

# Structural-temporal neighbourhoods of posts to characterize conversations in online forums

(Draft version: June 2, 2016)

Alberto Lumbreras  
Julien Velcin  
Laboratoire ERIC  
Université de Lyon, France  
alberto.lumbreras@univ-lyon2.fr  
julien.velcin@univ-lyon2.fr

Bertrand Jouve  
FRAMESPA/IMT  
Université de Toulouse, France  
jouve@univ-tlse2.fr

Marie Guégan  
Technicolor, France  
marie.guegan@technicolor.com

**Abstract**—Social network analysis provides tools to analyse the global and the local structure of graphs when these graphs represent users interactions. In the context of online conversations, usually represented by trees, there is a lack of tools to analyse the local structure of these conversations. In this paper, we introduced two novel definitions of neighbourhood adapted to conversation trees. We analyse them in two Reddit forum and, finally, we show that they can successfully detect different types of conversationalists in online forums.

## I. INTRODUCTION

The interactions between users in online forums are often modelled as complex networks. As such, they can be studied from many levels, each of which is represented by a graph where vertices and edges may represent different things. The most studied graph is a graph where vertices represent users and edges represent interactions between users (a post from one user replying to a post from another user). From the point of view of the community, some typically analysed properties are the degree distribution, the clustering coefficient, density, or the diameter of the graph. From the point of view of the individual users, analyses are typically focused either on the centrality of users, on the community structure or in blockmodeling (groups of users that tend to interact with the same other groups) [1]. Another common graph is a tree graph where vertices represent posts (or e-mails) and edges from one vertex to another indicate that the first is a reply to the second. Given the different discussion trees of a forum (a forest) the global analyses include depth distribution, branching factors, or even the time between two posts [2], [3], [4].

Studying discussion trees from the point of view of the users means, since trees are discussions, studying the user at a conversational level. Unfortunately, there is a lack of tools to perform an structural analysis of the conversation trees at a local level. In this paper, we propose to fill this gap through the concept of *structural-temporal neighbourhood*.

Our goal in this paper is two-fold: on the one hand, to illustrate how structural-temporal neighbourhoods can play the role of triads for conversation trees, in the sense that they show us the local structures or dynamics from which the bigger graph emerges. On the other hand, to show how structural-

temporal neighbourhoods can be used for the detection of different types of conversationalists in online forums or any other type of online discussion that is representable by a tree structure (e.g.: e-mails).

The remaining of the paper is as follows. We first discuss about the convenience of the classic neighbourhood definition in dynamic graphs such as discussion trees. Then we introduce our two new definitions of structural-temporal neighbourhoods. To illustrate the kind of neighbourhoods that we obtain in a real data, we analyse two Reddit forums. Finally, we apply our time-based neighbourhood definition to detect clusters of users that tend to appear in the same neighbourhoods. We close the paper with some suggestions of future research.

## II. DISCUSSION TREES

We represent a discussion thread by a tree graph  $G = (V, E)$  where  $V$  is a set of  $n$  vertices representing the posts (also called messages or comments) and  $E$  is a set of  $n - 1$  directed edges that say which posts replied to which post. Vertices have two attributes namely time and author. If for two given vertices  $v_i, v_j$  there exists an edge  $e = (v_i, v_j) \in E$  then  $v_j$  is called the *parent* of  $v_i$ , denoted as  $p(v_i) = v_j$ . The *root* of the tree is the only vertex with no parent, and corresponds to the post that starts the discussion. We say that two vertices  $v_i, v_k$  are *siblings* if and only if  $p(v_i) = p(v_k)$ . A vertex  $v_l$  is an ancestor of a vertex  $v_i$  and  $v_i$  is descendant of  $v_l$  if and only if  $v_l$  is in the path from  $v_i$  to the root. A *leaf* is a post with no replies. A *branch* is the path between a leaf and the root. Two vertices  $v_i$  and  $v_m$  are said to be *neighbours* if either  $p(v_i) = v_m$  or  $p(v_m) = v_i$ . The neighbours of a vertex  $i$  are its parent and its children.

## III. STRUCTURAL NEIGHBOURHOODS

Extending the classic definition of neighbourhood, we recall the definition of *structural neighbourhood*<sup>1</sup>:

**Definition 1:** Given a tree graph  $G$ , the *structure-based neighbourhood* of radius  $r$  of post  $i$ , denoted as  $\mathcal{N}_i(r)$ , is

<sup>1</sup>We coined the term *structural* to differentiate it from the other neighbourhoods discussed in this paper.

the induced graph composed by all the vertices that are at distance equal or less than  $r$  from post  $i$ .

In the context of discussion threads, this definition has two limitations. First, the decision on whether to include some post in the neighbourhood is only based on the structural distance, and therefore two posts that are at distance  $d \leq r$  are considered neighbours of  $i$  regardless of the time when they were written although, in conversations, time plays an important role. Another consequence of looking only at the structure is that the number of possible neighbourhoods within a radius  $r$  is infinite. This poses a problem when trying to categorize conversations since many conversations, while structurally different, can be considered semantically similar (e.g.: we might not want to put in different categories two structures representing, respectively, a post with 50 replies and another with 40).

