

# A multilevel clustering technique for community detection

Isa Inuwa-Dutse\*, Mark Liptrott, Ioannis Korkontzelos\*\*

*Department of Computer Science, Edge Hill University,  
St Helens Rd, Ormskirk, L39 4QP, Lancashire, United Kingdom*

---

## Abstract

A network is a composition of many communities, i.e., sets of nodes and edges with stronger relationships, with distinct and overlapping properties. Community detection is crucial for various reasons, such as serving as a functional unit of a network that captures local interactions among nodes. Communities come in various forms and types, ranging from biologically to technology-induced ones. As technology-induced communities, social media networks such as Twitter and Facebook connect a myriad of diverse users, leading to a highly connected and dynamic ecosystem. Although many algorithms have been proposed for detecting socially cohesive communities on Twitter, mining and related tasks remain challenging. This study presents a novel detection method based on a scalable framework to identify related communities in a network. We propose a multilevel clustering technique (*MCT*) that leverages structural and textual information to identify local communities termed *microcosms*. Experimental evaluation on benchmark models and datasets demonstrate the efficacy of the approach. This study contributes a new dimension for the detection of cohesive communities in social networks. The approach offers a better understanding and clarity toward describing how low-level communities evolve and behave on Twitter. From an application point of view, identifying such communities can better inform recommendation, among other benefits.

*Keywords:* Clustering, Multilevel clustering, Community detection, Twitter, Social networks

---

## 1. Introduction

A network comprises of many sub-networks or communities with distinct and overlapping properties. Networks exhibit varying degrees of organisations [1], and discovering the structure of various network forms has been investigated [2, 3, 4, 5]. As network size increases, so does the possibility of fragmentation [6, 7], leading to a decrease in the homogeneity of behaviour and attitude across groups [8]. Because similarity breeds attraction and interaction [9], network communities are defined by sets of nodes and edges with strong relationships. Communities are a fundamental organisation principle, especially in vast networks, allowing to analyse the structure and function of networks [3, 10]. Identifying local network structures: (a) provides a means for complex network analysis [11], for applications such as the detection of inter-related web-pages [12, 13], (b) allows to detect cliques [10] and facilitates intelligent recommendations [14] (c)

---

\*Corresponding author

\*\*Principal corresponding author

*Email addresses:* isadutse91@gmail.com (Isa Inuwa-Dutse), Mark.Liptrott@edgehill.ac.uk (Mark Liptrott), Ioannis.Korkontzelos@edgehill.ac.uk (Ioannis Korkontzelos)

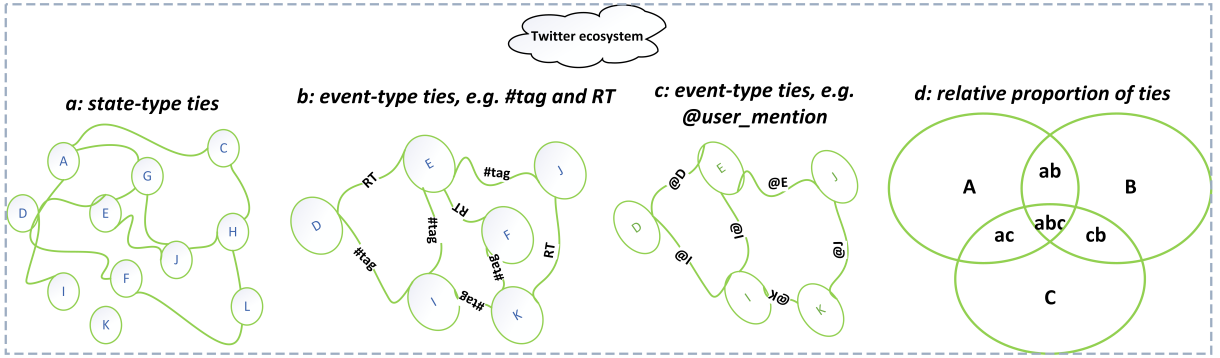


Figure 1: Examples of *event-type ties* ((a),(b) and (c)) on Twitter, allowing users to openly connect via: (a) unidirectional or directed means (e.g. friend or follower), bidirectional or undirected (among friends and followers) (b – c), and indirect or *transitory events*: *retweets*, *mentions* or *likes*. These flexible connections challenge cohesive community detection and contribute to the proliferation of spurious content. In (d),  $A = \{a, a_1, \dots, a_l\}$ ,  $B = \{b, b_1, \dots, b_m\}$ , and  $C = \{c, c_1, \dots, c_n\}$  denote users. Reciprocal ties (e.g.  $ac$ ) or transitive ties (e.g.  $abc$ ) are rare.

12 allows to discover organisational principles of networks [15, 16], and (d) helps in studying  
 13 social behaviour of users [17]. Examples of biological, social or technological networks where  
 14 community detection has been applied are: *protein–protein interaction networks* [18], *social*  
 15 *networks* [4, 10], *food webs* [11], *collaboration networks*, [19] and the *World Wide Web* [4].

16 The underlying difference across many network communities refers to the definition of  
 17 connections: some are deterministic, while some are just probabilistic and potentially non-  
 18 deterministic. Social media, e.g., *Twitter* and *Facebook*, connect a myriad of diverse users,  
 19 leading to a highly connected, dynamic ecosystem. The complexity and dynamism of this  
 20 ecosystem results in multiple interaction types at various layers of granularity and intensity:  
 21 global or local, positive or negative, influential or not, high or low-level. Such interactions  
 22 culminate in the formation of communities at various levels. Despite the proliferation of various  
 23 community detection methods, identifying socially cohesive communities on Twitter still poses  
 24 challenges. Communities with low presence are implicit and require extensive exploration  
 25 to understand the mechanism governing their behaviour [20]. Since social networks exhibit  
 26 properties from other networks [10], the limitations of existing approaches are due to:

27 *Methodological viewpoints and connection types.* Social network theorists hold two viewpoints  
 28 in investigating social relationships in a network: realist, based on a pre-conceived notion of the  
 29 existence of relationships, and nominalist, based on questions posed by the investigator [21].  
 30 Moreover, social ties are formed around *event-type ties*, a transitory connection that often results  
 31 in socially distant members. Such connections on Twitter include subscribing to trendy hashtags  
 32 or retweeting popular users. *State-type ties* are based on static, structural connections among  
 33 users, which suggest familiarity and trust [22]. Community detection on Twitter focusses mostly  
 34 on directed connections (*event-type ties*) based on the *realist's* approach. This is valid in many  
 35 networks, but can lead to many unrelated sets of users. We argue that the wealth of connection  
 36 forms on Twitter, shown in Figure 1, contribute to widespread spurious content and imply the  
 37 existence of less cohesive user communities.

38 *Proliferation and complexity of online content.* a rapid increase in network size increases the  
 39 likelihood of fragmentation [6, 7], which in turn decreases the homogeneity of behaviour and

40 attitude across groups [8]. With an average 139m daily users contributing to 500m content<sup>1</sup>,  
41 it is becoming more challenging to keep track of socially cohesive communities on Twitter.  
42 Furthermore, large scale and transitory content (mostly from influential users) often dominate the  
43 space leading to many forms of explicit communities [23]. Thus, basing a community detection  
44 task on transitory aspects of *metadata* such as popular hashtags or trending topics does not often  
45 reflect true connectivity [24], hence limiting the full realisation of the benefits in communities  
46 such as *cliquishness* and *local coherence*.  
47 This study attempts to address the identified challenges, to advance our knowledge concerning  
48 community detection problems.

### 49 1.1. A Multilevel Clustering Technique

50 A community detection paradigm involves prediction and quantification to identify a com-  
51 munity structure and relevant details about a network [25]. Predicting membership and assigning  
52 items to clusters is achieved using equivalence measures or scoring functions. Establishing  
53 the equivalence of network entities is achieved based on (a) equivalent units with the same  
54 connection pattern to the same neighbours or (b) equivalent units that have the same or a similar  
55 connection pattern to different neighbours [26]. Accordingly, communities are formed around  
56 two primary modalities or information sources: *network structure* and *node attributes*. Until  
57 recently, community detection methods relied on a single information source. Conventional  
58 methods such as *normalised cut* [27] and modularity [28] rely on the topological structure of  
59 networks. A bi-modal approach, based on network structure and the corresponding features or  
60 attributes of nodes as information sources, is becoming popular [29, 30, 31, 16]. According to  
61 Figure 1, connections on *Twitter* may manifest differently, such as *sharing a link*, *re-tweeting*  
62 (*RT*), using the same or similar *hashtags*, *user mention* (@) or *follower-ship*. Such connections  
63 are porous, allowing to connect with many diverse users and hindering the identification of  
64 cohesive groups. Noting that these eccentric connections patterns can mislead the detection of  
65 socially related users and encourage the propagation of spurious content, we propose a *multi-*  
66 *level clustering technique (MCT)* to identify socially cohesive user groups on Twitter, termed  
67 *microcosms*. No practical reasons prevent MCT to apply to other domains that involve network  
68 data. However, it would require minor amendments for platforms where a reciprocal tie is the  
69 default connection, e.g. Facebook. Failing to recognise Twitter’s eccentric topological structure  
70 would make the approach less generalisable. Focusing on Twitter leads to a more encompassing  
71 framework that can be mapped to other networks.

72 MCT measures similarity within a community of users using local and global information,  
73 modelling *structural* and intrinsic *textual* features jointly. In Figure 2(a) and (b) user communities  
74 exist based on *structural* and *content* or *textual* similarity, respectively. Users under the structural  
75 component, a form of a *state-type tie*, are related based on reciprocal ties, which are rare in  
76 Twitter, and the community is more cohesive than the community of users based on content  
77 or textual similarity, a form of an *event-type tie*. A more cohesive community is the one that  
78 recognises both structural and content similarity, in Figure 2(c). Intuitively, the degree of  
79 cohesiveness varies across different communities: a community based on both modalities is  
80 the most cohesive, followed by a community based on high structural similarity but low or no  
81 content or textual similarity. Finally, the least cohesive community exhibits high similarity in the

---

<sup>1</sup>See [www.omnicoreagency.com/twitter-statistics](http://www.omnicoreagency.com/twitter-statistics)

Table 1: Relevant notations and their corresponding descriptions

Notation	Description
$\mathcal{D}$	network data
$\mathcal{V}$ and $\mathcal{E}$	sets of nodes and edges, respectively
$m_{v_i}$	denotes the network of a user $v_i$
$fr_{v_i}$ and $fl_{v_i}$	denote sets of friends and followers of user $v_i$
$\kappa \in m_{v_i}$	set of reciprocal ties of user $v_i$
$\mathcal{A}_f$	set of all possible attributes of a node
$\mathcal{X}_f$	set of features inducing reciprocity
$a \succ b$	a binary relation between $a$ and $b$
$\exp x$ or $e^x$	base of a natural logarithm, $\ln$ , $\exp(\ln x) = x$ ; $e \approx 2.71828$
$d_i, p_i, q_i, u_i$	respective $i$ th component of vectors, $\mathbf{D}, \mathbf{P}, \mathbf{Q}, \mathbf{U}$
$\tau$	a predefined threshold for comparison, e.g. $\tau \geq 0.5$
$\mathcal{S}_r$ and $\mathcal{T}_r$	sets of structurally-related and textually-related nodes, respectively

82 textual component but low or no structural similarity. Groups of structurally similar nodes are  
 83 analysed by *spectral clustering*, which involves a series of methods ranging from adjacency and  
 84 affinity matrices to dimensionality reduction. The textual component complements the structural  
 85 aspect through a form of document-pivot clustering, in which weights are assigned to features in  
 86 the document according to a weighing scheme [32, 33, 34, 35].

87 *1.2. Contributions*

88 MCT relies on reciprocal ties, based on the assumption that combining structural and textual  
 89 features offers a more cohesive community representation. Our contributions are two-fold:

90 *A new dimension to the detection of cohesive communities.* The ability to follow anyone on  
 91 Twitter results in many unidirectional connections between socially unrelated users, affecting  
 92 clustering and the integrity of online content. To counter the challenging and time-consuming  
 93 task of collecting large scale reciprocal ties on Twitter, we proposed a strategy that returns  
 94 the likelihood of reciprocity among users. As a result, the detection of socially cohesive  
 95 communities is enhanced, providing a useful analysis tool and strengthening the validity of  
 96 online content. Moreover, by identifying communities of users with a strong cohesion, a well-  
 97 informed recommendation that recognises structural and textual similarities can be achieved.

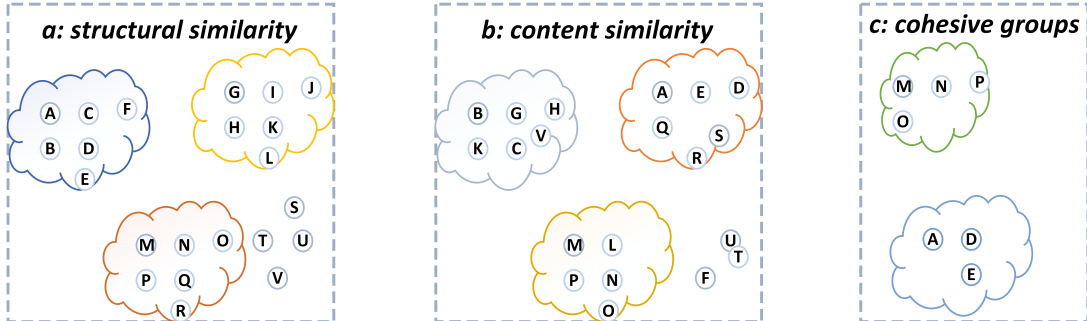


Figure 2: Node clustering in MCT in three stages, according to (a) structural similarity, (b) *content* or *textual* aspect similarity, and (c) joint *structural* and *content* similarity.

108 *A bi-modal community detection approach.* MCT addresses the problem of structurally unrelated  
109 users by adding a layer of social cohesion to existing community detection methods. Specifically,  
110 MCT advances existing techniques through: (a) an in-depth utilisation of a bi-modal source of  
111 information for community detection, (b) detection of network communities at various levels, (c)  
112 a robust and scalable community detection algorithm that is less prone to noise in the network  
113 data, and (d) an intuitive interpretation of the detected communities.

