

Meme-tracking and the Dynamics of the News Cycle

Jure Leskovec, Lars Backstrom and Jon Kleinberg
Computer science department
Cornell University

<http://memetracker.org>

CULTURAL INFORMATION

PRACTICE OR **IDEA** OR CONCEPT

THEORIES PRACTICES HABITS SONGS

NATURAL SELECTION

EXAMPLES MIGHT INCLUDE THOUGHTS IDEAS

CHARLES DARWIN'S IDEAS

SELF-PROPAGATING

SURVIVAL AND COMPETITION INFLUENCE THEM

MEME

Information and Media

- Intersection of news media, technology, and the political process
- From its early stages, a tension between **global effects** from the mass media and **local effects** carried by social structure

How does information transmitted by the media interact with the personal influence arising from social networks?

Fragmentation and Acceleration

- Internet, blogging, and social media:
 - Social media means the dichotomy between global and local influence is evaporating
 - Speed of media reporting and discussion has intensified: very rapid progression of stories, with no pauses
- The "24-hour news cycle":
 - Difficult to define, but associated with technological acceleration and a challenge to healthy civic discourse [Kovach-Rosenstiel '99]

Defining News Cycle (1)

- **Sept. 11, 2008 (New York Times):** *“Mr. McCain's increasingly aggressive campaign has sought to put Mr. Obama on the **defensive in each news cycle**, using any development at hand, like Mr. Obama's colloquial comment this week about ‘putting lipstick on a pig.’”*
- **Oct. 10, 2008 (New York Times):** *“Mr. McCain's traveling road show has veered from message to message **Each news cycle seems to bring another tactic** as the campaign appears to be trying anything and everything to see what might work.”*

Defining News Cycle (2)

- **Question:** Is the “news cycle” simply a metaphorical construct describing our perception of the news, or is it visible in data?
- And if it's visible, can we measure some of its basic properties?

Units of analysis?

- What basic “units” make up the news cycle?
 - Cascading hyper-links to articles: too fine-grained [Adar et al 04, Gruhl et al 04, Kumar et al 03, Leskovec et al]
 - Topics as probabilistic term mixtures: too coarse-grained [Blei-Lafferty 06, Wang-McCallum 06, Wang et al 07]
 - Named entities: too coarse-grained
Obama, McCain, Microsoft, Paris, Apple
 - Common sequence of words: too noisy
“I love you”, “web 2.0”, “Oh my God”, “Made in China”

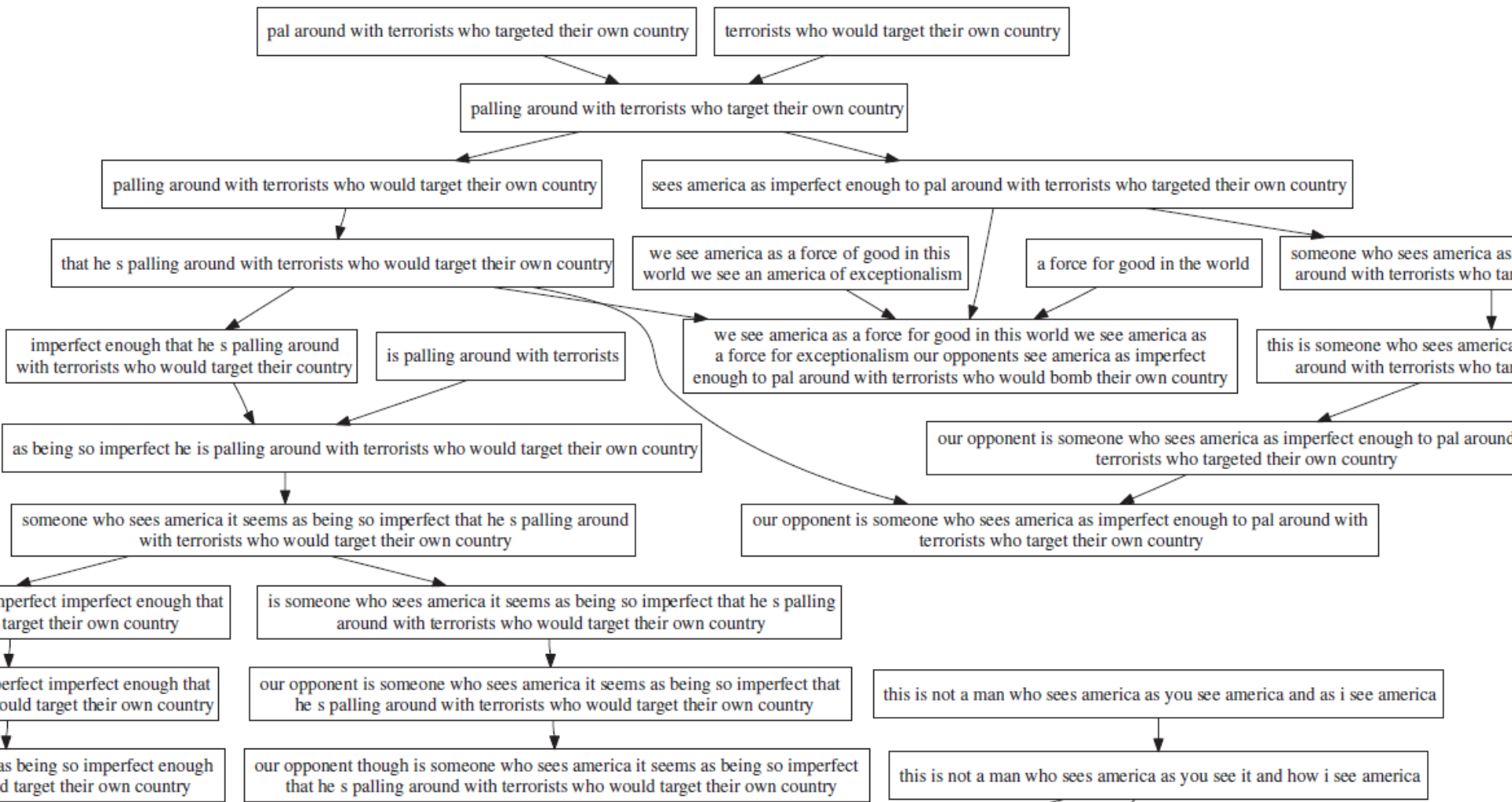
How to detect memes?

- Need **units** that:
 - correspond to aggregates of articles,
 - vary over the order of days,
 - and can be handled at terabyte scale
- **Plan:** identify text fragments, phrases, memes that travel relatively unchanged through many articles.
- **Idea: quoted phrases: “.*”**
 - are integral parts of journalistic practices
 - tend to follow iterations of a story as it evolves
 - are attributed to individuals and have time and location

Online media

- Data from Spinn3r on the 3 months leading up to the 2008 U.S. Presidential Election:
 - 1 million news articles and blog posts per day
 - Essentially a complete online media coverage:
 - 20,000 sites that are part of Google News
 - 1.6 million blogs
 - From August 1 to October 31 2008
 - 90 million documents from 1.65 million sites, 390GB
 - We extract 112 million quotes (phrases)

