

ChatGPT vs. Crowdsourcing vs. Experts: Annotating Open-Domain Conversations with Speech Functions

Anonymous ACL submission

Abstract

This paper deals with the task of annotating open-domain conversations with speech functions. We propose a method for annotating dialogs following the topic-oriented, multi-layered taxonomy of speech functions with the use of hierarchical guidelines. These guidelines comprise simple questions about the topic and speaker change, sentence types, pragmatic aspects of the utterance, and examples that aid untrained annotators in understanding the taxonomy. We compare the results of dialog annotation performed by experts, crowdsourcing workers, and ChatGPT. To improve the performance of ChatGPT, several experiments utilising different prompt engineering techniques were conducted. We demonstrate that in some cases large language models can achieve human-like performance following a multi-step tree-like annotation pipeline on complex discourse annotation, which is usually challenging and costly in terms of time and money when performed by humans.

1 Introduction

Analyzing discourse structure as a method of an abstract dialog representation is used in various NLP tasks: dialog management (Liang et al., 2020; Galitsky and Ilvovsky, 2017), dialog generation (Yang et al., 2022; Gu et al., 2021), dialog summarization (Chen et al., 2021), emotion recognition (Shou et al., 2022), etc. Discourse structure can be defined in a variety of ways, but in general, it is considered to be an interconnected system of linguistic features such as a topic, pragmatics, and semantics. One of the main goals of discourse analysis is to describe pragmatics of communicative actions performed by speakers, i.e., characterise interlocutors' intentions at a certain moment of dialog (Coulthard, 2014).

There are two basic theoretical approaches to this kind of a dialog representation: theory of dialog acts (DA theory) (Core and Allen, 1997) and

discourse annotation in the style of Segmented Discourse Representation theory (SDRT) (Lascarides and Asher, 2007) inheriting principles of Rhetorical Structures theory (RST) (Mann and Thompson, 1988). According to the SDRT style, firstly, a relation between two elementary discourse units (EDUs) needs to be defined and then characterized with a discourse class (for instance, Question-Answer, Clarification, etc.). While SDRT represents a dialog structure as a graph (Asher et al., 2016; Li et al., 2020), most of DA theory interpretations such as DAMSL (Allen and Core, 1997), SWBD-DAMSL (Jurafsky, 1997), MIDAS (Yu and Yu, 2019) describe it sequentially giving pragmatic characteristics to each EDU. In addition, most classes used in DA taxonomies do not represent pragmatic purposes but rather focus on semantics or grammar form of utterances within a dialog, using tags such as 'yes/no question', 'statement', 'positive answer'.

To represent the discourse structure of dialogs in a more advanced way, H. Bunt suggested Dialogue Annotation Markup Language (DiAML), a taxonomy including nine functional dimensions and 49 specific classes (Bunt et al., 2010, 2012). Even though DiAML is claimed to be an ISO standard for DA annotation, it is challenging to apply it to real-world problems for several reasons. First, DiAML supports multi-label annotation, i.e., several classes can be assigned to one EDU, which complicates automatic classification. Moreover, there is not enough labelled data to experiment with the taxonomy. One more taxonomy designed to represent a conversational structure on several levels is Dependency Dialogue Acts (DDA) (Cai et al., 2023). A combination of dialog acts and rhetorical relations in the SDRT style showed a potential of applying multi-layered and multi-dimensional approaches for analyzing discourse structure within conversations. However, because there is no annotated data with this taxonomy, it is not clear

whether it is applicable to automated tasks.

In this paper, we provide an extensive research on a topic-oriented, multi-dimensional and hierarchical taxonomy of speech functions introduced by D.Eggins and S.Slade (Eggins and Slade, 2004) as an alternative for abstract dialog representation. In addition, the taxonomy includes classes being too close to each other in terms of pragmatics that are not featured in other annotation schemes. We analysed how experts with backgrounds in linguistics and non-professional crowdsourcing workers can manage such a challenging classification. Furthermore, this paper explores the potential of using LLMs, ChatGPT in particular, for discourse structure annotation.

2 Taxonomy of Speech Functions

Taxonomy of speech functions introduced by S.Eggins and D.Slade is multi-layer, multi-dimensional and hierarchical, which allows to analyze dialog structure in a very consistent way (Eggins and Slade, 2004). This annotation scheme originally included 45 speech functions, but we reduced them to 32 labels. Unlike other multidimensional schemes, the taxonomy of speech functions supports single-label annotation. All functional dimensions are embedded in a single label, so it is easier to comprehend the discourse structure. While inheriting the principle of assigning one label to a specific EDU from DA theory, speech functions taxonomy also considers relationships between utterances following the SDRT style. The tag of a current label is determined in connection with the previous one, so it is important to take into account the utterances' previous context when assigning the correct label.

2.1 Functional Dimensions

S.Eggins and D.Slade's tag set consists of speech functions representing five different functional dimensions (Eggins and Slade, 2004). The dimensions are embedded in speech functions but distributed unevenly between tags: from two to five dimensions can be featured in one speech functions (see Figure 1).

- **Turn Management** denotes a speaker change at the current moment of conversation, which is represented in all speech functions except Opening moves defining a new topic. At this functional level, a *Sustain* label indicates that a speaker continues the conversation, whereas



Figure 1: Example of Speech Function Structure

a *React* label implies that a speaker changes or the same speaker reacts to previous utterances of an interlocutor.

