# Twitter Influencers Clustering and Personalized Recommendation

**Engage with influencers the right way.**

# Turn Data Into $$$?

Data          Popularity    →    $$$

Articles                           User

Web Traffic                        Advertiser

Social Media

# Turn Data Into $$$?

Data → Popularity → $$$

Articles

Web Traffic

Social Media
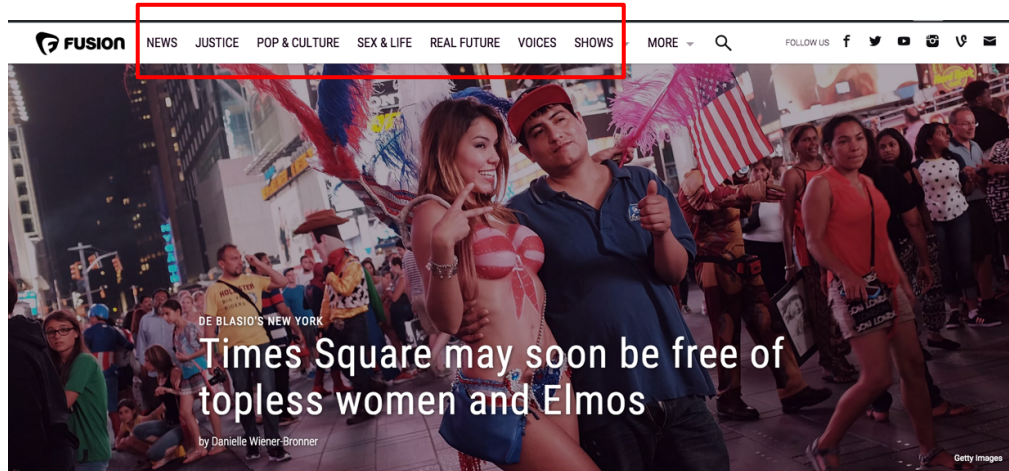
Give the right people the right content at the right times
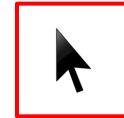
User

Advertiser

# Fusion Data Sources

# Mining Twitter

Who are influential to @thisisfusion:

- Participate in @thisisfusion tweets
  - retweet ⟵ Centrality Score
  - favorite
  - comment

  InfluenceFlow Score
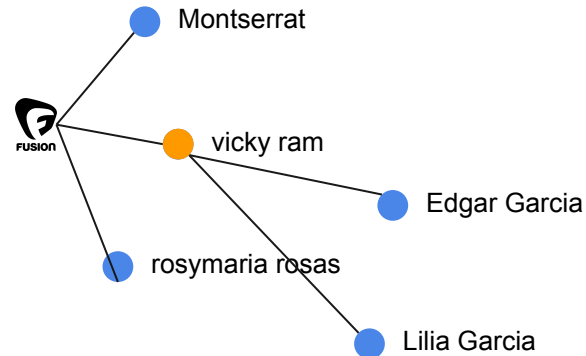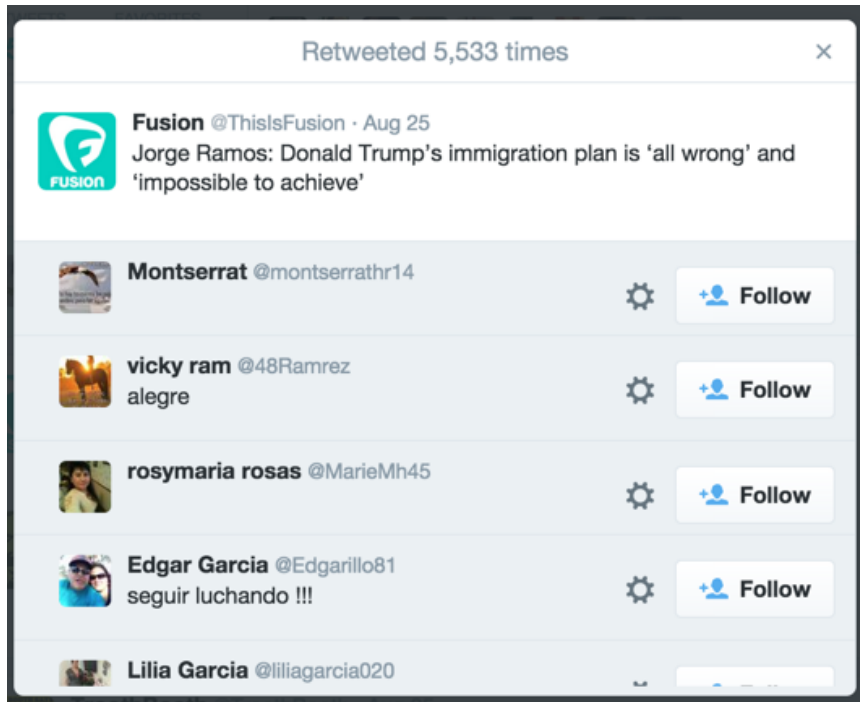
- Initiate tweets mentioning @thisisfusion
  - direct mention

