# Learning Communities in the Crowd: Characteristics of Content Related Interactions and Social Relationships in MOOC Discussion Forums

Alyssa Friend Wise, Yi Cui

*Learning Analytics Research Network*, New York University, 239 Greene Street, 2nd Floor, New York, NY 10003 USA alyssa.wise@nyu.edu (Alyssa Wise), yc65@nyu.edu (Yi Cui)

## ABSTRACT

This mixed method study used social network analysis (SNA) and inductive qualitative theme analysis to compare social relationships and the underlying interactions they represent in discussions related and unrelated to the learning of course content in a statistics MOOC. It additionally examined the impact of how social relationships are conceptualized (via network tie definition) on resultant network structures and properties. Using a previously developed natural language classifier, 817 threads containing 3,124 discussion posts from 567 forum participants were characterized as either related to the course content or not. Content, non-content, and overall interaction networks were constructed based on five different tie definitions: Direct Reply, Star, Direct Reply+Star, Limited Copresence, and Total Copresence. Results showed network properties were robust to differences in tie definition with the notable exception of Total Copresence. Comparison of content and non-content networks showed key differences at the network, community, and node (individual) levels. First, the two networks consisted of largely different people, with less than a third of the forum participants found in both networks. Second, participants in the content network and communities had higher average node degree and edge weight than the non-content network and communities. This indicates that participants in the content discussions had more repeated interactions with a larger number of peers. Analysis of the contributing threads helped to explain factors leading to some of these differences, showing the content discussions to be more diverse and complex in their communication purposes, conversation structures, and participants' interaction techniques. Within content discussions, the network of learners surrounding each of the two instructors showed distinct characteristics that appeared related to the instructor's facilitation approach: one instructor often used hints and leading questions instead of providing straightforward answers when responding to learners' questions; he / she also frequently used social presence indicators. Learners in the community surrounding this instructor had more repeated interactions with a larger number of other learners than those in the other instructor's community. Finally, a group of learners tightly connected to each other through content discussions showed nascent learning community-like characteristics.

## 1 INTRODUCTION

A commonly cited concern of Massive Open Online Courses (MOOCs) is a lack of social interaction as a valuable form of learning support (Rosé & Ferschke, 2016). Interaction is an important element of quality in online learning generally (Trentin, 2000) and of particular importance for learners to connect with and engage in a MOOC (Khalil & Ebner, 2013). While interaction is a worthy goal, the effectiveness of traditional student-to-instructor communication is challenged in MOOCs by the tremendous numbers of students involved and the diversity in learner backgrounds, needs, and intents (Jacobsen, 2017). Many MOOCs therefore rely on peer-to-peer communication as the primary vehicle for interaction (Kellogg, Booth, & Oliver, 2014), with online discussion forums serving as the central medium.

Despite the potential for peer interaction to support learning and improve student experiences, the actual benefits of MOOC discussion forums reaped thusfar are not clear. First, MOOC discussions are often

plagued by a host of problems, such as low levels of participation (Breslow, Pritchard, DeBoer, Stump, Ho, & Seaton, 2013), overwhelming quantity and disorganization of posts (McGuire, 2013), and a lack of responsivity between learners (Agrawal, Venkatraman, Leonard, & Paepcke, 2015). Second, examinations of the relationship between MOOC forum participation and learning outcomes have yielded mixed and contradictory findings (Jiang, Fitzhugh, & Warschauer, 2014; Santos, Klerkx, Duval, Gago, & Rodríguez, 2014). Third, models of exactly how MOOC forums are thought to play a role in learning are rarely specified. In general, studies seem to rely on one of two broad explanations: either the forums provide a primarily social purpose that supports motivation ("I'm in this with other people", e.g. Khalil & Ebner, 2013; Xiong, Li, Kornhaber, Suen, Pursel, & Goins, 2015) or they fulfill informational needs that can help improve knowledge structures ("I ask a question and get a useful answer", e.g., Jiang et al., 2014; Wang, Yang, Wen, Koedinger, & Rosé, 2015). However the detailed mechanisms of these processes, and the relationship between them, have not been explored in-depth.

This study seeks to address the gap by using mixed-methods to explore the interactions that occur and the social relationships that develop around discussion of the content of a MOOC on statistics. Specifically: (a) social network analysis (SNA) is used to compare the structural properties of relationships created through content and non-content discussions; and (b) inductive qualitative theme analysis is used to: (i) explain these differences through features of student dialogue; and (ii) probe in-depth the ways students and instructors engage together around course content. In addition, the study explores the important question of how social relationships in MOOC discussions are conceptualized and operationalized, examining whether the use of different tie definitions affects the findings. The contributions made are four-fold: first, this work deepens our understanding of MOOC forum learning processes; second, it draws connections between the interactional features of discussion engagement and structural properties of the social networks that result; third, it provides evidence demonstrating the importance of separately examining content and non-content discussions in MOOC forums; finally, it highlights the conceptual and empirical impact of the choice of tie definition in SNA studies of MOOC forums.

## 2 LITERATURE REVIEW

**2.1 Social Network Analysis (SNA) as a Tool for Studying Learning in MOOC Discussion Forums**
MOOC discussion forums are populated with numerous participants and large amount of activities. SNA (Scott & Carrington, 2011; Wellman & Berkowitz, 1988) is thought to be useful for investigating interactions in these environments due to its ability to extract patterns of connections between learners across the large volumes of posts present (Cho, Gay, Davidson, & Ingraffea, 2007; Yusof & Rahman, 2009). However to date, the results of using SNA to increase our understanding of *learning* through interactions in MOOC discussion forums have been underwhelming. Some MOOC studies have used SNA to characterize patterns of social relationship and how they are formed (Kellogg et al., 2014; Joksimović, Manataki, Gašević, Dawson, Kovanović, & De Kereki, 2016). The results of these studies did not indicate a clear pattern of homophily (the tendency of people to collect in groups with those similar to them); even had such a result been found, the existence of such structures does not offer direct insight into how learning occurs through them. Other studies have examined the relationship between social centrality and the linguistic characteristics of learner posts. For example Dowell et al. (2015) found that learners who were highly connected overall tended to use more narrative, conversational language with simple syntactic structures. However, learners who performed well in the course tended to use expository language with high referential cohesion. This raises the question of whether SNA applied globally to discussion corpora focuses attention sufficiently on learning per se. A third set of studies has focused directly on the relationship between social network properties and learning outcomes; but results have not offered a clear picture. For example, Jiang et al. (2014) found a significant correlation between measures of forum social centrality and final grade in an algebra MOOC, but no relationship between the same variables in a finance course. Joksimović et al. (2016) found that some measures of social centrality were significantly associated with completion and distinction in two offerings of a programming MOOC,

while others were useful in one course but not the other. Finally Houston, Brady, Narasimham, and Fisher (2017) found that social centrality was not a significant predictor of final grade for any of the three courses they examined after the number of threads contributed to was taken into account.

Put together, the collective insight from SNA into how learning occurs in MOOC forums has been relatively small thusfar. There are several possible explanations for this. First, it has been argued that as a structural approach which aggregates across interactions, SNA is not well equipped to offer insight into learning processes (c.f., Wise & Schwarz, 2017). However prior work has shown that networks of how people are connected through ideas over time (Suthers, 2015) and differences between distributed versus dominated patterns of communication (Brooks, Greer, & Gutwin, 2014) can offer insight into learning processes when these structures are connected back to features of the interactions that generated them. A second explanation arises from the great diversity of topics and purposes found in MOOC discussions (ranging from clarification of course content to logistical questions about assignments, and from sharing deep connections with the learning material to pure sociality, see Authors, Blinded; Stump et al., 2013). These different kinds of interactions play very different roles in the learning process, thus analyzing the social networks that result from them collectively compiles interactions with distinct characteristics which may confound relationships and conceal important patterns. A final explanation derives from a lack of consistency in how social networks for MOOCs are constructed; specifically differences arising from whether ties are defined based on replying directly to someone (e.g., Kellogg et al., 2014) or simply participating in the same overall discussion (e.g., Jiang et al., 2014).

Following this inspection of potential factors contributing to the limitation of MOOC discussion SNA research, the current study seeks to examine learning in a statistics MOOCs by using quantitative network analysis in combination with qualitative examination of the contributing threads, while probing the analytic importance of differentiating interaction types and selecting a tie definition. The justification of the mixed methods research approach and framing of how SNA and qualitative analysis are combined in the study is addressed in Section 3 (Study Framing). The following sections review the literature on the issues of differentiating interaction types and selecting a tie definition.

## 2.2 Differentiating Interactions in MOOC Discussions

As described above, discussions in MOOC forums involve highly diversified topics and purposes. To address the problem of amalgamating interactions of different nature, a small number of studies on MOOC discussion forum have recently begun to differentiate analysis of interactions based on what is being discussed.

One straightforward approach to doing so uses the presence of sub-forum designations for categorization. This work has shown that patterns of interaction vary across forum activities on different topics; notably there are key differences between sub-forums whose purpose directly relate to the learning of course content and those which do not. For example, Gillani and Eynon (2014) built networks based on copresence within a thread in a business MOOC for each of seven sub-forums: Readings, Lectures, Cases, Final Project, Course Material Feedback, Technical Feedback, and Study Groups. Sub-forum networks consisted of largely distinct groups of learners and showed different levels of participant persistence over time, with "Cases" (for discussing learning material) showing the highest level of persistence overall. In another study, Gillani, Yasseri, Eynon, and Hjorth (2014) built social networks for two successive offerings of a business MOOC based on thread copresence in eight sub-forums: Readings, Lectures, Cases, Final Project, Questions for Professor, Course Material Feedback, Technical Feedback, and Study Groups. For both offerings, the proportion of one-off ties (edges with a weight of one) differed across sub-forums: highest for "Feedback" (used for technical support) and lowest for "Cases" (for discussing learning material). In addition, densities of interaction differed, with "Cases" and "Final Project" (both used for working on course material) showing greater cohesiveness than "Study Groups" (used for locating study partners).

While differentiating interactions based on which sub-forum they take place in is straightforward technically, it is a non-optimal approach for two reasons. First, each MOOC sets different sub-forums, often defined quite narrowly in relation to the specific course; therefore the generalizability of findings based on such divisions is limited. Second, prior studies have shown that misplaced postings across sub-forums are very common in MOOCs (Rossi & Gnawali, 2014). Consequently, social networks built based on sub-forums may not always accurately reflect the nature of relationships formed in forum interactions. An alternative approach differentiates interactions based on what is actually discussed by looking at the language used in the posts. This requires setting up a scheme for categorization and coding posts accordingly. The advantage here is the ability to designate a useful set of groupings that is concise and generalizable (rather than minutely tailored to a particular course) and more accurate categorization based on the actual posts themselves. Poquet and Dawson (2016) adopted such an approach to study a MOOC on solar energy. They manually coded all discussions into five categories: "cognitive task" (conversations about material related to quizzes and assignment), "social task" (conversations about learner emotions about the assignments), "cognitive non-task" (conversations about the course topic not directly related to assignments), "social non-task" (purely social aspects), and "administrative / technical issues" (conversations about tools etc.). Constructing an undirected social network for the whole forum based on thread copresence, they found that these discussion topics did not significantly explain network formation. They suggest that this may be because different kinds of interactions play a role at different times in the course and because the network only included posting (not reading data). In addition, it is possible that the categories used here were overly refined. Good learning discussion that provide useful information and build community around the shared pursuit of knowledge may contain both comments about material related to quizzes and assignment, comments about the course topic not directly related to assignments and learner emotions about these topics. Thus in the current work we simplify the categorization to simply differentiate conversations which related to the course content (whether directly for an assignment or not) from those which are not (technical, logistical or pure sociality).

