Identifying Emergent Leadership in OSS Projects Based on Communication Styles

YUEKAI HUANG, Laboratory for Internet Software Technologies, Institute of Software Chinese Academy of Sciences; University of Chinese Academy of Sciences, Beijing, China

YE YANG*, School of Systems and Enterprises, Stevens Institute of Technology, Hoboken, NJ, USA

JUNJIE WAN G^* , Laboratory for Internet Software Technologies, State Key Laboratory of Computer Sciences, Institute of Software Chinese Academy of Sciences; University of Chinese Academy of Sciences, Beijing, China

WEI ZHENG, School of Business, Stevens Institute of Technology, Hoboken, NJ, USA

QING WANG*, Laboratory for Internet Software Technologies, State Key Laboratory of Computer Sciences, Institute of Software Chinese Academy of Sciences; University of Chinese Academy of Sciences, Beijing, China

In open source software (OSS) communities, existing leadership indicators are dominantly measured by code contribution or community influence. Recent studies on emergent leadership shed light on additional dimensions such as intellectual stimulation in collaborative communications. To that end, this paper proposes an automated approach, named iLead, to mine communication styles and identify emergent leadership behaviors in OSS communities, using issue comments data. We start with the construction of 6 categories of leadership behaviors based on existing leadership studies. Then, we manually label leadership behaviors in 10,000 issue comments from 10 OSS projects, and extract 304 heuristic linguistic patterns which represent different types of emergent leadership behaviors in flexible and concise manners. Next, an automated algorithm is developed to merge and consolidate different pattern sets extracted from multiple projects into a final pattern ranking list, which can be applied for the automatic leadership identification. The evaluation results show that iLead can achieve a median precision of 0.82 and recall of 0.78, outperforming ten machine/deep learning baselines. To demonstrate practical usefulness, we also conduct empirical analysis and human evaluation of the identified leadership behaviors from iLead. We argue that emergent leadership behaviors in issue discussion should be taken into consideration to broaden existing OSS leadership viewpoints. Practical insights on community building and leadership skill development are offered for OSS community and individual developers, respectively.

CCS Concepts: • Software and its engineering → Open source model.

Additional Key Words and Phrases: Leadership, Communication style, Linguistic pattern, Open source software

Authors' addresses: Yuekai Huang, huangyuekai18@mails.ucas.ac.cn, Laboratory for Internet Software Technologies, Institute of Software Chinese Academy of Sciences; University of Chinese Academy of Sciences, Beijing, China; Ye Yang, yyang4@stevens.edu, School of Systems and Enterprises, Stevens Institute of Technology, Hoboken, NJ, USA; Junjie Wang, junjie@iscas.ac.cn, Laboratory for Internet Software Technologies, State Key Laboratory of Computer Sciences, Institute of Software Chinese Academy of Sciences; University of Chinese Academy of Sciences, Beijing, China; Wei Zheng, wzheng11@stevens.edu, School of Business, Stevens Institute of Technology, Hoboken, NJ, USA; Qing Wang, wq@iscas.ac.cn, Laboratory for Internet Software Technologies, State Key Laboratory of Computer Sciences, Institute of Software Chinese Academy of Sciences; University of Chinese Academy of Sciences, Beijing, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

©~2021 Association for Computing Machinery.

Manuscript submitted to ACM

^{*}Corresponding author

ACM Reference Format:

1 INTRODUCTION

Steady influx of developer resource is key to healthy and sustainable OSS communities [4, 27]. The presence of community leaders, who both steer the development direction and motivate more developers to contribute is crucial to ensure a successful outcome of an OSS project [39]. In addition, it is also reported that developer retention is highly associated with the attention and treatment they receive in the project [31, 59, 70]. Therefore, it is an essential function for an OSS community to effectively monitor and appropriately recognize developer contribution in timely and fair manners.

In light of the importance of having leadership figures, many studies have investigated various developer leadership factors across OSS communities. Most rely on metrics representing individuals' code contribution or community influence. Example metrics of code contribution include code size (e.g., #lines of code added/removed [44], #commits [5, 16]), code structure (e.g., code instability [5]), code quality (e.g. #problems reported [44], #buggy commits [16]), or code topics [41]. Example metrics of community influence include the number of followers (e.g., #followers [39, 45]) and influence on other developer's motivations [53].

Most of today's OSS communities operate beneath an umbrella organization, while others are organized entirely independently, and yet others follow a strategy somewhere in between [21]. Different OSS organizational forms correspond with different leadership structure and coordination strategies, which complicates the measurement of leadership effectiveness. As detailed in Section 2, we find that there is mismatch between these developers' contribution in issue discussion and their contributions measured by other dimensions such as code contribution or community followers. We conjecture that in OSS communities where development is highly autonomous, and explicit coordination and control are loose or non-existent, those who know how to communicate effectively with others will contribute more in leading, coordinating, and influencing the development. This indicates the need of a new lens in identifying and recognizing OSS leadership behaviors corresponding to more effective and efficient issue resolution.

To bridge the gap, we start with constructing six categories of emergent leadership behaviors applicable to OSS projects from existing leadership models [2, 14, 39, 66, 67, 71]. These include proposal, redirection, confirmation, inquiry, operation, and volunteer (as show in Table 1). And then we propose an automated approach, named iLead, to mine communication styles and identify such emergent leadership behaviors in OSS communities, using issue comments data. The construction of iLead consists of three steps. First, we carry out the data labelling process to identify the category of leadership per issue comment. Second, we adapt the definition of linguistic patterns in previous studies [17, 40, 58] and manually extract the heuristic linguistic patterns to express different types of communication styles in flexible and concise manners. Third, based on the two previous steps, we develop an automatic algorithm to merge and consolidate different pattern sets extracted from multiple projects into a final pattern ranking list. When applying iLead for leadership identification, the final pattern ranking list can be used to automatically match new issue comments and identify corresponding leadership behaviors. The evaluation is conducted on 10 popular OSS projects with 10,000 comments. Results show that iLead can achieve a median precision of 0.82 and recall of 0.78, outperforming ten machine/deep learning baselines.

To demonstrate practical usefulness, we also conduct empirical analysis and human evaluation of the identified leadership behaviors with iLead. Results motivate the need for considering the emergent leadership behaviors in issue Manuscript submitted to ACM

discussion to broaden existing OSS leadership viewpoints. Practical insights on community building and leadership skill development for OSS community developers are also discussed. To summarize, the main contributions of this paper are as follows:

- The construction of 6 categories of emergent leadership behaviors in OSS projects from issue comments data, adapting existing leadership models into OSS context. This is the first taxonomy of the emergent leadership behaviors in issue discussions of OSS communities, which can motivate following-up studies in this direction.
- The development of iLead (with 304 defined linguistic patterns and pattern consolidation algorithm) to automatically identify emergent leadership behaviors in OSS projects from issue comments data. To our knowledge, this is the first study on automatic mining of emergent leadership behaviors in OSS issue discussion communications.
- The evaluation of iLead using 10,000 comments extracted from 10 popular OSS projects, achieving promising results; and empirical study for revisiting the leadership landscape in OSS projects.
- Publicly accessible dataset and source code¹ to facilitate the replication of our study and its application in other contexts.

2 BACKGROUND AND MOTIVATION

2.1 Leadership Models and Theories

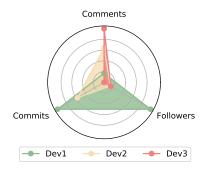
It is important to note the difference between "leadership" and "leader", in that leadership status is not necessarily based on an OSS community position or designated authority, i.e., leader status [21, 48, 62]. Instead, leaders emerge and earn their status through incremental influences and contributions to OSS communities. In the context of this study, we adopt the Yukl's definition of leadership as "influence exerted over other people to guide, structure, and facilitate relationships in a group" [65]. Existing work on leadership theories and styles can be largely categorized into two polarized yet important leadership styles: transformational and transnational, depending on followers' behavioral response [3]. Transformational leadership [12] requires leaders to work with followers to stimulate the enthusiasm of subordinates by creating a vision and motivation. Transactional leadership [29] is a kind of leadership behavior pays more attention to organizational management and performance. However, both styles are more focusing on traditional organizations with clear management hierarchy and team structure, which are not directly applicable to OSS communities.

In recent years, a number of new theories characterizing leadership in virtual teams were proposed, such as shared leadership [13, 65] and emergent leadership [64]. Shared leadership [50] is a leadership mechanism existing in a team with clear hierarchy, and it is a process in which individuals in the group influence each other and work together towards the goal. Emergent leadership [64] is a spontaneous leadership mechanism as an emergent phenomenon that develops over time through group processes.

