

MultiCAT: Multimodal Communication Annotations for Teams

Anonymous ACL submission

Abstract

Successful teamwork requires team members to understand each other and communicate effectively, managing multiple linguistic and paralinguistic tasks at once. Because of the potential for interrelatedness of these tasks, it is important to have the ability to make multiple types of predictions on the same dataset. Here, we introduce Multimodal Communication Annotations for Teams (MultiCAT), a speech- and text-based dataset consisting of audio recordings, automated and hand-corrected transcriptions. MultiCAT builds upon collected data for teams working collaboratively to save victims in a simulated search and rescue mission, and consists of annotations and benchmark results for the following tasks: (1) dialog act classification, (2) adjacency pair detection, (3) sentiment and emotion recognition, (4) closed-loop communication detection, and (5) phonetic entrainment detection. We posit that additional work on these tasks and their intersection will further improve understanding of team communication and its relation to team performance.

1 Introduction

Multimodal datasets are useful for a variety of machine learning tasks—e.g., automatic speech recognition (ASR) (Pratap et al., 2020), classifying speaker age (Sadjadi et al., 2016) and gender (Abouelenien et al., 2017), and classification of paralinguistic identifiers such as emotion and sentiment (Bagher Zadeh et al., 2018; Poria et al., 2019). Furthermore, spoken dialog datasets are useful for making predictions about dialog acts and group relationships (Shriberg et al., 2004). However, such datasets tend to focus on a single task or a small set of closely related tasks, thereby limiting the scope of research questions that they can be used to answer.

In this paper, we present *Multimodal Communication Annotations for Teams (MultiCAT)*, a novel speech- and text-based dataset that is annotated for

multiple distinct tasks of interest, as well as containing team-based measures of success to allow for examination of the interactions among tasks. MultiCAT consists of annotations for sentiment, emotion, dialog acts, adjacency pairs, phonetic entrainment, and closed loop communication for multiparty dialog in a collaborative search and rescue task. The primary contributions of this paper are the following:

(1) Dataset: We present a novel multiparty spoken dialog dataset with annotations for related paralinguistic and conversational classification and regression tasks. Notably, to our knowledge, this is the *first publicly available dataset for closed-loop communication detection*.

(2) Baseline models: We develop a set of baseline models for a number of related conversational and paralinguistic tasks and evaluate them on the dataset.

The rest of the paper is organized as follows. We provide a high-level summary of the dataset (§ 2), followed by sections that each focus on a single type of annotation (§ 3 – § 6), describing related work, the annotation procedure, and the benchmark results. Finally, we conclude in § 8.

2 Dataset

MultiCAT contains diverse annotations for text and audio data on individual- and conversation-based tasks. We annotate a subset of the ASIST Study 3 dataset (Huang et al., 2022b,a)—an existing dataset from a large-scale, remotely-conducted human-machine teaming experiment involving teams of three humans executing simulated urban search and rescue missions in a virtual Minecraft-based testbed. Each mission corresponds to an experimental trial, so we use the terms ‘mission’ and ‘trial’ interchangeably in this paper. Each team member has unique capabilities and information, ensuring that they must communicate with each

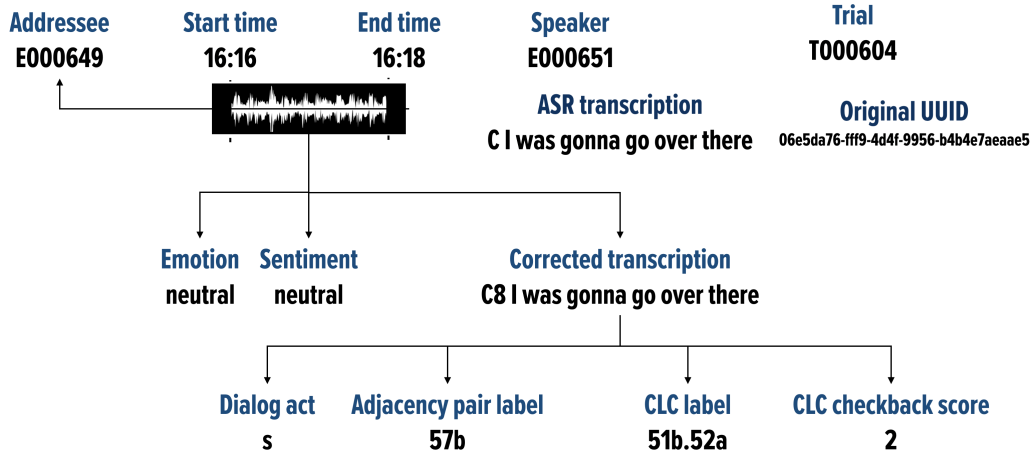


Figure 1: Organization of utterances and labels within the MultiCAT dataset, illustrated by example annotations for a single utterance. The figure also depicts the annotation flow—addressee, emotion, and sentiment annotation and transcript correction are based on the original audio recordings, followed by the corrected transcripts being used for the dialog act, adjacency pair, and CLC annotation tasks. For clarity, we omit IPU annotations in this figure.

other to achieve the best results. The goal of the mission is to maximize the team’s score, which is based on the number of victims identified, triaged, and moved to a safe zone within a 15-minute time limit. Participants are recruited from a pool of adults in the US who play Minecraft and speak English. Selected participant demographic details are provided in Table 8. The complete MultiCAT dataset is included in the supplementary material in the form of an SQLite3 database (multicat.db).

We annotate a subset of the data for sentiment, emotion, dialog acts, adjacency pairs, closed-loop communication events, utterance addressee, and interpausal unit boundaries (see Figure 1). In addition, we provide corrected gold transcriptions for the conversations, which originally had ASR-generated transcriptions. Details of the initial data collection, data annotation, and preprocessing procedures follow.

Data collection procedure Participants fill out a series of surveys related to their background with Minecraft, their leadership style, and sociological factors that may impact their performance in the study. They then participate in two separate missions with the same team, either on their own or with an AI advisor assisting them. Participants use their own computer for the task, and as such their setups may vary. Their speech is recorded on separate channels, with utterance-level transcriptions obtained in real time using Google’s enhanced phone call speech to text model.¹ Participants were

¹<https://cloud.google.com/speech-to-text/docs/enhanced-models>

compensated with either a \$35 Amazon gift card or course credit. If they were unable to complete the study due to technological issues, they were compensated at the rate of \$15 per hour, rounded up to the nearest hour.

Annotation procedure overview The starting point for data in MultiCAT is a set of utterance-aligned speech and text transcriptions. We trained three annotators (two graduate students and one undergraduate student) who completed annotation tasks that matched their knowledge and experience. The annotators were all native or highly proficient English speakers, and were paid the standard hourly student wage set by their respective universities. They underwent an iterative training procedure while working to achieve task-specific acceptable levels of agreement on a small portion of the data (the annotations from the training process is not included in the MultiCAT dataset); subsequent annotations were completed by one annotator each. The total numbers of items in MultiCAT with each label for each task are provided in Appendix C.

2.1 Dataset overview

The dataset is structured as follows. All utterances have a unique identifier (UUID) generated as part of the ASR transcription process, with the exception of a relatively small number of utterances (401) that were inserted as part of the manual transcript correction process—these can nevertheless be uniquely identified by combining their trial ID, participant ID, and start timestamp. Each item is associated with its speaker, as well as the mission in which

Quantity	Total					Annotation	# Trials	# Utts	
Trials	49		Mean	Min	Max	SD	Emotion	46	7731
Teams	25	Utts./spkr	151.0	42	287	53.5	Sentiment	46	7731
Speakers	73	Utts./trial	225.0	91	348	65.0	CLC	36	6544
Utterances	11024	Utts./spkr/trial	78.2	19	156	28.0	Gold transcript	45	4666
Word types	2607	Word types/utt.	8.7	1	74	8.2	Dialog act	45	10342
Word tokens	108475	Word tokens/utt.	9.9	1	118	10.8	APs	45	6846
							Entrainment	8	2896

(a) Totals

(b) Mean, minimum, maximum, and SD.

(c) Number of trials and annotated utterances for our annotation types.

Table 1: Highlights of the MultiCAT dataset. Not all utterances receive labels for all the tasks. AP, DA, and CLC tasks; only items with valid labels are counted here.

it was created and the start and end time of the utterance. Along with the task-specific labels, we also annotate instances of background noises.

A closer examination of the dataset (see Table 1 for details) reveals its particular benefits for the end user. The dataset contains a total of 11,024 utterances. Trials vary in amount of communication, ranging from 91 to 348 utterances. There is further variability in the amount of conversation attributed to an individual team member, with the number of utterances ranging from 19 to 156. This variability lends itself to an exploration of the dynamics of teamwork, different types of team members, and their relationships with team performance.

Differing numbers of trials were used for annotating different tasks due to small minority classes (emotion and sentiment annotation) and the difficulty of annotation (IPU boundary and addressee annotation). A detailed breakdown of which trials are annotated for which tasks can be found in Appendix D.

3 Dialog acts

Related work A dialog act (DA) is the communicative function underlying a speaker’s utterance (Bunt et al., 2020). While numerous annotated resources are available for DAs, their annotation schemes vary depending on their purpose, such as capturing domain-specific phenomena. The Switchboard Dialog Act (SwDA) (Jurafsky et al., 1997) and the Meeting Recorder Dialog Act (MRDA) (Shriberg et al., 2004) corpora are both based on naturally occurring conversations, and use the DAMSL (Core and Allen, 1997) tag-set with some modifications—an approach we adopt as well. While the SwDA corpus contains dyadic dialog, the MRDA dataset contains multi-party (defined as involving more than two interlocutors) dialog.

DailyDialog (Li et al., 2017) is a text-based

dataset using short human-written dyadic dialogs that follows Amanova et al. (2016). This dataset differs from ours in two notable ways. First, while DailyDialog contains annotations for only four DA labels, we use many more DA labels since we are interested in more fine-grained intentions. Second, the conversations in the DailyDialog corpus are more formal and less task-oriented compared to the conversations in our dataset that are naturalistic and occur in the context of a collaborative task. The STAC corpus (Asher et al., 2016) annotations capture the dialog structure in a multiparty setting. The communication occurs over a chat interface where the participants play a non-cooperative game with opposing goals. We capture the conversation flow by means of adjacency pairs.

Annotation procedure For our annotations of dialog acts (DAs), we used the framework from the MRDA dataset, which, like MultiCAT, consists of natural task-oriented human conversations. Under this framework, each utterance is annotated with a ‘general’ and zero or more ‘specific’ tags. Due to imperfect segmentation by the ASR system, our data contained single utterances that should have been split up into multiple utterances. To align the DA annotations with the rest of the annotation tasks while still letting an utterance have more than one DA label, we use the pipe symbol (|) to indicate segmentation. Finally, since inter-annotator agreement on the Accept (aa) and Acknowledgment (bk) tags was very low, we merged them into a single tag (aa). In total, there are 11 general tags and 38 specific tags². The inter-annotator agreement measured using Cohen’s κ is 0.6238 for the general DA category.

²We do not annotate for rising tone (rt), which is a non-DA tag.

Label set	C	Macro F1	Accuracy (%)
Fine-grained	50	30.75	63.24
Coarse-grained	5	42.15	93.92

Table 2: Macro F1 and accuracy for DA classification on fine-grained and top-level classes. |C| is the number of DA labels. The scores are the average of three random runs.

Adjacency pairs We also annotate the conversational structure in the dialogs using the conventions for adjacency pairs (APs) presented in MRDA (Dhillon et al., 2004). APs capture paired utterances such as question-answer, greeting-greeting, etc. An AP for a sequence of utterances is defined such that it contains two parts, each containing one or more utterances and uttered by different speakers (Levinson et al., 1983).

Baseline model We use He et al.’s (2021) baseline model, and include results for the 50 fine-grained and 5 basic labels³ on the corrected transcripts. Since this is a highly imbalanced dataset, we report the macro F1 score along with the accuracy in Table 2. See Appendix E for further details on the model training.