We hypothesize that content and non-content conversations involve different types of interaction that play different roles in the learning process. Content-related interactions are directly related to learning as students engage with the course content, while non-content interactions may also have value in supporting engagement and motivation but are more distal in their impact on learning. Prior work has shown both that these two kinds of conversations are characterized by very different linguistic features (Authors, Blinded) indicating different modes of interaction, and that these different ways of interacting may attract distinct collections of learners (Dowell et al., 2015). Furthermore, there is initial evidence that participation in content-related interactions is a better predictor of course performance than discussion participation overall (Authors, Blinded). This study will empirically examine if there are differences in the social networks formed around content and non-content discussions in a statistics MOOC, and whether any differences found can be explained by characteristics of the discussions that take place.

**2.3 Defining Social Ties**

As noted previously SNA studies on MOOCs have used varying definitions to construct social ties. For example, the studies cited in this paper thusfar have used copresence and (directed and undirected) direct reply to construct ties. Tie definition is critical for SNA studies because different ways of establishing ties carry different assumptions about the nature of the interaction in social networks that can have implications for network outcomes and their interpretation. This issue, however, has not been well addressed in the MOOC literature. The majority of studies simply establish their tie definition without explanation or rationale. Moreover, there has not yet been a clear articulation of a typology of different tie definitions, situations in which their use is appropriate, and the implications for results and interpretation.

Two classes of social networks have been commonly constructed based on discussion forum data. The first is based strictly on the reply relationship in a forum (i.e. who replies to whom). In this type of network, a tie is defined as "speaking to" someone. One approach that adopts this mechanism is Direct Reply, in which a tie is constructed only if there is a direct reply relationship between two nodes in the same thread, either between the thread starter and a reply post addressed to it, or between a reply post and its reply. Direct Reply only maps the reply relations between users, without making any assumption about others who may have been informed by a post but not replied directly to it. This is a straightforward approach to tie formation in online discussion and has been used by many studies (Joksimović et al., 2016; Kellogg et al., 2014). However, when making a post, learners do not always place it in the appropriate location. In addition, some discussion forums only support limited levels of posts meaning that replies to replies show up simply as additional replies to the prior post. Thus how accurately a Direct Reply network reflects the actual relations among learners is questionable.

To address such concerns, Zhu et al. (2016) proposed the Star tie definition for network extraction when investigating the relationships between course engagement, performance and social connectivity. Star also defines a tie as "speaking to" someone; but different than Direct Reply, Star considers all posts in the same thread being tied only to the thread starter. The rationale is that even if a reply post was not addressed directly to the starter, it was made in the context of the thread and should address the topic set by it, thus the tie is considered a traceable contact to the thread starter. The Star definition highlights the importance of thread starter; however, as it does not distinguish between different levels of replies, it overlooks connections formed between learners within the same thread.

A hybrid scheme that combines Direct Reply and Star to produce a more comprehensive network was introduced by Gruzd and Haythornthwaite (2008). Their work also explores alternative ways to operationalize ties through identifying posts that mention specific others by name or assigning decreasing weights to posts further away in the reference chain; however these approaches all still only consider the act of speaking in a threaded discussion. As a learner could access and be informed by multiple posts before making their own, whether they reply to them directly or not, social relationships could still be said to have formed among them through "listening to" each other. Therefore, methods that strictly follow reply relations leave out potential interactions that strengthen relationships between learners who have spoken on the same topic in the same thread, but not to each other directly.

Copresence tie construction approaches address the issue of interaction with posts that are not directly replied to by defining a connection as "being present" in the same part of a discussion. Thus a tie is created between two people as long as they participate in the same thread or subthread, even if they do not reply to one and other. This type of network formation mechanism represents the notion of online discussions as collective conversations rather than single streams of individual replies. Within the genre of copresence networks, the commonly used scheme is Total Copresence where any two nodes in the same thread are considered as having a tie (e.g., Poquet & Dawson, 2016). However, when this scheme is used to map interaction, the size of a thread can cause issues, especially when many distinct people are involved, as is the case in MOOCs. While it is reasonable to assume that a participant in a thread with a small number of replies has ties with all others in the same thread, this assumption becomes problematic when the number of replies and people involved is very large. To address this problem, we propose a variant on Total Copresence which sets a threshold number of posts in the same thread (or subthread) that a participant is assumed to read to create a measure of Limited Copresence. The five tie definitions are summarized in Table 1. We note that true measures of "listening" based on clickstream data are an exciting area for future work, but that such information is currently infrequently available due to a lack of systems that log discussion forum access at the level of individual post read.

**Table 1** Tie Definitions

| Tie Definitions | Index of | Tolerance for Misplaced Postings | Risk of Over Assumption |
|---|---|---|---|
| Direct Reply | Speaking | Very Low | Very Low |
| Star | Speaking | Low | Low |
| Direct Reply + Star | Speaking | Low | Low |
| Limited Copresence | Speaking + Listening | High | High |
| Total Copresence | Speaking + Listening | Very High | Very High |

Despite the existence of these two very different classes of tie definitions (and variations within each), MOOC SNA studies have rarely provided explicit justification for the choice of definition used. More concerningly, little work has examined the impact that such decisions may have on the resultant networks (c.f. Gruzd & Haythornthwaite, 2008), or the implications of this for the comparability of results across studies using different definitions. As an initial step towards remedying this deficiency, in this study we examine content and non-content networks based on each of the five definitions shown in Table 1 to investigate the robustness of network properties and interpretation to differences in tie definition.

## 3 STUDY FRAMING

This study seeks to contribute to our understanding of learning through discussions in MOOCs by examining the interactions that occur and the social relationships that develop in a MOOC on statistics. Our approach to doing so begins by using SNA to examine the social networks that develop based on content and non-content discussions. SNA also helps direct our attention to noteworthy communities and individuals within the larger network whose discussions merit the intensive time needed to study them qualitatively. Inductive analysis of these actual discussion interactions is then used to (a) help explain the network property differences between content and non-content discussions; and (b) probe in-depth the ways students and instructors engage together around course content. This work thus follows an explanatory mixed methods design (Creswell, Klassen, Plano Clark, & Smith, 2011) that leverages the advantages of quantitative approaches for working with large scale data and the benefits of qualitative approaches for providing insights into the details of learning interactions. In addition, within this overarching frame, we examine the impact of the choice of tie definition on network properties and interpretation.

The specific research questions asked are as follows:

RQ1: To what extent are the properties of MOOC discussion forum social networks constructed using different tie definitions comparable?
RQ2: Do social networks based on content and non-content MOOC discussions show distinct properties at the overall network, community and individual level?
RQ3: What characteristics of the interactions in content and non-content MOOC discussions may explain differences in social network properties?
RQ4: What are noteworthy aspects of interactions that occur around the course content in MOOC discussions and the social relationships that develop as a result?

## 4 PHASE ONE: NETWORK CONSTRUCTION AND COMPARISON

### 4.1 Methods

*4.1.1 Data Source*
This study used data from StatMed'14, a completed MOOC offered in 2014 on Stanford open-source platform Lagunita. The course is an introductory course on probability and statistics with a special focus on statistics in medical studies. The course provided a discussion forum for interaction in nine topic areas,

including General, Video, Homework, Course Material Feedback, External Resources, Tech Support, Introductions, Study Group, and Platform Feedback. Learners were invited to post questions and comments for response by peers, the TA and the instructor. Forum information provided in the dataset included the following: thread id; post id; user id; post position in thread (thread starting post, reply post, or reply to reply post); parent post; post text; post creation date and time; and number of votes post received. Thread titles were not included. The discussion forum was participated by 568 unique users. They made 817 thread starting posts, 1,277 reply posts, and 1,035 reply-to-reply posts in the forums. Of the 817 threads, 117 received no reply. Five reply posts that contained non-English language or only punctuation were removed, leaving a total corpus of 817 threads with 2,307 replies made by 567 users.

*4.1.2 Thread Classification*
The 817 threads were classified as either being content or non-content using a unigram and bigram model built on manually-coded starting posts from a prior offering of the course (Authors, Blinded). In previous work, the model demonstrated good reliability on StatMed'14 data for both thread starting and reply posts (accuracy > .81, kappa > .62) (Authors, Blinded; Authors, Blinded). This model was used in conjunction with (Blinded) method (Authors, Blinded) to categorize threads by comparing the classification of thread starting post and distribution of reply classifications. This additional step increases the estimation of classification accuracy to .88 (Authors, Blinded). Using this comprehensive characterization method, a total of 468 threads containing 1,446 replies were labeled as content and a total of 349 threads containing 861 replies were labeled as non-content.

*4.1.3 Network Participants*
The nodelist was extracted from discussion forum data using user id of posts. It was found that of the 567 forum users extracted from the cleaned data, 178 participated only in content threads, 232 participated only in non-content threads, and 157 participated in both kinds of threads. Thus the number of nodes for the overall, content, and non-content networks were 567, 335, and 389 respectively.

*4.1.4 Tie Extraction*
Edgelists were extracted using five different tie definitions for the overall, content, and non-content networks.
**Direct Reply:** The author of each post was connected with the author of its parent post; this represents the actual reply structure (see Figure 1a). Using this definition, a total of 2,307 ties were extracted from the overall network; after removing 286 self-loops, 2,021 ties representing 1,086 unique edges remained, including 1,249 content ties representing 625 unique edges and 772 non-content ties representing 551 unique edges.
**Star**: The author of each reply and reply-to-reply post was connected with the author of the thread starting post (see Figure 1b). Using this definition, a total of 2,307 ties were extracted from the overall network; after removing 502 self-loops, 1,805 ties representing 1,116 unique edges remained, including 1,092 content ties representing 625 unique edges and 713 non-content ties representing 558 unique edges.
**Direct Reply + Star**: Ties defined in both Direct Reply and Star were included but the same tie was never counted more than one time (see Figure 1c). Using this definition, a total of 3,339 ties were extracted from the overall network; after removing 683 self-loops, 2,656 ties representing 1,292 unique edges remained, including 1,697 content ties representing 747 unique edges and 959 non-content ties representing 643 unique edges.
**Limited Copresence**: All users in small threads (5 posts or fewer) were connected to each other; in larger threads users were connected to all other users in their sub-thread and the thread starter only (see Figure 1d). The threshold of 5 replies was picked because the histogram for threads with different numbers of posts showed a big drop between 5-post threads and 6-post threads with only 128 (16%) of the threads being larger than 5 posts. Users were connected to all other users in a sub-thread regardless of subthread size because there was no basis for further division; however of the 489 subthreads in larger threads, only 69 (14%) had more than four posts. Using a VBA script written for this definition, a total of 5,313 edges

were extracted from the overall network; after removing 1,066 self-loops, 4,247 ties representing 1,456 unique edges remained, including 2,879 content ties representing 848 unique edges and 1,368 non-content ties representing 724 unique edges.

**Total Copresence**: All authors in the same thread were connected with each other (see Figure 1e). Using a VBA script written for this definition, a total of 15,299 ties were extracted from the overall network; after removing 1,992 self-loops, 13,307 ties representing 5,578 unique edges remained, including 7,018 content ties representing 1,133 unique edges and 6,289 non-content ties representing 4,641 unique edges.
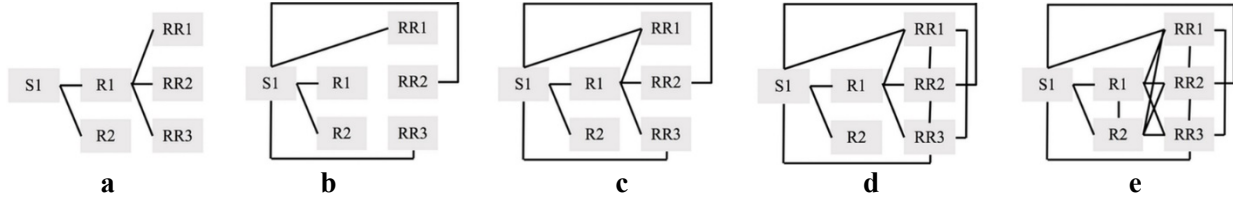


**Fig. 1** Ties based on five definitions (a) Direct Reply; (b) Star; (c) Direct Reply+Star; (d) Limited Copresence; (e) Total Copresence. S = thread starting post; R = reply post; RR = reply to reply post. Solid lines represent ties extracted using this definition.