2.2 OSS Community, Project and Developers

An OSS community is an ecosystem comprising of a set of closely-related OSS projects, which draws expertise and contributions from a pool of developers [33]. OSS developers usually assume certain roles by themselves according to their personal interest, rather than being assigned a task by someone else [46]. This freedom in assuming roles/responsibilities, while significantly differing from traditional organizations, necessarily nurtures the open innovation that drives the success of OSS. As examples of such open contribution, developers may invest effort on a specific OSS project instead of

¹https://github.com/20210827/iLead



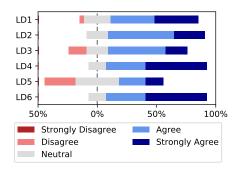


Fig. 1. Three example developers in atom

Fig. 2. Survey results

others, and he/she may choose to submit code contribution, or participate in issue discussions. As examples of developer impact measures, the number of commits and the number of followers typically reflect a developer's relative influence in the OSS community. Consequently, an OSS developer's role in an OSS project may vary significantly based on his/her interest, availability, and expertise. Nakakoji et. al. [46] proposed to classify OSS developers' roles into 8 types including project leader, core member, active developer, peripheral developer, bug fixer, bug reporter, reader, and passive user. The order implies the relative significance of their contribution. In addition, in measuring developer's contribution, they usually play greater weight on the developer's code contribution [25, 38], i.e., the number of commits, which is readily accessible from OSS platforms such as GitHub or SourceForge.

2.3 A Motivational Example

To explore differences among different leadership indicators, we compare three indicators, i.e., #commits, #followers, and #issue comments, across a set of contributing developers from the $atom^2$ OSS project. Figure 1 illustrates the comparison results of three example developers. It is clear that developer Dev1 (id:1476) has the highest #commits and #followers among all three, followed by Dev2 (id:7910250) then Dev3 (ID:4525388). The first two developers, i.e., Dev1 and Dev2, will likely be considered as having greater contributions than the third one following existing leadership indicators [5, 39, 44, 45] because they have more code contributions or have more followers. However, although developer Dev3 has less commits, he has actually published the most issue comments. Specifically, there are 267 comments made by Dev3, including the following examples:

- "Have you tried safe mode? Try -safe and see if the issue stops."
- "What is your atom version?"
- "Duplicate of #1667"

Through these communications, the developer *Dev3* either provides alternative solution to the issue reporter, or asks clarification questions, or redirects the issue to other similar ones. These behaviors are actually representative of emergent leadership and managerial responsibilities, which are essential to facilitate collaborative discussion and issue resolution.

To further examine the differences of these indicators at project level, we rank all developers in *atom* by the number of comments and the number of commits, respectively, and compare the top 25% of the developers in the two rankings.

²https://github.com/atom/atom

Note that the number of followers is defined at the overall GitHub, therefore, we did not include it for the project level comparison here. We find that large mismatch exists between these two lists. Specifically, 59% of developers on the two rankings are different, and those appearing in both rankings correspond with different orders as well. This indicates that existing code commit metric might significantly underestimate the influences and contributions of many developers from issue discussion aspects. Besides, although the number of comments can represent the frequency of developer' participation in issue discussion to a certain extent, it misses abundant semantic details reflecting underlying leadership behaviors, as illustrated in the examples above. Therefore, this study is motivated to develop an effective approach to mine and identify such communication associated leadership, to offer a complimentary perspective to existing studies.

As an example, an OSS community can integrate our leadership identification tool with their issue tracking system, to monitor different leadership behaviors at a particular frequency. The metrics data can enable informed decisions such as understanding leadership dynamics and influences, recommending leadership styles for stimulating issue discussion, reflecting on leadership variation, etc.

3 CATEGORIZING LEADERSHIP BEHAVIORS

This section presents a set of six categories of "emergent leadership" behaviors in collaborative issue discussion processes in OSS communities. Unlike previous studies on OSS leadership, we take a different approach in this study, by considering leadership as an influence process, not the totality of the behaviors of a particular developer. To that end, we aim at examining leadership behavior at individual comment level, instead of at the developer level (e.g. measuring the corpus of comments by a person). Such a conceptualization of leadership has two main advantages. First, it allows us to observe comment-level leadership behavior with finer granularity. Specifically, this enables us to examine immediate impact of each leadership behavior on other developers (e.g., as demonstrated by responding comments from others) and contextual factors surrounding leadership influence process, such as stimulating downstream discussions, and/or accelerating issue resolution. Second, using these comment-level leadership behaviors would also allow us to aggregate behaviors to the individual level. As OSS communities are constantly changing, more research is needed for aggregating comment-level metrics at the individual/team level, with respect to particular usage scenarios.

Existing leadership literature has offered many taxonomies of leadership behaviors. We build the six categories based on reviewing transformational leadership [2], Yukl's comprehensive taxonomy [66, 67], shared leadership[39], and emergent leadership[24]. The six leadership categories were adapted from 12 Yukl's comprehensive taxonomy [66, 67] according to the matching leadership functions in the OSS context, particularly for achieving the goal of efficient issue resolution. Table 1 summarizes these leadership categories with description and comment examples, along with the mapping with the 6 selected Yukl's categories [66, 67]. The other 6 Yukl's categories concerns more on change-oriented, hierarchy and external components of leadership, which are not included in this study due to less relevance to the OSS issue discussion context. More specifically, the leadership included in this study is as follows:

LD1 (Proposal). This category corresponds to the short-term planning category from Yukl's taxonomy, and refers to the idea proposal type of issue discussion, such as proposing ideas for resolving an issue, suggesting unexplored processes and resources to be explored. Issue discussion is an intelligence-intensive problem solving process, and LD1 discussions usually shape potential issue resolution plans, and accelerate the issue resolution cycles. In addition, we believe it is also consistent with the intellectual stimulation/inspiration category in transformational leadership [2, 47].

LD2 (Redirection). This category corresponds to the supporting category from Yukl's taxonomy, and refers to comments providing supportive information to redirect attentions to a new topic (e.g., concepts, processes, or resources) or to a more relevant information sources. Due to the open and interconnected nature of OSS communities, such topics

Manuscript submitted to ACM

Comment Type	Description	Examples	Yukl's Mapping	
LD1 (Proposal)	Propose a solution or alternative	1. Try install Bitcoin Core 0.17.	Planning	
		2. You could try to remove Homebrew's binutils.		
LD2 (Redirection)	Guide to other topics or places to discuss	1. This is a duplicate of #17576.	Supporting	
		2. This should be an issue in repo.		
LD3 (Confirmation)	Confirm an issue or opinion	1. I agree with @sipa.	Recognizing	
		2. Confirming that I got the same error on		
LD4 (Inquiry)	Ask for more information	1. Where did you get the code from?	Consulting	
		2. What version of Boost are you using?		
LD5 (Operation)	Suggest to close or reopen the issue report	Going to close due to lack of information.	Monitoring operation	
		2. It's just a workaround plz reopen this.		
LD6 (Volunteer)	Voluntary acceptance of a task	1. I would like to work on this.	Clarifying roles	
		2. I will open a PR fixing this soon.		

Table 1. Categories of emergent leadership behaviors in OSS projects

and information sources may be either internal (e.g., a related issue within the OSS project) or external (e.g., a site outside the OSS project).

LD3 (Confirmation). This category corresponds to the recognizing category from Yukl's taxonomy, and refers to the leadership function in confirming or reiterating an idea or opinion, which facilitates the consensus processes among members, through implicit voting or shaping shared perception among developers.

LD4 (**Inquiry**). This category corresponds to the consulting category from Yukl's taxonomy, and refers to cases where developers ask for more clarification question about the issue being reported/discussed. Example follow-up inquiries include asking for specific project version info or issue reproduce steps. Such issue comments typically provide some sort of confirmation to former opinions, and facilitate the issue discussion threads.

LD5 (Operation) and **LD6 (Volunteer)** represent two distinct types of administrative functions during issue life cycle, corresponding to the monitoring and clarifying roles categories, respectively, from Yukl's taxonomy. LD5 refers to comments related to either closing or reopening an issue, and LD6 relates to volunteering in undertaking some wrap-up bug fixing actions, respectively.

To validate the representativeness of the 6 comment-level leadership categories, we design and conduct a survey to gather feedback from OSS developers. The main question asked is whether they agree that these categories capture the relevant leadership behaviors which contribute to the issue discussion and resolution processes. We distribute the survey to the contributors of 10 popular projects (as shown in Table 2) on GitHub, and received 27 responses³. The results, as shown in Figure 2, indicates that an average of 71% of developers believe that these six categories of leadership provide a good representation of emergent leadership behaviors in OSS issue discussion communications. Furthermore, even the most controversial LD5 has no strong disagree, and the number of disagrees is only 26% (7/27). One possible reason for the disagreement on LD5 might be because developers tend to link leaders or leadership behaviors with more technical contribution, while LD5 mostly deals with process monitoring and technical focused.

4 APPROACH

When developers explain something similar, they are likely to use recurrent textual expressions in their comments [63]. Existing studies frequently rely on such expressions to extract linguistic patterns in order to facilitate context understanding and task automation [17, 40, 49, 58, 69]. In this section, we present a novel approach, iLead, which

³Due to space limit, we put the details of the survey on our project website.