#### IV. STRUCTURAL-TEMPORAL NEIGHBOURHOODS

We propose two refinements of the structure-based neighbourhood to take into account the order and time at which posts are written. In the order-based definition, a new radius is set to include at most the  $n$  closest posts in time. In the time-based definition, the new radius is set to include all posts (in the structural-neighbourhood) until the first point where the conversation slows down.

##### A. Order-based

*Definition 2:* Given an ordered tree graph  $G$ , the *order-based neighbourhood* of radius  $r$  and order  $n$  of vertex  $i$ , denoted as  $\mathcal{N}_i^O(r, n)$  is the induced subgraph from its structural neighbourhood  $\mathcal{N}_i(r)$  composed by the  $n$  vertices that are closest to  $i$  in time and for which there exists a path to  $i$  in  $\mathcal{N}_i(r, n)$ .

An example of order-based neighbourhood is given in Figure 1. This definition has two advantages over the *structural neighbourhood*. First, the temporal aspect of the conversation is better taken into account since the neighbourhood only includes posts that are not only near to  $i$  in the structure but also in time. Second, the size of the neighbourhood has an upper bound of  $\min(|\mathcal{N}_i(r)|, n)$  thus making the space of possible neighbourhood structures finite. The main limitation of this definition is that we have no *a priori* criteria to choose the proper parameter  $n$  other than making it small to capture only the local dynamic of the conversation around the post  $i$ .

##### B. Time-based

In order to take time explicitly into account, one might set fixed time-based boundaries for the neighbourhood and include only those posts whose timestamp  $t_j$  is at distance less than  $\tau$  from the ego post  $i$ ,  $|t_j - t_i| < \tau$ . However, the pace at which posts are added to the conversation may be very different between conversation threads (and also within a thread) and we have no *a priori* criteria for a proper choice of  $\tau$ .

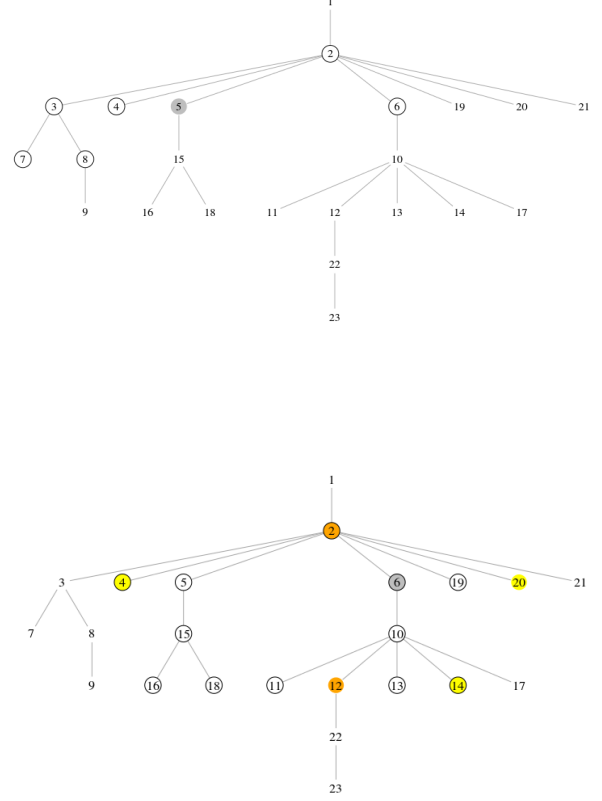


Fig. 1. Illustration of order-based (top) and time-based (bottom) neighbourhoods. Grey nodes represent the ego post. Nodes with circles are those in the neighbourhood. Numbers indicate the order of the post. The order-based neighbourhood has parameters  $r = 3$ ,  $n = 6$ . The time-based neighbourhood has parameter  $r = 3$ . In the time-based neighbourhood, horizontal changepoints (yellow) and vertical changepoints (orange) represent posts that are temporally far from their predecessors (siblings or parents) and therefore set the limits of the neighbourhood.

Rather than looking for a fixed time radius, we may instead decide it by looking at changes in the pace how new posts are added to the thread. In statistical analysis, a *changepoint* in a sequence  $x_1, \dots, x_n$  is a point that comes from a different probability distribution than its precedent values. If sequences are timestamps, they are monotonic increasing, therefore the changepoints will correspond to sudden pauses in the conversation. In the following, we will differentiate *horizontal changepoints*, that arise between siblings, from *vertical changepoints*, that arise within a branch. We say that a sequence posts with timestamps  $t_1, \dots, t_n$  belong to the same (vertical or horizontal) *local dynamic* if there is no changepoint  $t_i$  in the sequence such that  $1 < i \leq n$ . For the detection of changepoints we use the PELT algorithm [5]<sup>2</sup>. Now we can

<sup>2</sup>An implementation of this algorithm written by the authors themselves is available in the R library *changepoint*. We use their `cpt.meanvar` function.

introduce our second definition of neighbourhood.

*Definition 3:* Given a tree graph  $G$ , the *time-based neighbourhood* of radius  $r$  of vertex  $i$ , denoted as  $\mathcal{N}_i^T(r)$ , is the maximal subgraph of the structural neighbourhood  $\mathcal{N}_i(r)$  where all the vertices belong to the same vertical and horizontal local dynamic as  $i$ .

Note that we still have a radius parameter  $r$  to guarantee that the resulting structure remains local even if no changepoint is detected near the post  $i$ .