*4.1.5 Network Construction and Network Properties*

The edgelists and nodelists for each network were imported into Gephi 0.9.1. for Mac. Undirected weighted networks were constructed and visualized using the Force Atlas layout algorithm. For each of the fifteen networks (overall, content, non-content for each of the five tie definitions), the number of edges, average node degree (the average number of neighbors that a learner interacts with), average edge weight (the average number of times that a learner interacts with the same neighbor), and graph density (the ratio of the number of edges to the number of possible edges) were computed. Community detection was performed through modularity maximization using the Louvain method (resolution = 1.0). Randomization was used to improve decomposition with an acceptable tradeoff in computation time; ten runs were conducted for each of the fifteen networks. For twelve of the fifteen networks the resulting major communities were consistent. For the Limited Copresence overall and non-content networks, and the Total Copresence content network, the runs produced two different community structures, one with two of the key nodes (u1 and u417) in the same community and one with them in separate ones. Each structure was produced multiple times, necessitating a principled choice about which to interpret. The decision was made to use the structure with u1 and u417 in separate communities in order to allow for better examination of potential differences between them and to facilitate comparison of parallel communities across tie definitions.
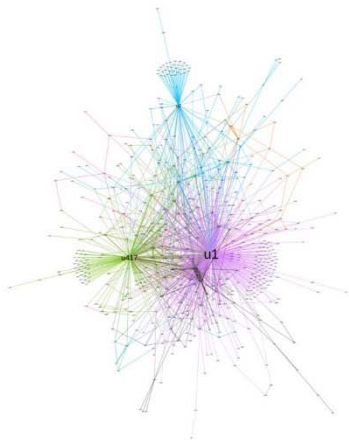
**4.2 Results and Discussion**

The properties of the fifteen network graphs are reported in Table 2 and the graphs are shown in Figure 2.
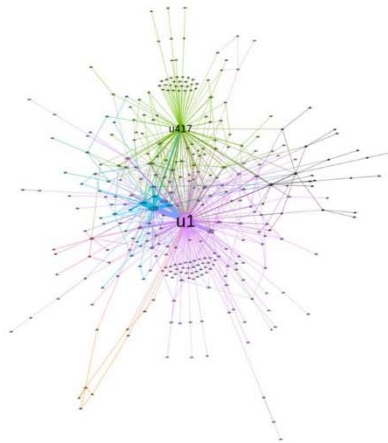
*4.2.1 Comparing Tie Definitions*

For overall, content, and non-content networks, there are clear trends across tie definitions in the number of edges with the more liberal tie definitions producing higher values (see Table 2). This shows an expected relationship between how ties are defined and the number of connections presumed to have taken place. Interestingly, while this trend is replicated in average node degree (more liberal tie definitions indicated a greater breadth of others with whom users were considered to have interacted), the average edge weights (indexing the strength of connection) based on the different definitions did not follow the same pattern: Total Copresence generated the *highest* average edge weight of all tie definitions for the content network, but it generated one of the *lowest* average edge weights of any definition for the non-content network. The network graphs extracted using the Total Copresence definition are dramatically different than those extracted using other definitions. Specifically, using the Direct Reply,
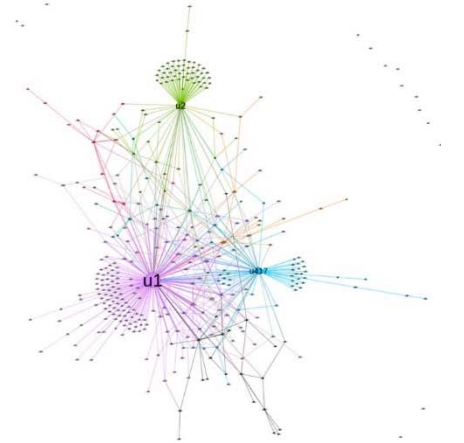
Star, Direct Reply+Star, and Limited Copresence tie definitions, all of the overall and non-content networks had three major communities, each dominated by a single node of high centrality (u1, u2, and u417, see Figure 2 a1-d1, a3-d3). In contrast, using Total Copresence definition, the overall and non-content networks extracted had only two major communities, one containing both u1 and u417, and the other a "balloon" of many similar-degree interconnected nodes (see Figure 2 e1 and e3). Examination of the post text contributed by u1 and u417 revealed them both to be members of the instructional team. Examination of the post-text data that contributed to the "balloon" showed that this was due to a single socializing thread started by learner u2 at the beginning of the course which received a total of 92 replies.
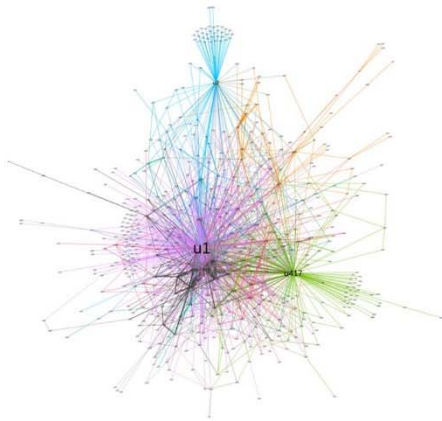
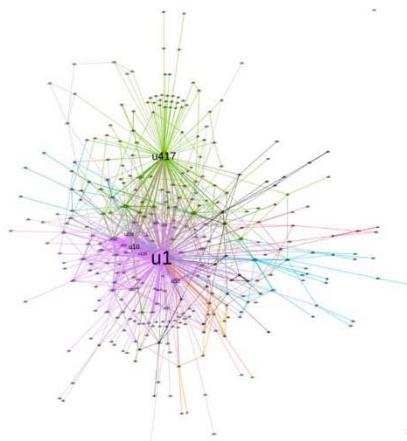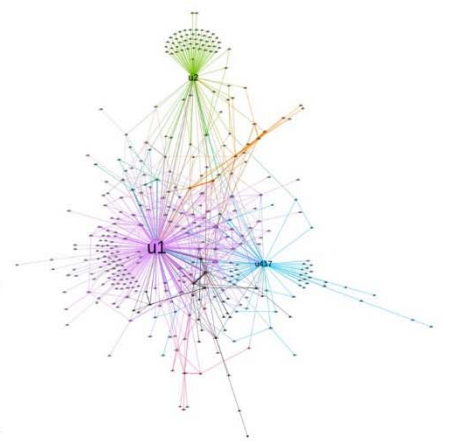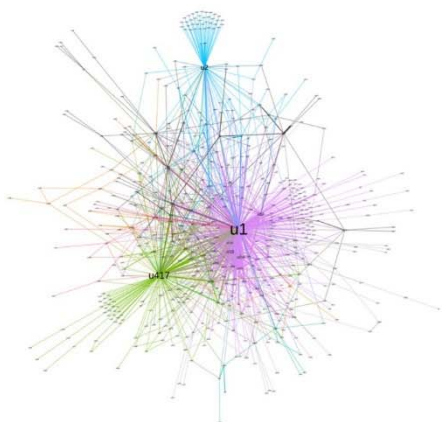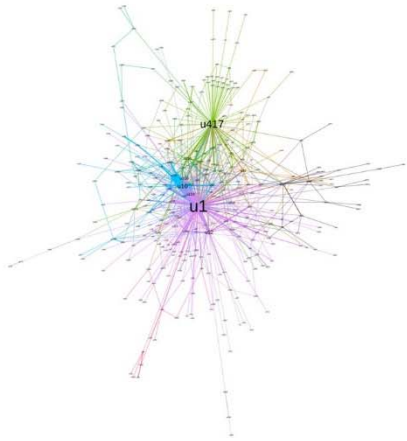**a1**

**a2**

**a3**

**b1**

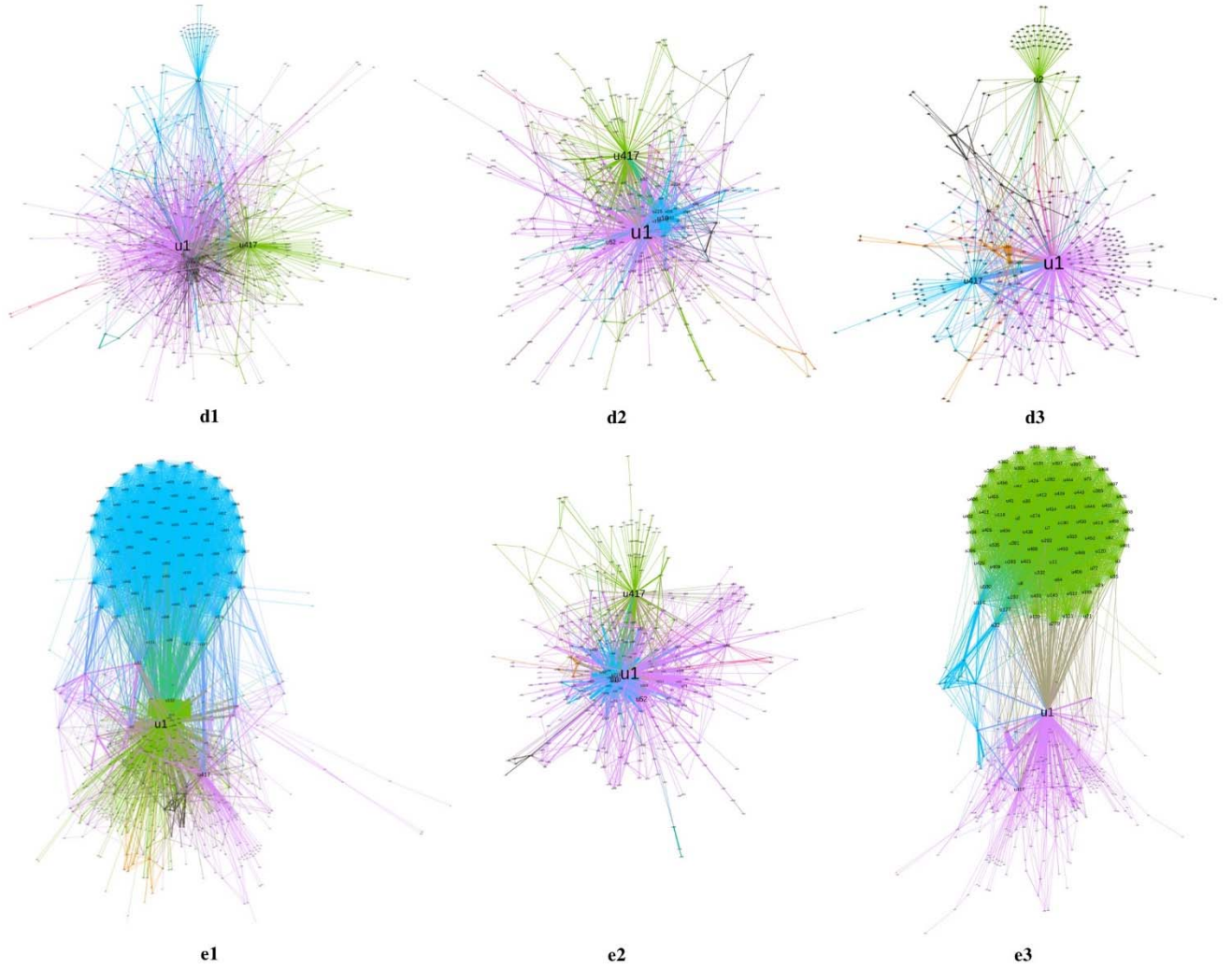**b2**

**b3**

**c1**

**c2**

**c3**

**Fig. 2.** Social networks constructed using (a) Direct Reply; (b) Star; (c) Direct Reply + Star; (d) Limited Copresence; (e) Total Copresence for (1) overall; (2) content; (3) non-content discussions. Only the primary components are shown in this figure. Node size represents degree. Color indicates community.