- **Topic Organisation** level denotes the beginning of the dialog or a topic shift, as well as the development of a topic. *Open moves* are used to indicate the start of a dialog or a new topic. Sustain moves include a *Continue* label that shows a progression of the current topic. The *Respond* label is embedded in Reaction moves to define classes that are more likely to end the dialog and do not contribute to the topic's development. Such classes encounter more passive responses in the form of answers, back channelling, and continuation of previous narration. *Rejoinder* labels, on the other hand, define more active development of the conversation topic that has an impact on the dialog flow.
- **Feedback** level is used to more accurately characterise moves of Reaction. *Confront* and *Support* labels indicate whether a speaker is challenging or supporting an interlocutor.
- **Communicative Acts** are used to specify groups of pragmatic purposes that are very close in terms of interpretation and united by the same functionality within conversation. For instance, three speech functions are used to prolong the narration produced by one speaker, but they're performed in three different ways in the dialog.
- **Pragmatic Purposes** level is the last one in hierarchical taxonomy of speech functions specifying speakers' intentions. This layer of annotation is considered to be the most challenging for annotation as those are very close classes in terms of pragmatics.

It is important to note that speech function taxonomy is flexible enough as there is a potential of

enriching the scheme with additional annotation layers indicating different features of utterances.

2.2 Levels of segmentation

(Bunt et al., 2012) defined EDUs as ‘functional segments’ and claimed that a speaker can perform several functions within one utterance. So, the boundaries of elementary discourse units are determined by communicative actions’ functions depending on a chosen taxonomy. As a taxonomy of speech functions is topic-oriented, the first level of segmentation is determined by a topic shift in the dialog. Utterances united by a specific topic compound a **discourse pattern** (see Figure 2). Every discourse pattern is segmented into **turns** defined by a speaker change that can include one or several **utterances**. In most cases, utterance boundaries coincide with sentence boundaries, but some speech functions demand a finer division or a combination of several sentences. Every utterance is actually a functional segment characterized by a particular speech function.

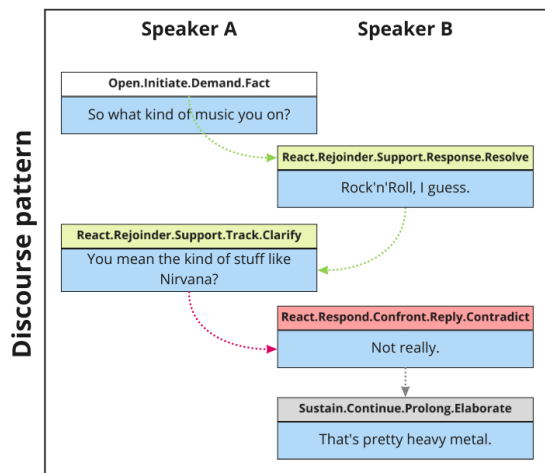


Figure 2: Discourse Pattern according to Speech Function Taxonomy

3 Human Annotation using Speech Function Taxonomy

The annotation of discourse structures or dialog acts is not a simple task as it requires either linguistic knowledge or trained workers (Yung et al., 2019). Additionally, understanding the speaker’s intentions in utterances can vary among individuals, further complicating the task. In this section, we compare the results of speech function annotation completed by experts with professional back-

grounds in linguistics and crowdsourcing assessors. To evaluate the agreement between the experts and between the assessors, we use Fleiss’ kappa that is an extension of Scott’s pi (π) for two coders. Fleiss’ kappa can deal with any number of annotators, where every item is not necessarily annotated by each annotator. It is the most commonly used method to evaluate taxonomy reliability in tasks related to discourse analysis. However, this method has the limitation of not considering the common mistakes of annotators. Therefore, we measured not only inter-annotator agreement but also three the most common metrics for multi-class classification tasks with imbalanced data — Macro F1, Weighted Precision and Weighted Recall, by comparing the workers’ annotations to the results of experts.

3.1 Tree-like Design of Annotation Instruction

To facilitate annotation, we designed a tree-like scheme comprised of a series of questions and their corresponding answer options that reproduces logic of a hierarchy of speech functions taxonomy. Due to multidimensional structure of speech functions, the path to each final label can be represented as a series of straightforward questions in form of instructions. This tree-like structure was used by both experts and annotators during annotation process.

3.2 Crowdsourcing Process

For crowdsourcing we used a platform for data annotation enabling project management and review cycles. When carrying out complex discourse annotation, the following two main problems are encountered:

- pragmatic classes are difficult to differentiate for annotators without a strong linguistic background;
- an issue of unreliable annotators who prioritize speed over accuracy.

To address the first issue, we used a tree-like design of guidelines rather than selecting one of 32 different speech functions. At each stage of annotation, a crowdsourcing worker answers a simple question with 2-4 possible options. An instruction with explanations and examples is attached to each question. Having answered all the questions in the chain, the annotator reaches the final label.

As for the second problem, we developed several mechanisms for tracking the quality of answers,

including (1) detecting the fast answers that are selected without reading instructions, (2) checking answer consistency across related questions, and (3) using trained classifiers to detect answers that do not match the expected annotation.

In addition, we developed a multi-level process of verification and examination to increase the quality of dialog annotation. The first stage involves both training and the exam process on a single dialog, with hints shown to crowdsourcing workers if they answer incorrectly. Workers who fail to achieve the appropriate quality can retry one more time. Those who pass the examination are selected for the main annotation pool. Each dialog is evaluated based on custom validation rules and control questions. If the dialog fails validation, annotators cannot continue the annotation.

3.3 Crowdsourcing vs. Experts

As the source of dialog data, we used DailyDialog (Li et al., 2017), a hand-crafted dataset of multi-turn human conversations about daily life. First, we splitted the utterances into EDUs. Second, three experts with at least B.A. in Linguistics annotated 64 dialogs (1030 utterances). Third, the same data was annotated via crowdsourcing with three non-professional workers annotating each dialog. We evaluated the results for 16 high-level cut labels and the complete taxonomy to identify the weak points of the established hierarchical guidelines (see Appendix B for an overview of taxonomy). We also examined cases of voting, in which the majority of annotators agreed on a tag. The cut labels were labeled with high accuracy by crowdsourcing workers, while the annotation of full tags was more challenging for non-experts, as proven by all metrics. Macro F1 value indicates that improving the quality of annotating low-level classes is necessary (see Table 3a). Fleiss’ Kappa revealed that differentiating tags with similar pragmatics is difficult not only for untrained workers but also for experts. Nonetheless, the chosen taxonomy is quite reliable, as Fleiss’ Kappa for experts’ annotation is more than 0.6, standing for substantial agreement (see Figure 3).