Algorithm 1 extracts time-based neighbourhoods of a given post according to the above definition. Since the changepoints only depend on the tree and not on the particular post we analyse, we previously detect the horizontal and vertical changepoints in the tree. In the case of multiple branches with some common posts, we consider that a common post is a vertical changepoint if it is a vertical changepoint in any of the branches. Once we have the changepoints, we can proceed with the algorithm. First, we extract the structural-neighbourhood. The time-based neighbourhood will be a subset of the later. Then, we look for horizontal and vertical changepoints, which mark the frontiers of the time-based neighbourhood. There are four possible cases:

- *A vertical changepoint in the ancestors:* If the changepoint is in the path between the post and the root, then the changepoint started the new local dynamic to which the ego post belongs. Thus, we remove the ancestors of the changepoint, but not the changepoint itself.
- *A vertical changepoint in the descendants:* If the changepoint is a descendant then it started a new different dynamic and therefore we must remove the descendants of the changepoint and the changepoint itself.
- *A horizontal changepoint either in the older siblings or in the ancestors:* If the changepoint is among the older siblings, then it started the horizontal dynamic to which the ego post belongs. Similarly, if the changepoint is among the ancestors, then every older sibling of the changepoint belongs to a previous local dynamic. In both cases, we remove the older siblings of the changepoint but not the changepoint itself.
- *A horizontal changepoint elsewhere:* In any other case, the horizontal changepoint starts a different local dynamic and therefore we remove its younger siblings and the changepoint.

Any time a vertex is removed we also remove its descendants.  $\mathcal{N}_i^T(r)$  is a connected induced subgraph of  $\mathcal{N}_i(r)$ . If there are no changepoints in the structural neighbourhood, then  $\mathcal{N}_i^T(r) = \mathcal{N}_i(r)$ .

An example of time-based neighbourhood is given in Figure 1 (right).

---

#### Algorithm 1 Extraction of time-based neighbourhood

---

**Input:** Posts tree  $G$ , vertical changepoints, horizontal changepoints, ego post  $i$ , radius  $r$   
**Output:**  $\mathcal{N}_i^T(r)$ : Time-based neighbourhood of  $i$  at radius  $r$   
 Compute structural neighbourhood  $\mathcal{N}_i(r)$   
 $\text{ancestors} \leftarrow \text{ancestors}(i)$  in  $\mathcal{N}_i(r)$   
 $\text{older\_siblings} \leftarrow \text{older\_siblings}(i)$  in  $\mathcal{N}_i(r)$   
 $\text{dump} \leftarrow \emptyset$   
**for**  $\text{bp} \in \text{vertical changepoints}$  **do**  
   **if**  $\text{bp} \in \text{ancestors}$  **then**  
      $\text{dump} \leftarrow \text{dump} \cup \text{ancestors}(\text{bp})$   
   **else**  
      $\text{dump} \leftarrow \text{dump} \cup \text{descendants}(\text{bp}) \cup \text{bp}$   
   **end if**  
**end for**  
**for**  $\text{bp} \in \text{horizontal changepoints}$  **do**  
   **if**  $\text{bp} \in (\text{older\_siblings} \cup \text{ancestors})$  **then**  
      $\text{dump} \leftarrow \text{dump} \cup \text{older\_siblings}(\text{bp})$   
   **else**  
      $\text{dump} \leftarrow \text{dump} \cup \text{younger\_siblings}(\text{bp}) \cup \text{bp}$   
   **end if**  
**end for**  
 $\mathcal{N}_i^T(r) \leftarrow \text{delete}(\text{dump\_posts} \cup \text{descendants}(\text{dump\_posts}))$   
 from  $\mathcal{N}_i(r)$

---

## V. NEIGHBOURHOOD COLOURING AND PRUNING

### A. Colouring

Even if the structure contains some important information about the type of conversation, there is still some ambiguity left, and two similar neighbourhoods can represent very different types of conversation. We can easily reduce this ambiguity by assigning colours to vertices, which allows us to identify some relevant property of the post.

In particular, we assign the following colours:

- *Red:* ego post.
- *Orange:* other posts written by the author of ego. It allows to identify re-entries (when the same author participates several times in the discussion [6])
- *Yellow:* parent of ego post and other posts written by the same author. It allows to identify, for instance, debates between the ego author (red and orange) and this other author (yellow). We have observed this phenomena in our data, and it is often the cause of long chains.
- *White:* root post. Differentiating the root post from the rest has been proven by previous research on online discussions. For instance, some types of users seem to get more replies than others when they initiate a thread ([7], [8]). Also, in terms of preferential attachment, root posts usually get more replies than the non-roots ([9], [10]).
- *Black:* none of the above.

Since some posts might correspond to several colours, we perform a sequential assignment of colours given by black, orange, red, yellow, white.

We do not claim this choice of labels to be of universal. Indeed, other labellings might also give interesting results since they would look at conversations from new points of view (e.g.: a colour for leaf vertices, colours according to the post length, or colours to represent the sentiment of the post).

### B. Pruning

At this point, we can still find structures the only difference of which is that one of them has more leafs hanging from some of the nodes. Thus, we prune every neighbourhood by leaving a maximum of two consecutive siblings of the same colour and where neither of them has any children. The reason of setting the limit to two is that it is the minimum necessary to distinguish between *zero*, *one* and *more than one* consecutive occurrences. In other words, we consider that the difference between one and zero replies to a post is relevant, but that five and six replies do not make any difference worth being represented. Moreover, this choice encourage smaller structures, which allows the computation of isomorphisms (Section VI-A) even in real time.

