

Analyzing the Dynamics of Community Formation using Brokering Activities

Matthias Trier, Annette Bobrik

Technical University Berlin, Germany
{trier, bobrik}@sysedv.cs.tu-berlin.de

Abstract. Understanding structures and processes of virtual communication networks can help to improve knowledge sharing and collaboration in a corporate setting. A widespread research method in that domain is Social Network Analysis. However, SNA only considers a summarized picture of the final structures of virtual community networks. It does not focus the understanding of the actors' underlying dynamic processes of structural evolution.

To overcome these shortcomings, we propose a methodology of dynamic network analysis based on small time windows and animations of network evolution. Based on that method we introduce the measure Brokering Activity to identify persons which have actively contributed to community formation. In a study of corporate e-mail traffic, we compare this dynamic measure with current measures to show that it uncovers important networking agents, previously ignored. By plotting Brokering Activity over time, further insights about the dynamics of networking can be achieved.

Introduction

During the past years, a growing attention on electronic collaboration and group formation among internet users but also among employees in knowledge related work contexts could be recognized. Indicators are the intensive discussion of the role of social software and web2.0 but also of corporate electronic communities of practice and knowledge management (e.g. Wasko and Faraj, 2000). This development invoked increased interest in observing, visualizing, analyzing, and even 'measuring' the structures of such networks.

Next to approaches based on simple activity logging (e.g. Cothrel, 2000), the rapid and regular advance in social network research provides a vast body of related measurements and methodologies (Wasserman and Faust, 1994). The according field of Social Network Analysis (SNA) is defined as a framework for the analysis of structured social relationships (Wasserman and Faust, 1994), which in the organizational context can reflect role-based authority relationships of formal organizational structures, informal structures based on communication, information exchange, or affection (Tichy et al., 1979).

The main hypothesis of SNA is that human behavior is influenced by structural properties (e.g. restrictions). This shifts the focus towards observing relationships between actors and the actors' embeddedness in a complex relationship network. This relationship network and the individual relationships have influential structural conditions. According to Wellman (1997), such social networks are virtually present whenever a group of people interacts electronically. That enables the systematic examination of such computer-networked communities.

According measurements within the domain of SNA can either include composition variables, i.e. the number and properties of actors, or structural variables, i.e. the properties of relationships. Some of the most important factors for evaluating actor networks are network size, relationship strength, network roles (broker, gatekeeper, pulsetaker, hub, isolate, transmitter, receiver, carrier), degree (activity, prominence, symmetry or reciprocity), betweenness and centrality, density, and diameter.

Although existing measurements for Social Network Analysis help to evaluate a variety of properties of larger networks of virtual communication, the method is only taking into consideration a snapshot of the final state of a network at time $t=T$. Such structural measures emphasize the perspective, which Leenders (1996) labels 'contagion'. Here, networks are regarded as the independent variable and actor attributes (e.g. behavior) as dependent. Leenders asserts, that viewing the network structure as the independent and changing variable influenced by actor behavior is far less addressed and calls this perspective 'selection'. Emirbayer (1997) also recognizes a 'structural bias' in the sociological conception of the social world. He differentiates between substances or processes, i.e. static 'things' or dynamic relationships. A complementing methodology to understand the corresponding processes in emerging communication networks still remains a challenging field of research. It could answer questions like how they evolved over time and how their organic properties can be fostered and utilized. According to Doreian and Stokman (1996, p.2), this may be due to the fact, that structures are easier to observe and 'social network processes may seem more elusive for formal model building.

The authors define a social process as a series of events involving relationships that generate (specific) network structure. Studying network processes therefore requires the use of time, i.e. temporally ordered information in addition to descriptions of network structures as summarized information.

Research of this type is hence focused on network change. Here, Moody et al. (2005) generally differentiate two forms of analysis: One approach plots network summary statistics as line graphs over time and the other is examining separate images of the network at each point in time. Such images are often difficult to interpret, since it is difficult to identify the sequence linking node position in one frame to position in the next.