# Mining One Tweet


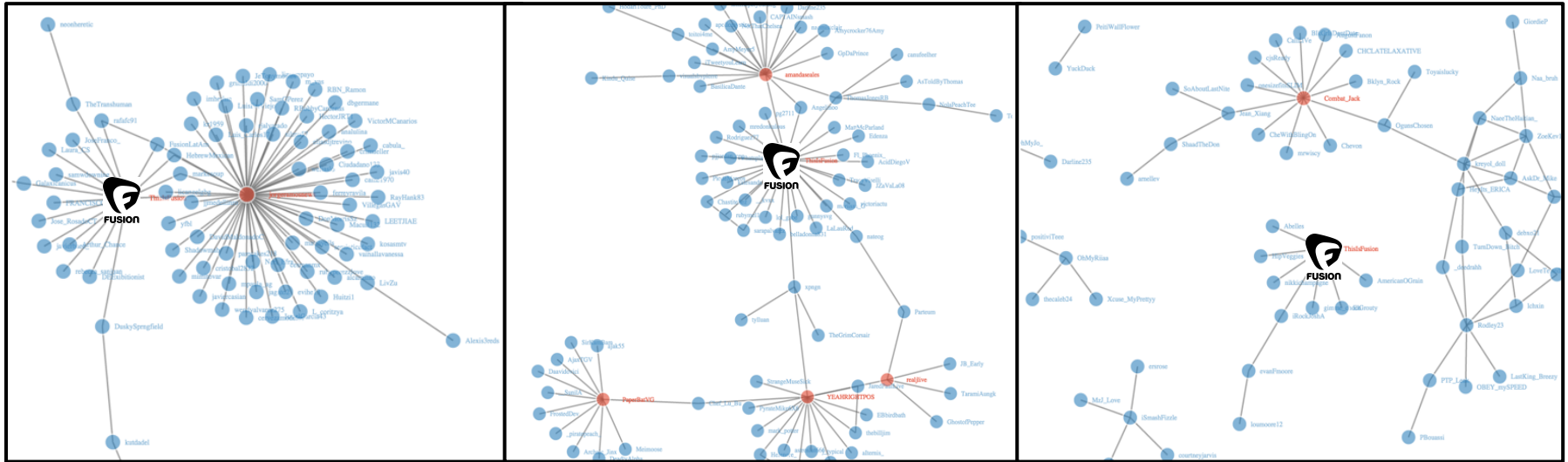
Centrality:
How important a user is within a retweet network.

Degree-centrality: No. of followers
Eigenvector-centrality: No. of important followers

# Mining One Tweet

How does each tweet transmit?

# Mining Fusion Community

Not limited to specific tweet, InfluenceFlow Score captures overall information flow in Fusion community:

| | User | Latin.America | follower | * | mention | = | score |
|---|------|---------------|----------|---|---------|---|-------|
| 1 | jorgeramosnews | 1 | 1439674 | | 29 | | 41750546 |
| 2 | rafafc91 | 1 | 301 | | 62 | | 18662 |
| 3 | FusionLatAm | 1 | 142 | | 125 | | 17750 |
| 4 | TheTranshuman | 1 | 1646 | | 5 | | 8230 |
| 5 | Arthur_Chance | 1 | 1538 | | 4 | | 6152 |
| 6 | Laura_CS | 1 | 1070 | | 3 | | 3210 |
| 7 | yfbl | 1 | 2485 | | 1 | | 2485 |

InfluenceFlow Score

Include:
No. of Retweets
No. of Direct mentions
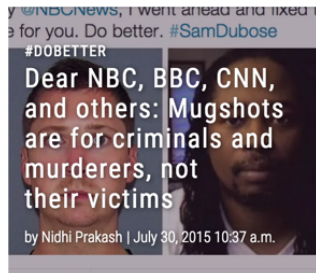
# Right People: Superfans

- Collect superfan data from every tweet
  - Scores
    - Centrality - influence within one tweet
    - InfluenceFlow - overall influence
  - Interests
    - Twitter Hashtags
    - Article Sections (news/justice/voices...)
    - Article Topics (drugs/transgender/mexico...)
    - ...

