



Project Acronym: **CATALYST**
Project Full Title: **Collective Applied Intelligence and Analytics for Social Innovation**
Grant Agreement: **6611188**
Project Duration: **24 months (Oct. 2013 - Sept. 2015)**

Masters of Networks 3 – Track 1 report
What makes of community a community?

Date: **March 10-11, 2015**
Location: **Rome, Italy**
Acknowledgement: **All participants of the Masters of Networks 3 event**



This project has received funding from the European Union's Seventh Framework Programme for research, technological development and demonstration under grant agreement n°6611188

Investigating the *Matera-Lote4* twitter community

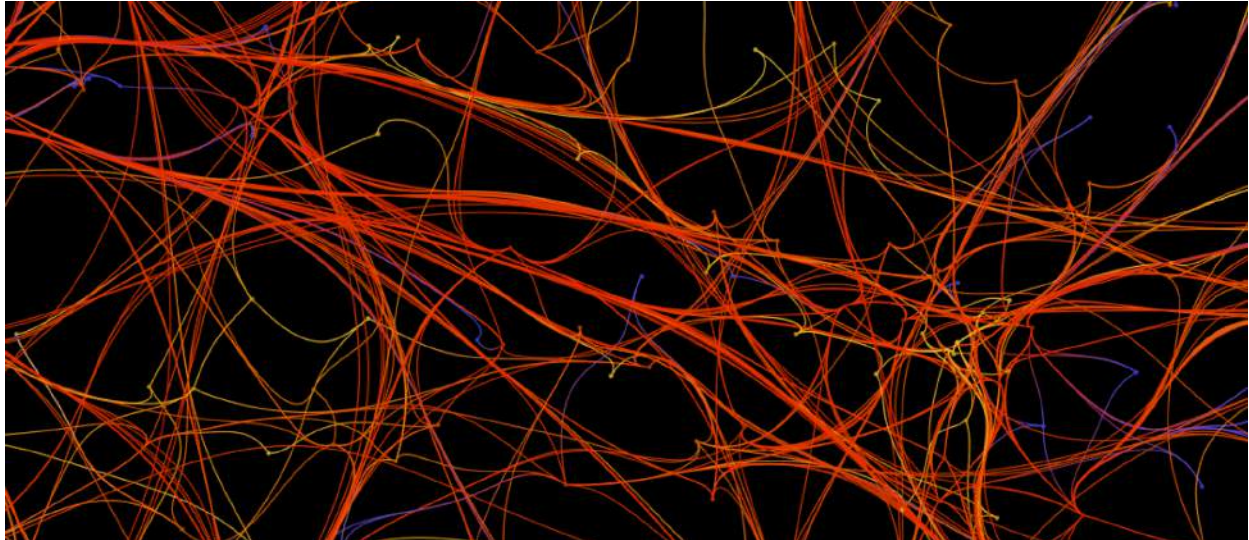


Figure 1 - A view at heart of the twitter nebula

When organizing an event, we can reach out using Twitter, but what is at the heart of these Twitter discussions and on what's the (new?) definition of a community in the twitter space?

We try to describe in this document a study of an analysis of Twitter, with a top-to-bottom approach: from the top-level general questions around the idea of a community we drill down to the bottom-level data gathering, analysis, and finish with visualization of the data.

Through this document, discussions with experts has redefined the notion of community along three different axes and analyse our twitter communities around these axes.

The document presents first the highest level questions we have wondered on this Twitter community, then the data gathered, followed by the ground baselines on which we've looked at communities, and finally the analysis

1. The question

So we've gathered many experts, our Masters of Networks (Community Managers Lee-Sang Huang, Noemi Salantiu, Laura Manconi, Rosa Strube and Collective Intelligence Researchers Marta Arniani, Yannis Treffot, Benoit Gregoire and Network Scientist Benjamin Renoust).

The first idea was to find the many questions that we can build around the Lote4 twitter community. The main idea was to identify first if and how a Twitter “conversation” around a hashtag can form a community?

One of our assumptions is that we can find that there are lots of isolated components in a Twitter hashtag stream, with people not really calling out to each other, whereas “tight” online community gives rise to a giant component that most of the nodes are connected to. Can we confirm that from a data perspective?

Another assumption is that people could easily form subgroups investigating specifics, and how we can find traces of, or understand the content gathering these subgroups?

2. The data & context

To that purpose, we have analyzed twitter data made available during the MoN3 event. But before talking about the data itself, we may mention some contextual information on how this data has been captured.

This data represents Tweets collected between 18/10/2014 and 23/02/2015, this data has been gathered from the search query “*#lote4 OR edgeryders OR unmonastery OR Matera*”.

Lote4 stands for Living On The Edge, a conference organised by the Edgeryders global community which took place in Matera (Italy) between the 23rd and the 26th of October 2014. The unMonastery was an artist and hacker residency program that Edgeryders ran in Matera during most of 2014; the Lote4 conference happened in the unMonastery building and with the help of the unMonastery events. The search string was expected to catch all tweets around the event taking place in Matera and their follow up over subsequent months.

Organizing the event, we expect authors such as Edgeryders and Matera2019 to be sort of moderators of the event, and to engage other twitter users from their own networks in driving conversations among the different participants.

The collection is composed of about 20k Tweets written by 7000 people involving another 1000 additional people (via mentions or RTs).

The data has information on who sends a Tweet, eventually whom the Tweet is sent to, who is mentioned in the Tweet, the date and which hashtag has been used.

3. The notion of community

We need to step back a bit here and question ourselves on what makes a community a community. Social network scientists such as myself have preconceived established notions of

what makes a community in terms of data analysis, but they all end up being empirical and somewhat fitting well in the boots of data analysis. Two main definitions might be recalled, the first one would come from Manski, for whom the group effect builds from gathering people alike (translated in data processing this would mean similarity of attributes). The second definition is used as a support to compute the Newman's modularity, states that a community has much more relationships between members within itself than with other members from outside the community. Discussions led by our panel of experts has pushed even further the different definitions of the notion of community.

We've tried to bring different perspectives on the notion of community, by asking these questions "what does it mean to be a community? what does it mean to belong to a community? what does it mean to look at a community?".

We gathered many answers which faceted a bit the notion of "community", into three main categories, and we've also discussed some other interesting characteristics of communities.

Awareness	Exchanging/discussion	Action
<i>Sense of belonging/endeavor</i> <i>inner sense of belonging</i> <i>sharing of interest</i> <i>some commonalities</i> <i>publish on similar twitter hashtag</i> <i>gather around specific goal</i> <i>share content</i>	<i>People talking to one another</i> <i>exchange in the community (both ways)</i>	<i>follow the same people / sign petition</i> <i>actually meet / community of practice</i> <i>people who do more than what they need to/have to behavior,</i> <i>can be negative engagement</i>
Other characteristics		
<i>groups = set of people</i> <i>transversal to existing organization</i> <i>classes of communities hierarchical</i> <i>over time can fuse or divide</i> <i>somebody who's not in the community (the rest of the world)</i>		

Table 1 - Summary of different characteristics a community can have

Our experts have extracted 3 levels that define a community from this point:

- the awareness
- the exchanges/discussions
- the action/actual engagement

4. Community in Twitter

The next step is to reconnect these notions of a community with evidences we can find in the data. In other words, what does “awareness”, “exchanges” and “engagement” mean in the context of Twitter publications?

In the world of Twitter, someone’s awareness can be measured by the semantics these people use, i.e. the hashtags they use in publishing so we can measure how often these keywords appear, the number of other users who relate to the same semantics, and the presence of our users and their posts on different platforms.

Sharing and exchanging between users is the basic purpose of such a micro blogging platform. This type of interaction can be represented in the world of Twitter by mentioning somebody or replying to somebody: it never means that an actual engaging conversation is going on, but it initiates potential collaboration.

Other measures can be of interest in tracing the interactions within a community (well, when the community is defined already). The number of connections (following, followers) of a user, their amount and frequency of posting, and the impact of the posts: how do others in the community endorse the posts? does it generate spin offs? all within the community, and out of the community? How to measure this impact? of each post? of each individual?

The network of practice in the material engages people in meeting and actually doing things together, working toward a common goal. We could find traces of physical presence, at events for example, of people from the geocoding, the hashtags they use, when related to an event, or via cross platforms activities such as FourSquare. Unfortunately these indications are not really reliable when confined to the sole Twitter information. Engagement on Twitter can take different shapes, it can mean reciprocal interactions, with the production of content and maybe some spin off actions. One reliable action on Twitter is the construction of actual conversation between people, meaning people replying to one another, reciprocally not only broadcasting information, or commenting on shared interest but real conversations.

Among the other characteristics of a community discussed, the most interesting would probably be the influence of time on the group evolution (fusing/dividing), but we’ll keep these aspects for a different analysis.

Notice that we have yet taken the “RT” or “retweet” relationship out of the picture as it is a versatile information. This action is the easiest and most represented action of the Twitter universe. The act of retweeting can bear two different meanings. It first helps showing your interest, people retweeting similar posts are show similar awareness, but it can be either

positive or negative engagement. The number of retweets can weigh actually the tweets, because when a tweet has been retweeted a lot can be considered as “impactful” or just popular.

5. The analysis

5.1 General analysis of the twitter data

So the dataset is composed of ~20,000 gathered from the query string “#lote4 OR edgeryders OR unmonastery OR Matera” between 18/10/2014 and 23/02/2015, published by ~7,000 different authors, involving 8,000 people including mentions and replies, and referring to over 6000 hashtags.

Here is the production of tweets, we can clearly see a few activity peaks around the event, first during the period of preparation, before the actual event, then during the event.

The tweets also peak around half November, and more activity can be noticed between the end of December and the beginning of January.

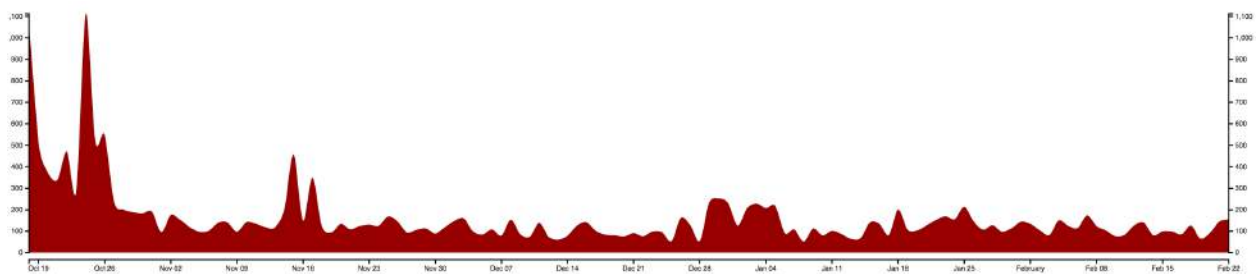


Figure 2 - The production of tweets over the period of capture

Among the 7000 authors, only 42 have produced over 50 tweets in the period of time, and 18 users have only retweeted information, about 850 more have published over 5 tweets in this period of time (and actually 200 are only retweeting information).

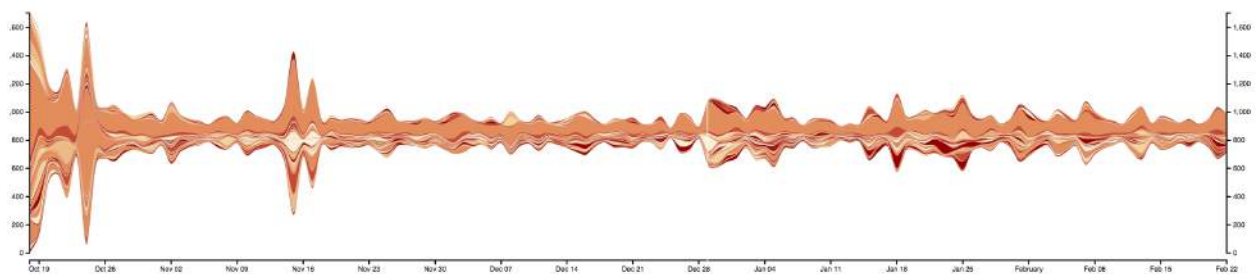


Figure 3- Occurrence of hashtags over time. Of course, #matera is the most occurrent over time in the dataset. Spikes correspond partly to #btwic2014 & #lote4 (before Oct. 26) #saleritana, #basilicata #matera2019 #mendicino (from Nov. 02) #labuonascuola #vivoazzurro #under21 #matera2019 (around Nov. 16) #capodanno #genova #neve #matera2019... (end of Dec, early Jan.)...

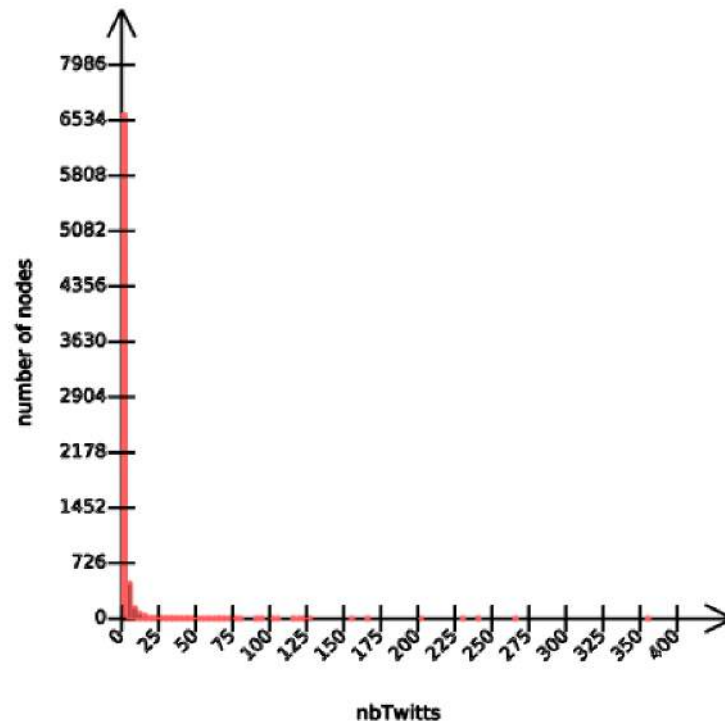


Figure 4 - Distribution of users (nodes) per number of tweets produced

So retweets generate a background noise and we'll keep them apart for a secondary analysis. When we remove the RTs, we can consider a total of 4500 twitterers replying and mentioning each other.

The network they compose is very disconnected, and half users captured here discuss in small groups of at most 10 people, producing each very little tweets. However the other half users (around 2100) are involved in a gigantic twitter conversation.

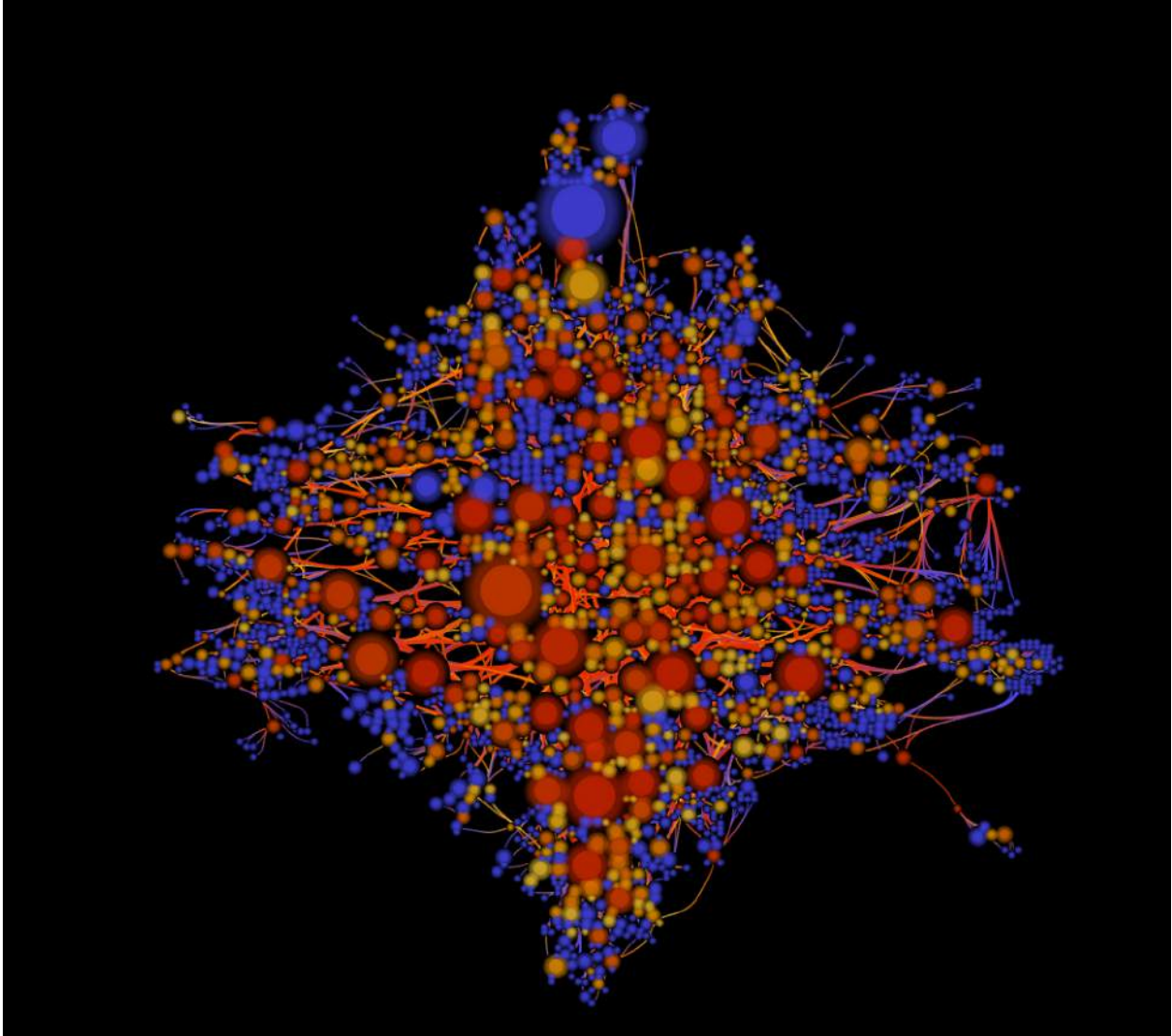


Figure 5 - The main connected component of twitterers, each node is a twitterer, each link a reply or a mention between two twitterers

2100 twitterers discussing about Matera, Lote4 and many other things. The size of a node is the number of tweets produced in the collection, the color of a node is its centrality. An edge means a direct reply and/or mention between two users.

5.2 Communities in Twitter

Drilling down to the heart of the community

Following the previously defined criteria, we've tried to define how is this community composed around the twitter hashtags, mentions and replies. Because we're looking for the strongest evidences of "communal" behaviors between twitter users, we advanced quick towards traces of

engagement between users. We have therefore considered first the “reply” relationship between users, and we’ve drilled down to only 500 of the 2100 users who are actually replying to each other in a big conversation (1600 are involved in small conversations).

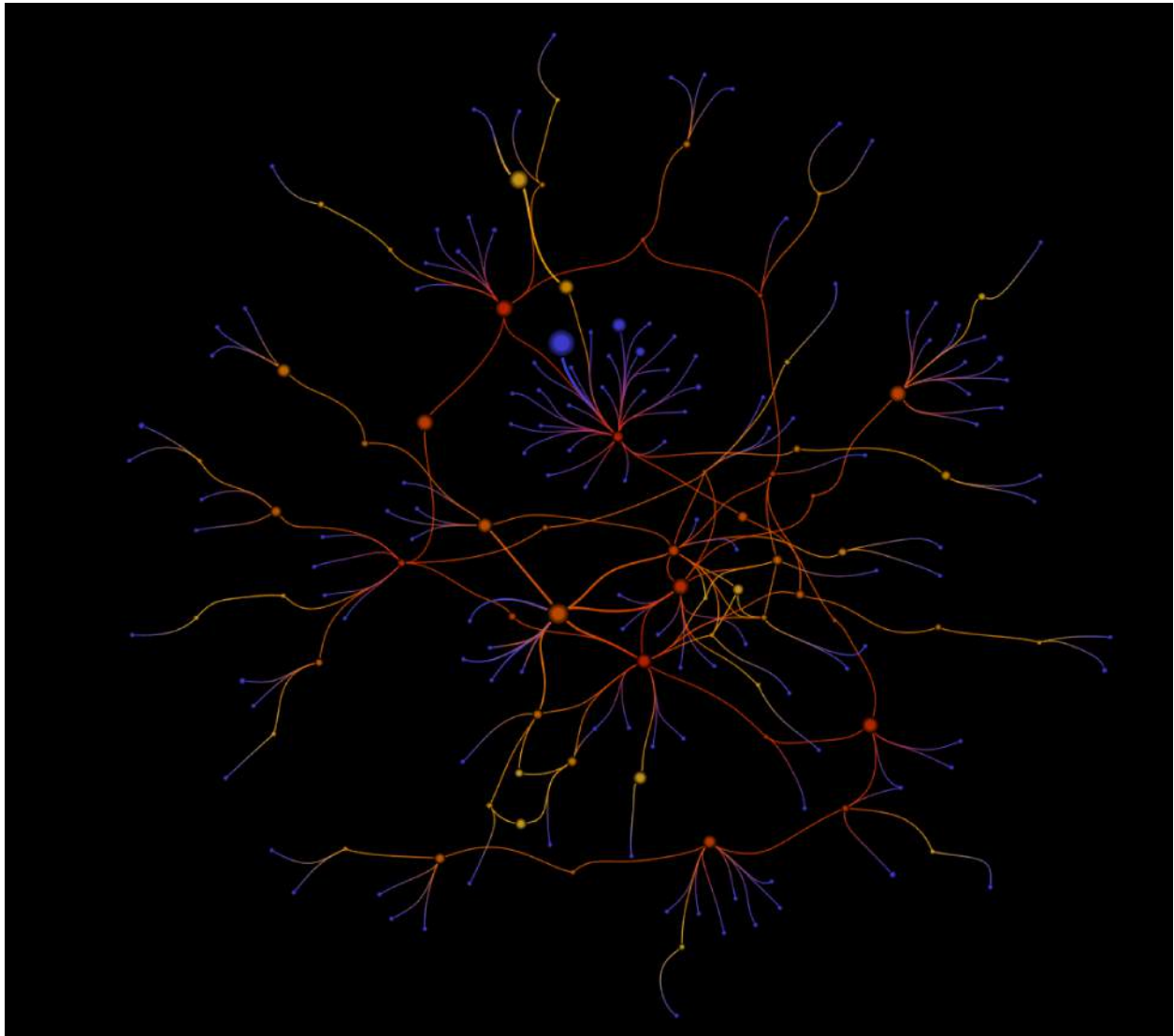


Figure 6 - 500 users replying to each other. We can notice the arborescent structure of some nodes.