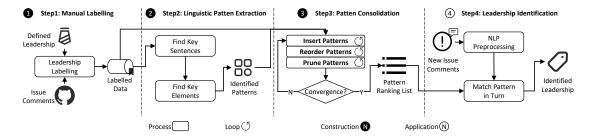


Fig. 3. iLead overview

employs a set of heuristic linguistic patterns to automatically mine developers' issue comments, and identify the emergent leadership behaviors (as described in Section 3).

Figure 3 illustrates the overview of the iLead approach, consisting of four steps including: 1) Manual labelling of issue comments, in order to establish ground truth of leadership labels (Section 4.1); 2) Extracting linguistic patterns specifying the leadership behaviors from the comment in labelled data (Section 4.2); 3) Developing a pattern consolidating algorithm to reduce the penalty due to potential pattern over-fitting and over-loading (Section 4.3); 4) For new issue comments, applying the pattern ranking list output from the previous step (*Step 3*) to identify the leadership behaviors of the issue comments (Section 4.4).

4.1 Manual Labelling

In this step, we will mark each issue comment with a corresponding leadership label (i.e., LD1, LD2, ..., or LD6), as introduced in Section 3. Specifically, the first three authors individually label the same issue comments, then discuss and merge the results on each individual issue comment. The merged leadership labels serve as the ground truth in the later steps, i.e., pattern consolidation and evaluation.

To ensure the validity of the manual labelling outcomes, we adopt a labelling process in the existing work [51], and the process includes four activities, i.e., defining a labelling schema, independently labelling by authors, group reviewing of individual labels, and building a final consensus corpus.

Specifically, the labelling schema is based on the six categories of emergent leadership behaviors, as mentioned in Section 3, plus an additional category (N) for comments not containing any leadership behaviors. A group discussion is conducted to ensure that everyone involved in labelling have a common understanding of these six leadership behaviors in open source issue discussion context.

Then, each comment will be independently labelled by the first three authors and produce the label of a corresponding leadership behavior, or N indicating no leadership behavior. The average Cohen's Kappa is 0.75 between each pair of annotators, indicating substantial agreement among the annotators.

Next, a group meeting is held to review individual labelling results. When a discrepancy arises, each author involved in labelling provides rationale for his/her choice, and the group then discusses and reaches an agreement.

Finally, after several rounds of discussions, agreements are reached for every comment and we get the final consensus corpus.

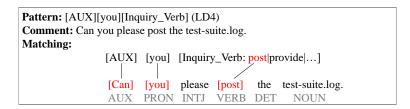


Fig. 4. Example of leadership identification with pattern

4.2 Linguistic Pattern Extraction

The second step is to manually extract linguistic patterns from the textual content of issue comment, which can be used for predicting each of the six leadership behaviors. For all comments with leadership labels (i.e., LD1 to LD6), we first extract the key sentence in the comments which determines it belonging to certain leadership behavior. Take the following LD4 (Inquiry) comment as an example "Can you provide more information? I've just done a build of master using the 64 bit Windows instructions and it's working correctly." The key sentence is "Can you provide more information?", since the other sentences help little in identifying the leadership behaviors. After that, the more structured part of the sentence like "...can...you...provide..." will be identified and the linguistic pattern will be generated base on it. In this example, a pattern containing three elements about "can", "you" and "provide" will be extracted.

The patterns generated from the above process is relatively large, if only concrete word/phrase elements are considered. For example, "Can you provide more information", "Could you give us the steps to reproduce", "Would you mind uploading the dockerfile" and "Are you sure you're not running master" all have similar structures and belong to LD4 (Inquiry). Since the words are different, we get four patterns: "can you provide", "could you give", "would you mind uploading" and "are you sure". This will lead to a huge pattern set, which is difficult to maintain.

In order to reduce the number of patterns, we follow the existing research [40, 58] by considering not only words as elements of a patterns, but also the Part-of-Speech(POS) tag of the words, and use dictionary to represent a group of words that have similar usage. For a pattern with *Dictionary*, it will be used as a template which means each word in the dictionary can be applied to the corresponding positions and formulates a specific pattern. As a result, the four LD4 patterns in the previous paragraph can be merge into "[AUX(POS)] [you(word)] [Inquiry_Verb(Dictionary)]". This not only reduces the total number of patterns, but also enables more sophisticated pattern matching. Figure 4 shows how this pattern can be applied in leadership identification.

Two of the authors participate in the process and employ the open card sorting [54] during this process. Specifically, we do not pre-define the relevant patterns and dictionaries. Each author extracts patterns independently on the collected data and organize their corresponding pattern set. Then group meetings are conducted to discuss and merge the pattern sets. After that, the merged pattern set will be maintained by the authors together and each author can add new patterns to the set or refine existing patterns in their subsequent pattern extraction. Finally, we get a set of patterns that can be used to match comment and identify leadership. For a complete list of all linguistic patterns, please refer to our project website on GitHub, and a total of 304 patterns are extracted for the follow-up leadership identification (see details in Table 3).

4.3 Pattern Consolidation

In this step, we develop an iterative pattern consolidation algorithm to merge patterns extracted from individual projects, monitor and tune the performance change of iLead against the newly merged pattern set in each iteration. This step is essential due to two main reasons: (1) patterns extracted from individual projects are at the risk of pattern over-fitting; (2) in cases that multiple patterns can be simultaneously applicable, appropriate priorities need to be assigned in order to avoid complexity and ambiguity caused by overloaded patterns. To address the pattern over-fitting and overloading issues, we develop an algorithm to automate the pattern merge and consolidation process. This algorithm consists of three sequential tasks (i.e., insert, reorder, and prune), and stops when the performance of iLead converges to a satisfactory level. We will elaborate each of the three tasks next.

4.3.1 *Insert Patterns*. This step takes two input: one is a target pattern ranking list (i.e., initially an empty list), and the other is a new pattern set to be inserted into the target pattern ranking list.

Specifically, for each new pattern, the algorithm will iteratively examine the impact of merging it into the target pattern ranking list, and only accept it when it positively contributes to the prediction performance of iLead. Following an existing study [40, 58, 69], we choose the F1-Score (mentioned in 5.4) as the performance metric. Specifically, the algorithm will probe every candidate position in the pattern ranking list for inserting a new pattern, calculate the performance of iLead in each probing, record the best performance and the corresponding optimal insertion position. As a result, a new pattern will only be inserted into the merged list if it positively contributes to the performance of iLead, and the optimal position is the one corresponds to the best performance. Otherwise, the new pattern will be discarded and will not be merged into the pattern ranking list.

4.3.2 **Reorder Patterns**. We observe that, after the previous merging process, some issue comments may be incorrectly identified according to a loosely relevant pattern with a higher priority, while a more relevant pattern is ranked lower in the list and never be applied. To address this, the algorithm reorders the ranking of patterns by examining these incorrectly labelled comments, to promote the overall performance on the whole dataset. The ranking update is conducted iteratively and dynamically in terms of each incorrectly labelled comment by current pattern ranking list.

In detail, for an incorrectly labelled comment (the comment mislabelled by the current pattern ranking list), we define the highest ranking pattern which correctly labels it as *realwinner*, while the patterns which incorrectly match it and higher than *realwinner* as the *fakewinners*. If *fakewinners* list is not empty, algorithm will try the ranking update. Specifically, algorithm will insert the *realwinner* ahead of the *fakewinners* in the pattern ranking list, and examine the leadership labelling performance on the whole dataset (i.e., F1-Score in Section 5.4). If the performance increases with this new pattern ranking list, the update is accepted and the current pattern ranking list will change to the new pattern ranking list. Then the algorithm iterates with the next incorrectly labelled comment found by the current pattern ranking list until it remains unchanged. As a result, we obtain the updated pattern ranking list which will be used next.

4.3.3 **Prune Patterns**. During this task, the algorithm will identify excessive patterns whose existance negatively contributes to the performance of iLead. Similar to the previous task, the pattern pruning is conducted iteratively and dynamically by examining each incorrectly labelled comment.

In detail, for each incorrectly labelled comment, the algorithm first detects the *realwinner* (i.e., the highestly ranked pattern which is able to correctly label the comment), as well as all *fakewinners* (i.e., the patterns which incorrectly match the comment and higher than *realwinner*). After that, the algorithm will remove the *fakewinners* in the pattern

ranking list, and accept the update if the performance (i.e., F1-Score in Section 5.4) increases. Then the algorithm iterates with the next incorrectly labelled comment found by the refined pattern ranking list until it remains unchanged.

After that, we will obtain the pattern ranking list and if the patterns from a new project need to be added, the algorithm will start executing from *Insert Patterns* again and use this list as the input (instead of the original empty list).