Early studies of network change compared only few points in time to evaluate progress (e.g. Freeman, 1984). Primary focus has been on the (cumulative) state of a network and the contained variations of individual measures, including intra-pair attractiveness, popularity, reciprocity, or transitivity (e.g. Doreian et al., 1996). Mostly, these studies were computing averages for the complete network and have not considered the role of actor clusters or the individual actors (also cf. Moody et al., 2005). Further, they had to deal with the problem of attrition in the sample.

Related approaches analyze the structural changes of networks after (disruptive) events or concentrate on the transitions in network structure between points in time (Hammer, 1980). A stochastic approach is pursued by Snijders (2001), who examines statistical models of network change (for a more comprehensive overview about different approaches, cf. Doreian and Stockman, 1996).

Only recently, Moody et al. (2005) introduced their related activities and started a methodical exploration of dynamic network visualization. In discussing the benefits of dynamic analysis of social networks, the authors also see that the issue of identifying important nodes is dependent on the longitudinal life cycle of the social network. For example, they note that “understanding of the betweenness of these centre nodes changes once the temporal nature of the network is revealed” (Moody et al., 2005, p.1218). Much like the approaches demonstrated in the next sections, the authors assert that a “static network pattern often *emerges* through a set of temporal interactions”. Further they stress the importance of dynamic visualization for understanding ‘change’ in a network and propose various new terms, like ‘dance’ or ‘pulse’ to describe dynamic patterns of node behavior.

Research Objective

With improved means of capturing large sets of longitudinal communication data from virtual communities, novel means of analysis, visualization, and measurement can be developed to improve the understanding of not only the structure of the final network and its general transition, but also of the actual processes of its formation driven by the community's members.

To approach this objective, we introduce and discuss a method which combines dynamic analysis of incremental network changes $\Delta t_{12} = t_2 - t_1$ ($t_1 < t_2$, $t_1, t_2 \in [0; T]$) (also compare for Gloor et al., 2004) and animation of the according network evolution. We take into account that after each specified period in time, messages and according relationships are added or have decayed. In the analysis the initial relationship link is formed, if one person contacts a second person and that in turn replies.

Based on this dynamic method, we introduce and discuss a networking measure, which we term Brokering Activity (BA). It shifts the observation perspective (contagion) from a passive analysis of an actor's network position towards an understanding of the active role of the involved persons in the process of network formation (selection).

Taking a sample of one year of corporate e-mail communication, we study the differences between this approach of dynamic network measurement and the results of measuring static community actor's message volume, betweenness, and degree centrality. We want to identify, if dynamic measures indicate actor's as being 'important' for the overall network formation, which would have not been identified by using static SNA. The according research questions are:

Research Question 1: How can nodes be identified with dynamic network analysis methods, which made significant contribution to network formation but are missed out by means of static SNA?

Research Question 2: What new insights can be derived about the formation of a community, when we look at the actions which fuel network evolution and relate them to static volume and SNA indicators?

Towards methods for dynamic network analysis

Looking at the dynamic evolution of a network over time provides the chance to move towards a whole new set of novel process oriented measures, which augment the existing set of structural SNA measures. Such novel metrics can take into consideration the iterative changes of the network structures occurring during the sample period. They can further be

broken down into personal activities. By this we can develop better means to identify important people, based on their actual contribution to the overall community structure.

To allow for insights into dynamic network evolution, we utilize available information about the time at which messages occur in order to decompose the dataset into a multitude of individual time windows. All information which is outside the time window is either not yet being included or has expired. Another important parameter is the step size (e.g. fetch the messages of the next n days at once). If the time window is set to be larger than the step size, then time windows start to overlap. For example, a time window of 30 days and a step size of 1 day are resulting in network slices from day 1 to day 30, day 2 to day 31, day 3 to day 32 etc.

For the analysis in this contribution we assume a time window and a step size of one day to capture added activity on a daily basis. On the next day, the previous day's activities have already expired and are thus removed from computation. The according time window thus contains each active node's additions to the network and we have the chance to move away from summarized pictures towards observing daily network forming activities of participating agents.

Defining the concept of dynamic Brokering Activity

Based on the dynamic analysis method above, we can establish a novel measure, which we have termed Brokering Activity (BA). For that, we assume that an activity of an actor is beneficial if it resulted in improved network connectedness, or to be more specific, in reduced path length between the network's nodes. This is equivalent to observing how an actor creates shortcuts in the network.