Now, looking for the strongest ties, we want to subset even further these conversations to identify actual traces of reciprocal conversations, i.e. replies over replies.

Only 21 people are actually taking part of conversations involving more than just a triplet of users, and this group is actually divided into 2 disconnected subgroups.

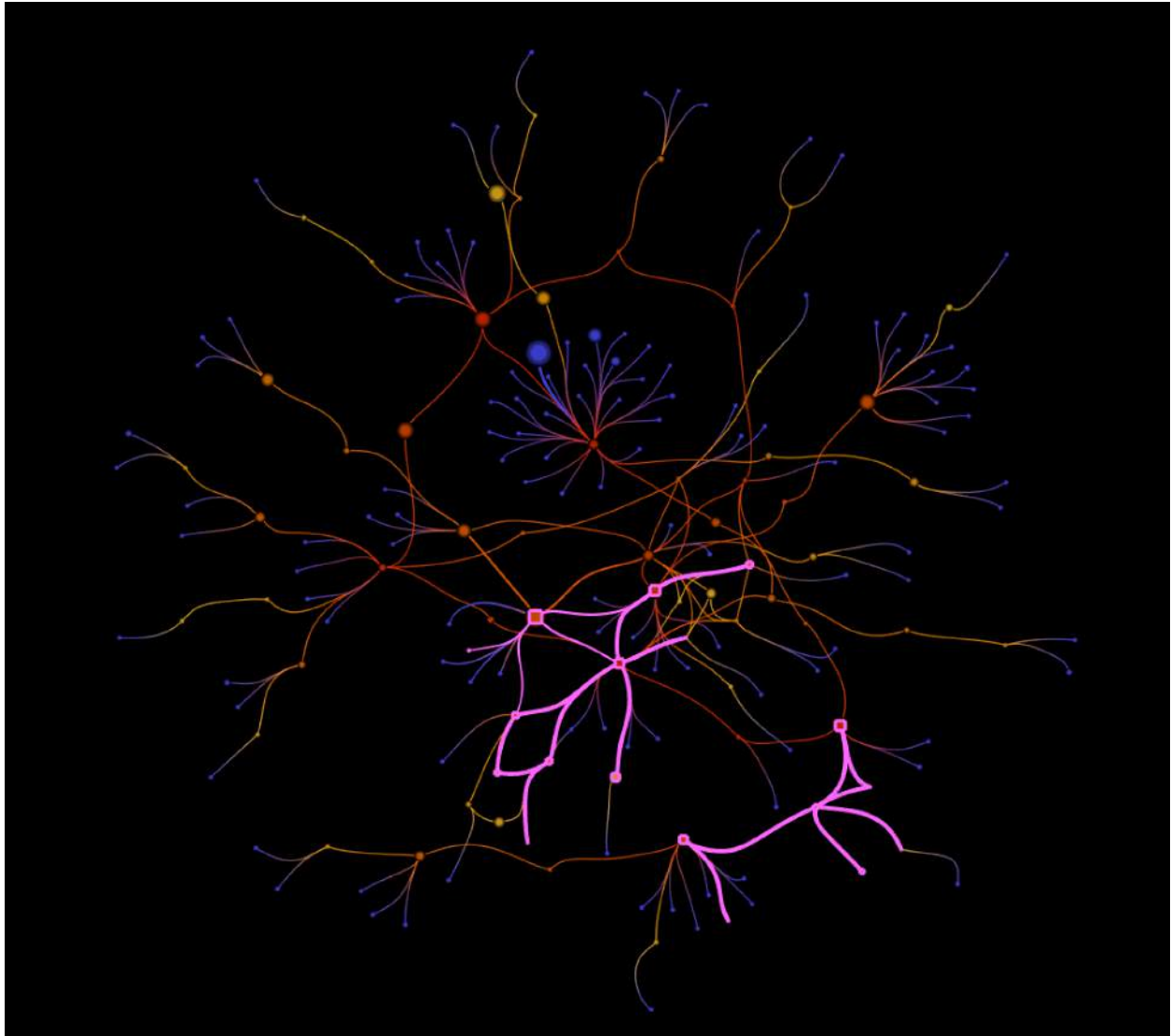


Figure 7 - The 21 core individuals involved in reciprocal conversations.

One subgroup is focused on Edgeryders, producing altogether 490 tweets, and the other group focused on travel companies and tourism, producing 171 tweets.

A quick search on how it distributes over time tells us that those two conversations happen at two different timings, the Edgeryders community happens mostly in the early period and the second group is focused later.

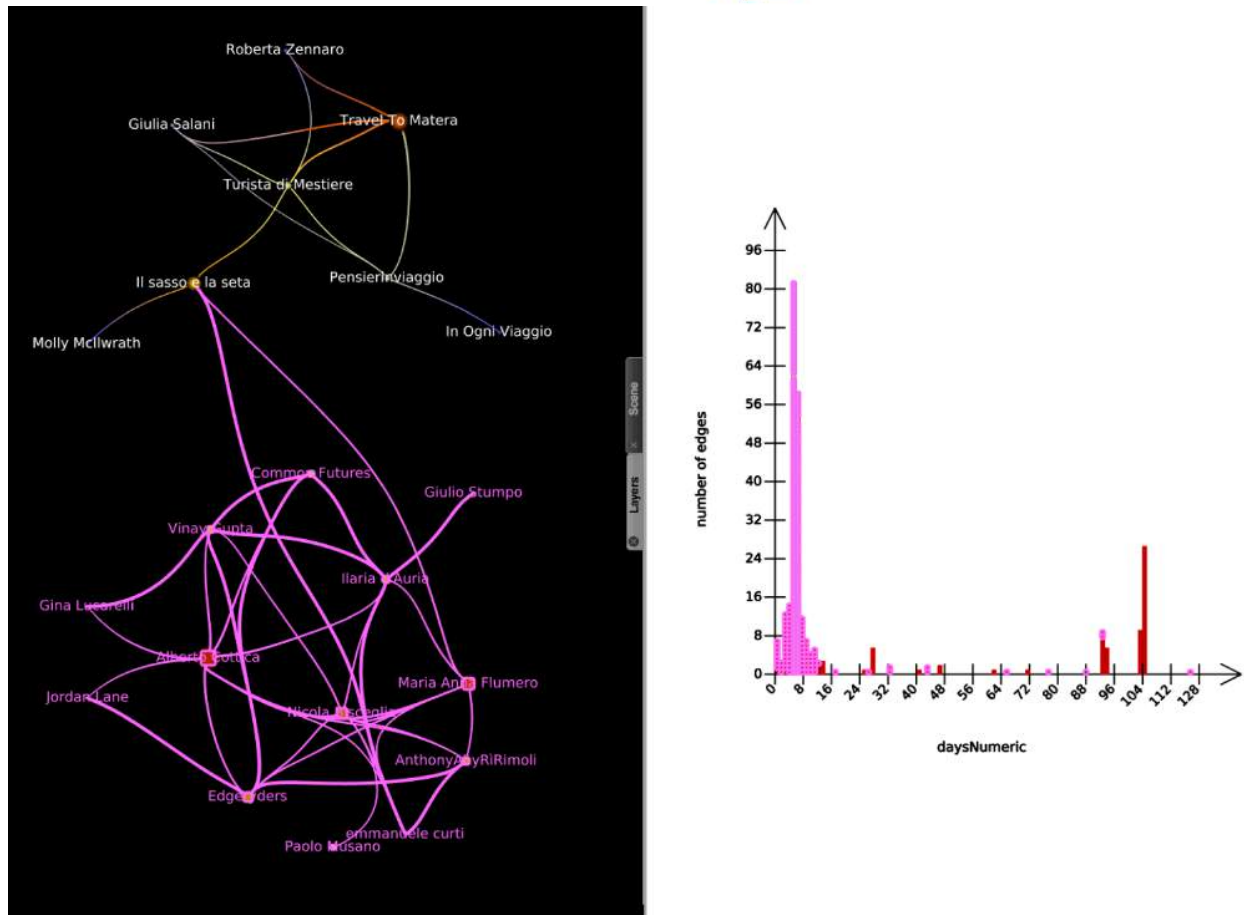


Figure 8 - Distribution of the tweets related to Edgeryders among time (selection in pink).

Going back to notion of community

Now, we may wonder how the two rather central communities have been brought close together, and how do they interact together? We can step back and look at how the different twitterers do at mentioning each other.

1500 people actually mention each other in a connected way. Among these 1500 people, only 100 form a core of reciprocal mentions between each other's, *i.e.* people acknowledging each other.

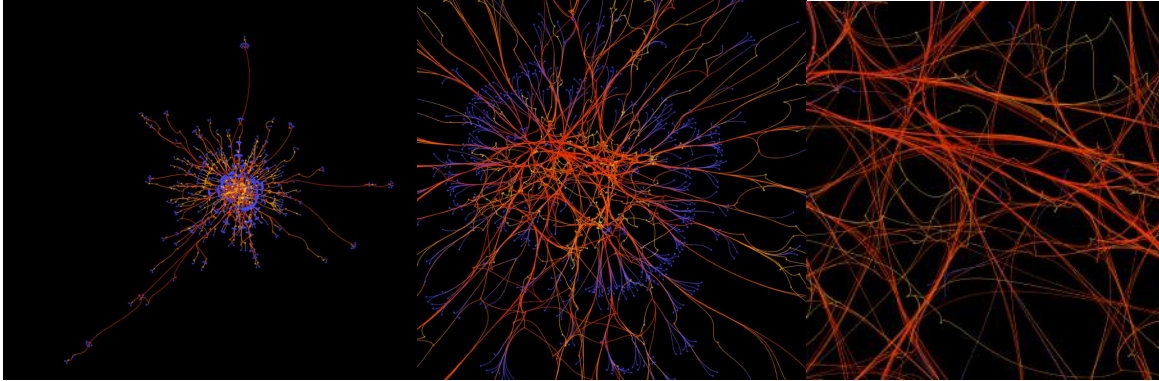


Figure 9 - Two consecutive mesmerizing zooms on 1500 people mentioning one another in an intricate conversation.

Now it is interesting to see how the core two communities we have previously identified collaborate with one another. We can observe that the two sub-communities do not acknowledge reciprocally each other.

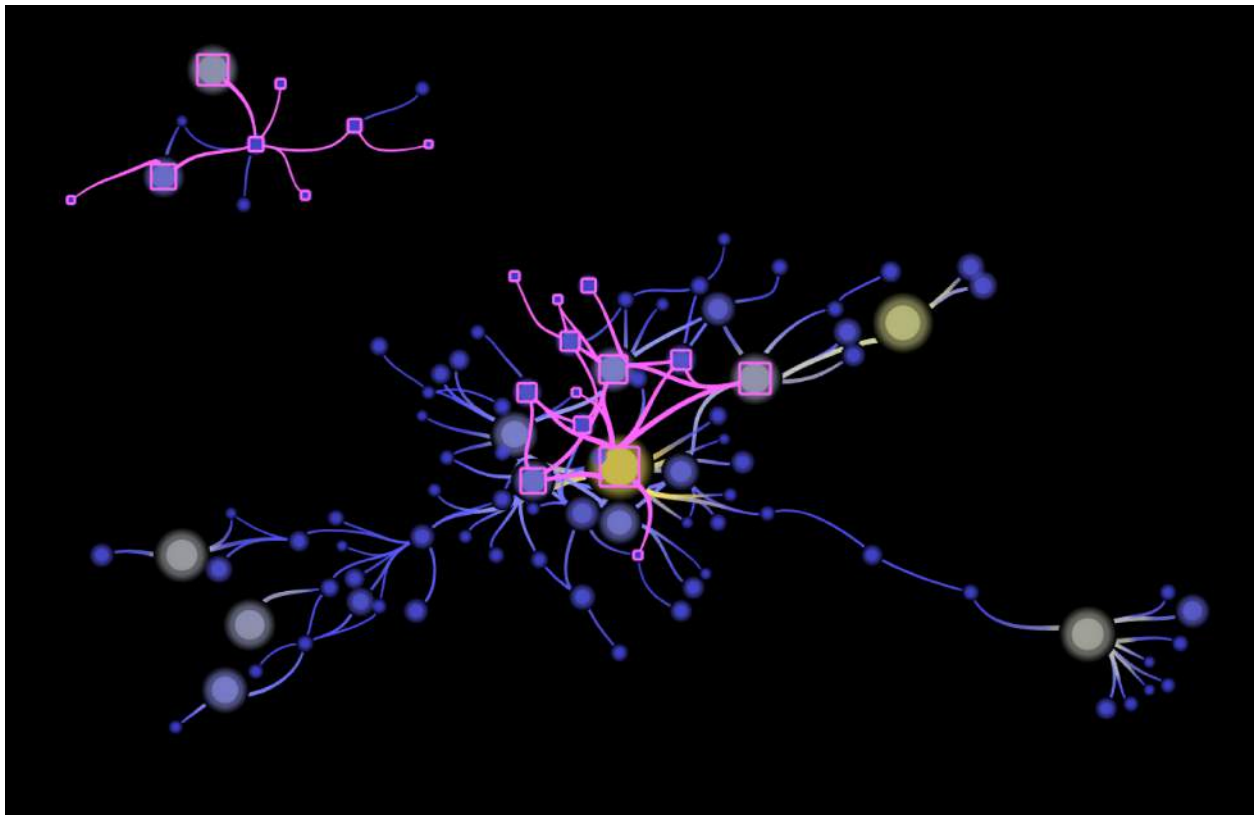


Figure 10 - Highlight (in pink) of the two communities in the connected components of the “Mention” network

However we can also identify the bridge elements between both communities by stepping back in the bigger “mention” component.

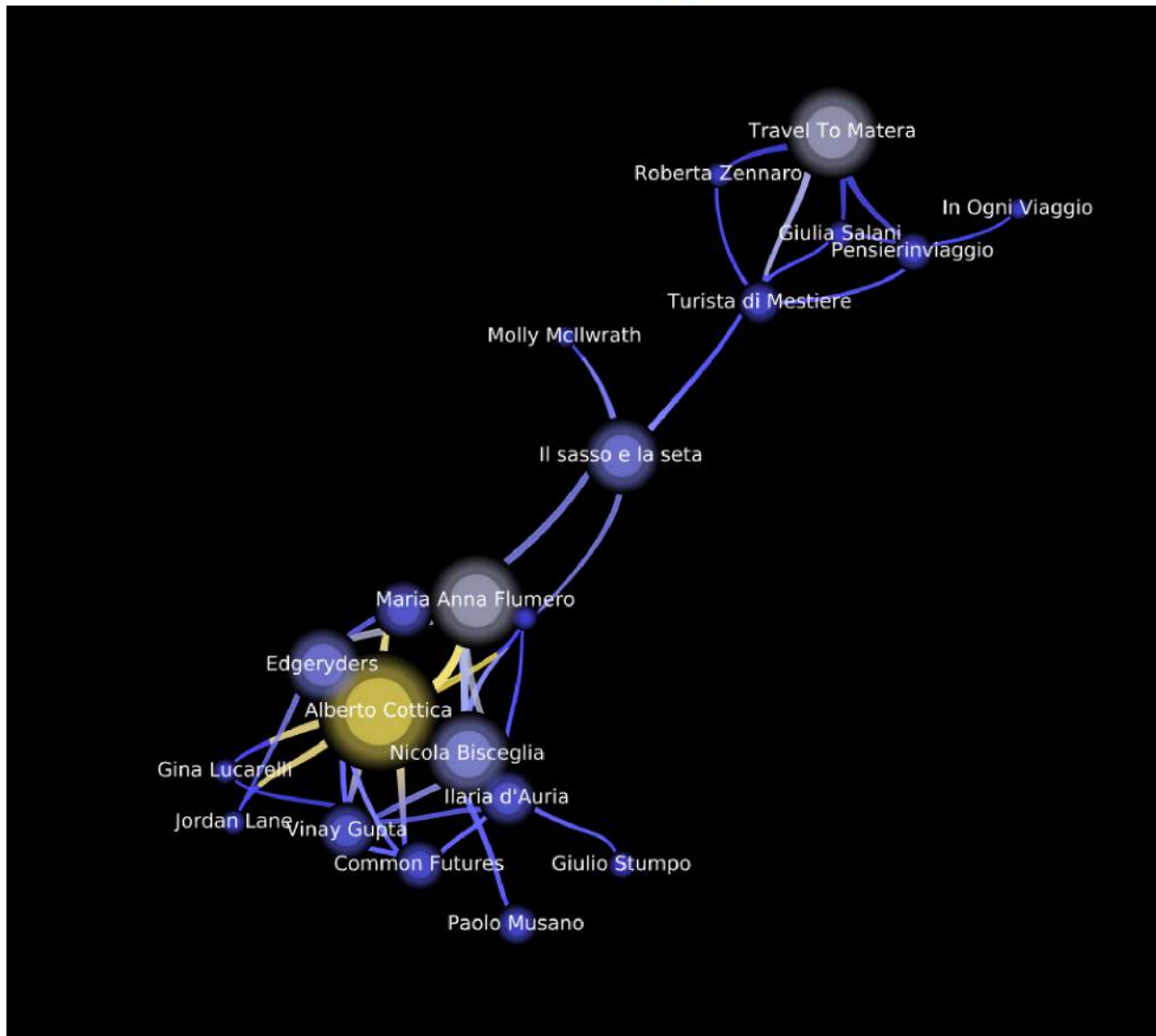


Figure 11 - How the two separated components are connected (extract from the “mention” connected component)

Conversing

Now that we have focused on how *people* connect (or not) together in the community. We can focus on *what* they interact about.

The idea is to compare the semantic space in which the individual exchange when they discuss together. Is it different from how they mention each other? and how? and what brings them together.

To do so, we have built a different network, it's actually a network in which links materialize the hashtags exchanged between two users. It has the same flat topology as the previous network of people, but it is rich of the semantics that people use when they converse.

Using this model, we can capture the schemes of conversations and see how hashtags bring the two sub-communities together.