## VI. APPLICATION TO REDDIT FORUMS

Our data consists of two different forums of Reddit, namely the Podemos dataset<sup>3</sup> and the Game of Thrones dataset<sup>4</sup>. The Podemos forum was conceived in March 2014 as a tool for internal democracy, and forum members used it to debate ideological and organizational principles that were later formalized in their first party congress hold in Madrid on October 18th and 19th, 2014. Nowadays, its members use it mainly to share and discuss about political news. The Game of Thrones forum is a casual discussion forum about the Game of Thrones TV series. The Podemos dataset contains 75,000 posts corresponding to 4,262 threads written by 9,160 users between April 25th and October 20th, 2014. The Game of Thrones dataset contains 75,000 posts corresponding to 5,862 threads written by 18,907 users.

### A. Neighbourhood extraction

As a previous step for all further analysis, we extracted the structure-based ( $r = 2$ ), the order-based ( $r = 2, n = 4$ ) and the time-based ( $r = 2$ ) neighbourhood around every post for the two forums.

For each method (structure-based, order-based or time-based) we extract the neighbourhood around every post of a forum as follows. A new label is assigned to every neighbourhood structure the first time it is detected. For every post, we extract its neighbourhood, colour and prune it (Sections III and IV). Then we check whether it is isomorphic to some of the  $n$  neighbourhoods that we have already seen. If this is the case, then the current neighbourhood is given the same label than the neighbourhood that we already saw. Otherwise we create a new label for the current neighbourhood. The set of neighbourhoods found and their labels constitutes a *dictionary*.

Once the six extractions are finished, we merge the dictionaries into a global one and re-label the neighbourhoods by frequency (1 for the most frequent and so forth). Figure 5 shows some frequent neighbourhoods that will be discussed later.

In the worst case, where every post has a unique neighbourhood, every extracted neighbourhood should be compared against every neighbourhood that already appears in the dictionary before concluding that it needs a new entry. In this case, the computational cost of this operation is  $\mathcal{O}(n^2)$  where  $n$  is the number of posts. In practice, however, the bigger the number of posts processed, the less likely we go through all the dictionary without finding an isomorphism. Moreover, if we keep the dictionary sorted by frequency, a large number of neighbourhoods will find an isomorphism among the first neighbours of the dictionary.

### B. Comparing neighbourhood methods

In this section we compare different aspects of the three methods, namely (a) the size and frequency distribution of the neighbours for each method, (b) the discrepancies between methods when they extract the neighbourhood of a same post.

1) *Size and frequency distribution*: We detected in the whole dataset 3,791 different structure-based neighbourhoods, 174 different order-based and 1,488 different time-based, which means that the structure-based captures a much wider set of structures. Of course, this is not necessarily an advantage since these extra structures may be just spurious and only occur a very small number of times. Indeed, Figure 2 shows that around half of the time-based neighbourhoods (41%) and the structure-based neighbourhoods (60%) occur only once, while the number of neighbourhoods that appear more than 100 times (among 150,000 posts) is similar for the three types of neighbourhood. To cover 95% of the posts we need 623 structure-based neighbourhoods, 40 order-based and 155 time-based.

The main reason why the frequency of the order-based neighbourhoods decays faster is that the most frequent neighbourhood (id 1) is much more frequent (upper-left red point in the left plot of figure 2) than the most frequent time-based neighbourhood. This is the main handicap of the order-based neighbourhood since, as we will see later, it avoids finding richer structures. Table VI-B1 shows the most frequent neighbourhoods for the three methods in each forum. The results are quite sensitive to the choice of the neighbourhood. Note that the order-based detects less differences among the most frequent neighbourhoods in both forums. Besides, it seems that when the time-based detects the neighbourhood 4, the structure-based detects the neighbourhood 7. This illustrates the effect of the changepoint: the presence of an horizontal changepoint in 7 (many posts replying to the root) converts it, from the point of view of the time-based neighbourhood, to a 4, which indicates that the ego has at most two siblings in the same period of time.

2) *Discrepancies*: The divergent number of non-isomorphic neighbourhoods between structure-based, order-based and