4.4 Leadership Identification

After completing the construction of pattern ranking list, we can use it to identify leadership behaviors. Specifically, for a new issue comment, iLead first preprocesses it to get the words and the Part-of-Speech (POS) tags. Then, it matches the preprocessed data with element(s) specified in a given pattern. We employ the following two distance constraints between elements, adapted from existing study [57], to promote the performance of leadership identification. The first constraint is that the latter matched element cannot appear before the prior matched element, since the pattern is extracted by traversing the words in sequence. The second constraint is that the distance between two matched elements should be no more than three, which is to avoid the mismatch of patterns since the leadership behaviors are usually expressed in a concise way. iLead will automatically match the comment with the patterns in the list in order, and following existing study [58], we choose the highest priority strategy which means the leadership of the comment is identified by the first matched pattern.

5 EXPERIMENT DESIGN

5.1 Research Questions

To evaluate the proposed approach, we answer the following three research questions:

- RQ1 (Performance Convergence) When can the generated linguistic patterns reach the state of convergence?
- RQ2 (Performance Evaluation) How effective is iLead in leadership identification on new projects?
- RQ3 (Baseline Comparison) Can iLead outperform other techniques in leadership identification?

5.2 Subject Projects

We collect 10 popular open source software projects from GitHub for evaluation. The projects are selected with the following criteria: (1) Popularity: with 10k+ stars and 1.5K+ forks; (2) Activeness: regularly updated within 1 month of the data collection date; (3) Diversity: representing different OSS organizational structures from various domains. A brief summary of the experimental projects are summarized in Table 2.

For the experimental projects, we use the GitHub's REST APIs to crawl the issues and related comments⁴. We only crawl the comments of closed issues since they represent the complete issue life-cycle. In addition, since it is our interest to compare the emergent leadership behaviors with traditional leadership indicators, we also extract two indicators at the developer level: the number of commits and the number of followers. For each project, we randomly sample a portion of issues and obtain 1,000 related comments (in Table 2), and conduct manual labelling, as introduced in Section 4.1. The overview for the ground truth of the leadership labels is shown in Table 2.

5.3 Data Pre-processing

We observe that some developers would refer to the previous comments in their own comments. The referenced descriptions may contain another developer's leadership, and should be removed to reduce potential noise leading to

⁴The data is crawled between July 2020 and February 2021.

Leadership labels Basic information ID Project Link Domain #Comments #Issues LD1 LD2 LD3 LD4 LD5 LD6 TotalLD P1 bitcoin [8] digital currency 57k 28,504 4.997 95 98 123 132 55 18 521 P2 450 sklearn machine learning 47k 40,432 6.287 81 47 63 73 105 81 55 Р3 ember.js [22] js framework 22k 32,225 5,833 69 54 132 77 89 43 464 P4 brew [30] package manager 29k 14,269 2.694 74 55 89 85 34 38 375 P5 atom text editor 56k 83,004 15,621 50 255 141 42 550 56 6 P6 bokeh 74 79 [10] visualization 15k 21,432 4.573 76 65 60 56 410 P7 entity framework 55,820 11,596 515 efcore [19] 11k 87 47 72 140 162 P8 knex [37] SQL query 15k 10,337 2,177 28 87 104 56 39 480 166 P9 roslvn compiler 15k 76,862 14,722 123 74 37 411 P10 solidity [23] program language 12k 11,567 2,513 44 52 84 90 46 28 344

Table 2. Subject projects

incorrect matching. In this study, we implement this via removing the node which is used to highlight the reference descriptions in the comment's HTML text.

In addition, we use Stanford natural language processing tool Stanza [52] to process the text in lowcase, tokenization, part of speech (POS) tagging and lemmatization, so as to obtain the data applied to the construction process of iLead. We also use Urlextract⁵ to annotate the URL contained in the comments, which is the element in certain patterns.

5.4 Evaluation Metrics

We employ three commonly-used metrics to evaluate leadership identification performance, i.e., precision, recall and F1-Score.

Precision measures how precise a classifier is in predicting the comments with specific category of leadership.

Recall measures the ability of classifier to find the comments with specific category of leadership.

F1-Score is the harmonic mean of precision and recall.

For the precision, recall, and F1-Score of multiple categories of leadership (as shown in Figure 6), we calculate these metrics of each category and use the average value.

5.5 Experiment Setup

5.5.1 Experiment for answering RQ1. We start with the set of linguistic patterns established from one project, and then employ an iterative process as follows: 1) add another set of linguistic patterns established from a new project; 2) apply the pattern consolidation algorithm from Section 4.3 to generate a pattern ranking list; 3) input the pattern ranking list to iLead for leadership identification; 4) evaluate the performance of iLead using the previously defined metrics. Throughout the iterative process, we continuously monitor the performance change of iLead, and stop the process when the performance converges (i.e., when the additional performance gain is less than 0.01). By the convergence time, we stop to incorporate new patterns, and finalize the pattern ranking list in iLead as the default pattern list used for further experiment and evaluation.

5.5.2 **Experiment for answering RQ2**. To evaluate the performance of iLead, we use the remaining projects (projects unused in RQ1) as test sets and measure the identification performance of iLead on each test project. In addition, we will

⁵https://github.com/lipoja/URLExtract

	Pa	ttern extra	ction resul	ts			
Projects	P1	P1-P2	P1-P3	P1-P4	P1-P5		
#Patterns	198	250	287	300	304		
#Add	-	55	41	13	4		
#Delete	-	2	4	0	0		
#Change	3	5	1	0	0		
Leadership identification performance							
Precision	0.81	0.82	0.83	0.83	0.83		
Recall	0.74	0.82	0.84	0.85	0.85		
F1-Score	0.76	0.82	0.83	0.84	0.84		

Table 3. Performance change with iterations (RQ1)

also examine the performance of iLead across the six leadership categories to obtain more comprehensive performance results.

5.5.3 Experiment for answering RQ3. We employ commonly-used machine learning based and deep learning based text identification approaches as baselines to compare with iLead. Specifically, we use Term Frequency and Inverse Document Frequency (TF-IDF) [43] to vectorize textual descriptions of comments, and use the term vectors as input for identification. For selection of machine learning classifiers, we evaluate four commonly-used machine learning models: Logistic Regression (LR) [6], Support Vector Machine (SVM) [15], Decision Tree (DT) [11], and Random Forest (RF) [28], and due to the limited space, we only show the best classifier (SVM). For the deep learning based classifiers, we utilize TextCNN [36], BiLSTM [42], CNN and RNN with Bert pre-training model [18].

In addition, we further explore whether our extracted patterns are helpful to the machine/deep learning models, following the existing studies [40, 58, 69]. In detail, we enrich the model inputs with the knowledge about pattern matching, which will be consider as a boolean vector, i.e., if a comment match with certain patterns, the corresponding position of the vector is set to 1, otherwise it is 0. We also run the above mentioned machine/deep learning models on the combined vectors, and they are short for <Model_Name-P> (e.g., TextCNN-P). In total, we have 10 baselines for comparison.

6 RESULTS AND ANALYSIS

6.1 Answering RQ1. Performance Convergence

As introduced in Section 5.5, We start with the pattern set extracted from one project, and iteratively add a pattern set extracted from a new project, consolidate into the pattern ranking list following the algorithm described in Section 4.3. In each iteration, we record the number of patterns added, deleted and changed, as well as the performance of iLead. We find that the performance of iLead converges after five iterations. Table 3 summarizes the results from each of the five iterations.

Across the five iterations, we can see that the number of patterns increase from 198 to 304, and the accuracy increases from 0.81 to 0.85. We also observe that different projects share common linguistic patterns in issue comments, as seen by the smaller number of newly added patterns in each iteration. In addition, the most pattern changes caused by newly added patterns is in the second iteration, when adding the pattern set of *sklearn*, corresponding to 55 added, 2 deleted, and 4 changed patterns. Though there are new patterns added in the later two iterations as well, however, their Manuscript submitted to ACM

Projects AVG bokeh efcore knex roslyn solidity Precision 0.80 0.83 0.82 0.82 0.84 0.82 Recall 0.77 0.84 0.77 0.81 0.78 0.79 F1-Score 0.77 0.83 0.79 0.81 0.80 0.80

Table 4. Performance on unfitted projects (RQ2)

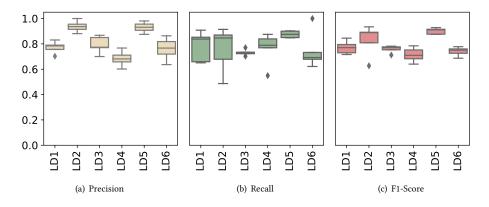


Fig. 5. Performance on six leadership categories (RQ2)

added impact on the performance is diminishing. When the fifth project is added, only very slight changes are observed in the pattern list and the three metrics. This means that the information from the fifth project does not bring much contribution to the pattern extraction, consolidation and performance increase. We believe that the results can support the conclusion on the sufficient state of convergence.

In the remaining experiments, the set of patterns and their rankings obtained from the fifth iteration are finalized in iLead, and used to investigate other research questions.