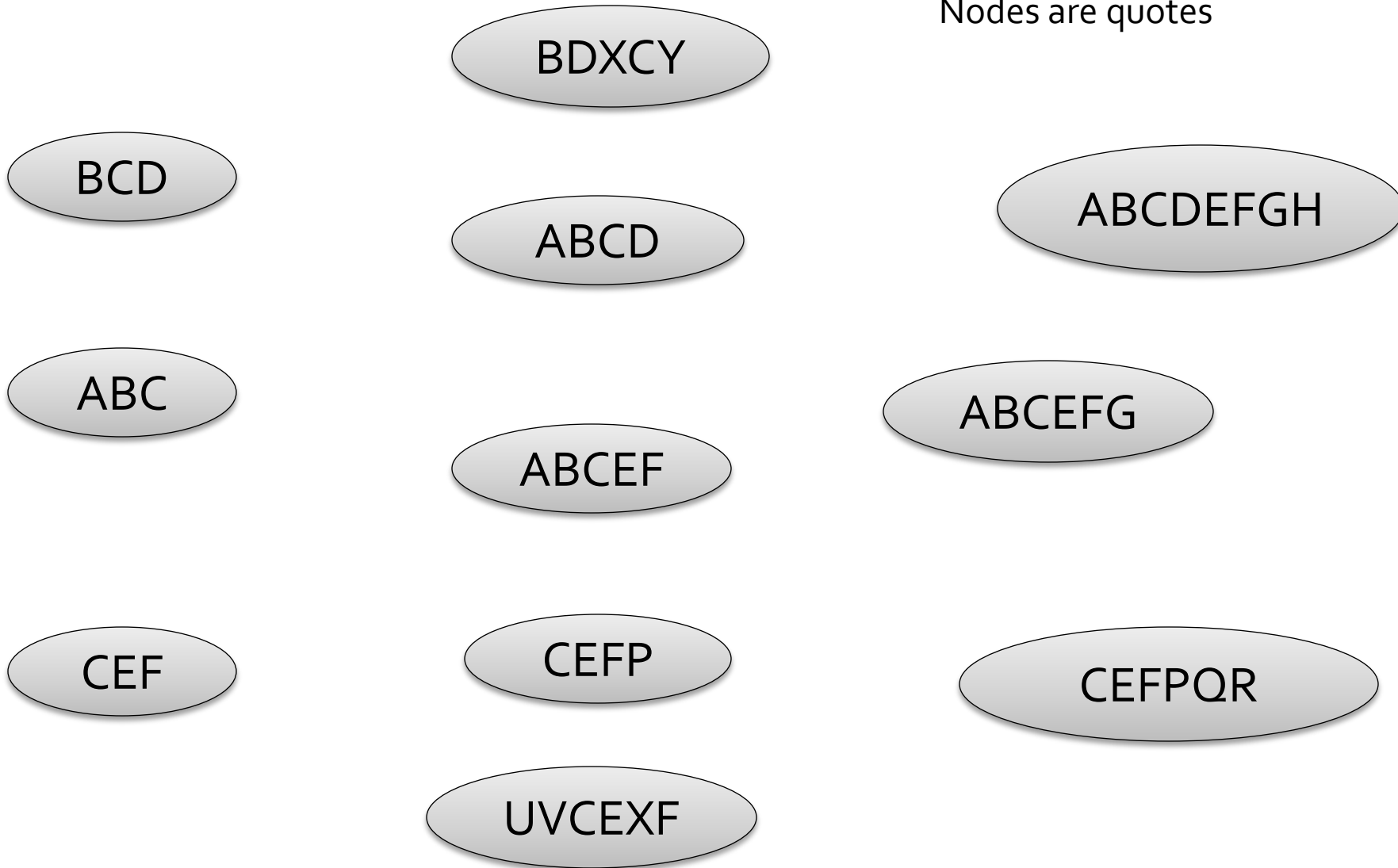
Challenge: Quotes Mutate... A lot!



Our opponent is someone who sees America, it seems, as being so imperfect, imperfect enough that he's palling around with terrorists who would target their own country.