14

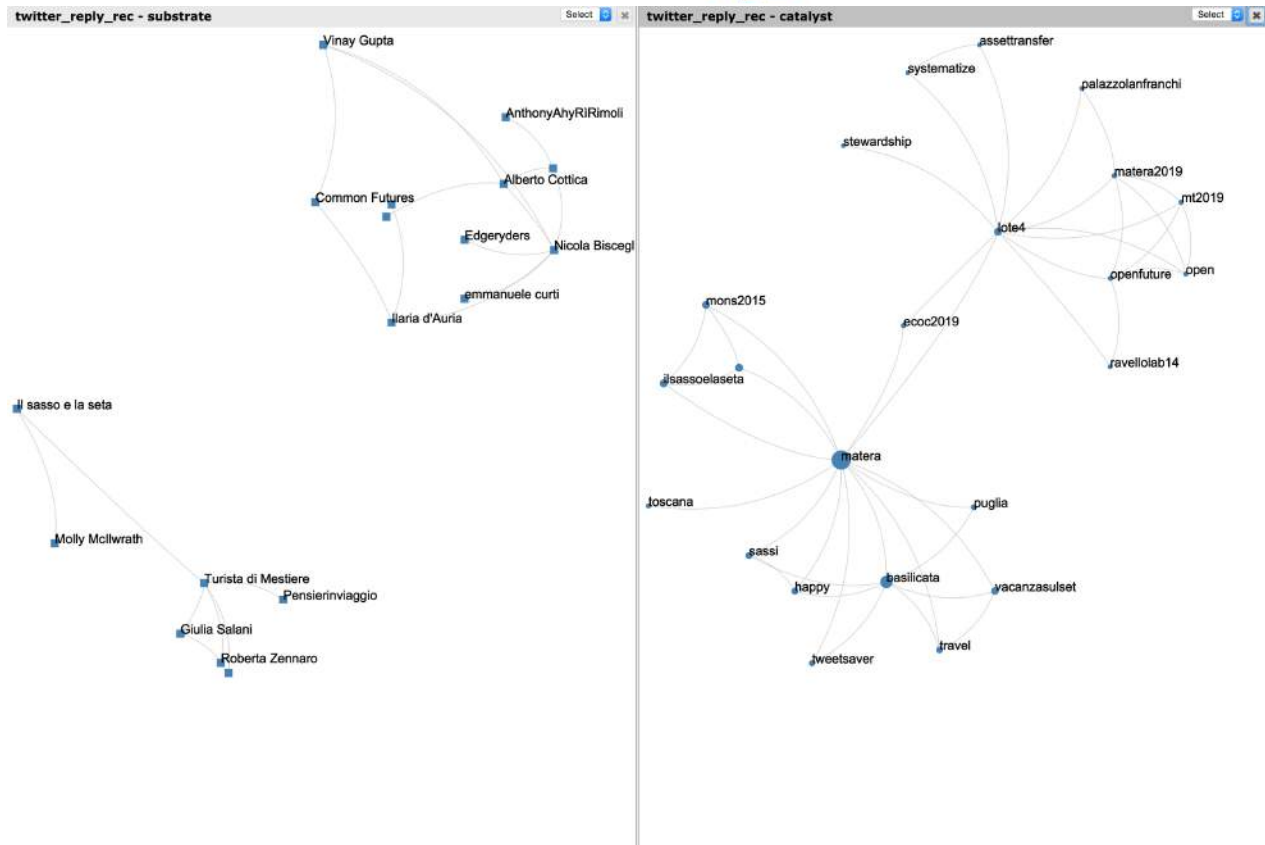


Figure 13 - Communities of people replying to each other. We can see the two disconnected group of people on the left, aligned with their corresponding hashtags on the right (people on top, converse about topics on top). The frontier is very clear. People on top discuss of *lote4*, and people on bottom of *Matera* in terms of travel.

However by looking at how these two groups are mentioning their members together, the frontiere still exists of course, but it is by far more blurry, and *Matera2019* seems to be an interesting gathering point between these communities.

By the way, we have analysed that, even if *#matera* is the heaviest most occurring hashtags in our dataset, *#lote4* is by far more often more co-occurring with many other hashtags putting *#lote4* as the most influent hashtag in terms of group cohesion. In other words, the sub-community centered on Edgeryders is more cohesive because people discuss more often about the same focused topics. Could it be an interesting side-effect of an effective moderation? or could it be that semantic cohesion makes a community really a community?

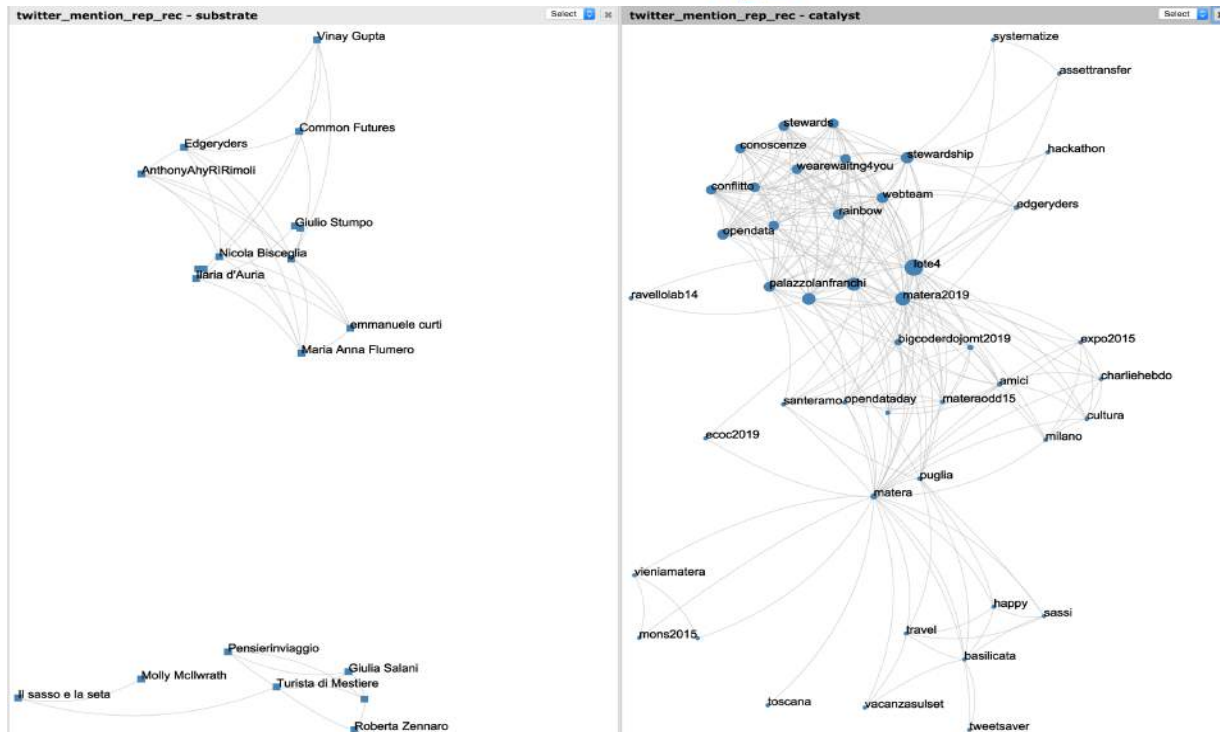


Figure 13 - Communities of people mentioning each other. We can see the two disconnected group of people on the left, aligned with their corresponding hashtags on the right (people on top, converse about topics on top). The frontier is more fuzzy and a core set of topics (on top) relate very much to #lote4 whereas #matera is very diffuse (on bottom)

Conclusion

Of course, this analysis is only the result of a 2-day workshop, and we would wish to push it further to have a complete understanding of the structure of the community. But it is a nice example how we can dig effective elements of discussions, topics of interests, and central people, at the heart of a very noisy spread Twitter conversation.

Also, considering the context of the event it would make sense to filter the data after, the end of November. Later tweets refer to Matera or to the unMonastery, but at this point the two terms no longer refer to a unity. There also could be many other questions, focused on data properties but also wider openings: How would we compared and define that with another community, such as Imagination4People's mailing list / communities? What is the specific role of these [put-your-list-here] individuals? What were their focus of attention during [put-your-timeframe-here]? etc. etc.

Now that we know all these tools and metrics are available, the most challenging task would be to set a (*Cartesian*) *method* to systematize and build integrated tools for the analysis of such a Twitter community. Here are some examples ways of applying this to the real world:

- Showing the graphs on events to demonstrate which were the main topics tweeted during the day and which people were most active - this implies that the graphs are developed rather quickly
- You can identify people who are especially interested in specific topics, at least if mentioning the topic in # is a good indicator for this
- Really cool would be if this could be used instantly on a twitter wall during events
- By reducing the overall number of members of a community to the very few (18 in this case) really active ones, you know who you would need to approach for future events, discussions etc. > find the champions
- As it comes down to few topic # on our graph, we might be able to encounter new emerging topics or other unexpected stuff

Development

All the data is made open from the Masters of Networks 3 website. This analysis has been processed during the event, Tulip 4.6 (tulip.labri.fr) has been used to process the CSV data, and build the initial networks, detangler (detangler.labri.fr:31497) has been used for the paired semantic analysis, and a little bit of d3 (d3js.org) to process the time series (just a bit after the event).

Reportback

Retweets: not a vector of community, a weak link

~8000 people / ~20.000 Tweets

Matera is both a hashtag and a physical community, the name of a town

Reply : is the strongest (is the strongest social connection)

Mention : is in the middle (mutual respect)

Retweet is the weakest (not necessarily an endorsement)

8000 people have tweeted or retweeted something hashtagged #LOTE4

4500 people total attached to #LOTE4

2000 people have mentioned/replied each other (within the whole Twitter space) in a connected "community" (component)

1500 people have built a network of mentions (the core community)

About 100 people have reciprocal mentions

200 people have replied to another

18 people have engaged in long conversations (more than just 1 reply)

These 18 people form 2 separate subgroups

2 different anchors for the subgroups - what they discuss and what they talk about

Our visualization shows how these 2 different subgroups interrelate.

one is #Matera. The other is #Lote4

The common connector between these two groups is actually #Matera, not #Lote4

No clear division between the topics of conversation, but #LOTE is clearly the focus of discussion for 1 group of people. But #Matera is more the bridge between people. One

Next Step: we only have 18 people at the center of the conversation, it would be great to be able to do a qualitative content analysis of what the tweets were actually talking about.

LOTE4 is the most common term, but Matera is the co-occurrent (the bridging term)

Open questions:

- How quickly can you get the data you need from Twitter and make the graphs out of it?
- It could be, that the connecting topic of #Matera has been used by both communities at different times. For ensuring that this topic was really an interaction point, we need to consider the time line as well.

Ways of applying this to the real world:

- Showing the graphs on events to demonstrate which were the main topics tweeted during the day and which people were most active (e.g. next CAPS meeting) - this implies that the graphs are developed rather quickly
- You can identify people who are especially interested in specific topics, at least if mentioning the topic in # is a good indicator for this
- Really cool would be if this could be used instantly on a twitter wall during events
- By reducing the overall number of members of a community to the very few (18 in this case) really active ones, you know who you would need to approach for future events, discussions etc. > find the champions
- As it comes down to few topic # on our graph, we might be able to encounter new emerging topics or other unexpected stuff

Track 1: What makes a community a community?

This builds on Alberto's question: is a Twitter conversation around a hashtag "a community"?

I suppose a way to investigate this would be to compare the interaction networks generated by relatively close-knit communities (like Imagination4People's mailing lists, or Edgeryders) to those generated by hashtags.

What do you see? How could you describe why these two are different, if they are?

For example, you could find that there are lots of isolated components in a Twitter hashtag stream, with people not really calling out to each other, whereas "tight" online community give rise to a giant component that most of the nodes are connected to.

Or not. Anyway, this is an interesting question, and people could easily form subgroups investigating specifics: for example, comparing "loosely knit" and "tightly knit" communities from the point of view of network modularity, or centralization, or clustering.

Remember: the presence of community managers means that we have an independent qualitative assessment on the tightness of each community.

We've uploaded the gexf files to this google drive folder.

<https://drive.google.com/folderview?id=0B-sizRV1qoRXfk5vTINyVTAwZmVfY25PdFJpM3dxSFFwUVILVUpETko3QjhHamZfNEk5Ukk&usp=sharing>

Questions:

Goals (twitter):

- 1- Connect them to one another
- 2- Have them actually attend event

What is a community manager?

Innovation ipa:

People have different roles

General:

Some have specific moderator personas, others interact under their own name, like any other user.

Definition of a community

People talking to one another

Sense of belonging/endeavor

inner sense of belonging

sharing of interest

some commonalities

twitter hashtag

exchange in the community (both ways)

notion of time (punctual or lasts)

somebody who's not in the community (the rest of the world)

people who do more than what they need to/have to

behavior, follow the same people / sign petition

actually meet / community of practice

transversal to organization

groups = set of people

gather around specific goal

sense of belonging

share content

can be negative engagement

classes of communities

hierarchical/over time

can fuse or divide

3 classes that define a community

- interaction / awareness (manski) (hashtag mentions)
- discussion / sharing (username mention - one direction)
- practice / action (reciprocal interaction - conversation on twitter - back and forth)

1) Awareness / Interest

- a) Semantics (Hashtags, key words occurrence)
- b) Number of users who relates to the semantics
- c) Cross-Platform presence

2) Interaction

- a) Reciprocity
- b) Number of connections // user popularity
- c) Frequency of posting
- d) Impact of each post: endorsement, capacity to generate spin offs
- e) Relationship between the size and number of interaction / number of posts
- f) Independency from the community from the moderator

3) Action / Practice

- a) Occurrence of hashtags during events
- b) Cross-Platforms activities
- c) *Production of content (how to measure?)*
- d) *Spin off actions (how to measure?)*

Opening questions. How would we compared and define that with Imagination4People's mailing list or ImaginatoriPA communities? I suppose a way to investigate this would be to compare the interaction networks generated by relatively close-knit communities (like Imagination4People's mailing lists, or Edgeryders) to those generated by hashtags.

What do you see? How could you describe why these two are different, if they are?

Comparison, comparing "loosely knit" and "tightly knit" communities from the point of view of network modularity, or centralization, or clustering.

What could be an evidence of the role of community managers in these communities?

Shared folder for pictures:

<https://drive.google.com/folderview?id=0B3qYgN-9xWcofmRQWnUxYzBWbzZrSk5qdnc1Z1ZWWDdCQW1qN3dXeFc0NkxaRFNKTjZiZXc&usp=sharing>

Raw data (.csv.zip), Tulip data available (.tlpx), detangler (<http://detangler.labri.fr:31497>) data available (.json) here:

Now what's in it?

We've computed from the raw data the network of twitters who are mentioning and / or replying together (excluded RTs yet as not informative enough - the RT metadata has actually to be reconstructed to be relevant)

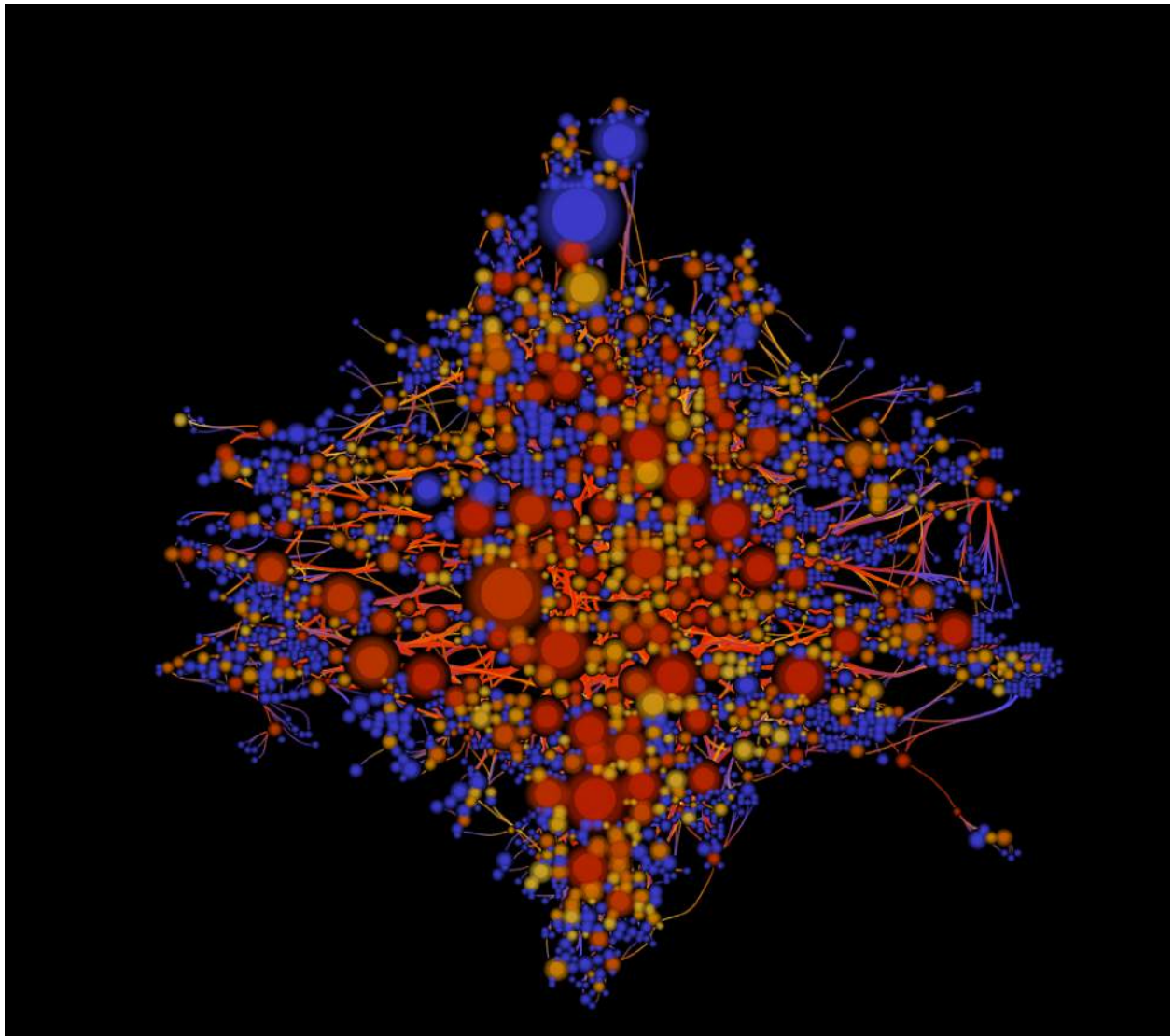
Visual encoding:

Nodes are authors

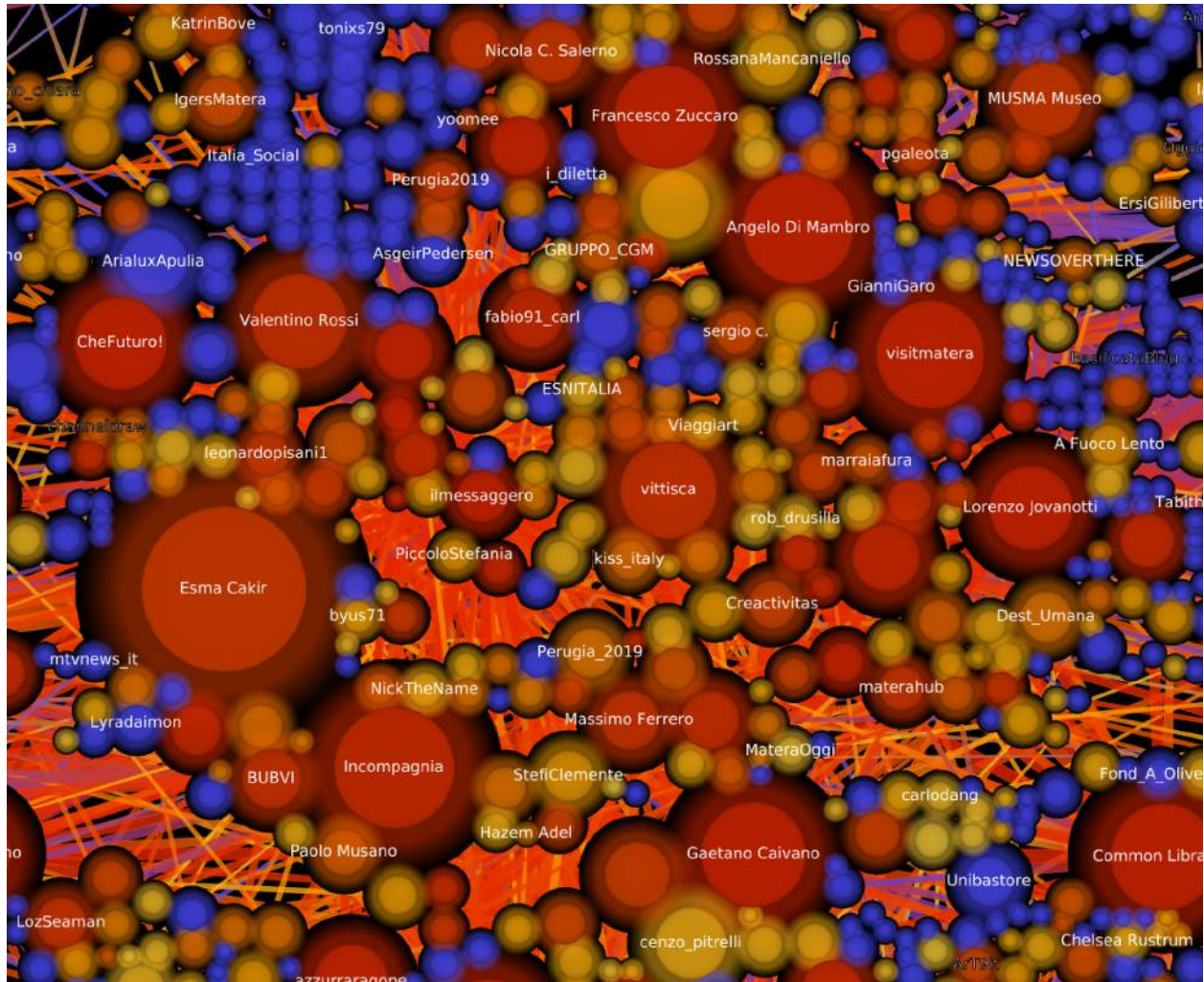
Node size is mapped to the number of twitts produced