4 Closed-loop communication

Related work Good teamwork processes enable teams to perform beyond the sum of their parts (Roberts et al., 2021). Closed-loop communication (CLC) has been proposed in the team science literature as one of the coordinating mechanisms for effective teamwork (Salas et al., 2005). This communication strategy has been implemented in military contexts to reduce the frequency of communication breakdowns in teams (Burke et al., 2004), and is being explored in the context of health-care as well (Parush et al., 2011). CLC has been shown to be correlated with improved outcomes in both simulations (Diaz and Dawson, 2020) and the real world (Härgestam et al., 2013; El-Shafy et al., 2018), with studies suggesting that high-performing teams tend to display CLC more often than low-performing teams (Bowers et al., 1998), and that deviations from CLC can lead to information loss (Parush et al., 2011) and degraded task performance (Lieber et al., 2022). These findings suggest the utility of developing methods to automatically detect deviations from CLC proto-

cols in real-time, in order to provide appropriate interventions—e.g., an AI agent that informs the team in a timely manner when there is a communication breakdown.

Automated CLC detection is a relatively understudied task. Rosser et al. (2019) developed an NLP-based method to identify CLC and found positive relationships between the outputs of their algorithm and annotations performed by a trained human annotator. However, we were not able to find further details on their method or dataset. Winner et al. (2022) assess the usability of a ‘Team Dynamics Measurement System’ (TDMS) prototype, which implements a measure of CLC that relies solely on communication flow data (e.g., interlocutor identity, utterance timing, and turn-taking patterns), while ignoring the actual content of the utterances. Robinson et al. (2023b) improve upon the flow-based measure by incorporating keyword analysis to analyze the content of the utterances. The dataset used for both of these studies (Robinson et al., 2023a) is not publicly available, limiting our ability to compare our work to theirs.

Though varying definitions of CLC can be found in the literature (Diaz and Dawson, 2020; Salik and Ashurst, 2022; Salas et al., 2005; Marzuki et al., 2019; Härgestam et al., 2013), most definitions of what we refer to as a CLC ‘event’ include the following three sub-events occurring in sequence:

1. *Call-out*: Interlocutor I_1 shares information with/gives an instruction to interlocutor I_2 (Butcher, 2018).
2. *Check-back*: I_2 confirms their understanding of the information/instruction by repeating it back to I_1 .
3. *Closing*: I_1 confirms that I_2 has received and understood the information or performed the desired action.

To the best of our knowledge, MultiCAT is the first publicly available dataset for studying closed-loop communication (CLC). Most existing CLC research is conducted by watching videos and recording only the parts that researchers are interested in (e.g., CLC categories (Marzuki et al., 2019) and task completion time (El-Shafy et al., 2018)) without transcribing the entire communication.

Annotation procedure Annotators were trained to identify and label CLC sub-events and score the quality of check-backs on a scale of 1–3, as detailed

³The 5 basic tags are Statement, Filler, Backchannel, Disruption, and Question.

in Table 3. We used a, b, and c to denote call-outs, check-backs, and closings, respectively, to partially align our CLC labels with the labels for AP components. The inter-annotator agreement calculated using Krippendorff’s α was 0.676. We deemed this an acceptable level of agreement given the challenging nature of this annotation task, which involves a nontrivial amount of subjective interpretation, dealing with ambiguity, and keeping large windows of utterances in the annotator’s working memory.

Baseline Model We use a three-stage approach to identifying CLC events.

1. In the first step, we construct TF-IDF feature vectors from lemmatized versions of the utterances, which are then used as inputs to a logistic regression model that predicts whether or not an utterance corresponds to a call-out sub-event (i.e., a).
2. For each utterance that is labeled as a call-out, we examine the next three utterances following that utterance that are from a speaker other than the source of the call-out utterance. For each of the call-outs and their three candidate check-back pairs, we use a RoBERTa-based sequence classification model fine-tuned on MultiCAT to predict whether the candidate utterances check back to the call-out utterance (i.e., b).
3. Given the rarity of ‘closing’ sub-events, and the resultant analytical complexity, we combine subevent sequences ab and abc into a CLC event category, contrasting it against isolated call-outs classified as ‘open loop events’. This pragmatic categorization is consistent with the prevalence of two-stage CLC events in real-world scenarios noted by Robinson et al. (2023b) and Marzuki et al. (2019).

We aggregated the labels from prior steps to classify the overall CLC event status into three categories: closed-loop event, open-loop event, and non-CLC event. For every utterance, if a call-out sub-event is detected, and if at least one check-back is detected within the next three utterances from speakers other than the original speaker, we conclude that this call-out is ‘closed’ and a CLC event has occurred. Conversely, if no check-back is detected among the three candidate utterances for that call-out, then the call-out by itself

forms an open-loop event. Non-CLC events are categorized as situations where the initial call-out is not detected at all.

Results for all three stages are provided in Table 4, and details on our model training are provided in Appendix E.

5 Sentiment/Emotion recognition

Previous work Datasets for sentiment and emotion have largely been annotated for one or both tasks, but not others. GEMEP (Bänzinger et al., 2012) and IEMOCAP (Busso et al., 2008) contain a total of 10 actors each simulating a range of emotions. Both contain high-quality recordings but are relatively small corpora. RAVDESS (Livingstone and Russo, 2018) likewise contains actors simulating emotion, with an additional annotation for the intensity of the emotion. The YouTube dataset (Morency et al., 2011) contains 47 videos of single speakers, with utterances annotated for sentiment. Similarly, ICT-MMMO (Wöllmer et al., 2013) contains single-speaker data annotated for sentiment, with each item being relatively long.

The Multimodal Emotion Lines Dataset (MELD) (Poria et al., 2019) consists of conversations from the TV show *Friends* and is annotated for Ekman’s universal emotions (Ekman, 1992) and positive, negative, or neutral sentiment. Likewise, the CMU Multimodal Opinion Sentiment and Emotion Intensity (CMU-MOSEI) dataset (Bagher Zadeh et al., 2018) is annotated for both tasks, with seven sentiment labels ranging from strong negative to strong positive. CMU-MOSEI uses monologue data from YouTube. The DailyDialog dataset—mentioned earlier in this section—is also annotated for Ekman’s universal emotions. While all of these datasets contain annotation types that have some overlap with those present in MultiCAT, none contain the range we present here.

Annotation procedure Two annotators were trained to identify the opinions of the speaker towards the subject (sentiment) and the affect shown by the speaker (emotion) during an utterance, by listening to it in context. Inter-annotator agreement was calculated using Cohen’s κ ; annotators achieved an agreement score of .886 for sentiment and .830 for emotion.

We use the same set of emotions as MELD and DailyDialog, namely Ekman’s universal emotions (Ekman, 1992)—anger, disgust, fear, joy, sadness,

Criteria	Example	Score
<i>Incomplete/inaccurate:</i> The recipient did not confirm their understanding of the information or instruction.	<i>Okay.</i>	1
<i>Partially complete and accurate:</i> The recipient partially confirmed their understanding of the information or instruction.	<i>Okay, I am on my way.</i>	2
<i>Complete and accurate, with all key information repeated:</i> The recipient fully confirmed their understanding of the information or instruction.	<i>Okay, I am on my way to B4 to clear the rubble.</i>	3

Table 3: Rubric for evaluating checkbacks in closed-loop communication events. The middle column shows examples of replies to the hypothetical call-out: “Engineer, can you clear the rubble room B4?”

Stage	Accuracy	F_1
Call-out detection	.771	.789
Check-back detection	.761	.432
Complete CLC event detection	.514	.447

Table 4: Results for CLC detection baseline approach. For the complete CLC event detection stage, we report a weighted F_1 score due to the very small number of ‘closing’ sub-events in the data.

and surprise—along with neutral. Sentiment labels are positive, negative, and neutral.

Baseline model We use a multitask sentiment and emotion classifier based on the model created by Culnan et al. (2021). This model uses low-level acoustic features from the Interspeech 13 feature set created for tasks including emotion and social cues (Schuller et al., 2013) extracted with openS-MILE (Eyben et al., 2010). We use 768-d word embeddings generated with BERT (Devlin et al., 2019) model bert-base-uncased as text features. Text is fed through a bidirectional LSTM, while acoustic features are averaged and fed through feed forward layers. The output of these two components are then concatenated and fed through two feed forward layers to reduce their dimension to 100. Finally, the output of these feed forward layers is passed to task-specific heads to make sentiment and emotion predictions.

The model is pretrained on data from MELD and CMU-MOSI. CMU-MOSI contains sentiment labels from strong negative to strong positive, so we collapse over negative and positive label types to get the same three classes of interest as in MultiCAT. We report F1 for each class and overall macro F1 over a single run due to significant class imbalances. These results are shown in Table 5.

We find that our multitask sentiment and emotion prediction model is more successful at prediction of

sentiment than emotion, with better performance for majority classes than minority classes. In the case of emotion prediction, difficulty arises from two very small minority classes, anger (total support of 18) and disgust (total support of 25).

6 Entrainment detection

Entrainment is the adaption of verbal and non-verbal actions by conversation partners to more closely resemble one another (Borrie and Liss, 2014). It facilitates effective turn taking, builds rapport, and aids in communicating positive sentiments. Correlations between entrainment and desired social outcomes have been reported in cooperative games (Yu et al., 2019; Levitan et al., 2015), patient-therapist relations (Nasir et al., 2020; Borrie et al., 2019), effective communication in study groups (Friedberg et al., 2012), and romantic success (Ireland et al., 2011). Besides English, research has been extended to Hebrew (Weise et al., 2022), Russian (Kachkovskaia et al., 2020; Menshikova et al., 2020), Slovak, Spanish, and Chinese (Levitan et al., 2015).

The study of entrainment faces many challenges. Many popular corpora have relatively a modest number of teams. For example, the Columbia Games Corpus⁴ and the Brooklyn Multi-Interaction Corpus (Weise et al., 2022) have 12 each. Some are also restricted due to being sensitive in nature, such as the Suicide Risk Assessment Corpus (Baucom et al., 2014) and the Couples Therapy Corpus (Christensen et al., 2004), or prohibitively expensive to obtain, such as the Fisher Corpus (Cieri et al., 2004).

Previous studies have heavily relied on pristine recording conditions with professional recording equipment and manual preparation of an acoustic-

⁴<http://www.cs.columbia.edu/speech/games-corpus/>

	Sentiment				Emotion							
	Neg.	Neut.	Pos.	All	Anger	Disgust	Fear	Joy	Neut.	Sadness	Surprise	All
Maj. Class	0.00	72.76	0.00	24.25	0.00	0.00	0.00	0.00	87.97	0.00	0.00	12.57
Strat.	15.02	51.47	28.05	31.51	5.41	0.00	3.23	4.18	77.50	5.56	3.70	14.22
Baseline	43.54	62.66	49.79	51.99	0.00	9.30	16.19	20.11	76.86	30.53	29.20	26.03
Support	370	1310	611	2291	18	25	70	154	1799	145	80	2291

Table 5: Results of baseline model on MultiCAT’s sentiment and emotion test partition. Number of items per class and overall are shown in the bottom row of the table. Overall results reported are Macro F1. Model ‘Maj. Class’ represents a classifier that predicts only the majority class; ‘Strat.’ represents a classifier predicting results using the probability of each class appearing from the training set; ‘Baseline’ is our baseline model for sentiment and emotion prediction. Neg=Negative, Neut=Neutral, Pos=Positive

prosodic feature set, restricting entrainment-specific datasets to laboratory conditions. In contrast, MultiCAT is based on data collected in more realistic conditions, where researchers exert limited control over recording channels, environments, and participant interactions. MultiCAT also enables the analysis of entrainment in short-lived, randomly formed teams in which the teammates do not know each other beforehand.

Annotation procedure Previous research on vocalic entrainment has concentrated on dyadic interactions with balanced turn-taking and responses directed at one intended listener. However, the distribution of utterances in a multi-party (i.e., more than two interlocutors) conversation is less likely to be balanced than in a dyadic conversation. Additionally, in a multi-party conversation, utterances could be aimed at the group as a whole, rather than one intended listener. Thus, there is a need to identify speaker dyads and separate them from utterances with no specific intended listener.

We identify the subset of utterances in three-member trials in which there is a single intended addressee to find dyadic interactions within a multi-party conversation. For this annotation task, 4 teams, i.e., 8 trials, were randomly selected. Annotators completed this annotation task in Praat (Boersma, 2001), using each speaker’s individual audio stream (in order to avoid speaker overlap), gold transcriptions, and Praat textgrids. The data for one of the eight trials (T000605) is missing audio data for one speaker, thus yielding data for 8 trials and 11 unique speakers, with an average duration of 28.1 minutes.