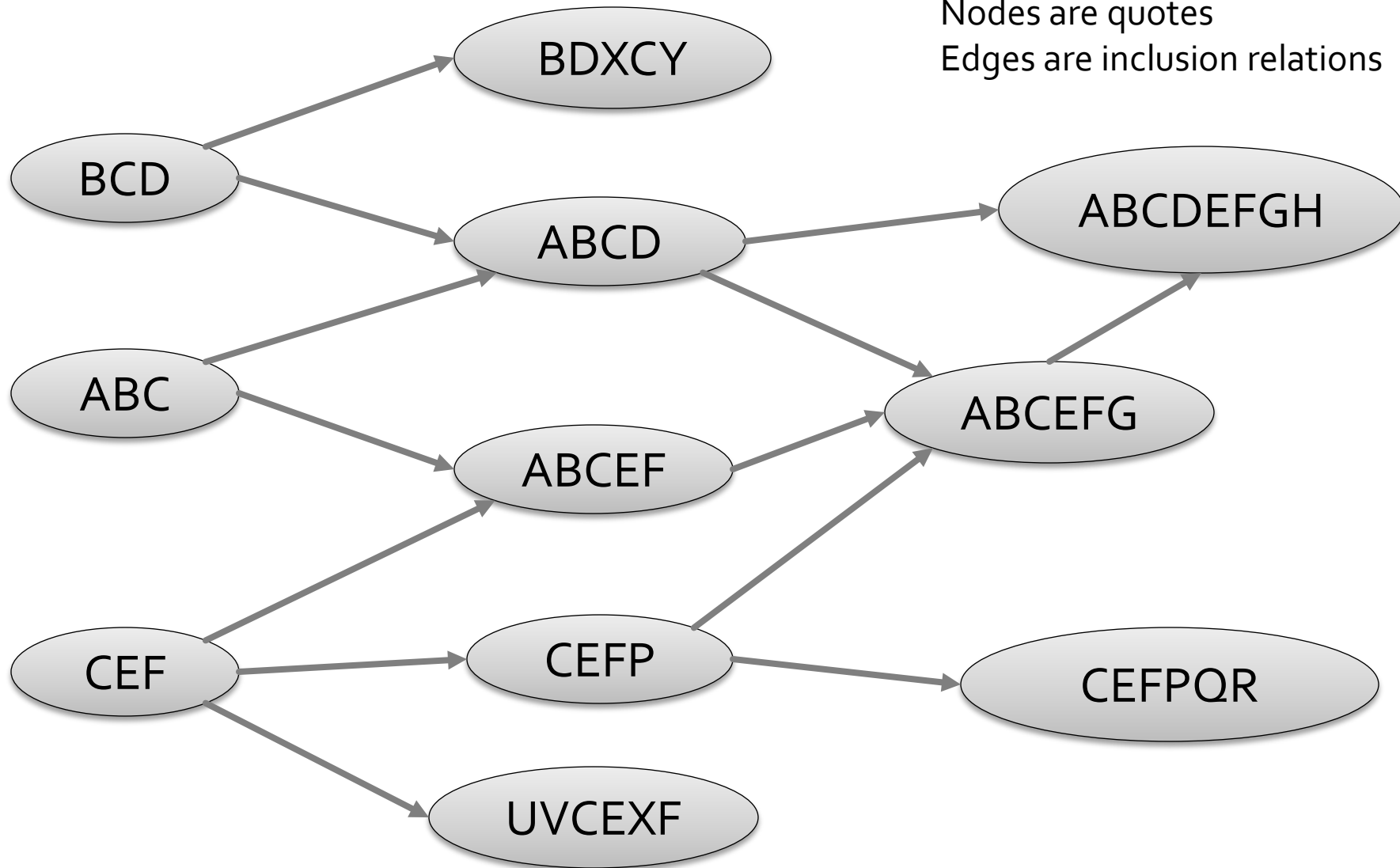
Creating clusters of Mutations

Nodes are quotes



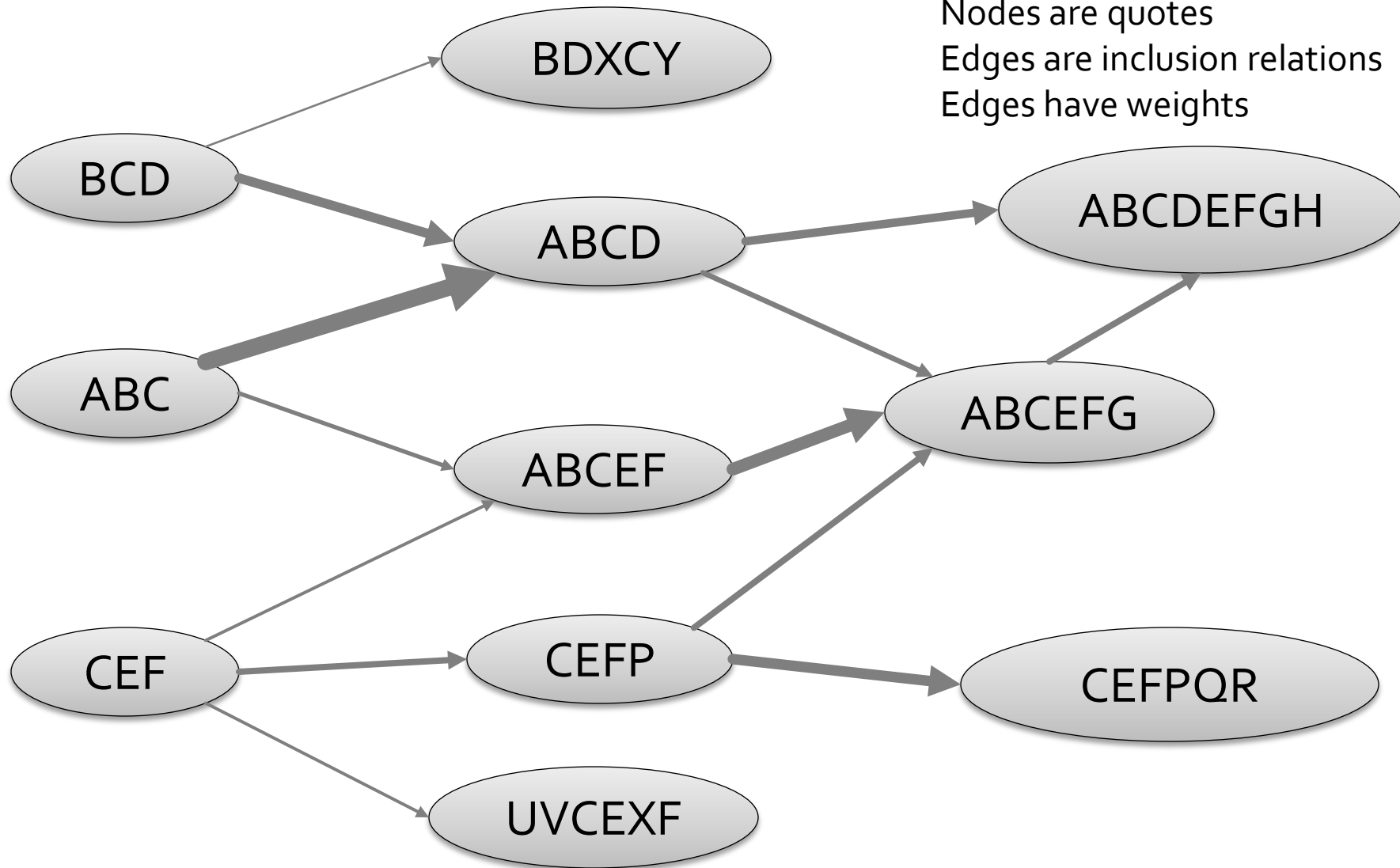
Creating clusters of Mutations

Nodes are quotes
Edges are inclusion relations



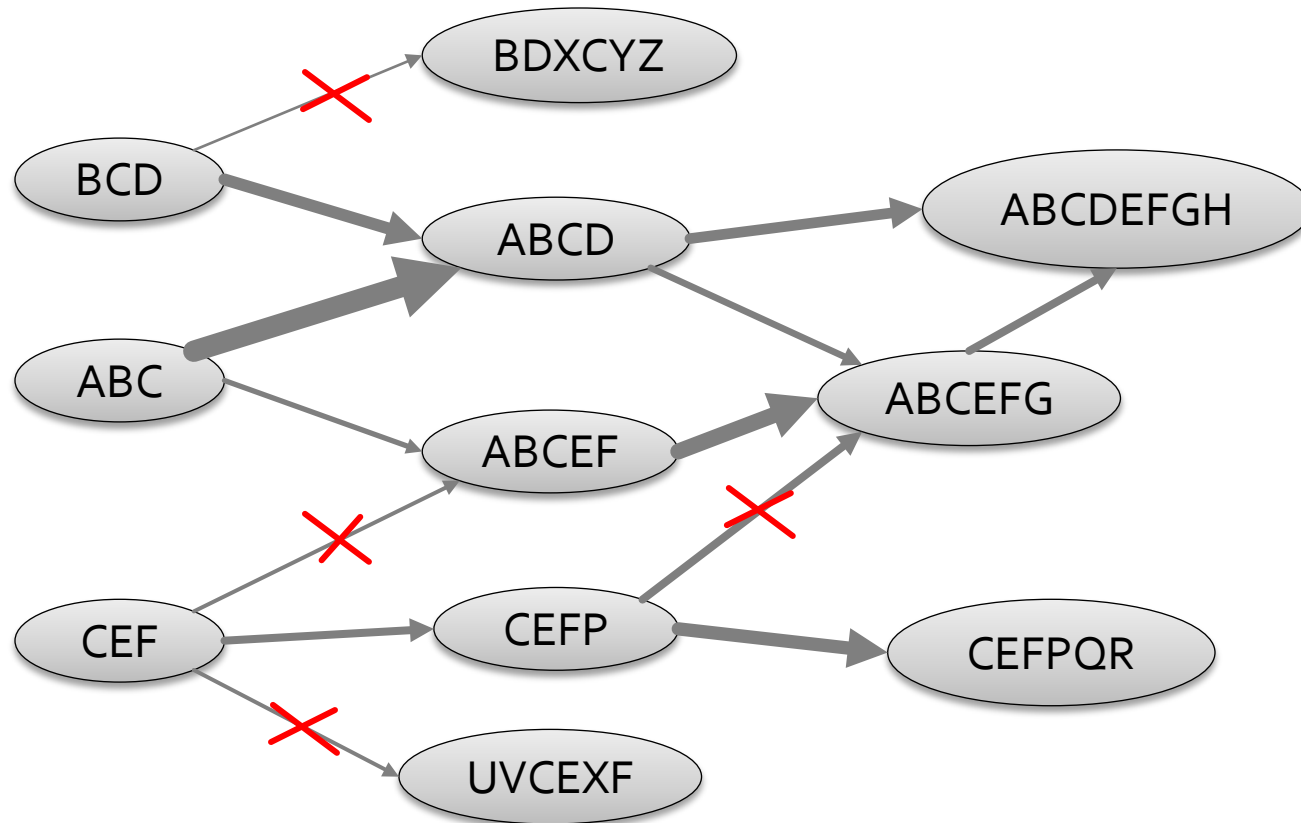
Creating clusters of Mutations

Nodes are quotes
Edges are inclusion relations
Edges have weights



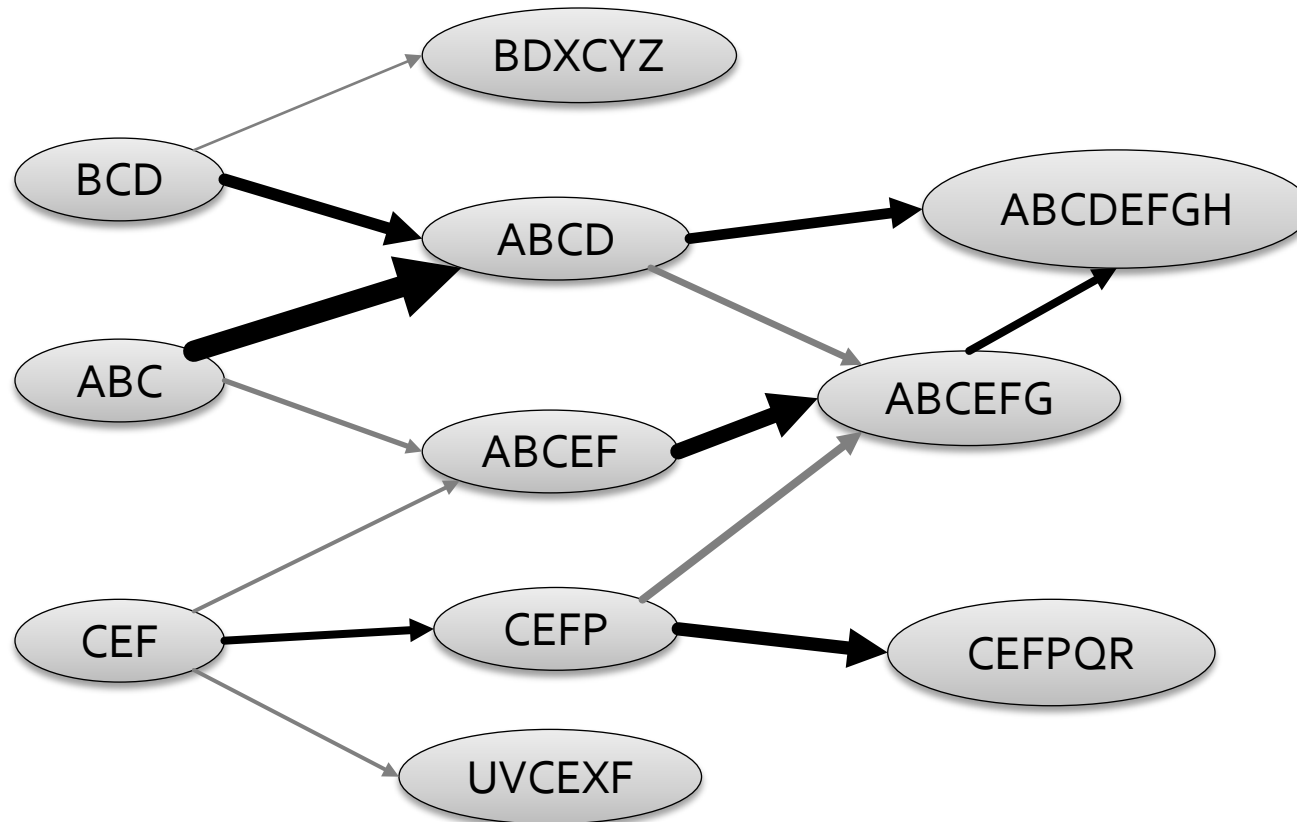