6.2 Answering RQ2. Performance Evaluation

Table 4 presents the performance of iLead applying to the five new projects. The results show that iLead has an average precision of 82%, and an average recall of 79% across the five projects. In addition, the results are quite close to those from Table 3, which indicates the performance stability of iLead, as well as the generality of the pattern set and their relative ranking in iLead. It is also noticeable that the performance is slightly higher than the inter-rater agreement level of 0.75 during the manual annotation process (Section 4.1). A possible reason is associated with the group consensus based on majority vote. More specifically, some inter-raters discrepancy involves long comments containing multiple sentences, in which two or more sentences insides a long issue comment may correspond to different leadership labels. We employed group discussion and majority vote to keep 1 dominant leadership label. This handling might lead to the performance of iLead slightly higher than the interrater's initial agreement level. Nevertheless, we observe a slight decrease in precision and recall, compared with the performance in Table 3. One possible reason is that these five new projects might contain new patterns in expressing leadership behaviors, which necessitates further tuning or addition investigation to include more patterns to improve the performance.

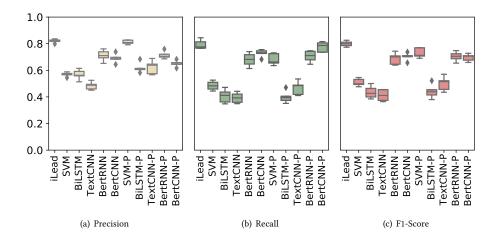


Fig. 6. Performance comparison with baselines (RQ3)

We also examine the performance of iLead in identifying each leadership category on the five new projects, as shown in Figure 5. From the comparison, we can see that LD2 (Redirection) and LD5 (Operation) are among the categories with highest performance, meaning relatively easier to identify. These two categories both reach a median precision over 0.9, median recall over 0.8. The lowest performance is observed in LD4 (Inquiry) and LD6 (Volunteer), yet with a satisfactory median of 0.70 F1-Score.

6.3 Answering RQ3. Baseline Comparison

Figure 6 demonstrates the comparison of the performance of iLead with the ten baselines. We can see that iLead outperforms all the baselines in all the investigated metrics, e.g., the median F1-Score of iLead is 13% (i.e., (0.80-0.71)/0.71) higher than the best baseline *SVM-P*. Furthermore, we conduct Mann-Whitney U test between the performance of iLead and the baselines. Results show that iLead significantly (p-value <0.01) outperform all the baselines in F1-Score, and significantly outperform most of the baselines in precision (except SVM-P) and recall (except BertCNN-P). This implies the effectiveness of our extracted patterns as well as the pattern consolidation algorithm in iLead.

Among the baselines, the machine learning classifier with combined features, e.g., *SVM-P*, is better than its counterpart only with textual features. For example, the median F1-Score achieved by *SVM-P* is 42% (i.e., (0.71-0.50)/0.50) higher than its counterpart *SVM*. This implies our designed patterns are useful even when integrated into the machine learning models.

Nevertheless, they still underperform our proposed approach, because our proposed approach has well-designed mechanism to optimize the patterns and their rankings, while the machine learners consider less of the application order of the patterns. We also observe that the deep learning based approaches do not achieve a satisfactory result as anticipated. This might because training an efficient deep learning models requires a large amount of labelled data, yet the dataset utilized in our experiment is far from that amount. Besides, there are many low frequency words which distinguish a certain category of leadership behaviors from others. But the deep learning models are not good at capturing these low frequency words.



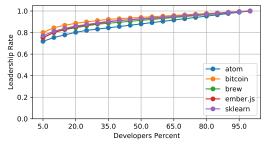


Fig. 7. Distribution of leadership

Fig. 8. Cumulative distribution of developers' leadership

7 REVISITING OSS LEADERSHIP LANDSCAPE

This section describes a large-scale empirical evaluation with iLead on OSS projects. We apply iLead on all issue comments of the 5 fitted projects (from which the patterns are extracted) as shown in Table 2, to automatically identify the leadership behaviors, and explore: 1) the distribution of leadership behaviors; 2) the correlation between our emergent leadership behaviors and existing leadership indicators. We also conduct a human evaluation to verify the usefulness of iLead.

7.1 Distribution of Leadership Behaviors

Figure 7 illustrates the distribution of the identified emergent leadership behaviors on the five projects. It is observed that: (1) Overall, leadership-associated comments account for 39% of all comments on average (ranging from 37% to 44%); (2) The Top-3 dominant leadership categories are LD1 (Proposal), LD3 (Confirmation), and LD4 (Inquiry), corresponding to an average of 7%, 11%, and 12% of all comments, respectively; and (3) The other three types, i.e., LD2 (Redirection), LD5 (Operation) and LD6 (Volunteer), are observed at less frequency, corresponding to an average of 3%, 2%, and 3% of all comments, respectively. These results provide quantitative evidences on developers leadership contribution in issue discussion, i.e., through brainstorming ideas, information inquiry, and consensus-based decision making.

Figure 8 illustrates the pareto curves of individual developers' leadership-related actions in issue discussion, across the five projects. Two outstanding results are noted: (1) On the one hand, an average of 11% (ranging from 5% to 20%) developers contribute to more than 80% of leadership behaviors; (2) On the other hand, the majority (i.e., 61%) developers (ranging from 40% to 75%) are associated with merely less than 10% of leadership behaviors. This implies the potential need for more research to study OSS leadership through this new lens, e.g., supporting OSS leadership evolution and individual/team leadership coaching.

It is conceivable that leadership indicators may be correlated by the number of comments, corresponding to the common perception that when a developer publishes more comments, they are practicing more leadership. Indeed, we can observe the Top-11% highly rated developers contribute 80%+ leadership behaviors on average, as shown in Fig. 8. Particularly, for the Top-5%, the correlation with the comment counts is 0.89. However, the focus of this study is be able to identify and cultivate emergent leadership behaviors among the majority middle-level developers to become more effective contributors and leaders. To explore whether the proposed emergent leadership indicators can offer more insights on such middle-level developers, we extracted two emergent leadership indicators for the middle 60% of developers ranked by the total number of comments. The two emergent leadership indicators are the leadership count and the leadership percentage (i.e., leadership count/total comments) across all six leadership categories (i.e.,

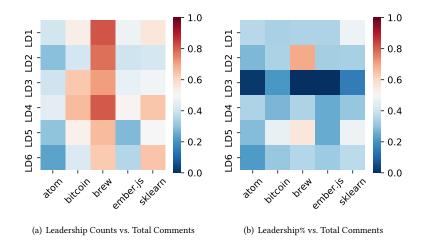


Fig. 9. Correlation between leadership indicators and total comments at individual developer level

LD1 LD6). Figure 9 shows the correlation results between the two new indicators of leadership and the traditional leadership metric of total comments at individual developer level, using data from 5 OSS projects. From Figure 9(a), the high correlation is only observed on the *brew* project, but not on the other 4 projects.

As shown in Figure 9(b), the correlation between the percentage of leadership comments and the total comments is much lower, suggesting that developers making the same amount of comment contribution might correspond to different degree of emergent leadership behaviors. It can be concluded that from both charts, the number of comments does not necessarily mean more leadership, which indicates that our proposed leadership indicators provide a novel lens to offer new insights to study OSS leadership.

7.2 Comparison among Leadership Indicators

iLead is committed to identify the emergent leadership of developers and complement the existing viewpoints of developer contributions. To that end, we conduct a comparison analysis to measure OSS developers using different leadership metrics. For each developer, we aggregate the number of his/her comments corresponding with a specific leadership category. Then we conduct a correlation analysis to compare individual developers' emergent leadership behaviors with two existing leadership indicators, i.e. code contribution and community followers.

Figure 10(a) shows the correlation between the number of code commits and the number of identified leadership behaviors across the six emergent leadership categories. In general, the correlation between the number of leadership behaviors and the number of commits is low. This again confirms that emergent leaders are not always the leading code contributors. We can observe a relatively higher correlation between these two indicators in the *brew* project. This may be because the *brew* project entrusts its daily management activities to a third-party service organization[56], and the developers focus more on the code contribution and issue discussion, i.e., resulting a higher correlation between these two indicators. Moreover, Figure 10(b) shows a much lower correlation between the number of followers and the number of leadership behaviors than that from Figure 10(a). This confirms one of frequent findings in existing leadership literature [34, 35], which refers that leaders are not depending on how much followers they have, but how Manuscript submitted to ACM

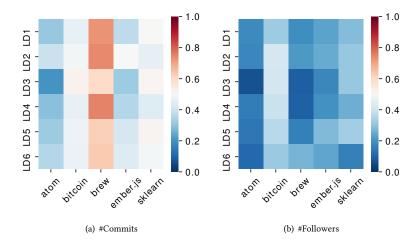


Fig. 10. Correlation between existing contribution indicators and #leadership

many leaders they create. In OSS community, there are a variety of motivating factors for following or not following, such as user's interest in learning, socializing, obtaining updates, or easy access to others, etc. [9].