The use of speech function taxonomy implies a noticeable class imbalance, with certain speech functions occurring more frequently than others (see confusion matrix 6a in Appendix A). Classes that have a limited number of examples are Rebound, Re-challenge, Refute, etc. Certain classes

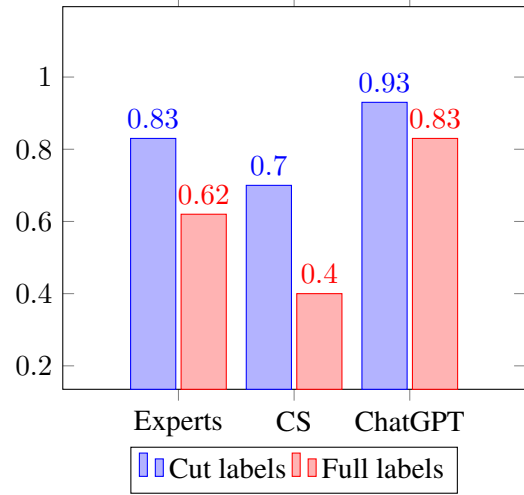


Figure 3: Inter-annotator Agreement (CS - crowdsourcing)

are well-defined and easily distinguishable, including Open.Attend, Register, Resolve, Clarify, and Open.Demand.Fact. However, the classes of Extend, Enhance, and Elaborate are challenging to distinguish accurately because they are very close in terms of pragmatics.

4 Large language models’ heyday

In the recent years, the paradigm of training and using NLP models has been undergoing significant changes. With the advance of Large Language Models (LLMs), the focus has shifted from the previously dominating “pre-train, fine-tune” procedure to “pre-train, prompt, and predict” (Liu et al., 2023), where an LLM is applied to downstream tasks directly. In this case, textual prompts are used to guide the models’ behaviour and achieve the desired output without additional fine-tuning. Scaling up LLMs to billions of parameters leads to significantly improved results in terms of few-shot and zero-shot prompting (Brown et al., 2020; Wei et al., 2021, *i.a*). However, as the objective of training most LLMs is not following the instructions but simply predicting the next token, they may fail to perform the task. One solution is fine-tuning LLMs using Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017) to align its behaviour in accordance with the trainers’ values and needs (Ouyang et al., 2022). An example of such model is ChatGPT (OpenAI, 2022) that has shown state-of-the-art or comparable performance on a number of NLP tasks in few-shot or zero-shot setting and provoked a tide of research

articles testing its capabilities in areas ranging from coding and bug-fixing (Tian et al., 2023; Kashefi and Mukerji, 2023; Sobania et al., 2023, *i.a*) to medical applications (Nori et al., 2023; Kung et al., 2023, *i.a*).

4.1 ChatGPT for annotation

There have been claims that ChatGPT outperforms crowdsourcing workers on a number of annotation tasks while being significantly cheaper. The tasks in question included annotation of relevance, stance, topics, and frames detection (Gillard et al., 2023); political affiliation classification of tweets (Törnberg, 2023); hate speech detection (Huang et al., 2023; Li et al., 2023; Zhu et al., 2023); sentiment analysis and bot detection (Zhu et al., 2023). In the above-listed works, the approach to obtaining the final label was straightforwardly simple. With one prompt containing textual instruction and a datapoint, the model had either to answer a question about the datapoint, assigning a label to it, or score the probability of the datapoint belonging to some class.

However, to the best of our knowledge, there have been no attempts to apply LLMs to more complex annotation tasks that deal with tens of labels and require multi-step reasoning. The purpose of this work is to test whether LLMs such as ChatGPT may achieve human-like performance on such annotation tasks.

5 Experiments

The annotation task in question required careful instruction preparation even for human annotators as opposed to simpler tasks such as sentiment classification, bot detection, etc. Thus, the process of creating the best prompt for an LLM is also a challenging and multi-step process. We conduct a number of experiments in order to find the best way to use ChatGPT for complex discourse annotation tasks. In all cases, the `system_message` we used while querying ChatGPT API was “You are a professional linguist annotator who has to perform a discourse annotation task”. The `user_message` varied for different experiments. See Figure 5 for an example of `user_message`.

To reduce the number of API calls and thus the time and the cost of the annotation, we also used automatic methods other than ChatGPT on some steps of the annotation. For example, in all our experiments we used Topic Shift Classifier to de-

tect the beginning of a new topic in a dialog. It is worth noting that ChatGPT did not perform well in this particular task. The Topic Shift Classifier was trained using the DeepPavlov (Burtsev et al., 2018) library utilizing a double sequence binary classifier model based on `roberta-large-mnli`, with two sequential utterances as input. The true labels indicate topic change in the utterances. The following hyper-parameters were used to train the model: learning rate = $2e-5$, optimizer = AdamW, input max length = 128. To successfully train the model, we used the early-stopping technique. The classifier was able to transfer the knowledge acquired during pre-training on `mnli` to the related problem of shift identification by using a pre-trained model (Kononov et al., 2020; Gulyaev et al., 2020). Due to the limitations in funding and a large number of experiments, they were carried out on 12 dialogs, 189 utterances only (approximately 1/5 of the final corpus).

5.1 Choosing the best annotation scheme

First, we compare three approaches to automatic discourse annotation using ChatGPT:

- Direct annotation – providing an full list of labels to choose from;
- Step-by-step scheme with intermediate labels;
- Complex tree-like scheme with intermediate labels and yes-no questions prevailing on each step.

5.1.1 Direct annotation scheme

The most straightforward approach is providing the final labels, their description and 2 examples for each to the model as they are. However, even at this step we chose to distinguish between 6 *Open* speech functions – the ones that begin the dialog or a new topic in the dialog – and 27 *React/Sustain* speech functions via a preliminary classification step. Here, the pipeline consists of two steps. See Figure 4a for an overview.