If for example some person directly connects to a contact, which has previously been at a path length of three steps (two nodes in between), the actor effects positive change in the overall structure of the network, i.e. he shortens pathways across the network for all nodes and hence allows surrounding and indirectly connected nodes to 'move closer'. This in turn increases the probability of direct connections between neighbors (i.e. via triadic closure) and generates shorter paths for better information dissemination in the community. To measure BA, we need to identify and count such action patterns, where the observed actors actively 'reach out' and connect to other nodes.

Formally, we define the Brokering Activity (BA) as the number of new connections or shorter paths between other nodes generated in one time

window by a node's activity. We thus count added messages after every day. For that, the algorithm needs to eliminate the actor's node i from the network slice (showing the current day) and recalculate shortest paths between the remaining nodes j to k (with $j > k$ and $i \neq j \neq k$). Comparing shortest path length among all nodes without actor i versus with actor i (and thus with its activities) gives information about the improvement of connections between the nodes caused by actor i 's activity on the observed day. All paths which increase in step size (or which became infinite) after removing actor i 's actions imply that i 's actions have a positive impact. The according formal equation is:

$$BA = \sum_{j>}^g \sum_k^g b(p'_{jk} > p_{jk}) \text{ with } i \neq j \neq k$$

p_{jk} shortest path between j and k in the path matrix of i

p'_{jk} shortest path between j and k in the path matrix without i

$b(p'_{jk} > p_{jk}) \in [0;1]$ Boolean value

g number of nodes in the network

As additional information we compute a subset of BA which we term Originating Brokering Activity (OBA): It counts how many new connections have been created by actor i 's activities after all time windows of a dataset. In analogy to BA, actor i is removed and shortest paths of all remaining nodes j to k are recalculated. All paths which are now of infinite path length imply that node i 's activity in that time window originated (already decayed or completely new) relationships and thus account for i 's OBA score.

$$OBA = \sum_{j>}^g \sum_k^g b(p'_{jk} > p_{jk}) \text{ with } i \neq j \neq k \text{ and } p'_{jk} = \infty$$

p_{jk} shortest path between j and k in the path matrix of i

p'_{jk} shortest path between j and k in the path matrix without i

$b(p'_{jk} > p_{jk}) \in [0;1]$ Boolean value

g number of nodes in the network

Brokering Activity can be illustrated using the example of figure 1. Within a time window of one day, an additional structure consisting out of 6 nodes has formed. To indicate the contribution of node 2, it is eliminated from the network. The path matrix without node 2 reveals, that 5 connec-

tions are longer (including infinite length) than before, thus $BA = 5$. Four connections within the network would not be possible without the node's activity on that day, therefore $OBA = 4$. If this procedure is being executed for all time windows of the overall dataset and each node's BA and OBA are accumulated, we can generate an overall estimate of its role for establishing the connections within the network and thus its role for network formation.