For each trial, annotators identified the boundaries of a stream of audio separated by a pause of 50ms or more, also known as an inter-pausal unit (IPUs). Next, they mapped the audio in each

IPU to the corresponding text from the transcript (an utterance can have one or more IPUs), and identified the addressee of each IPU. The addressee labels had four possibilities—an identifier for each of the three participants, or ‘all’ to indicate a general response or an unknown audience. Annotators achieved a Cohen’s κ score of .478. We deem this acceptable due to the complexity of the task. IPUs with a specific addressee comprise 27.42% of the total number utterances per trial on average (SD = 25.43%).

Baseline We replicate the baseline model used in Nasir et al. (2020) for assessment of their unsupervised model, using the same training corpus, acoustic feature set and hyperparameters. First, 80% of the utterances from the Fisher Corpus English Part 1 (LDC2004S13) (Cieri et al., 2004) (total 5850 spontaneous telephonic dyadic conversations, 10 minutes each) are randomly chosen. A feed-forward neural network encoder-decoder model is used to encode entrainable information from a given utterance and predict the next turn, which is compared to its referent (that is the real ‘next turn’) to compute the loss function of the model.

In order to verify if this model is able to detect entrainment in a multi-party system, we use the verification measures from Nasir et al. (2020), in which the model classifies conversations as ‘real’ (all pairs of adjacent utterances are in order) or ‘fake’ (turns scrambled so that the entrainment information is not preserved) when presented with sample conversations from the test set. We report mean results over 30 runs.

First, dyadic interactions are extracted using the addressee labels for each of the 8 trials ($8 \times 3 = 24$ possibilities). This process yields 11 interactions, a number lower than the expected number (23) due to the fact that not all participants are judged to

have addressed both their team mates. Turn-level acoustic features are then extracted and processed to function as a test set for the model.

The classification accuracy for the MultiCAT entrainment set was 51.86%. This is a much lower score than observed for the two-party Fisher test set and Suicide Corpus in [Nasir et al. \(2020\)](#) (72.10% and 70.44% respectively). This could be due to two factors. First, the increase in the number of interlocutors from a two-party to a multi-party system increases the complexity of detecting entrainment. Second, the differences recording conditions for the training corpus and the MultiCat corpus (controlled vs real-world) pose a challenge to detecting vocalic entrainment, an effect that is sensitive to recording conditions. We choose to report these results because to the best of our knowledge, there are no current benchmarks for unsupervised multi-party entrainment detection.

7 Relationships between label types and team performance

This section presents an exploratory analysis (i.e., for the purpose of generating hypotheses rather than testing them) of the relationships among annotation types and between annotations and team performance. Performance is measured by scores achieved by a team in a single mission. Each team participated in two missions, so each may have two scores associated with it.

Mission scores were calculated based on the number of victims saved in the simulated search and rescue mission, with 10 points awarded for two types of victims, and 50 points awarded for a third victim type. The trials in MultiCAT have final scores ranging from 190 to 890, with a mean of 609.6 and standard deviation of 140.2.

7.1 Relationships among label types

We calculated chi square tests of independence for crosstabs of classes in our tasks. This test shows that AP and emotion labels have a significant relationship ($\chi^2(12)=186.99, p < .001$). Likewise, there is a significant relationship between AP labels and sentiment labels ($\chi^2(4) = 543.49, p < 0.001$)

We compare DA labels with other labels by examining only the general DA label. A chi-squared test of independence for DA labels and sentiment shows a significant relationship between the two ($\chi^2(54) = 395.74, p < .001$). Comparing DA labels with emotion labels demonstrates that there is

	AP	CLC	DA	Sent/Emo	All
RMSE	159.43	142.97	301.14	196.81	204.62
MAE	126.92	112.89	243.38	165.39	175.51

Table 6: Results of basic model predictions for team score on annotations from MultiCAT. Results shown for a linear regression model.

likewise a relationship between the two ($\chi^2(168) = 287.78, p < 0.001$). Sentiment and CLC also show a significant relationship ($\chi^2(12) = 406.95, p < .001$), as do emotion and CLC labels ($\chi^2(36) = 262.04, p < .001$).

7.2 Comparing annotations with outcomes

To examine the relationship between our annotations and the outcomes of individual missions, we generate feature vectors to feed into a simple linear regression model. From the DA annotations we select the general tags as features; if a single item contains multiple general tags, we retain them all. From AP and CLC annotations, we take the general utterance type (i.e. a or b for AP, or a, b, or c for CLC), removing associated number. We use all sentiment and emotion classes. As an additional feature set, we combine all of these feature vectors. We apply min-max scaling to our features prior to using them with a basic model to generate predicted scores.

We use 5-fold cross-validation to predict score for each of the missions containing DA, AP, CLC, sentiment, and emotion labels. Our results are shown in [Table 6](#). We anticipate that creating more sophisticated models to include multiple annotation types will demonstrate additional benefit to an end user over focusing on a single task for downstream tasks related to team performance.

8 Conclusion

In this paper, we have presented MultiCAT, a novel dataset annotated for six separate computational tasks that may be studied individually or in concert to make assessments about team performance. We examine relationships among the annotated tasks and with team scores. For a subset of our tasks, we provide baseline models to be used for comparison in future research. We demonstrate that MultiCAT has use both for individual tasks and for downstream tasks involving multiple annotation types.

9 Limitations

As with any novel dataset, MultiCAT has its limitations. Firstly, data is only in English, largely from native speakers of American English. Conclusions drawn from and patterns found in this dataset may not generalize to other languages or populations.

Additionally, because natural language does not have an equal distribution of items from all dialog act classes, for example, and because each emotion does not appear with equal frequency, datasets consisting of conversations of unconstrained natural language that are created for these tasks will be inherently imbalanced. This is true of MultiCAT, as well. This limitation necessarily affects models seeking to make good predictions about minority classes, as there may be few examples of a given minority class. We believe that acknowledging these limitations in future research will help avoid the risks of overgeneralizing results to other populations and making assumptions about patterns of data in non-English languages.

10 Ethics Statement

In this work, we annotated a subset of the publicly available ASIST Study 3 dataset (Huang et al., 2022b). Our use of the dataset is consistent with its terms of use (CC0 1.0).

Both the collection of the ASIST Study 3 dataset and our analysis of it were approved by IRBs. Participants in the ASIST Study 3 dataset were voluntary participants who signed informed consent forms and were aware of any risks of harm associated with their participation.

The dataset collection process and conditions are described in § 2. The group of annotators was comprised of three graduate students and one undergraduate student. All annotators were compensated fairly for their time in accordance with the standard hourly wages set by their respective departments (in the case of graduate students) or their university (in the case of the undergraduate student).

The characteristics of the dataset are provided in Appendix B. We provide information about the compute resources required for model training in Appendix E.

Intended use If our technology functions as intended, it could be deployed as part of social AI agents embedded in human-machine teams—these agents would be able to understand the affective states of their human teammates, as well as social

dynamics within the team.

Failure modes Failure modes of our technology involve incorrect predictions. It is conceivable (in the context of human-machine teaming) that deteriorated outcomes may result from ineffective human-machine teaming that occurs due to a social AI agent’s inability to understand their human teammates.

Misuse potential It is also conceivable that malicious actors may endow AI agents with the ability to infer sentiment, emotion, team dynamics, etc. in order to perform social engineering for nefarious purposes.

Collecting data from users We are not proposing a system to collect data from users in this paper.

Potential harm to vulnerable populations To our knowledge, the possible harms we have identified are not likely to fall disproportionately on populations that already experience marginalization or otherwise vulnerable.

References

- Mohamed Abouelenien, Verónica Pérez-Rosas, Rada Mihalcea, and Mihai Burzo. 2017. Multimodal gender detection. In *Proceedings of the 19th ACM international conference on multimodal interaction*, pages 302–311.
- Dilafruz Amanova, Volha Petukhova, and Dietrich Klakow. 2016. *Creating annotated dialogue resources: Cross-domain dialogue act classification*. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 111–117, Portorož, Slovenia. European Language Resources Association (ELRA).
- Nicholas Asher, Julie Hunter, Mathieu Morey, Benamara Farah, and Stergos Afantenos. 2016. *Discourse structure and dialogue acts in multiparty dialogue: the STAC corpus*. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 2721–2727, Portorož, Slovenia. European Language Resources Association (ELRA).
- AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2018. *Multimodal language analysis in the wild: CMU-MOSEI dataset and interpretable dynamic fusion graph*. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, Melbourne, Australia. Association for Computational Linguistics.
- Tanja Bänzinger, Marcello Mortarillo, and Klaus R Scherer. 2012. *Introducing the geneva multimodal*

732	expression corpus for experimental research on	
733	emotion perception. <i>Emotion (Washington, D.C.)</i> ,	
734	12(5):1161–1179.	
735	BR Baucom, AO Crenshaw, CJ Bryan, TA Clemans,	
736	TO Bruce, and MD Rudd. 2014. Patient and clinician	
737	vocally encoded emotional arousal as predictors of	
738	response to brief interventions for suicidality. <i>Brief</i>	
739	<i>Cognitive Behavioral Interventions to Reduce Sui-</i>	
740	<i>cide Attempts in Military Personnel. Association for</i>	
741	<i>Behavioral and Cognitive Therapies.</i>	
742	Paul Boersma. 2001. Praat, a system for doing phonetics	
743	by computer. <i>Glott International</i> , 5(9/10):341–345.	
744	Stephanie A Borrie, Tyson S Barrett, Megan M Willi,	
745	and Visar Berisha. 2019. Syncing up for a good	
746	conversation: A clinically meaningful methodology	
747	for capturing conversational entrainment in the speech	
748	domain. <i>Journal of Speech, Language, and Hearing</i>	
749	<i>Research</i> , 62(2):283–296.	
750	Stephanie A Borrie and Julie M Liss. 2014. Rhythm as	
751	a coordinating device: Entrainment with disordered	
752	speech. <i>Journal of Speech, Language, and Hearing</i>	
753	<i>Research</i> , 57(3):815–824.	
754	Clint A. Bowers, Florian Jentsch, Eduardo Salas, and	
755	Curt C. Braun. 1998. Analyzing communication se-	
756	quences for team training needs assessment. <i>Human</i>	
757	<i>Factors</i> , 40:672+. Article.	
758	Harry Bunt, Volha Petukhova, Emer Gilmartin, Cather-	
759	ine Pelachaud, Alex Fang, Simon Keizer, and Laurent	
760	Prévot. 2020. The ISO standard for dialogue act	
761	annotation, second edition. In <i>Proceedings of the</i>	
762	<i>Twelfth Language Resources and Evaluation Confer-</i>	
763	<i>ence</i> , pages 549–558, Marseille, France. European	
764	Language Resources Association.	
765	C S Burke, Eduardo Salas, K Wilson-Donnelly, and	
766	H Priest. 2004. How to turn a team of experts into	
767	an expert medical team: guidance from the aviation	
768	and military communities. <i>BMJ Quality & Safety</i> ,	
769	13(suppl 1):i96–i104.	
770	Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe	
771	Kazemzadeh, Emily Mower, Samuel Kim, Jean-	
772	nette N Chang, Sungbok Lee, and Shrikanth S	
773	Narayanan. 2008. Iemocap: Interactive emotional	
774	dyadic motion capture database. <i>Language resources</i>	
775	<i>and evaluation</i> , 42(4):335–359.	
776	Brad W Butcher. 2018. Leadership and crisis manage-	
777	ment. <i>Rapid Response System: A Practical Guide</i> ,	
778	19.	
779	Andrew Christensen, David C Atkins, Sara Berns, Jen-	
780	nifer Wheeler, Donald H Baucom, and Lorelei E	
781	Simpson. 2004. Traditional versus integrative behav-	
782	ioral couple therapy for significantly and chronically	
783	distressed married couples. <i>Journal of consulting</i>	
784	<i>and clinical psychology</i> , 72(2):176.	
785	Christopher Cieri, David Miller, and Kevin Walker. 2004.	
786	The fisher corpus: A resource for the next generations	
787	of speech-to-text. In <i>LREC</i> , volume 4, pages 69–71.	
	Mark G Core and James Allen. 1997. Coding dialogs	788
	with the DAMSL annotation scheme. In <i>AAAI Fall</i>	789
	<i>Symposium on Communicative Action in Humans and</i>	790
	<i>Machines</i> , volume 56, pages 28–35. Boston, MA.	791
	John Culnan, Seongjin Park, Meghavarshini Krish-	792
	naswamy, and Rebecca Sharp. 2021. Me, myself,	793
	and ire: Effects of automatic transcription quality on	794
	emotion, sarcasm, and personality detection. In <i>Pro-</i>	795
	<i>ceedings of the Eleventh Workshop on Computational</i>	796
	<i>Approaches to Subjectivity, Sentiment and Social Me-</i>	797
	<i>dia Analysis</i> , pages 250–256, Online. Association for	798
	Computational Linguistics.	799
	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	800
	Kristina Toutanova. 2019. BERT: Pre-training of	801
	deep bidirectional transformers for language under-	802
	standing. In <i>Proceedings of the 2019 Conference of</i>	803
	<i>the North American Chapter of the Association for</i>	804
	<i>Computational Linguistics: Human Language Tech-</i>	805
	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	806
	4171–4186, Minneapolis, Minnesota. Association for	807
	Computational Linguistics.	808
	Rajdip Dhillon, Sonali Bhagat, Hannah Carvey, and	809
	Elizabeth Shriberg. 2004. Meeting recorder project:	810
	Dialog act labeling guide. Technical report, Interna-	811
	tional Computer Science Institute, Berkley CA.	812
	Maria Carmen G. Diaz and Kimberly Dawson. 2020. Im-	813
	pact of simulation-based closed-loop communication	814
	training on medical errors in a pediatric emergency	815
	department. <i>American journal of medical quality</i> ,	816
	35(6):474–478.	817
	Paul Ekman. 1992. An argument for basic emotions.	818
	<i>Cognition & emotion</i> , 6(3-4):169–200.	819
	Ibrahim Abd El-Shafy, Jennifer Delgado, Meredith Ak-	820
	erman, Francesca Bullaro, Nathan A. M. Christopher-	821
	son, and Jose M. Prince. 2018. Closed-loop com-	822
	munication improves task completion in pediatric	823
	trauma resuscitation. <i>Journal of surgical education</i> ,	824
	75(1):58–64.	825
	Florian Eyben, Martin Wöllmer, and Björn Schuller.	826
	2010. Opensmile: the munich versatile and fast open-	827
	source audio feature extractor. In <i>Proceedings of the</i>	828
	<i>18th ACM international conference on Multimedia</i> ,	829
	pages 1459–1462.	830
	Heather Friedberg, Diane Litman, and Susannah BF	831
	Paletz. 2012. Lexical entrainment and success in	832
	student engineering groups. In <i>2012 IEEE spoken</i>	833
	<i>language technology workshop (SLT)</i> , pages 404–409.	834
	IEEE.	835
	Maria Härgestam, Marie Lindkvist, Christine Brulin,	836
	Maritha Jacobsson, and Magnus Hultin. 2013. Com-	837
	munication in interdisciplinary teams: exploring	838
	closed-loop communication during in situ trauma	839
	team training. <i>BMJ open</i> , 3(10):e003525.	840
	Charles R. Harris, K. Jarrod Millman, Stéfan J. van der	841
	Walt, Ralf Gommers, Pauli Virtanen, David Cour-	842
	napeau, Eric Wieser, Julian Taylor, Sebastian Berg,	843