5.1.2 Step-by-step annotation scheme

Here, the annotation process was broken down into smaller steps. The pipeline consisted of 2-5 steps depending on the outcome of each step. In the end, the model once again had to choose between several final labels (from 4 to 12). See Figure 4b for an overview.

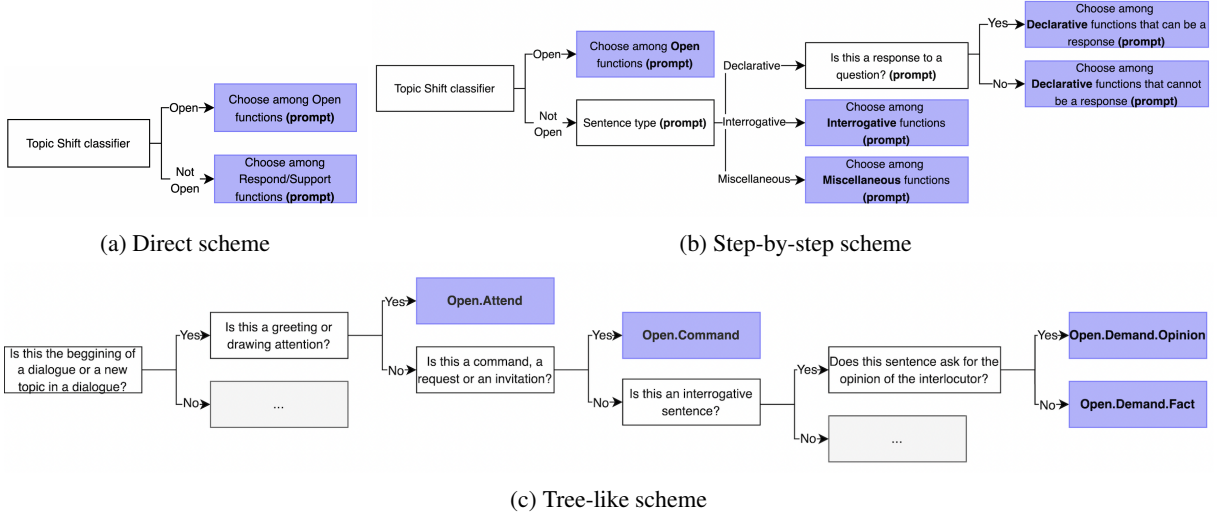


Figure 4: Experiment pipelines

```

TASK: This is part of the dialog is between 2
speakers. Answer QUESTION about CURRENT UTTERANCE.
You must analyze relations between CURRENT UTTERANCE
and PREVIOUS CONTEXT.

PREVIOUS CONTEXT:
speaker_1: Hey!
speaker_1: I heard you'd annotated a corpus of 1000
utterances in just an hour!
speaker_1: Is that true?
CURRENT UTTERANCE:
speaker_2: Well, technically, I made ChatGPT do that.

QUESTION: Can this utterance be an answer to the
previous speaker's question?
POSSIBLE ANSWERS: Yes, No
You must always select an option. Provide only one
answer without explanation.
ANSWER (Yes or No):

```

Figure 5: An example of user_message.

5.1.3 Tree-like annotation scheme

In this experiment, we used a complex tree-like annotation pipeline that was primarily designed to facilitate human crowdsourcing annotation process. As breaking the task of selecting one of many labels into smaller sub-tasks of a tree-like structure with simpler questions on each step is used to improve performance of humans on complex discourse annotation tasks (Scholman et al., 2016), we speculate that the same holds true for annotation via ChatGPT. Additionally, novel research suggests that making the model follow a number of tree-like structured prompts may greatly improve its performance (as applied to sudoku puzzles in (Yao et al., 2023)). The major difference from the Step-by-step annotation scheme is that the Tree-like annotation scheme favours prompts containing yes-no questions over prompts asking to select one option out of many. As a results, the scheme is much more

complex than the ones described before, with 2-12 steps to be completed before reaching the final label. However, the majority of questions are extremely simplified, guiding the model to the final label via a series of yes-no questions. For an example of how some final labels can be reached, see Figure 4c¹.

5.1.4 Results

Naturally, with more detailed schemes and simpler questions on each step, the model achieved better results. As Table 3a demonstrates, Macro F1 is significantly lower than Weighted Recall and Weighted Precision for complex schemes, Step-by-step and Tree-like annotation. The Speech Function annotation scheme is deemed to produce imbalanced data classes due to its nature – some classes are by definition more common and some are rare. Thus, the difference between higher Weighted Recall and Precision demonstrate that we were able to classify more common categories well as those categories have a greater influence on weighted metrics. On the opposite, as Macro F1 treats all classes equally regardless of their size, lower Macro F1 in all schemes shows that the model’s performance consistently deteriorates on smaller classes.

Even though Weighted Precision is higher for less complex Step-by-step scheme, we can say that with Tree-like scheme the model performed the task better as higher Macro F1 demonstrates that it was better at distinguishing between smaller classes.

¹The entire scheme is too large to be included into the article. However, the link to our GitHub where it is published will be provided after the blind review stage.

	Weighted Recall	Weighted Precision	Macro F1
Direct annotation	0.23	0.33	0.28
Step-by-step scheme	0.57	0.75	0.31
Tree-like scheme	0.62	0.67	0.43

Table 1: Evaluation of annotation by ChatGPT using different annotation methods (on a subset of dialogs)

5.2 Hyperparameter tuning

While examining the results of the annotation in Subsection 5.1, we noticed that in some cases the model was confused by the names of the classes it had to choose from on the last step of the annotation. For example, when asked to choose from labels Check, Confirm, Clarify, and Probe, the model tended to ignore the instruction that Check is only used to get the previous speaker to repeat something, and overuse this label (see Appendix B for detailed definitions of each label). When asked to provide an explanation of its choice, the model would produce explanations based on the semantics of the word Check, e.g. “The speaker wanted to check what the previous speaker thinks”. Thus, we decided to check if the performance improves if the final labels are masked, replacing the speech function name with a number and leaving the definitions and instructions intact.