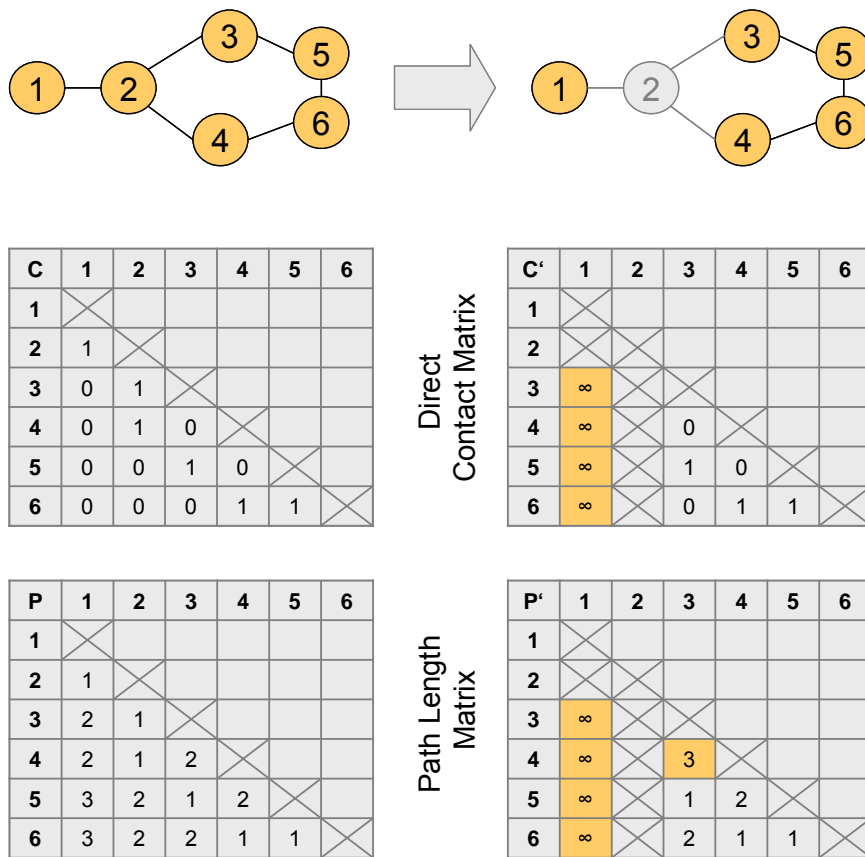


Fig. 1. Computing BA and OBA. Without node 2, five paths would be longer (hence less efficient) or impossible ($1 > 3..6$; $3 > 4$). This can be implied by the two path length matrices on the bottom. BA accumulates to 5. Four of those connections would not be connected at all ($1 > 3..6$). Thus $OBA = 4$

Software-based analysis and visualization methods

Although the computation of a node's contribution to network formation is not very complicated, a massive set of computations for a large number of time-windows, nodes, and shortest paths is required to evaluate the contribution of each participant for the final network. Hence, the above methods of measuring dynamic network properties should be implemented as a software algorithm. For our research we were able to work with a dynamic network analysis tool called Commetrix. It allows for computing time window measures and additionally provides very sophisticated functionality for animating the community evolution to visually demonstrate the actors' activities. This helps to actually represent and visually trace change in a network and adds additional insight to the quantitative results. Figure 2 shows an example of the e-mail network's cumulative evolution.

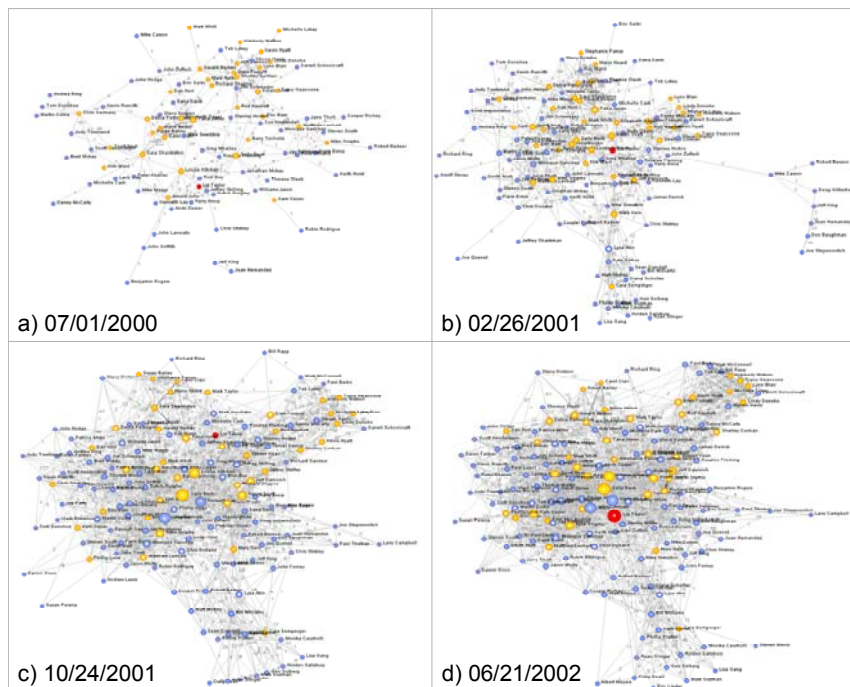


Fig. 2. The evolution of the most central author's position in the corporate e-mail network of Enron. The observed (most central) author is marked red, the size is representing the nodes' degree. Orange nodes represent members of the 'core group' of active people. The final static picture of the ENRON e-mail corpus as would be produced by an SNA Sociomatrix is shown on the bottom right. The according animation is available at <http://www.commetrix.de/enron/>