Quote clustering: DAG partitioning

- **Objective:** in directed acyclic graph (approx. quote inclusion), **delete min total edge weight** s.t. **each connected component has a single “sink” node**



Quote clustering: DAG partitioning

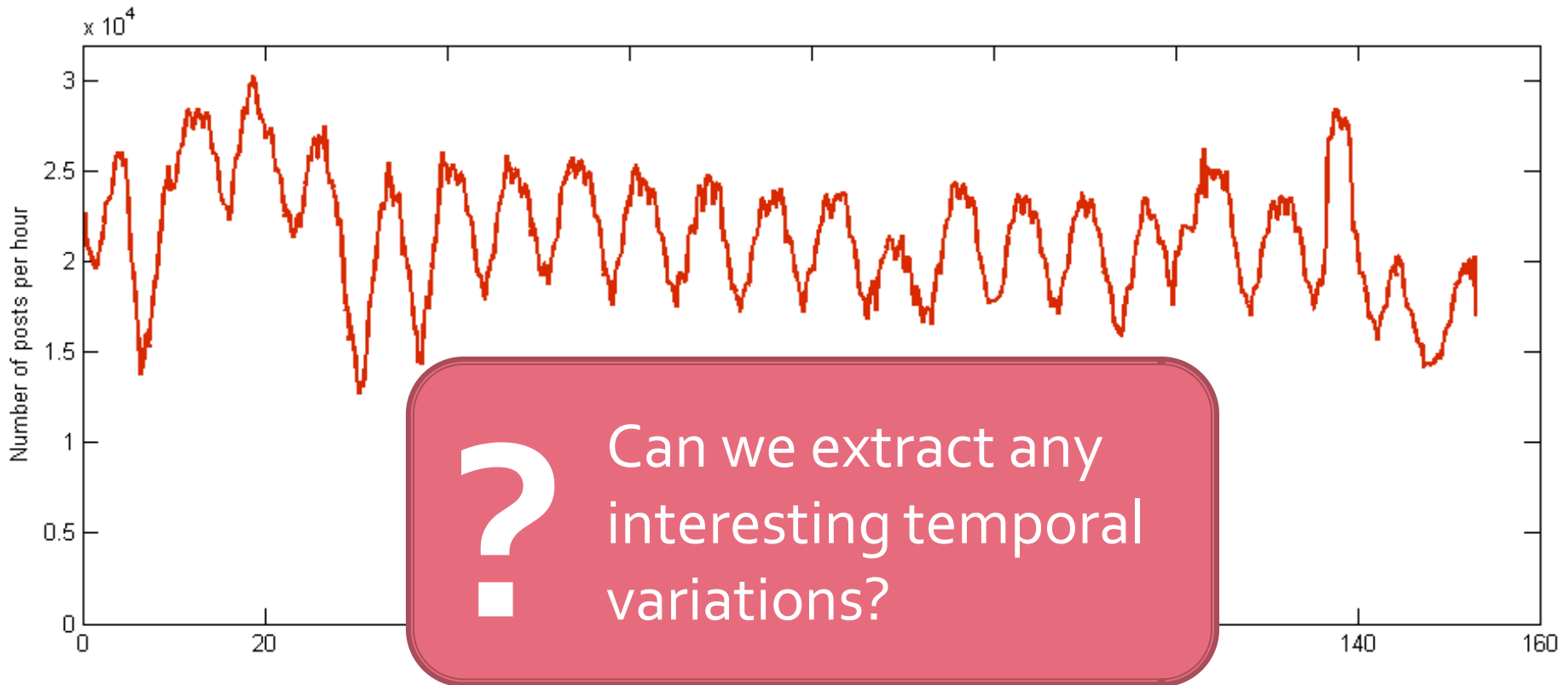
- **Observation:** enough to know node's parent
- **Heuristic:** proceed top down and assign node to strongest cluster