We also experimented with model temperature (0.0, 0.5, 0.9), a hyperparameter that controls the randomness of the generated content. Higher temperature, meaning higher randomness and diversity, turned out to work best.

Another feature that we tested was a modification of zero-shot Chain-of-Thought prompting as described in (Kojima et al., 2022). Here, the model was asked to provide an answer in the following format: “Reasoning: (your reasoning). The final answer: (your final answer)”. However, in our case, generating reasoning and grounding the final answer in it did not improve the quality.

Finally, we experimented with the size of the context window (1, 3, 5), i.e., the number of previous utterances provided to the model. The longer context seems to confuse the model, as the windows of sizes 1 and 3 performed better. For a detailed

overview of metrics for all experiments, refer to Table 2.

Overall, there has been no significant difference in performance between the models with different hyperparameters. the best performing option turned out to be the model with temperature = 0.9, masked labels, context window = 1, and no reasoning.

5.3 Evaluation

Then, we used ChatGPT to annotate the same corpus that was previously annotated by experts and crowdsourcing workers (described in 3.3). For that, we selected the best set of hyperparameters (as shown in 5.2): tree-like scheme with masking, temperature = 0.9, and context window = 1. The results are shown in Table 3b.

As can be seen, ChatGPT performed well on a subset of 12 dialogs (see Table 2), but on the entire dataset, it performs noticeably worse for full and cut tags. We also tried to employ the voting method when utilizing ChatGPT, similar to what was done with crowdsourcing annotation, to enhance the reliability of the annotation. We ran the annotation pipeline three times, counted the votes and got the results that are also shown in Table 3b. As can be seen from the table, the implementation of voting had minimal impact on the results. ChatGPT consistently provided answers, as indicated by the Fleiss Kappa scores of 0.83 for full tags and 0.93 for cut tags, representing an almost perfect level of agreement and model consistency, despite temperature being set to 0.9 (meaning more diverse responses).

The lower quality of the annotation by ChatGPT compared to crowdsourcing can be explained by two main reasons (see Figure 6b in Appendix A). Firstly, distinguishing between close subclasses such as Extend/Enhance/Elaborate is challenging, even for humans, and it appears to be even more difficult for ChatGPT. Additionally, ChatGPT struggles with differentiating between Acknowledge/Affirm/Agree. Secondly, ChatGPT not only has difficulties in distinguishing among subclasses, but it also frequently confuses Resolve (detailed answer) with Replies (positive and negative answers). Furthermore, it often misclassifies Extend as Affirm or Agree. In general, the difference in metrics between 12 and 64 dialogs can be explained by the individuality and complexity of each dialog, with some being significantly more complicated than

Experiment	Weighted Recall	Weighted Precision	Macro F1
No masking; context=1; t=0.9	0.62	0.67	0.43
Masking; context=1; t=0.9	0.61	0.72	0.43
Masking; context=1; t=0.0	0.58	0.69	0.41
Masking; context=1; t=0.5	0.58	0.69	0.4
Masking; context=1; t=0.9; reasoning	0.58	0.67	0.42
Masking; context=3; t=0.9	0.59	0.72	0.41
Masking; context=5; t=0.9	0.61	0.67	0.42

Table 2: Evaluation of annotation by ChatGPT using Tree-like scheme (on a subset of dialogs)

	Weighted Recall	Weighted Precision	Macro F1
Full tags	0.56	0.67	0.44
Full tags & voting	0.6	0.71	0.46
Cut labels	0.81	0.82	0.54
Cut labels & voting	0.84	0.86	0.59

(a) Crowdsourcers

	Weighted Recall	Weighted Precision	Macro F1
Full tags	0.41	0.59	0.34
Full tags & voting	0.42	0.6	0.33
Cut labels	0.74	0.78	0.5
Cut labels & voting	0.73	0.77	0.49

(b) ChatGPT

Table 3: Evaluation of final annotation by ChatGPT and crowdsourcing workers as compared to expert annotation (all dialogs)

others.

As for cost, annotation with ChatGPT varies depending of a tree length for a particular dialog from 0.03\$ to 0.07\$ while crowdsourcing workers need to be paid from 0.12\$ to 0.22\$ for one dialog annotation. So, ChatGPT can be used as a silver standard of annotation instead of crowdsourcing results, which would reduce the time and money spent on experts' post-annotation. However, working with such abstract annotation classes, it is still important to rely on non-expert annotators to make the taxonomy easy to comprehend.

6 Conclusion and Future Work

We conducted several experiments on the annotation of casual conversations with speech function taxonomy performed by experts in linguistics, crowdsourcing workers, and ChatGPT. In this paper, we took a closer look at the problems of defining multilayer taxonomies in real dialogs and, furthermore, explored whether it is possible to differentiate between those classes when annotating. Experiments with ChatGPT have demonstrated the potential of using LLMs for linguistic annotation with accuracy that is close to crowdsourcing workers' performance on some dialogs. Even though guid-

ing the model across a tree-like structure of instructions to reach the final label seems to be promising, it still falls short of non-expert performance on such tasks and does not let the researchers explore variations in how non-experts understand discourse structures. With the inevitable future improvement in LLMs' performance, we can expect automatic discourse annotation to serve as gold standard.

It is important to mention that a significant drawback of the method we propose is the necessity of expert involvement in writing prompts and structuring them the right way. However, with LLMs, this process turned out to be extremely similar to the process of writing instructions for non-expert crowdsourcing workers and should thus pose no difficulty to a discourse researcher.

Possible areas for the future work are: 1) trying out other instruction-based models; 2) conducting a more comprehensive selection of hyperparameters; 3) adding criticism steps to the current pipeline, allowing the model to reflect on its answers and to correct itself (Kim et al., 2023).