# Live Demo

# Live Demo

# Live Demo

# **Right Content?**

Is there an algorithm to suggest personalized recommendations to superfans?

In Progress:

- Text mining on all superfans timeline content
- Clustering texts on vectors from AlchemyAPI
- Observe distributions of superfan scores and interests in different clusters
- Content-based recommendation (match vectorized user timeline text and article text using cosine similarity)

# Right Content?

| | Unnamed: 0 | Academic | AcademicInstitution | Accommodation | Actor | AdministrativeDivision | Airline | Airport | AirportOperator |
|---|---|---|---|---|---|---|---|---|---|
| 0 | LaurenLaCapra | 0.000000 | 0.392141 | 0 | 0.650576 | 1.810971 | 0 | 0 | 0.000000 |
| 1 | Sexyman469 | 0.307363 | 0.000000 | 0 | 0.000000 | 1.506483 | 0 | 0 | 0.000000 |
| 2 | alonso_529 | 0.296390 | 0.000000 | 0 | 0.202315 | 1.155570 | 0 | 0 | 0.000000 |
| 3 | IvankaOC | 0.000000 | 0.000000 | 0 | 1.598437 | 1.180828 | 0 | 0 | 0.000000 |
| 4 | AngryYuca | 0.423798 | 0.000000 | 0 | 3.961219 | 1.176333 | 0 | 0 | 0.000000 |
| 5 | TooTurntNelly | 0.000000 | 0.000000 | 0 | 0.000000 | 0.309068 | 0 | 0 | 0.000000 |
| 6 | natashalennard | 0.000000 | 0.613226 | 0 | 0.471438 | 1.602814 | 0 | 0 | 0.000000 |
| 7 | paul_boyd | 0.296297 | 0.000000 | 0 | 0.491203 | 1.323246 | 0 | 0 | 0.000000 |
| 8 | SirKamBam | 0.000000 | 0.000000 | 0 | 0.444928 | 0.906519 | 0 | 0 | 0.000000 |
| 9 | SteampunkMuppet | 0.000000 | 0.000000 | 0 | 0.876238 | 0.618926 | 0 | 0 | 0.000000 |
| 10 | _piratepeach_ | 0.000000 | 0.435684 | 0 | 0.000000 | 0.970973 | 0 | 0 | 0.000000 |
| 11 | hatterfan | 0.000000 | 0.000000 | 0 | 0.000000 | 3.599519 | 0 | 0 | 0.000000 |
| 12 | drseid | 0.378882 | 0.000000 | 0 | 0.000000 | 1.216191 | 0 | 0 | 0.000000 |
| 14 | Biibekk | 0.000000 | 0.000000 | 0 | 0.299256 | 0.973483 | 0 | 0 | 0.000000 |
| 15 | bequarius | 0.460585 | 0.000000 | 0 | 0.594332 | 0.459007 | 0 | 0 | 0.000000 |
| 16 | AceHoffman | 0.349410 | 0.406325 | 0 | 0.000000 | 1.329949 | 0 | 0 | 0.000000 |

# Limits

- Time Limit (3-week-project)
    - Mainly focused on identifying and classifying superfans
    - Didn't include time analysis


- Data Limit (Twitter API limitations)
    - Rate limit - takes 1 to 2 days to get basic data needed for 20 tweets
    - Query limit - can't retrieve more than 3,200 tweets and 100 retweets, which sacrificed the flexibility of popularity analysis
    - Lack of infrastructure support - data consistency issue occurred a lot

# Make It Scalable

# About Us

This is an open source project collaborated between NYC Data Science Academy and Fusion media (an ABC-Univision joint-venture).

**--NYC Data Science Academy**

Fangzhou Cheng (PM, Data Science Fellow)

Shu Yan (Data Science Fellow)

Alex Singal (Data Science Fellow)

**--Fusion**

Noppanit Charassinvichai (Data Engineer)