Links are relationship between authors (reply to or mention)

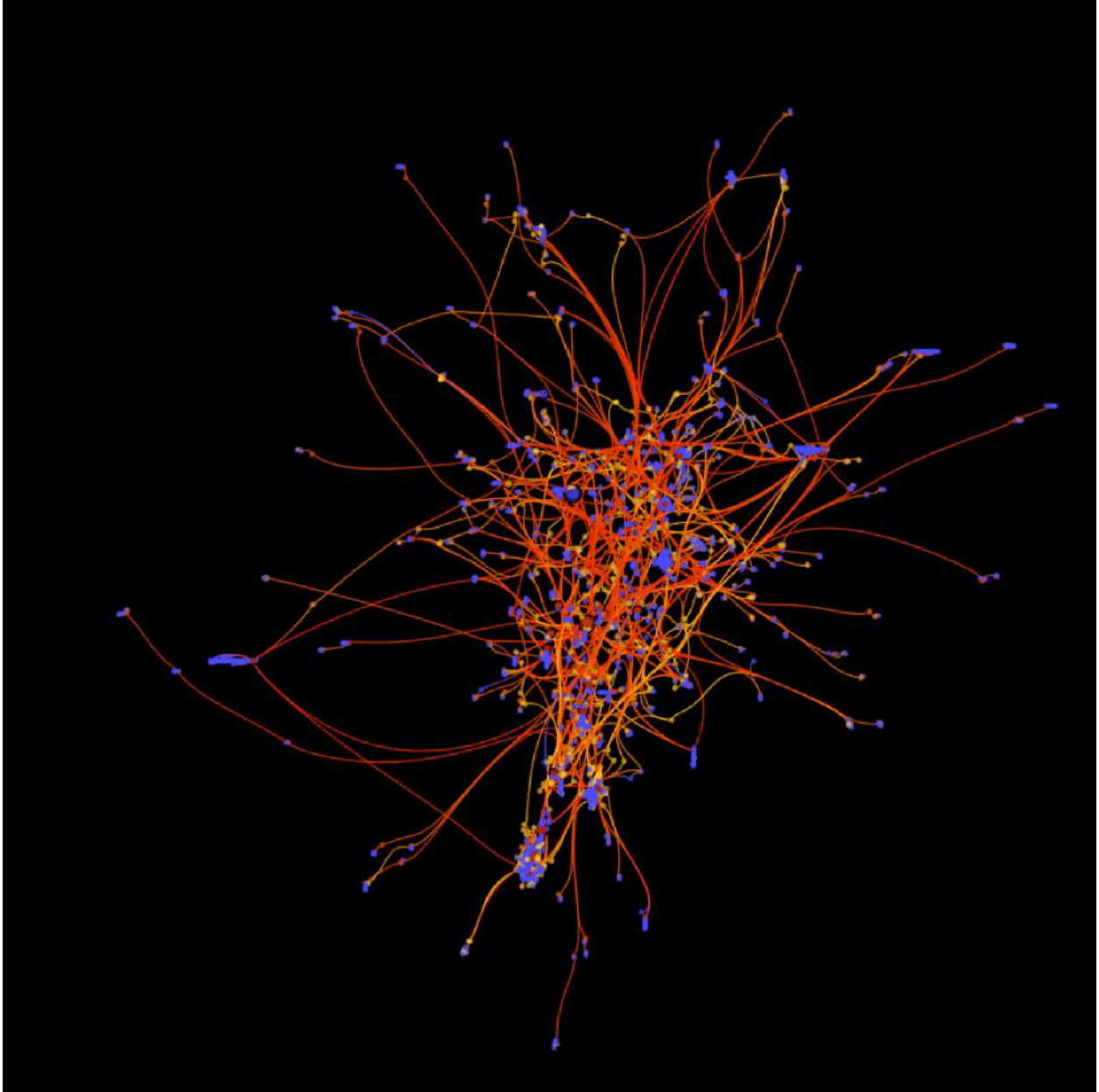
Node color is centrality (Edge color is interpolated from nodes)



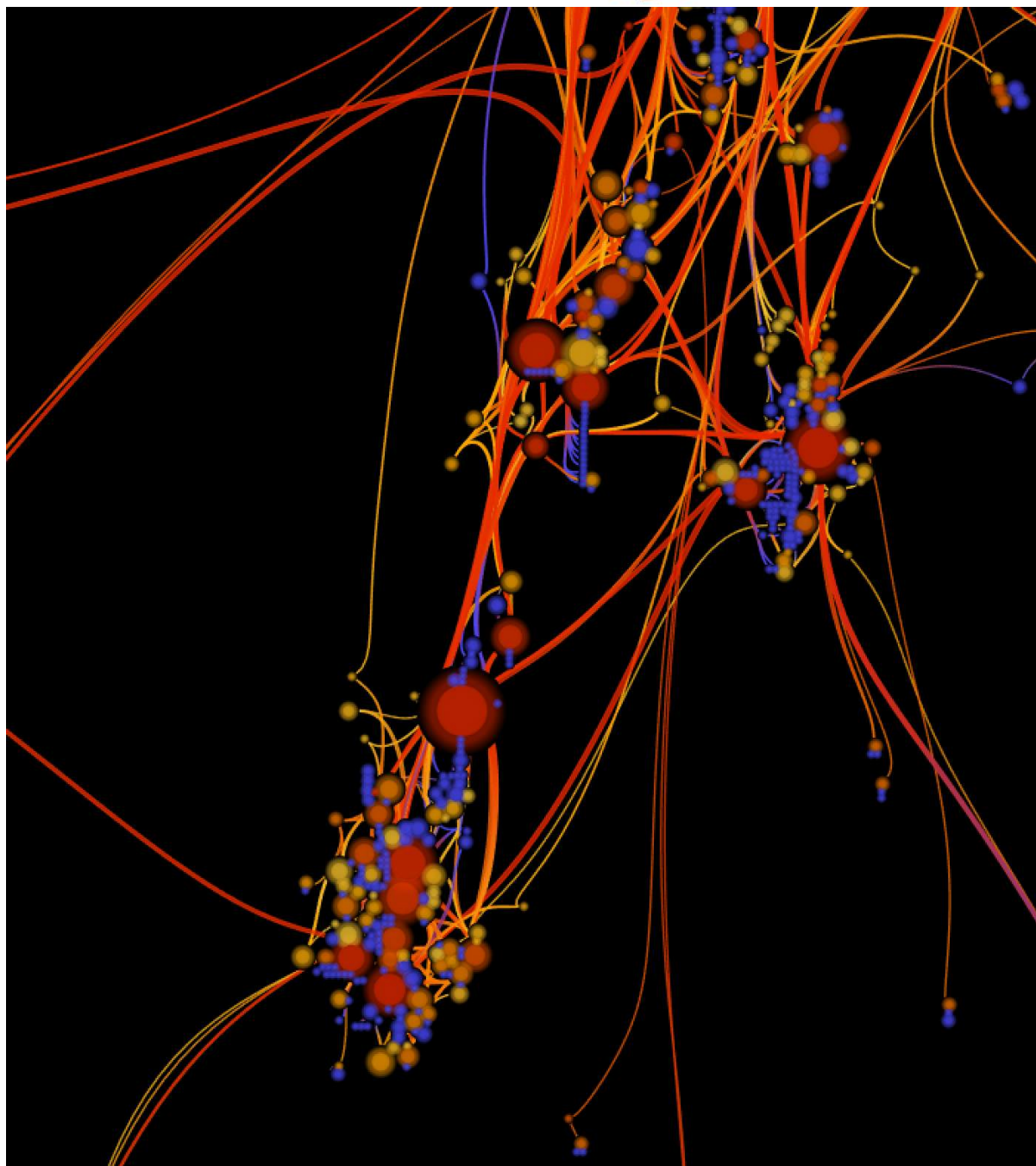
Obviously this graph is dense and hard to process, here are a few focus

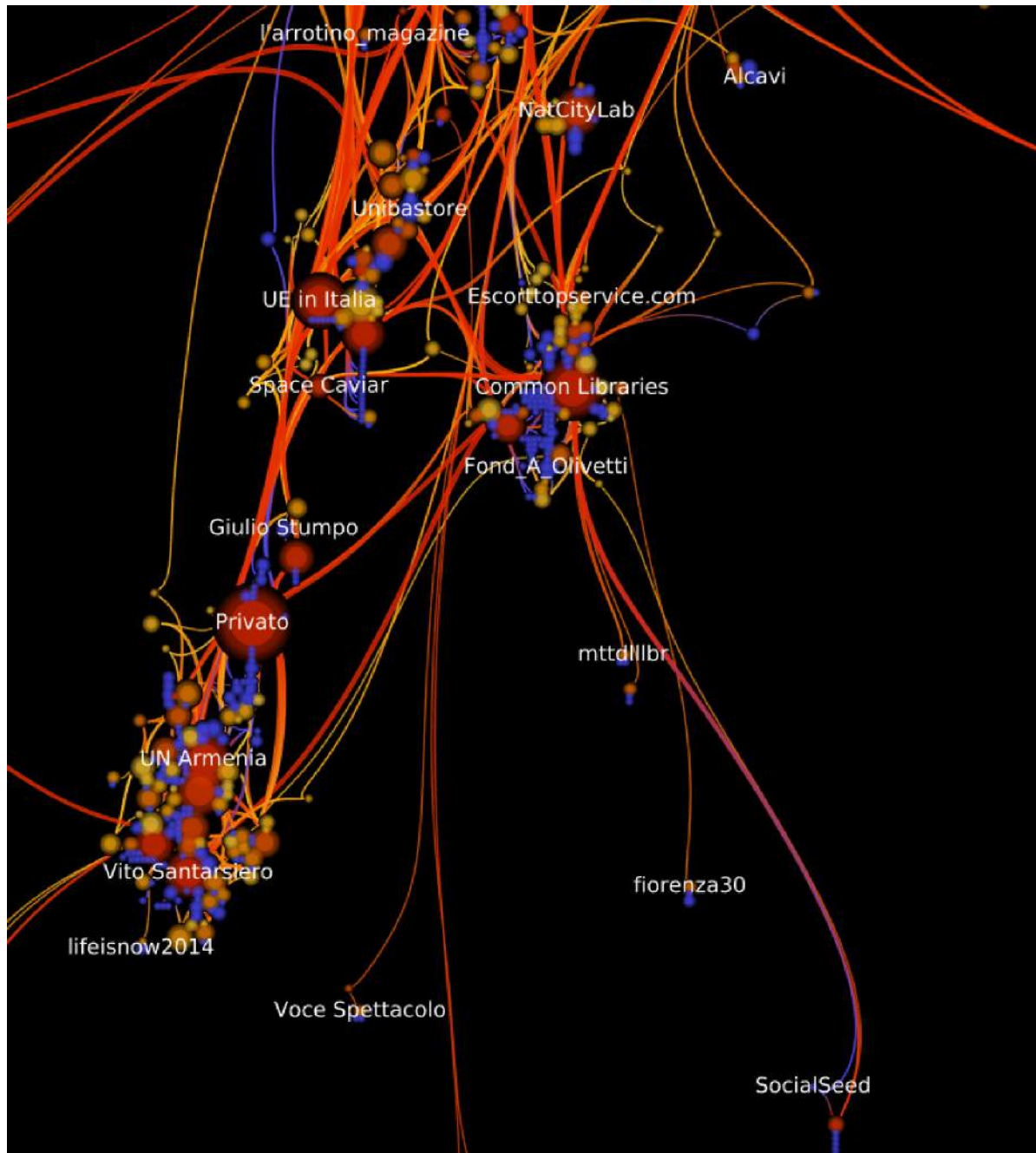


Now we can try to visually clusterize the network



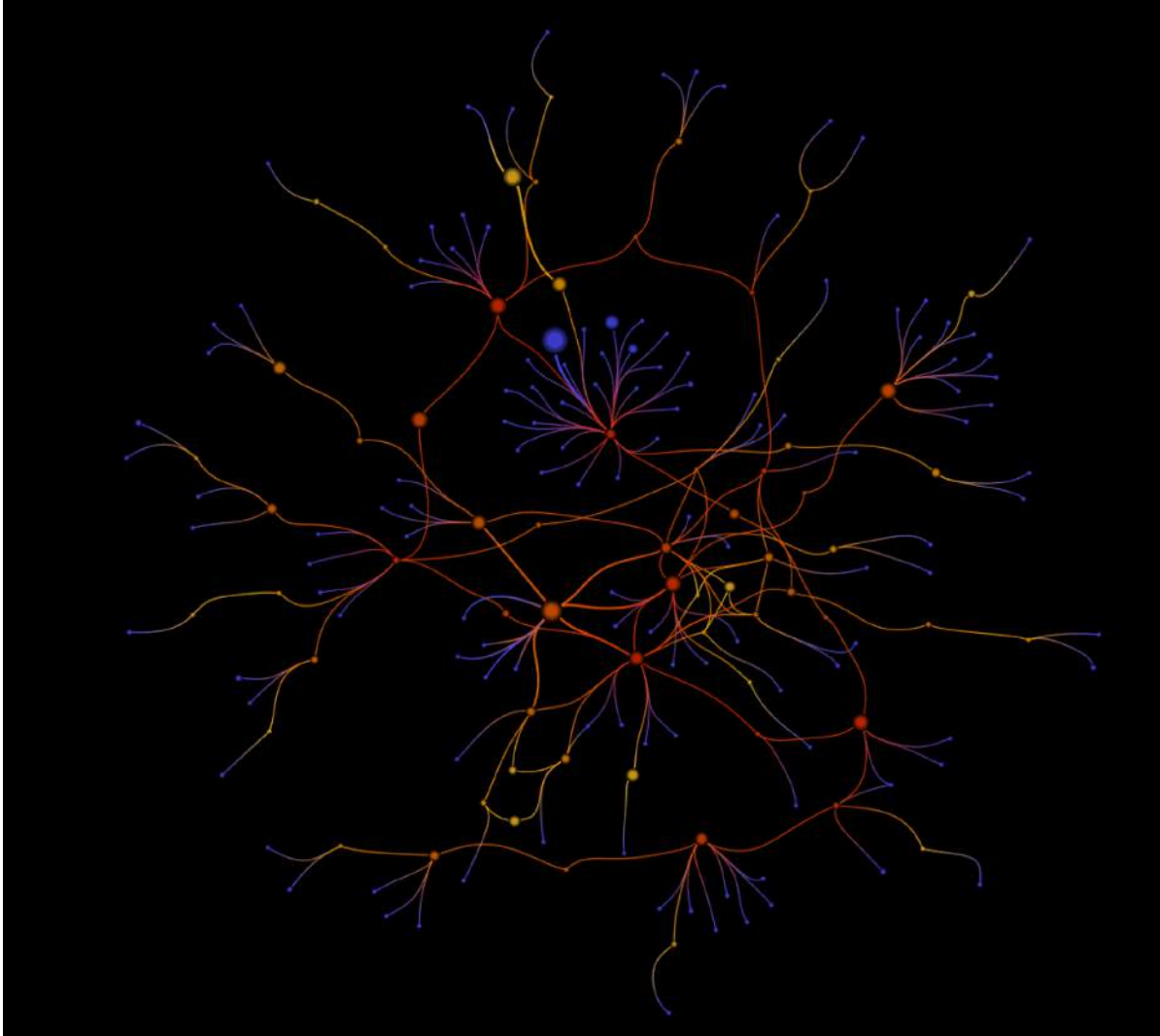
We can see some subgroups appearing

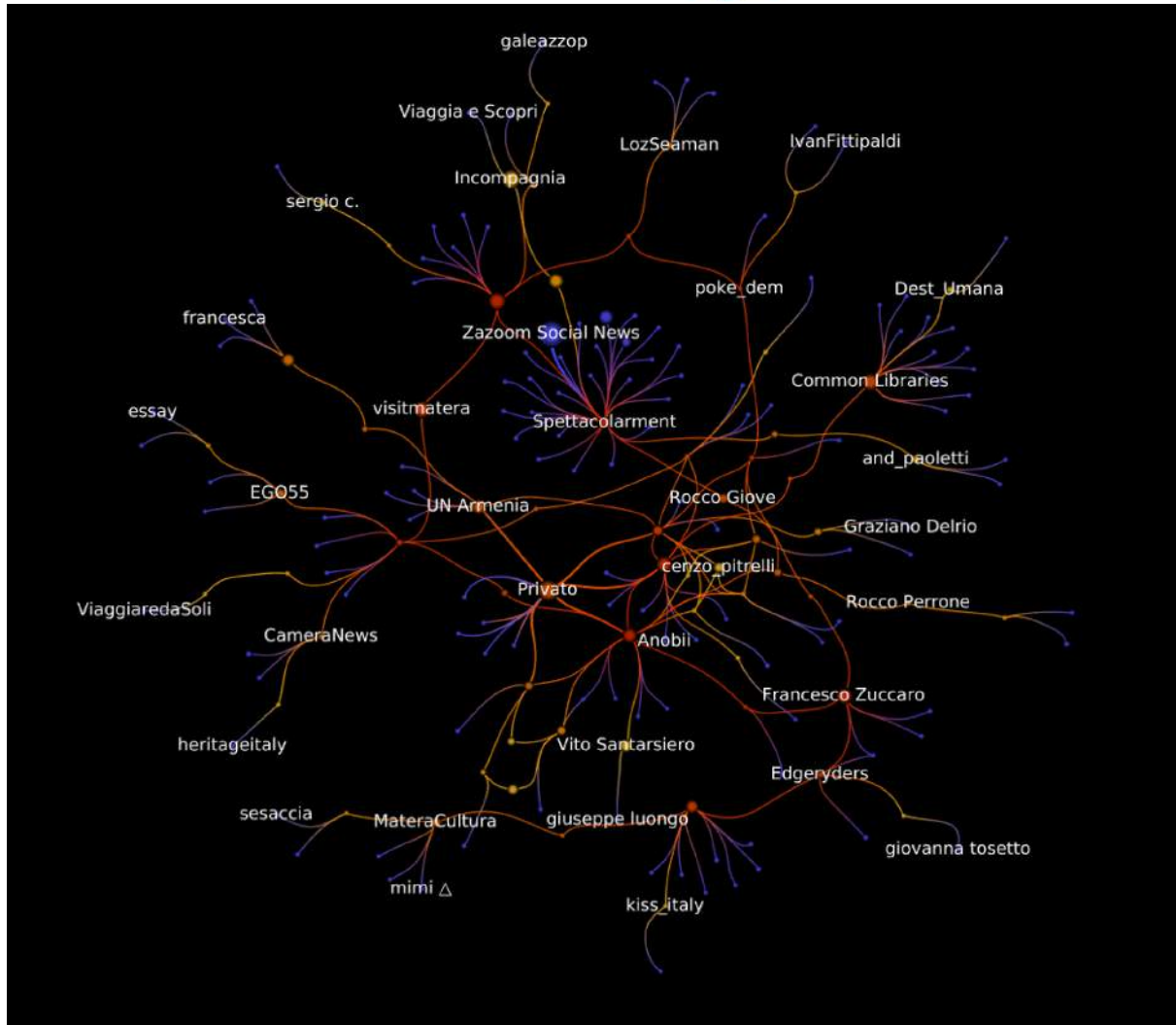




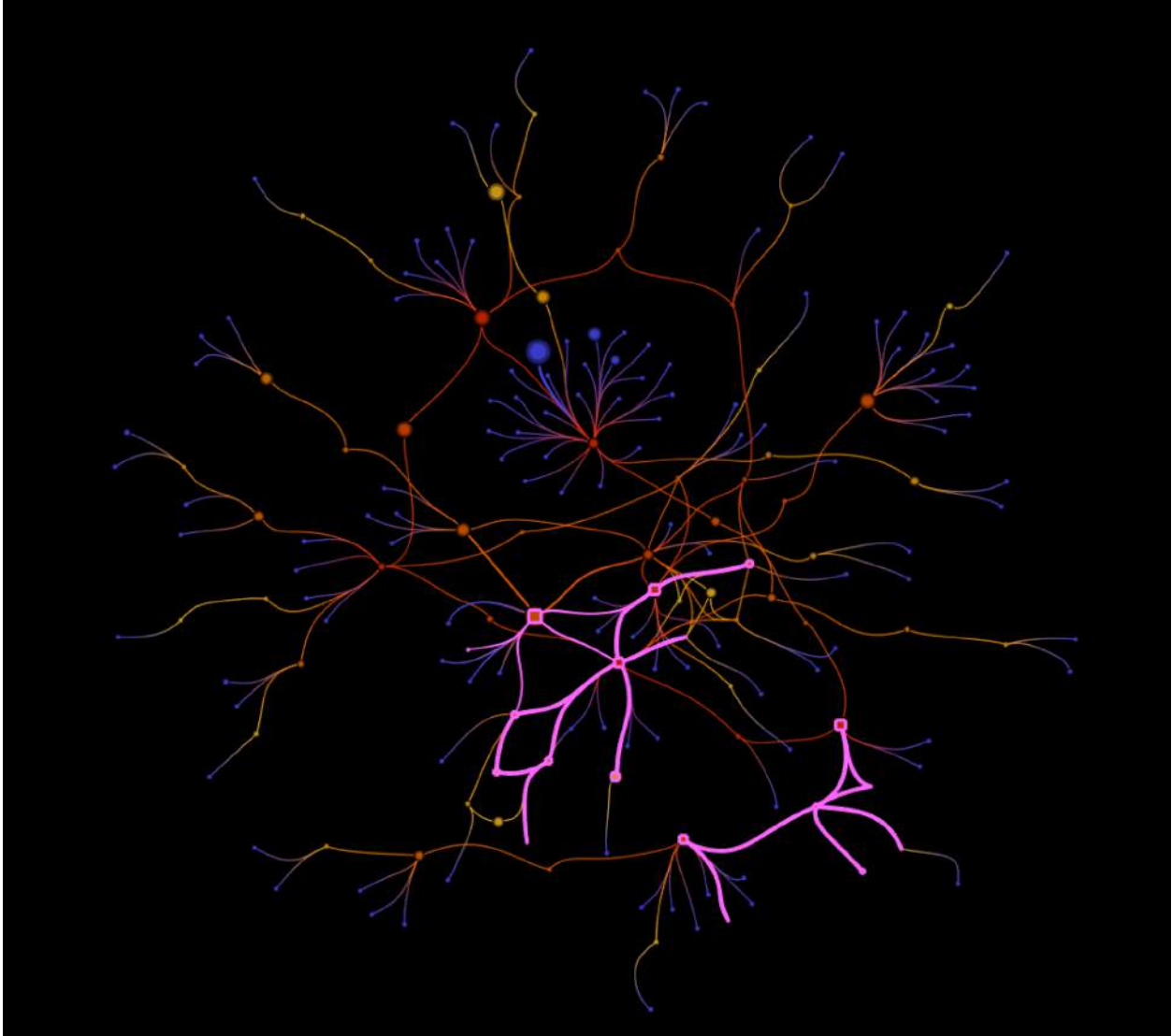
So we have two different types of relationships, we'll try to find cues of actual community behavior from them.