114 The remaining of this paper is structured as follows. Section 2 provides background details  
115 and related work. Section 3 formulates the problem and describes the MCT framework. We  
116 describe the experimentation process in Section 4 and discuss the main observations in Section 5.  
117 Finally, Section 6 concludes the study and provides some considerations for future work. For  
118 ease of referencing, Table 1 summarises the notations used.

## 109 2. Background

110 Humans can effortlessly abstract complex phenomena, but efficiently automating the process  
111 is daunting partly due to the multidimensional nature of clustering data [36]. In this section, we  
112 review relevant topics and studies associated with clustering and community detection tasks.

### 113 2.1. Network and community structure

114 A network comprises of heterogeneous nodes connected via edges. The topological structure  
115 of networks and other quantities related to them are useful in understanding complex networks  
116 across numerous domains [4]. Various levels of relationship forms in networks have been  
117 analysed, from the structure of microscopic organisms to complex networks, such as the internet  
118 [2, 3, 4]. Complex networks were once considered to as random and the classic random graph  
119 model [37] was the standard analysis tool until regular patterns in various networks were  
120 discovered, e.g., via statistical analysis. Fundamentally, network complexity [38] is defined  
121 by: (a) *Clustering coefficient* quantifies the probability of a node to be clustered, assuming that  
122 users with common friends are likely to know each other. (b) *Degree distribution* quantifies the  
123 probabilistic distribution of nodes. (c) *Small-worldliness* is a network property associated with  
124 short path-length, i.e., many structured short-range connections and few random long-range ones,  
125 and network diameter that is exponentially less than its size [3].

126 In *social media*, communication happen at various layers of granularity and intensity: global  
127 or local, positive or negative, influential or not. In contrast with the early unidirectional *two-step*  
128 communication model, where few users serve as intermediaries between mass communication  
129 and the public [39], the design of social media allows users to generate and consume information.  
130 On social media, communication follows the influence network model, enabling multi-way flow,  
131 where users can simultaneously generate and consume information [38]. Twitter is dominated by  
132 influential users, logically dividing a clique of content pushers and consumers, resembling the  
133 two-step flow model [39]. This division is strengthened by strategies, such as content promotion,  
134 that entice users to engage more, and to follow or add friends. Using these strategies social  
135 media users can increase their network of friends, generating more value to the platforms. A  
136 social media network is the synthesis of many user communities, and identifying these structures  
137 is a vital task. Because members of a community are highly similar to each other and less or not  
138 at all similar to members of other communities [40, 1, 41], a *community structure* has densely  
139 connected node groups and sparser connections to other communities [42]. Thus, community

140 identification involves prediction and quantification tasks to detect the relevant structures and  
141 their characteristics [25]. Selecting an effective similarity measure is crucial as it allows a  
142 clustering algorithm to identify groups and affects the *signal-to-noise-ratio* within the instance  
143 matrix [43].

## 144 2.2. Related work

145 Network partitioning has attracted interest from various domains of expertise, hence diverse  
146 strategies have been put forward to identify relevant communities embedded in a network.

147 *Clustering and Community Detection.* Often clustering and community detection are used  
148 interchangeably in the literature. Clustering mostly focuses on a single modality, e.g., using  
149 node attributes to group network objects, whereas community detection focuses on network  
150 structure as a function of connectivity involving social interaction. As a form of dimensionality  
151 reduction, clustering entails unsupervised partitioning a network into groups of related objects  
152 using a domain-specific scoring function and maximising in-group similarity. There are two  
153 principal lines of research in this direction: graph partitioning and hierarchical modelling [41].  
154 We follow this classification, as it reflects the approach in this study. Methods can also be divided  
155 in dimensionality reduction based ones, and graph partitioning (hierarchical or not [44]) and  
156 hierarchical ones [45].

### 157 2.2.1. Graph-based and hierarchical methods

158 Graph-based clustering assumes that a community structure exists in the network and attempts  
159 to discover it using specific techniques. Graph partitioning divides the network into predefined  
160 node groups and suits applications where the number and size of groups are known, e.g., in  
161 parallelisation of computing processors. The approach may involve hierarchical agglomerative  
162 clustering [46] following a random walk model [37], or based on modularity [42] optimisation,  
163 such as in the Louvain detection algorithm [47]. The clustering method can be based on *iterative*  
164 *bisection*, which divides the network optimally into two parts and repeats until the required  
165 number of partitions is reached [28]. The modularity measure measures community strength  
166 and detects groups, assuming that community structures correspond to an interesting statistical  
167 arrangement of edges. Positive values indicate the presence of community structures, i.e., that  
168 nodes within a community are more tightly connected than by chance [28]. The modularity  
169 value of real networks ranges from .3 to .7. The higher the score, the more cohesive the  
170 community structure [41]. Predefining the maximum bisection size is required, which may affect  
171 performance. Metrics such as *betweenness* or *shortest loop edges*, are central to the operation of  
172 algorithms that process graphs to detect groups of similar nodes [48].

173 Hierarchical modelling follows a different technical approach from graph-based clustering.  
174 Assuming that there are natural subgroups in a network, hierarchical clustering utilises a similar-  
175 ity measure, such as Euclidean distance or Pearson correlation to analyse the network [2]. In  
176 particular, pairwise node similarities are computed and nodes are iteratively and deterministi-  
177 cally assigned to clusters. Commonly, similar clusters are iteratively merged into larger ones  
178 [45]. Furthermore, categorisation based on a generative or model-based and discriminative  
179 or similarity-based is used in the literature. Model-based or generative clustering algorithms,  
180 e.g., Latent Dirichlet Allocation (LDA) [49, 29, 50], are a form of Expectation Maximisation  
181 (EM) that aim to learn a generative model from data segments, where each model represents a  
182 cluster [51]. The EM-based models estimate the maximum likelihood of data-points to belong

183 to a cluster and is suitable for incomplete data. On the other hand, similarity-based clustering  
184 algorithms are based on optimising a scoring function that is used to compute pairwise similarity  
185 between data-points. This form of clustering follows hierarchical agglomerative clustering or  
186 block modelling.

### 187 2.2.2. *Multiview and bi-modal clustering*

188 Multiview and bi-modal techniques aim to improve clustering performance using multiple  
189 independent data sources; thus, multiview clustering relies on data that can be split into indepen-  
190 dent sub-features or attributes [52, 53]. For instance, a *web page* can be described by its textual  
191 content and pages that link to it [54, 55]. The advantages of multiview clustering over its single  
192 view counterpart has been investigated using algorithms based on *K-means* and *Expectation*  
193 *Maximisation* [54].

194 Bi-modal clustering technique is based on the fact that network communities are formed  
195 around two primary modalities or information sources: *network structure* and *node attributes*. In  
196 many cases, structural and textual aspects evolve simultaneously and communities are discovered  
197 according to the nodes' similarity degree vis-à-vis those two aspects. Until recently, studies  
198 mostly focused on one aspect, not both [16, 29, 31, 56, 57]. A study closely related to our  
199 approach, proposes a generative model for networks with node attributes [16]. However, the  
200 depth of the features, especially the nodes' attributes, is shallow and the node attribute (*hashtag*) is  
201 insufficient to analyse the depth of similarity between network entities in a complex environment  
202 such as Twitter in which the structural component is not fully captured due to reliance on directed  
203 edges. The connected k-centre approach employs both structural and attribute information for a  
204 given network partition [56]. The problem is NP-hard, leading to many heuristics. Similarly,  
205 SA-cluster method combines structural and attributes' similarities for community detection by  
206 partitioning a network into cohesive k-clusters with structural and attribute information using a  
207 distance metric to estimate pairwise node similarity or closeness [57].

208 Conventional methods, such as normalised cut [27] and modularity [28], are based on topo-  
209 logical structures. However, many networks come with incomplete information, e.g., a *terrorist*  
210 *network* or food web [30]; thus, community detection in networks with edge uncertainty or  
211 incomplete information is getting traction. Inferring links in incomplete networks is challenging,  
212 because the information is usually localised within a small, linked group. The full wealth of  
213 data has been used to learn a generalisable distance metric to complete the missing information  
214 [30]. However, this approach is too complicated and does not account for the breadth required  
215 in *textual* aspects in networks with many transient connections, such as Twitter<sup>2</sup>. The MCT is  
216 a two-stage clustering technique that recognises different modalities as information sources;  
217 it incorporates multiview aspects at various levels, structural and textual, using independent  
218 features.

## 219 3. MCT framework

220 Noting that nodes in a community are highly similar and edges among communities are  
221 infrequent, community detection is usually formalised to identify network partitions that satisfy  
222 specific requirements. The problem focuses on detecting smaller groups with high similarities,  
223 using a joint similarity function that considers global and local information as a two-stage process

---

<sup>2</sup>As shown in Figure 1, Twitter communities are formed based on many factors.

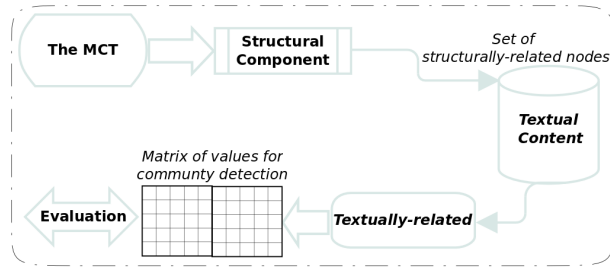


Figure 3: Execution pipeline of the MCT framework - The structural component processes a collection of structurally similar nodes, promoting group formation among them. Then, *textual analysis* of structurally related node content identifies groups according to discussion topics. MCT combines *state-type* and *event-type* ties.

224 comprising of structural and textual components (shown in Figure 2). Figure 3 shows a block  
 225 diagram of the stages in the MCT framework.

### 226 3.1. Structural component

227 The structural component is based on dyadic, pairwise edge between two users, and transitive  
 228 ties, which are the basic forms of establishing reciprocal ties in social networks. We aim  
 229 to identify groups of users with true reciprocal relationships at dyadic and transitive level.  
 230 Transitivity expresses the social preference to be friends with a *friend-of-a-friend* and has been  
 231 characterised as a peculiar network feature [3]. Transitive ties are synonymous to *Simmelian ties*,  
 232 strong social relationships among three or more individuals, which are vital in understanding a  
 233 network’s social tie structure [58]. Our approach assumes that community detection or clustering  
 234 methods that take reciprocal ties into account offer a more cohesive community representation.  
 235 Our analysis of Twitter datasets concluded that true reciprocal ties are rare. Thus, we use a  
 236 method that strengthens the possibility of finding Twitter users with reciprocal ties. A user  
 237 with many reciprocated ties can represent a *microcosms*, allowing to analyse a user group as a  
 238 unit. Research in social science suggest that users compare themselves with one another and  
 239 adopt similar behaviour with users similar to them [9]. Homophily on Twitter can be interpreted  
 240 as a reciprocal relationship among users. Noting this insight and the inspiration drawn from  
 241 social homophily, we argue that users with similar profiles are more likely to connect on Twitter.  
 242 Therefore, structural equivalence is mapped to a *state-type tie* to infer structural similarity  
 243 according to the node’s attributes. Figure 4 shows features that contribute in finding structurally  
 244 related nodes.

#### 245 3.1.1. Modelling structural clusters

246 *Definitions.* This section begins with the definitions of relevant concepts and terms in the  
 247 implementation. Table 1 provides a summary of all relevant notations used in the study.

- 248 - *Network data*  $\mathcal{D}$  consists of sets of nodes,  $v_1, v_2, \dots, v_m \in \mathcal{V}$ , and edges,  $e_1, e_2, \dots, e_n \in \mathcal{E}$ .  
 249 Each node is described by its structural and textual features, as shown in Figure 4.
- 250 - *Dyadic and transitive ties:* a relation,  $\succ$ , between two nodes  $v_i, v_j \in \mathcal{D}$  is dyadic<sup>3</sup> if  $v_i$   
 251 follows  $v_j$  and vice versa, i.e.  $v_i \succ v_j \Leftrightarrow \forall v_i, v_j \in \mathcal{D}$ . In this context,  $v_i$  follows  $v_j$  is a  
 252 directed relationship; if  $v_j$  follows  $v_i$  back, it is undirected; see some examples in Figure 5.

<sup>3</sup>Dyadic tie, pairwise, 2-star or binary relations are used interchangeable in this study.



253 A transitive or *Simmelian tie* is a social relationship within groups of three or more. A  
 254 binary relation,  $\succ$ , over a set  $\mathcal{D}$  is *transitive*:  $v_i \succ v_j, v_j \succ v_k \Leftrightarrow v_i \succ v_k, \forall v_i, v_j, v_k \in \mathcal{D}$ .

255 - *Prediction features and reciprocity*: To identify structurally related nodes, we use features,  
 256 such as *indegree* (*ind*), the number of incoming edges to a node; *outdegree* (*out*), the  
 257 number of outgoing edges from a node; and *category* (*cat*), indicating if a node is *verified*  
 258 or *unverified*. Account verification ascertains legitimacy of the account holder.  $\mathcal{A}_f$   
 259 denotes the set of all possible node features, which can be used to extract feature subsets,  
 260  $\mathcal{X}_f = \{ind, out, cat\} \subset \mathcal{A}_f$ , from a profile, as in Figure 4. The features of a node pair,  
 261  $v_i, v_j$ , are denoted by:  $\mathcal{X}_{f_{v_i}} = \{ind_{v_i}, out_{v_i}, cat_{v_i}\}$  and  $\mathcal{X}_{f_{v_j}} = \{ind_{v_j}, out_{v_j}, cat_{v_j}\}$  and  
 262 are used for training models that predict the likelihood of *reciprocity*.

263 *Nodes reciprocity and constant error*. We predict node sets with dyadic ties on Twitter using the  
 264 approach in [59], and we use them for clustering. The formulation assumes that reciprocity is  
 265 based on the features that can induce friendship, in Figure 4(b). Nodes reciprocity hypothesises  
 266 that the decision to reciprocate or establish friendship correlates with the idea of homophily and  
 267 structural equivalence. Reciprocal ties are predicted based on these concepts between node pairs.