844	Nathaniel J. Smith, Robert Kern, Matt Picus, Stephan	Christopher S. Lieber, Yancy Vance Paredes, Aaron	901
845	Hoyer, Marten H. van Kerkwijk, Matthew Brett, Al-	Zhen Yang Teo, and Nancy J. Cooke. 2022. Analysis	902
846	lan Haldane, Jaime Fernández del Río, Mark Wiebe,	of Voice Transmissions of Air Traffic Controllers	903
847	Pearu Peterson, Pierre Gérard-Marchant, Kevin Shep-	in the Context of Closed Loop Communication De-	904
848	pard, Tyler Reddy, Warren Weckesser, Hameer Ab-	viation and its Relationship to Loss of Separation.	905
849	basi, Christoph Gohlke, and Travis E. Oliphant.	<i>Proceedings of the Human Factors and Ergonomics</i>	906
850	2020. Array programming with NumPy. <i>Nature</i> ,	<i>Society Annual Meeting</i> , 66(1):672–676.	907
851	585(7825):357–362.		
852	Zihao He, Leili Tavabi, Kristina Lerman, and Moham-	Steven R. Livingstone and Frank A. Russo. 2018. The	908
853	mad Soleymani. 2021. Speaker turn modeling for	ryerson audio-visual database of emotional speech	909
854	dialogue act classification. In <i>Findings of the Associ-</i>	and song (ravdess): A dynamic, multimodal set of	910
855	<i>ation for Computational Linguistics: EMNLP 2021</i> ,	facial and vocal expressions in north american english.	911
856	pages 2150–2157, Punta Cana, Dominican Republic.	<i>PloS one</i> , 13(5):e0196391–e0196391.	912
857	Association for Computational Linguistics.		
858	Lixiao Huang, Jared Freeman, Nancy Cooke, John	Ilya Loshchilov and Frank Hutter. 2019. Decoupled	913
859	Colonna-Romano, Matthew D Wood, Verica	weight decay regularization. In <i>International Confer-</i>	914
860	Buchanan, and Stephen J Kaufman. 2022a. Exer-	<i>cence on Learning Representations.</i>	915
861	cises for artificial social intelligence in minecraft		
862	search and rescue for teams.	Ernisa Marzuki, Hannah Rohde, Chris Cummins, Holly	916
863	Lixiao Huang, Jared Freeman, Nancy Cooke, John “JCR”	Branigan, Gareth Clegg, Anna Crawford, and Lisa	917
864	Colonna-Romano, Matt Wood, Verica Buchanan, and	MacInnes. 2019. Closed-loop communication during	918
865	Stephen Kaufman. 2022b. Artificial Social Intelli-	out-of-hospital resuscitation: Are the loops really	919
866	gence for Successful Teams (ASIST) Study 3.	closed? <i>Communication and Medicine</i> , 16(1):54–66.	920
867	Molly E Ireland, Richard B Slatcher, Paul W Eastwick,	Alla Menshikova, Daniil Kocharov, and Tatiana	921
868	Lauren E Scissors, Eli J Finkel, and James W Pen-	Kachkovskaia. 2020. Phonetic entrainment in coop-	922
869	nebaker. 2011. Language style matching predicts	erative dialogues: A case of Russian. In <i>Proceedings</i>	923
870	relationship initiation and stability. <i>Psychological</i>	<i>of Interspeech 2020</i> , pages 4148–4152.	924
871	<i>science</i> , 22(1):39–44.		
872	Dan Jurafsky, Elizabeth Shriberg, and Debra Bi-	Louis-Philippe Morency, Rada Mihalcea, and Payal	925
873	asca. 1997. Switchboard SWBD-DAMSL Shallow-	Doshi. 2011. Towards multimodal sentiment analysis:	926
874	Discourse-Function Annotation Coders Manual,	Harvesting opinions from the web. In <i>Proceedings</i>	927
875	Draft 13. Technical report, University of Colorado at	<i>of the 13th international conference on multimodal</i>	928
876	Boulder and +SRI International.	<i>interfaces</i> , pages 169–176.	929
877	Tatiana Kachkovskaia, Tatiana Chukaeva, Vera Evdoki-	Md Nasir, Brian Baucom, Craig Bryan, Shrikanth	930
878	mov, Pavel Kholiavin, Natalia Kriakina, Daniil	Narayanan, and Panayiotis Georgiou. 2020. Model-	931
879	Kocharov, Anna Mamushina, Alla Menshikova, and	ing vocal entrainment in conversational speech using	932
880	Svetlana Zimina. 2020. SibLing corpus of Russian	deep unsupervised learning. <i>IEEE Transactions on</i>	933
881	dialogue speech designed for research on speech	<i>Affective Computing.</i>	934
882	entrainment. In <i>Proceedings of the Twelfth Lan-</i>	Avi Parush, Chelsea Kramer, Tara Foster-Hunt, Kathryn	935
883	<i>guage Resources and Evaluation Conference</i> , pages	Momtahan, Aren Hunter, and Benjamin Sohmer.	936
884	6556–6561, Marseille, France. European Language	2011. Communication and team situation aware-	937
885	Resources Association.	ness in the or: Implications for augmentative infor-	938
886	Stephen C Levinson, Stephen C Levinson, and S Levin-	mation display. <i>Journal of Biomedical Informatics</i> ,	939
887	son. 1983. <i>Pragmatics.</i> Cambridge University Press.	44(3):477–485. Biomedical Complexity and Error.	940
888	Rivka Levitan, Štefan Beňuš, Agustín Gravano, and Julia	Adam Paszke, Sam Gross, Francisco Massa, Adam	941
889	Hirschberg. 2015. Acoustic-prosodic entrainment	Lerer, James Bradbury, Gregory Chanan, Trevor	942
890	in Slovak, Spanish, English and Chinese: A cross-	Killeen, Zeming Lin, Natalia Gimelshein, Luca	943
891	linguistic comparison. In <i>Proceedings of the 16th</i>	Antiga, et al. 2019. Pytorch: An imperative style,	944
892	<i>Annual Meeting of the Special Interest Group on</i>	high-performance deep learning library. <i>Advances in</i>	945
893	<i>Discourse and Dialogue</i> , pages 325–334.	<i>neural information processing systems</i> , 32.	946
894	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang	Soujanya Poria, Devamanyu Hazarika, Navonil Ma-	947
895	Cao, and Shuzi Niu. 2017. DailyDialog: A manually	jumder, Gautam Naik, Erik Cambria, and Rada Mi-	948
896	labelled multi-turn dialogue dataset. In <i>Proceedings</i>	halcea. 2019. MELD: A multimodal multi-party	949
897	<i>of the Eighth International Joint Conference on Natu-</i>	dataset for emotion recognition in conversations. In	950
898	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	<i>Proceedings of the 57th Conference of the Association</i>	951
899	pages 986–995, Taipei, Taiwan. Asian Federation of	<i>for Computational Linguistics, ACL 2019, Florence,</i>	952
900	Natural Language Processing.	<i>Italy, July 28- August 2, 2019, Volume 1: Long Pa-</i>	953
		<i>pers</i> , pages 527–536. Association for Computational	954
		Linguistics.	955