The Commetrix software consists of two elements (for a more detailed description, cf. Trier, 2005). First, a data model together with the according mining algorithms captures as much information from individual electronic discourses as possible in a systematic and standardized way to prepare the data for subsequent analysis. The primary elements are authors, their messages, and the relational information, i.e. how the authors and their messages relate to other authors' comments. Second, there are visual model specifications, which utilize the underlying data to enable insights into the complex structures, activities, and contents of electronic collaboration via node size, node color, edge length, edge color, rings, orientation, etc.

The most important feature for this paper's research objective is a special algorithm called time filter. It allows for filtering out a set of authors, which are active in a specified time period of the discourse. Obviously, this supports the computation of the time windows and the according added network structures as required for the BA measure. Our BA formula was implemented as a prototype study and we generated statistical data outputs and according animations (accessible at <http://www.commetrix.de/enron/>) of the brokering activities within the network.

Data Source

For demonstration and evaluation of the above method we will take a sub-sample of a corporate e-mail network of Enron managers.

For the process of dataset preparation it is important to note, that for general semi-automated practical application, pre-filtering approaches need to be employed. They eliminate unrelated messages, which would otherwise count as valuable messages and thus would affect the results of BA. According filters must ensure that only a subset of messages with the desired length and content (i.e. related to a list of defined glossary words) is considered. However, we regard this as a configuration option of the analyst during the process of dataset preparation. As far as the current Enron dataset is concerned, all such unrelated messages have already been eliminated by its providers. We decided to include short messages as they also contribute to the overall relationship development.

We consider a set of 4526 messages between 112 authors between from 01/04/2000 to 21/30/2000. The authors formed 394 relationships with the average relationship strength of 19 messages. The network's density is 6.34 percent. The diameter only amounts to a path length of 9 steps. The core group of active people, which together accumulate 80 percent of overall message volume, has a share of 14 percent. The maximum degree

is 23 contacts. The most active author has sent 786 and received 587 messages from the other contacts in the sample. The dataset provides information about the senders, receivers, time stamp, contents, and organizational hierarchies.

The time period has been analyzed with conventional SNA measures like betweenness and degree centrality. Additionally message volume has been counted for selected nodes. The dataset has then been disaggregated into time slices of one day length in order to apply the measurement method for counting brokering activities (BA) as introduced above. The resulting set of ‘important’ actors is being compared with the initial picture from conventional static social network visualization in order to identify differences.

Data Analysis and Results

Message volume, degree (number of node *i*’s direct contacts), degree centrality (the percentage of node *i*’s direct contacts versus all contacts), betweenness (measured as the percentage of shortest paths which run through node *i*), static BA (where we have set the time window to equal the overall duration *T* of the sample, i.e. one year) and dynamic BA (time window set to 1 day as introduced above) have been analyzed for the dataset. We computed the top five actors of the network for each of the indicators and compared the ranking lists. The resulting differences are shown in table 1.

Table 1. Comparing measurement results of static methods (time window = overall sample length = *T*) versus dynamic BA (time window = 1 day)

Message Volume		Degree and Degree Centrality		Betweenness		BA static		BA dynamic	
Node	Val.	Node	Val.	Node	Val.	Node	Val.	Node	Val.
1173	1373	1173	23, 20.7%	1173	17.7%	1173	830	1173	2782
1192	1168	1192	21, 18.9%	1264	16.6%	1202	728	1157	1284
1271	902	1271	21, 18.9%	1261	12.8%	1244	635	1174	940
1175	597	1175	19, 17.1%	1202	11.7%	1264	570	1264	894
1198	523	1198	17, 15.3%	1244	10.2%	1174	304	1263	799