268 Consider the sets of nodes,  $\mathcal{V}$ , and edges,  $\mathcal{E}$ . The likelihood of *reciprocity*,  $p(R_{v_i, v_j}), \forall v_i, v_j \in$   
 269  $\mathcal{V}$  involved in the computation of reciprocal units (see Section 4.2.1) is described by Eq. 1 to  
 270 3, leading to *reciprocal-communities*  $\mathcal{C}_{rc}$ . The features of a node pair,  $v_i, v_j$ , for comparison  
 271 are denoted by:  $\mathcal{X}_{f_{v_i}}$  and  $\mathcal{X}_{f_{v_j}}$ , such that the ratio of the attributes, e.g. *ind* or *out*, is a real  
 272 value quantity given by:  $\frac{ind_{v_i}}{ind_{v_j}} \in \mathbb{R} \quad \forall f_r \in \mathcal{X}_{f_{v_i}, v_j}$ . If the comparison evaluates to a value in  
 273  $[0.75, 1.25]$ , the pairs are assumed to have similar attributes (1), otherwise dissimilar attributes  
 274 (0). The interval is to allow extra freedom for minor discrepancies between the features. For  
 275 instance, if the ratio equals 1.0, the pair has identical attributes, which is useful in analysing  
 276 aspects of homophily and social equivalence. The binary feature comparison values are used to

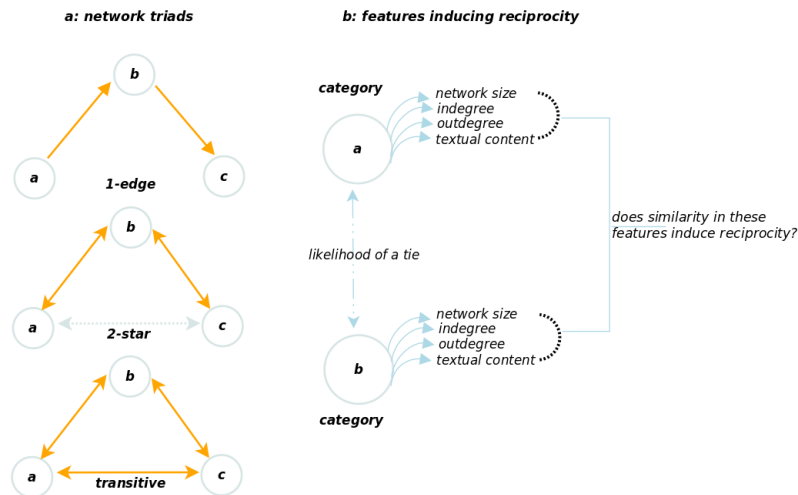


Figure 4: (a) Possible social ties in a network triad - Each node is associated to a set of nodes with a directed or reciprocal tie. (b) An example dyad and the features responsible for tie formation between users on Twitter. A probability score is assigned to each feature, to discover their inter-dependencies and enable reciprocal ties.

277 compute the overall similarity using the *Jaccard Index*,  $J$ :

$$J(\mathcal{X}_{f_{v_i}}, \mathcal{X}_{f_{v_j}}) = \frac{|\mathcal{X}_{f_{v_i}} \cap \mathcal{X}_{f_{v_j}}|}{|\mathcal{X}_{f_{v_i}} \cup \mathcal{X}_{f_{v_j}}|} \quad (1)$$

278 Moreover, modelling reciprocity or the response to a *friendship* request is associated with a  
 279 decision error, that can be quantified based on *response probability*. Response probability applied  
 280 to cases where a feature set for a decision involves a constant probability of error (Eq. 2) in the  
 281 choice [60]. Thus, the probability of *reciprocity* between pairs based on the similarities in their  
 282 features  $J(v_i, v_j)$  (Eq. 1) is given by [61]:

$$\epsilon_{v_i, v_j} = \frac{1}{\zeta \times (1 + \log(J(v_i, v_j) + \zeta))} \quad (2)$$

283 The *constant error term*,  $\zeta$ , is assigned the value of 1/3 and the final relation is given by:

$$p(R_{v_i, v_j}) = \frac{1}{1 + \exp \varphi} \quad (3)$$

284 where  $\varphi = -\log(\epsilon_{v_i, v_j} + J(v_i, v_j)) \times (\epsilon_{v_i, v_j} + J(v_i, v_j))$ . Any node pair  $v_i, v_j$  are structurally  
 285 similar or related  $\mathcal{S}_r$ , if their degree of reciprocity (Eq. 3) is greater than a predefined threshold  $\tau$ .  
 286 It follows that  $\mathcal{S}_r : \forall v_i \in \mathcal{S}_r \exists v_j$ , such that  $p(R_{v_i, v_j}) \geq \tau$ , where  $\mathcal{S}_r \subset \mathcal{D}$ .

287 *Collection of structurally related nodes.* Eq. 3 allows to identify as many node sets with a high  
 288 likelihood of establishing reciprocal ties, thus adding a layer of social cohesion to the MCT  
 289 framework. Identifying groups of structurally related nodes begins with a high-level aggregation  
 290 of nodes according to network size (for *network-communities*) and reciprocity (for *reciprocal-*  
 291 *communities*). This led to the question *what does it mean for nodes to be structurally related?*  
 292 As a simplistic example, consider a finite set  $\mathcal{V}_{13}$  that contains 13 nodes:  $\{v_1, \dots, v_{13}\} \in \mathcal{V}_{13}$ .  
 293 After executing *Algorithm f-sim* (Algorithm 1), which predicts the likelihood of reciprocity, the  
 294 following pairs of nodes are structurally-similar or related<sup>4</sup>:  $v_1 \sim v_2, v_1 \sim v_3, v_1 \sim v_5, v_2 \sim$   
 295  $v_4, v_2 \sim v_5, v_2 \sim v_9, v_3 \sim v_{11}$ . Accordingly, three structurally related communities can be  
 296 identified:  $c_1 = \{v_1, v_2, v_3, v_5\}$ ,  $c_2 = \{v_2, v_4, v_5, v_9\}$  and  $c_3 = \{v_3, v_9\}$ . The communities can  
 297 be expressed in a matrix form for spectral analysis. Matrix entries can be based on states, such  
 298 as the *reciprocity potential* of nodes defined as the ratio of *outdegree* over *network size*.

299 *Spectral analysis.* Since structurally related nodes can be easily transformed into a graph, we  
 300 apply *spectral analysis* to identify clusters and enabling *Sociometry*<sup>5</sup>, a means to measure social  
 301 relationships [62]. Spectral clustering involves operations ranging from the construction of

<sup>4</sup>The symbol ' $\sim$ ' is used to denote structural similarity between pairs.

<sup>5</sup>Metrics such as structural hole, homophily, and centrality metrics – degree, closeness, betweenness

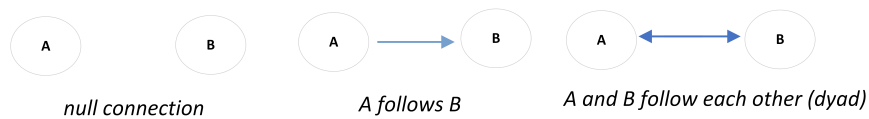


Figure 5: Examples of possible relationships between pairs

---

**Algorithm 1** : *Algorithm f-sim* returns the likelihood of reciprocity between pairs.

---

```

1: Initialisation:  $\{\} \leftarrow \mathcal{S}_r; \{\} \leftarrow \mathcal{S}_u$ 
2: Input: a finite collection of network data  $\mathcal{D}$ 
3: while  $\mathcal{D} \neq \emptyset$  do
4:    $\forall v_i, v_j \in \mathcal{D}$ , compute  $p(R_{v_i, v_j})$  using Eq. 3            $\triangleright v_i \neq v_j$ 
5:   if  $p(R_{v_i, v_j}) \geq \tau$  then                                $\triangleright \tau$ , a predefined threshold
6:      $\mathcal{S}_r \leftarrow (v_i, v_j)$                                   $\triangleright$  structurally related
7:   else
8:      $\mathcal{S}_u \leftarrow (v_i, v_j)$                                 $\triangleright$  structurally unrelated
9:   end if
10: end while
11: Output:
12:  $\mathcal{S}_r, \mathcal{S}_u, \mathcal{M}_{A_{i,j}}^{n \times n}$                               $\triangleright$  adjacency matrix  $\mathcal{M}_{A_{i,j}}^{n \times n}$ 

```

---

302 adjacency or affinity matrices to clustering in a reduced dimension [63]. We construct the  
303 *adjacency*, *similarity* and *degree* matrices based on ground-truth data and Eq. 3. The adjacency  
304 matrix  $\mathcal{M}_{A_{i,j}}$  (Eq. 4) encodes the structural similarities among node pairs. Presence or absence  
305 of reciprocity is marked with 1 or 0, respectively.

$$\mathcal{M}_{A_{i,j}} = \begin{cases} 1 & \text{if } p(R_{v_i, v_j}) \geq \tau \\ 0 & \text{if otherwise} \end{cases} \quad (4)$$

306 The *degree matrix*,  $\mathcal{M}_D$ , is a diagonal matrix obtained by summing the entries in Eq. 4 across  
307 the rows, in which entry  $i, i$  denotes the degree or number of edges to node  $i$ . Thus, each entry in  
308 the diagonal  $d_i$  (Eq. 5), of matrix  $\mathcal{M}_D$  is defined by:

$$d_i = \sum_{\{j | (i,j) \in E\}} p(R_{v_i, v_j}) \geq \tau \quad (5)$$

309 Subtracting the adjacency matrix,  $\mathcal{M}_{A_{i,j}}$ , from the degree matrix,  $\mathcal{M}_D$ , gives the graph *Laplacian*,  
310  $\mathcal{M}_L = \mathcal{M}_D - \mathcal{M}_A$ , whose *eigenvectors* and *eigenvalues* offer informative features for clustering.  
311 Diagonal entries in Eq. 6 denote the degree of nodes, off-diagonal negative entries ( $-p(R_{v_i, v_j})$ )  
312 represent edges between node pairs and zero entries signify no edges.

$$\mathcal{M}_{L_{i,j}} = \begin{cases} d_i & \text{if } i = j \\ -p(R_{v_i, v_j}) & \text{if } p(R_{v_i, v_j}) \geq \tau \\ 0 & \text{if otherwise} \end{cases} \quad \triangleright \text{edge}(i, j) \text{ exists} \quad (6)$$

### 313 3.1.2. Structural optimisation

314 Given a set of structurally related nodes, hidden local communities can be uncovered via  
315 matrix decomposition on the following matrices of interactions and corresponding dimensions:

- 316 -  $\mathcal{M}_{c_{ns}} \mapsto n \times p$ : a matrix of nodes according to a concept, e.g., network size
- 317 -  $\mathcal{M}_{c_{vr}} \mapsto n \times k$ : a matrix of nodes according to reciprocity
- 318 -  $\mathcal{M}_{c_{nr}} \mapsto p \times k$ : a matrix of high-level and local communities

319 The network-communities matrix,  $\mathcal{M}_{C_{ns}}^{n \times p}$ , is decomposed into its approximate constituents:

$$\mathcal{M}_{C_{ns}} \approx \mathcal{M}_{vr} \mathcal{M}_{C_{nr}}^T \quad (7)$$

320 Interpretability is desirable in the MCT framework, hence, the decomposition follows a non-  
 321 negative matrix factorisation (NMF) scheme [64]. NMF provides an intuitive factorisation, in  
 322 which non-negative constraints are imposed on the optimisation parameters [65].

$$\min_{\mathcal{M}_{vr}, \mathcal{M}_{C_{nr}}} \|\mathcal{M}_{C_{ns}} - \mathcal{M}_{vr} \mathcal{M}_{C_{nr}}^T\|_F^2, \text{ subject to } \mathcal{M}_{vr}, \mathcal{M}_{C_{nr}} \geq 0 \quad (8)$$

323 We simplify the notation of matrices as follows:  $\mathcal{M}_{C_{ns}} \mapsto D$ ,  $\mathcal{M}_{vr} \mapsto P = [p_{is}]$ ,  $\mathcal{M}_{C_{nr}} \mapsto$   
 324  $Q = [q_{js}]$ . The formulation in Eq. 8 allows to consider *Lagrangian relaxation* to optimise the  
 325 squared Frobenius norm ( $\|\cdot\|_F^2$ ) of the matrix. Consequently, NMF's non-negative constraints  
 326 are relaxed by introducing the *Lagrangian multipliers*, two new parameters ( $\alpha, \beta \leq 0$ ), to the  
 327 corresponding entries of the optimisation parameters ( $P, Q$ ). Accordingly, the objective function  
 328  $M_{sr}$  is expressed as a *minmax* problem, that requires a simultaneous minimisation over  $P, Q$  and  
 329 maximisation<sup>6</sup> over all applicable values of  $\alpha$  and  $\beta$ :

$$M_{sr} = \|D - PQ^T\|_F^2 + \sum_{i=1}^n \sum_{s=1}^k p_{is} \alpha_{is} + \sum_{j=1}^p \sum_{s=1}^k q_{js} \beta_{js} \quad (9)$$

330 The optimisation starts by computing the gradient of the relaxed Lagrangian, with respect to the  
 331 first aspect of the *minmax* (i.e. minimisation) optimisation variables. Although  $\alpha$  and  $\beta$  offers  
 332 a degree of flexibility (that comes at a cost), optimal solution requires optimisation conditions  
 333 to be based on  $P, Q$  only. We apply a handy technique based on Karush-Kuhn-Tucker (KKT)  
 334 optimality condition [66] to eliminate the Lagrangian multipliers. The KKT condition suggests  
 335 that  $p_{is} \alpha_{is} = q_{js} \beta_{js} = 0$ . Non-negative random values in  $(0, 1]$  are iteratively assigned to the  
 336 parameters  $P$  and  $Q$  based on the following update rule (after simplifying equations as shown in  
 337 the Appendix):

$$p_{is} \leftarrow \frac{(DQ)_{is} p_{is}}{(PQ^T Q)_{is}}, \quad q_{js} \leftarrow \frac{(D^T P)_{js} q_{js}}{(QP^T P)_{js}}, \quad \forall i \in \{1, \dots, n\}, j \in \{1, \dots, n\}, s \in \{1, \dots, k\} \quad (10)$$

338 The process of updating  $P$  and  $Q$  involves comparing their values to the original matrix  $D$ ,  
 339 aiming in minimising the difference or *error*. Parameters  $p_{is}$  and  $q_{js}$  are iteratively updated until  
 340 convergence. A successful iterative update process ensures that the underlying matrices exhibit  
 341 strong correlations among their respective entries.