A quote cluster

Quoted text	Volume
the fundamentals of our economy are strong	3654
the fundamentals of the economy are strong	988
fundamentals of our economy are strong	645
fundamentals of the economy are strong	557
if john mccain hadn't said that the fundamentals of our economy are strong on the day of one of our nation's worst financial crises the claim that he invented the blackberry would have been the most preposterous thing said all week	224
fundamentals of the economy	172
the fundamentals of the economy are sound	119
i promise you we will never put america in this position again we will clean up wall street	83
the fundamentals of our economy are sound	81
clean up wall street	78
our economy i think still the fundamentals of our economy are strong	75
fundamentals of the economy are sound	72
the fundamentals of our economy are strong but these are very very difficult times and i promise you we will never put america in this position again	68
the economy is in crisis	66
these are very very difficult times	63
the fundamentals of our economy are strong but these are very very difficult times	62
do you still think the fundamentals of our economy are strong genius	62
our economy i think still the fundamentals of our economy are strong but these are very very difficult times	60
mccain's first response to this crisis was to say that the fundamentals of our economy are strong then he admitted it was a crisis and then he proposed a commission which is just washington-speak for i'll get back to you later	55
i still believe the fundamentals of our economy are strong	53
i think still the fundamentals of our economy are strong	50
cut taxes for 95 percent of all working families	50

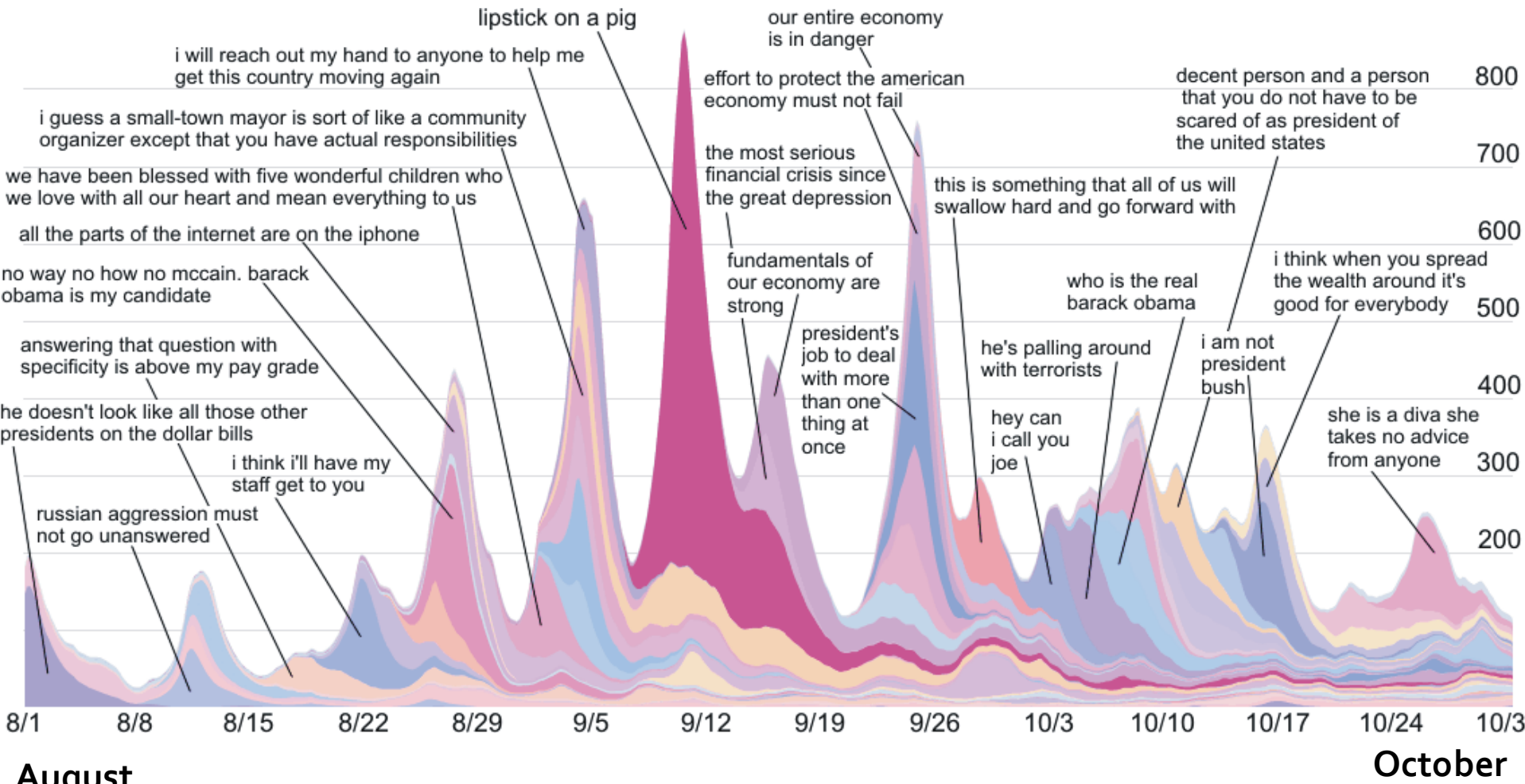
Articles/phrases over time



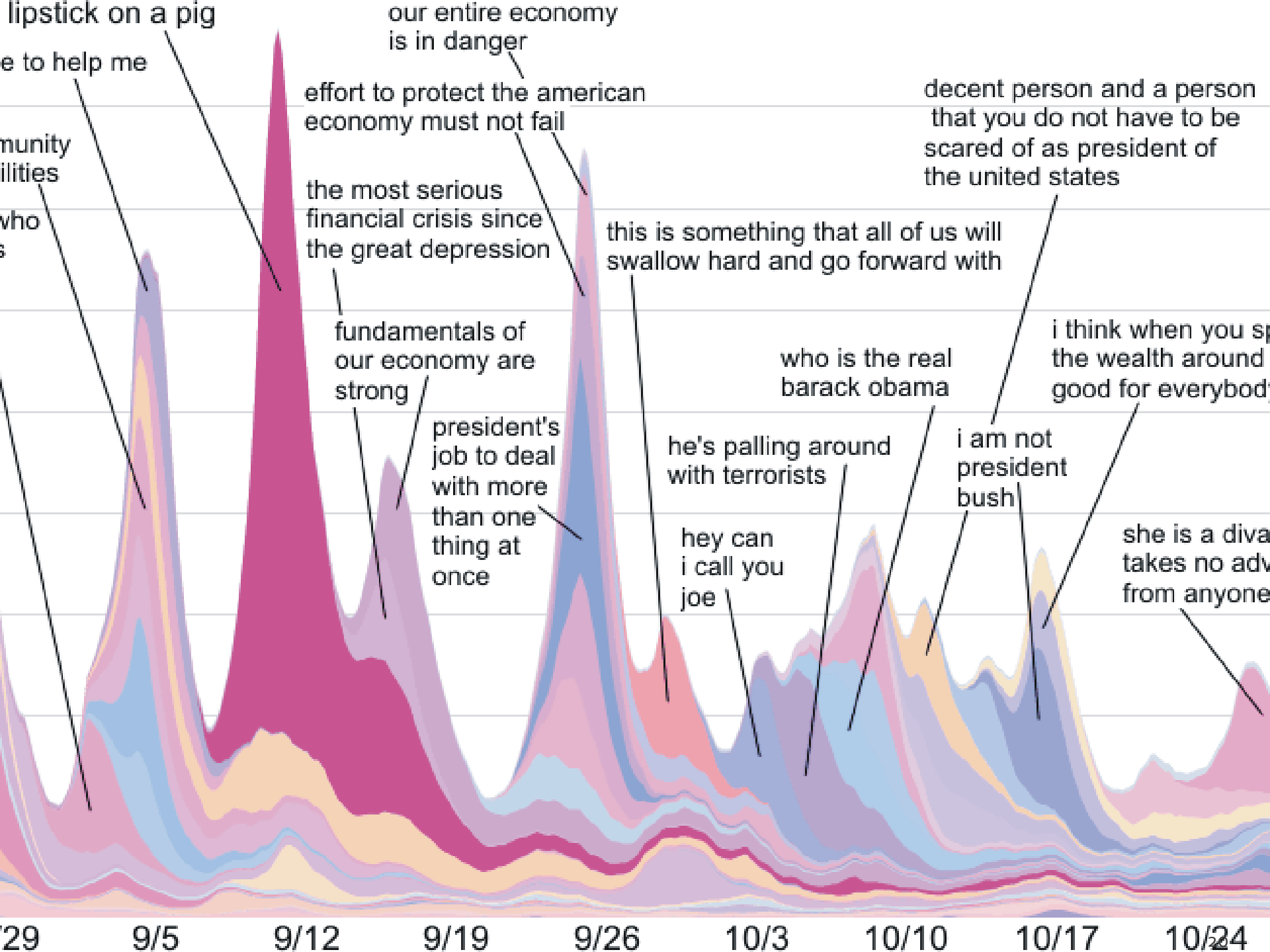
... is periodic, has no trends.