956	Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert. 2020. MIs: A large-scale multilingual dataset for speech research. <i>arXiv preprint arXiv:2012.03411</i> .	1012
957		1013
958		1014
959		1015
960	Aaron P. J. Roberts, Leonie V. Webster, Paul M. Salmon, Rhona Flin, Eduardo Salas, Nancy J. Cooke, Gemma J. M. Read, and Neville A. Stanton. 2021. <i>State of science: Models and methods for understanding and enhancing teams and teamwork in complex sociotechnical systems</i> . <i>Ergonomics</i> , pages 1–45.	1016
961		1017
962		
963		1018
964		1019
965		1020
966	F. Eric Robinson, Lt Col Sarah Huffman, Lt Col Daniel Bevington, DeAnne French, Clayton Rothwell, LTC Christopher Stucky, Marissa Tharp, and Ashton Hughies. 2023a. <i>Team coordination style is an adaptive, emergent property of interactions between critical care air transport team personnel</i> . <i>Air medical journal</i> , 42(3):174–183. ObjectType-Article-1.	1021
967		1022
968		1023
969		
970		1024
971		1025
972		1026
973	Frank E Robinson, David Grimm, Dain Horning, Jamie C Gorman, Jennifer Winner, and Christopher Wiese. 2023b. <i>Using natural language processing to develop an automated measure of closed loop communication among critical care air transport teams</i> .	1027
974		1028
975		
976		1029
977		1030
978	Alexandra Rosser, Sarah Sullivan, Ryan Thompson, and Hee Soo Jung. 2019. <i>1774: Automated natural language processing of closed-loop communication in trauma resuscitations</i> . <i>Critical care medicine</i> , 47(1 Suppl 1):860–860.	1031
979		1032
980		1033
981		
982		1034
983	Seyed Omid Sadjadi, Sriram Ganapathy, and Jason W Pelecanos. 2016. Speaker age estimation on conversational telephone speech using senone posterior based i-vectors. In <i>2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 5040–5044. IEEE.	1035
984		1036
985		1037
986		1038
987		
988		1039
989	Eduardo Salas, Dana E. Sims, and C. Shawn Burke. 2005. <i>Is there a “big five” in teamwork?</i> <i>Small Group Research</i> , 36(5):555–599.	1040
990		1041
991		1042
992	Irim Salik and John V. Ashurst. 2022. <i>Closed Loop Communication Training in Medical Simulation</i> . Stat-Pearls Publishing, Treasure Island (FL).	1043
993		1044
994		1045
995	Björn Schuller, Stefan Steidl, Anton Batliner, Alessandro Vinciarelli, Klaus Scherer, Fabien Ringeval, Mohamed Chetouani, Felix Weninger, Florian Eyben, Erik Marchi, et al. 2013. The INTERSPEECH 2013 computational paralinguistics challenge: Social signals, conflict, emotion, autism. In <i>Proceedings INTERSPEECH 2013, 14th Annual Conference of the International Speech Communication Association, Lyon, France</i> .	1046
996		1047
997		1048
998		1049
999		1050
1000		1051
1001		1052
1002		1053
1003		1054
1004	Elizabeth Shriberg, Rajdip Dhillon, Sonali Bhagat, Jeremy Ang, and Hannah Carvey. 2004. <i>The ICSI meeting recorder dialog act (MRDA) corpus</i> . In <i>Proceedings of the SIGDIAL 2004 Workshop, The 5th Annual Meeting of the Special Interest Group on Discourse and Dialogue, April 30- May 1, 2004, Cambridge, Massachusetts, USA</i> , pages 97–100. The Association for Computer Linguistics.	1055
1005		1056
1006		1057
1007		1058
1008		1059
1009		1060
1010		1061
1011		1062
	Andreas Weise, Matthew McNeill, and Rivka Levitan. 2022. <i>The Brooklyn Multi-Interaction Corpus for Analyzing Variation in Entrainment Behavior</i> . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 1721–1731, Marseille, France. European Language Resources Association.	
	Jennifer Winner, Jayde King, Jamie Gorman, and David Grimm. 2022. <i>Team coordination dynamics measurement in enroute care training: Defining requirements and usability study</i> . <i>Proceedings of the International Symposium on Human Factors and Ergonomics in Health Care</i> , 11(1):21–25.	
	Martin Wöllmer, Felix Weninger, Tobias Knaup, Björn Schuller, Congkai Sun, Kenji Sagae, and Louis-Philippe Morency. 2013. Youtube movie reviews: Sentiment analysis in an audio-visual context. <i>IEEE Intelligent Systems</i> , 28(3):46–53.	
	Mingzhi Yu, Diane Litman, and Susannah Paletz. 2019. Investigating the relationship between multi-party linguistic entrainment, team characteristics and the perception of team social outcomes. In <i>The Thirty-Second International Flairs Conference</i> .	
	A Introduction	
	In these appendices we provide additional details on the dataset, the model training, and the annotation procedures (Appendix G , Appendix H , Appendix I , Appendix J , Appendix K).	
	B Data Statement	
	B.1 Curation Rationale	
	The ASIST Study 3 dataset contains data from eight experimental conditions: (i) teams with no advisor, (ii) teams with human advisors, and (iii) teams with one of six AI advisors (i.e., six conditions). Of these, we opted to exclude trials with human advisors for two reasons: (i) unlike with the actual study participants, we did not have source-separated audio streams for the human advisors, who were experimental confederates, and (ii) we believed that there would be some level of phonetic entrainment between the participants in the ‘human-advisor’ condition and their human advisor, which would introduce an additional confounding variable into our analysis of phonetic entrainment. For the trials involving AI advisors, we sampled trials relatively equally across all six AI advisors. We sampled at the team level, so sampling an additional team for a given AI advisor results in two additional trials for that AI advisor (since each team completes two Minecraft missions).	
	We exclude trials that were for the purpose of training participants on how to perform the task.	

Advisor	# of Trials
None	31
ASI-CMURI-TA1	2
ASI-CRA-TA1	2
ASI-DOLL-TA1	2
ASI-SIFT-TA1	2
ASI-UAZ-TA1	2
ASI-USC-TA1	2

Table 7: Number of trials annotated for each advisor condition.

We disfavor—but do not completely exclude—trials with data quality issues (e.g. trials that are missing utterances due to technical issues with the audio capture setup). For trials in which the audio capture for one or more speakers failed due to technical issues, we were still able to annotate dialog acts, sentiment and emotion, but were unable to annotate for CLC events and entrainment.

B.2 Speaker Demographic

Speaker demographics are provided in Table 8.

B.3 Annotator Demographic

Annotator demographics are provided in Table 9.

B.4 Speech Situation, Recording Quality

The audio recordings were conducted as part of a remote experiment that took place in 2022. Spoken, synchronous participant dialog was captured using the participants’ own computers, often with background noises (which we try to annotate). The dialog was spontaneous, arising in the context of the collaborative virtual search-and-rescue task being performed by the participants. The intended audience for the speakers are their teammates that are performing the search-and-rescue task with them at the moment.

B.5 Database contents

The entirety of the MultiCAT dataset is provided through a single SQLite3 database (multicat.db in the supplementary material for the paper). The entity-relation diagram showing the structure of the database (tables, foreign key relationships, etc.) is shown in Figure 2.

C Items per class in MultiCAT

Tables 10, 11, 12, 13, 14, and 15 show the number of items per class in each task within MultiCAT. Note that some tasks allow multiple labels for a single utterance, so the number of items for a particular

class in a task do not add up to the number of utterances annotated for that task.

D Breakdown of annotations by team and trial

The breakdown of annotations in MultiCAT by team and trial are shown in Table 16. Different tasks had different goals and different levels of complexity, so trials that were ideal for some were not always ideal for all annotation types. For entrainment detection annotation, teams with two missions composed of clear audio files were selected. For sentiment and emotion annotation, extra trials were selected with the goal of increasing examples of small minority classes.

E Model training details

Below are the details of parameters, computational resources used and specifics of our training procedures for our baseline models.

E.1 DA classification

The training, validation, and test splits we used are shown in Table 17. We use version 1.13.1+cu117 of the PyTorch library (Paszke et al., 2019). The learning rate is set to 10^{-4} . The AdamW optimizer (Loshchilov and Hutter, 2019) is used with a decay of 10^{-5} . We train for a maximum of 100 epochs with early stopping after no improvement on the validation set for 10 epochs. The model has around 127M parameters, and takes ≈ 23 minutes to train. All experiments are performed on a single NVIDIA RTX A6000 GPU.

E.2 CLC detection

For the logistic regression model, we use as the training set the following 25 trials: T000603, T000604, T000607, T000608, T000613, T000627, T000628, T000631, T000632, T000633, T000634, T000635, T000636, T000637, T000638, T000713, T000714, T000715, T000716, T000719, T000720, T000723, T000724, T000729, T000730.

For the check-back detection step, we used the following 20 trials as the training set: T000603, T000604, T000627, T000628, T000631, T000632, T000635, T000636, T000637, T000638, T000713, T000714, T000715, T000716, T000719, T000720, T000723, T000724, T000729, T000730, and the following 5 trials as the validation set: T000607, T000608, T000613, T000633, T000634.

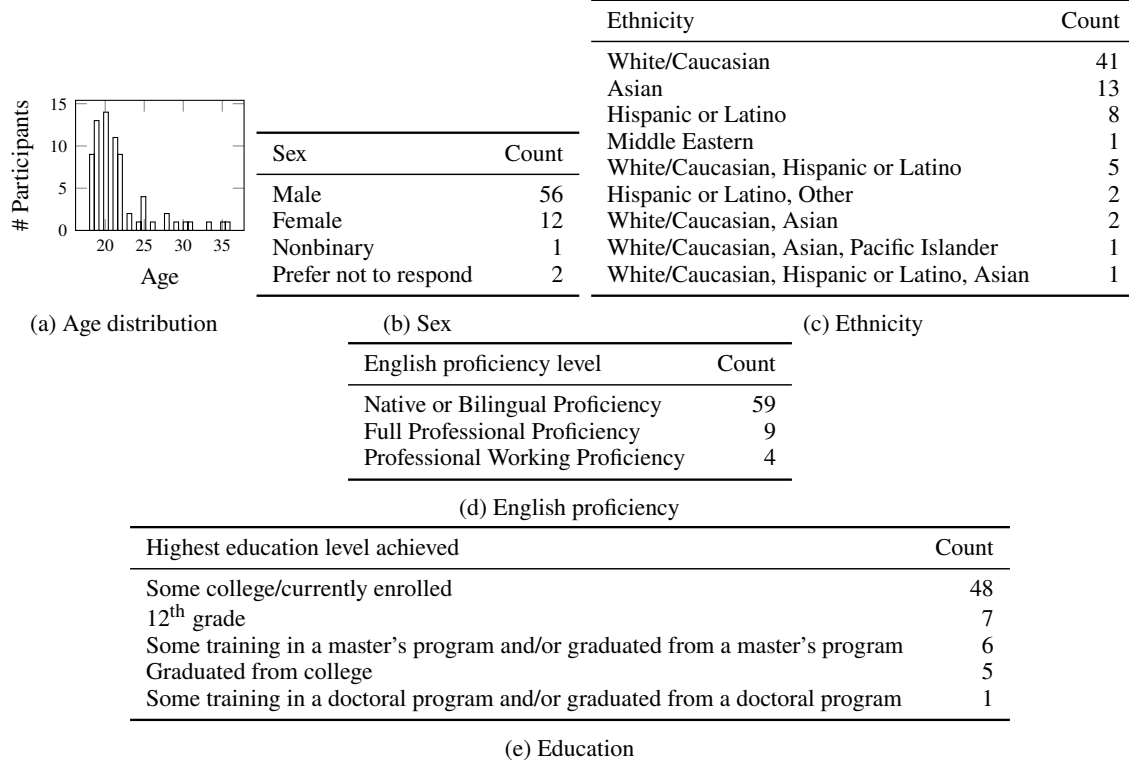


Table 8: Aggregated speaker demographic data for selected dimensions.

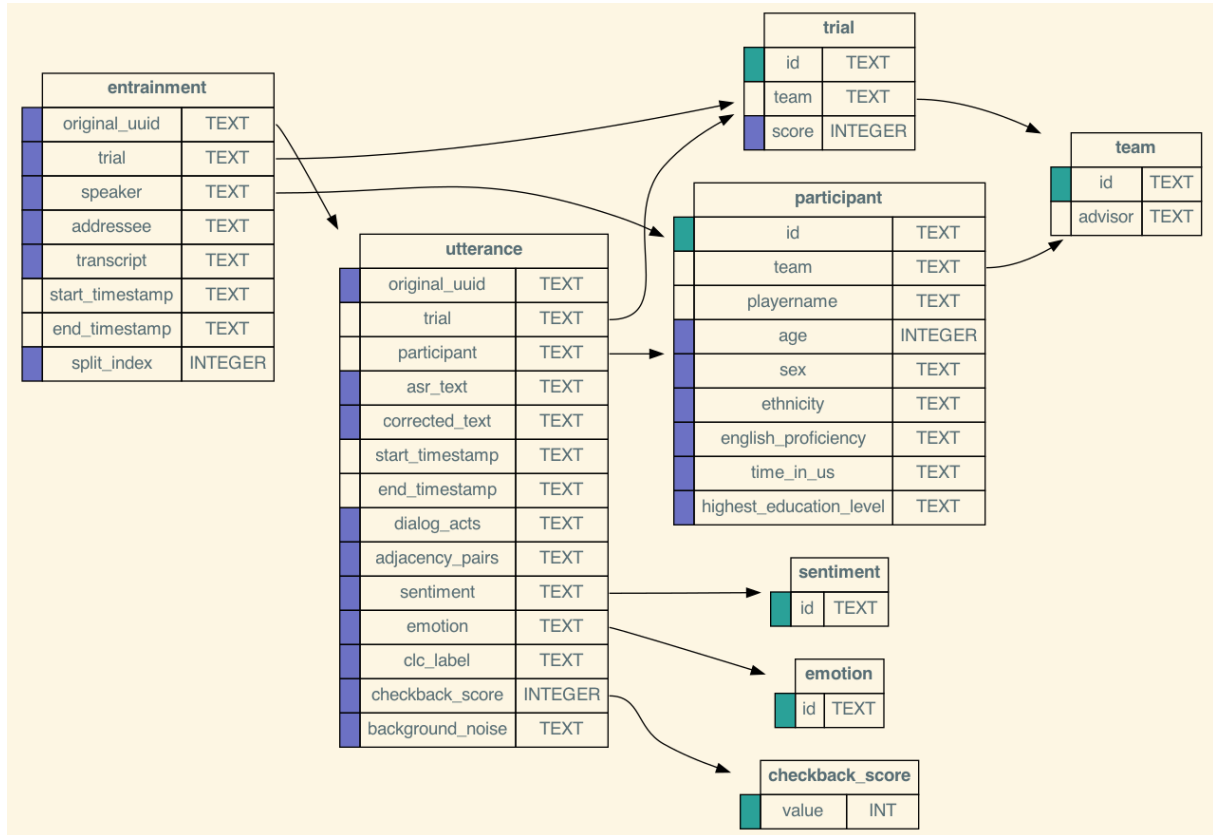


Figure 2: Entity-relation diagram for the MultiCAT database.

