# Discovering the Functions of Language in Online Forums

**Youmna Ismaeil, Oana Balalau, Paramita Mirza**
Max Planck Institute for Informatics
{yismaeil, obalalau, paramita}@mpi-inf.mpg.de

## Abstract

In this work, we revisit the functions of language proposed by linguist Roman Jakobson and we highlight their potential in analyzing online forum conversations. We investigate the relation between functions and other properties of comments, such as controversiality. We propose and evaluate a semi-supervised framework for predicting the functions of Reddit comments. To accommodate further research, we release a corpus of 165K comments annotated with their functions of language.

## 1 Introduction

Understanding human conversations has long been an active area of research and has become even more important with the pervasiveness of intelligent assistants in our daily life. A vast amount of work has been dedicated to *speech act* (also referred to as *dialogue act* or *discourse act*) categorization for the purpose of characterizing the discourse of conversations or discussions. Speech acts focus on the addresser's intent in using language and were first introduced by Austin (1975). One of the most influential subsequent work by Searle (1976) focused on the addresser's intent in using language and proposed five categories for speech acts: *representatives*, *directives*, *commissives*, *expressives*, and *declarations*.

With the rise of the internet and online communication, recent works focus on utilizing dialogue acts for analyzing emails, online forums and live chats (Zhang et al., 2017; Joty and Hoque, 2016; Jeong et al., 2009; Forsyth, 2007; Wu et al., 2002).

However, even though they employed sets of dialogue acts based on the Dialogue Act Markup in Several Layers (DAMSL) scheme (Core and Allen, 1997), each work proposed different subsets to annotate the data with, tailored for each specific purpose. Zhang et al. (2017) proposed 9
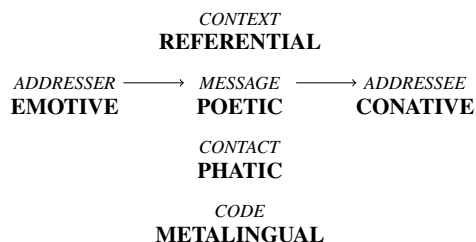


Figure 1: Functions of language (in bold), each focuses on different elements of communication (in italic).

| X: | I am so happy, our paper got accepted! | Emotive, Referential |
| Y: | Seriously?? Congratz. | Emotive, Phatic |
| X: | Well, the pen is mightier than the sword. | Poetic |
| Y: | What do you mean? | Metalingual |
| X: | It's a proverb, meaning to persuade reviewers with words instead of force. | Metalingual |
| Y: | I see. | Phatic |
|  | Could you send me your paper? | Conative |
| X: | Sure, it's here: <link> | Referential |
| Y: | Thanks! | Phatic |

Table 1: An example of a discussion where all functions of language are present.

speech acts for characterizing comments in Reddit; Joty and Hoque (2016) utilized 5 coarser classes from 12 acts used in Jeong et al. (2009); while Forsyth (2007) and Wu et al. (2002) defined 15 dialogue acts based on Stolcke et al. (2000). The lack of formalism and the diversity of taxonomies make it difficult to compare different annotated datasets. It is also not clear if the proposed acts cover all kinds of utterances in various conversation types. For instance, Zhang et al. (2017) labelled comments expressing disgust or anger as *negative reaction*, however, the counterpart *positive reaction* is not available as a label. Meanwhile, Joty and Hoque (2016) acknowledge only certain positive reactions labelled as *politeness*.

In this work, we address these issues by adopt-

ing the theory of *language functions* proposed by Jakobson (1960). One key difference between speech acts and language functions is that the former describes the meaning of utterances, while the latter also explains *why* utterances have different meanings, through the dynamic relationship between the elements of communication and their functions. Hence, we argue that the functions of language are a more comprehensive framework for capturing the discourse of human communication. Jakobson's model distinguishes six elements that are *necessary* for communication to occur: a **message** sent by an **addresser** to an **addressee** requires *(i)* a certain **context** to be understood, *(ii)* a common **code**, i.e., common definitions understood by both addresser and addressee, as well as *(iii)* a **contact**, i.e., a physical and psychological connection enabling both addresser and addressee to stay in the communication. Figure 1 illustrates the communicative functions of language in relation to those elements. In Table 1, we present an anecdotal example where all functions are expressed in the conversation.