### 342 3.2. Textual component

343 The textual component applies a form of document-pivot clustering based on weighted  
 344 features [32, 33, 34, 35]. Due to the volume of tweets and their short length, it is difficult to  
 345 gain a broad perspective about topics. Hence, using a single tweet may not provide sufficient  
 346 information. To understand the discussion topics and the degree of similarity among tweets,  
 347 a fixed number of tweets is extracted from nodes in the structurally related sets,  $\mathcal{S}_r$ . We

<sup>6</sup>This is needed because the Lagrangian multipliers are initialised with negative values.

348 utilise Latent Dirichlet Allocation (LDA), which has been previously applied for similar tasks  
 349 [67, 50, 68]. LDA is a probabilistic generative model that assigns word distributions to *topics* and  
 350 topic distributions to *documents* in a corpus, so that the documents represent random mixtures  
 351 over latent topics [49]. In this study, tweets collected from each node  $v_i \in \mathcal{S}_r$  define a corpus  
 352  $\mathcal{T}_{v_i}$ , whose overall theme is analysed for comparison with other nodes.

### 353 3.2.1. Modelling textual clusters

354 Identifying textually related nodes,  $\mathcal{T}_r$ , starts by aggregating a finite collection of textual  
 355 content,  $\mathcal{T}$ , from each node  $v_i$ . The collection of  $k$  tweets produced by node  $v_i$  over time,  
 356 i.e.  $\{t_{i1}, t_{i2}, t_{i3}, \dots, t_{ik}\} \in \mathcal{T}_{v_i}$  consists of a set of  $m$  n-gram features  $\{f_{i1}, f_{i2}, f_{i3}, \dots, f_{im}\} \in t_i \in$   
 357  $\mathcal{T}_{v_i}$ . Each stream of tweets is preprocessed to extract shingles<sup>7</sup> for transformation following the  
 358 *term-frequency-inverse document frequency (tf-idf)* weighting scheme [69]. The *tf-idf* vector of  
 359 a tweet,  $\mathbf{v}_i$ , can be normalised or not; we apply the  $L_2$  - *norm* given by:  $\mathbf{v}_i = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|_2}$ .

360 *Similar documents.* The collections of tweets for any node pair are:  $\mathcal{T}_{v_i}, \mathcal{T}_{v_j} \in \mathcal{T}_r \subset \mathcal{S}_r$ . The  
 361 aggregation scheme ensures that each node has a distinct fingerprint for comparison with others.  
 362 We train the LDA model so that each tweets in the corpus has a finite distribution over topics,  
 363 and topics have distributions over words. The distribution of each tweet, dubbed *anchor tweet*,  
 364 is used to compare with other tweets to locate the most similar tweets and generate relevant  
 365 matrices. Because LDA-based comparison relies on the probability distributions of tweets, we  
 366 apply Jensen-Shannon Divergence (JSD), a useful statistical metric, to measure tweet similarity  
 367 as the degree of divergence in the respective distributions. Unlike Kullback-Leibler Divergence,  
 368 *JSD* is symmetric, which is crucial in the task of comparing tweets, since similarity should be the  
 369 same irrespective of the order, i.e.,  $X \mapsto Y$  or  $Y \mapsto X$  be equal. For example, given two discrete  
 370 distributions  $X$  and  $Y$ , *JSD* is defined as:

$$JSD(X||Y) = \frac{1}{2}D(X||\mu) + \frac{1}{2}D(Y||\mu) \quad (11)$$

371 The *JSD* distance measure ( $JSD_{dist}$ ) is obtained by squaring its *divergence relation*:

$$JSD_{dist} = \sqrt{(JSD(X||Y))} \quad (12)$$

372 It follows that any pairs of tweets,  $t_i$  and  $t_j$ , are textually-similar or related,  $\mathcal{T}_r$ , if their similarity  
 373 degree,  $\phi$ , is greater than a predefined threshold  $\tau$ . Thus,  $\forall t_i \in \mathcal{T}_r \exists t_j : \phi(t_i, t_j) \geq \tau, \mathcal{T}_r \in \mathcal{S}_r$ .  
 374 Because a finite collection of tweets is extracted from each node in  $\mathcal{S}_r$ , each  $t_i \in \mathcal{T}_r$  consists of  
 375 a node and its set of tweets. LDA outputs come in a dense  $d \times t$  matrix  $\mathcal{M}_{lda}^{d \times t}$ , consisting of  $d$   
 376 tweets and their corresponding  $t$  topics. Moreover, two matrices apply to the textual component:  
 377 (a)  $\mathcal{M}_{vt} \mapsto m \times q$ : matrix of  $m$  nodes and top  $q$  topics, and (b)  $\mathcal{M}_{va} \mapsto m \times m$ : affinity matrix  
 378 of nodes according to topic similarity. Consequently, node communities are formed around  
 379 common discussion topics and the goal to cluster them according to topical similarities, as in  
 380  $tr(\mathcal{M}_{va} \mathcal{M}_{vt} \mathcal{M}_{vt}^T)$ . Algorithm 2 describes how to obtain the textually related clusters.

<sup>7</sup>Shingles are obtained by removing stopwords and other non-content bearing terms in a tweet.

---

**Algorithm 2** : *Algorithm text-sim identifies textually related clusters*


---

```

1: Initialisation:  $\{\} \leftarrow \mathcal{T}_r; \{\} \leftarrow \mathcal{T}_u$ 
2: Input: collection of structurally related nodes  $\mathcal{S}_r$ 
3:  $\forall v_i \in \mathcal{S}_r$ , get  $k$  texts  $g(\mathcal{T}_{v_i})$  ▷  $g(\mathcal{T}_{v_i})$  set of  $k$  texts of  $v_i$ 
4:  $\mathbf{x}_i \leftarrow t_i \in g(\mathcal{T}_{v_i})$  ▷ get texts vectors  $\mathbf{x}_i$ 
5: truncate  $\mathbf{x}_i$  ▷ retain  $b$  top terms in vector  $\mathbf{x}_i$ 
6:  $m(\mathcal{T}_{v_i}) = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|_2}$  ▷ mean of  $L_2$  normalised  $\mathbf{x}_i$ 
7:  $LDA(m(\mathcal{T}_{v_i}))$  ▷ invoke the LDA on  $m(\mathcal{T}_{v_i})$ 
8:  $\mathcal{T}_{sim}(\mathcal{T}_{v_i}, \mathcal{T}_{v_j}) = JS_{dist}(\mathcal{T}_{v_i} || \mathcal{T}_{v_j})$  ▷ get similar texts using Eq. 12
9: if  $\mathcal{T}_{sim}(\mathcal{T}_{v_i}, \mathcal{T}_{v_j}) \geq \tau$  then
10:   update  $\mathcal{T}_r$  ▷ textually related
11: else
12:   update  $\mathcal{T}_u$  ▷ textually-unrelated
13: end if
14: Output:
15:  $\mathcal{T}_r, \mathcal{T}_u, \mathcal{M}_{ta}^{m \times m}$  ▷ affinity matrix  $\mathcal{M}_{ta}^{m \times m}$ 

```

---

### 3.3. Microcosm detection algorithm

381 The problem of discovering community structures is modelled as a multilevel clustering  
382 task, in which nodes are grouped according to scoring functions. Using the affinity matrix based  
383 on  $\mathcal{S}_r$  or  $\mathcal{T}_r$ , various algorithms can be used to identify relevant partitions by optimising the  
384 separate and joint similarities of  $\mathcal{S}_r$  and  $\mathcal{T}_r$ :  $\psi(\mathcal{S}_r, \mathcal{T}_r) = \phi(\mathcal{S}_r) + \phi(\mathcal{T}_r)$ . A cohesive community  
385 is a collection of nodes,  $\mathcal{V}$ , with high degree of similarity structurally and textually. Thus, the  
386 microcosm detection problem can be formally defined as:  
387

388 given a collection of network data  $\mathcal{D}$ , defined by sets of nodes,  $\mathcal{V}$  and edges,  $\mathcal{E}$ , for  
389 each node  $v_i \in \mathcal{V}$  consisting of sets of *structural* and *textual* features, the goal is to  
390 identify a collection of highly cohesive sub-networks  $\mathcal{P}$ :

$$\mathcal{P} : \forall v_i \in \mathcal{S}_r \exists v_j : p(R_{v_i, v_j}) \geq \tau \text{ and } \forall t_i \in \mathcal{T}_r \subset \mathcal{S}_r \exists t_j : \phi(t_i, t_j) \geq \tau, \mathcal{P} \subset \mathcal{D}$$

391 The above formulation means that for all node pair in the partition  $\mathcal{P}$ , both the *structural* and  
392 *textual* similarities are greater than their respective threshold  $\tau$ <sup>8</sup>.

393 *Community of related nodes.* With  $\mathcal{S}_f$  and  $\mathcal{T}_f$  denoting feature sets of structurally and textually  
394 related nodes,  $\mathcal{M}_{S_f}^{n \times n}$  and  $\mathcal{M}_{T_f}^{m \times m}$  define an adjacency matrix of the structural component and  
395 an affinity matrix based on the textual similarity component, respectively. Therefore, for each  
396 matrix there exist community sets ( $\mathcal{K} \in \mathbb{R}^{n \times k}$ ,  $\mathcal{Q} \in \mathbb{R}^{m \times q}$ ), such that  $\{k_1, \dots, k_n\} \in \mathcal{K}$  denotes  
397 possible communities in  $\mathcal{M}_{S_f}$  and  $\{q_1, \dots, q_m\} \in \mathcal{Q}$  denotes communities in  $\mathcal{M}_{T_f}$ . For a matrix  
398 of reciprocal relationships  $R \subseteq \mathcal{V} \times \mathcal{V}$ <sup>9</sup> and the associated network data  $\mathcal{D}$  ( $\mathcal{V}, R \in \mathcal{D}$ ), there  
399 exist numerous communities  $\{c_1, \dots, c_k\} \in \mathcal{C}$  such that  $\emptyset \subset c_i \subseteq \mathcal{V}$  and  $\mathcal{C}$  denotes a community  
400 set. With any pair of similar nodes denoted by  $v_i \sim v_j \iff \exists c_i \in \mathcal{C} : v_i, v_j \in c_i$ , a more  
401 socially cohesive node community is formed by identifying overlapping nodes in both  $\mathcal{K}$  and  $\mathcal{Q}$ ,  
402 through a repetitive partition optimisation. Accordingly, the MCT framework contributes into

---

<sup>8</sup>For all experiments,  $\tau \geq 0.5$ , i.e. pairs are considered similar (1) if  $\tau \geq 0.5$ , otherwise dissimilar (0).

<sup>9</sup> $R$  applies to both  $\mathcal{M}_{S_f}$  and  $\mathcal{M}_{T_f}$ .

---

**Algorithm 3** : *Algorithm MCT* identifies local communities known as *microcosms* in a network.

---

1:	<b>Initialisation:</b> $\{\} \leftarrow \mathcal{S}_r, \{\} \leftarrow \mathcal{S}_u, \{\} \leftarrow \mathcal{T}_r, \{\} \leftarrow \mathcal{T}_u$	
2:	<b>Input:</b> a collection of network data $\mathcal{D}$	
3:	structural-component:	$\triangleright$ invoke <i>f-sim</i> (alg. 1)
4:	$\text{f-sim}(\mathcal{D}) \mapsto \{\mathcal{S}_r, \mathcal{S}_u\}, \mathcal{M}_{A_{i,j}}^{n \times n}$	$\triangleright$ alg. 1 output
5:	textual-component:	$\triangleright$ invoke <i>text-sim</i> (alg. 2)
6:	$\forall v_i \in \mathcal{S}_r$ get $k$ tweets	$\triangleright$ set of texts $\mathcal{T}_{v_i}$
7:	$\text{text-sim}(\mathcal{S}_r) \mapsto \{\mathcal{T}_r, \mathcal{T}_u\}, \mathcal{M}_{ta}^{m \times m}$	$\triangleright$ alg. 2 output
8:	compare all topics ( $\mathcal{T}_{v_i} \in \mathcal{S}_r$ ) using Eq. 11	$\triangleright$ affinity matrix
9:	local clusters:	
10:	$\psi(\mathcal{S}_r, \mathcal{T}_r)$	$\triangleright \mathcal{S}_r, \mathcal{T}_r \geq \tau$
11:	<b>Output:</b>	
12:	$\mathcal{C}_{c_i}^{m \times p}$	$\triangleright$ local communities

---

403 two operational categories: (1) *optimising matrices of values* and (2) *optimising intra-cluster*  
404 *similarity*. The MCT can be considered as a multivariate function, made up by structural and  
405 textual components, allowing to define an objective function that maximises the overall joint  
406 similarity.

### 407 3.3.1. *Optimising matrices of values*

408 Recall that the set of textually related nodes  $\mathcal{T}_r$  is a subset of the structurally related nodes  $\mathcal{S}_r$   
409 ( $\mathcal{M}_{c_{vr}}$ ), i.e  $\mathcal{T}_r \subseteq \mathcal{S}_r$ . Since the optimisation goal is to *maximise*  $\mathcal{T}_r$  ( $\mathcal{M}_{vt}$ ), the two are equated  
410 under the constraint:  $\mathcal{M}_{vt} = \mathcal{M}_{c_{vr}}$ , such that  $\mathcal{M}_{vt} - \mathcal{M}_{c_{vr}} = 0$ . Noting the constraint, the  
411 simplified representation used in Eq. 9 also applies to  $\mathcal{M}_{vt}$ , given by  $\mathcal{M}_{vt} = \mathcal{M}_{c_{vr}} = P$ , to  
412 achieve the maximum values possible by determining the extremum of the function. Thus, the  
413 goal is to maximise the joint models under the constrained function according to Eq. 13<sup>10</sup>:

$$\psi_{st} = \|D - PQ^T\|_F^2 + \sum_{i=1}^n \sum_{s=1}^k p_{is} \alpha_{is} + \sum_{j=1}^p \sum_{s=1}^k q_{js} \beta_{js} - \lambda Tr(\mathcal{M}_{vt}^T \mathcal{M}_{va} \mathcal{M}_{vt}) \quad (13)$$

### 414 3.3.2. *Optimising intra-cluster similarity*

415 Intra-cluster similarity optimisation is similar to the approach in Section 3.3.1 through  
416 the use of value matrices, but different objective function. The approach in Eq. 9 and the  
417 corresponding update rule (Eq. 10) are based on a matrix factorisation, which poses challenges  
418 with respect to exact or one-one mapping to the textually related clusters ( $\mathcal{T}_r$ ). We know that  
419 the two are related at a higher level, since  $\mathcal{T}_r \subset \mathcal{S}_r$ , but the details about the shared clusters  
420 are not fully established. To address this challenge, we propose the following approach based  
421 on the node similarity. Information about similar nodes is stored in the nodes' affinity matrix  
422 ( $\mathcal{M}_{va}^{n \times n}$ ), in which the magnitude of pairwise similarity decides entries in the matrix. Nodes are  
423 assigned to communities based on their degree of similarity denoted by  $\mathcal{M}_{C_{vr}}^{n \times k}$  (of  $n$  nodes and  $k$   
424 reciprocal-communities). For example,  $\{c_{vr_1}, \dots, c_{vr_k}\} \in \mathcal{C}_{vr}$  represents a set of nodes-reciprocal  
425 communities and membership in a cluster is qualified by Eq. 1 - 3. With a higher probability of  
426 forming a tie for nodes in the same cluster, community detection is based on optimising the joint

---

<sup>10</sup>The proportionality constant  $\lambda$  in Eq. 13 denotes a Lagrange multiplier.