The first relationship, the strongest, is the "Reply to" relationship. If we filter out the other edges, we get (we just kept the biggest component for the sake of clarity) :

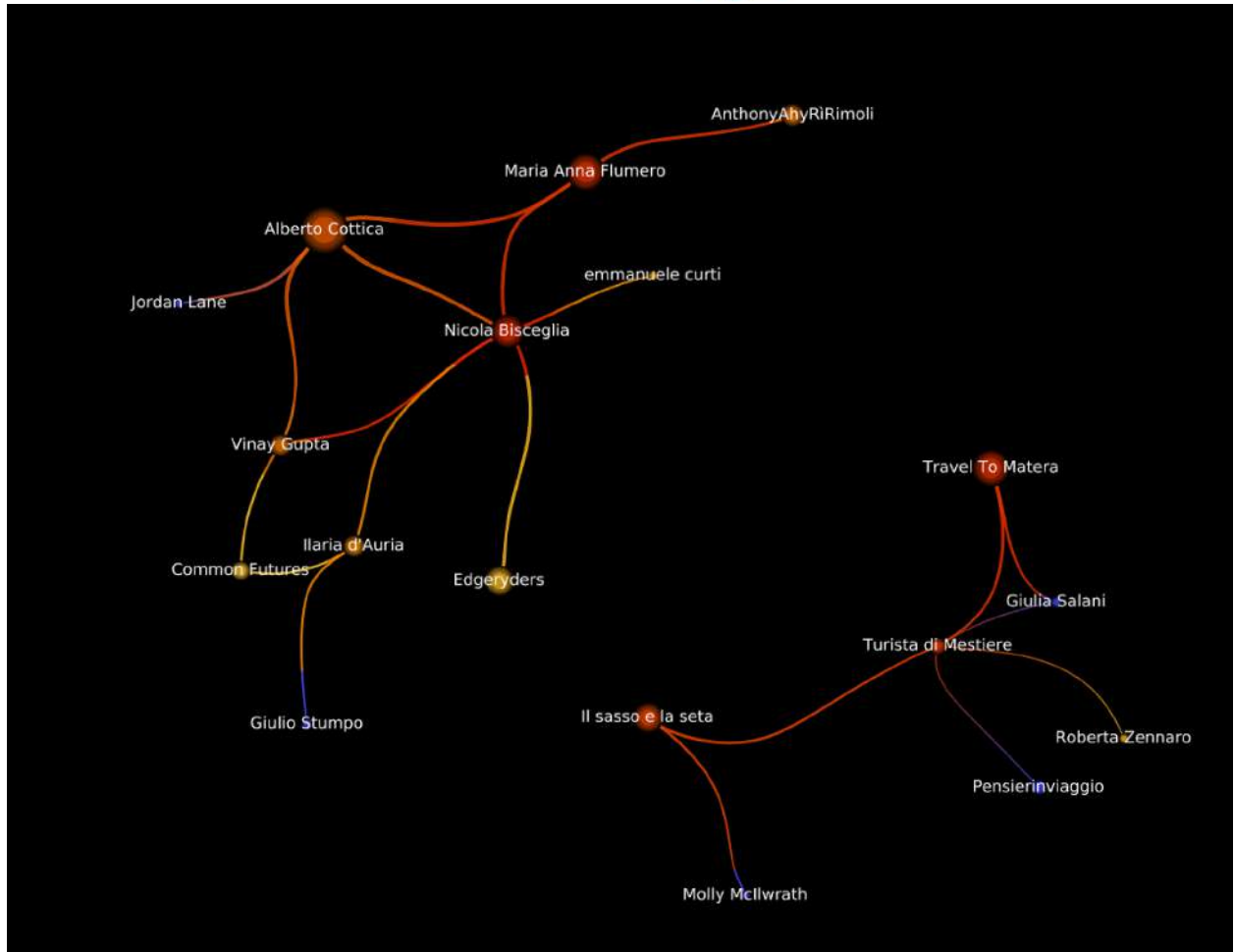




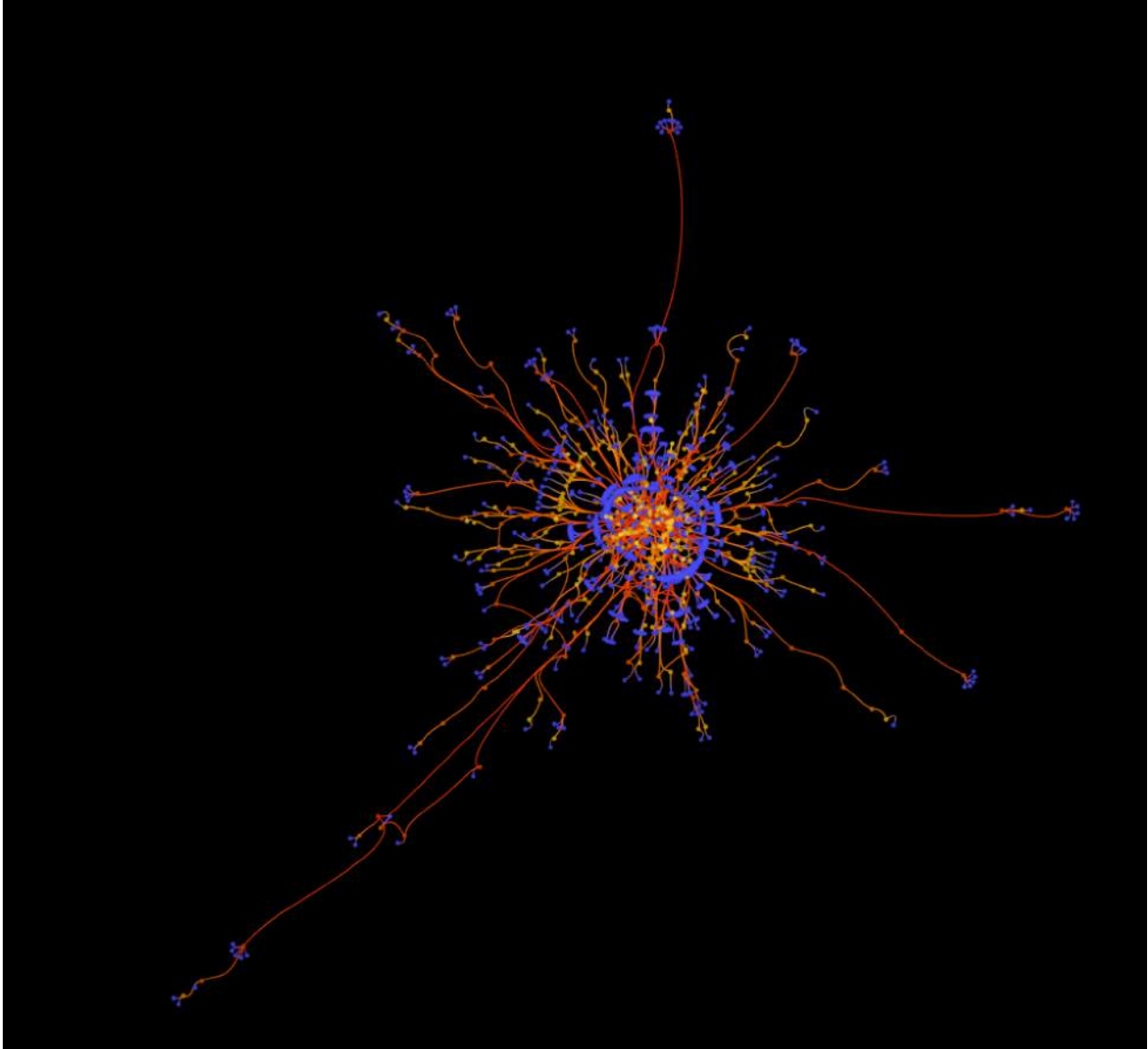
We expect the strongest relationships to be mutual, so we subset to reciprocal relationships



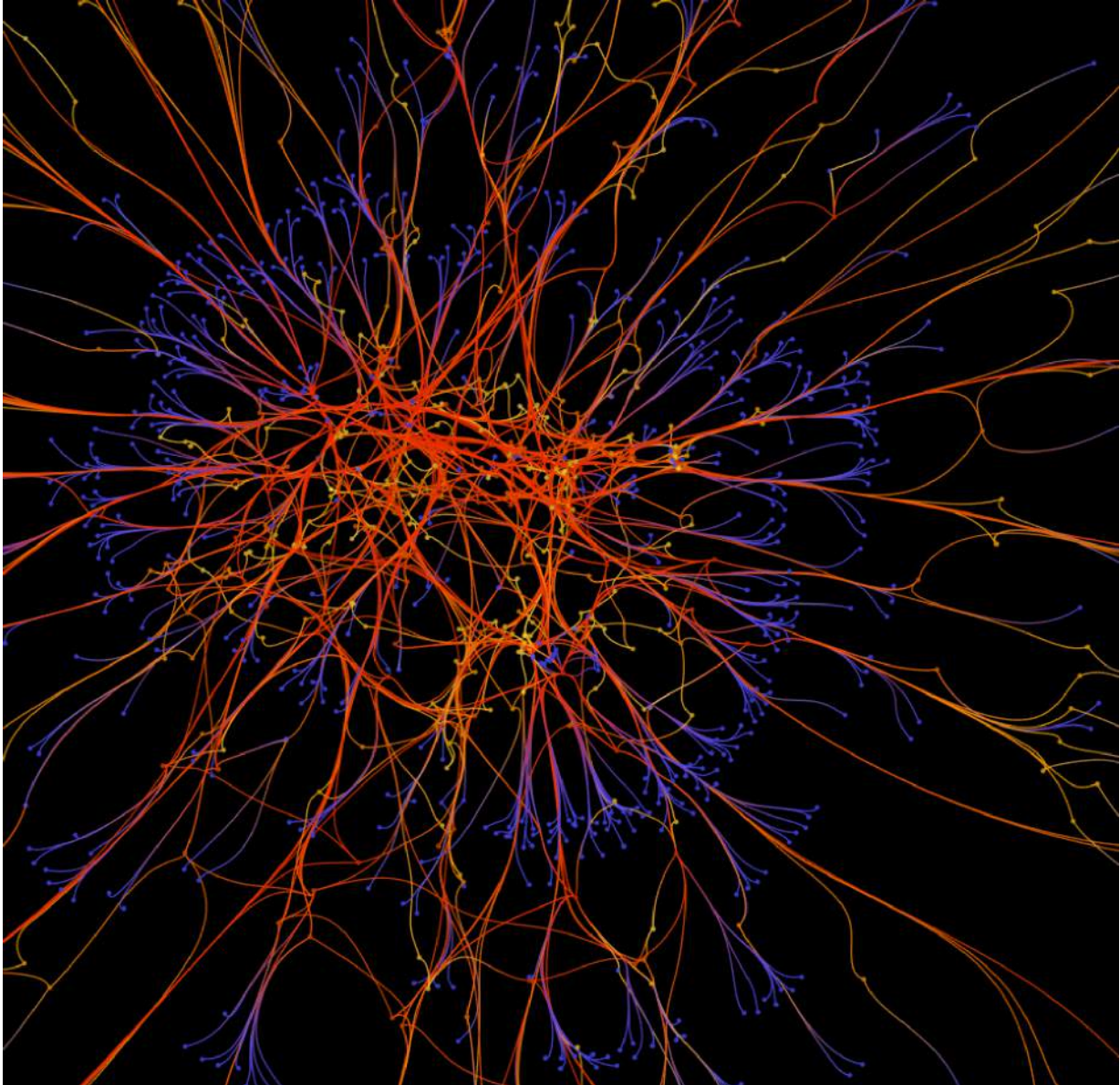
Interestingly it takes the shape of two disconnected subgroups:

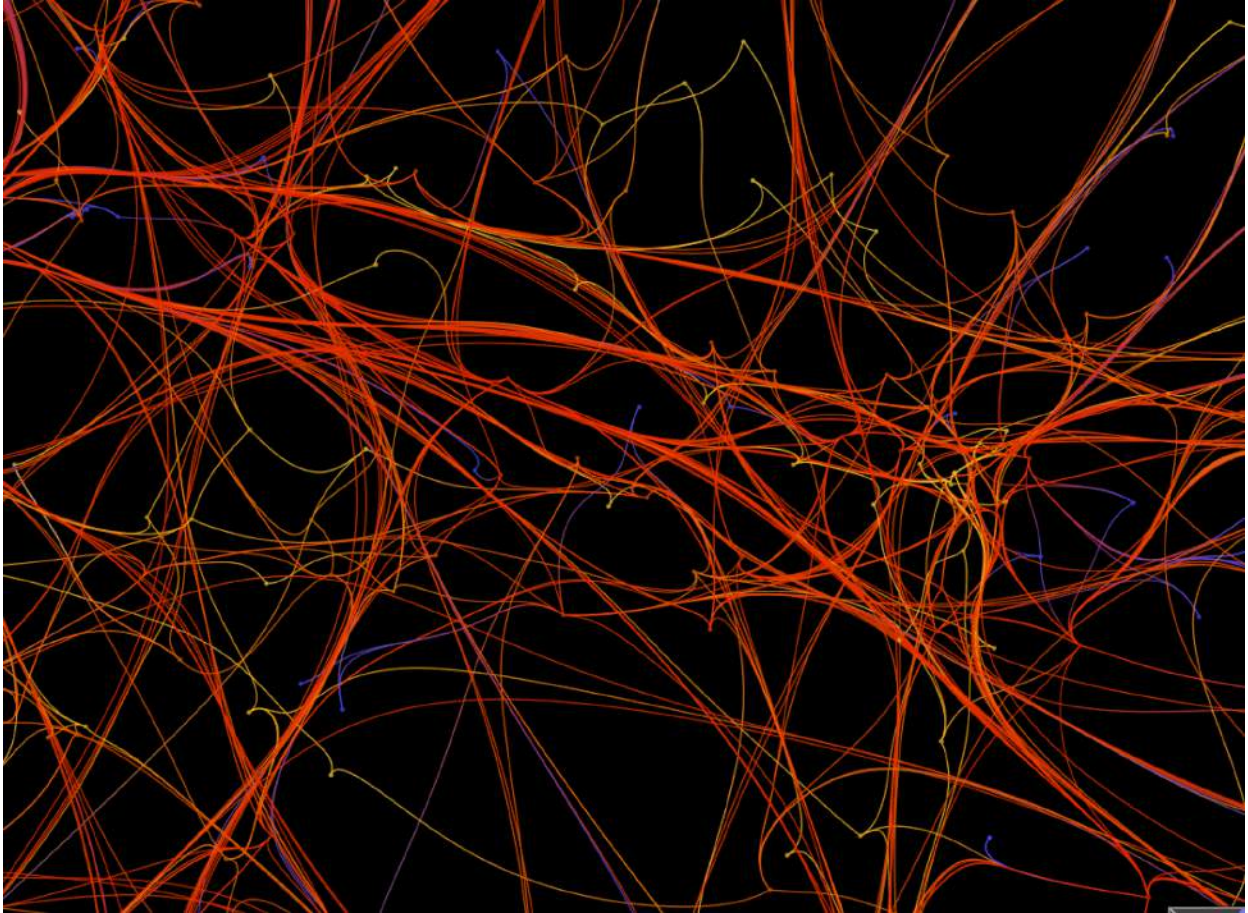


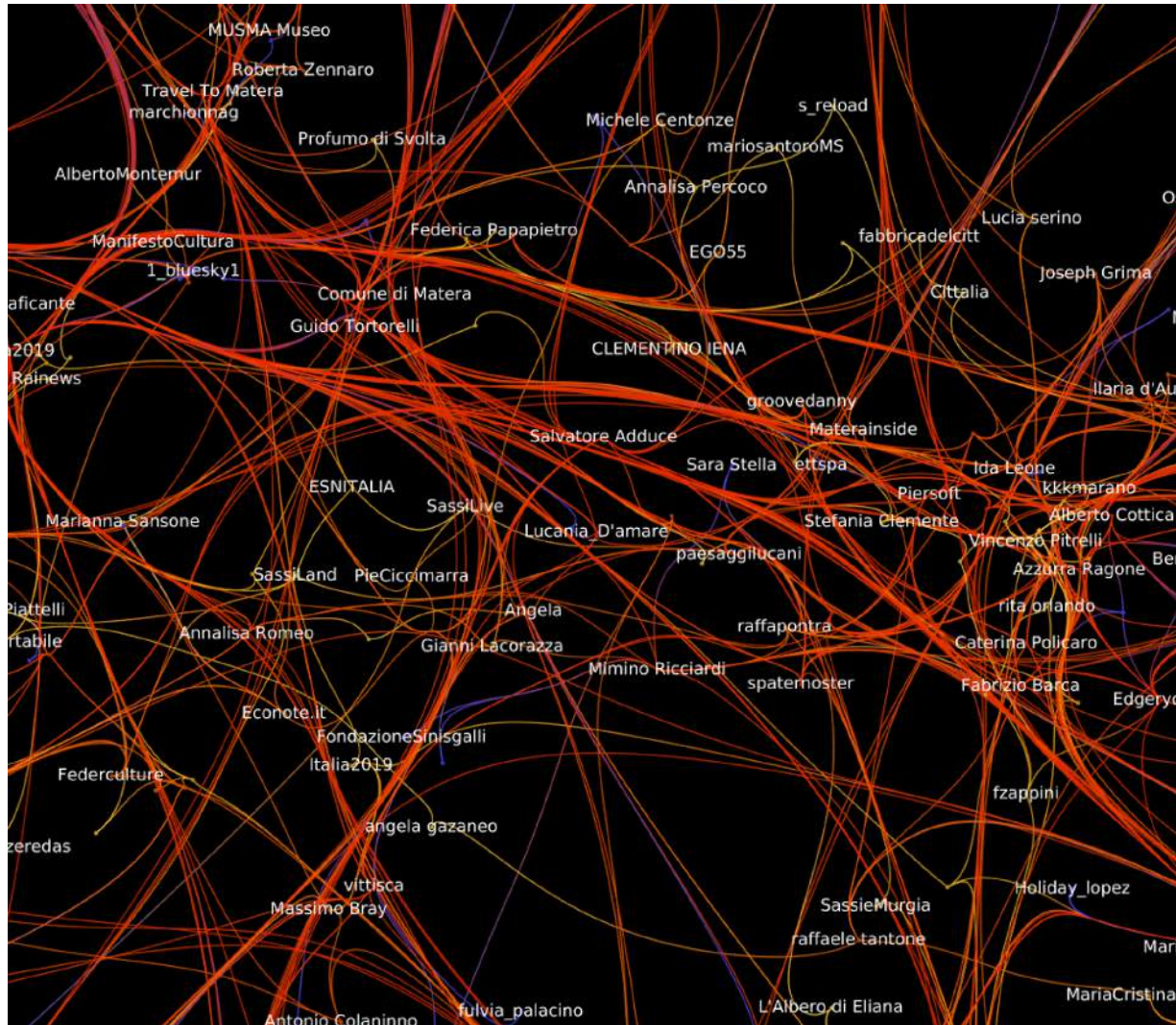
Now a more shallow relationship is the mention:

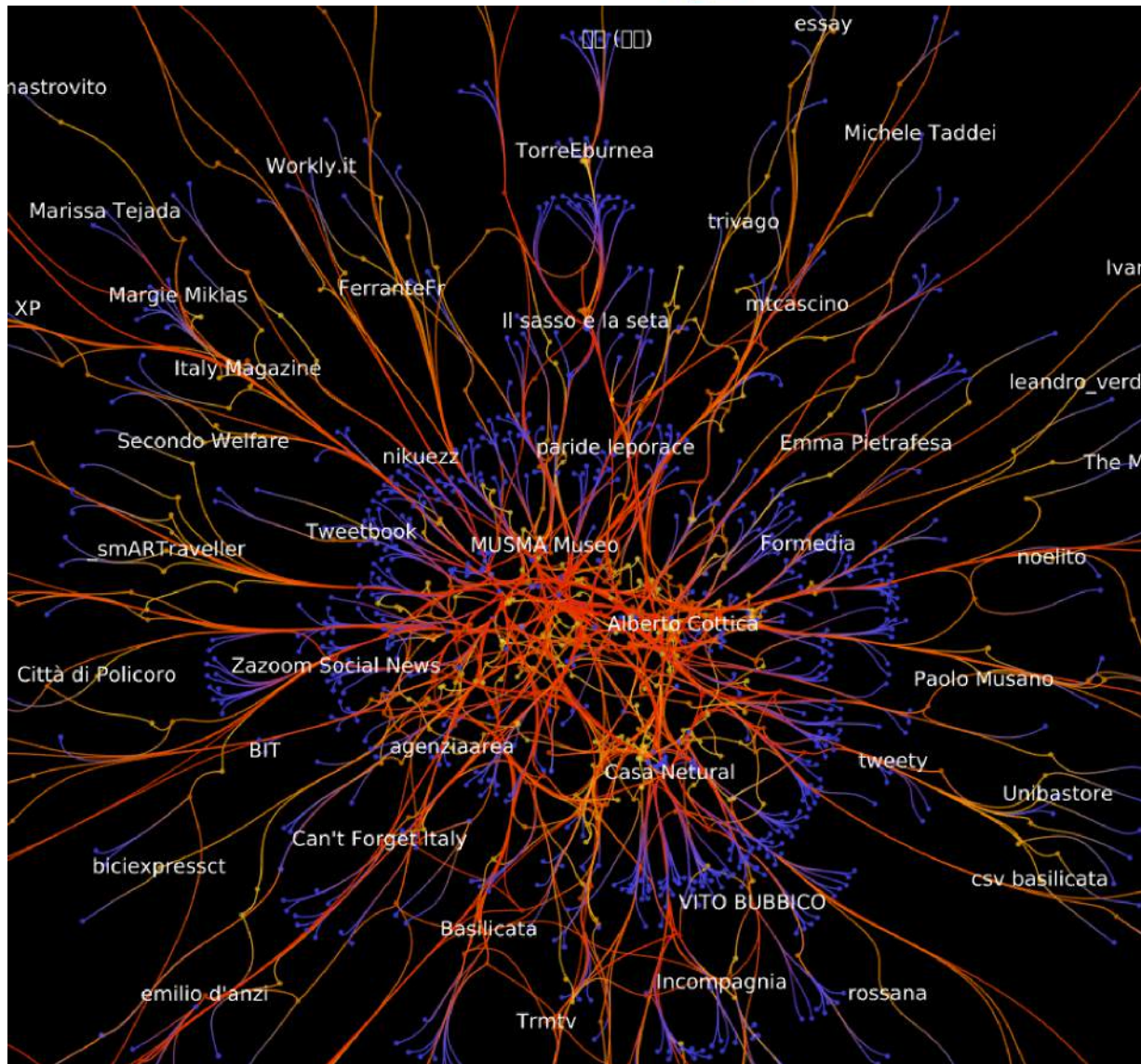


yet again the network is very centralized around a core of undistinct active actors