“Bandwidth” of the online media is constant

Cluster volume over time



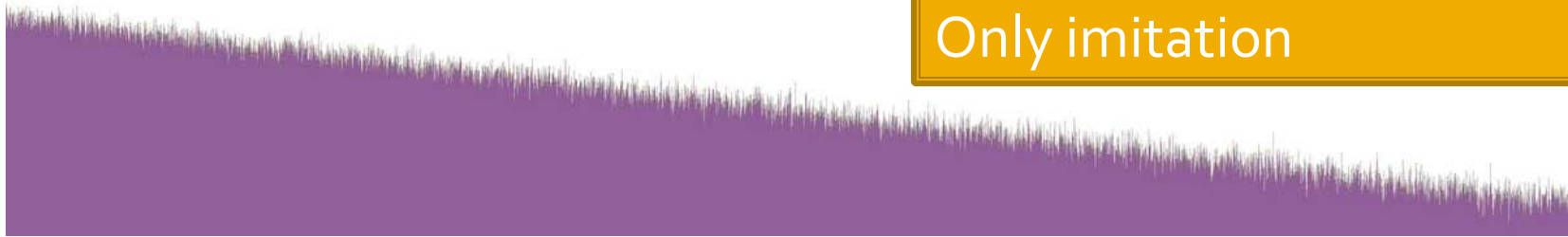
Volume over time of top 50 largest total volume quote clusters



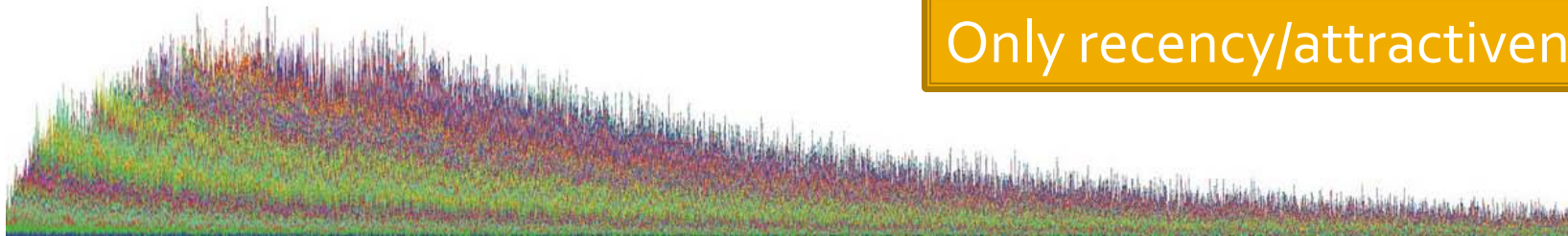
Modeling the temporal variation

- What ingredients are essential to qualitatively reproduce the observed dynamics?
 - Temporal variation has potential connections with natural processes
 - Species competing for resources in an ecosystem.
 - Biological systems synchronize to favor small number of individuals [Lacker-Peskin 1981]
- N news sources, one new story per time step. Source's choice of what to cover controlled by:
 - **Imitation**: increasing in number of sources covering story
 - **Recency**: decreasing in time since story's appearance
 - **Attractiveness**: prefer more interesting stories