**Table 2** Network measures of five overall networks.

| | Overall (N=567) | | | | Content (N=335) | | | | Non-content (N=389) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | # of edges | Avg node degree (SD) | Avg edge weight (SD) | Graph density | # of edges | Avg node degree (SD) | Avg edge weight (SD) | Graph density | # of edges | Avg node degree (SD) | Avg edge weight (SD) | Graph density |
| DR | 1,086 | 3.83 (14.54) | 1.86 (3.38) | 0.007 | 625 | 3.73 (11.33) | 2.00 (3.62) | 0.011 | 551 | 2.83 (10.47) | 1.40 (1.16) | 0.007 |
| S | 1,116 | 3.94 (12.92) | 1.62 (2.50) | 0.007 | 625 | 3.73 (9.59) | 1.75 (2.89) | 0.011 | 558 | 2.87 (9.50) | 1.28 (0.79) | 0.007 |
| DR+S | 1,292 | 4.56 (15.48) | 2.06 (4.36) | 0.008 | 747 | 4.46 (12.04) | 2.27 (4.90) | 0.013 | 643 | 3.31 (11.20) | 1.49 (1.41) | 0.009 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LC | 1,456 | 5.14 (16.69) | 2.92 (9.49) | 0.009 | 848 | 5.06 (13.20) | 3.40 (11.10) | 0.015 | 724 | 3.72 (12.09) | 1.89 (2.75) | 0.010 |
| TC | 5,578 | 19.68 (36.63) | 2.39 (14.98) | 0.035 | 1,133 | 6.76 (15.64) | 6.19 (31.73) | 0.020 | 4,641 | 23.86 (38.06) | 1.36 (1.53) | 0.061 |

DR = Direct Reply; S = Star; DR+S = Direct Reply+Star; LC = Limited Copresence; TC = Total Copresence

When this thread was removed from the non-content network and the properties were recalculated, the average node degree was 4.40 (SD = 12.70), still the highest among all definitions, and the average edge weight was 1.85 (SD = 1.96), slightly lower than that of Limited Copresence. This suggests the Total Copresence tie definition can be problematic because one or two very large threads can dramatically inflate average node degree and deflate average edge weight dramatically.

Looking globally, differences between Direct Reply and Star definitions were smaller than expected, both in terms of network properties (see Table 2) and graph appearance (see Figure 2). Direct Reply+Star and Limited Copresence, although based on different conceptions of tie definition, appeared to produce networks with similar properties (see Table 2). These findings suggest that social network analysis is more robust to tie definition than was expected, with the exception of Total Copresence.

*4.2.2 Comparing Content and Non-Content Networks*
Content and non-content networks were found to have several distinct characteristics. First, both average node degree and average edge weight were higher in the content network than in the non-content network for all tie definitions except Total Copresence (see Table 2). In Total Copresence, the non-content network had a dramatically higher average node degree and lower average edge weight than the content network. However, the average node degree was inflated and the average edge weight deflated by the thread with 93 posts from 88 participants in the non-content network. When this thread was removed, the recalculated average node degree on the Total Copresence non-content network became lower (4.40, SD = 12.70) than that of the content network. The recalculated average edge weight rose somewhat from the original value (1.85, SD = 1.96) but was still lower than that of the content network. These cross-definition trends indicate that learners interact with more people and have more repeated interactions with the same people in content discussions than in non-content discussions.

Second, the two kinds of networks showed different community patterns (see Figure 2). All of the content networks contained two major communities, each dominated by one of the instructors (u1, u417) with substantial connections between the two communities. In contrast, all of the non-content networks (except Total Copresence) showed three major communities (centered around instructors u1 and u417 and learner u2) with fewer links between them. Even when the 93 post thread was removed from the non-content network, this pattern still held, with another learner-only community taking the place of the balloon.

**5 PHASE TWO: IN-DEPTH COMPARISON OF COMMUNITIES AND INDIVIDUALS IN CONTENT AND NON-CONTENT NETWORKS**

**5.1 Methods**

The second phase of analysis compares the interaction characteristics of communities and individuals in the content and non-content networks of the MOOC discussions. Based on the analysis in phase one, it was noted that networks constructed using all three reply-based tie definitions and the Limited Copresence tie definition had similar properties. In phase two, the Limited Copresence was used to construct the networks, as conceptually it encapsulates both speaking and listening activities but avoids the disproportionate influence of large threads. It is noted that these results are comparable to those obtained via Direct Reply.

*5.1.1 Identification of Major Communities*
The distribution of community size showed that in both the content and non-content networks, there was a sharp drop in community size after the first few communities followed by a large number of micro-communities each containing less than 5% of the total participants. The current study focuses on examining the major communities. It can be interesting to investigate the smaller ones in the future.

All communities that contained more than 5% of the total participants in either the content or non-content network were extracted for further investigation. From the content network three communities were extracted; from the non-content network, four communities were extracted. Communities are referred to in shorthand as C1(I), C2(I), C3 and N1(I), N2(I), N3, N4, where (I) indicates that the community included an instructor.

*5.1.2 Characterization of Major Communities*
For each community the number of nodes, number of edges, average node degree, average edge weight, and graph density were calculated.

All threads that contributed to the connections in each community were extracted and examined. Threads were analyzed intact to maintain the context of interaction for interpretation. For each community, the following were calculated: number of contributing threads, average thread length (number of posts per thread), average number of subthreads per thread, and average subthread length (number of posts per subthread).

Threads were then analyzed qualitatively. Given the large and varied number of threads for C1(I), C2(I), and N1(I), stratified random sampling was conducted to produce similarly sized samples for each community. Strata were formed based on number of posts in threads; then threads were randomly sampled within each stratum as follows: 41 threads for C1(I); 35 threads for C2(I); 36 threads for N1(I) (exact numbers differ due to variations in thread length). Threads contributing to the other communities were small enough to be analyzed in their entirety: 30 threads in C3; 38 threads in N2(I); 2 threads in N3; 11 threads in N4. Threads were manually checked to verify if they had been properly categorized as content / non-content using the (Blinded) method. Only 9 out of 193 threads were determined to be miscategorized (6 originally labeled as content and 3 non-content). These threads were not included in subsequent analysis.

The posts in each thread were then analyzed qualitatively using inductive thematic analysis following the constant comparative method (Auerbach & Silverstein, 2003; Gibson & Brown, 2009). The purpose such analysis was to identify emergent themes and patterns through probing the characteristics, similarities and differences between interactions in the communities from the content and non-content networks. Threads were first color coded by participant id to aid interpretation and then a description was written for each thread that includes: (a) the expressed purpose for the thread initiation; (b) the overall dynamics of replies that resulted; (c) any other purposes for the thread that emerged through discussion; and (d) if the purpose(s) were eventually fulfilled. In addition, predominant interaction techniques (e.g., providing factual information, sharing personal understanding, asking leading questions) and social presence indicators (e.g., greetings, addressing people by name, see Authors, Blinded) were noted for each thread. These descriptions were then examined for emergent themes across threads contributing to each community. One researcher conducted the initial analysis and interpretation; descriptions and proposed themes were examined and discussed with a second researcher leading to revision and refinement. Care was taken to ensure the trustworthiness of analysis via triangulation, referential adequacy, and creation of an audit-trail (Guba, 1981).

*5.1.3 Identification and Characterization of Central Learners*

Degree was calculated for all learners in the content and non-content networks. Learners that ranked in the top 10 for any of the two networks were selected for further investigation, including examination of which network(s) they participated in (content, non-content, both) and whether they also ranked in the top 10 list for the other network.

*5.1.4 Characterization of Instructors in Content Communities*

To investigate the instructors' participation in learning activities directly related to the course content, the post texts contributed by the instructors in the content threads were examined qualitatively using the same inductive thematic analysis method described in Section 5.1.2. The purpose was to identify themes and patterns for the instructor's participation techniques (e.g. providing straight forward answers, giving hints, asking leading questions) and their use of social presence cues, through probing the characteristics, similarities, and differences between the instructors' participation.

**5.2 Results and Discussion**

*5.2.1 Communities*
**Network Structure**
Several differences between content and non-content networks were found at the community level. First, there was a difference in the size of instructor communities: C1(I) and C2(I) together contained 77% of all nodes in the content network; while N1(I) and N2(I) together contained 55% of all nodes in the non-content network. Thus the content network was more dominated by interactions that happened around the instructors.

Second, there were differences in the structural properties for communities from the two networks. For the instructor communities, average node degree in C1(I) and C2(I) were both higher than those in the corresponding instructor communities in the non-content network (N1(I) and N2(I), see Table 3). These differences indicated that these learners interacted more broadly with different peers. The structures of all the communities, however, was still a hub-and-spoke centered around the instructor (see Figure 3). For the learner-only communities, C3 had twice the number of edges to the comparably-sized N4, leading to a higher average node degree and density; average node degree in C3 was also greater than that in the larger-sized N3. Differences in patterns of connectivity were also seen in the network structures between C3 (highly-connected web), N3 (hub-and-spoke), and N4 (elongated chain) (see Figure 3).

In addition, looking at the strength of connection, the average edge weight in C3 was dramatically higher than that in all other communities (see Table 3), indicating that learners had substantially more repeated interactions with the same peers. Edge weights were highest among a central clique of five learners (see Figure 3).

**Table 3** Network structures of major communities (ties based on Limited Copresence).

|  | C1(I) | C2(I) | C3 | N1(I) | N2(I) | N3 | N4 |
|---|---|---|---|---|---|---|---|
| # of nodes | 184 | 75 | 23 | 168 | 47 | 62 | 23 |
| # of edges | 400 | 105 | 57 | 315 | 55 | 71 | 28 |
| Graph density | 0.024 | 0.038 | 0.225 | 0.022 | 0.051 | 0.038 | 0.111 |

| Avg node degree (SD) | 4.35 (11.06) | 2.80 (7.56) | 4.96 (4.08) | 3.75 (11.18) | 2.34 (6.03) | 2.29 (7.44) | 2.44 (2.39) |
|---|---|---|---|---|---|---|---|
| Avg edge weight (SD) | 2.23 (3.21) | 1.83 (1.72) | 14.67 (32.61) | 2.11 (2.48) | 1.20 (0.44) | 1.06 (0.29) | 1.89 (1.63) |
| Instructor degree | 145 | 67 | NA | 144 | 43 | NA | NA |



**Fig. 3.** Comparison of (a) C1(I); (b) C2(I); (c) C3; (d) N1(I); (e) N2(I); (f) N3; (g) N4 (ties based on Limited Copresence).

**Thread structure**

Differences were found in the number of threads contributing to the communities from content and non-content networks. All content communities had more contributing threads than the corresponding non-content communities (see Table 4). This difference indicates that interactions in the content communities were more distributed across the discussion forum. Moreover, differences were also found in thread and subthread size. Threads contributing to C3 had more than twice the number of posts with relatively longer sub-threads compared to N3 and N4; this trend was also present, though less dramatic, in the instructor communities (see Tables 4).

**Table 4** Structure of threads contributing to major communities (ties based on Limited Copresence).

|  | C1(I) | C2(I) | C3 | N1(I) | N2(I) | N3* | N4 |
|---|---|---|---|---|---|---|---|
| # of contributing threads | 162 | 70 | 30 | 137 | 38 | 2 | 11 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Avg thread length (SD) | 5.26 (5.74) | 3.37 (1.44) | 13.40 (14.85) | 4.07 (2.75) | 2.68 (0.76) | 6.00, 93.00 | 6.00 (5.75) |
| Avg # of subthreads per thread (SD) | 2.07 (1.49) | 1.40 (0.55) | 2.73 (2.05) | 1.77 (1.17) | 1.34 (0.62) | 5.00, 82.00 | 3.00 (3.41) |
| Avg subthread length (SD) | 2.64 (1.29) | 2.54 (1.02) | 4.80 (2.76) | 2.48 (1.22) | 2.18 (0.72) | 1.20, 1.13 | 2.57 (1.51) |

\* This community had only two contributing threads of very different characters, thus results are reported for each thread separately.