There is a limited body of work exploring this scheme for content analysis. Bonini and Sellas (2014) use the functions of language to study the behavior of public radio channels on social media. Morrison and Martens (2018) incorporate the phatic function in a dialog system that would follow social norms. We believe we are the first to investigate Jakobson's functions of language for characterizing online forum discussions.

## 2 Jakobson's Functions of Language

We expand the original definitions of language functions (Jakobson, 1960) with further clarifications from the literature, as well as examples of how each function can be used to characterize messages on online forums.

**Referential.** The referential function, which is the most frequent one in communication, is marked by a reference to the context of the discussion, which can be a situation, a person, or an object. The message is used to transmit information and the words most often carry literal definitions (denotative). Instances of referential messages include observations, opinions, and factual information.

*Examples:* factual information (*"Trump won the election"*), opinions (*"He has a shot"*).

**Poetic.** The poetic function appears when the focus is on the message, marked in conversations by the use of figurative language. Figurative language concerns the use of metaphor, alliteration, onomatopoeia, idioms, irony and oxymorons, among others. Particularly in online forums, users often use slang expressions, which can be considered as poetic as well.

*Examples:* slang (*"Thanks Obama[1]"*), onomatopoeia (*"ding ding ding"*).

**Emotive.** The emotive function reflects the attitude or mood of the addresser towards the information being communicated. The message can be perceived as conveying emotion, such as anger, anticipation, joy and sadness. Emotive messages focus more attention on the addresser and less on the information being sent. Despite the absence of emotional tone and nonverbal cues, people can distinguish emotions in a text-based communication (Hancock et al., 2007).

*Examples:* emotions are often expressed using emojis or slang such as *"lol"* or *"omg"*, as well as words bearing strong sentiment (*"what a horrible human being"*).

**Conative.** The conative function is marked by a focus on the addressee. A conative message would make the addressee react by performing a verbal act (e.g., answering a question), a psychological act (e.g., changing a conviction), or a physical act (e.g., closing a door). More precisely, messages have a conative function if they represent orders, demands, advice, or wishes, among others.

*Examples:* demands (*"link please"*, *"Vote for Bernie Sanders"*), warnings (*"don't count on it"*).

**Phatic.** Sometimes referred to as back-channel or small talk, the phatic function serves the purpose of preserving the physical and psychological contact between speakers. The physical contact is related to the physical environment in which the conversation takes place and in the case of online forums, this will be a reference to the platform, e.g., *"happy cake day![2]"*. The psychological contact refers to the personal relation between speakers and the involvement in the conversation.

*Examples:* involvement in the conversation (*"I see"*), agreement and disagreement between

---

[1] https://www.urbandictionary.com/define.php?term=Thanks%20Obama

[2] Reddit "Cake Day" is the yearly anniversary of when a user signed up on Reddit.

speakers (*"good point"*, *"I don't think so"*).

**Metalingual.** The metalingual function corresponds to clarifications regarding the concepts used in the conversation, which can be related to the language used (as the common code) or the system/environment where the communication takes place. The metalingual function is often indicated by linguistic cues such as *"what is a"* or *"what do you mean by"*. The metalingual function appears when we need definitions, as well as ambiguity resolution. *Examples:* clarifying the vocabulary *"what is a noob?"*, or more general concepts *"what does the Supreme Court do?"*.

**Relations between Functions.** Messages will generally have more than one function of language. Jakobson (1960) highlights the relation between **referential** and **poetic** functions. The author argues that a poetic message will make referential information ambiguous, however, it will not completely discard it. Klinkenberg (2000), on the other hand, justifies the relation between **conative** and **referential** functions. The transfer of information between the addresser and the addressee might determine a change in the behavior of the addressee.