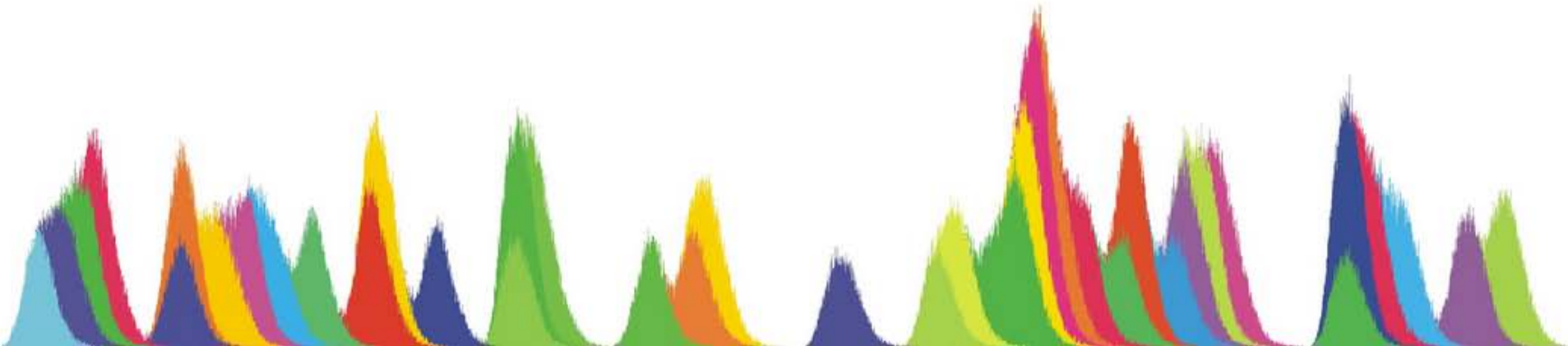
Modeling the temporal variation



Only imitation

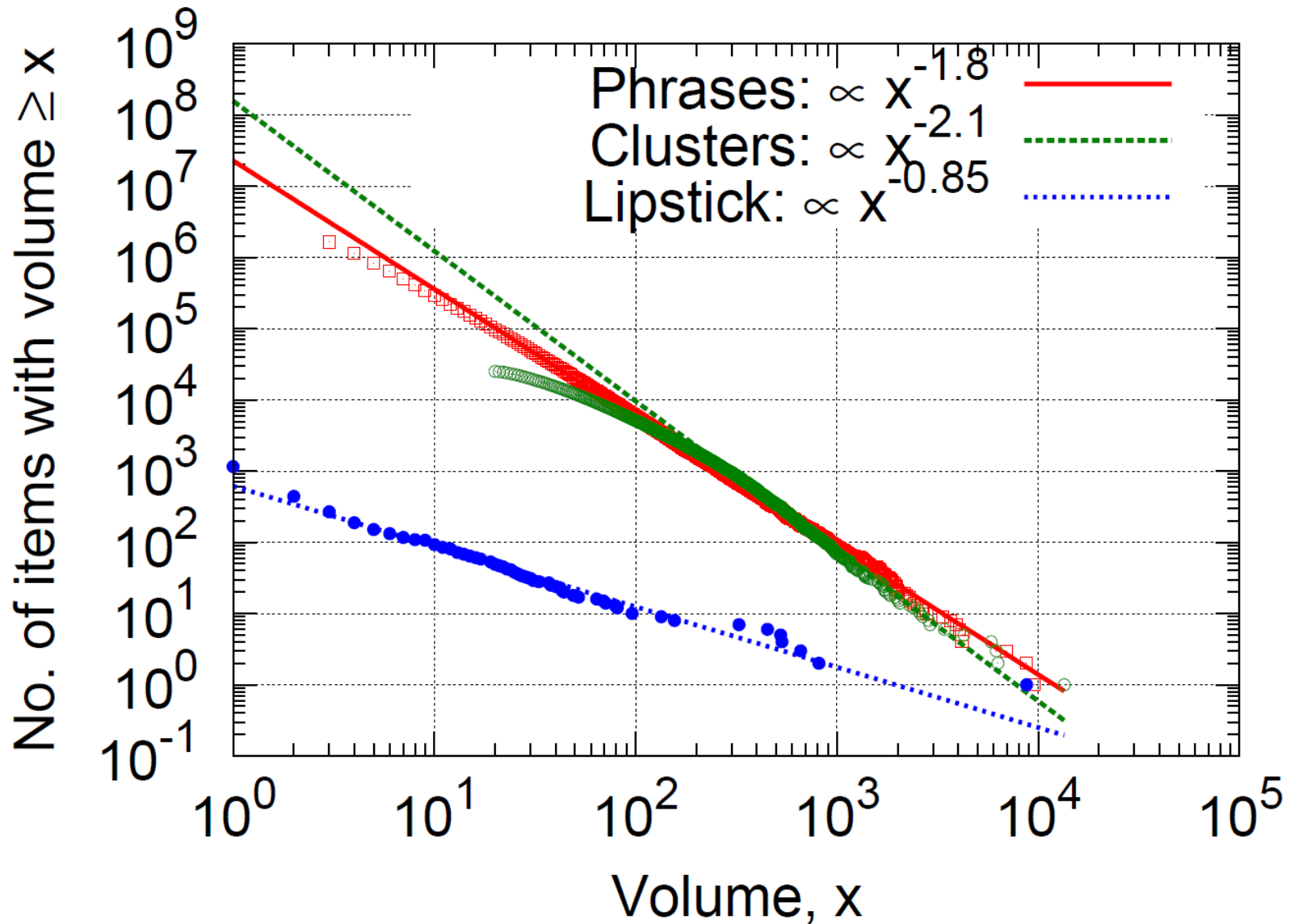


Only recency/attractiveness



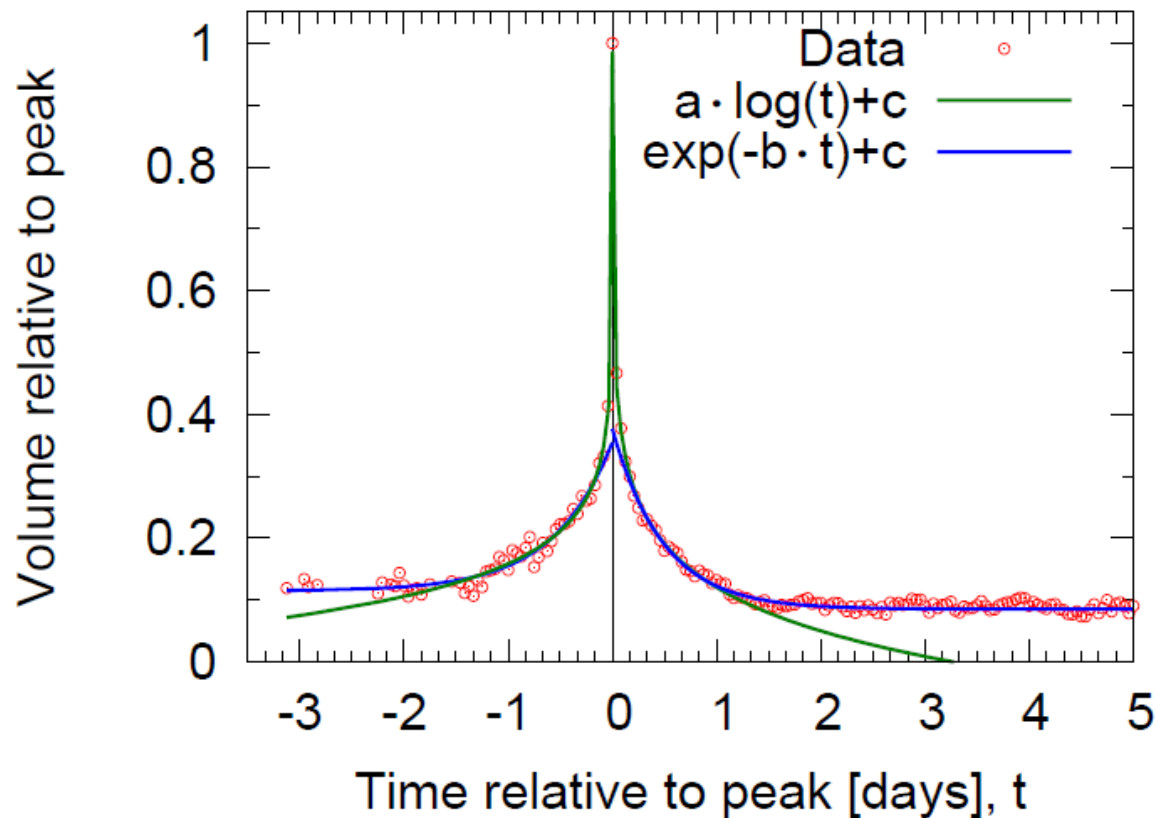
Imitation & Recency

Volume of phrases

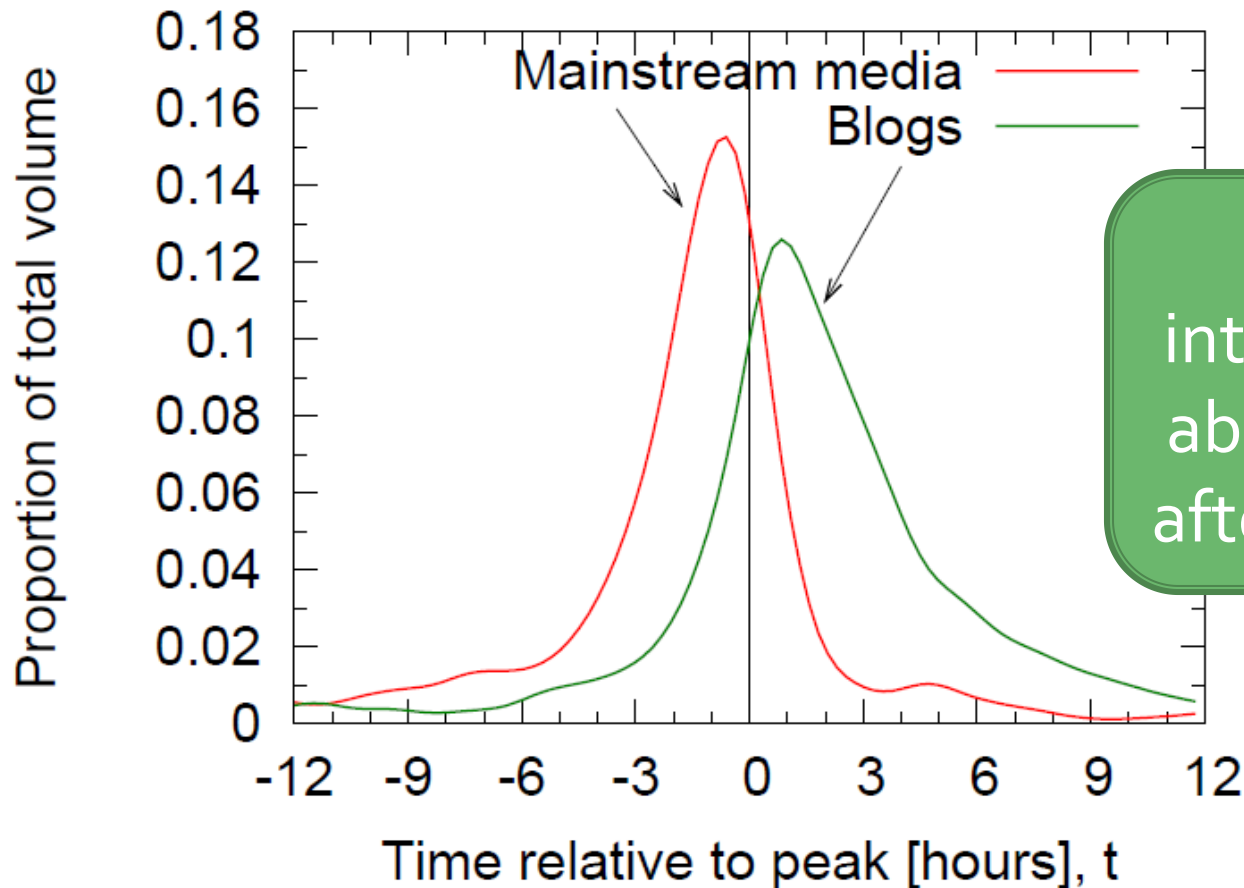


Interaction of News and Blogs

- Can study typical quote cluster volume curve
- Peak behaves like a delta function. Phrases are very short lived.



Interaction of News and Blogs



- Using Google News we label:
 - Mainstream media: 20,000 sites (44% vol.)
 - Blog (everything else): 1.6 million sites (56% vol.)

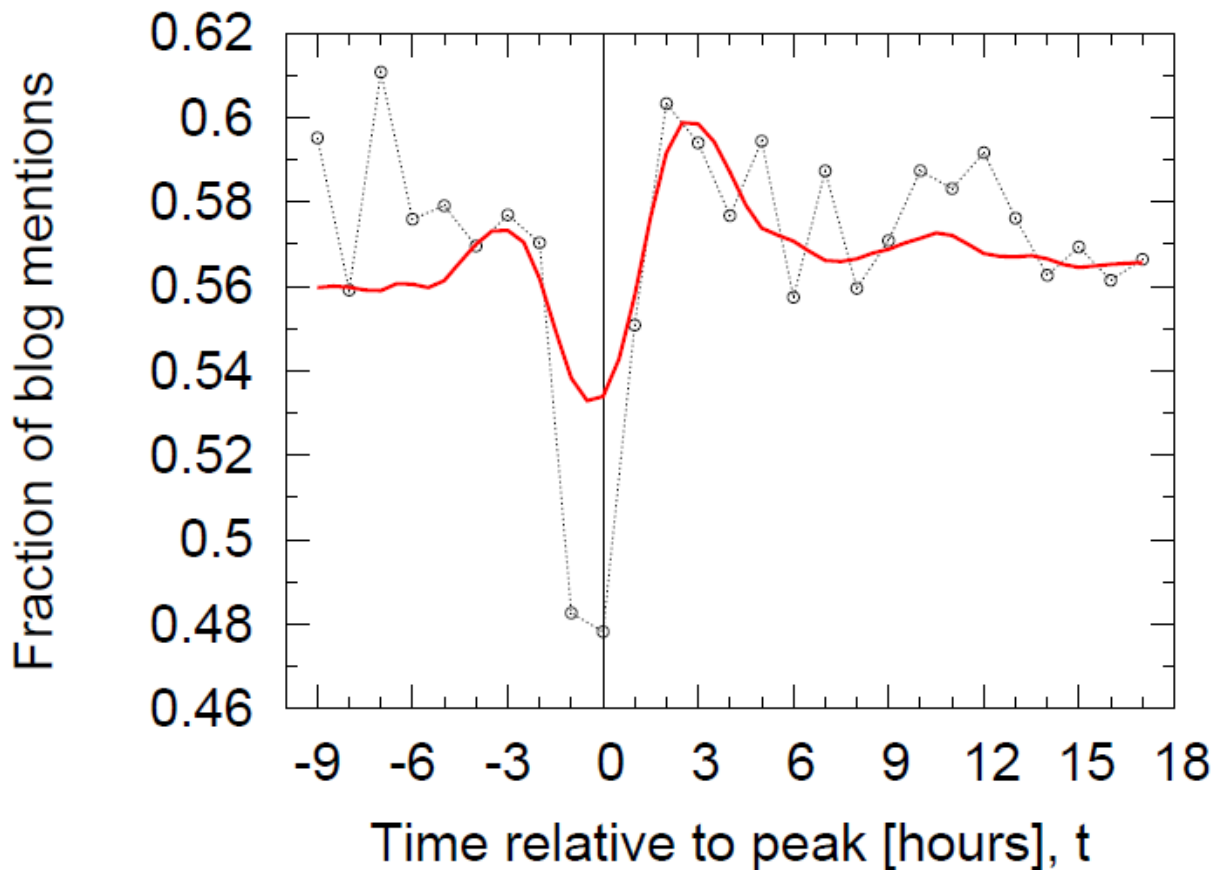
How quickly sites mention quotes?

- Can classify individual sources by their typical timing relative to the peak aggregate intensity

	Rank	Lag [h]	Reported	Site
Professional blogs	1	-26.5	42	hotair.com
	2	-23	33	talkingpointsmemo.com
	4	-19.5	56	politicalticker.blogs.cnn.com
	5	-18	73	huffingtonpost.com
	6	-17	49	digg.com
	7	-16	89	breitbart.com
	8	-15	31	thepoliticalcarnival.blogspot.com
	9	-15	32	talkleft.com
	10	-14.5	34	dailykos.com
	News media	30	-11	32
34		-11	72	cnn.com
40		-10.5	78	washingtonpost.com
48		-10	53	online.wsj.com
49		-10	54	ap.org

Interaction of News and Blogs

- Can study “oscillation” of attention between news and media



Stories catalyzed by blogs

- Can formulate queries for different temporal “signatures”: e.g., stories catalyzed by blogs:

$[x; y; t]$ -query: between x and y frac. of total quote volume (f_b) occurred on blogs at least t days before overall the peak

M	f_b	Phrase
2,141	.30	Well uh you know I think that whether you're looking at it from a theological perspective or uh a scientific perspective uh answering that question with specificity uh you know is uh above my pay grade.
826	.18	A changing environment will affect Alaska more than any other state because of our location I'm not one though who would attribute it to being man-made.

In total 3.5% of phrases migrate from blogs to media

Conclusion & Further questions

- A framework for tracking memes through the news, to quantify the dynamics of the news cycle.
- Demo + Data:
<http://memetracker.org>
- Many further questions:
 - Which elements of the news cycle do we miss?
 - Can this analysis of memes help identify dynamics of polarization? (cf. [Adamic-Glance, 2005])
 - How are these memes actually spreading among people?

THANKS!

Demo + Data:

<http://memetracker.org>