---

**Algorithm 4** : *Algorithm MCT-2* identifies local communities known as *microcosms* in a network.

---

```

1: Input: a collection of network data  $\mathcal{D}$ 
2: structural-component:
3:   f-sim( $\mathcal{D}$ )  $\mapsto$   $\{\mathcal{S}_r, \mathcal{S}_u\}, \mathcal{M}_{sa}^{n \times n}$ 
4:   textual-component:
5:      $\forall v_i \in \mathcal{S}_r$  get k tweets
6:     text-sim( $\mathcal{S}_r$ )  $\mapsto$   $\{\mathcal{T}_r, \mathcal{T}_u\}, \mathcal{M}_{ta}^{n \times n}$ 
7:     compare all topics( $\mathcal{T}_{v_i} \in \mathcal{S}_r$ ) using Eq. 11
8: Clusters initialisation:
9:   select four random seed nodes:  $v_i, v_j, v_k, v_l \in \mathcal{S}_r$  and  $v_i, v_j, v_k, v_l \in \mathcal{T}_r$ 
10:  compute pairwise similarities among  $v_i, v_j, v_k, v_l$  using  $\psi_{s_r, t_r}(v_i, v_j)$ 
11: if  $\mathcal{T}_{sim}(\mathcal{T}_{v_i}, \mathcal{T}_{v_j}) \geq \tau$  then
12:   create single cluster  $\mathcal{C}_{ij}$ 
13: else
14:   create two clusters  $\mathcal{C}_i, \mathcal{C}_j$ 
15: end if
16: repeat 9 – 15 until  $|\mathcal{C}_{ik}|_{k=1}^M = M$   $\triangleright$  maximum clusters  $M$ 
17: Assign nodes to clusters:
18:    $\forall v_i \in \mathcal{S}_r$  compute similarity with cluster's mean
19:    $max_{\phi(v_i, \mu_{\mathcal{C}_i})}$   $\triangleright$  assign  $v_i$  to the most similar  $\mu_{\mathcal{C}_i}$ 
20: update cluster's mean:  $\mu_{\mathcal{C}_i} \leftarrow \mu_{\mathcal{C}_i}$ 
21: Output:
22:   local communities

```

---

427 similarities of  $\mathcal{S}_r$  and  $\mathcal{T}_r$ :

$$\psi_{st}(v_i, v_j) = (\lambda) \cdot \mathcal{S}_r(v_i, v_j) + (1 - \lambda) \cdot \mathcal{T}_r(v_i, v_j) \quad (14)$$

428 The goal of Eq. 14 is to maximise the joint similarity between  $\mathcal{S}_r$  and  $\mathcal{T}_r$  according to an  
429 aggregation criterion inspired by [70], based on the similarity scores between pairs and a user-  
430 defined balancing parameter<sup>11</sup>  $\lambda$ , with values in  $(0, 1)$ . We follow the approach in [71] to find  
431 optimum value for the  $\lambda$ . Algorithm 4 describes how nodes are assigned to relevant clusters until  
432 the stopping criterion, a user-defined integer  $M$  signifying the desired number of clusters, is  
433 reached.

## 434 4. Experimentation

435 This section presents our experimentation to evaluate the MCT against other existing methods.

### 436 4.1. Datasets

437 We utilise the following diverse datasets for the experimentation.

#### 438 4.1.1. Ground-truth and predicted data

439 Unlike previous studies in which datasets from various social networks were collected  
440 [31, 72, 73], this study focuses on nodes with reciprocal, not directed, ties. The reciprocal

---

<sup>11</sup>Note that this is different from the one used in optimisation based on matrices of values.



---

**Algorithm 5** *Algorithm search-dyads* profiles users with directed and undirected ties on Twitter

---

```
1: Initialisation:  $1 - edge \rightarrow \{\}, dyads \rightarrow \{\}$ 
2: Input: begin with an arbitrary set of seed users, say  $k$ 
3: while  $k \neq \emptyset$  do
4:    $\forall v_i \in k$ , get sets of friends  $fr_{v_i}$ , followers  $fl_{v_i}$ ,  $fr_{v_i}, fl_{v_i} \in m_{v_i}$ ;  $m_{v_i}$  denotes  $v_i$  network
5:    $\forall v_j \in fr_{v_i}$ , retrieve the sets  $fr_{v_j}$  and  $fl_{v_j}$ ,  $fr_{v_j}, fl_{v_j} \in m'_{v_j}$ ;  $m'_{v_j}$  denotes  $v_j$  network
6:   if  $v_i \in fr_{v_j}$  then
7:      $v_i \sim v_j$   $\triangleright$  both follows one another
8:     update dyads
9:   else
10:     $v_i$  follows  $v_j$ 
11:    update 1-edge
12:   end if
13: end while
```

---

441 collection consists of *dyadic* and *transitive* datasets, which were collected using Twitter’s  
442 Application Programming Interface (API) according to Algorithm 5. The process returns a  
443 collection of *tweet objects*, a complex object with many descriptive fields, which allows to  
444 extract structural and textual components for analysis. The collection begins with a search on the  
445 network profile of each from a finite set of seed users<sup>12</sup>, or a network composition  $m_{v_i}$ , consisting  
446 of lists of friends  $fr_{v_i}$  and followers  $fl_{v_i}$ , to determine user pairs that follow each other. The set  
447 of reciprocal pairs is denoted by  $\kappa \in m_{v_i}$  and the transitive dataset is a scaled-version of dyadic  
448 data.

449 In addition to the collection of nodes with actual pairwise ties (denoted as G-pTie in Table 4.1),  
450 the ground-truth dataset also consists of public data associated with COVID-19 outbreak (G-  
451 pMention) related to aspects of scepticism and myths about the pandemic [75]. The data contains  
452 two broad categories: information put forward by credible sources, such as the World Health  
453 Organisation (WHO), and information from users dismissing WHO’s guidelines on combating  
454 the pandemic. The dataset consists of interaction information about users who mention each  
455 other. Nodes with frequent mentioning are highly likely to be in the same community. For the  
456 dataset consisting of predicted pairwise ties (P-pTie), a reciprocal tie exists between  $v_i$  and  $v_j$  if  
457  $p(R_{v_i, v_j}) \geq \tau$ , otherwise just a directed tie. In Table 2, SND1 refers to synthetic network data  
458 generated based on LFR approach (see Section 4.3.2 for details).

#### 459 4.1.2. Public datasets

460 To reinforce evaluation and generalisation, we use the following collections of publicly  
461 available datasets. The datasets consist of real-world networks commonly used for community  
462 detection. Essentially, the following datasets have been used: Zachary’s karate club [76], dolphin  
463 social network [77], political blog dataset [78] and Ego-network, consisting Facebook and  
464 Twitter datasets [79]. The Facebook data contains anonymised ‘*circles*’ or ‘*friends lists*’, and  
465 *node features (profiles)*. Each node has *node ids*, sets of *connections* or *edges*, and *anonymised*  
466 *features* encoding information about its *circle*. The Facebook data allows to explore communities  
467 using each user’s network circle in terms of size and diversity of membership. Moreover, the

---

<sup>12</sup>Seed users are verified or unverified accounts devoid of spammers or social bots were collected by SPD filtering [74]. A ‘*list*’ on Twitter allows a user to store a set of preferred users and can be used to obtain relevant information.

468 collection consists of *synthetic data*, which is based on the approach proposed in [71] to generate  
 469 synthetic network with known parameters. The synthetic nature of the networks makes it possible  
 470 to explore the parameters' space for the best community structure in the network. Table 2 shows  
 471 basic statistics of datasets used in this study.

#### 472 4.2. Meta-analysis

473 Owing to the prevalence of unreciprocated and event-type ties on Twitter, we conjecture  
 474 that mining tasks, such as community detection, are less effective and more challenging. In  
 475 this section, our goal is to apply a pragmatic approach that provides a statistical analysis of  
 476 relevant metrics in the datasets to identify strongly correlated node attributes (Figure 4(b)) with  
 477 reciprocity among nodes. The empirical cumulative distribution function (ECDF) gives the  
 478 probability of a quantity evaluated at arbitrary points. We use it to analyse observations, such as  
 479 the variation of dyads or Simmelian ties across user categories or network size.

Table 2: A summary of microcosms detection datasets. V and E denote the node and edge size, respectively. G-pTie and P-pTie denote groundtruth and predicted sets of users with pairwise connectivity; G-Mention denotes collection of users with pairwise mentioning;  $\mu_{deg.}$  refers to the average degree in each data category.

Category	V	E	Description
G-pTie	18973	15538	$\mu_{deg.} = 1.6379$
P-pTie	15038	1298998	$\mu_{deg.} = 172.7621$
G-pMention	514	259	$\mu_{deg.} = 1.0078$
Karate club	34	78	Consists of 2 groundtruth communities
Dolphin	62	—	Consists of 2 groundtruth communities
Pol. Blog	1224	—	Consists of 2 groundtruth communities
ego-Facebook	4039	88234	$\mu_{deg.} = 43.6910$
ego-Twitter	81,306	1768149	—
SND1	1000	—	—

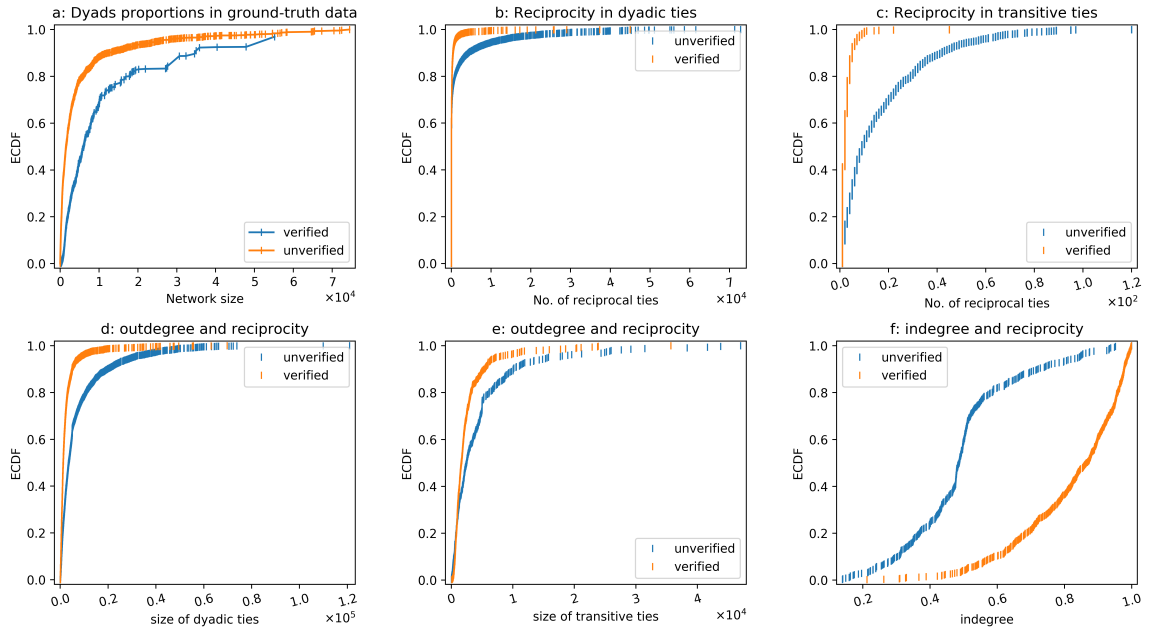


Figure 6: Sub-figures (a), (b) and (c) show the reciprocity effect on nodes with many dyadic relationships. Sub-figures (d), (e) and (f) show outdegree-reciprocity and indegree-reciprocity relationship in the ground-truth data.

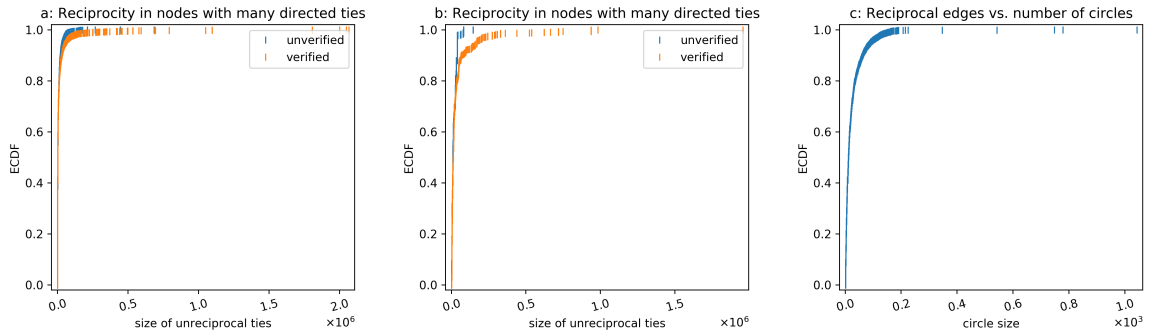


Figure 7: Relationship between the number of directed ties and reciprocity across different nodes in the datasets.