Specification	Value
Age	23–33 years
Gender	Female (3), Male (2)
Race/ethnicity	East Asian (1), South Asian (2), Middle Eastern (1)
Native language	Korean (1), Tamil (1), Hindi (1), English (2), Persian (1), Sindhu/Urdu (1)
Socioeconomic status	Middle class (4), upper middle class (1)
Training	Master’s degree in linguistics (2), PhD student in linguistics (1), No direct linguistics training, but work in NLP (1), Took linguistics courses (1)

Table 9: Annotator demographics

The detection of the call-out step with the logistic regression model takes 0.1 second to train.

We adopted the Transformer-based RoBERTa-base model for the detection of the check-back step. The learning rate is set to 5×10^{-5} , the model is trained with a batch size of 16 for 3 epochs. This model takes approximately 30 minutes to train.

The CLC detection experiments are performed on a Apple M1 CPU.

E.3 Sentiment and emotion classification

We train our sentiment and emotion baseline on a high performance computing environment with a Tesla V100S-PCIE-32GB GPU. For the sentiment and emotion classification tasks, we use the same training, validation, and test splits as in Table 17, except for including an additional trial (T000614) in the validation split.

We train this model using version 2.2.0+cu121 of PyTorch. Our baseline model contains 1,904,690 parameters. Our best hyperparameter settings are a learning rate of 10^{-3} with an Adam optimizer with a weight decay of 10^{-4} .

We perform a limited grid search over our pre-training corpora, then fine-tune with MultiCAT data on the best of these. The model takes approximately 15 minutes to train and 2 minutes for fine-tuning.

E.4 Entrainment identification

We train our entrainment model using PyTorch version 1.9.0+cu111 with torchaudio version 0.7.0⁵ and NumPy version 1.22.4 (Harris et al., 2020). This is done on an NVIDIA A100-PCIE-40GB GPU. We use the same hyperparameters identified as best in Nasir et al. (2020). Training of the model

Class	Count
2	19
%	92
%-	123
%—	125
aa	1858
aap	10
am	14
ar	58
arp	1
b	39
ba	227
bc	6
bd	17
br	46
bs	17
bsc	94
bu	113
cc	1201
co	889
cs	251
d	206
df	233
e	449
fa	121
fe	152
ft	140
fw	1
g	58
j	44
m	136
na	263
nd	45
ng	32
no	43
qo	9
qr	52
qw	308
qy	808
r	44
s	6033
t1	141
x	116
z	264

Table 10: Items per class for DA classification

⁵<https://pytorch.org/audio/stable/index.html>

Class	Count
Neutral	4081
Positive	2436
Negative	1214

Table 11: Items per class for sentiment analysis

Class	Count
Neutral	5977
Joy	571
Sadness	452
Fear	319
Surprise	280
Anger	66
Disgust	66

Table 12: Items per class for emotion prediction.

Class	Count
a	4115
b	4473

Table 13: Items per class for adjacency pair identification.

Class	Count
a	3671
b	2767
b	386

Table 14: Items per class for CLC detection.

Class	Count
Addressee	2896

Table 15: Items per class for entrainment detection.

on the Fisher corpus took an average of 70 minutes, and testing on MultiCAT takes 3.5 minutes (for all 30 iterations).

F Software

The code used to generate the database and the results in the paper will be added to the supplementary material for the camera-ready version upon paper acceptance.

G ASR transcript correction guidelines

Basic Setup The data should be in CSV format with one column for ASR and one column of corrected transcripts. The annotator is expected to listen to the full audio and read the ASR transcripts, whenever there are any discrepancies, those should be corrected and entered only in the corrected transcripts column.

Segmentation The segmentation of speaker utterances as done by ASR is not to be changed. For example, even if the annotator feels utterance B should come before utterance A, they should not change the order of the utterances.

Missing Utterances At times the ASR fails to pick up on small utterances, especially those that are just a few words long. In that case, a new row should be inserted in the CSV file and the text of the utterance should be manually entered. The field for the ASR transcript should be left empty. The annotator should also enter the speaker name and start and end timestamps.

Relative Order of New Utterance The utterance should be inserted based on the start timestamp and its relative order with the already present utterances.

Noise Picked up by ASR When ASR picks up noise as an utterance, a special character of hyphen "-" should be added as the corrected transcript.

H DA annotation guidelines

H.1 MRDA Framework

Our annotations follow the same guidelines as that of the ICSI MRDA corpus. The manual for MRDA contains detailed examples and definitions of different tags. This manual further builds on the MRDA manual (Dhillon et al., 2004) and addresses special cases we encountered when annotating MultiCAT.

Team	Trial	SentEmo	CLC	DA	AP	Entrainment
TM000201	T000602	✓				
TM000202	T000603	✓	✓	✓	✓	✓
TM000202	T000604	✓	✓	✓	✓	✓
TM000203	T000605	✓	✓	✓	✓	✓
TM000203	T000606	✓	✓	✓	✓	✓
TM000204	T000607	✓	✓	✓	✓	
TM000204	T000608	✓	✓	✓	✓	
TM000205	T000609	✓		✓	✓	
TM000205	T000610	✓		✓	✓	
TM000206	T000611	✓		✓	✓	
TM000206	T000612	✓		✓	✓	
TM000207	T000613	✓	✓	✓	✓	
TM000207	T000614	✓				
TM000210	T000619	✓				
TM000210	T000620	✓		✓	✓	
TM000211	T000621	✓				
TM000211	T000622	✓		✓	✓	
TM000212	T000623	✓	✓	✓	✓	
TM000212	T000624	✓		✓	✓	
TM000213	T000625	✓		✓	✓	
TM000213	T000626	✓		✓	✓	
TM000214	T000627		✓	✓	✓	
TM000214	T000628	✓	✓	✓	✓	
TM000216	T000631	✓	✓	✓	✓	
TM000216	T000632	✓	✓	✓	✓	
TM000217	T000633	✓	✓	✓	✓	
TM000217	T000634	✓	✓	✓	✓	
TM000218	T000635	✓	✓	✓	✓	
TM000218	T000636	✓	✓	✓	✓	
TM000219	T000637	✓	✓	✓	✓	
TM000219	T000638	✓	✓	✓	✓	
TM000236	T000671	✓	✓	✓	✓	
TM000236	T000672	✓		✓	✓	
TM000252	T000703		✓	✓	✓	
TM000252	T000704		✓	✓		
TM000257	T000713	✓	✓	✓	✓	
TM000257	T000714	✓	✓	✓	✓	
TM000258	T000715	✓	✓	✓	✓	
TM000258	T000716	✓	✓	✓	✓	
TM000260	T000719	✓	✓	✓	✓	✓
TM000260	T000720	✓	✓	✓	✓	✓
TM000262	T000723	✓	✓	✓	✓	✓
TM000262	T000724	✓	✓	✓	✓	✓
TM000264	T000727	✓	✓	✓	✓	
TM000264	T000728	✓	✓	✓	✓	
TM000265	T000729	✓	✓	✓	✓	
TM000265	T000730	✓	✓	✓	✓	
TM000269	T000737	✓	✓	✓	✓	
TM000269	T000738	✓	✓	✓	✓	

Table 16: A list of all trials with the team that trial represents indicating which types of annotation each trial contains.

Split	# of trials	Trial IDs
Train	28	T000603, T000604, T000611, T000612, T000620, T000622, T000623, T000624, T000627, T000628, T000631, T000632, T000635, T000636, T000637, T000638, T000703, T000704, T000713, T000714, T000715, T000716, T000719, T000720, T000723, T000724, T000729, T000730
Validation	5	T000613, T000607, T000608, T000633, T000634
Test	12	T000605, T000606, T000671, T000672, T000625, T000626, T000727, T000728, T000737, T000738, T000609, T000610

Table 17: Train, validation, and test split composition for the DA classification and AP detection tasks.

H.2 Questions

Discontinuous Question When speaker A asks a question but they get interrupted by speaker B. after the interruption, speaker A goes on to finish the question. Two scenarios can arise.

- Speaker B answered the question, in this case the subsequent utterances by speaker A would be marked with statement general tag and elaboration specific tag. Since speaker A's intent behind the latter utterances is not to elicit an answer. Check page 34 of MRDA manual for a similar use case.
- Speaker B does not answer the question, the rest of speaker A utterances completing the question would get the same question tag(s).

H.3 Segmentation with Pipe

Floor Mechanisms (FM) <fg>, <fh>, <h> at the start or end of an utterance can be ignored. No need to pipe separate an utterance or include the FM tag in the label.

Short Response For tags <aa> and <ar> at the start or end of an utterance, make the response tag as part of a single combined utterance tag. That is, the general tag will be shared by the whole utterance.

Different General Tags with Pipe Pipe should be used for cases where segments of the utterance require different tags and cannot be merged into one label because of different general tags. The pipe would then be added to both the utterance and the label.

Utterance	DA
Oh you do? So you probably discard	qh s^cs

Table 18: An example illustrating the use of pipe bar to annotate an utterance for multiple general tags.

H.4 Acknowledgment <bk> & Accept <aa>

<bk> and <aa> tags have been merged into a single tag - <aa>.

H.5 <df> and <e> for a Single Utterance

The tag <df> can be assigned to a single utterance without having to associate it with a previous utterance. The same is not true for <e>. <e> tag can only be assigned in relation to some previous utterance.

Special case of <df> and <e> in same utterance

If an utterance were to be segmented to assign <df> tag while some portion has already been assigned the <e> tag, the <df> and <e> tags can be merged under the same general tag (if after pipe <df> was to receive the same general tag as well)

Speaker	Utterance	DA
A	So yeah I would move.	s^cs
B	Um.	h
A	down to Breaker's Bridge and shore it up, cause I don't think there's anything we can do.	s^df^e

Table 19: <df> and <e> can occur in the same utterance but <e> still has to be in relation to a prior utterance of the same speaker.

H.6 Commitment <cc> in Present Actions

In MultiCAT data, players often verbalize the action they are carrying out at the present moment, any such actions should also be considered as <cc>.

Utterance	DA
yep on my way.	s^aa^cc

Table 20: <cc> for present actions.

I Sentiment/emotion annotation guidelines

One task to complete during this summer's annotation effort is the annotation of utterances for sentiment and emotion. This document discusses the method that should be used when annotating each.

I.1 Key terminology

I.1.1 Utterance

For purposes of this task, we define the term **utterance** as a single unit transcribed by Google's ASR. In some cases, this will correspond to a single sentence without a pause; in others, this may actually be composed of more than one sentence. Occasionally, a single sentence is even split into two utterances by the ASR.

I.1.2 Emotion

Emotion in this task refers to the discrete emotion shown by a speaker during an utterance. The emotion is selected from the set of labels described in section 3 below.

I.1.3 Sentiment

Sentiment in this task refers to the feelings a speaker shows towards the topic of an utterance. The sentiment may be positive, negative, or neutral. Sentiment labeling is discussed in section 4 below.

I.2 Basic annotation procedure

You will be asked to make your annotations using spreadsheets and while accessing the full audio files for a mission. Below is the annotation procedure that we will be following.

I.2.1 Materials needed

To complete this annotation task, you will need a spreadsheet containing each of the corrected/uncorrected utterances (which should be provided to you) with empty columns where you will enter your annotation labels, as well as the corresponding audio files.

You should select a quiet place to work and use headphones to ensure that you can clearly hear the entirety of the audio.