Acknowledgements

References

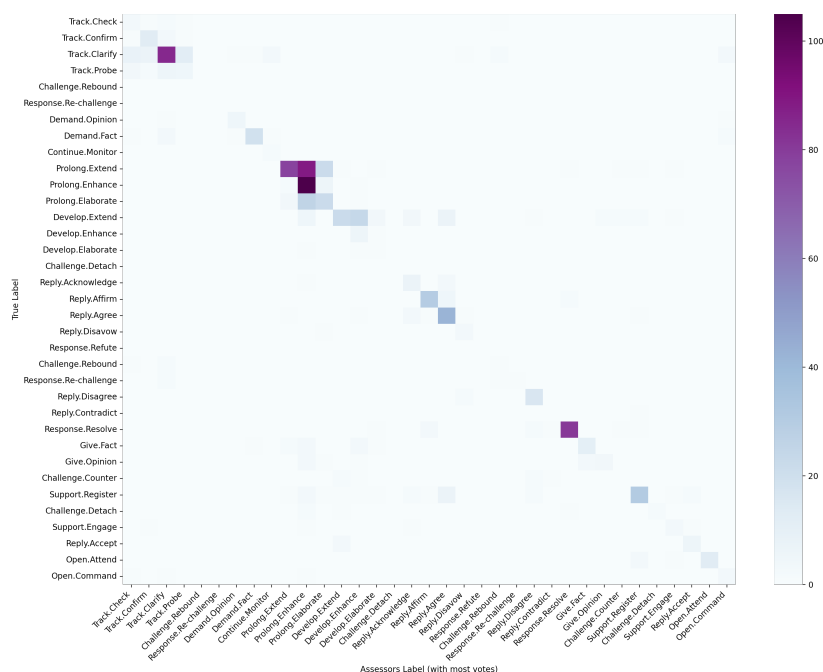
- James Allen and Mark Core. 1997. Damsl: Dialogue act markup in several layers (draft 2.1). In *Technical Report, Multiparty Discourse Group, Discourse Resource Initiative*.
- Nicholas Asher, Julie Hunter, Mathieu Morey, Farah Benamara, and Stergos Afantenos. 2016. Discourse structure and dialogue acts in multiparty dialogue: the stac corpus. In *10th International Conference on Language Resources and Evaluation (LREC 2016)*, pages 2721–2727.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Harry Bunt, Jan Alexandersson, Jean Carletta, Jae-Woong Choe, Alex Chengyu Fang, Kiyong Lee, Volha Petukhova, Andrei Popescu-Belis, Laurent Romary, Claudia Soria, et al. 2010. Towards an iso standard for dialogue act annotation. In *Seventh conference on International Language Resources and Evaluation (LREC’10)*.
- Harry Bunt, Jan Alexandersson, Jae-Woong Choe, Alex C Fang, Koiti Hasida, Volha Petukhova, Andrei Popescu-Belis, and David Traum. 2012. Iso 24617-2: A semantically-based standard for dialogue annotation. Technical report, University of Southern California Los Angeles.
- Mikhail Burtsev, Alexander Seliverstov, Rafael Airapetyan, Mikhail Arkhipov, Dilyara Baymurzina, Nickolay Bushkov, Olga Gureenkova, Taras Khakhulin, Yuri Kuratov, Denis Kuznetsov, and Vasily Konovalov. 2018. [DeepPavlov: Open-source library for dialogue systems](#). In *NIPS*.
- Jon Z Cai, Brendan King, Margaret Perkoff, Shiran Dudy, Jie Cao, Marie Grace, Natalia Wojcik, Ananya Ganesh, James H Martin, Martha Palmer, et al. 2023. Dependency dialogue acts—annotation scheme and case study. *arXiv preprint arXiv:2302.12944*.
- Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. 2021. Dialogsum: A real-life scenario dialogue summarization dataset. *arXiv preprint arXiv:2105.06762*.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30.
- 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, pages 28–35. Boston, MA.

- Malcolm Coulthard. 2014. *An introduction to discourse analysis*. Routledge.
- Suzanne Eggins and Diana Slade. 2004. *Analysing casual conversation*. Equinox Publishing Ltd.
- Boris Galitsky and Dmitry Ilvovsky. 2017. Chatbot with a discourse structure-driven dialogue management. In *Proceedings of the Software Demonstrations of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pages 87–90.
- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks. *arXiv preprint arXiv:2303.15056*.
- Xiaodong Gu, Kang Min Yoo, and Jung-Woo Ha. 2021. Dialogbert: Discourse-aware response generation via learning to recover and rank utterances. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12911–12919.
- Pavel Gulyaev, Eugenia Elistratova, Vasily Konovalov, Yuri Kuratov, Leonid Pugachev, and Mikhail Burtsev. 2020. [Goal-oriented multi-task bert-based dialogue state tracker](#).
- Fan Huang, Haewoon Kwak, and Jisun An. 2023. Is chatgpt better than human annotators? potential and limitations of chatgpt in explaining implicit hate speech. *arXiv preprint arXiv:2302.07736*.
- Dan Jurafsky. 1997. Switchboard swbd-damsl shallow-discourse-function annotation coders manual. [www.dcs.shef.ac.uk/nlp/amities/files/bib/fics-tr-97-02.pdf](#).
- Ali Kashefi and Tapan Mukerji. 2023. Chatgpt for programming numerical methods. *Journal of Machine Learning for Modeling and Computing*.
- Geunwoo Kim, Pierre Baldi, and Stephen McAleer. 2023. Language models can solve computer tasks. *ArXiv*, abs/2303.17491.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*.
- Vasily Konovalov, Pavel Gulyaev, Alexey Sorokin, Yuri Kuratov, and Mikhail Burtsev. 2020. [Exploring the bert cross-lingual transfer for reading comprehension](#). In *Komp’juternaja Lingvistika i Intellektual’nye Tehnologii*, pages 445–453.
- Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. 2023. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. *PLoS digital health*, 2(2):e0000198.