As a first result, it can be observed, that BA is delivering a different set of important actors in the electronic community. To better compare the effects of BA with the other measurements, we have computed an artificial static BA. Here, we have set the measured time window to equal the sample period. This implies, that there is only one time window of length $[0;T]$. For example, eliminating the actor 1173 would result in 830 longer and thus less efficient connections (including infinite length). The resulting measure indicates the improvement for the final network structure caused by node i 's position in the network, i.e. how many paths would be longer or absent if node i would not have existed? To some extent this relates to the betweenness, as if an actor would be positioned on shortest paths between other nodes and would be removed he would cause the paths to increase in length or become eliminated. However, BA is focusing on nodes, and examines how many nodes would be unconnected or connected less directly if node i would not have existed. In other words BA is eliminating a node from the network and calculates how the network decomposes.

The actual novel list of actors who were important for forming the network during its evolvement is provided in the column on the right hand side of table 1. The dynamic measure emphasized 3 actors, which have not previously been recognized by betweenness and even 4 new nodes compared with message volume and degree.

We can see that nodes 1157 and 1263 could not be identified with static computations, i.e. by only considering the final network structure. They only come into focus after the community's formation processes are analyzed. Figure 3 shows a comparison of static and dynamic BA in the resulting network visualization. Black nodes were active in networking (top five in BA ranking). Their node size shows, how important they were when measured with the respective static message volume, degree, betweenness, or static BA. Small black nodes thus imply, that they are ignored by the conventional approaches.

The result can be interpreted as follows: The now identified persons have done much for forming the network. They reached out to other nodes and improved their links among each other. Still, in the end, those nodes have no exclusive network position anymore in terms of betweenness or static BA. This mainly results from the activities of their surrounding nodes, which indeed connected well and slowly made the observed node less important for the structure.

The ranking computed with static analysis emphasizes nodes, which are important within the final structure of the network, e.g. nodes which connect to small peripheral clusters (1202, 1244, 1264), or nodes, which have a high betweenness because shortest paths are dependent on them (1173). Hence, static analysis identifies obvious nodes. Compared to that, dynamic

analysis also identifies other nodes (e.g. 1157). Those do not seem to establish critical positions but at specific periods in time have been important for forming dense network structures. Compared to betweenness, dynamic analysis does not overemphasize nodes with connections to small peripheral clusters if the nodes do not frequently work to maintain that connection (1202, 1244), which could be an indicator of inferior importance.