#### 4.2.1. Proportion of reciprocal units

Noting the flexibility of connections and the rarity of reciprocal links on Twitter, large scale dyadic ties are rare and difficult to locate. Using Algorithm 5, we collected *directed* or *1-edge* and *undirected* data, and examined the *network topology* of each category and its utility in the detection of local communities. In Figure 6, there is a high proportion of reciprocity in unverified users in comparison to the verified counterpart. The reciprocity ability slightly decreases with increasing network size of the user, which can be attributed to the difficulty in keeping track of and responding to all followership requests. Sub-figures 6(a) and (b) show the relationship between reciprocity and the number of reciprocal ties. While there is higher reciprocity in the unverified category, the verified category shows almost 100% reciprocity with a relatively small network size. Sub-figure 6(c) shows similar behaviour in transitive ties, but is more evident in the unverified users category. Similarly, sub-figures 6(d),(e) and (f) show the relationship between outdegree and reciprocity and indegree and reciprocity in the ground-truth data. The behaviour resembles an inverse relationship in which reciprocity decreases with increasing outdegree (sub-figures 6(d) and (e)). Sub-figure 6(f) shows an almost linear relationship between indegree and reciprocity, especially among the unverified category. In the verified category, the effect is low and seems to shoot once the network size increases (vis-a-vis indegree/number of followers). There is an instant reciprocity in the unverified users, which can be explained by suggesting that the users are interested in expanding their network. Figure 7 shows the relationship between the number of directed ties and reciprocity across user categories in the data. The results demonstrate that verified users have many directed nodes or unreciprocated ties but with less reciprocity. This observation holds for nodes with many dyadic and transitive ties in the data.

#### 4.3. Evaluation

To ascertain the efficacy and relevance of the study’s output, evaluation entails thorough analyses and comparison with relevant baselines drawn from the literature. Quantitative analyses of experimentation on various datasets using the baseline algorithms is involved. Other forms of evaluation are specific to the structural and textual levels of the MCT strategy. The evaluation process aims to: (a) investigate the effect of utilising structurally-related nodes in identifying local communities in social networks, (b) compare structurally-related clusters with textually-related clusters, and (c) evaluate the performance of MCT in comparison with baseline models.

### 510 4.3.1. Evaluation metrics

511 This section discusses quantitative measures for validating the performance of MCT and  
512 baseline methods. Because of the multilevel approach, the metrics are suitable for evaluating  
513 network structure (structural clusters) and textual-clusters (roughly considered as labels).

514 *Clustering coefficient and Community cohesion.* Clustering coefficient,  $C_{coeff}$ , is used to  
515 quantify the clustering tendency of a given node in relation to other nodes within a network  
516 [38]. Computing  $C_{coeff}$  requires:  $edges = \frac{k_i(k_i-1)}{2}$  and  $C_{coeff_i} = \frac{2E_i}{k_i(k_i-1)}$  where  $i, k_i, E_i$  denote  
517 a network node, the number of edges connecting  $i$  to  $k_i$  other nodes in the network, and the  
518 actual number of existing edges between  $k_i$  nodes, respectively. The ratio  $E_i \propto \frac{k_i(k_i-1)}{2}$  defines  
519 the clustering coefficient of a node. *Community cohesion* demonstrates the level of connectivity  
520 within a community and is captured by measuring the degree of cohesiveness. Due to the presence  
521 of a strong connectivity among nodes, a well-connected community is intuitively difficult to  
522 divide into sub-communities [80]. Any useful metric that reveals the degree of cohesion can be  
523 used to evaluate cohesiveness, i.e., if the community is well-connected and difficult to partition.  
524 In this study, cohesiveness is measured by the degree of similarities  $\mathcal{S}_r$  and  $\mathcal{T}_r$ . We compute the  
525 *average degree* ( $\mu_{deg}$ ), defined as the average node degree to other member nodes [81]. Moreover,  
526 we use the *accuracy metric*, i.e., the fraction of predicted labels to the total number of data points.

527 *Modularity and NMI.* Modularity,  $\mathbf{Q}$ , measures the strength of communities, as the number of  
528 edges falling within groups minus the expected number in an equivalent network with edges  
529 placed at random [42]. Usually,  $\mathbf{Q} > 0$  signifies the possible presence of a community structure  
530 and the higher the values the better [41]. Normalised Mutual Information (NMI) is another  
531 statistical tool to evaluate the quality of network clusters [82]. NMI measures the degree of  
532 agreement between network partitions, based on the assumption that each node in a community,  
533  $v_i \in \mathcal{V}$ , is associated with both the *true community* and the *predicted community*, such that  
534  $l_{v,p} = i$  defines the predicted community  $i$  of a node [83]. Furthermore, we apply *Rand* and  
535 *Jaccard* similarity metrics, which are based on tracking both correctly and incorrectly classified  
536 pairs of nodes, especially in groundtruth datasets.

### 537 4.3.2. Baseline Models

538 For evaluation, *MCT* is applied alongside the following detection algorithms with different  
539 modes of operation on the datasets described in Section 4.1 to identify local community structures.

540 *Girvan-Neuman (G-N) and Label propagation (LP).* The G-N algorithm assumes that a commu-  
541 nity detection algorithm can naturally detect divisions among vertices without external influence  
542 or imposed restrictions on the divisions [5]. Accordingly, Girvan and Newman [84] proposed  
543 the iterative G-N algorithm that progressively removes network edges based on betweenness,  
544 a metric to quantify traffic flow among nodes. Each node’s betweenness score dictates which  
545 edge to remove. The most critical nodes are likely to experience high traffic flow, hence will  
546 possibly create a bottleneck. The LP algorithm is an iterative clustering method that converts  
547 unlabelled data to labelled given an initial seed of labelled data. Labelling involves a repetitive  
548 random node reshuffling and tagging with the most frequent label among its neighbours until  
549 convergence [85]. The labelled data information is then propagated across the whole network.

550 *Synthetic Network Model.* This is achieved using the widely used approach, or LFR model,  
 551 proposed in [86] to generate synthetic networks with planted partitions or community structures.  
 552 For a given network  $G$  generated via the LFR, the following basic model’s parameters are  
 553 defined:  $\gamma, \beta, \bar{d}, \hat{\mu}$  denoting exponents of the power-law degree distribution, community size  
 554 distribution, mean degree and mixing parameter, respectively. Accordingly, the model ensures  
 555 that nodes’ degrees are sampled independently whose distribution exhibit power-law behaviour  
 556 and the mixing parameter,  $\hat{\mu}$ , to distribute nodes’ indegree and outdegree such that  $1 - \hat{\mu}$  and  
 557  $\hat{\mu}$  denote the proportions of edges shared with nodes in the same and different communities,  
 558 respectively. The SND1 network in Table 2 is generated based on the LFR approach. The  
 559 network consists of 1000 nodes,  $\gamma = 1.5, \bar{d} = 15, \gamma_C = 0.8, C_{min} = 30, C_{max} = 300$ , and the  
 560 mixing parameter,  $\hat{\mu}$ , sampled from 0.1, 0.01, 0.2, 0.3, 0.5, 0.7, 0.9, 1. Because the parameters  
 561 pertaining the network and the embedded community structure are known, relevant community  
 562 detection methods should be able to detect or identify values (especially for the community) that  
 563 approximate such parameters.

564 The Planted Partition Model (PPM) is a form of likelihood optimisation algorithms that  
 565 are commonly used for community detection task. Due to their mathematical efficacy, many  
 566 algorithms are defined based on relevant assumptions about the underlying structure in the  
 567 network. Under this approach, a network is a composition of communities, which are used to infer  
 568 the network [87]. The PPM relies on the community membership of nodes to probabilistically  
 569 decide whether any pairs of nodes are connected. We apply variant of the PPM (degree-corrected  
 570 planted partition model [88]) and an extended version of the LFR model proposed in [71] as part  
 571 of the evaluation.

#### 572 4.4. Detection of community structure

573 In this section, we focus on the detection of community structures using our proposed  
 574 method<sup>13</sup>, introduced in Section 3.3, and the baseline models, described in Section 4.3.2. The  
 575 detection process consist of four steps: (1) retrieve a set of nodes with reciprocal ties on Twitter,  
 576 (2) compute the similarity proportion between pairs, using Algorithm 1, (3) compare prediction  
 577 accuracy using the ground-truth, and (4) perform clustering for community detection.

##### 578 4.4.1. Effectiveness of tie prediction

579 Using Algorithm 1, which computes the similarity between the corresponding features of  
 580 any pairs of nodes, we report its efficacy in the prediction pipeline. Due to the availability  
 581 of empirical data, the effectiveness of the model is quantified with respect to the degree of  
 582 conformity with the ground-truth data. This is vital because the tie prediction segment is not  
 583 relevant if it does not add value to the overall detection framework. The accuracy of the prediction  
 584 is obtained by computing the ratio of predicted reciprocal ties to true reciprocal ties. The best  
 585 result achieved is .608 accuracy; depending on the threshold  $\tau$ , the accuracy may be lower or  
 586 higher. Sub-figure 8(c) shows possible values of  $\tau$  and the corresponding accuracy.

##### 587 4.4.2. Community structure

588 We examine how the use of a collection of structurally-related nodes affects community  
 589 detection, and compare performance. For the experiments, we apply the proposed method,

---

<sup>13</sup>relevant datasets and implementation code available at <https://github.com/ijdutse/mct>

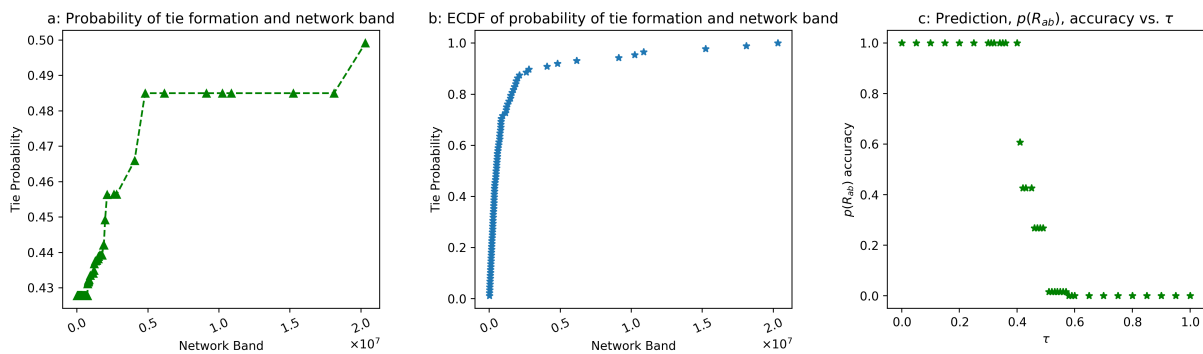


Figure 8: Sub-figures *a* and *b* show the probability of tie formation as a function of *network size*; there is a high chance of reciprocating a tie among users in a network band of  $0.5 \times 10^7$ . Sub-figure *c* shows the prediction accuracy versus the threshold value in *f-sim* (Algorithm 1). The prediction accuracy is almost 100% when the value of the threshold is low; conversely, the accuracy is almost 0% when the threshold value is very high. A *switch-point* can be observed toward the midpoint in which the accuracy is just above 60% at a threshold value of about 0.41. With additional features, the prediction can be improved. For instance, the inclusion of a description feature led to a significant improvement; however, it requires training on a large corpus to obtain the embedding of terms in text.

590 (MCT), and the baselines, G-N [84] and LP [85]. Table 3 shows the results of applying the  
 591 community detection algorithms on the data according to the evaluation metrics described in  
 592 Section 4.3.1. Although all the algorithms detected community structures, there are quantitative  
 593 variations among the outcomes. Our analysis is along the following dimensions.

594 *Effect of datasets:* All algorithms perform best on the ground-truth data, followed by *ego-*  
 595 *Facebook*, then *predicted*, and worst on *ego-Twitter*. The *ego-Facebook* data consists of nodes  
 596 with reciprocal ties, but the textual feature set is small, making it less complex than the other  
 597 datasets. We consider *homophily* and *structural equivalence* as precursors of communities, in  
 598 which nodes with similar profiles or social status are more likely to interact and establish a  
 599 small community. For instance, sub-figures 8(a) and (b) show homophily as a form of structural  
 600 equivalence based on network size and indegree for examining the probability the formation of  
 601 an edge. Sub-figures 8(a) and (b), depicts a behaviour that resembles an inverse relationship:  
 602 increase in network size results in decrease in reciprocity.

603 *Effect of models:* Table 3 also demonstrates the performance of each model. The *MCT* results  
 604 indicate a more localised structure noting the magnitude of *Q*, *NMI* and the number of detected  
 605 communities (*#DC*) with respect to the *ground-truth* data. We attribute the improvement to the  
 606 use of in-depth structural features that introduce a connectivity layer. *MCT* explores the data  
 607 for community structures at local and global level through a high-level grouping of nodes into

Table 3: Results of experiments on three datasets for community detection using algorithms based on structural properties. G-N: Girvan–Neuman, LP: Label Propagation, MCT: Multilevel Clustering Technique, #DC: Number of Detected Communities

Dataset	G-N			LP			MCT		
	Metric		#DC	Metric		#DC	Metric		#DC
	Q	NMI		Q	NMI		Q	NMI	
G-pTie	.908	.794	308	.77	.602	1319	.915	.791	263
ego-Twitter	.334	.197	1431	.215	.131	2131	.307	.230	1131
ego-Facebook	.522	.590	1037	.421	.304	1780	.503	.372	1845
P-pTie	.473	.311	1107	.360	.267	2071	.601	.472	985

608 communities, according to the network size and the recognition of bi-modal information sources.

609 For the parameterised approach, groundtruth datasets form the basis of the evaluation.  
 610 Therefore, we discuss the results obtained in Table 4 using Rand and Jaccard scores as the  
 611 evaluation metrics. Generally, the results in the Table indicate good performance, especially with  
 612 respect to the Rand score and on small datasets such as the Karate club. The values associated  
 613 with the MCT signify better performance across all the datasets. However, there are instances  
 614 where the algorithm’s performance lags behind. For instance, the PPM and ILFR perform better  
 615 on the Karate and Pol. Blogs datasets, which we attribute to the small size nature of the datasets.  
 616 Moreover, there is a significant improvement on performance on the synthetic datasets, i.e. SND1.  
 617 This is expected since the network consists of well defined community structure. A common trait  
 618 among the algorithms is that they perform poorly based on the Jaccard index, suggesting that the  
 619 metric is somewhat strict or further optimisation is needed.