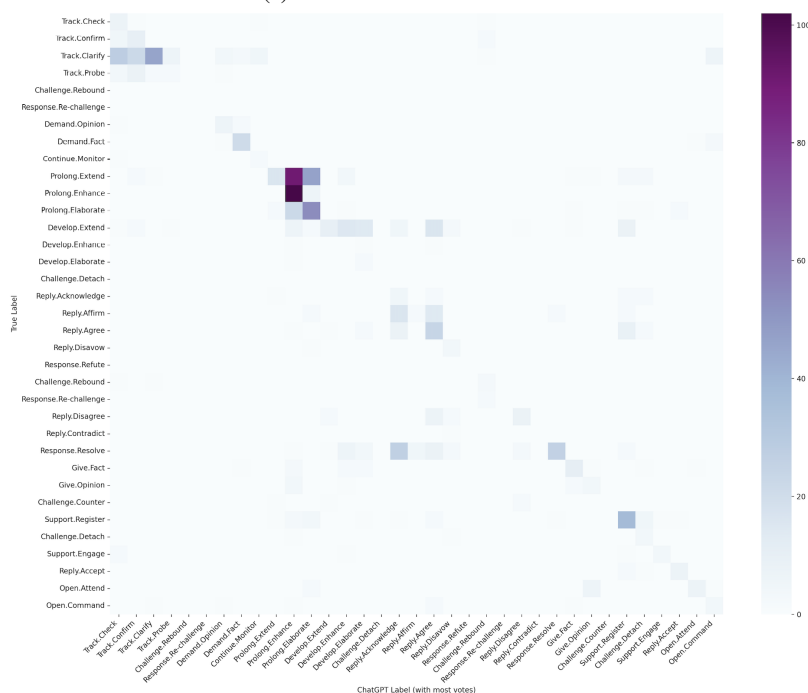
718	Alex Lascarides and Nicholas Asher. 2007. Segmented	Dominik Sobania, Martin Briesch, Carol Hanna, and	772
719	discourse representation theory: Dynamic semantics	Justyna Petke. 2023. An analysis of the automatic	773
720	with discourse structure. <i>Computing meaning</i> , pages	bug fixing performance of chatgpt. <i>arXiv preprint</i>	774
721	87–124.	<i>arXiv:2301.08653</i> .	775
722	Jiaqi Li, Ming Liu, Min-Yen Kan, Zihao Zheng, Zekun	Haoye Tian, Weiqi Lu, Tsz On Li, Xunzhu Tang, Shing-	776
723	Wang, Wenqiang Lei, Ting Liu, and Bing Qin. 2020.	Chi Cheung, Jacques Klein, and Tegawendé F Bis-	777
724	Molweni: A challenge multiparty dialogues-based	syandé. 2023. Is chatgpt the ultimate program-	778
725	machine reading comprehension dataset with dis-	ming assistant—how far is it? <i>arXiv preprint</i>	779
726	course structure. <i>arXiv preprint arXiv:2004.05080</i> .	<i>arXiv:2304.11938</i> .	780
727	Lingyao Li, Lizhou Fan, Shubham Atreja, and Libby	Petter Törnberg. 2023. Chatgpt-4 outperforms experts	781
728	Hemphill. 2023. "hot" chatgpt: The promise of chat-	and crowd workers in annotating political twitter	782
729	gpt in detecting and discriminating hateful, offensive,	messages with zero-shot learning. <i>arXiv preprint</i>	783
730	and toxic comments on social media. <i>arXiv preprint</i>	<i>arXiv:2304.06588</i> .	784
731	<i>arXiv:2304.10619</i> .		
732	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin	785
733	Cao, and Shuzi Niu. 2017. Dailydialog: A manually	Guu, Adams Wei Yu, Brian Lester, Nan Du, An-	786
734	labelled multi-turn dialogue dataset. <i>arXiv preprint</i>	drew M Dai, and Quoc V Le. 2021. Finetuned lan-	787
735	<i>arXiv:1710.03957</i> .	guage models are zero-shot learners. <i>arXiv preprint</i>	788
		<i>arXiv:2109.01652</i> .	789
736	Kaihui Liang, Austin Chau, Yu Li, Xueyuan Lu, Dian	Yang Yang, Juan Cao, Yujun Wen, and Pengzhou Zhang.	790
737	Yu, Mingyang Zhou, Ishan Jain, Sam Davidson, Josh	2022. Multiturn dialogue generation by modeling	791
738	Arnold, Minh Nguyen, et al. 2020. Gunrock 2.0: A	sentence-level and discourse-level contexts. <i>Scien-</i>	792
739	user adaptive social conversational system. <i>arXiv</i>	<i>tific Reports</i> , 12(1):20349.	793
740	<i>preprint arXiv:2011.08906</i> .		
741	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	794
742	Hiroaki Hayashi, and Graham Neubig. 2023. Pre-	Thomas L. Griffiths, Yuan Cao, and Karthik	795
743	train, prompt, and predict: A systematic survey of	Narasimhan. 2023. Tree of thoughts: Deliberate	796
744	prompting methods in natural language processing.	problem solving with large language models.	797
745	<i>ACM Computing Surveys</i> , 55(9):1–35.		
746	William C Mann and Sandra A Thompson. 1988.	Dian Yu and Zhou Yu. 2019. Midas: A dialog act an-	798
747	Rhetorical structure theory: Toward a functional the-	notation scheme for open domain human machine spo-	799
748	ory of text organization. <i>Text-interdisciplinary Jour-</i>	ken conversations. <i>arXiv preprint arXiv:1908.10023</i> .	800
749	<i>nal for the Study of Discourse</i> , 8(3):243–281.		
750	Harsha Nori, Nicholas King, Scott Mayer McKinney,	Frances Yung, Vera Demberg, and Merel Scholman.	801
751	Dean Carignan, and Eric Horvitz. 2023. Capabili-	2019. Crowdsourcing discourse relation annotations	802
752	ties of gpt-4 on medical challenge problems. <i>arXiv</i>	by a two-step connective insertion task. In <i>Proceed-</i>	803
753	<i>preprint arXiv:2303.13375</i> .	<i>ings of the 13th Linguistic Annotation Workshop</i> ,	804
		pages 16–25.	805
754	OpenAI. 2022. Introducing chatgpt . Accessed on May	Yiming Zhu, Peixian Zhang, Ehsan-Ul Haq, Pan Hui,	806
755	13, 2023.	and Gareth Tyson. 2023. Can chatgpt reproduce	807
		human-generated labels? a study of social computing	808
		tasks. <i>arXiv preprint arXiv:2304.10145</i> .	809
756	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,		
757	Carroll Wainwright, Pamela Mishkin, Chong Zhang,		
758	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.		
759	2022. Training language models to follow instruc-		
760	tions with human feedback. <i>Advances in Neural</i>		
761	<i>Information Processing Systems</i> , 35:27730–27744.		
762	Merel CJ Scholman, Jacqueline Evers-Vermeul, Ted JM		
763	Sanders, et al. 2016. A step-wise approach to dis-		
764	course annotation: Towards a reliable categoriza-		
765	tion of coherence relations. <i>Dialogue & Discourse</i> ,		
766	7(2):1–28.		
767	Yuntao Shou, Tao Meng, Wei Ai, Sihan Yang, and Keqin		
768	Li. 2022. Conversational emotion recognition stud-		
769	ies based on graph convolutional neural networks		
770	and a dependent syntactic analysis. <i>Neurocomputing</i> ,		
771	501:629–639.		