We think there are a couple of reasons for such low correlation. First, OSS developers may have different roles and preferences in different projects. For example, some prefer coding, while others participate more in the discussion. This suggests the need for a more dynamic and diverse lens to measure OSS leadership. Second, this reveals further opportunities for in-depth analysis to demonstrate the correlation between emergent leadership behaviors with existing leadership indicators, e.g., through data segmentation, time-lag analysis, etc. The low correlation between the identified emergent leadership behaviors and existing leadership indicators suggest that the identification of emergent leadership helps to broaden existing code contribution and followers viewpoints to embrace and recognize more comprehensive perspectives.

7.3 Leadership Influence

As discussed earlier, we consider leadership as an influence process and measure leadership at individual comment level. Different leadership behaviors might have different impacts and the most obvious effect is on other discussions. To explore the influence of each leadership, we extracted 10 features encompassing two categories for each issue from the 5 projects. The first category focuses on examining the influence on process stimulation, and the second category focuses on examining issue resolution efficiency. We set an effect window of 24 hours, to extract process stimulation features reflecting later issue discussions after a leadership comment is posted.

- other_commenter: The number of commenter who is not the issue reporter or the author of this comment.
- comment_num: The number of comments.
- reporter_response: The number of comments posted by the issue reporter.
- self_response: The number of comments posted by the author of this comment.
- other_response: The number of comments posted by other commenters.
- ld_num: The number of leadership comments.

Feature Category LD1 LD2 LD3 LD4 LD5 LD6 (-)*** (+)*** (-)*** other_commenter Process $(-)^*$ # (-)*** (+)*** (-)*** comment num Process (+) # (-)*** reporter_response Process (+)** (-)*** (-)** (+)** (-)*** (+)** (-)*** (+)** self_repsonse Process other_response Process $(-)^*$ (-)* (+)**(-)** # ld num Process (-)**(+)** (+)** (-)*** $(+)^{***}$ (+)*** (-)*** (-)** ld_types Process # (-)*** $(+)^{***}$ (-)*** word_divergence (+) Process # # time_from_start Resolution $(-)^*$ (-)**(+)** (+)*** $(+)^{**}$ time_to_close Resolution (-)** (+)*** (-)*** (-) (+)*** $(+)^*$

Table 5. Leadership hypothesis testing results

- $p \ge 0.05$, p<0.05, p<0.01, p<0.01
- ld_types: The number of leadership types.
- word divergence: set(words) / len(words)
- time_from_start: Time between the comment and the the opening of issue.
- time_to_close: Time between comment and issue closing.

Table 5 shows the hypothesis testing results of the mean features' values of all 6 leadership indicators (Compared with non-leadership) at issue level, using Mann-Whitney U test. We use the positive signs to represent positive/increasing influences and use the negative signs to represent the negative/decreasing influences. Consistent results are observed across the 5 projects, for space consideration, Table 5 only shows the results on *ember.js*. Results from other projects can be found on the study website ⁶.

We conclude: (1) LD1 (Proposal) tends to be able to trigger more discussions between the issue reporter and the LD1 publisher, and it may promote the issue to be solved; (2) Comments after LD2 (Redirection) and LD5 (Operation) are less active. This is intuitive and indicates the issue originator's question is solved (LD2), and the issue has been closed (LD5), respectively, and both may promote the issue to be closed; (3) After the release of LD3 (Confirmation), the number of leadership increased, but the number of comments remained unchanged, indicating that it can promote the emergence of leadership, and it may increase the time for discussion; (4) The discussions after LD4 (Inquiry) are more active; (5) LD6 (Volunteer) has no difference in most features, but in comparison, it is far from the opening and closing of the issue. We observed that it is generally in the late stage of the issue and will start the discussion of issue repair, so it takes longer to close the issue.

In summary, it can be concluded that LD3, LD4 play important role in stimulating issue discussion process, and LD1, LD2 contributes to bring speedy issue resolution.

7.4 Human Evaluation

To evaluate the usefulness of iLead, we recruited 9 developers from three OSS communities (i.e., *bitcoin*, *atom*, *ember.js*) as evaluator. Specifically: (1) use iLead to identify leadership behaviors in these three projects, and then derive a ranked list of developers according to the frequency of leadership behaviors; (2) create two subsets of developers: one with the Top 20 developers (used as study group), and the other with developers ranking 21-40 (used as comparison group); (3)

⁶https://github.com/20210827/iLead

randomly select 10 developers from each group, and ask the human evaluators (from the same OSS project) to vote (Yes or No) for the 20 randomly selected developers, based on whether they believe the names on the list play important role(s) in driving issue resolution progress. Note that the developers are provided in random order (without ranking information) to the human evaluators to avoid introducing potential bias.

The results show that the random developers from the first subset (i.e., Top-20) account for 74% of overall votes, while those from the second subset receives 26% of overall votes. This indicates that, human evaluators' votes are mostly consistent with the leadership status identified by iLead.

8 DISCUSSION

8.1 Practical Insights for OSS Communities

Our findings suggest that leadership in OSS projects should encompass the key facets of emergent leadership behaviors in issue discussion. Based on the results from this study, we provide the following practical insights for OSS community.

First, while developers with large code submissions are naturally considered as leaders/core contributors to the project, OSS projects may not necessarily pay attention to emergent leadership of other developers, who serve important influential roles in guiding/facilitating issue resolution. In addition to recognize project leaders based on their code contribution, OSS projects should also consider developers emergent leadership behaviors in their interpersonal communications, e.g. issue comments and live chat.

Second, with support from iLead, OSS community can automate the following tasks accordingly: (1) to identify developers with the most influential emergent leadership behaviors according to the LD1 (Proposal) category; (2) to identify developers coordinating the most dominant brainstorming and exploration tasks according to the LD3 (Confirmation) and LD4 (Inquiry) categories; and (3) to identify developers facilitating the management and operation tasks such as redirecting a conversation, issue closing and volunteering, according to LD2 (Redirection), LD5 (Operation) and LD6 (Volunteer). With these knowledge, there is also great potential to appropriately automate and streamline such issue discussion and operational tasks, in order to accelerate issue resolution cycles.

8.2 Practical Insights for OSS Developers

We also observe several findings that may shed light to practices related to leadership skill development. First, developers will benefit from learning about different roles of emergent leadership. To emerge and to be accepted as a leader in an OSS project, a developer needs to know how to communicate effectively with others via textual messages. As the core construct of iLead, the linguistic pattern set embodies strong verbal-cues correlated with the emergent leadership. These can be used to coach newcomer developers in effectively inquiring about an issue, collaboratively confirming a candidate solution, or clearly proposing an alternative option, etc.

Second, while intuitively, it may seem that OSS leaders should play greater roles in LD1 (Proposal), since it is the most influential leadership behaviors, we actually observe something slightly different. We observe no much difference between junior and senior OSS developers in three types of leadership behaviors, i.e., LD1 (Proposal), LD4(Inquiry), and LD6(Volunteer). In addition, junior developers tend to conduct more LD3 (Confirmation), and less LD2 (Redirection) and LD5 (Operation); and senior developers conduct more LD2 (Redirection) and LD5 (Operation), and less LD3 (Confirmation). It is easy to understand that LD2 and LD5 rely on comprehensive knowledge about a specific OSS projects. Therefore, we encourage new or in-experienced developers to focus more on LD1, LD3, LD4, and LD6 when looking for opportunities to make contributions to a new OSS project.

8.3 Threats to Validity

The first threat concerns the generality of the proposed approach. The iLead is only evaluated using 10 popular OSS projects, which might not be representative of other projects. However, the projects are from various domains, and the performance convergence experiment is conducted to ensure the comprehensiveness and effectiveness of the finalized linguistic rules.

The second threat comes from the representativeness of the six categories of emergent leadership behaviors. One co-author's area of research is in leadership, and the six leadership categories are carefully identified based on reviewing and examining both traditional leadership styles and emergent leadership studies for virtual teams. In addition, the leadership survey results also show that this provides a meaningful representation of emergent leadership behaviors in OSS issue discussion communications.

The third threat is associated with the results of manual labelling. We follow an iterative process[51], and conduct the labelling by three authors independently, and common consensus are reached after several rounds of discussion.

The last threat is associated with the human evaluation since the number of participants is small. We cannot guarantee that the votes from the human evaluators are fair. To mitigate this threat, we ask the evaluators to vote on 10 randomly selected developers from two subsets, and use the total votes for comparison, to mitigate individual bias.

9 RELATED WORK

Linguistic Patterns. Previous researches noticed that people tend to use recurrent expressions to illustrate similar things, and leverage such phenomenon, existing researches constructed the linguistic patterns for automating various software engineering tasks. Zhao et al.[69] conducted a manual analysis of 980 performance issue reports and extracted 80 linguistic patterns; they then combined these patterns with machine/deep learning algorithms to automatically detect the performance issue reports. Lin et al.[40] manually analyzed 388 API-related sentences and defined 157 linguistic patterns, and used these patterns to classify the opinions about API on Q&A website. Shi et al.[58] proposed an linguistic patterns based approach to automatically understand the feature request. They manually defined a group of patterns and leverage them to generate fuzzy rules which can automatically identify feature request. Panichella et al.[49] manually inspected 500 app reviews and defined 246 patterns to automatically classify the app reviews into categories relevant to software maintenance and evolution. Motivated by these studies, this paper utilized linguistic patterns to identify the leadership behaviors.