<sup>3</sup><https://www.reddit.com/r/podemos>

<sup>4</sup><https://www.reddit.com/r/gameofthrones>

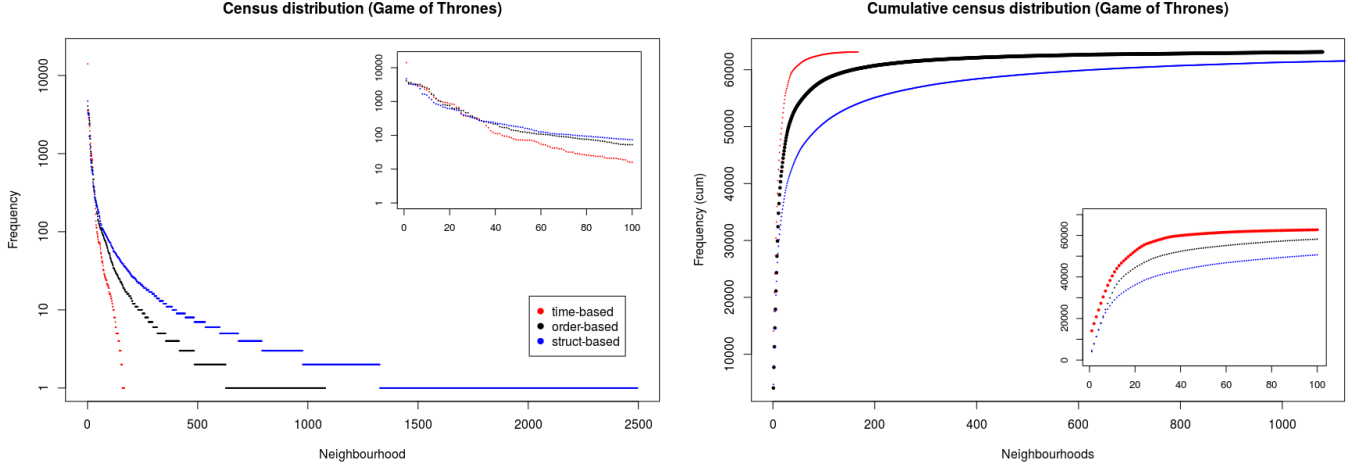


Fig. 2. Comparison of census distributions with order-based and time-based neighbourhoods.

struct-based		order-based		time-based	
POD	GOT	POD	GOT	POD	GOT
7 (9.6)	7 (7.4)	1 (25.2)	1 (25.2)	4 (9.0)	4 (6.4)
3 (16.5)	9 (12.6)	3 (32.1)	3 (32.1)	8 (17.2)	8 (12.2)
2 (23.1)	3 (17.8)	2 (38.8)	2 (38.8)	2 (25.3)	2 (17.9)
20 (27.8)	11 (22.9)	13 (45.3)	13 (45.3)	3 (31.2)	9 (23.2)
4 (31.1)	5 (27.8)	4 (50.1)	4 (50.1)	10 (36.6)	14 (28.4)
13 (34.1)	2 (32.3)	18 (54.7)	18 (54.7)	14 (41.8)	11 (33.5)
6 (36.9)	12 (36.1)	6 (58.9)	6 (58.9)	6 (45.8)	10 (38.5)
5 (39.4)	27 (38.8)	12 (65.5)	10 (62.8)	15 (48.9)	5 (43.1)
15 (41.9)	4 (41.4)	16 (68.5)	12 (65.5)	19 (51.5)	6 (47.4)
12 (43.7)	6 (43.9)	5 (70.6)	16 (68.1)	5 (53.9)	3 (51.3)
21 (45.3)	21 (45.9)	17 (72.8)	5 (70.7)	7 (55.9)	25 (55.1)
30 (46.9)	15 (47.7)	23 (75.5)	17 (72.8)	16 (57.6)	15 (57.8)
16 (48.4)	13 (49.2)	31 (76.1)	23 (74.5)	36 (59.4)	12 (60.2)
39 (49.8)	20 (50.6)	19 (77.5)	31 (76.1)	25 (61.2)	7 (62.2)
44 (51.2)	30 (51.9)	22 (79.0)	19 (77.5)	21 (62.6)	21 (63.9)
8 (52.4)	22 (53.1)	24 (80.3)	22 (78.9)	34 (63.9)	19 (65.5)
22 (53.3)	43 (54.2)	26 (81.5)	24 (80.3)	38 (65.1)	17 (66.7)
51 (54.2)	8 (55.2)	29 (82.7)	26 (81.5)	23 (66.4)	36 (68.0)

TABLE I  
MOST FREQUENT NEIGHBOURHOODS (AND CUMMULATIVE FREQUENCY)  
FOR THE THREE METHODS IF PODEMOS (POD) AND GAME OF THRONES  
(GOT) FORUMS.

times-based methods suggests that some neighbourhoods, for instance, time-based, might be systematically mapped to the same neighbourhoods in order-based. Or similarly, that some structure-based neighbourhoods are systematically mapped into the same time-based neighbourhoods. Figure 3 gives a general overview of these mappings. The left figure shows, for every post, its time-based neighbourhood and its order-based neighbourhood, while the right figure shows the time-based and the structure-based. The points in the diagonal correspond to posts that have been assigned the same neighbourhood by the two compared methods.

The diagonal of the time-based / order-based plot contains 36% of the posts. The diagonal of the structure-based / time-based plot contains 40% of the posts. For the order-based / struct-based (not shown) there are 41% of the posts in the

diagonal. In general, the structures with a bigger agreement are among the most frequent. This is probably due to the fact that a high amount of threads are small and therefore there is less room for discrepancy. There is a 30.75% of the posts that are assigned the same neighbourhood structure by the three methods.

Rather than looking at where the methods agree, it is even more interesting to see where they disagree, since this will reveal whether one of the methods is incapable of detecting some interesting conversational structure. In Figure 3 some horizontal lines can be spotted specially in the lower part of the plots. These correspond to posts with different time-based neighbourhoods that are assigned the same order-based neighbourhood (left), and different time-based neighbourhoods that are assigned the same order-based neighbourhood (right). Looking for the order-based neighbourhoods that have a correspondence with a larger set of time-based neighbourhoods, we detect that these correspond to neighbourhoods 1 (77 different time-based) and 50 (57 different time-based), 17 and 18. Looking for the time-based neighbourhoods that have a correspondence with a larger set of structure-based neighbourhoods, we detect that these are neighbourhoods 10, 4, 8, 19, 14.