A Confusion matrices comparing crowdsourced/ChatGPT annotation with true labels

810



(a) Crowdsourced annotation



(b) ChatGPT annotation

B Speech Functions list

811

Cut labels	Full labels	Definition
Open.Demand.Fact	Open.Demand.Fact	Demanding factual information.
Open.Demand. Opinion	Open.Demand.Opinion	Demanding judgment or evaluative information from the interlocutor.
Open.Give.Fact	Open.Give.Fact	Providing factual information.
Open.Give.Opinion	Open.Give.Opinion	Providing judgment or evaluative information.

Open.Command	Open.Command	Making a request, an invitation or command to start a dialog or discussion of a new topic.
Open.Attend	Open.Attend	These are usually greetings.
React.Rejoinder. Confront.Response	React.Rejoinder.Confront. Response.Re-challenge	Offering an alternative position, often an interrogative sentence.
React.Rejoinder. Support.Track	React.Rejoinder.Support.Track. Probe	Requesting a confirmation of the information necessary to make clear the previous speaker's statement.
	React.Rejoinder.Support.Track. Check	Getting the previous speaker to repeat an element or the entire statement that the speaker has not heard or understood.
	React.Rejoinder.Support.Track. Clarify	Asking a question to get additional information on the current topic of the conversation. Requesting to clarify the information already mentioned in the dialog.
	React.Rejoinder.Support.Track. Confirm	Asking for a confirmation of the information received.
Sustain.Continue. Prolong	Sustain.Continue.Prolong. Extend	Adding supplementary or contradictory information to the previous statement.
	Sustain.Continue.Prolong. Enhance	Adding details to the previous statement, adding information about time, place, reason, etc.
	Sustain.Continue.Prolong. Elaborate	Clarifying / rephrasing the previous statement or giving examples to it.
React.Rejoinder. Confront.Challenge. Rebound	React.Rejoinder.Confront. Challenge. Rebound	Questioning the relevance, reliability of the previous statement, most often an interrogative sentence.
React.Respond. Support.Reply	React.Respond.Support.Reply. Affirm	A positive answer to a question or confirmation of the information provided. Yes/its synonyms or affirmation.
	React.Respond.Support.Reply. Acknowledge	Indicating knowledge or understanding of the information provided.
	React.Respond.Support.Reply. Agree	Agreement with the information provided. In most cases, the information that the speaker agrees with is new to him. Yes/its synonyms or affirmation.
React.Respond. Support.Develop	React.Respond.Support.Develop. Extend	Adding supplementary or contradictory information to the previous statement.
	React.Respond.Support.Develop. Enhance	Adding details to the previous statement, adding information about time, place, reason, etc.
	React.Respond.Support.Develop. Elaborate	Clarifying / rephrasing the previous statement or giving examples to it. A declarative sentence or phrase (may include for example, I mean, like).
React.Respond. Confront.Reply	React.Respond.Confront.Reply. Disagree	Negative answer to a question or denial of a statement. No, negative sentence.
	React.Respond.Confront.Reply. Contradict	Refuting previous information. No, sentence with opposite polarity. If the previous sentence is negative, then this sentence is positive, and vice versa.
	React.Respond.Confront.Reply. Disavow	Denial of knowledge or understanding of information.
Sustain.Continue. Monitor	Sustain.Continue.Monitor	Checking the involvement of the listener or trying to pass on the role of speaker to them.
Sustain.Continue. Command	Sustain.Continue.Command	Making a request, an invitation or command to start a dialog or discussion of a new topic.
React.Respond. Support.Register	React.Respond.Support.Register	A manifestation of emotions or a display of attention to the interlocutor.
React.Respond. Support.Engage	React.Respond.Support.Engage	Drawing attention or a response to a greeting.
React.Respond. Support.Reply. Accept	React.Respond.Support.Reply. Accept	Expressing gratitude.
React.Rejoinder. Support.Response. Resolve	React.Rejoinder.Support. Response.Resolve	The response provides the information requested in the question.
React.Respond. Command	React.Respond.Command	Making a request, an invitation or command to start a dialog or discussion of a new topic.
React.Rejoinder. Confront.Challenge. Detach	React.Rejoinder.Confront. Challenge.Detach	Terminating the dialog.