## 620 5. Discussion

621 In this section, we discuss some significant observations from the study.

622 *Impact of reciprocal units and text aggregation for clustering.* One of the assumptions of this  
 623 is that recognising a set of reciprocal units for community detection offers a more cohesive  
 624 community representation. Since small groups allow modular analysis of social networks  
 625 [89, 90], we examined reciprocal ties, dyadic and Simmelian, as the basic units of relational  
 626 interaction on Twitter. However, Twitter’s flexible and eccentric connections entangle locating  
 627 nodes with reciprocal links. Structural similarity allows to organise nodes into connected clusters  
 628 and simplifying community detection. Structurally similar nodes are more likely to connect and  
 629 belong to the same community. The high volume and small size of tweets make comparisons of  
 630 discussions context challenging. Because a single tweet may not yield enough information about  
 631 a discussion, we need to balance between quantity and quality. We collected a finite set of tweets  
 632 from each node  $v_i$  that defines a user corpus  $\mathcal{T}_{v_i}$ , and computed its overall theme to compare  
 633 with other nodes. Textually-related nodes  $\mathcal{T}_r$  are identified by a topic modelling technique that  
 634 compares the similarities of the discussion topics of structurally similar nodes.

635 *Improving social cohesion in the detection task.* Online content increases rapidly in volume and  
 636 complexity and is dominated by influential users. These facts make the detection of socially  
 637 cohesive groups on Twitter challenging. With respect to *sociometry*, the formation of a social tie  
 638 can be based on *event-type* or *state-type ties*. The size of a network and the size of communities

Table 4: Results of experiments on three datasets for community detection using relevant algorithms and evaluation metrics.

	Model Metric	PPM		LFR		MCT		ILFR	
		Rand	J	Rand	J	Rand	J	Rand	J
Datasets	G-pMention	0.501	0.005	0.501	0.005	0.662	0.169	0.501	0.005
	Karate club	0.734	0.475	0.734	0.475	0.701	0.461	0.720	0.425
	Dolphin	0.651	0.379	0.581	0.255	0.696	0.340	0.536	0.167
	Pol. Blob	0.903	0.821	0.878	0.775	0.565	0.487	0.581	0.166
	G-pTie	0.621	0.012	0.621	0.012	0.629	0.210	0.621	0.012
	P-pTie	0.652	0.001	0.652	0.001	0.696	0.237	0.652	0.001
	SND1	0.846	0.431	0.795	0.395	0.879	0.510	0.601	0.317
	ego-Facebook	0.683	0.317	0.597	0.301	0.697	0.332	0.579	0.298

639 are almost linearly correlated. Similarly, the size of a network is inversely correlated with its  
640 degree of homogeneity. The degree of interaction is higher among structurally similar users.  
641 Often, users that discuss with and mention each other are engaged in reciprocal ties, showing  
642 strong social cohesion. Based on the idea of social *homophily*, users with many reciprocated ties  
643 are crucial in analysing socially cohesive groups. Figure 1 shows that communities on Twitter  
644 can be formed in many ways. A bi-modality approach differs across networks with respect  
645 to the depth of the features associated with the structural and textual modalities [29, 31, 16].  
646 Bi-modalities, e.g., network structure, features and attributes of nodes, lead to better and more  
647 interpretable community detection results. In Twitter, the structural component is not fully  
648 captured as it relies on directed connections. *MCT* exploits the usability of features in the  
649 detection of a local community through the impact analyses of both modalities, especially the  
650 structural one. We have shown that a structural component is useful in community detection  
651 and has minimal practical requirements. *MCT* offers a compact way to find and represent  
652 co-occurring users or user groups, allowing to explore local and global clustering requirements.

## 653 6. Conclusion

654 Many natural networks exhibit a certain degree of organisation, in which node groups form  
655 tightly connected units called communities. Community detection allows to understand the  
656 network structure and extract useful information. Detecting socially cohesive communities  
657 on Twitter is still challenging. While many methods have been proposed, they often discover  
658 disparate communities, likely to be socially unrelated. We observed that the topology of  
659 eccentric connections contributes to the detection of socially unrelated users and encourages the  
660 propagation of spurious content. Consequently, we propose a *multilevel clustering technique*  
661 (*MCT*) to identify socially cohesive user groups, i.e. *microcosms*, on Twitter.

662 The proposed *MCT* framework, jointly modelling structural and intrinsic textual features,  
663 contributes toward a methodological paradigm for cohesive community detection in a dynamic  
664 and heterogeneous social media. This is important because until recently, community detection  
665 algorithms focused on single modality, e.g. using node attributes or connectivity. Recent studies  
666 that combine information modalities are limited in capturing the nuances and intricate connection  
667 structure in platforms, such as Twitter. To improve the identification of socially cohesive  
668 communities, *MCT* offers a scalable detection strategy. The approach addresses the problem of  
669 structurally unrelated users, by adding a layer of social cohesion to existing community detection  
670 methods. In summary, *MCT* contributes: (1) a systematic exposition of community detection or  
671 clustering algorithms, (2) an in-depth utilisation of the bi-modality for community detection, and  
672 (3) detection of network communities at various levels.

673 A note on the proposed method's complexity is in order here. When the network size is  
674 huge, it is challenging to authoritatively specify when a given community detection algorithm  
675 will converge. Thus, we rely on a single iteration to analyse the algorithm's complexity, which  
676 will provide insights to its future performance. Let us assume that the execution complexity of  
677 a basic parameterised algorithm is  $f(C)$ , then the term  $O(f(C) \times s \times r \times m)$ , where  $s$  is the  
678 number of comparisons in deciding the next cluster,  $r$  is the size size and  $m$  the number runs.  
679 Execution wise, the complexity of the algorithm is relatively low. However, it tends to increase  
680 with growing data-points, hence the needs for further improvement in future work.



681 **Acknowledgements**

682 The third author has participated in this research work as part of the *TYPHON* Project,  
 683 which has received funding from the European Union’s Horizon 2020 Research and Innovation  
 684 Programme under grant agreement No. 780251.

685 **Appendix A: Supplementary information**

686 *Structural Communities: optimisation and interpretability.* Recall that the *network-communities*  
 687  $(\mathcal{M}_{C_{ns}}^{n \times p})$  matrix is decomposed into its approximate constituents given by Eq. 7, i.e.  $\mathcal{M}_{C_{ns}} \approx$   
 688  $\mathcal{M}_{c_{vr}} \mathcal{M}_{c_{nr}}^T$  and the optimisation function (Eq. 8) given by:  $\min_{\mathcal{M}_{c_{vr}}, \mathcal{M}_{c_{nr}}} \|\mathcal{M}_{C_{ns}} - \mathcal{M}_{c_{vr}} \mathcal{M}_{c_{nr}}^T\|_F^2$   
 689 subject to  $\mathcal{M}_{c_{vr}}, \mathcal{M}_{c_{nr}} \geq 0$ . The following conventions are used to represent the matrices:  
 690  $\mathcal{M}_{C_{ns}} \mapsto D$ ,  $\mathcal{M}_{c_{vr}} \mapsto P = [p_{is}]$ ,  $\mathcal{M}_{c_{nr}} \mapsto Q = [q_{js}]$ . We follow the *NMF* scheme [64] in the  
 691 *modelling of structural communities*.

*Iterative Update.* In response to the additional parameters ( $\alpha, \beta$  with values  $\leq 0$ ,) induced by  
 the *Lagrangian relaxation*, the objective function  $M_{sr}$  is given by the following equation:

$$M_{sr} = \|D - PQ^T\|_F^2 + \sum_{i=1}^n \sum_{s=1}^k p_{is} \alpha_{is} + \sum_{j=1}^d \sum_{s=1}^k q_{js} \beta_{js}$$

692 To solve the optimisation problem, the process begins with computing the gradient of the La-  
 693 grangian relaxation with respect to the first aspect of the *minmax* (i.e. minimisation) optimisation  
 694 variables. To achieve an optimal solution, the optimisation condition needs to be based on  $P, Q$   
 695 only. Hence, to eliminate the introduced *Lagrangian multipliers*, the *KKT optimality condition*,  
 696 which suggests that  $p_{is} \alpha_{is} = 0$  and  $q_{js} \beta_{js} = 0$ , is applied. We then solve for the optimisation  
 697 parameters as follows.

$$\begin{aligned} \|D - PQ^T\|_F^2 &= (D - PQ^T)^T (D - PQ) \\ &= (D^T - P^T Q) (D - PQ) \\ &= \underbrace{D^T D}_1 - \underbrace{D^T P Q}_2 - \underbrace{Q P^T D}_3 + \underbrace{Q P^T P Q}_4 \end{aligned} \quad (15)$$

698 In Eq. 15, the second term (2) and third term (3) are equal and the fourth term (4) can be  
 699 expressed in a quadratic form depending on the parameter of interest (for minimisation). Thus,

$$M_{sr} = D^T D - 2QP^T D + P^2 Q^T Q + \sum_{i=1}^n \sum_{s=1}^k p_{is} \alpha_{is} + \sum_{j=1}^d \sum_{s=1}^k q_{js} \beta_{js} \quad (16)$$

700 From Eq. 16, the partial differentiation with respect to  $P$  gives:

$$\frac{\partial}{\partial p_{is}} M_{sr} = -(2DQ)_{is} + (2PQ^T Q)_{is} + \alpha_{is}$$

divide by 2 and equate to zero

$$= -(DQ)_{is} + (PQ^T Q)_{is} + \alpha_{is} = 0$$

to eliminate the relaxation parameters multiply with  $p_{is}$  throughout

$$= -(DQ)_{is}p_{is} + (PQ^T Q)_{is}p_{is} + \alpha_{is}p_{is} = 0 \quad (17)$$

the term  $\alpha_{is}p_{is}$  equates to 0 according to KKT optimality, thus

$$(PQ^T Q)_{is}p_{is} = (DQ)_{is}p_{is}$$

the update rule:

$$p_{is} = \frac{(DQ)_{is}p_{is}}{(PQ^T Q)_{is}}$$

701 The last term or expression in Eq. 17 is the update rule for the parameter  $P$ . A similar process  
702 applies to the parameter  $Q$ :

$$M_{sr} = D^T D - 2QP^T D + P^T P Q^2 + \sum_{i=1}^n \sum_{s=1}^k p_{is} \alpha_{is} + \sum_{j=1}^d \sum_{s=1}^k q_{js} \beta_{js} \quad (18)$$

703 The partial derivative with respect to  $Q$  is given by the following:

$$\frac{\partial}{\partial q_{js}} M_{sr} = -(2D^T P)_{js} + (2QP^T P)_{js} + \beta_{js}$$

divide by 2 and equate to zero

$$= -(D^T P)_{js} + (QP^T P)_{js} + \beta_{js} = 0$$

to eliminate the relaxation parameters multiply with  $q_{js}$  throughout

$$= -(D^T P)_{js}q_{js} + (QP^T P)_{js}q_{js} + \beta_{js}q_{js} = 0 \quad (19)$$

the term  $\beta_{js}q_{js}$  equates to 0 according to KKT optimality, thus

$$(QP^T P)_{js}q_{js} = (D^T P)_{js}q_{js}$$

the update rule:

$$q_{js} = \frac{(D^T P)_{js}q_{js}}{(QP^T P)_{js}}$$

704 The last term or expression in Eq. 19 is the update rule for the parameter  $Q$ . The process  
705 of updating  $P, Q$  involves comparing their values to the original matrix  $D$ , and the goal is to  
706 minimise the difference or *error*. The iterative update of the parameters ( $p_{is}$  and  $q_{js}$ ) continues  
707 until convergence.

## 708 References

709 [1] A. Lancichinetti, S. Fortunato, J. Kertesz, Detecting the overlapping and hierarchical  
710 community structure in complex networks, *New journal of physics* 11 (2009) 033015.

- 711 [2] J. Scott, Social network analysis, *Sociology* 22 (1988) 109–127.
- 712 [3] D. J. Watts, S. H. Strogatz, Collective dynamics of ‘small-world’ networks, *nature* 393  
713 (1998) 440.
- 714 [4] R. Albert, A.-L. Barabási, Statistical mechanics of complex networks, *Reviews of modern*  
715 *physics* 74 (2002) 47.
- 716 [5] M. E. Newman, M. Girvan, Finding and evaluating community structure in networks,  
717 *Physical review E* 69 (2004) 026113.
- 718 [6] B. Berelson, G. A. Steiner, *Human behavior: An inventory of scientific findings.* (1964).
- 719 [7] M. E. Shaw, *Group dynamics: The psychology of small group behavior* (1971).
- 720 [8] M. Granovetter, Problems of explanation in economic sociology, *Networks and organiza-*  
721 *tions: Structure, form, and action* (1992) 25–56.
- 722 [9] D. J. Brass, K. D. Butterfield, B. C. Skaggs, Relationships and unethical behavior: A social  
723 network perspective, *Academy of management review* 23 (1998) 14–31.
- 724 [10] M. E. Newman, J. Park, Why social networks are different from other types of networks,  
725 *Physical review E* 68 (2003) 036122.
- 726 [11] R. J. Williams, N. D. Martinez, Simple rules yield complex food webs, *Nature* 404 (2000)  
727 180.
- 728 [12] G. W. Flake, S. Lawrence, C. L. Giles, F. M. Coetzee, Self-organization and identification  
729 of web communities, *Computer* (2002) 66–71.
- 730 [13] S. Papadopoulos, Y. Kompatsiaris, A. Vakali, P. Spyridonos, Community detection in social  
731 media, *Data Mining and Knowledge Discovery* 24 (2012) 515–554.
- 732 [14] C. Cao, Q. Ni, Y. Zhai, An improved collaborative filtering recommendation algorithm  
733 based on community detection in social networks, in: *Proceedings of the 2015 annual*  
734 *conference on genetic and evolutionary computation, ACM, 2015*, pp. 1–8.
- 735 [15] M. E. Newman, Properties of highly clustered networks, *Physical Review E* 68 (2003)  
736 026121.
- 737 [16] J. Yang, J. McAuley, J. Leskovec, Community detection in networks with node attributes,  
738 in: *2013 IEEE 13th International Conference on Data Mining, IEEE, 2013*, pp. 1151–1156.
- 739 [17] V. Arnaboldi, A. Guazzini, A. Passarella, Egocentric online social networks: Analysis of  
740 key features and prediction of tie strength in facebook, *Computer Communications* 36  
741 (2013) 1130–1144.
- 742 [18] N. J. Krogan, G. Cagney, H. Yu, G. Zhong, X. Guo, A. Ignatchenko, J. Li, S. Pu, N. Datta,  
743 A. P. Tikuisis, et al., Global landscape of protein complexes in the yeast *saccharomyces*  
744 *cerevisiae*, *Nature* 440 (2006) 637.