In retrospective, these results are not surprising. What we see is simply that structures that are often captured by one method collapse into another single structure when another method with a smaller space of structures is used. Yet, the three methods mostly agree when they find a frequent neighbourhood. But this might be due to the fact that simpler motifs tend to be in shorter threads, where there is no room for two neighbourhood methods to detect a different structure.

It is clear from this analysis that the order-based does not do a good job on capturing all the richness of the conversational structures due to the dominance of neighbourhood 1. There are two other reasons to choose the time-neighbourhood over the structure based. First, the time-based has less spurious

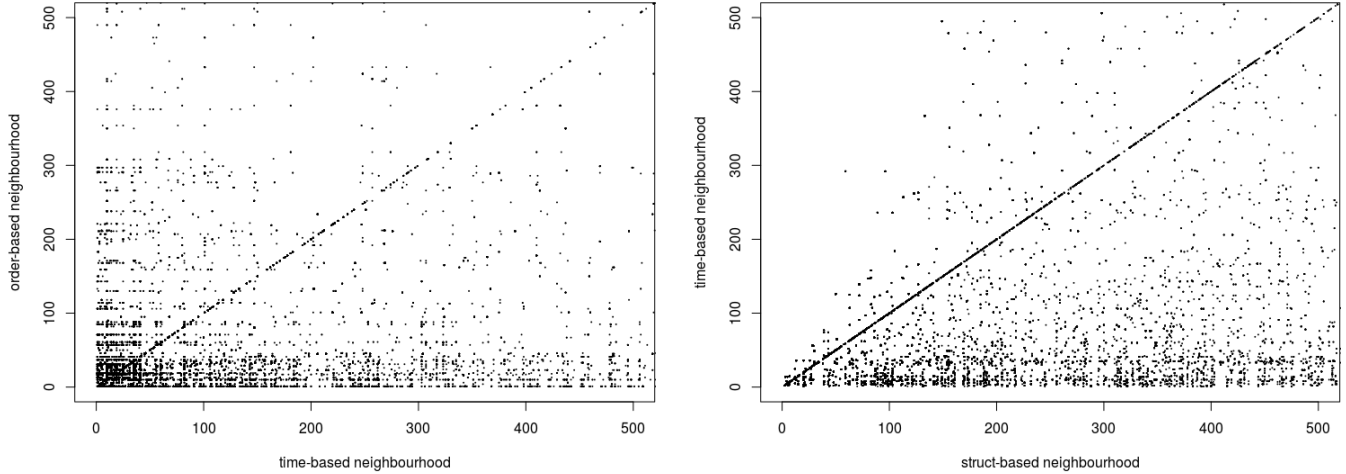


Fig. 3. Neighbourhood of posts from the point of view of order-based and time-based. Neighbourhoods with the same label in the order-based and the time-based axes are isomorphic. The points in the diagonal correspond to posts to whom the two methods assign exactly the same neighbourhood. Horizontal rows of points correspond to posts that are assigned many different time-based neighbourhoods but the same single order-based neighbourhood.

structures, covering 95% of posts with 155 neighbourhoods, against the 623 neighbours necessary to cover that amount with the structure-based. Second, the time-based neighbourhood has more grounded criteria to decide where the neighbourhood ends, while the structure-based has no criteria at all.

### C. Conversation-based clustering of users

In this section, we apply the time-based neighbourhoods to cluster users in the Podemos dataset based on the neighbourhoods of their posts. Our intuition is that some users like participating in some kind of discussion rather than other. Certainly, most part of the information necessary to understand the nature of a discussion in its textual content. However, due to the huge diversity of topics, vocabulary, and the difficulty of current algorithms to capture the language subtleties such as humour, irony, or context, we turn our attention towards the structure of the discussions, which can also contain some information. We work under the hypothesis that the structural-neighbourhoods in which a user post in embedded reflect the kind of conversation in that part of the thread.

Our dataset contains a set of 100 *active users* who wrote more than 100 posts. For those users, we create an initial feature matrix  $U \times N$  where  $U$  is the number of users and  $N$  is the number of neighbourhoods in the dictionary, and where the position  $(u, n)$  is a counter of the number of times that a post written by user  $u$  has a neighbourhood isomorphic to  $n$ . We drop those feature columns that are zero for every user (these are neighbourhoods seen only around posts of users with low-level activity). To make the feature vector of a user independent on the number of posts, we transform the counts into percentages. And since some features have much higher percentages than others for most users, we scale and

normalize the matrix so that every feature has mean 0 and variance 1. To avoid non-significant scores, we remove also the feature column corresponding to neighbourhoods that have a frequency less than 50 among the active users.

We use k-means to find the clusters, though one can use any other clustering method. Since the plot over the Within-Cluster Sum of Squares for  $k = 1, \dots, 20$  did not show any clear elbow we chose  $k = 3$  clusters so that clusters are easier to interpret. We did the clustering over a feature matrix with order-based neighbourhoods and another feature matrix with time-based neighbourhoods.