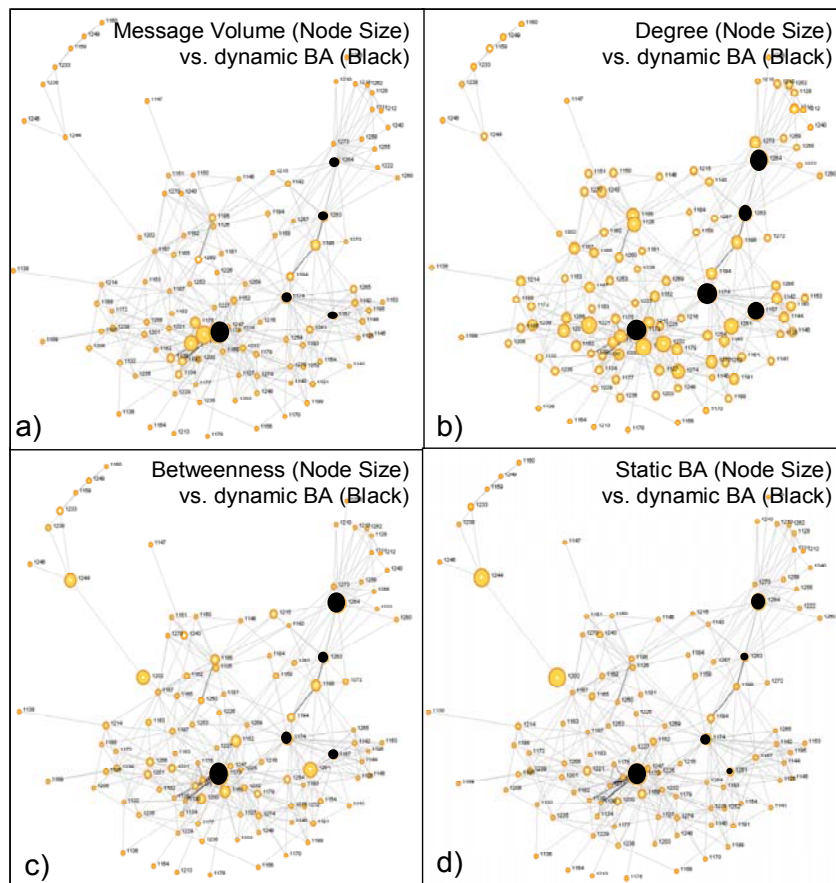


Fig. 3. Comparing results of static measures with dynamic BA. Node size represents the static measures message volume, degree, centrality, and static BA. The black nodes are the top five nodes in the BA ranking. Small black nodes thus imply that the respective measure misses these active nodes' value for network evolution

To visualize how the contribution of nodes affected the network, in figure 4, the brokering activity of three active nodes is shown over time. Generally, we found that Brokering Activity increases throughout the year.

The network grows simultaneously. The three nodes show different activity patterns. Node 1173 shows relatively constant Brokering Activity in the whole sampling period. The impact on the network is relatively small, though. Node 1262 and 1157 only appear as brokers in the last third of the year, however, their activities affect much more connections between surrounding nodes. Actors 1262 and 1157 are participating in the improvement of a few large subnetworks, whereas 1157 is improving many small subnetworks during the sampling period. The reason for the limited reach of that actor's networking activity and the resulting missing formation of larger networks remains to be examined.

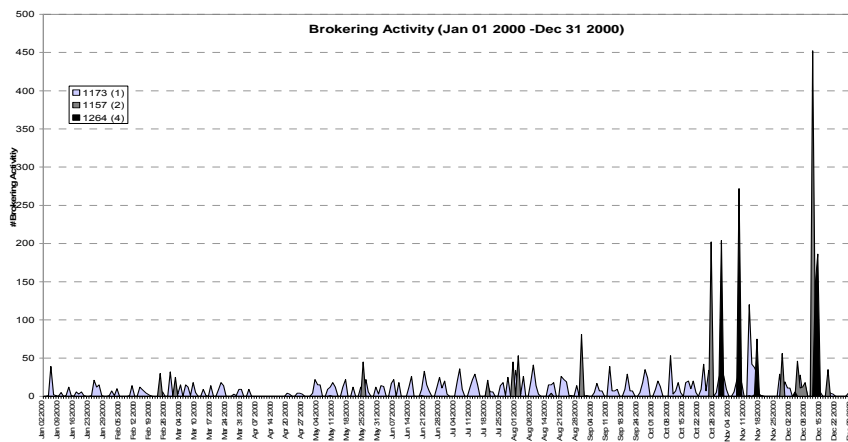


Fig. 4. Comparing brokering activity of three actors. Node 1173 is continuously acting as mediator but connects only few nodes, his broker role has only a small reach. 1157 and 1264 show no significant activity in the first half of the year, but had far reaching brokering impacts in the second half

Discussion and Conclusion

The network study introduced in this paper shows that static measures of Social Network Analysis are not able to indicate all important agents for the network's formation processes. They rely on a static snapshot of the final network structure and ignore much information in their analysis. Eventually, the resulting measures are very ambiguous. Large message volume or a large number of contacts does not mean that the actor was important for the formation of the community. Rather, many actors remain undiscovered although they contributed much to the final network structure but have not many contacts or messages. Moreover, nodes with high between-

ness may be important for a network's final structure, but this measure has a negative connotation as nodes identified as important are indispensable to some extent and are critical intermediaries (or bottlenecks) for information dissemination. It can not be recognized, whether they have actively contributed to get into this exclusive position or whether they are located at exclusive paths to only unimportant peripheral nodes with only few communication acts.

These reasons motivated the examination of the novel measure Brokering Activity (BA) based on a dynamic network analysis and visualization method. The measure identifies important actors by counting their activities that resulted in shortening the network's paths (i.e. creating short-cuts) and thus moved participating nodes closer to each other. Such processes are eventually the essence of forming digital communities with dense relationship networks over time.