I.2.2 Procedure

For this task, you should have the transcript and label spreadsheet open while listening to the audio. If you cannot look at the transcript and listen to the audio at the same time, you should read the

transcript for each single utterance immediately before listening to that utterance.

For the sake of consistency, we will be using **uncorrected** transcripts for this task. This means that the words may not form a logical sentence, and at times may be difficult to understand. When this happens, do your best to pay attention to the words in the recording (as these should make sense) and use these to help inform your decisions.

You will need to download the transcripts and the relevant audio files from kraken. The transcripts may be found in the following location: /home/tomcat/annotations/transcriptions. The audio files may be found in: /home/tomcat/annotations/wav. Some of these transcript files may contain corrected transcripts; however, you should focus on the uncorrected transcripts (the column labeled 'utt' or 'utterance').

Select a transcription and the corresponding audio; open the transcription to take up at least half of your screen, ensuring that you can see the entirety of each transcribed utterance that is within the window.

After listening to a single utterance, pause the recording, then enter the emotion label and the sentiment label into the corresponding cells in the spreadsheet. You may then play the recording again and examine the next utterance.

I.3 Emotion task

The first of the two annotations that you will be completing as you go through the files is the emotion task. For this task, you will need to decide which of a set of emotions is the best label for each individual utterance, as defined above. The set of labels used in this task and examples of annotations for each appear below.

I.3.1 Emotion labels

While there are several methods for capturing emotional information from audio, we are using a set composed of Ekman's universal emotions + a neutral label. This label set is:

1. **anger**: the speaker is angry, upset, and reveals this through words, tone or both.
2. **disgust**: the speaker is disgusted; in this dataset, disgust frequently appears when a player walks into the same trap room more than once, when someone is having a little bit of trouble with the controls, or when any sort

1365	of glitch occurs. This emotion label is more	• If one emotion seems much stronger than	1410
1366	like frustration than anger.	the other, choose the stronger emotion	1411
1367	3. fear : the speaker is afraid of something.	• If one emotion dominates the utterance,	1412
1368		choose the dominant emotion	1413
1369	4. joy : the speaker is happy, having a good time,	• otherwise (assuming equal parts of each	1414
1370	or otherwise enjoying something. This emo-	of two emotions):	1415
1371	tion frequently occurs at the end of missions	(a) If one emotion is fear and the other	1416
1372	immediately after time has run out, though	is anything else, choose fear	1417
1373	some speakers show moments of joy through-	(b) If one emotion is sadness and the	1418
	out the mission.	other is anything but fear, choose sad-	1419
1374	5. neutral : (no clear emotion)—the speaker	ness	1420
1375	doesn't demonstrate any emotions; they may	(c) If one emotion is anger and the other	1421
1376	be explaining something or providing infor-	is not fear or sadness, choose anger	1422
1377	mation about their movements to their team.	(d) If one emotion is disgust and the other	1423
1378	This sort of neutral language is very common	is joy or surprise, choose disgust	1424
1379	in the ASIST data.	(e) If one emotion is joy and the other is	1425
		surprise, choose joy	1426
1380	6. sadness : the speaker is sad or disappointed, of-	• If there are ever three emotions in one	1427
1381	ten because something has happened that they	utterance, follow the points above to make	1428
1382	did not want to have happen (like repeatedly	your decision about which to select	1429
1383	entering a trap room), or because something		
1384	hasn't happened that they wanted to see hap-		
1385	pen (e.g. the number of victims saved is lower		
1386	than they had hoped).		
1387	7. surprise : something surprising has happened,		
1388	the speaker is suddenly given new unexpected		
1389	information or corrected about something they		
1390	thought they knew but that turned out to be		
1391	incorrect.		
1392	Each utterance should be given a single label.		
1393	This label may be based on the words that the		
1394	participant produces, the way in which they speak,		
1395	or both.		
1396	I.3.2 How to decide which emotion label to		
1397	select		
1398	Determining which label to use is often straight-		
1399	forward; sometimes, however, you may not be sure		
1400	of which label to assign an utterance. In general,		
1401	follow these rules:		
1402	1. If an utterance contains no obvious emotional		
1403	information, give it a label of neutral		
1404	2. If most of an utterance contains no obvious		
1405	emotional information, but one part of it does		
1406	contain emotion, provide the label of the non-		
1407	neutral emotion demonstrated		
1408	3. If an utterance contains two emotions, do the		
1409	following:		
		I.3.3 Examples of emotion annotations	1430
		"Okay can you make sure you mark it?" Said with	1431
		a neutral tone, this would be given the label neutral.	1432
		The speaker is making a request of another player.	1433
		"Oh shoot that's the wrong one" The participant	1434
		suddenly realized they have gone to the wrong	1435
		location. This should be given the label surprise.	1436
		"and then wacky fun little update guys both of	1437
		our C zones are blocked right now" While the ASR	1438
		transcription isn't perfectly accurate, this speaker is	1439
		indicating that they are stuck in a room. With the	1440
		intonation from the audio, we can tell that 'wacky	1441
		fun little update' is sarcastic, so this utterance	1442
		should be given the label disgust.	1443
		"shit" This speaker just shouted this word out,	1444
		showing that they were feeling mad, this would be	1445
		given the label anger.	1446
		"guys I'm starting to think we're not going to get	1447
		everyone" This speaker is disappointed that their	1448
		performance is not as good as the team had hoped.	1449
		This would be given the label sadness.	1450
		"I was like 3 seconds away oh I died" At the end	1451
		of the game, the speaker has not managed to save the	1452
		last victim they were carrying. Then the game ends	1453
		by showing the speaker's character dying. Without	1454
		the audio, it may seem as though this person is	1455
		disgusted, angry, or surprised, but they are in fact	1456
		laughing and having fun, while being surprised by	1457
		the event. This could have been labeled either joy	1458

or surprise, so following the guidelines above, we select label joy.

“Ah, what’s happening?” The mission has ended and the screen has suddenly changed, but the speaker thinks they have done something wrong somehow. They show both surprise and fear, so using the guidelines above, we select the label fear.

“oh geez now she’s been a red turn its meeting throws a 720” While the ASR is not quite right, this person is annoyed at an aspect of the mission that they have no control over (their speed). This could show surprise, disgust, or anger, so using the guidelines above we select anger.

I.4 Sentiment

The second annotation task that you will complete while going through these files is sentiment annotation. For this task, you will assign each item a sentiment label according to the sentiment expressed in the statement. For this task, as with the above, you will want to pay attention to both what is said and how it is said.

I.4.1 Sentiment labels

Sentiment: the content/meaning of each utterance should be marked as one of the following.

1. **positive:** the utterance refers to a subject that the speaker feels positively about.
2. **neutral:** the utterance does not reveal positive or negative sentiment; this is generally the case with instructions, updates, descriptions of players’ movements and when speakers provide general information.
3. **negative:** the utterance refers to a subject that the speaker feels negatively about.

I.4.2 How to decide which sentiment label to select

Because there are only three sentiment labels to select from, it is much less likely that you will have to make difficult decisions about which to choose.

1. If there is no indication of either positive or negative sentiment, choose the neutral label
2. If any part of the utterance demonstrates positive or negative sentiment, select that sentiment, even if the majority of the utterance is neutral

3. If both positive and negative sentiment are shown in equal amounts in the same utterance, select the negative label
4. Politeness does not convey any information other than politeness. Thus, select neutral label
5. ‘Okay’ should be labeled depending on tone and pitch
 - negative: sarcasm, annoying situation
 - neutral: gap filler
 - positive: other than the aforementioned

There is a correlation between sentiment labels and emotion labels (e.g. ‘happy’ utterances would tend to also have a positive sentiment), although there is not an exact mapping of sentiments onto emotions (e.g. ‘surprise’ could be positive or negative). The vast majority of the utterances seem to be neutral in both emotion and sentiment, and that’s okay. One of the recordings I listened to only had one utterance that showed a non-neutral emotion/sentiment value (the last utterance, actually).

Sometimes, however, the emotion a participant shows is NOT the same as the sentiment they express. For example, sometimes someone expresses joy through their tone, but the words they are saying actually indicate a negative sentiment (e.g. they are having fun playing the game, but they say ‘We did really poorly this round!’).

I.4.3 Examples of sentiment annotations

“It might actually be best to start in the middle and then work our way either left or right because the middle is where we spawn” This speaker is giving suggestions on what they think is the correct way to organize their movements during a mission that is just starting. They are neutral in their tone. This should be labeled neutral.

“Okay engineer to enter so critical in here yeah” The ASR has not given an accurate transcription here, but we can see that most of the words themselves seem neutral. However, with the speaker’s tone, we see that they feel positively about the event taking place at the end (where a critical victim is found), so this would be labeled positive.

“Other that sorry that’s the one you know it’s not okay so we got that b there’s two critical Zone here speak out that one but” The ASR is again not quite accurate, but we can see that this person does not seem to feel positively about the room that

1551	they have just entered. Using this knowledge, plus	• Automatically filled textgrids (one per audio	1594
1552	phrases like ‘sorry’ and ‘it’s not okay’, this would	file) with two tiers, ‘silences’ and ‘addressee’.	1595
1553	be labeled as negative.	The ‘silences’ tier will have two types of au-	1596
		tomatically detected labels: ‘silence’ (which	1597
1554	J Entrainment annotation guidelines	is the label for non-speech sounds as well as	1598
		silences), and ‘sound’ (for speech).	1599
1555	In this annotation task, we search for the intended	You should select a quiet place to work and	1600
1556	listener of a given spoken unit. You task is to listen	use headphones to ensure that you can clearly	1601
1557	to the audio, read the transcripts for every utterance	hear the entirety of the audio.	1602
1558	in the recording, find the inter-pausal units within		
1559	each utterance, and ascertain who the inter-pausal	J.2.2 Procedure	1603
1560	unit is aimed at.	For this task, keep the transcript open on any spread-	1604
		sheet reader, along with the audio and Praat textgrid	1605
1561	J.1 Key terminology	open on Praat.	1606
1562	J.1.1 Utterance		
1563	A section of the spoken interaction that the auto-	1. Download the transcripts, textgrids and the	1607
1564	matic transcription service has detected as a unit of	relevant audio files from kraken. The tran-	1608
1565	speech.	scripts may be found in the following location:	1609
		‘/home/tomcat/annotations/transcriptions’,	1610
1566	J.1.2 Vocal Entrainment	and the audio and textgrids in	1611
1567	Vocal Entrainment is the shift in vocalic features	‘/home/tomcat/annotations/wav’.	1612
1568	(such as fundamental frequency) of a speaker in		
1569	order to resemble their conversation partner.	2. On Praat, move your cursor to the first chunk	1613
		where the experiment participant is speaking.	1614
1570	J.1.3 Inter-pausal Unit (IPU)		
1571	A stream of audio separated by a pause of 50ms or	3. Listen until you hear the speaker pausing, and	1615
1572	more. This can be a whole or part of an utterance.	check if the pause is over 50 ms. You can	1616
		see the length of the selected audio above the	1617
1573	J.2 Basic annotation procedure	waveform, or by clicking on ‘Query’ > ‘Get	1618
1574	For this task, you will be working to assess and	length of selection’ in the menu on the top	1619
1575	correct the IPU boundaries on a automatically filled	left corner of the screen. If the pause is less	1620
1576	Praat textgrid. For each IPU you correct and finalize,	than 50 ms, continue listening until you hear	1621
1577	you will add the corresponding transcription in	a pause.	1622
1578	the ‘silences’ tier from the transcript spreadsheet		
1579	provided. Finally, you will identify the intended	4. If you see a longer pause, make sure the start	1623
1580	addressee of every IPUs and annotate for it in	and end of the speech has boundaries on both	1624
1581	the ‘addressee’ tier. Your final submission is a	the ‘silences’ and ‘addressee’ tiers. Drag the	1625
1582	corrected textgrid with labels in the ‘silences’ and	boundaries until they enclose the speech and	1626
1583	the ‘addressee’ tiers.	move them as close to the speech chunk as	1627
1584	You will be asked to make your annotations	possible.	1628
1585	using spreadsheets and the audio files from the		
1586	individual recording channels for each player in	5. Ensure that the silences on each side of the	1629
1587	given a mission. The procedure is outlines in the	speech chunk have the automatically generated	1630
1588	‘Procedure’ section below.	label ‘silence’.	1631
1589	J.2.1 Materials and technology needed		
1590	• Praat software.	6. From the spreadsheet, copy and paste the	1632
1591		chunk of the transcript that matches the words	1633
1592	• The spreadsheet containing the corrected ut-	you hear into the ‘silences’ tier. These words	1634
	terances for a given trial.	may be just a portion of the utterance in the	1635
1593		cell. The rest may belong to the following	1636
	• The corresponding audio files.	IPU.	1637
		7. Identify the addressee of the IPU. You can	1638
		determine this from the context of the conver-	1639
		sation. For example, the speaker could have	1640