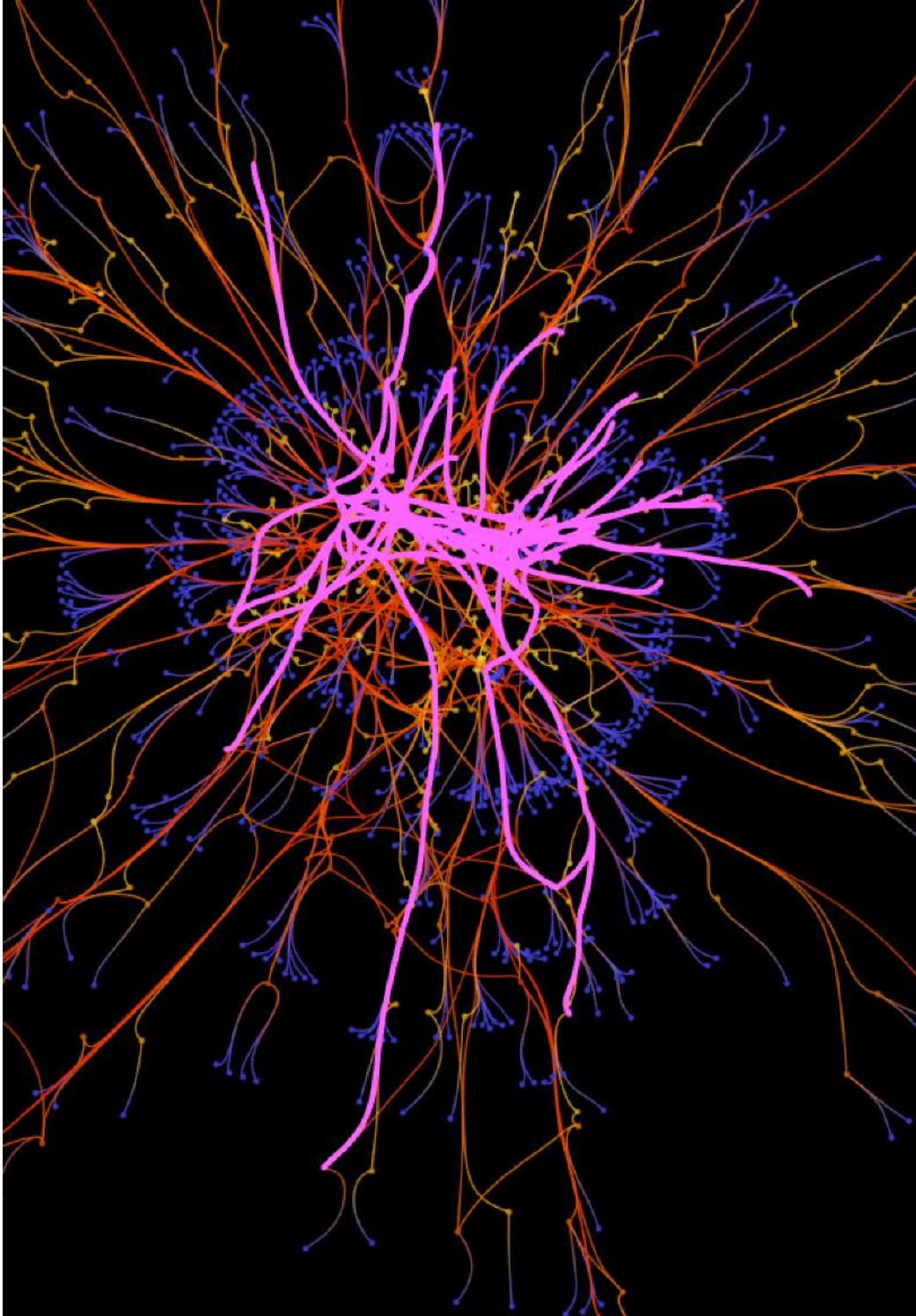


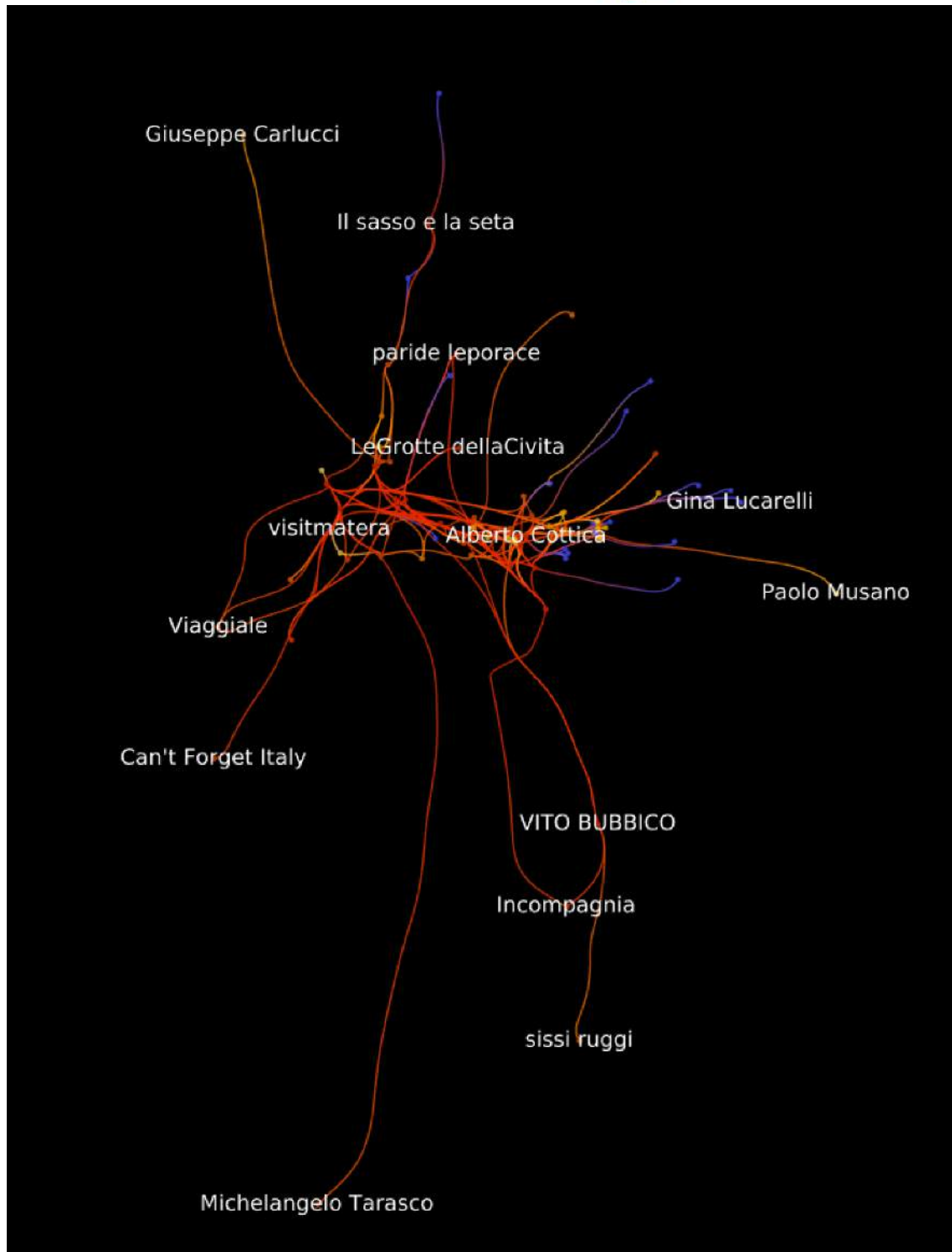






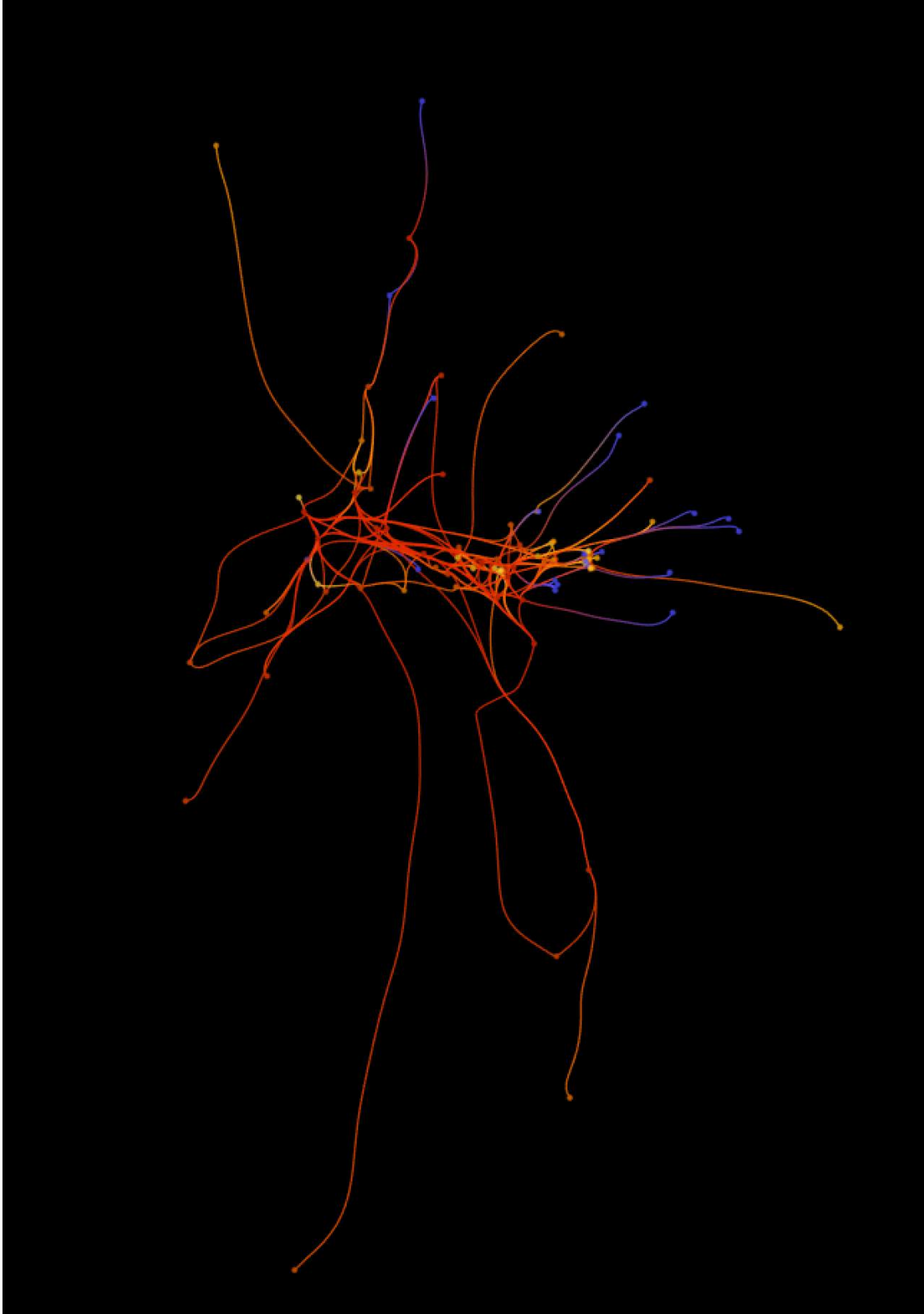
The bilateral mentions

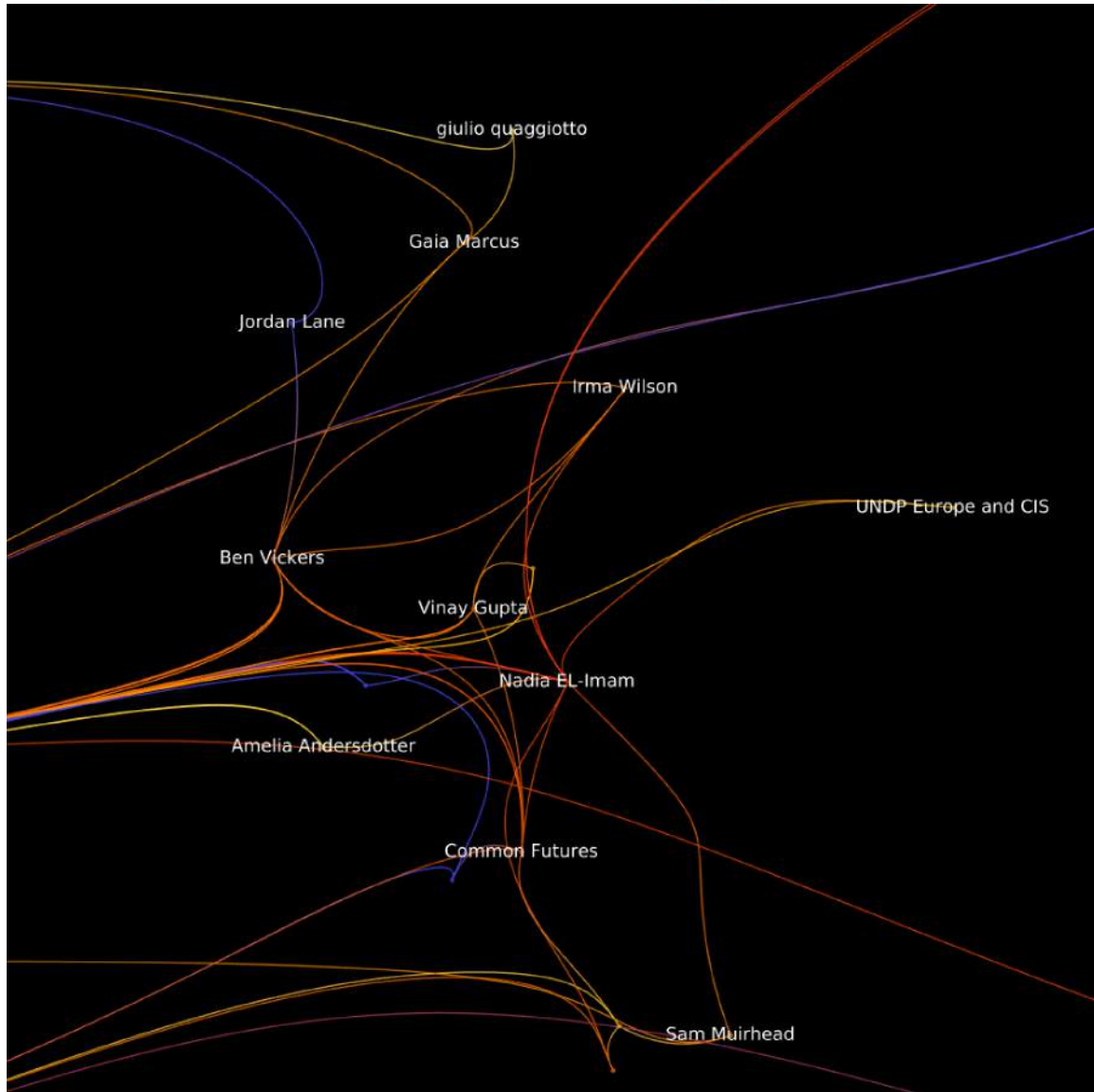




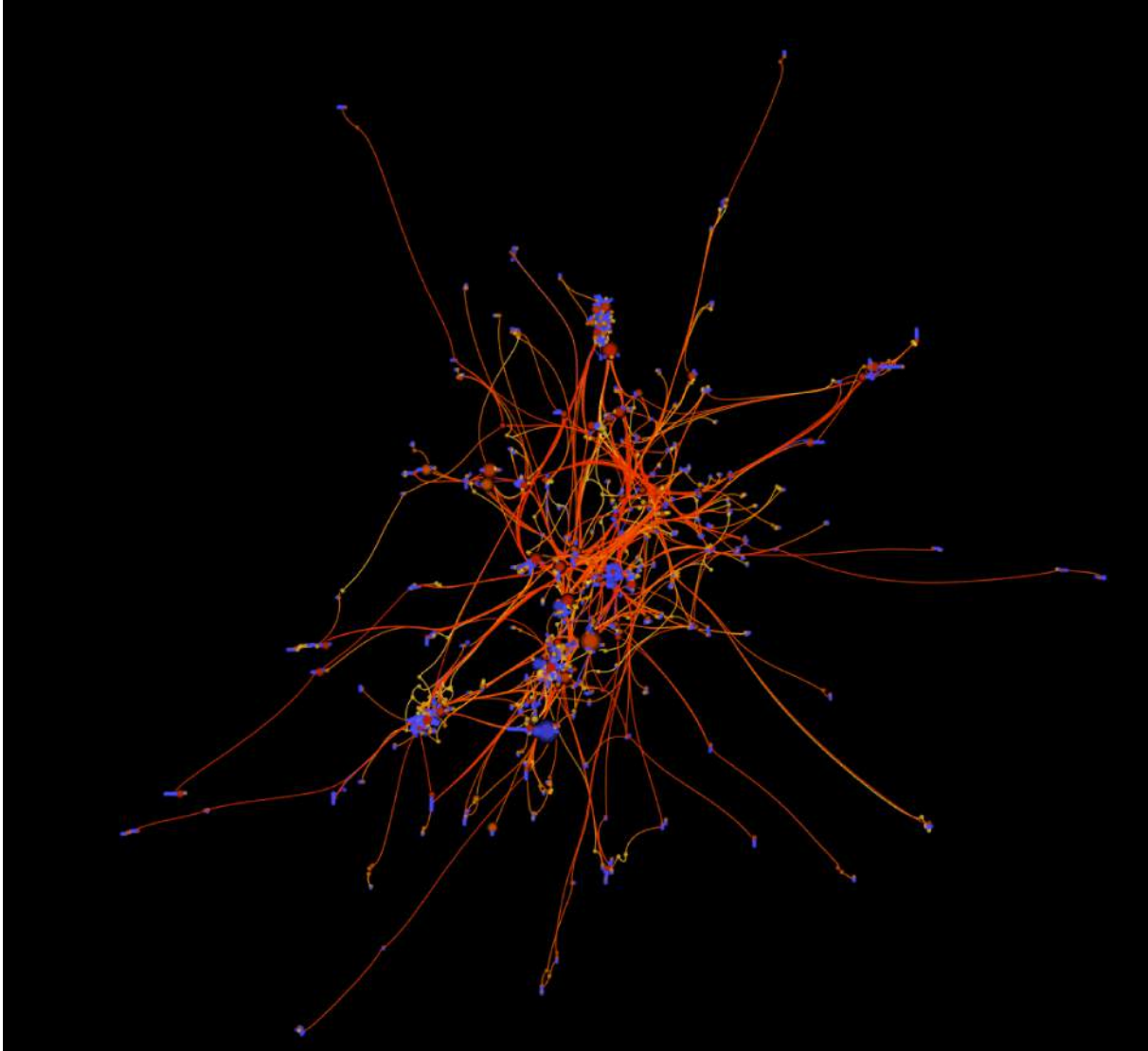
with some focus

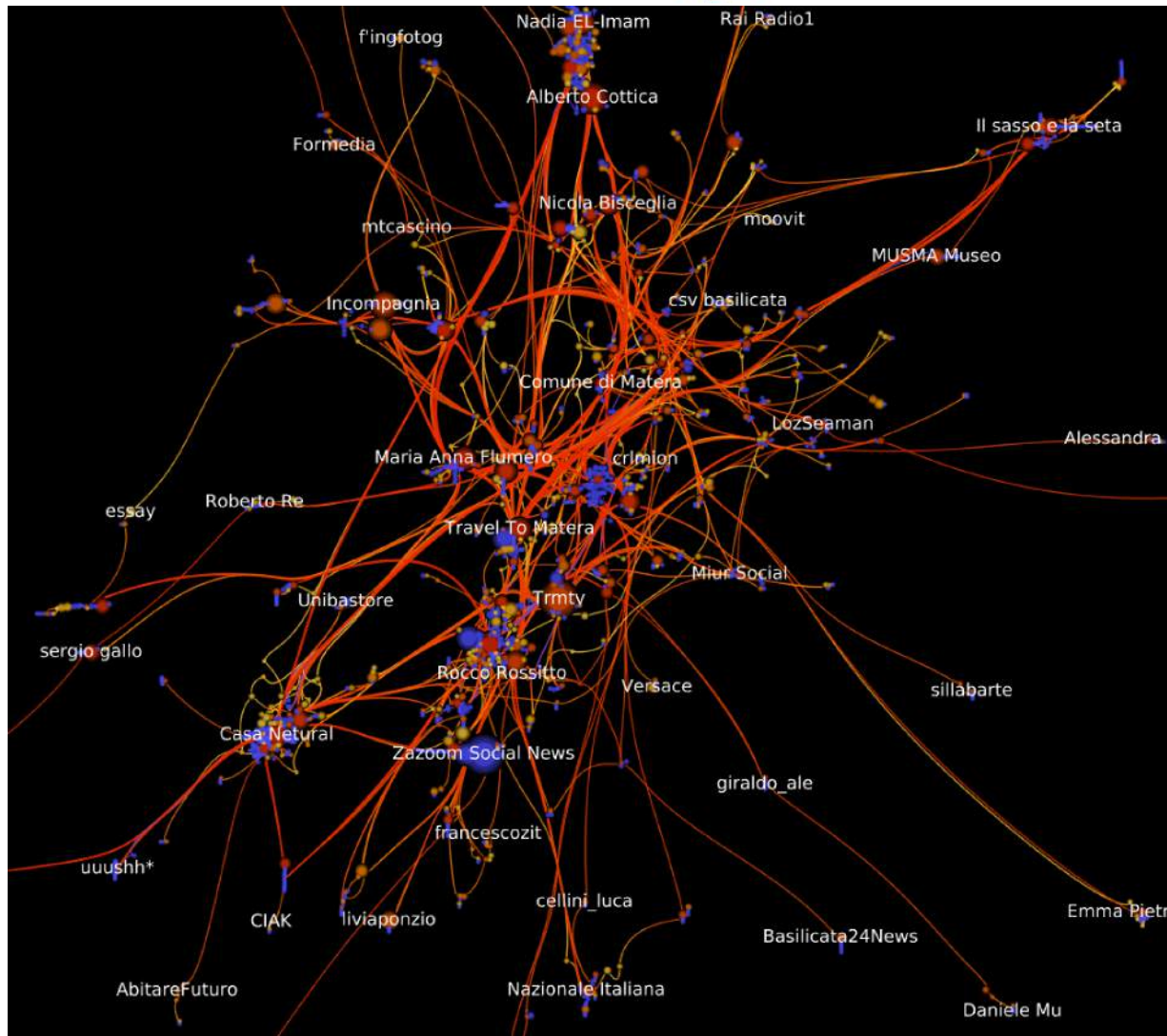


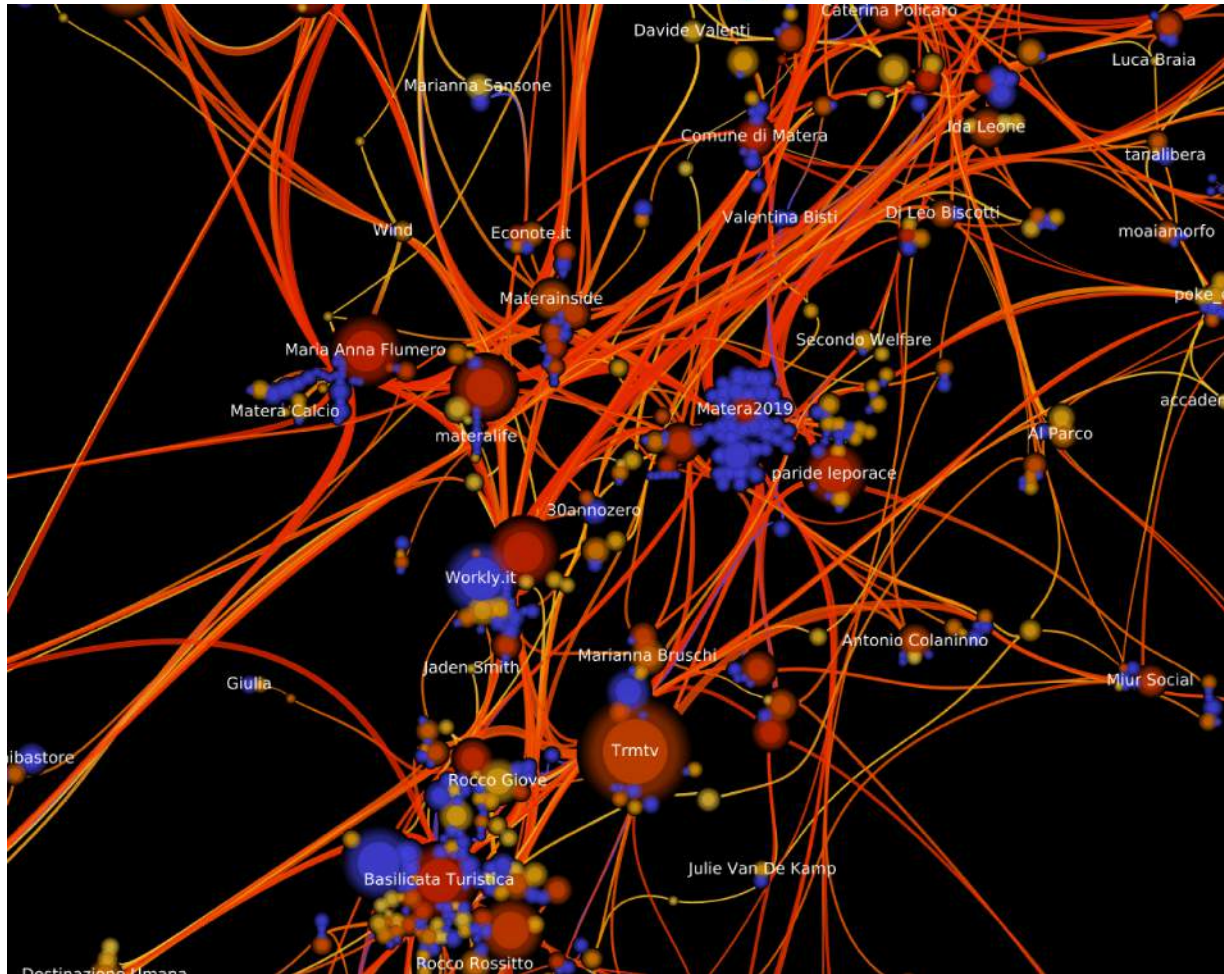




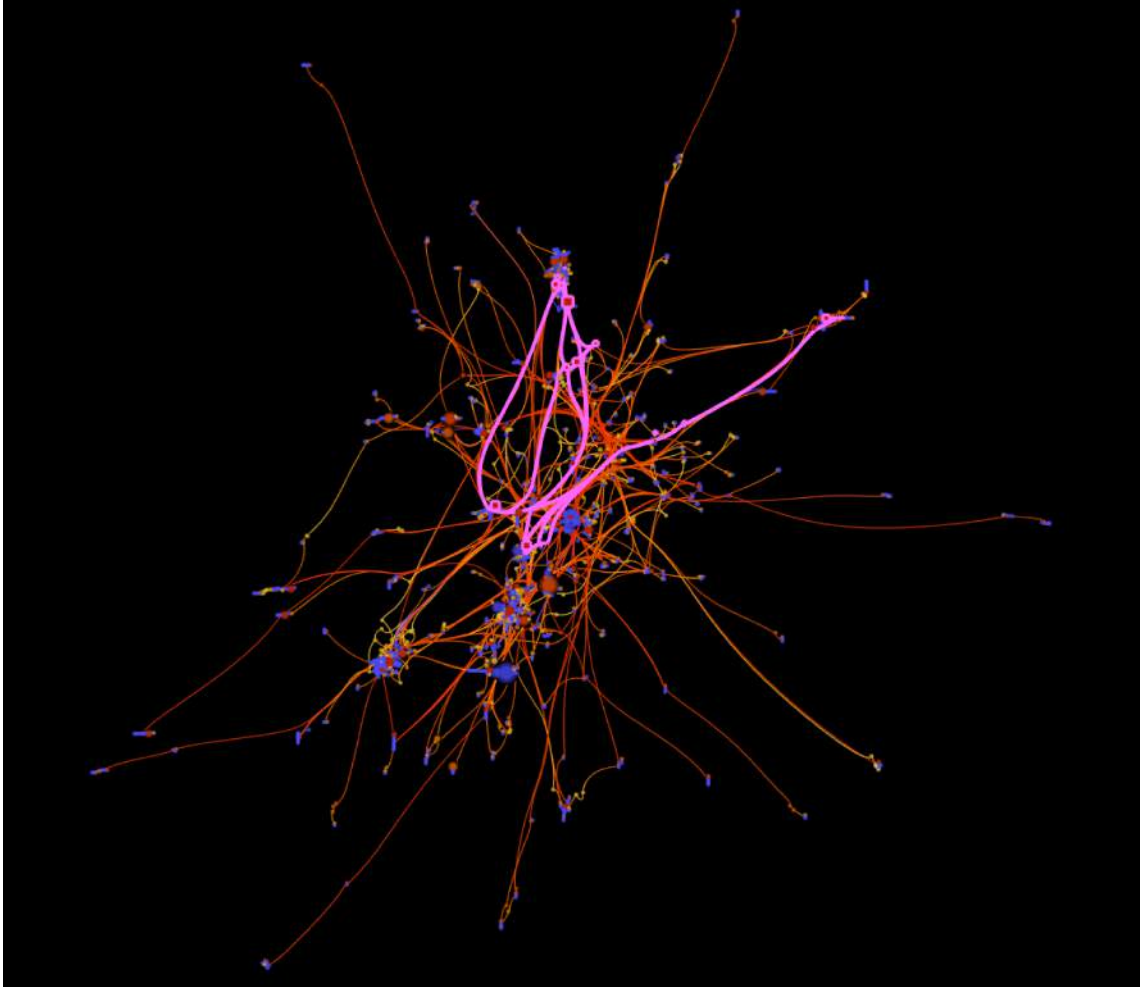
We can actually also clusterize this network

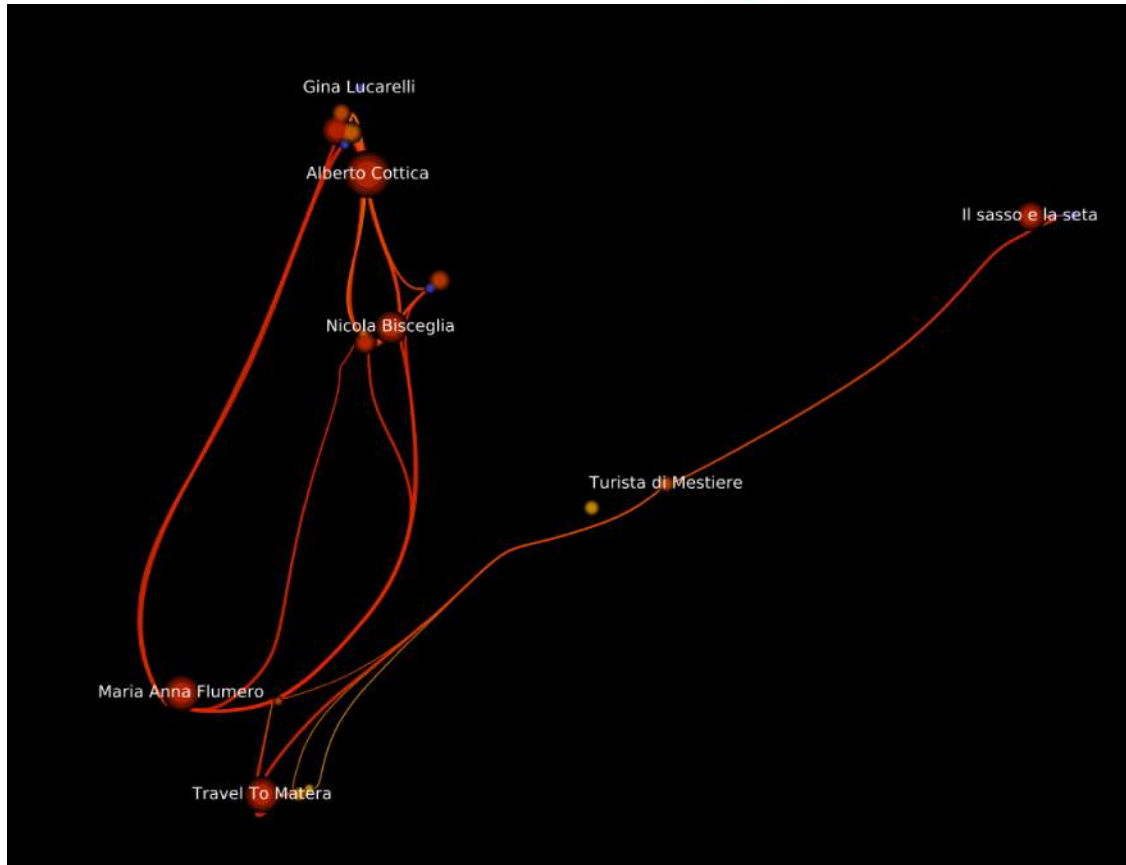






And we may want to look at how the people replying bilaterally do in mentions:



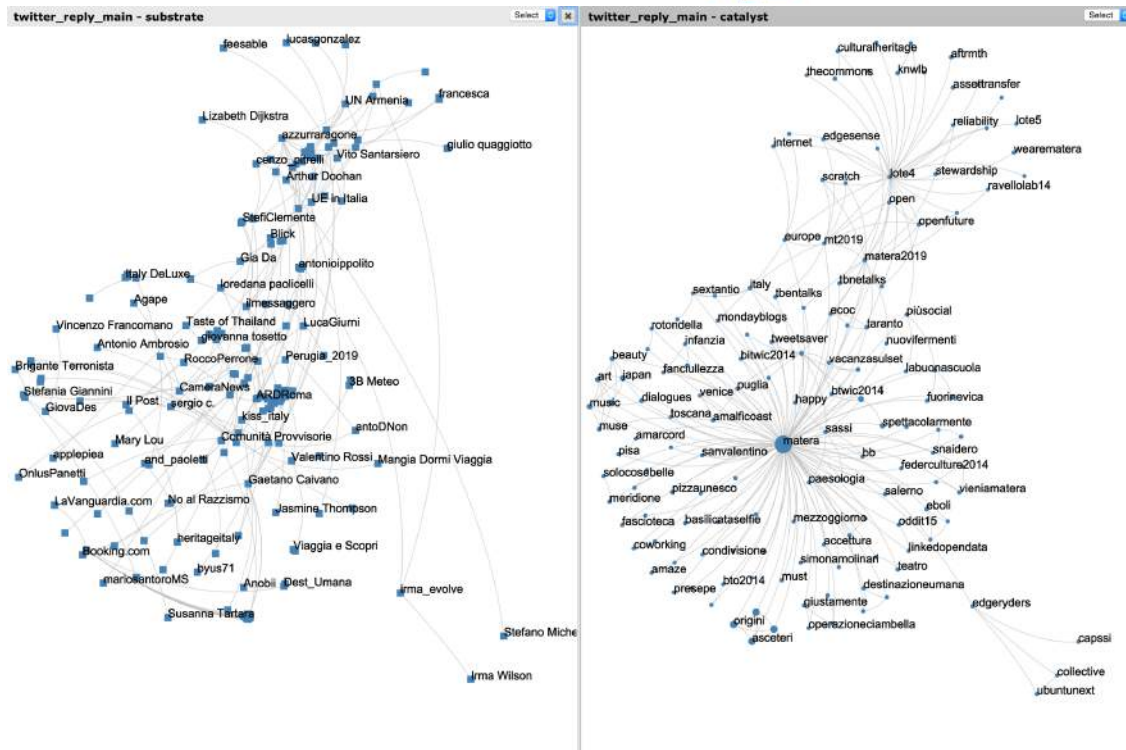


Interestingly they are also disconnected.