Since the set of feature dimension is relatively large, we show in Figure 4 how the users in a cluster are distributed for each feature. The features shown are the ones that have an bigger absolute value (in mean) for each cluster (above the percentile 0.95).

- Green cluster: these users are prominent in conversations 8, 4, 15, which correspond to replies to the root in threads where the root post attracted some attention.
- Red cluster: these users tend to appear in 20, 3, 19, 14, 13, and 18, which are almost all cascades, sometimes between only two users (20, 18).
- Blue cluster: these users seems to avoid cascades (3, 19) while they do not limit themselves to reply to the root, but instead the participate deeper in the discussion (6).

## VII. CONCLUSIONS

The goal of is paper was to shed some light on what would be a proper method to characterise discussions online forums. For that, we analysed three different kinds of neighbourhood applicable to tree structures by which conversation threads are represented. We named this three types of neighbourhood structure-based, order-based and time-based. While the first

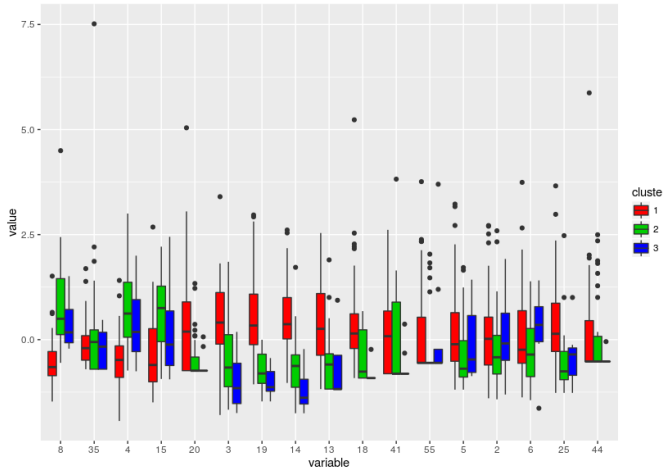


Fig. 4. Whiskers plots in the most relevant features (time-based)

one is a mere extension of the classical concept of neighbourhood in a graph, the aim of the other two is to explicitly consider time when deciding which are the neighbours of a given post.

We showed the limitations of the order-based to capture the structure of conversations. We think this finding is relevant because, in some contexts, order is more relevant than real time, and assuming that posts (or whatever vertices represent) arrive in homogeneous intervals gives excellent results [10].

In the context of online discussions, time is critical and must be considered in the definition of the neighbourhood motifs. Our choice has been to use *change-points*, variations in the pace of the discussion, to decide where a neighbourhood ends. Besides explicitly considering time, this method gives in practice relatively small neighbourhoods, which guarantees its scalability since isomorphisms are considerably fast to compute in small graphs.

We also showed that, thanks to the use of change-points to set limits of the neighbourhoods, the space of neighbourhoods is considerably reduced with respect to the structural neighbourhood.

Finally, we used this neighbourhoods to characterise users in terms of the structure of the conversations they participate in and showed that, indeed, there are different types of structural conversationalists.

The concept of structural-temporal neighbourhood opens the door to some interesting paths of research. The space of neighbourhoods is still quite large and, besides using pruning strategies more aggressive, it would be interesting to analyse which neighbourhoods, despite of being non-isomorphic, may represent the same type of conversation. Such analyse would allow to merge some of the neighbourhoods resulting in an even more reduced and meaningful space. Also, other kind of neighbourhoods can also be studied to analyse different phenomena. For instance, we might take only those neighbours which are descendants of the ego post. Such kind of neighbourhood may be used to classify users in terms in the

reaction they trigger in the discussion. The inverse technique, that is, take all the neighbours that are not descendants of the ego, may help to understand the kind of discussion that a users tends to be attracted to. This might be useful for posts recommendations or even to improve current generative models that try to reproduce the way how a thread grows [10], [11].

## REFERENCES

- [1] A. McCallum, X. Wang, and A. Corrada-Emmanuel, "Topic and role discovery in social networks with experiments on enron and academic email," *Journal of Artificial Intelligence Research*, vol. 30, no. 1, pp. 249–272, 2007. [Online]. Available: <http://www.aaai.org/Papers/JAIR/Vol30/JAIR-3007.pdf>
- [2] R. Bhatt and K. Barman, "Global Dynamics of Online Group Conversations," in *Proceedings of the Sixth International AAAI Conference on Weblogs and Social Media*, 2012.
- [3] N. Gaumont, T. Viard, R. Fournier-S'niehotta, Q. Wang, and M. Latapy, "Analysis of the temporal and structural features of threads in a mailing-list," in *Workshop on Complex Networks CompleNet*, 2016. [Online]. Available: <http://arxiv.org/pdf/1512.05002v1.pdf>
- [4] R. Dorat, M. Latapy, B. Conein, and N. Auray, "Multi-level analysis of an interaction network between individuals in a mailing-list," *Annales des télécommunications*, vol. 62, no. 3-4, pp. 325–349, 2007. [Online]. Available: <http://link.springer.com/article/10.1007/BF03253264>
- [5] R. Killick, P. Fearnhead, and I. Eckely, "Optimal detection of change-points with a linear computational cost," *Journal of the American Statistical Association*, pp. 1590–1598, 2012.
- [6] L. Backstrom, J. Kleinberg, L. Lee, and C. Danescu-Niculescu-Mizil, "Characterizing and Curating Conversation Threads: Expansion, Focus, Volume, Re-entry," in *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*, 2013, pp. 13–22.
- [7] I. Himmelboim, E. Gleave, and M. Smith, "Discussion catalysts in online political discussions: Content importers and conversation starters," *Journal of Computer-Mediated Communication*, vol. 14, no. 4, pp. 771–789, jul 2009.
- [8] A. Lumbreras, J. Lanagan, J. Velcin, and B. Jouve, "Analyse des rôles dans les communautés virtuelles : définitions et premières expérimentations sur IMDb," in *Modèles et Analyses Réseau : Approches Mathématiques et Informatiques (MARAMI)*, 2013, pp. 1–12.
- [9] V. Gómez, H. J. Kappen, and A. Kaltenbrunner, "Modeling the structure and evolution of discussion cascades," in *Proceedings of the 22nd ACM conference on Hypertext and hypermedia*, 2010, pp. 181–190.
- [10] V. Gómez, H. J. Kappen, N. Litvak, and A. Kaltenbrunner, "A likelihood-based framework for the analysis of discussion threads," *World Wide Web*, p. 31, apr 2012.
- [11] R. Kumar, M. Mahdian, and M. McGlohon, "Dynamics of Conversations," in *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2010, pp. 553–562.