**Table 5** Average number of post made by participants in contributing threads for major communities (ties based on Limited Copresence).

| | C1(I) | C2(I) | C3 | N1(I) | N2(I) | N3 | N4 |
|---|---|---|---|---|---|---|---|
| Average post per participant (SD) | 1.39 (0.58) | 1.29 (0.31) | 2.41 (1.62) | 1.30 (0.43) | 1.10 (0.19) | 1.13 (0.07) | 1.42 (0.53) |

Second, participants in the threads contributing to C3 made more posts per thread than those who posted to contributing threads for other communities (see Table 5). This finding indicates that participants in C3 revisited the threads more often.

**Thread content**
During the qualitative analysis, three themes emerged associated with the differences in information seeking and giving behaviors in content and non-content communities.

***Theme 1: Discussion questions and responses in the content network were more complex.***
The kinds of questions asked in content versus non-content threads were different, which may explain differences in thread and subthread length. Non-content starting posts often asked for straightforward factual information which could be easily provided without extended conversation (i.e. short threads). For example:

> *From community N4*
> *U48: I included the % sign by mistake in the answer box.… Is there any possibility to counteract this problem? Please...*
> *U32: I remember in the introduction to the course, there have been a note that you can't amend your answer after submission.…*

> *From community N2(I)*
> *U143: Will it be possible to get the lecture slides being used in the background?*
> *U417: They are available under the overview/teaser section for each unit.*

In contrast, starting posts in content threads commonly asked for help with problem-solving or understanding complicated concepts which required multiple rounds of back-and-forth comments to resolve (i.e. long threads). In these conversations, it was common for participants to use diverse interaction techniques such as (1) paraphrasing or clarifying the question, (2) giving explanation, examples, and comparisons, (3) asking follow-up questions, and (4) using leading questions. These techniques were used both by instructors and some learners. For example:

> *From community C1(I)*
> *U10: Hello [instructor u1] and classmates.… I do not know how to parse out the twin data or understand the cases / controls and who was exposed in order to do the statistical test I think is*

*correct. I wanted to make a 2X2 table, and obviously I don't understand the 1X4 table in the question.*
*U1: Hello [u10], do you know what type of test is best to analyze this type of information (twin studies)? … What is the best way to set up a 2X2 table from twin pairs? Hint: the +/- next to the case and control twin pairs are the exposures.*
*U10: OK, I thought about it some more. Does this look right? Am I thinking of the right test? (shared a link to his/her solution)*
*U1: Dear [u10], big hint: google "twin pair 2X2 table"…. They have an example of how this is set up! … Also check out module 4…. I think once you have done this, you can easily solve the problem!*
*U152: Hi [instructor u1], does this look correct to you? (shared a link to his/her solution)*
*U1: Nice job! Do you understand why the table should be this way?*
*U152: Took me a while but yeah, now I do! Thanks!*

**From community C3**
*U209: Can anybody help me with question 10 of unit 4? Do we have to consider the mean = proportion = 112/200 = 0.56?*
*U167: Good morning, the question states that you should use the normal approximation to the binomial. I would go back and look at slides 139-141 for this unit and double check the equation that you're using for the mean. The mean is not a proportion, it is =n* p. The problem wording gives you both of the values for those variables, You just need to plug them in!*
*…*
*U209: Thanks, but I'm still confused. Don't we have to use the statistics of proportion here? 112/200 =0.56 and if I'm using the formula mean= n*p, and X = 112, then the z score is coming to zero. Does that make any sense?*
*U502: P of flip a coin is 0.5, X=112, mean(u)=n*p=0.5*200. You can calculate SD using sigma^2=np(1-p) and z=(x-u)/sd, and use Standard Normal Distribution Table.*
*U10: Hint: Don't forget to take the square root in the denominator when solving for the SD (a mistake I made!).*

***Theme 2: Conversation forms in the content network were more complicated.***
The qualitative analysis also revealed differences in the form of conversation between content and non-content communities. Non-content discussions usually involved a single topic; when sub-threads were used, the structure did not seem to indicate anything substantive about sub-topics of the conversation. In contrast, many content discussions contained multiple subthreads devoted to different subtopics. This may help explain the overall greater length of threads contributing to the content communities.

Three variations of structures for content conversations emerged during the qualitative analysis.

First, many starting posts in C3 were expansive in nature, either making a general call for discussion or asking about several homework questions together. For example,

> *U10: Has anyone begun on this homework? I'd love to start some threads like we had going in homework 7. Thanks!*

Responses to such starting posts consequently often broke into separate sub-threads focused on specific questions or topics.

Second, the starting posts of some content threads raised specific questions that could be addressed in a single subthread, but the participants initiated new subthreads to address the topic with different

perspectives / approaches. For example, the following excerpt shows a starting post with four replies that each initiated a new subthread:

> ***From community C3***
> **U288: (Starting post)** *In the modules it discussed clinical significance is largely based on what is being studied and what clinical affect it has…. For the list of relative risks on this question we don't have any of that information so it is difficult to determine at what cut-off the risk ratios are clinically significant. Am I missing something?*
> **U10: (Subthread 1)** *Hi [u288], I think that we have the information we need to make an informed decision here. … a value of 1 means no effect so if something has an effect above and below 1, it isn't clinically significant...*
> **U536: (Subthread 2)** *I don't really think there are cut-off values, but if the relative risk is extremely close to one, there is just a tiny extra risk which is probably clinically insignificant.*
> **U410: (Subthread 3)** *The way I approached the problem was this: forget about everything we've done in this class. Forget about confidence intervals. Just look at the main result…. The goal with the tests is to reject the null hypothesis. Which results can 'confidently' reject the null?*
> **U52: (Subthread 4)** *Probably you have to look to the effect of the Drug…. Here it would be clinical significant, however statistical significant does not care about this point.…*

Third, in some content conversations, after the specific questions raised in the starting posts were addressed, the participants initiated new subthreads to ask follow-up or other relevant questions. For example, the following excerpt shows a starting post with the one reply starting a subthread to address the question and a later one starting a subthread on a related topic:

> ***From community C2(I)***
> **U258: (Starting post)** *Hello, at approximately 6:50 into the video, the p value is declared to be 4%. Can anyone explain how this was calculated? Thanks.*
> **U197: (Subthread 1)** *You can calculate it using the binomial probability …. Another way is by using the Z scores …. Thus, the two-tailed probability is 4%.*
> *…*
> **U52: (Subthread 2)** *How did you calculate standard error 11.1?*

### Theme 3: Interactions in community C3 had unique characteristics.

Two distinct characteristics were found in community C3. These also found, but to a lesser extent, in C1(I), and rarely seen in any of the other communities.

First, in community C3, learners often revisited the conversations. In other communities, once the answer/solution was provided, the interaction usually stopped. In contrast, in many of the threads in C3, and some in C1(I), the thread initiator revisited their threads to interact with peers who responded to them. In addition, after the initial questions had been answered, learners often made subsequent posts to provide help for others, which often led to multiple rounds of talks. This may explain the higher average number of posts made by participants in the threads contributing to C3. For example:

> **U110:** *Hello [u10], For Question 8…what model was used to generate this beta coefficient? My logic…. In the equation I should quantify the breast feeding score…. For 9-11 I saw [the instructor's] videos …. Do you think I am on the right track? For 3 I actually made a pseudo data …. Do you think I am using the right logic? Best Wishes.*
> **U10:** *[u110], Q8: Are you asking if Y = visual reception and then what is the equation to estimate Y …. For the equation problems…I went back to the slides …. Q9-11 … I would just say to review an example …. Q3 Probably module 3 is the best one to use for this question…when I was*

*considering the question, I looked closely at .… I think you can answer the question by looking at the graph. Good luck and feel free to ask me more questions :).*
*U110: Thank you [u10]. :) … we can only identify categorical variables as binary if they have been dealt as binary right?… Best Wishes.*
*U10: Hi [u110], We create dummy codes when there are … but I believe that a binary variable is.… Actually - I think I just reworded what you said!*
*U110: So [u10], does Linear Regression always handle categorical predictors as binary?*
*U10: [u110]… we have to create the new variables if they have more than 2 categories …. I think it is also important, here - for the homework.…*
*U110: [u10] Thank s for the advice. I am using your advice and Community 5 to make a decision.*

Second, many excerpts from the C3 communications show signs of a budding learning community. In many cases, participants expressed a desire to work with others in the course and used indicators of social presence (such as including a greeting, sharing feelings and using personal pronouns):

*U299: Hi everybody, I am enjoying this course so much, I can't wait to apply my newly acquired stat skills to my work. However, I ran into an unexpected complication and would like to summon your help (instructors and participants) to sort it out…*

Moreover, a conversation among the learners connected by high-weighted edges in C3 towards the end of the course suggests the development of interpersonal bonds through participating in learning-related discussions:

*U10: I am done! I missed one we didn't even discuss …, but anyways, OVER! And so is the exam... sigh... this was the toughest MOOC I have taken - grading wise.*
*U225: Congrats [u10]! Yes, it has been hard, but fun, and we learned an awful lot, right?*
*U110: Great! Everyone it was a pleasure to work with you. Thank you.…*
*U10: YES [u225]! And [u110] - the test was scary - I thought of my discussion board friends often!!*
*U216: Thanks, thanks so much to [u10], [u152], [u110], [u225] and everybody who helped us to understand this beautiful course! And in my case also for writing many posts, I see I have improved my English skills and my statistics vocabulary!!!*
*U225: [u10], [u216], [u152], [u110], [u515] and everyone, your discussions helped me so much. I was always a few days behind you in homework - glad I was able to catch up in the last weeks and participate a little bit.…*

In contrast, the very large thread in N3 initiated by u2 for learners to make self-introductions did not fulfill its potential for building social connections. While the thread was joined by 87 other participants, many of whom expressed interest in interacting with others, most participants just made a monologue self-introduction and never participated in this thread again:

*U32: Hi everyone, I'm a doctoral student … from Egypt. My specialty is .… I'm looking forward to complete this course, hopefully with distinction. Hope to be friend with all of you guys. Best wishes for all.*

Only 5 of 88 participants posted in this thread more than once. This indicates that minimal social interaction actually took place in this (purely) socially-oriented thread. Furthermore, over half of the contributors to this thread (56%, N=49) never took part in any discussions about the course content and 42% (N=37) never participated in any discussion at all ever again.

*5.2.2 Learners with high degree*
There were 17 distinct learners in the top 10 lists (ranked by degree) for the content and non-content networks. Except for u10, u216, and u21 who appeared on the high degree lists for both networks, the remaining 14 learners had high degree in one network or the other but not both (see Table 6). This indicates that top players in the two networks were largely different people and that even if participating in both content and non-content threads, learners who were highly connected in one network were not necessarily highly connected in the other.

**Table 6** Top 10 learners ranked by degree centrality in content and non-content networks (ties based on Limited Copresence).

| Rank | Content network User (Degree) | Non-content network User (Degree) |
|---|---|---|
| 1 | u10 (54) | u2 (87) |
| 2 | u52 (46) | u73 (22) |
| 3 | u216 (27) | u79 (22) |
| 4 | u32 (26) | u225 (16) |
| 5 | u110 (24) | u21 (15) |
| 6 | u236 (24) | u56 (15) |
| 7 | u152 (21) | u216 (15) |
| 8 | u60 (19) | u10 (14) |
| 9 | u46 (19) | u23 (13) |
| 10 | u21 (19) | u30 (13) |

Shade indicates user appears on both lists.