## 3 Functions of Language on Reddit

### 3.1 Dataset

In order to have a diverse tone of comments in long discussions, we consider the *Politics*[3] subreddit, a popular forum for political U.S. news. We retrieved $10.6M$ comments on the *Politics* subreddit for the year 2016, the year of the presidential election, using the Reddit API.

In this work, we focus on *short comments*, as they are challenging for existing automatic content analysis tools such as topic models. However, they often carry clear language functions on their own and can be easily distinguished by humans. We consider short comments to be the ones that consist of at most two syntactic phrases or chunks[4], e.g. "I see your point" (NP-VP), obtaining $165K$ comments. After removing punctuation and converting the text to lowercase, we have a final dataset of $4,482$ distinct utterances, which we will refer to as **messages**. Each message might represent several comments and be used in different contexts. Our intuition for removing punctua-

---

https://www.reddit.com/r/politics/
[4]We used OSU Twitter NLP by Ritter et al. (2011).

tion and uppercase is that the additional meaning can be added using simple rules. For example, an exclamation mark or text in all caps may suggest surprise or anger.

**Manual Annotation.** We set aside 920 (420 most frequent and 500 randomly selected) messages from the $4,482$ messages to be manually annotated. Each message was annotated by three human annotators, who are trained with the descriptions and examples of functions of language as have been explained in Section 2. We observed that almost all messages strongly express, and hence, annotated with at most two language functions, as we focus only on short comments. When a message is very ambiguous, the annotators were encouraged to give the label *unclear*. A message receives as a final label a function of language $f$ if that function is assigned by at least two annotators. The Krippendorff's alpha agreement score among annotators is $0.565$. We remark that the agreement score which is comparable with results reported for speech act labeling on Reddit (cf. Table 1 in Zhang et al. (2017)). Out of the 920 labelled messages, the annotators disagree on 67 messages, i.e. no label is voted more than once, which we exclude from our final dataset. We also removed 10 messages that are consistently labeled *unclear* by three annotators, leaving 843 labeled messages in our final dataset used for analysis, and later for experiments with automatic methods. The final label distribution is as follows: 352 referential, 288 phatic, 147 emotive, 104 poetic, 71 conative and 16 metalingual.

### 3.2 Analysis

We now analyze the properties of the annotated messages in relation to the functions of language. For each distinct message, we first retrieve the initial comments containing it. For example, the text *"thank you"* appears in $2,292$ comments, with different letter cases and punctuation. A comment has several properties, including *author*, *controversiality* and *parent comment*. A comment receives the tag *controversial* when it has a significant amount of votes, and these votes are roughly equally split between upvotes and downvotes. The parent comment is the comment to which the current comment is replying. From the parent-child relation of comments, we infer the *number of replies* of a given comment.

**Controversiality.** We first investigate which language functions commonly follow controversial comments. Our intuition is that emotive (e.g., expressing surprise or anger) and conative messages (e.g., asking users to behave or to provide evidence) will be written frequently in response to controversial topics or controversial users. The percentage of controversial parent comments per function shown in Table 2 confirms our intuition that conative comments are used to reply to controversial content more often than the other comments. Emotive messages are also written more frequently in reply to controversies, as well as referential messages. For the latter, the user may bring more information about the topic, either to approve or disprove the parent comment.

| Function | % controversial parents | % receive reply |
|---|---|---|
| *referential* | **15.87%** | **36.74%** |
| *poetic* | 12.86% | 18.17% |
| *emotive* | **16.59%** | 15.78% |
| *conative* | **20.14%** | **51.03%** |
| *phatic* | 11.43% | 14.61% |
| *metalingual* | 13.41% | **26.51%** |

Table 2: Analysis per language function.