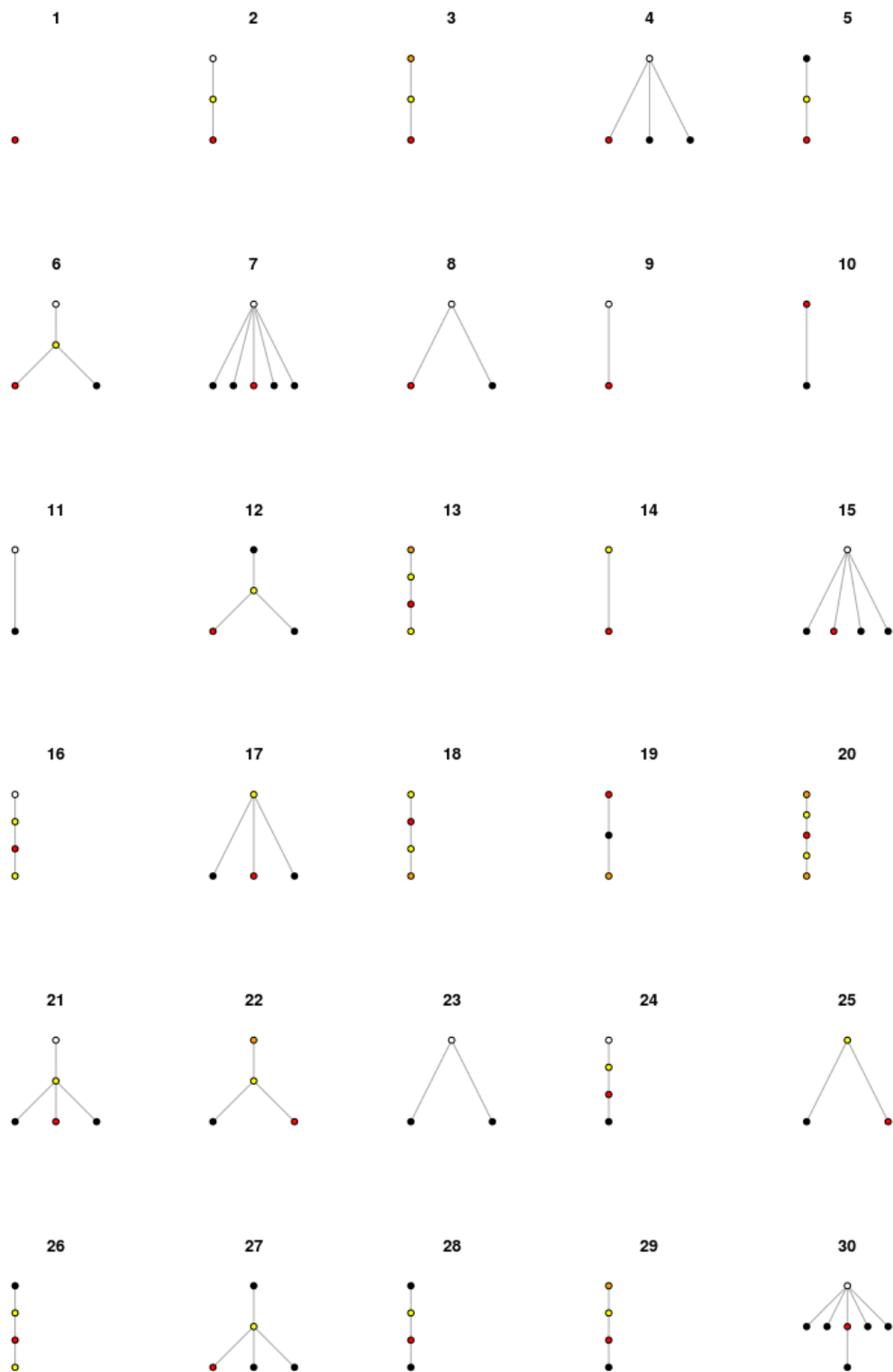


Fig. 5. Neighbourhoods reference list