*5.2.3 Instructors in the Content Communities*
**Network structure**
Social network analysis revealed differences in interactions around the two instructors. First, C1(I) had 2.5 times as many nodes and 4 times as many edges as C2(I) (see Table 3), indicating that the community formed around u1 contained more people and interactions than that formed around u417. In addition, the degree of u1 was more than twice that of u417 (see Table 3), indicating that u1 interacted directly with many more learners than u417. Second, differences were also found in the community graphs: in both communities, a large proportion of learners were only connected to the instructor and not to any learner in the community (see Figure 3a and 3b), indicating that a large number of learners only had interactions with the instructor. Comparatively, there were more interconnections among learners in C1(I) than in C2(I). This is also indicated by the finding that the average node degree in C1(I) was 1.5 that for C2(I) (see Table 3). These characteristics indicate participants in C1(I) interacted with more other learners than those in C2(I). Moreover, the average edge weight in C1(I) was somewhat higher than that in C2(I) (see Table 3), indicating that participants in the former community had more repeated interactions with the same other learners than those in the latter community.

**Participation pattern**
Posts made by the two instructors revealed distinct characteristics of instructors' forum activities. First, the two instructors had different patterns of posting. Instructor u1 made 353 posts in 240 content threads, including 216 replies to thread starting posts and 137 replies within subthreads. This indicates that this instructor not only revisited the threads that he/she had participated in, but also commented on other learners' replies to learner-initiated threads. In contrast, instructor u417 made 121 posts in content

threads, all which were direct replies to starting posts. This indicates that u417 only addressed the subset of all leaners that initiated threads.

**Communication techniques**
Qualitative analysis showed the instructors also used distinct communication techniques. Instructor u1 often tried to encourage and help learners to work out the answer or solution themselves. For example, he/she often gave hints instead of answering a question outright:

> *U1: Looks like you are making great progress! ... You are correct about dealing with categorical data, and the observations are certainly correlated! Think about it again using the hint and let me know if you have any other questions?"*

Instructor u1 also used leading questions to help learners work through problems and figure out solutions themselves:

> *U1: That is correct - Nice! So how would you use this to solve the question?*

In addition, u1 used a variety of social presence indicators such as greetings and addressing learner's by name in his/her messages:

> *U1: Hello [learner' s name], do you think you could clarify your example data set a little more?*

In contrast, instructor u417 tended to provide straightforward answers or instructions to address learners' questions and used social presence cues infrequently. For example:

> *U417: A bell shape is not necessary. You could have a 'bimodal' distribution (with two distinct peaks in the distribution) where the two groups do not follow a bell shape.*

## 6 GENERAL DISCUSSION

### 6.1 Differences in Content and Non-Content Networks and Activities
In this study, we found the content and non-content discussions and networks to show distinct characteristics. First, in the content network, participants interacted with more people and developed stronger ties with the same people than in the non-content network. This expands findings reported by Gillani et al. (2014) for the networks formed in two successive offerings of a business MOOC (using the Total Copresence tie definition). In that study the sub-forum for discussing course materials contained the strongest social connections (low proportion of one-off ties) while the sub-forum for conversations related to technical support and expressions of gratitude contained weaker connections (high proportion of one-off ties).

Second, content conversations seemed to be of greater depth than non-content ones, indicated by the topics involved and the greater number of contributing posts. Previous studies have suggested that content related activities are more useful for learning (Romero, López, Luna, & Ventura, 2013; Wang et al., 2015), which may be related to the deeper discussions in the content threads, the stronger connections formed during the discussions, or both. For instance, Romero et al. (2013) found that various features of learners' content-related contributions (including outdegree centrality) were more useful than features of non-content contributions for predicting learning outcome in a computer science course. Wang et al. (2015) found that the quantity of students' on-task discourse was a significant predictor of their learning gains in an introductory psychology MOOC. Although causal relations cannot yet be claimed, the association between learning outcomes and content-related interactions is a clear direction for future study.

Differences in participants' behavior in the content and non-content networks may be explained in two ways. First, the differences may have to do with who chooses to participate in which network. In this study the two networks were populated by largely distinct people. Only 28% of all forum participants engaged in both content and non-content discussions. A similar observation was made by Gillani and Eynon (2014) who found that in a business MOOC where different sub-forums were designated for content and non-content discussions, cross-sub-forum participation was rare. In addition, in the current study learners who gained high degree centrality in the content and non-content networks were largely distinct people. Therefore, it is possible that the differences observed in the two networks were to some degree due to intrinsic differences in the participants, such as their motivation for participation or learning goals. A similar result was found by Dowell et al. (2015) in a MOOC on infrastructures where learners with higher performance and or higher social centrality were associated with different kinds of linguistic features. Thus it is possible that achievement-focused learners were more likely to engage in learning-related discussions while others were more interested in purely social connections in the forums. Second, the topics of discussions in the two contexts may have affected participant's behavior. The content discussions involved more complicated issues, such as problem-solving or understanding complicated concepts, which often took several rounds of talks to resolve and involved contributions from multiple learners. In this process the learners not only interacted with the same people repeatedly, but were exposed to opportunities to interact with additional other peers. In contrast, the non-content discussions usually involved short exchanges of factual information, which were accomplished straightforwardly with input from a smaller number of learners. The hypothesis that the topic of conversation drove some of the differences observed is supported by our finding that learners who gained high degree centrality in both networks made longer posts and participated in a greater number of threads in the content discussions than in the non-content ones.

### 6.1.1 Implication: Learning in MOOC Discussions as an Integrated Process
Prior work on MOOC discussion forums has considered their role either in fulfilling a social purpose that supports motivation or an informational one that helps students develop understanding (e.g., Jiang et al., 2014; Xiong et al., 2015), but the specific mechanisms through which this occurs have not yet been examined. This study contributes to our knowledge of these processes by documenting the extended rounds of explanations, clarifications and follow-up questions through which questions are answered and social connections are built. There are two aspects of this finding that are noteworthy. First, prior work on question answering in MOOC forums has focused on the provision of answers either as a problem of *commodity* (pointing a learner to where the information for answering the question can be found, e.g., Agrawal et al., 2015) or one of *source* (identifying which learners are well-equipped to provide answers, Yang, Adamson, and Rosé, 2014). This work suggests that there is value in taking the perspective of question answering as a *process* in which understanding is developed as learners engage in a variety of interactions such as (a) clarifying what is being asked, (b) giving explanation, examples, and comparisons, (c) raising follow-up questions that arise based on the conversation, and (d) using leading questions to help others figure out answers themselves. Second, the social purpose of MOOC discussions has often been discussed in isolation from the informational one; that is the motivation that results from knowing there are other people involved in the course is disconnected from whether a learner actually engages with others in the course about the content. This study found both that rich interpersonal connection-building occurred in the context of discussing course content, and that purely social threads don't always foster it. The implication of this finding is that the development of content understanding and social connections need to be considered together, rather than separately.

### 6.1.2 Implication: Encouraging Content-Related Discussions
The finding that content discussions contained deeper and wider interactions can inform MOOC course design. Courses that aim to promote learner connections are well advised to go beyond simply providing forums as an open space for interaction and to think consciously about how to invite and support

particular kinds of content-related conversations. One way to do this could be to embed "challenge questions" in the video lectures explicitly designed to stimulate learner discussion. This could be supported with a guide to support the process of question answering through the processes described in 6.1.1. Another approach could encourage learners to use their distinct expertise to support each other. For example, in a management course, video lectures and readings could include examples from different industries and learners from the relevant fields could be invited to take on facilitation roles. In addition, course designers can consider how to help learners connect with others about specific aspects of the content; for example, building on Yang et al.'s (2014) question recommendation system (which matched learners with people likely to be able to answer their question), a tool could be developed to suggest connections between students with similar interests and provide them with prompts and process through which to engage.

### 6.1.3 Implication: Differentiating Discussions for a Clearer View of Interaction and Relationships
These findings support the value of content-based differentiation of discussions in MOOC research. Such differentiation (separating content and non-content discussions prior to analysis) has multiple benefits. First, as the two kinds of discussions can be participated in by largely different people, content-based differentiation can help researchers to identify more precisely particular kinds of individuals or subgroups. For instance, for studies that want to identify community TAs (Papadopoulos, Sritanyaratana, & Klemmer, 2014) or influential learners for disseminating instructor information (Jiang, Zhang, Liu, & Li, 2015), it is useful to distinguish learners with high social status in the content network from those who are only prominent in the non-content network. Similarly, marginalized learners in the two networks may have different needs for support and intervention (Brugha & Restoule, 2016). Second, as content-related and non-content activities have been found to associate differently with retention and learning outcomes (Kuh, 2002; Romero et al., 2013; Wang et al., 2015), content-based differentiation can enable researchers to engineer variables that are most useful for their research purpose. For instance, network-level and node-level properties in the content network are expected to be more useful for studying relationships between peer connections and learning outcomes (Wise & Cui, 2018).

## 6.2 The Role of Instructor Participation in MOOC Discussions

Instructors' participation in discussion forums is perceived as an important factor contributing to quality in online learning by both instructors and learners (Dennen, Aubteen Darabi, & Smith, 2007; Hew, 2015). However, the ways in which it is important and the expected impact on the discussions and resultant learning is less agreed upon. Some people see an instructor's presence in the forum as important for student motivation (Baker, 2010), while others see it primarily as a source of quality information (Jiang et al., 2015). One common argument claims that instructor participation should stimulate both the quantity and quality of student contributions (Bangert, 2008; Dennen, 2005); however the evidence supporting this is mixed. While some studies have found evidence that instructor involvement is positively associated with the quality and quantity of student contributions (Brinton, Chiang, Jain, Lam, Liu, & Wong, 2014; Dennen, 2005), others reported null or negative associations (Mazzolini & Maddison, 2007; Tomkin & Charlevoix, 2014) in which instructor involvement may actually shut down conversation (for example the authority of the instructor may be interpreted to provide a definitive answer that cannot be questioned or further explored).

This study contributed two findings about the relationship between instructor and learner participation. First, learner interactions in the MOOC were overwhelmingly instructor-centric: the two instructor communities in the content network contained 77% of the total participants and 60% of the connections in the network. Both instructor communities showed hub-and-spoke structure with the instructor being the only central node and relatively limited learner-learner interactions. A similar hub-and-spoke structure (constructed via ties of direct reply) was reported in Brooks et al. (2014), showing the dominant role in communication an instructor can unintentionally play due to their role and status in the course.

However, the current study also showed that the two instructors engaged very differently in the forums, and this was associated with corresponding differences in learners' participation and the resultant network structure. This indicates that to understand the impact of instructor activity on learner participation in discussion forums, it is important to consider not just if or how much the instructor participates, but the ways in which they do so. In this study, instructor u417 only responded to thread starting posts, provided straightforward answers, and seldom used social presence cues; the associated network community showed few learner-learner connections. In contrast, instructor u1 responded to posts on all levels, helped learners to work out the answers by providing hints and asking leading questions, and frequently used social presence cues; learners in this community showed stronger ties with a greater number of peers. We hypothesize that the difference in learner-learner connections is explained by the fact that u1's facilitation resulted in more opportunities for learners to interact with other peers in prolonged discussions. This may help to explain Mazzolini and Maddison's (2007) finding that instructor participation did not stimulate learner discussion, as almost 70% of the instructors' forum postings in that study consisted of providing answers to questions. Similarly, in Tomkin and Charlevoix (2014), instructor forum participation was focused on motivating learners through providing positive feedback, but comments did not create a substantive need (such as a question) for learners to respond. While direct benefits on learning outcomes were not assessed in this study, we note that u1's activities and the associated learner activity align more closely with constructivist theories of learning and associated pedagogies that see the role of the instructor as 'guide-on-the-side" who supports students in learning how to figure out answers.