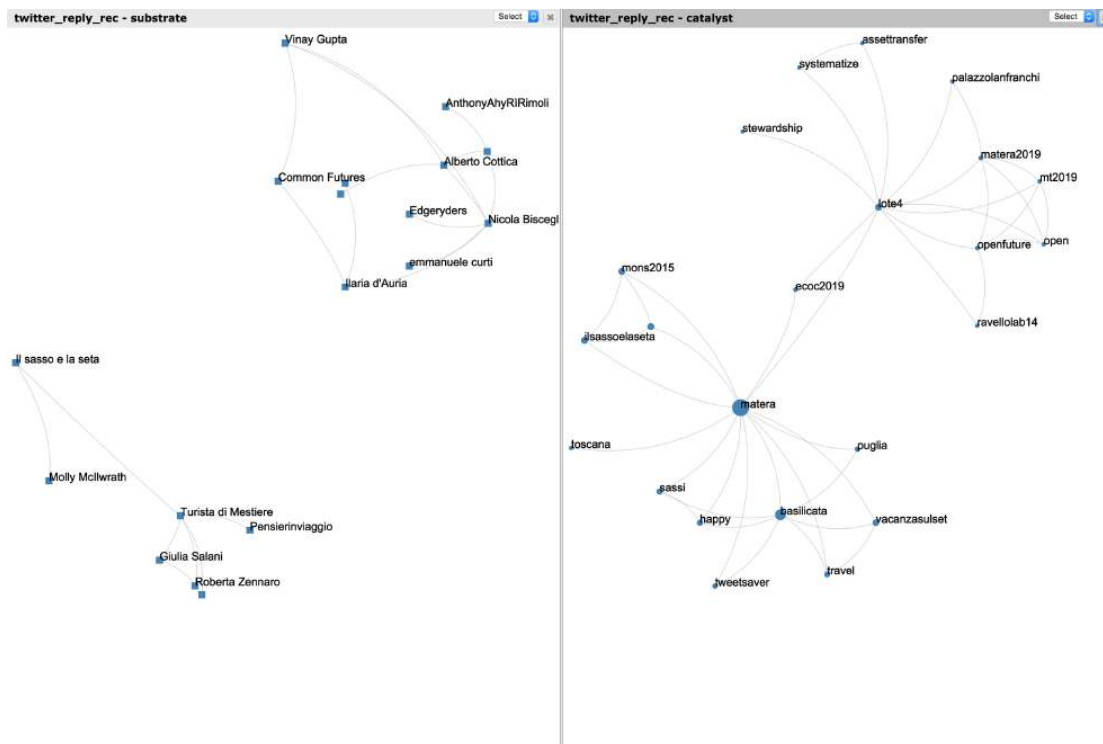
Now a focus on the semantics behind, the interest:

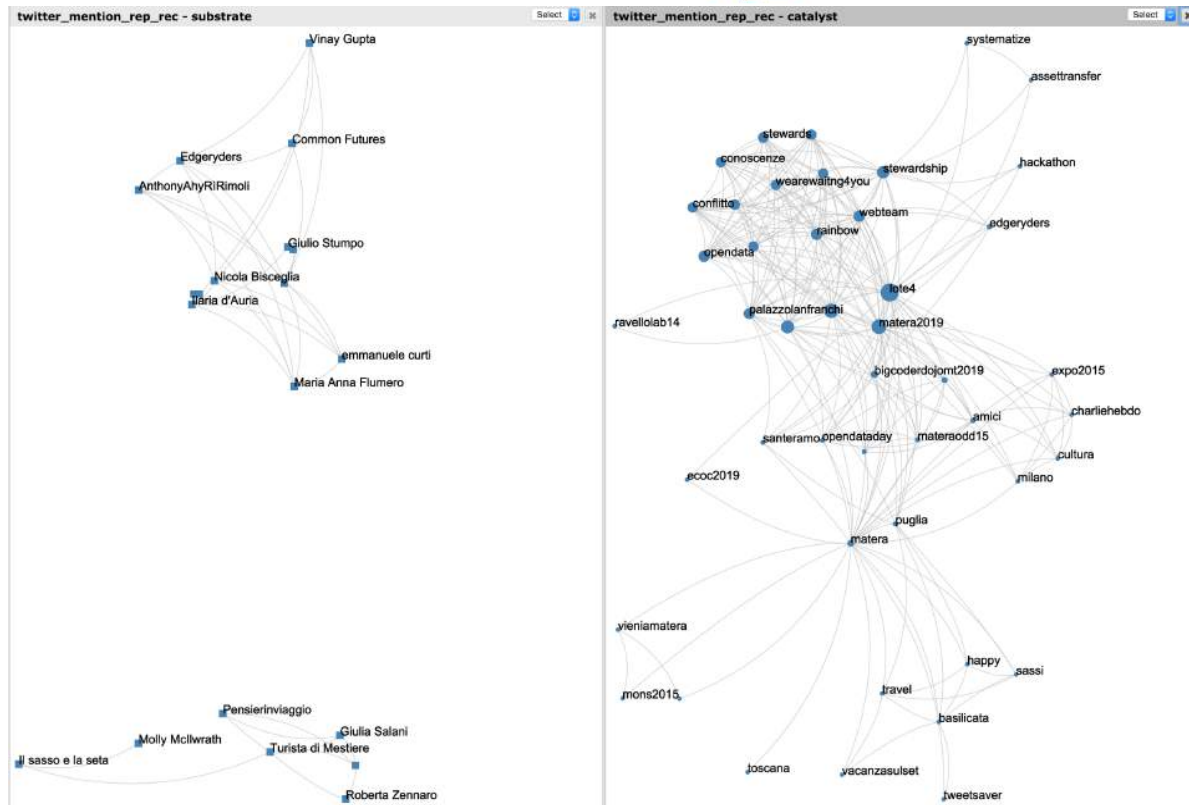
This data shows how hashtags are distributed over mentions and replies

This one displays what people discuss about in direct replies (twitter_reply_main.json)

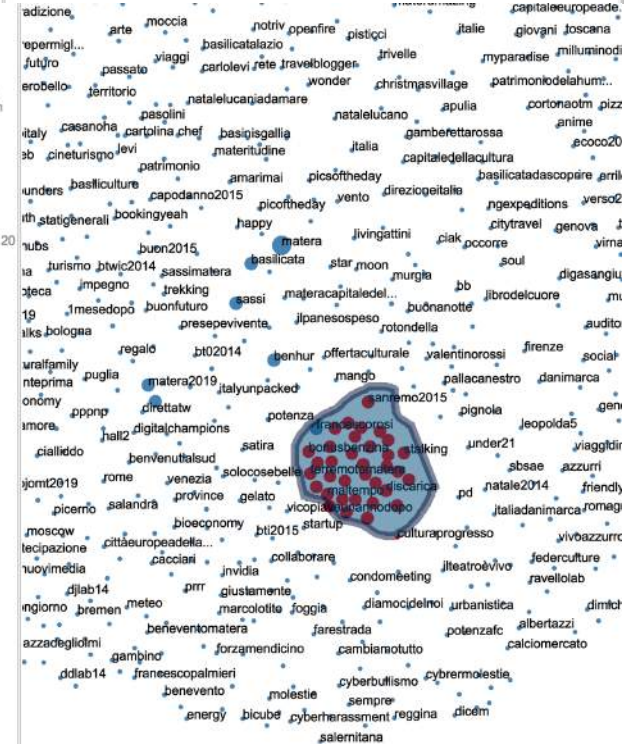


Here limited to the set of reciprocal interactions:

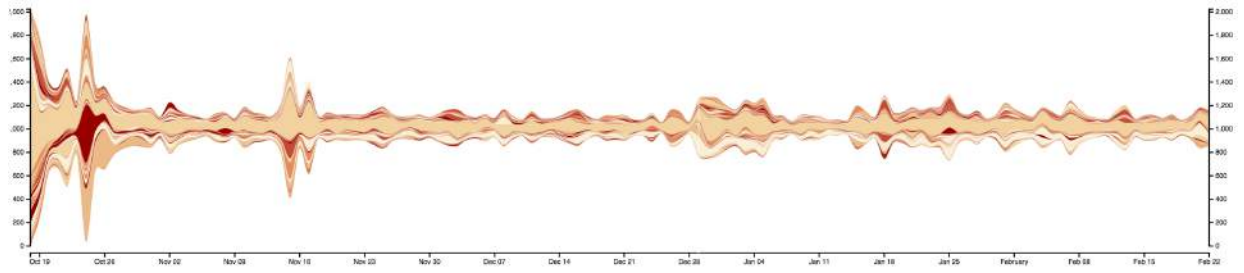




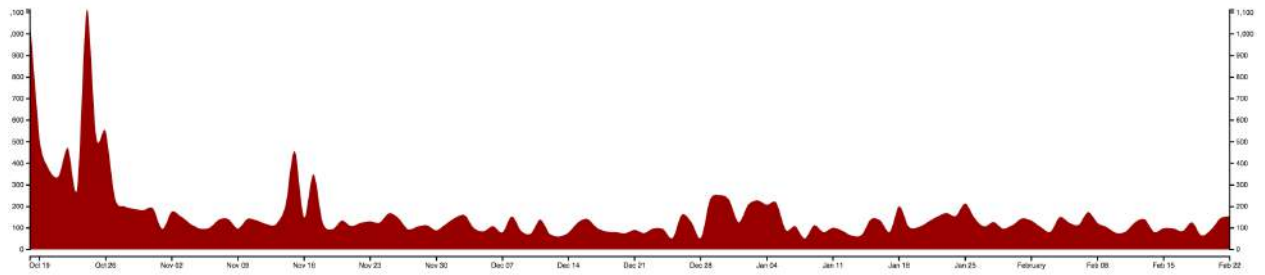
Just for the fun, here is the whole set of mentions



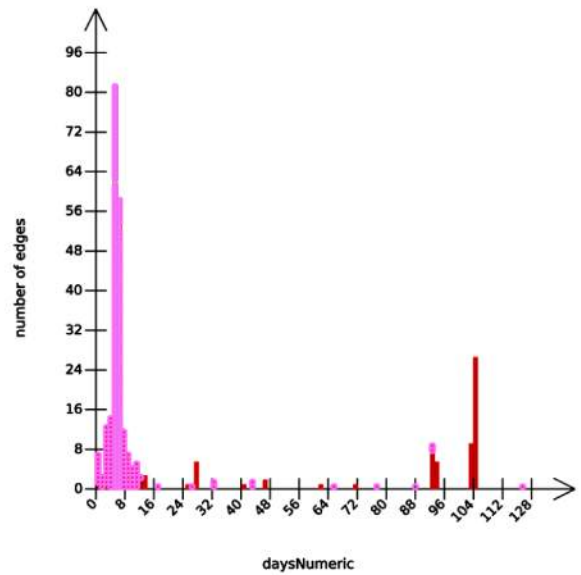
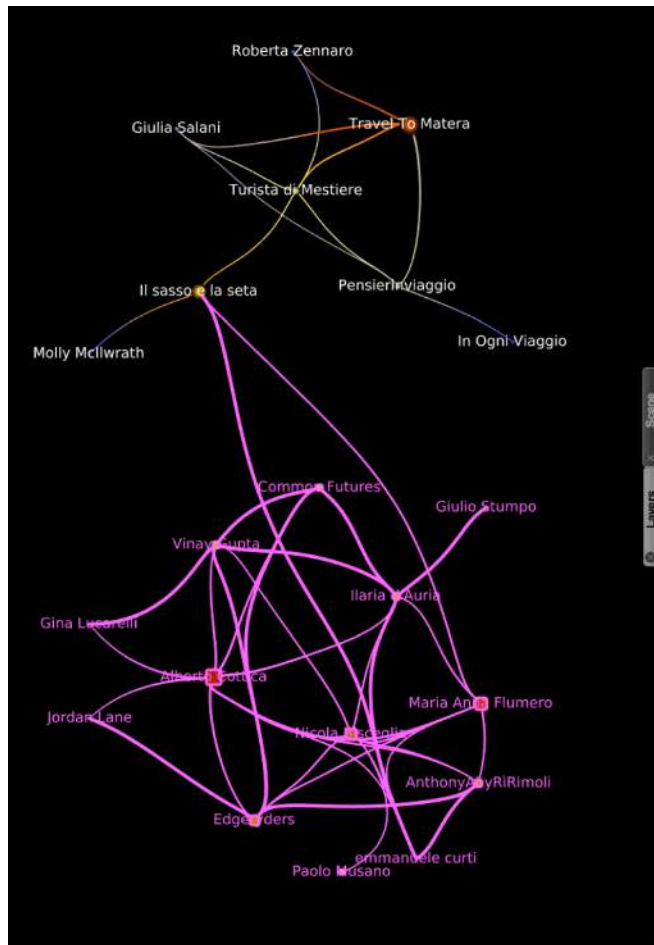
51



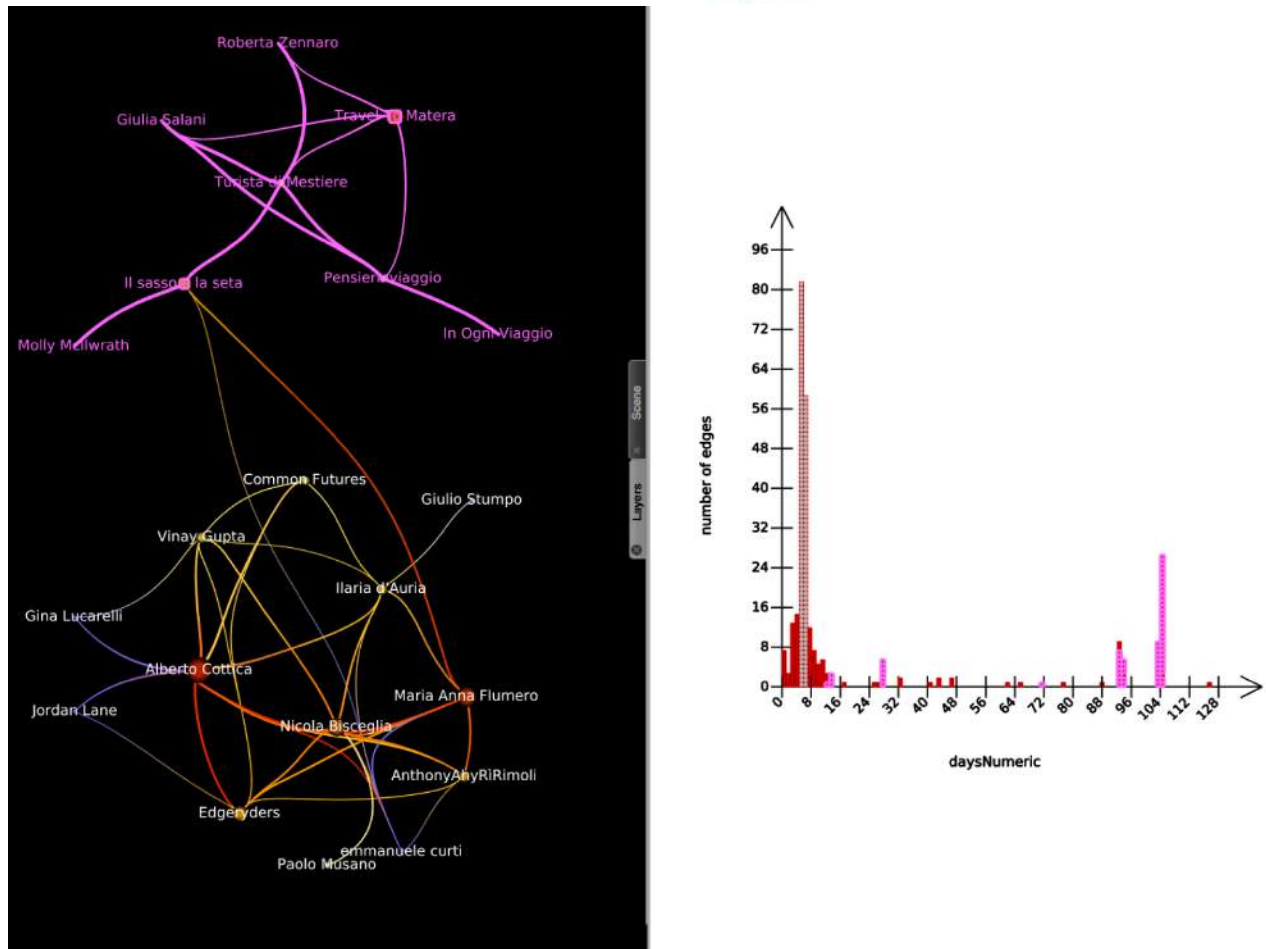
all words



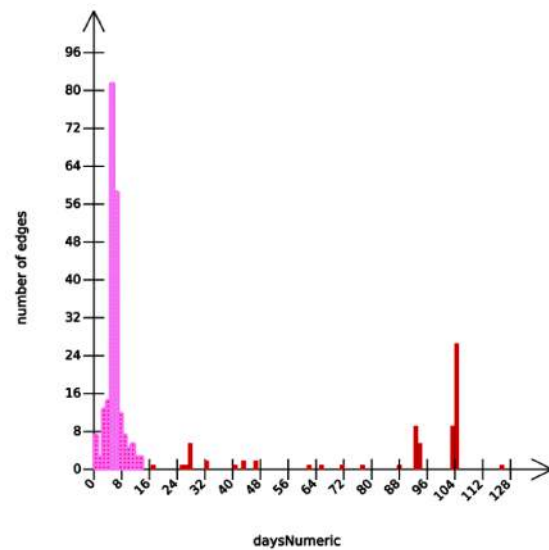
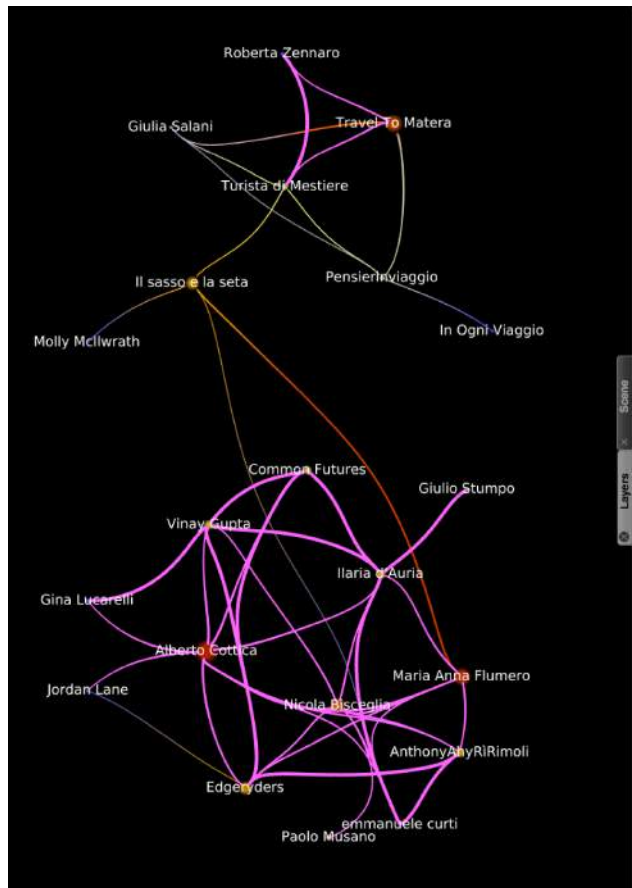
All tweets



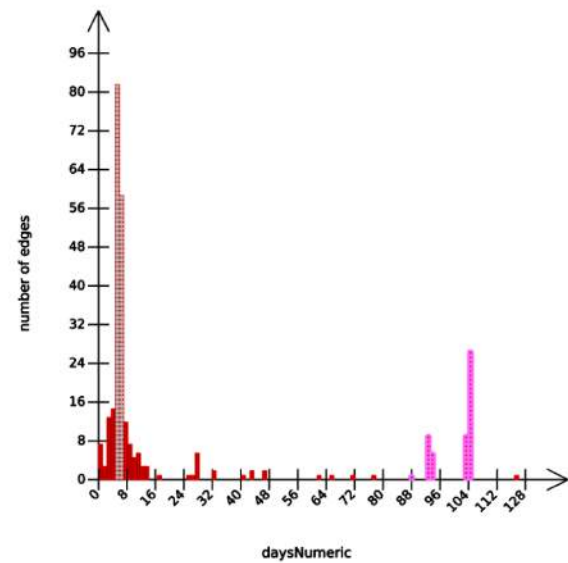
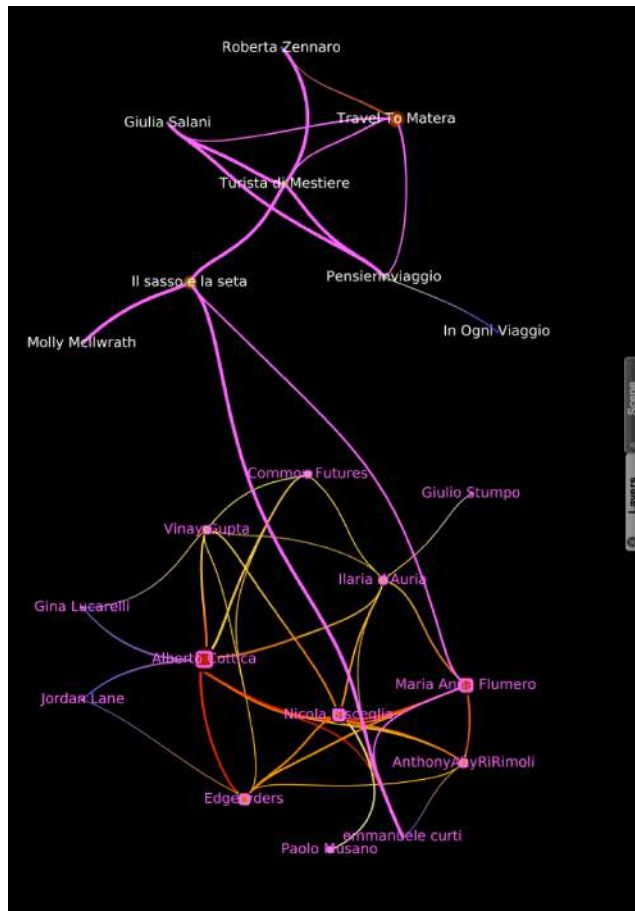
Distribution of the tweets of the core communities corresponding to the edgeryders community (selection from users)



Distribution of the tweets of the core communities corresponding to the “travel_matera” community (selection from users)



Distribution of the tweets of the core communities corresponding to the early peak (selection from time)



Distribution of the tweets of the core communities corresponding to the late peak (selection from time)