Leadership in OSS Community. The leadership behaviors have long been investigated in business and management domain, and have attracted increasing attention in software engineering domain in recent years. Giuri et al. [26] characterized the difference in individual and project-specific characteristics between project leaders and other project members. Bian et al. [7] explored the correlation between leader characteristics and success of OSS projects from the behavioral, structural and cognitive dimensions. Li et al. [39] investigated the relationship between leadership characteristics and the motivation to contribute. Jensen et al. [32] explored the collaborative efforts, the leadership and control structures in OSS community. Tan et al. and Zhang et al. [60, 61, 68] investigated how developers and companies collaborate and communicate with each other during OSS development. By comparison, this paper proposes an automatic approach to identify leadership behaviors in OSS projects, with which one can conduct the leadership related investigation to facilitate the OSS development.

10 CONCLUSION

To facilitate the developing and sustaining health growth of the OSS communities, this paper explores an innovative perspective to look at OSS leadership, and proposes an automated approach, iLead, to mine communication styles and identify emergent leadership behaviors in OSS projects using issue comments data. Based on the constructed six categories of leadership and extracted heuristic linguistic patterns, iLead employs an automated algorithm to consolidate the pattern ranking list for leadership identification. The evaluation results demonstrate the effectiveness of the lightweight heuristic rule-based approach employed in iLead, outperforming 10 Machine Learning baseline models. Applying the proposed iLead to 5 OSS projects, we conduct comparison of identified leadership with existing leadership indicators and offer practical insights on community building and leadership skill development. This study provides a new lens for examining leadership and its influence in OSS projects. With the support of iLead, one can look at how leadership develops, at who tends to become a leader, how leadership evolves in a developer's career, etc.

Future researches include: 1) further evaluate and improve the performance of iLead using more OSS projects; 2) investigate the impact of various emergent leadership behaviors; 3) examine the relationship between emergent leadership and temporal, contextual structure of OSS communities; 4) develop intelligent support for facilitating leadership skill development of OSS developers.

REFERENCES

- [1] atom. 2021. Atom. https://github.com/atom/atom.
- [2] Bernard M Bass and Bruce J Avolio. 1996. Multifactor leadership questionnaire. Western Journal of Nursing Research (1996).
- [3] Bernard M Bass, Bruce J Avolio, Dong I Jung, and Yair Berson. 2003. Predicting unit performance by assessing transformational and transactional leadership. *Journal of applied psychology* 88, 2 (2003), 207.
- [4] Hoda Baytiyeh and Jay Pfaffman. 2010. Volunteers in Wikipedia: Why the Community Matters. J. Educ. Technol. Soc. 13, 2 (2010), 128–140. http://www.ifets.info/abstract.php?art_id=1046
- [5] Xu Ben, Shen Beijun, and Yang Weicheng. 2013. Mining developer contribution in open source software using visualization techniques. In 2013 Third International Conference on Intelligent System Design and Engineering Applications. IEEE, 934–937.
- [6] Joseph Berkson. 1944. Application of the logistic function to bio-assay. Journal of the American statistical association 39, 227 (1944), 357–365.
- [7] Yiyang Bian, Wen Mu, and J. Leon Zhao. 2018. Online Leadership for Open Source Project Success: Evidence from the GitHub Blockchain Projects. In 22nd Pacific Asia Conference on Information Systems, PACIS 2018, Yokohama, Japan, June 26-30, 2018. 189.
- [8] bitcoin. 2021. Bitcoin. https://github.com/bitcoin/bitcoin.
- [9] Kelly Blincoe, Jyoti Sheoran, Sean Goggins, Eva Petakovic, and Daniela Damian. 2016. Understanding the popular users: Following, affiliation influence and leadership on GitHub. Information and Software Technology 70 (2016), 30–39.
- [10] bokeh. 2021. Bokeh. https://github.com/bokeh/bokeh.
- [11] Leo Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. 1984. Classification and Regression Trees. Wadsworth.
- [12] James M Burns et al. 1978. Leadership.
- [13] Jay B Carson, Paul E Tesluk, and Jennifer A Marrone. 2007. Shared leadership in teams: An investigation of antecedent conditions and performance. Academy of management Journal 50, 5 (2007), 1217–1234.
- [14] Traci A Carte, Laku Chidambaram, and Aaron Becker. 2006. Emergent leadership in self-managed virtual teams. *Group Decision and Negotiation* 15, 4 (2006), 323–343.
- [15] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. Machine learning 20, 3 (1995), 273-297.
- [16] Daniel Alencar da Costa, Uirá Kulesza, Eduardo Aranha, and Roberta Coelho. 2014. Unveiling developers contributions behind code commits: an exploratory study. In Proceedings of the 29th Annual ACM Symposium on Applied Computing. 1152–1157.
- [17] Alberto Rodrigues da Silva. 2017. Linguistic Patterns and Linguistic Styles for Requirements Specification (I): An Application Case with the Rigorous RSL/Business-Level Language. In Proceedings of the 22nd European Conference on Pattern Languages of Programs, EuroPLoP 2017, Irsee, Germany, July 12-16, 2017. ACM, 22:1-22:27. https://doi.org/10.1145/3147704.3147728
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. https://doi.org/10.18653/v1/n19-1423
- [19] dotnet. 2021. Efcore. https://github.com/dotnet/efcore.

- [20] dotnet. 2021. Roslyn. https://github.com/dotnet/roslyn.
- [21] Remo Eckert, Matthias Stuermer, and Thomas Myrach. 2019. Alone or Together? Inter-organizational affiliations of open source communities. Journal of systems and software 149 (2019), 250–262.
- [22] emberjs. 2021. Ember.js. https://github.com/emberjs/ember.js.
- [23] ethereum. 2021. Solidity. https://github.com/ethereum/solidity.
- [24] Susan Ferebee and James Davis. 2012. Emergent leadership, persuasion, and trust in virtual leaderless groups. The Exchange 1, 1 (2012).
- [25] GitHub. 2021. GitHub: Viewing a project's contributors. https://docs.github.com/en/github/visualizing-repository-data-with-graphs/accessing-basic-repository-data/viewing-a-projects-contributors.
- [26] Paola Giuri, Francesco Rullani, and Salvatore Torrisi. 2008. Explaining leadership in virtual teams: The case of open source software. Inf. Econ. Policy 20, 4 (2008), 305–315.
- [27] Alexander Hars and Shaosong Ou. 2002. Working for Free? Motivations for Participating in Open-Source Projects. Int. J. Electron. Commer. 6, 3 (2002), 25–39. https://doi.org/10.1080/10864415.2002.11044241
- [28] Tin Kam Ho. 1995. Random decision forests. In Proceedings of 3rd international conference on document analysis and recognition, Vol. 1. IEEE, 278-282.
- [29] Edwin P Hollander. 1978. Leadership Dynamics: A Transactional Perspective. Technical Report. STATE UNIV OF NEW YORK AT BUFFALO DEPT OF PSYCHOLOGY.
- [30] Homebrew. 2021. Brew. https://github.com/Homebrew/brew.
- [31] Carlos Jensen, Scott King, and Victor Kuechler. 2011. Joining Free/Open Source Software Communities: An Analysis of Newbies' First Interactions on Project Mailing Lists. In 44th Hawaii International International Conference on Systems Science (HICSS-44 2011), Proceedings, 4-7 January 2011, Koloa, Kauai, HI, USA. IEEE Computer Society, 1–10. https://doi.org/10.1109/HICSS.2011.264
- [32] Chris Jensen and Walt Scacchi. 2005. Collaboration, Leadership, Control, and Conflict Negotiation and the Netbeans.org Open Source Software Development Community. In 38th Hawaii International Conference on System Sciences (HICSS-38 2005), CD-ROM / Abstracts Proceedings, 3-6 January 2005, Big Island, HI, USA.
- [33] Corey Jergensen, Anita Sarma, and Patrick Wagstrom. 2011. The onion patch: migration in open source ecosystems. In Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering. 70–80.
- [34] Susanne Kean, ELAINE HAYCOCK-STUART, Sarah Baggaley, and Maggie Carson. 2011. Followers and the co-construction of leadership. Journal of Nursing Management 19, 4 (2011), 507–516.
- [35] Robert E Kelley. 1988. In praise of followers. Harvard Business Review Case Services.
- [36] Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, Alessandro Moschitti, Bo Pang, and Walter Daelemans (Eds.). ACL, 1746–1751. https://doi.org/10.3115/v1/d14-1181
- [37] knex. 2021. Knex. https://github.com/knex/knex.
- [38] Amanda Lee and Jeffrey C. Carver. 2017. Are One-Time Contributors Different? A Comparison to Core and Periphery Developers in FLOSS Repositories. In 2017 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM 2017, Toronto, ON, Canada, November 9-10, 2017, Ayse Bener, Burak Turhan, and Stefan Biffl (Eds.). IEEE Computer Society, 1-10. https://doi.org/10.1109/ESEM.2017.7
- [39] Yan Li, Chuan-Hoo Tan, and Hock-Hai Teo. 2012. Leadership characteristics and developers' motivation in open source software development. Information & Management 49, 5 (2012), 257–267.
- [40] Bin Lin, Fiorella Zampetti, Gabriele Bavota, Massimiliano Di Penta, and Michele Lanza. 2019. Pattern-based mining of opinions in Q&A websites. In Proceedings of the 41st International Conference on Software Engineering, ICSE 2019, Montreal, QC, Canada, May 25-31, 2019, Joanne M. Atlee, Tevfik Bultan. and Ion Whittle (Eds.). IEEE / ACM. 548-559. https://doi.org/10.1109/ICSE.2019.00066
- [41] Erik Linstead, Paul Rigor, Sushil Bajracharya, Cristina Lopes, and Pierre Baldi. 2007. Mining eclipse developer contributions via author-topic models. In Fourth International Workshop on Mining Software Repositories (MSR'07: ICSE Workshops 2007). IEEE, 30–30.
- [42] Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. 2016. Recurrent Neural Network for Text Classification with Multi-Task Learning. In Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016, Subbarao Kambhampati (Ed.). IJCAI/AAAI Press, 2873–2879. http://www.ijcai.org/Abstract/16/408
- [43] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. Introduction to information retrieval. Cambridge University Press. https://doi.org/10.1017/CBO9780511809071
- [44] Audris Mockus, Roy T Fielding, and James Herbsleb. 2000. A case study of open source software development: the Apache server. In Proceedings of the 22nd international conference on Software engineering. 263–272.
- [45] Mahdi Moqri, Xiaowei Mei, Liangfei Qiu, and Subhajyoti Bandyopadhyay. 2018. Effect of "following" on contributions to open source communities. Journal of Management Information Systems 35, 4 (2018), 1188–1217.
- [46] Kumiyo Nakakoji, Yasuhiro Yamamoto, Yoshiyuki Nishinaka, Kouichi Kishida, and Yunwen Ye. 2002. Evolution patterns of open-source software systems and communities. In Proceedings of the international workshop on Principles of software evolution. 76–85.
- [47] Derrick J Neufeld and Haoyue Gu. 2019. Leadership Emergence and Impact on Open Source Software Project Success: A Comparative Case Study. In Academy of Management Proceedings, Vol. 2019. Academy of Management Briarcliff Manor, NY 10510, 11698.
- [48] Peter G Northouse. 2021. Leadership: Theory and practice. Sage publications.