A longitudinal plot of BA can further show, if the observed node came into its position via many actions with small impact or only a few interventions with a large impact. A plot of all brokering activities simultaneously shows a proxy of general networking activity over time.

As most SNA measures provide more insights when used in combination, we suggest relating BA results with conventional static measures to achieve a diverse and complementary set of analytical results. Every indicator captures different aspects of importance: Degree shows how many contacts a node has established, message volume shows the intensity of a node's communication and relationships but counts only messages to existing neighbors. It can hence not indicate crucial activities for further community formation. Both measures do not uncover, how good a node utilizes its existing structure to integrate the overall network. A person does not have to be important, if it knows the right people to have a large impact on the overall network formation. Thus it can be more advisable to promote a few contacts, than to maintain many contacts with low reach into the network. The measure Brokering Activity thus emphasizes that being an intermediary in a network is as important as activities and contacts. It shows how far the influence of an actor's activity reaches and how 'good' his direct contacts are.

In future research, we want to extend our analysis of Brokering Activity from counting (Boolean) activities towards weighing them according to their reach. This will answer the question of how much individuals improved the network to become densely knit. Further, we will extend our studies to examining the relationship between BA and the node's position in the network pattern (e.g. peripheral or central in star-networks or linear

networks). Finally, we want to approach the issue, that the current addition of messages assumes that only context-related messages are in the dataset.

References

- Cothrel, J.P. 2000. Measuring the success of an online community. *Strategy & Leadership* 28(2000)2, p.17-21, <http://www.participate.com/research/StrategyLeadership.pdf>, Accessed: 2002-11-25.
- Doreian, P., and Stockman, F.N. 1996. The Dynamics and Evolution of Social Networks. In: *Evolution of Social Networks*, edited by P. Doreian and Frans N. Stokman. New York: Gordon & Breach, p. 1–17.
- Emirbayer, M. 1997. Manifesto for Relational Sociology. *American Journal of Sociology* 103:281–317.
- Freemann, L.C. 1984. The impact of computer based communication on the social structure of an emerging scientific speciality. *Social Networks*, 6: 201-221.
- Gloor, P., Laubacher, R., Zhao, Y., and Dynes, S. 2004. Temporal Visualization and Analysis of Social Networks, NAACSOS Conference, June 27 - 29, Pittsburgh PA, North American Association for Computational Social and Organizational Science.
- Hammer, M. 1980. Predictability of social connections over time. *Social Networks*, 2: 165-180.
- Leenders, R.T.A.J. 1996. Longitudinal Behavior of Network Structure and Actor Attributes: Modeling Interdependence of Contagion and Selection. In: *Evolution of Social Networks*, edited by P. Doreian and F.N. Stokman. New York: Gordon & Breach, p. 165–84.
- Moody, J., McFarland, D., Bender-DeMoll, S. 2005. Dynamic Network Visualization. *American Journal of Sociology*. *AJS* Volume 110 Number 4 (January 2005), p.1206–41
- Snijders, T.A.B. 2001. The Statistical Evaluation of Social Network Dynamics. In: *Sociological Methodology Dynamics*, edited by M. Sobel and M. Becker, Basil Blackwell, Boston and London, p. 361–95.
- Tichy, N.M., Tushman, M.L., and Fombrun, C. 1979. Social Network Analysis for Organizations. *Academy of Management Review*, 4(1979)4, p. 507-519.
- Trier, M. 2005. IT-supported visualization of knowledge community structures. *Proceedings of the 38th Hawaii International Conference on System Sciences*. Los Alamitos: IEEE Press.
- Wasko, M., and Faraj, S. It Is What One Does: Why People Participate and Help Others in Electronic Communities of Practice. *Journal of Strategic Information Systems* (9:2-3) 2000, p. 155-173.
- Wasserman, S. and Faust, K. 1994. *Social Network Analysis: Methods and Applications*. Cambridge University Press: Cambridge, 1994.
- Wellman, B. 1997. An Electric Group is Virtually a Social Network. *Culture of the Internet*, p. 179-205.