**Replies.** We also examine which language functions are often followed by at least one reply comment, shown in Table 2. The findings corroborate with the definitions of the language functions, as conative comments, which put the focus on the addressee, often receive replies. Referential and metalingual are the other functions that often receive replies, since they bring more information and naturally prolong the discussion. Meanwhile, the opposite is observed for emotive and phatic comments. Emotive comments such as *"lol"*, *"haha"* representing joy usually require no follow-up in verbal conversations. Poetic comments receive also relatively few replies. Drew and Holt (1998) found that figures of speech are used as transitions between topics or as ending remarks on a topic. This phenomenon may also be present on Reddit. On the subreddit Politics, users initiate discussions via an article or video, hence, phatic messages mostly express involvement (e.g., *"I see"*) or (dis)agreement (e.g., *"good point"*), which require no replies and serve the role of ending the conversation.

**Applications.** Given the definitions of language functions and the previous observations, we illustrate two use cases for the analysis of language functions in online forums. First, they can be used in combination with other features for the automatic classification of comments or threads as *controversial*, in the absence of sufficient voting activity. Controversial messages require the immediate attention of moderators as they might contain hate speech or false information. Secondly, understanding conversational patterns related to language functions (e.g., a conative message asks for a referential reply, while an emotive message calls for a thoughtful and empathetic response) are beneficial for building smarter chatbots.

### 3.3 Semi-Supervised Inference of Labels

Annotating posts on social media with functions of language, or any semantic or discourse label in general, is a time-consuming and labor-intensive task. To overcome this challenge, we investigate the utility of a graph-based semi-supervised label propagation framework with the Modified Adsorption (MAD) algorithm (Talukdar and Pereira, 2010), which makes predictions by taking into consideration both labeled and unlabeled data. MAD was shown to perform the best when compared with other semi-supervised frameworks, such as the Label Propagation (LP-ZGL) algorithm and the Adsorption algorithm (Talukdar and Pereira, 2010). MAD computes a soft assignment of labels of the nodes in a graph, allowing multi-label classification. Graph-based semi-supervised learning is widely used by the NLP community, particularly for tasks where acquiring annotated data is expensive, such as semantic parsing (Das and Smith, 2011).

We construct a graph in the following way: the set of $4,482$ messages (e.g., *"thank you"*, *"thanks"*) is considered as the set of nodes and an edge is added between a message and its $k$-nearest neighbor in the embedding space. Cosine similarity between two messages *embeddings*[5] is assigned as the weight of the edge. We experiment with different values for $k$ and we find that the algorithm performs best for $k = 4$. For evaluation, we used 5-folds cross-validation on the annotated dataset (843 messages) and we report precision, recall, and $F1$-score per function of language in Table 4. Note that we exclude metalingual comments, as they were not sufficient for propagating

---

[5] Pre-trained embeddings from Google's Universal Sentence Encoder (Cer et al., 2018).

| referential | phatic | emotive | poetic | conative |
|---|---|---|---|---|
| *she has a credible claim* | *absolutely agree* | *shameful* | *feeltheburn* | *mind explaining why* |
| *what time was that* | *I'm sorry* | *barely even human* | *inverted triple bern* | *keep fooling yourself* |
| *it's all marketing* | *I upvoted you* | *epic simply epic* | *duality of man* | *dude relax* |

Table 3: Example of predictions of our semi-supervised approach.

| Function | Precision | Recall | F1 |
|---|---|---|---|
| *referential* | 0.893 | 0.888 | 0.891 |
| *phatic* | 0.905 | 0.889 | 0.897 |
| *emotive* | 0.868 | 0.838 | 0.853 |
| *poetic* | 0.680 | 0.798 | 0.734 |
| *conative* | 0.822 | 0.890 | 0.855 |
| **Average** | **0.834** | **0.861** | **0.847** |

Table 4: Precision, recall and $F1$ score per function.

the labels. However, we hypothesize that one sure way to identify metalingual function is by looking at the presence of linguistic cues such as *"what is a"* or *"what do you mean by"*. Apart from metalingual, our approach performs well on all functions yielding around $0.734 - 0.897$ $F1$-scores. The lowest $F1$-score is coined by poetic functions due to the high variation of figurative language. Users can make comparisons, metaphors or puns, among others, making the task at hand challenging and deserving of a focused effort.