- 745 [19] M. A. Nascimento, J. Sander, J. Pound, Analysis of sigmod's co-authorship graph, *ACM*  
746 *Sigmod record* 32 (2003) 8–10.
- 747 [20] G. Palla, A.-L. Barabási, T. Vicsek, Quantifying social group evolution, *Nature* 446 (2007)  
748 664.
- 749 [21] E. O. Laumann, P. V. Marsden, D. Prensky, The boundary specification problem in network  
750 analysis, *Research methods in social network analysis* 61 (1989) 87.
- 751 [22] S. P. Borgatti, D. S. Halgin, On network theory, *Organization science* 22 (2011) 1168–1181.
- 752 [23] H. Kwak, C. Lee, H. Park, S. Moon, What is twitter, a social network or a news media?,  
753 in: *Proceedings of the 19th international conference on World wide web*, AcM, 2010, pp.  
754 591–600.
- 755 [24] C. Wilson, B. Boe, A. Sala, K. P. Puttaswamy, B. Y. Zhao, User interactions in social  
756 networks and their implications, in: *Proceedings of the 4th ACM European conference on*  
757 *Computer systems*, AcM, 2009, pp. 205–218.
- 758 [25] J. Chen, O. R. Zaiane, R. Goebel, Detecting communities in large networks by iterative  
759 local expansion, in: *2009 International Conference on Computational Aspects of Social*  
760 *Networks*, IEEE, 2009, pp. 105–112.
- 761 [26] P. Doreian, V. Batagelj, A. Ferligoj, Positional analyses of sociometric data, *Models and*  
762 *methods in social network analysis* 77 (2005) 77–96.
- 763 [27] J. Shi, J. Malik, Normalized cuts and image segmentation, *Departmental Papers (CIS)*  
764 (2000) 107.
- 765 [28] M. E. Newman, Modularity and community structure in networks, *Proceedings of the*  
766 *national academy of sciences* 103 (2006) 8577–8582.
- 767 [29] R. Balasubramanyan, W. W. Cohen, Block-lda: Jointly modeling entity-annotated text and  
768 entity-entity links, in: *Proceedings of the 2011 SIAM International Conference on Data*  
769 *Mining*, SIAM, 2011, pp. 450–461.
- 770 [30] W. Lin, X. Kong, P. S. Yu, Q. Wu, Y. Jia, C. Li, Community detection in incomplete  
771 information networks, in: *Proceedings of the 21st international conference on World Wide*  
772 *Web*, ACM, 2012, pp. 341–350.
- 773 [31] J. Leskovec, J. J. McAuley, Learning to discover social circles in ego networks, in:  
774 *Proceedings of NIPS*, 2012, pp. 539–547.
- 775 [32] J. Allan, R. Papka, V. Lavrenko, On-line new event detection and tracking., in: *Sigir*,  
776 volume 98, Citeseer, 1998, pp. 37–45.
- 777 [33] Y. Yang, A study of thresholding strategies for text categorization, in: *Proceedings of*  
778 *the 24th annual international ACM SIGIR conference on Research and development in*  
779 *information retrieval*, ACM, 2001, pp. 137–145.

- 780 [34] T. Brants, F. Chen, A. Farahat, A system for new event detection, in: Proceedings of  
781 the 26th annual international ACM SIGIR conference on Research and development in  
782 informaion retrieval, ACM, 2003, pp. 330–337.
- 783 [35] G. P. C. Fung, J. X. Yu, P. S. Yu, H. Lu, Parameter free bursty events detection in text  
784 streams, in: Proceedings of the 31st international conference on Very large data bases,  
785 VLDB Endowment, 2005, pp. 181–192.
- 786 [36] C. M. Bishop, Pattern recognition and machine learning, springer, 2006.
- 787 [37] P. Erdős, A. Rényi, On the evolution of random graphs, Publ. Math. Inst. Hung. Acad. Sci  
788 5 (1960) 17–60.
- 789 [38] D. J. Watts, P. S. Dodds, Influentials, networks, and public opinion formation, Journal of  
790 consumer research 34 (2007) 441–458.
- 791 [39] E. Katz, P. F. Lazarsfeld, E. Roper, Personal influence: The part played by people in the  
792 flow of mass communications, Routledge, 2017.
- 793 [40] H. Sundaram, Y.-R. Lin, M. De Choudhury, A. Kelliher, Understanding community  
794 dynamics in online social networks: a multidisciplinary review, IEEE Signal Processing  
795 Magazine 29 (2012) 33–40.
- 796 [41] M. E. Newman, Detecting community structure in networks, The European Physical  
797 Journal B 38 (2004) 321–330.
- 798 [42] M. E. Newman, Fast algorithm for detecting community structure in networks, Physical  
799 review E 69 (2004) 066133.
- 800 [43] D. J. Lawson, D. Falush, Population identification using genetic data, Annual review of  
801 genomics and human genetics 13 (2012) 337–361.
- 802 [44] C. D. Manning, C. D. Manning, H. Schütze, Foundations of statistical natural language  
803 processing, MIT press, 1999.
- 804 [45] C. Aggarwal, K. Subbian, Evolutionary network analysis: A survey, ACM Computing  
805 Surveys (CSUR) 47 (2014) 10.
- 806 [46] P. Pons, M. Latapy, Computing communities in large networks using random walks., J.  
807 Graph Algorithms Appl. 10 (2006) 191–218.
- 808 [47] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, E. Lefebvre, Fast unfolding of communities  
809 in large networks, Journal of statistical mechanics: theory and experiment 2008 (2008)  
810 P10008.
- 811 [48] A. Pothen, H. D. Simon, K.-P. Liou, Partitioning sparse matrices with eigenvectors of  
812 graphs, SIAM journal on matrix analysis and applications 11 (1990) 430–452.
- 813 [49] D. M. Blei, A. Y. Ng, M. I. Jordan, Latent dirichlet allocation, Journal of machine Learning  
814 research 3 (2003) 993–1022.

- 815 [50] X. Yan, J. Guo, Y. Lan, X. Cheng, A biterm topic model for short texts, in: Proceedings of  
816 the 22nd international conference on World Wide Web, ACM, 2013, pp. 1445–1456.
- 817 [51] P. Berkhin, A survey of clustering data mining techniques, in: Grouping multidimensional  
818 data, Springer, 2006, pp. 25–71.
- 819 [52] K. Chaudhuri, S. M. Kakade, K. Livescu, K. Sridharan, Multi-view clustering via canonical  
820 correlation analysis, in: Proceedings of the 26th annual international conference on machine  
821 learning, ACM, 2009, pp. 129–136.
- 822 [53] J. Liu, C. Wang, J. Gao, J. Han, Multi-view clustering via joint nonnegative matrix  
823 factorization, in: Proceedings of the 2013 SIAM International Conference on Data Mining,  
824 SIAM, 2013, pp. 252–260.
- 825 [54] S. Bickel, T. Scheffer, Multi-view clustering., in: ICDM, volume 4, 2004, pp. 19–26.
- 826 [55] G. Chao, S. Sun, J. Bi, A survey on multi-view clustering, arXiv preprint arXiv:1712.06246  
827 (2017).
- 828 [56] M. Ester, R. Ge, B. J. Gao, Z. Hu, B. Ben-Moshe, Joint cluster analysis of attribute data  
829 and relationship data: the connected k-center problem, in: Proceedings of the 2006 SIAM  
830 International Conference on Data Mining, SIAM, 2006, pp. 246–257.
- 831 [57] Y. Zhou, H. Cheng, J. X. Yu, Graph clustering based on structural/attribute similarities,  
832 Proceedings of the VLDB Endowment 2 (2009) 718–729.
- 833 [58] M. S. Granovetter, The strength of weak ties, in: Social networks, Elsevier, 1977, pp.  
834 347–367.
- 835 [59] I. Inuwa-Dutse, M. Liptrott, Y. Korkontzelos, Analysis and prediction of dyads in twitter,  
836 in: International Conference on Applications of Natural Language to Information Systems,  
837 Springer, 2019, pp. 303–311.
- 838 [60] A. Marley, M. Regenwetter, Choice, preference, and utility: Probabilistic and deterministic  
839 representations, New handbook of mathematical psychology 1 (2016) 374–453.
- 840 [61] I. Inuwa-Dutse, M. Liptrott, I. Korkontzelos, Simmelian ties on twitter: empirical analysis  
841 and prediction, in: 2019 Sixth International Conference on Social Networks Analysis,  
842 Management and Security (SNAMS), IEEE, 2019, pp. xx–xx.
- 843 [62] S. Wasserman, K. Faust, Social network analysis: Methods and applications, volume 8,  
844 Cambridge university press, 1994.
- 845 [63] J. Han, J. Pei, M. Kamber, Data mining: concepts and techniques, Elsevier, 2011.
- 846 [64] D. D. Lee, H. S. Seung, Learning the parts of objects by non-negative matrix factorization,  
847 Nature 401 (1999) 788.
- 848 [65] C. C. Aggarwal, Machine learning for text, Springer, 2018.

- 849 [66] D. P. Bertsekas, Nonlinear programming, *Journal of the Operational Research Society* 48  
850 (1997) 334–334.
- 851 [67] E. M. Airoldi, D. M. Blei, S. E. Fienberg, E. P. Xing, Mixed membership stochastic  
852 blockmodels, *Journal of machine learning research* 9 (2008) 1981–2014.
- 853 [68] P. Yali, Y. Jian, L. Shaopeng, L. Jing, A biterm-based dirichlet process topic model for  
854 short texts, in: *3rd International Conference on Computer Science and Service System*,  
855 Atlantis Press, 2014.
- 856 [69] G. Salton, C. Buckley, Term-weighting approaches in automatic text retrieval, *Information  
857 processing & management* 24 (1988) 513–523.
- 858 [70] C. C. Aggarwal, K. Subbian, Event detection in social streams, in: *Proceedings of the 2012  
859 SIAM international conference on data mining*, SIAM, 2012, pp. 624–635.
- 860 [71] L. Prokhorenkova, Using synthetic networks for parameter tuning in community detection,  
861 in: *International Workshop on Algorithms and Models for the Web-Graph*, Springer, 2019,  
862 pp. 1–15.
- 863 [72] T. Yoshida, Toward finding hidden communities based on user profile, *Journal of Intelligent  
864 Information Systems* 40 (2013) 189–209.
- 865 [73] J. Yang, J. Leskovec, Defining and evaluating network communities based on ground-truth,  
866 *Knowledge and Information Systems* 42 (2015) 181–213.
- 867 [74] I. Inuwa-Dutse, M. Liptrott, I. Korkontzelos, Detection of spam-posting accounts on twitter,  
868 *Neurocomputing* 315 (2018) 496–511.
- 869 [75] I. Inuwa-Dutse, I. Korkontzelos, A curated collection of covid-19 online datasets, *arXiv  
870 preprint arXiv:2007.09703* (2020).
- 871 [76] W. W. Zachary, An information flow model for conflict and fission in small groups, *Journal  
872 of anthropological research* 33 (1977) 452–473.
- 873 [77] D. Lusseau, K. Schneider, O. J. Boisseau, P. Haase, E. Sloaten, S. M. Dawson, The  
874 bottlenose dolphin community of doubtful sound features a large proportion of long-lasting  
875 associations, *Behavioral Ecology and Sociobiology* 54 (2003) 396–405.
- 876 [78] L. A. Adamic, N. Glance, The political blogosphere and the 2004 us election: divided  
877 they blog, in: *Proceedings of the 3rd international workshop on Link discovery*, 2005, pp.  
878 36–43.
- 879 [79] J. Leskovec, A. Krevl, SNAP Datasets: Stanford large network dataset collection, [http:  
880 //snap.stanford.edu/data](http://snap.stanford.edu/data), 2014.
- 881 [80] J. Leskovec, K. J. Lang, M. Mahoney, Empirical comparison of algorithms for network  
882 community detection, in: *Proceedings of the 19th international conference on World wide  
883 web*, ACM, 2010, pp. 631–640.

- 884 [81] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, D. Parisi, Defining and identifying  
885 communities in networks, *Proceedings of the national academy of sciences* 101 (2004)  
886 2658–2663.
- 887 [82] L. Danon, A. Diaz-Guilera, J. Duch, A. Arenas, Comparing community structure identifi-  
888 cation, *Journal of Statistical Mechanics: Theory and Experiment* 2005 (2005) P09008.
- 889 [83] A. L. Fred, A. K. Jain, Data clustering using evidence accumulation, in: *Object recognition*  
890 *supported by user interaction for service robots*, volume 4, IEEE, 2002, pp. 276–280.
- 891 [84] M. Girvan, M. E. Newman, Community structure in social and biological networks,  
892 *Proceedings of the national academy of sciences* 99 (2002) 7821–7826.
- 893 [85] X. Zhu, Z. Ghahramani, Learning from labeled and unlabeled data with label propagation,  
894 Technical Report, Citeseer, 2002.
- 895 [86] A. Lancichinetti, S. Fortunato, F. Radicchi, Benchmark graphs for testing community  
896 detection algorithms, *Physical review E* 78 (2008) 046110.
- 897 [87] P. J. Bickel, A. Chen, A nonparametric view of network models and newman–girvan  
898 and other modularities, *Proceedings of the National Academy of Sciences* 106 (2009)  
899 21068–21073.
- 900 [88] B. Karrer, M. E. Newman, Stochastic blockmodels and community structure in networks,  
901 *Physical review E* 83 (2011) 016107.
- 902 [89] L. C. Freeman, Some antecedents of social network analysis, *Connections* 19 (1996)  
903 39–42.
- 904 [90] R. I. Dunbar, The social brain hypothesis, *Evolutionary Anthropology: Issues, News, and*  
905 *Reviews: Issues, News, and Reviews* 6 (1998) 178–190.