- [49] Sebastiano Panichella, Andrea Di Sorbo, Emitza Guzman, Corrado Aaron Visaggio, Gerardo Canfora, and Harald C. Gall. 2015. How can i improve my app? Classifying user reviews for software maintenance and evolution. In 2015 IEEE International Conference on Software Maintenance and Evolution, ICSME 2015, Bremen, Germany, September 29 - October 1, 2015, Rainer Koschke, Jens Krinke, and Martin P. Robillard (Eds.). IEEE Computer Society, 281–290. https://doi.org/10.1109/ICSM.2015.7332474
- [50] Craig L Pearce and Jay A Conger. 2002. Shared leadership: Reframing the hows and whys of leadership. Sage Publications.
- [51] James Pustejovsky and Amber Stubbs. 2012. Natural Language Annotation for Machine Learning: A guide to corpus-building for applications. "O'Reilly Media, Inc.".
- [52] Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python Natural Language Processing Toolkit for Many Human Languages. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations.
- [53] Jeffrey A Roberts, Il-Horn Hann, and Sandra A Slaughter. 2006. Understanding the motivations, participation, and performance of open source software developers: A longitudinal study of the Apache projects. Management science 52, 7 (2006), 984–999.
- [54] Gordon Rugg and Peter McGeorge. 2005. The sorting techniques: a tutorial paper on card sorts, picture sorts and item sorts. Expert Syst. J. Knowl. Eng. 22, 3 (2005), 94–107. https://doi.org/10.1111/j.1468-0394.2005.00300.x
- [55] scikit learn. 2021. Scikit-learn. https://github.com/scikit-learn/scikit-learn.
- [56] sfconservancy. 2021. Software Freedom Conservancy: Member Project Services. https://sfconservancy.org/projects/services/.
- [57] Dou Shen, Jian-Tao Sun, Qiang Yang, and Zheng Chen. 2006. Text classification improved through multigram models. In Proceedings of the 2006 ACM CIKM International Conference on Information and Knowledge Management, Arlington, Virginia, USA, November 6-11, 2006, Philip S. Yu, Vassilis J. Tsotras, Edward A. Fox, and Bing Liu (Eds.). ACM, 672–681. https://doi.org/10.1145/1183614.1183710
- [58] Lin Shi, Celia Chen, Qing Wang, Shoubin Li, and Barry W. Boehm. 2017. Understanding feature requests by leveraging fuzzy method and linguistic analysis. In Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering, ASE 2017, Urbana, IL, USA, October 30 November 03, 2017, Grigore Rosu, Massimiliano Di Penta, and Tien N. Nguyen (Eds.). IEEE Computer Society, 440–450. https://doi.org/10.1109/ASE.2017.8115656
- [59] Igor Steinmacher, Igor Scaliante Wiese, Ana Paula Chaves, and Marco Aurélio Gerosa. 2012. Newcomers Withdrawal in Open Source Software Projects: Analysis of Hadoop Common Project. In 2012 Brazilian Symposium on Collaborative Systems, Sao Paulo, Brazil, October 15-18, 2012. IEEE Computer Society, 65-74. https://doi.org/10.1109/SBSC.2012.16
- [60] Xin Tan and Minghui Zhou. 2019. How to Communicate when Submitting Patches: An Empirical Study of the Linux Kernel. Proc. ACM Hum. Comput. Interact. 3, CSCW (2019), 108:1–108:26.
- [61] Xin Tan, Minghui Zhou, and Brian Fitzgerald. 2020. Scaling open source communities: an empirical study of the Linux kernel. In ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June - 19 July, 2020. 1222–1234. https://doi.org/10.1145/3377811.3380920
- [62] Virginia J Vanderslice. 1988. Separating leadership from leaders: An assessment of the effect of leader and follower roles in organizations. *Human relations* 41, 9 (1988), 677–696.
- [63] Terry Winograd, Fernando Flores, and Fernando F Flores. 1986. Understanding computers and cognition: A new foundation for design. Intellect Books.
- [64] Youngjin Yoo and Maryam Alavi. 2004. Emergent leadership in virtual teams: what do emergent leaders do? *Information and organization* 14, 1 (2004), 27–58.
- [65] Gary Yukl. 1981. Leadership in Organizations, 9/e. Pearson Education India.
- [66] Gary Yukl. 2012. Effective leadership behavior: What we know and what questions need more attention. Academy of Management perspectives 26, 4 (2012), 66–85.
- [67] Gary Yukl, Angela Gordon, and Tom Taber. 2002. A hierarchical taxonomy of leadership behavior: Integrating a half century of behavior research. Journal of leadership & organizational studies 9, 1 (2002), 15–32.
- [68] Yuxia Zhang, Minghui Zhou, Klaas-Jan Stol, Jianyu Wu, and Zhi Jin. 2020. How do companies collaborate in open source ecosystems?: an empirical study of OpenStack. In ICSE '20: 42nd International Conference on Software Engineering, Seoul, South Korea, 27 June 19 July, 2020. 1196–1208.
- [69] Yutong Zhao, Lu Xiao, Pouria Babvey, Lei Sun, Sunny Wong, Angel A. Martinez, and Xiao Wang. 2020. Automatically identifying performance issue reports with heuristic linguistic patterns. In ESEC/FSE '20: 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Virtual Event, USA, November 8-13, 2020, Prem Devanbu, Myra B. Cohen, and Thomas Zimmermann (Eds.). ACM, 964-975. https://doi.org/10.1145/3368089.3409674
- [70] Minghui Zhou and Audris Mockus. 2012. What make long term contributors: Willingness and opportunity in OSS community. In 34th International Conference on Software Engineering, ICSE 2012, June 2-9, 2012, Zurich, Switzerland, Martin Glinz, Gail C. Murphy, and Mauro Pezzè (Eds.). IEEE Computer Society, 518–528. https://doi.org/10.1109/ICSE.2012.6227164
- [71] Haiyi Zhu, Robert Kraut, and Aniket Kittur. 2012. Effectiveness of shared leadership in online communities. In Proceedings of the ACM 2012 conference on Computer Supported Cooperative Work. 407–416.