*6.2.1 Implications: Informing and Supporting Instructor Participation with Network Analytics*
The finding that instructors' approaches to forum participation are associated with different social network patterns offers an exciting opportunity to use timely learning analytics to provide actionable feedback to instructors (Ifenthaler, Adcock, Erlandson, Gosper, Greiff, & Pirnay-Dummer, 2014). Literature shows that instructors do not always make accurate judgements about their teaching activities. For instance, Mazzolini and Maddison (2007) found although most instructors indicated that they often combined answers with follow-up questions to promote constructivist learning in the forums, only 12% of their postings actually contained follow-up questions. The benefit of providing learning analytics to instructors was shown in Brooks et al. (2014): when the instructor saw a network diagram indicating that the discussion forum activities were highly instructor centric, she changed the assignment to be more problem-solving based and modified the evaluation criteria in a subsequent offering to successfully promote increased learner-learner interactions. Similarly, the findings in the current study about the two instructors' participation approaches and the characteristics of learner-learner interactions developed around them can be shown to the instructors so that they are better able to align their participation approaches with pedagogical goals.

**6.3 Finding Community in the Crowd**

The ability of MOOC discussion forums to realize effective peer support and collaborative learning has not yet been conclusively established. While some studies claim MOOC forums have the potential to foster social networks and facilitate peer-connections (Kellogg et al., 2014), others claim that MOOC forum participants are dispersed crowds rather than communities of learners, evidenced by findings that in the modularized and short-lived discussion groups, learners do not move from peripheral participation to playing important roles in supporting each other's learning (Gillani & Eynon, 2014).

This study found that learner-learner interaction in the MOOC forums was outweighed by instructor-centric interaction. However, there was one group of learners in the content network that exhibited evidence of a budding community of practice (Wenger, 1998). First, members in this community had a common interest in participating in learning-related collaboration. As revealed by the qualitative examination of their post texts, not only were their communications primarily content-related, but the

participants also valued the collaborative discussions with their peers and felt they learned from them. Second, the social connections in this community were more cohesive than in any other community. The members had strong ties with a larger number of peers and used social presence cues frequently, such as greetings and addressing each other by name. Moreover, this community also had a distributed core consisting of multiple central learners who were connected via strong ties and had interactions with a large number of peers. Similar distributed core structures have been identified by Kellogg et al. (2014) in the discussion forums in two MOOCs on digital learning and mathematics learning. Distributed core structures indicate robust and sustainable learner-learner connections, which can be valuable for retention and peer support in MOOCs (Kellogg et al., 2014). Futures studies can investigate the development of distributed core structures over time with the goal of figuring out how to nurture them.

This study did not investigate the influence of core learners on others' forum activity; however, prior research on learners with similar characteristics offer avenues for future study. For instance, learners at the core of the learner-only community in the current study were all superposters (quantity of forum contributions ranked among top 5% in the course). In prior work Huang, Dasgupta, Ghosh, Manning, and Sanders (2014) found that superposters' prolific forum behaviors were associated with not just greater quantity but also better quality activity.  Future study can investigate the influence of associated central network members on group learning and social connections in MOOC discussion forums.

### *6.3.1 Implication: Facilitating Learner Community*
Participants in MOOCs are often seen as large and dispersed crowds (Gillani et al., 2014). Indeed, the findings from this study for the purely social "self-introductions" thread indicate a large crowd that gathered and dispersed quickly, many never to participate again. In contrast, however, the small group of strongly connected learners identified in this study indicates that communities can form in MOOC discussion forums. This group showed rich interpersonal connection-building as learners both discussed substantive content and expressed emotions around this shared pursuit. Moreover, the finding that the social connections formed between the core learners in this community were stronger and more robust than those formed around the instructors indicates the potential to enhance peer support in MOOCs by facilitating community development. While prior work to support MOOC forums has focused on matching information seekers with people who have the expertise to provide a specific answer (Chandrasekaran, Kan, Tan, & Ragupathi, 2015; Yang et al., 2014), cultivating a learner community requires a different set of approaches that provide a route to ongoing peer support and social connections. Various strategies can be adopted to cultivate content-related learner communities, building on the set of techniques offered to support content-related conversations in Section 6.1. In addition to the already mentioned strategies, MOOCs can offer optional activities with a sustained group that involve group collaboration or sharing of personal experience (Shackelford & Maxwell, 2012); instructors can model and encourage the use of discussion interaction techniques, such as asking leading questions and using social presence cues (Rovai, 2000); learners can be invited to volunteer as peer facilitators for units that they are interested in (Zydney, 2014); and instructors can use analytics to monitor the dynamics of learner interaction in order to phase out his/her facilitation as learner communities become robust and self-sustaining (Brooks et al., 2014).

## 6.4 Impact of Tie Definition

In addition to the contributions made to our understanding of learning through discussion in MOOCs, and the value of differentiating content-based discussions, this study also made methodological contributions related to tie definition.

First, it was found that network formation was relatively robust to tie definition; with the important exception of Total Copresence, the networks constructed did not differ dramatically in structure or properties. Interestingly, although Direct Reply / Star and Limited Copresence are based on different

assumptions about speaking and listening activities, the resultant networks were quite similar. It is worth investigating to what extent this is specific to the data used in this study (either due to the MOOC forum that only allowed two levels of replies or just the actual patterns of posting that occurred). This might usefully be done with simulated data designed to have various different properties (Hewitt, 2005). In addition, it would also be interesting to examine how different limits for the maximum size of threads in which learners are assumed to have read all posts impacts the structure of the resultant Limited Copresence network.

Second, Direct Reply and Total Copresence are the two most commonly used tie definitions in SNA studies on MOOC discussion forums (Gillani & Eynon, 2014; Gillani et al., 2014; Jiang et al., 2014; Joksimović et al., 2016; Kellogg et al., 2014) but were found to produce dramatically different networks in terms of structure and properties. This may help to partially explain the discrepancies in the MOOC literature regarding SNA results and highlights the critical importance in both selecting and sharing tie construction decisions in MOOC SNA research. Third, Total Copresence was found to be particularly sensitive to very large threads, reflected by the inflated average node degree and deflated average edge weight. This (a) must be taken into account and specifically screened for if this definition is to be used; and (b) suggests that studies using copresence and reply-based tie definitions may not always be directly comparable.

### 6.4.1 Implications: Aligning Tie Definition with the Nature of Interaction and Relationship

Many kinds of relationships can be studied using social network analysis, such as similarity, social relations, interactions and flows (Borgatti, Mehra, Brass, & Labianca, 2009). The network relationships in MOOC discussion forums have been studied from different perspectives, such as interest (Poquet & Dawson, 2016), reciprocity (Kellogg et al., 2014), and information transmission (Jiang et al., 2015). The nature of the relationship to be examined should determine the appropriate tie definition. For instance, when studying reciprocity in posting behavior, reply-based definitions can be appropriate as social connection is considered as one participant "speaking" to another (Kellogg et al., 2014). However, for connections between forum participants related to information flow or interest, a copresence definition is more appropriate as it takes into account both "speaking" and "listening". Moreover, when choosing a tie definition, it is also important to explicitly acknowledge its strengths and limitations. For instance, Direct Reply traces visible interactions among participants and has a low risk for constructing false connections, however researchers have noted that the limited number of posting levels in MOOC forums (such as EdX and Coursera) may impact the resulted structural information extracted from threads (Chandrasekaran et al., 2015; Rossi & Gnawali, 2014). In addition, this definition is also more vulnerable to misplaced posts caused by participants' liberal or unintentional use of the posting structure. In contrast, Total Copresence constructs learners' connections based on their presence in the same threaded discussion and thus can capture "invisible" connections. However, the risk of overestimating connections should be accounted for due to the fact that as the discussion evolves, it may change direction or involve new topics that not all participants in the thread are interested in (Stump et al., 2013).

At a more general level, the effects of tie definition found here highlight the overall importance and impact of the many micro decisions researchers make during SNA. Another example of a decision made in this study that could be examined in the future is the choice of community detection algorithm. In this work we used the Louvain method for calculating modularity (which identifies communities as non-overlapping clusters) because we wanted to identify and examine distinctions in patterns of connections among learner-only and instructor-included communities. Given other analysis goals, community detection measures that allow for either sparse or dense overlap in community membership (Yang & Leskovec, 2013) might be appropriate. Importantly, operational choices around community detection, tie definition and content-based differentiation all contain presumptions about the interactions and relationships from the learning environment being represented. It is important to reflect on and be

transparent about the impact of such decisions on the findings and their interpretation so as to allow meaningful comparison across studies.

## 7 LIMITATIONS AND FUTURE STUDY

The scope of this study is limited to examining the end-of-course networks in one statistics MOOC. Additional work on MOOCs in other domain areas (e.g., the humanities), using other pedagogies (e.g., constructivism), and those specifying other purposes for the discussion forums (e.g., structured group-work) is needed to determine the extent to which the current findings generalize more broadly. In addition, SNA literature has suggested the usefulness of network refining techniques such as filtering nodes and edges and assigning differentiated weight to edges based on a ranking of interaction types (Dowell et al., 2015; Gillani et al., 2014). Depending on the tie definition used, edges could be assigned different weights based on position in the thread (e.g., in a Direct Reply tie construction starter-reply connections could be weighted more heavily than reply-reply ones) or scaled based on the overall thread size (e.g., in Total Copresence the weight of a tie between two people decreases in inverse proportion to the size of the thread). Moreover, future work can use dynamic network analysis techniques such as exponential random graph models (Shumate & Palazzolo, 2010) to investigate the evolution of content networks over time and statistically test differences in network properties. Finally, in addition to the directions already mentioned in the discussion section, future work needs to test the extent to which differentiating discussions based on their content and network measures contribute additional explanatory power to the prediction of learning outcomes above and beyond overall participation counts.

## 8 CONCLUSION

Online discussion forums are commonly provided in MOOCs as a medium for learning support and interaction, but their actual value in fulfilling that purpose is not clear. This work provided insight into specific differences in the way learners interacted in content and non-content discussions in a statistics MOOC and the resultant structures of social relationships that developed. In doing so, it contributes to a nascent understanding of MOOC discussion learning processes by: (a) documenting the extended rounds of explanations, clarifications and follow-up questions through which social connections are built; (b) identifying instructor practices associated with greater learner-learner interactions; and (c) establishing evidence to support possibilities for learner community within the crowd. In addition, drawing connections between network structures and specific discussion practices lays the groundwork for the development of social network analytics that can usefully inform ongoing discussion activity. Finally, the study makes methodological contributions by: (a) providing empirical evidence that demonstrates the importance of the separate examination of content and non-content discussions; and (b) drawing attention to the importance of tie definition and potential problems associated with the Total Copresence definition.

## REFERENCES

Agrawal, A., Venkatraman, J., Leonard, S., & Paepcke, A. (2015). YouEDU: Addressing confusion in MOOC discussion forums by recommending instructional video clips. In *Proceedings of the 8th International Conference on Education Data Mining* (pp. 297-304). ACM, New York, NY, USA.

Auerbach, C., & Silverstein, L. B. (2003). *Qualitative data: An introduction to coding and analysis*. NYU press.

Bangert, A. (2008). The influence of social presence and teaching presence on the quality of online critical inquiry. *Journal of Computing in Higher Education, 20*(1), 34-61. http://doi.org/10.1007/BF03033431

Baker, C. (2010). The impact of instructor immediacy and presence for online student affective learning, cognition, and motivation. *The Journal of Educators Online, 7*(1). http://doi.org/10.9743/JEO.2010.1.2

Borgatti, S., Mehra, A., Brass, D., & Labianca, G. (2009). Network analysis in the social sciences. *Science, 323*(5916), 892-895.

Breslow, L., Pritchard, D. E., DeBoer, J., Stump, G. S., Ho, A. D., & Seaton, D. T. (2013). Studying learning in the worldwide classroom research into edX's first MOOC. *Research & Practice in Assessment, 8*, 13-25.

Brinton, C. G., Chiang, M., Jain, S., Lam, H., Liu, Z., & Wong, F. M. F. (2014). Learning about social learning in MOOCs: From statistical analysis to generative model. *IEEE Transactions on Learning Technologies, 7*(4), 346-359.

