)

3. Has the Needle Moved? Social Media (Influential Followers on Twitter)



3. Has the Needle Moved? Social Media Campaign (Facebook Fanpage Posts, 428)



- Impact Assessment: FORD Foundation, grant 0125-6162.
- Sudan: National Science Foundation (NSF) IGERT 9972762, the Army Research Institute (ARI) W91WAW07C0063, the Army Research Laboratory (ARL/CTA) DAAD19-01- 2-0009, the Air Force Office of Scientific Research (AFOSR) MURI FA9550-05-1-0388, the Office of Naval Research (ONR) MURI N00014-08-11186, and a Siebel Scholarship.

References

- Lazer, D., Pentland, A. S., Adamic, L., Aral, S., Barabasi, A. L., Brewer, D., . . . Gutmann, M. (2009). Life in the network: the coming age of computational social science. *Science*, 323(5915), 721-723.
- On movie impact assessment:
 - Diesner J, Aleyasen A, Kim J, Mishra S, Soltani S (2013) Using Socio-Semantic Network Analysis for Assessing the Impact of Documentaries. WIN (Workshop on Information in Networks) 2013, New York, NY
- On Sudan project:
 - Diesner J, Carley KM, Tamabyong L (2012) Mapping socio-cultural networks of Sudan from open-source, large-scale text data. *Computational and Mathematical Organization Theory (CMOT)* 18(3), 328-339.
 - Diesner J (2012) Uncovering and Managing the Impact of Methodological Choices for the Computational Construction of Socio-Technical Networks from Texts. Technical Report CMU-ISR-12-101.

29

Thank you!

Q&A

- For questions, comments, feedback, follow-up:
Jana Diesner
jdiesner@illinois.edu
Phone: (412) 519 7576
Web: <http://people.lis.illinois.edu/~jdiesner>

30