**Qualitative Analysis.** We show in Table 3 examples of the predictions for the unlabelled dataset made by our approach. Even for the difficult task of identifying figurative language, MAD can make good predictions.

## 4 Conclusion

This paper revisits the functions of language introduced by Jakobson (1960) and investigates their potential in analyzing online forum conversations, specifically political discussions on Reddit. We highlight interesting relations between comments, their properties, and the language functions they express. In addition, we present a graph-based semi-supervised approach for automatic annotation of language functions.

**Dataset.** For further research in this area, we release a corpus[6] of $165K$ comment IDs labeled with their functions of language.

---

[6]https://github.com/nyxpho/jakobson

## References

John Langshaw Austin. 1975. *How to do things with words*. Oxford university press.

Tiziano Bonini and Toni Sellas. 2014. Twitter as a public service medium? A content analysis of the Twitter use made by Radio RAI and RNE. *Communication & Society*, 27(2):125–146.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.

Mark G Core and James Allen. 1997. Coding dialogs with the DAMSL annotation scheme. In *AAAI fall symposium on communicative action in humans and machines*, volume 56. Boston, MA.

Dipanjan Das and Noah A. Smith. 2011. Semi-supervised frame-semantic parsing for unknown predicates. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 1435–1444. Association for Computational Linguistics.

Paul Drew and Elizabeth Holt. 1998. Figures of speech: Figurative expressions and the management of topic transition in conversation. *Language in Society*, 27(04):495–522.

Eric N. Forsyth. 2007. Improving automated lexical and discourse analysis of online chat dialog. In *Masters thesis, Naval Postgraduate School*.

Jeffrey T. Hancock, Christopher Landrigan, and Courtney Silver. 2007. Expressing emotion in text-based communication. page 929.

Roman Jakobson. 1960. Linguistics and poetics. In *Style in language*, pages 350–377. MA: MIT Press.

Minwoo Jeong, Chin-Yew Lin, and Gary Geunbae Lee. 2009. Semi-supervised speech act recognition in emails and forums. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing*, pages 1250–1259, Singapore. Association for Computational Linguistics.

Shafiq Joty and Enamul Hoque. 2016. Speech Act Modeling of Written Asynchronous Conversations with Task-Specific Embeddings and Conditional Structured Models. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1746–1756.

Jean-Marie Klinkenberg. 2000. *Précis de sémiotique générale*. Le Seuil.

Hannah Morrison and Chris Martens. 2018. "How Was Your Weekend?" A Generative Model of Phatic Conversation. pages 74–79.

Alan Ritter, Sam Clark, Mausam, and Oren Etzioni. 2011. Named entity recognition in tweets: An experimental study. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 1524–1534. Association for Computational Linguistics.

John R Searle. 1976. A classification of illocutionary acts. *Language in society*, 5(1):1–23.

Andreas Stolcke, Noah Coccaro, Rebecca Bates, Paul Taylor, Carol Van Ess-Dykema, Klaus Ries, Elizabeth Shriberg, Daniel Jurafsky, Rachel Martin, and Marie Meteer. 2000. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational Linguistics*, 26(3).

Partha Pratim Talukdar and Fernando Pereira. 2010. Experiments in graph-based semi-supervised learning methods for class-instance acquisition. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1473–1481. Association for Computational Linguistics.

Tianhao Wu, Faisal M. Khan, Todd A. Fisher, Lori A. Shuler, and William M. Pottenger. 2002. Posting act tagging using transformation-based learning. In *Foundations of Data Mining and knowledge Discovery*.

Amy Zhang, Bryan Culbertson, and Praveen Paritosh. 2017. Characterizing online discussion using coarse discourse sequences. In *International AAAI Conference on Web and Social Media (ICWSM'17)*.