Brooks, C., Greer, J., & Gutwin, C. (2014). The data-assisted approach to building intelligent technology-enhanced learning environments. In Larusson, J. A. & White, B. (Eds), *Learning Analytics* (pp. 123-156). Springer New York.

Brugha, M., & Restoule, J. P. (2016). Examining the learning networks of a MOOC. In ElAtia, S., Ipperciel, D., & Zaïane, O. R (Eds.), *Data Mining and Learning Analytics: Applications in Educational Research* (pp. 121-138). Wiley.

Chandrasekaran, M. K., Kan, M. Y., Tan, B. C., & Ragupathi, K. (2015). Learning instructor intervention from MOOC forums: Early results and issues. In *Proceedings of the 8th International Conference on Education Data Mining* (pp. 218-225). ACM, New York, NY, USA.

Cho, H., Gay, G., Davidson, B., & Ingraffea, A. (2007). Social networks, communication styles, and learning performance in a CSCL community. *Computers & Education, 49*(2), 309-329.

Creswell, J. W., Klassen, A. C., Plano Clark, V. L., & Smith, K. C. (2011). Best practices for mixed methods research in the health sciences. *Bethesda (Maryland): National Institutes of Health*, 2094-2103.

Dennen, V. P. (2005). From message posting to learning dialogues: Factors affecting learner participation in asynchronous discussion. *Distance Education*, *26*(1), 127-148.

Dennen, V. P., Aubteen Darabi, A., & Smith, L. J. (2007). Instructor–learner interaction in online courses: The relative perceived importance of particular instructor actions on performance and satisfaction. *Distance Education*, *28*(1), 65-79.

Dowell, N., Skrypnyk, O., Joksimović, S., Graesser, A. C., Dawson, S., Gašević, D., Vries, P. d., Hennis, T., & Kovanović, V. (2015). Modeling learners' social centrality and performance through language and discourse. In *Proceedings of the 8th International Conference on Educational Data Mining* (pp. 250-257). ACM, New York, NY, USA.

Gibson, W., & Brown, A. (2009). *Working with qualitative data*. Sage.

Gillani, N., & Eynon, R. (2014). Communication patterns in massively open online courses. *The Internet and Higher Education*, 23, 18-26.

Gillani, N., Yasseri, T., Eynon, R., & Hjorth, I. (2014). Structural limitations of learning in a crowd: Communication vulnerability and information diffusion in MOOCs. *Nature Scientific Reports, 4*. http://doi.org/10.1038/srep06447

Gruzd, A.A., & Haythornthwaite, C. (2008). Automated discovery and analysis of social networks from threaded discussions. In *Proceedings of the International Network of Social Network Analysts 2008*. Retrieved September 28, 2016, from http://hdl.handle.net/10150/105081.

Guba, E. G. (1981). Criteria for assessing the trustworthiness of naturalistic inquiries. *Educational Technology Research and Development*, *29*(2), 75-91.

Hew, K. F. (2015). Student perceptions of peer versus instructor facilitation of asynchronous online discussions: Further findings from three cases. *Instructional Science*, *43*(1), 19-38.

Hewitt, J. (2005). Toward an understanding of how threads die in asynchronous computer conferences. *The Journal of the Learning Sciences*, *14*(4), 567-589.

Houston, S. L., Brady, K., Narasimham, G., & Fisher, D. (2017). Pass the idea please: The relationship between network position, direct engagement, and course performance in MOOCs. In *Proceedings of the 4th (2017) ACM Conference on Learning@ Scale* (pp. 295-298). ACM, New York, NY, USA.

Huang, J., Dasgupta, A., Ghosh, A., Manning, J., & Sanders, M. (2014). Superposter behavior in MOOC forums. In *Proceedings of the First ACM Conference on Learning@ Scale Conference* (pp. 117-126). ACM.

Ifenthaler, D., Adcock, A. B., Erlandson, B. E., Gosper, M., Greiff, S., & Pirnay-Dummer, P. (2014). Challenges for education in a connected world: Digital learning, data rich environments, and computer-based assessment - Introduction to the inaugural special issue of technology, knowledge and learning. *Technology, Knowledge and Learning*, *19*(1-2), 121.

Jacobsen, D. Y. (2017). Dropping out or dropping in? A Connectivist approach to understanding participants' strategies in an e-learning MOOC pilot. *Technology, Knowledge and Learning*. http://doi.org/10.1007/s10758-017-9298-z

Jiang, S., Fitzhugh, S. M., & Warschauer, M. (2014). Social positioning and performance in MOOCs. In *Proceedings of Graph-Based Educational Data Mining Workshop at the 7th International Conference on Educational Data Mining* (pp. 55-58). CEUR-WS.

Jiang, Z., Zhang, Y., Liu, C., & Li, X. (2015). Influence analysis by heterogeneous network in MOOC forums: What can we discover?. In *Proceedings of the 8th International Conference on Education Data Mining* (pp. 242-249). ACM, New York, NY, USA.

Joksimović, S., Manataki, A., Gašević, D., Dawson, S., Kovanović, V., & De Kereki, I. F. (2016). Translating network position into performance: Importance of centrality in different network configurations. In *Proceedings of the 6th International Conference on Learning Analytics & Knowledge* (pp. 314-323). ACM New York, NY, USA. http://doi.org/10.1145/2883851.2883928

Kellogg, S., Booth, S., & Oliver, K. (2014). A social network perspective on peer supported learning in MOOCs for educators. *The International Review of Research in Open and Distributed Learning*, *15*, 5. http://doi.org/10.19173/irrodl.v15i5.1852

Khalil, H., & Ebner, M. (2013). "How satisfied are you with your MOOC?" - A research study on interaction in huge online courses. In *Proceedings of EdMedia 2013* (pp. 830-839). AACE.

Kuh, G. (2002). From promise to progress: How colleges and universities are using student engagement results to improve collegiate quality. *National Survey of Student Engagement Annual Report*. Bloomington, IN: Indiana University.

Mazzolini, M., & Maddison, S. (2007). When to jump in: The role of the instructor in online discussion forums. *Computers & Education*, *49*(2), 193-213.

McGuire, R. (2013). Building a sense of community in MOOCs. *Campus Technology*, *26*(12), 31-33.

Papadopoulos, K., Sritanyaratana, L., & Klemmer, S. R. (2014). Community TAs scale high-touch learning, provide student-staff brokering, and build esprit de corps. In *Proceedings of the First ACM Conference on Learning@ Scale Conference* (pp. 163-164). ACM.

Poquet, L., & Dawson, S. (2016). Untangling MOOC learner networks. In *Proceedings of the 6th International Conference on Learning Analytics & Knowledge* (pp. 208-212). ACM New York, NY, USA. http://doi.org/10.1145/2883851.2883919.

Romero, C., López, M. I., Luna, J. M., & Ventura, S. (2013). Predicting students' final performance from participation in on-line discussion forums. *Computers & Education*, *68*, 458-472.

Rosé, C. P., & Ferschke, O. (2016). Technology support for discussion based learning: From computer supported collaborative learning to the future of Massive Open Online Courses. *International Journal of Artificial Intelligence in Education*, *26*(2), 660-678.

Rossi, L.A., & Gnawali, O. (2014). Language independent analysis and classification of discussion threads in Coursera MOOC forums. In *Proceedings of 2014 IEEE 15th International Conference on Information Reuse and Integration* (pp. 654-661). IEEE. http://doi.org/10.1109/IRI.2014.7051952.

Rovai, A. P. (2000). Building and sustaining community in asynchronous learning networks. *The Internet and Higher Education*, *3*(4), 285-297.

Santos, J.L., Klerkx, J., Duval, E., Gago, D., & Rodríguez, L. (2014). Success, activity and drop-outs in MOOCs: An exploratory study on the UNED COMA courses. In *Proceedings of the 4th International Conference on Learning Analytics & Knowledge* (pp. 98-102). ACM New York, NY, USA. http://doi.org/ 10.1145/2567574.2567627

Scott, J., & Carrington, P. J. (2011). *The SAGE handbook of social network analysis*. SAGE publications.

Shackelford, J. L., & Maxwell, M. (2012). Sense of community in graduate online education: Contribution of learner to learner interaction. *The International Review of Research in Open and Distributed Learning*, *13*(4), 228-249.

Shumate, M., & Palazzolo, E. T. (2010). Exponential random graph (p*) models as a method for social network analysis in communication research. *Communication Methods and Measures*, *4*(4), 341-371.

Stump, G. S., DeBoer, J., Whittinghill, J., & Breslow, L. (2013). Development of a framework to classify MOOC discussion forum posts: Methodology and challenges. In *Proceedings of NIPS 2013 Workshop on Data Driven Education* (pp. 1-20). NIPS Foundation.

Suthers, D. D. (2015). From contingencies to network-level phenomena: Multilevel analysis of activity and actors in heterogeneous networked learning environments. In *Proceedings of the Fifth International Conference on Learning Analytics and Knowledge* (pp. 368–377). ACM.

Tomkin, J. H., & Charlevoix, D. (2014). Do professors matter?: Using an a/b test to evaluate the impact of instructor involvement on MOOC student outcomes. In *Proceedings of the First ACM Conference on Learning@ Scale Conference* (pp. 71-78). ACM.

Trentin, G. (2000). The quality-interactivity relationship in distance education. *Educational Technology, 40*(1), 17-27.

Wang, X., Yang, D., Wen, M., Koedinger, K., & Rosé, C. P. (2015). Investigating how student's cognitive behavior in MOOC discussion forums affect learning gains. In *Proceedings of the 8th International Conference on Educational Data Mining*. International Educational Data Mining Society. Retrieved December 6, 2016, from http://files.eric.ed.gov/fulltext/ED560568.pdf.

Wellman, B., & Berkowitz, S. D. (Eds.). (1988). *Social structures: A network approach* (Vol. 2). CUP Archive.

Xiong, Y., Li, H., Kornhaber, M. L., Suen, H. K., Pursel, B., & Goins, D. D. (2015). Examining the relations among student motivation, engagement, and retention in a MOOC: A structural equation modeling approach. *Global Education Review*, *2*(3), 22-33.

Wenger, E. (1998). *Communities of practice: Learning, meaning, and identity*. Cambridge, UK: Cambridge University Press.

Wise, A. F., & Cui, Y. (2018). Unpacking the relationship between discussion forum participation and learning in MOOCs: Content is key. In *Proceedings of the 8th International Learning Analytics & Knowledge Conference*. ACM.

Wise, A. F. & Schwarz, B. S. (2017). Visions of CSCL: Eight provocations for the future of the field. *International Journal of Computer-Supported Collaborative Learning 12*(4), 1-45.

Yang, D., Adamson, D., & Rosé, C. P. (2014). Question recommendation with constraints for massive open online courses. In *Proceedings of the 8th ACM Conference on Recommender Systems* (pp. 49-56). ACM, New York, NY, USA.

Yang, J., & Leskovec, J. (2013). Overlapping community detection at scale: A nonnegative matrix factorization approach. In *Proceedings of the Sixth International Conference on Web Search and Data Mining* (pp. 587-596). ACM.

Yusof, N., & Rahman, A. A. (2009). Students' interactions in online asynchronous discussion forum: A social network analysis. In *Proceedings of 2009 International Conference on Education Technology and Computer* (pp. 25-29). IEEE.

Zhu, M., Bergner, Y., Zhang, Y., Baker, R., Wang, Y., & Paquette, L. (2016). Longitudinal engagement, performance, and social connectivity: A MOOC case study using exponential random graph

models. In *Proceedings of the 6th International Conference on Learning Analytics & Knowledge* (pp. 223-230). ACM, New York, NY, USA. http://doi.org/10.1145/288385.

Zydney, J. (2014). Strategies for creating a community of inquiry through online asynchronous discussions. *Journal of Online Learning and Teaching*, *10*(1), 153-165.