1641	called out to a specific player. Or the IPU	(or zooming out, as seen in Figure 5 on the textgrid,	1689
1642	could be part of an answer to a question asked	you can see that both the previous utterances did	1690
1643	in a previous utterance.	not have a specific addressee (thus labelled ‘all’).	1691
1644	8. Add an addressee label in the ‘addressee’ tier.	Based on the context, we will mark this IPU as ‘all’	1692
1645	You have four options. If you identify a dis-	in the ‘addressee’ tier on the textgrid.	1693
1646	distinct addressee, annotate with the name of	This completed the annotation task for this IPU,	1694
1647	any one Minecraft roles played by the players	and we can scroll to the next one.	1695
1648	(‘engineer’, ‘transporter’, ‘medic’).		
1649	9. Or, if you can’t identify a specific addressee,	K CLC annotation guidelines	1696
1650	or if the IPU is directed at the experimenter,	This document discusses the method of annotating	1697
1651	simply mark it as ‘all’.	closed-loop communication events in multi-party	1698
1652		dialogues.	1699
1653	10. Continue scrolling through the IPU’s until you	K.1 Definition of Closed-Loop	1700
1654	have corrected, transcribed and addressee-	Communication	1701
1655	identified each IPU. Save your annotated	In team communication, especially in emergency	1702
1656	textgrid frequently.	situations, there’s a standard scheme of communi-	1703
1657	J.2.3 An example for IPU detection	cation, called Closed-loop communication. Closed-	1704
1658	Figure 3 has a Praat window open with the	loop communication aims to achieve safe commu-	1705
1659	waveform (top), spectrogram (middle), as well	nication by reducing the risk of miscommunication	1706
1660	as the textgrid (bottom) containing the auto-	and ensuring clear communication. Closed-loop	1707
1661	matically detected voice activity for the files	communication is usually trained and adopted in	1708
1662	‘HSRData_ClientAudio_Trial-T000719_Team-	high-stakes team environments like Crew Resource	1709
1663	TM000260_Member-E000888_CondBtwn-	Management, medical surgery teams, and emer-	1710
1664	ASI-UAZ-TA1_CondWin-na_Vers-1.wav’ and	gency departments. In our Minecraft games which	1711
1665	‘HSRData_ClientAudio_Trial-T000719_Team-	simulate the urban search and rescue scenario, the	1712
1666	TM000260_Member-E000888_CondBtwn-ASI-	appearance of Closed-loop communication is con-	1713
1667	UAZ-TA1_CondWin-na_Vers-1.TextGrid’. The	sidered a good approach to team communication,	1714
1668	view shows the audio divided into chunks of sound	although the participants of the game are not trained	1715
1669	and silence (labelled in the first tier). In reality, this	in doing so.	1716
1670	is one inter-pausal unit in which the consonants	Closed-loop communication includes three	1717
1671	have been incorrectly labelled as silences by the	phases:	1718
1672	automatic speech detector. Our first task is to	Call-out The sender initiates a message.	1719
1673	correct the IPU boundaries and add the transcript	Check-back The receiver acknowledges the mes-	1720
1674	corresponding to it.	sage, usually by paraphrasing or repeating the	1721
1675	First, we remove the unwanted boundaries and	main information of the message.	1722
1676	labels such that only the initial and final boundaries	Closing-of-the-loop The sender verifies that the	1723
1677	remain. Next, we adjust the start and end boundaries	message has been received and interpreted	1724
1678	until they enclose only speech. Finally, we add the	correctly.	1725
1679	text from the transcription spreadsheet. The end	Table 21 is an example of closed-loop communi-	1726
1680	result should look like Figure 4 .	cation. The detection of Closed-loop communica-	1727
1681	J.2.4 An example for addressee identification	tion will be triggered by recognizing the Call-out	1728
1682	Using the same IPU as the above section, we now	phase, and then searching for the Check-back phase,	1729
1683	move on to identifying the speaker and their ad-	and finally the closing-of-the-loop phase. There	1730
1684	ressee. First, we look in the transcript spreadsheet	might be situations where only a sender calls out	1731
1685	for utterances preceding the IPU of interest, and	but no one checks back to the sender, or there’re	1732
1686	who was the speaker. In the example, the utterances	call-out and check-back but no final acknowledg-	1733
1687	preceding ‘this is transporter there’s a critical vic-	ment to close the loop. We have different labels	1734
1688	tim in A4’ (‘this is’ and ‘three’) are also uttered by	for the three phases. Table 22 is a list of common	1735
	the same speaker (‘transporter’). By scrolling back	semantic types of the CLC phases.	1736

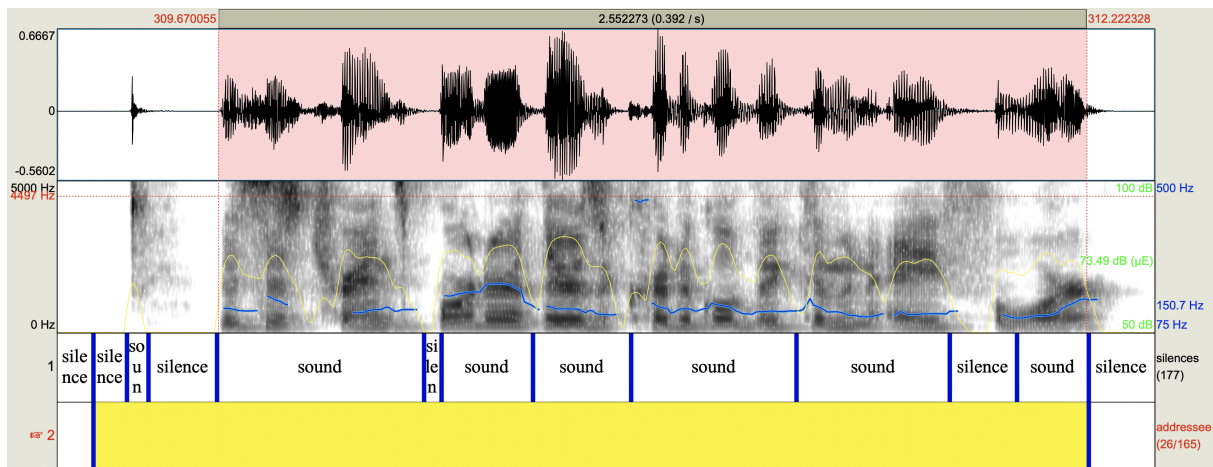


Figure 3: Original textgrid

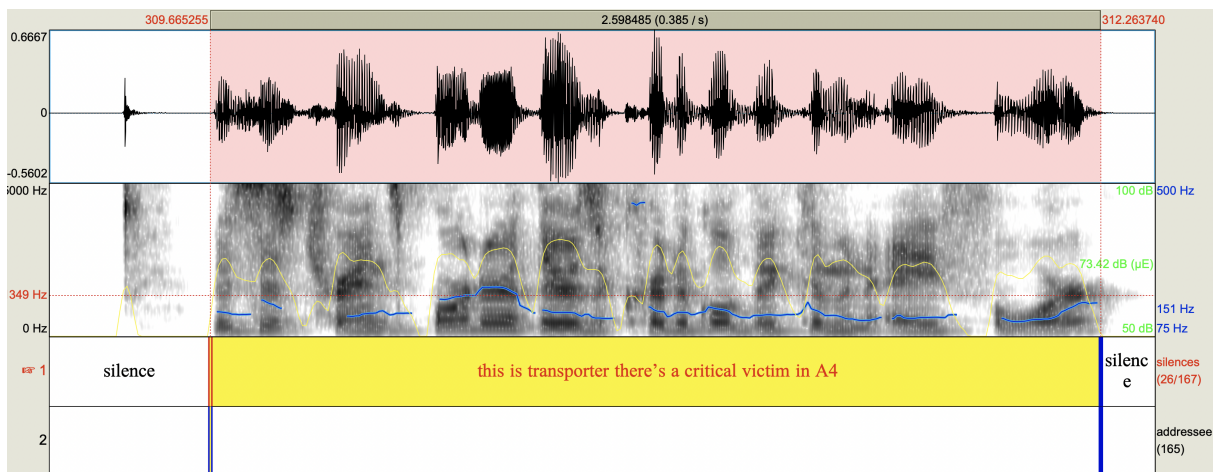


Figure 4: Textgrid with IPU boundaries and transcript

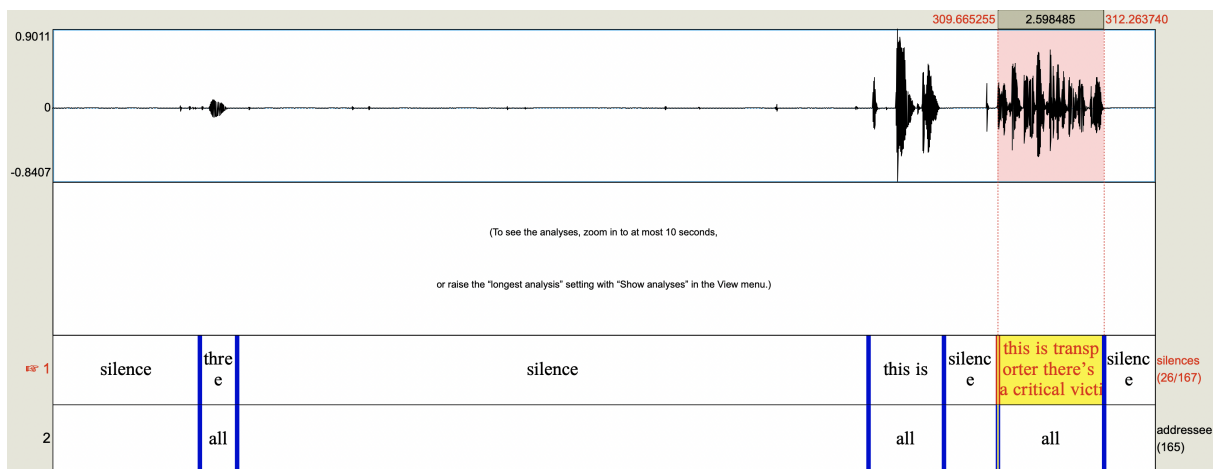


Figure 5: Textgrid with IPU boundaries and transcript

Role	Utterance	CLC Phase
Green	This is Green. I'm finishing this side, blue, could you check the central?	Call-out
Blue	This is Blue. I'll go check the central.	Check-back
Green	Thank you, Blue.	Closing-of-the-loop

Table 21: An example of the closed-loop communication

K.2 Labels and Scores

The transcripts of utterances are saved in CSV files. The annotations are in columns: CLC_Label, Checkback_Score.

At the beginning of each trial, there are several pre-game chatting utterances, which happen before players enter the scene and they were communicating with each other about team strategies. At the end of each trial, there're also several post-game utterances after the game session ends. We will not include those in our CLC annotation.

The three phases of the CLC are labeled with letters *a*, *b*, and *c*:

- Call-out: *a*
- Check-back: *b*
- Closing-of-the-loop: *c*

We follow the MRDA (Multi-Dimensional Annotation) framework for annotating adjacency pairs and adapt it to our CLC annotation with the format:

<CLC number><CLC phase>.<CLC number><CLC phase>-<nth speaker>[+...]

The <CLC number> is the index number of CLC events, which helps us keep linking call-outs and their follow-up check-backs and closing-of-the-loops, especially when they are several utterances away from the call-outs. The <CLC phase> are *a*, *b*, and *c* phases for each CLC event. The <nth speaker> is useful when there're multiple check-outs for one call-out, and the [+...] suffix is used to note a continued CLC phase from the same speaker, which usually happens when a sentence is cut off into more than one utterances. For example:

8a.9a indicates two call-out events in one utterance, see table 23.

a+/a++ indicates continued call-out events by the same speaker, see table 24.

b+/b++ indicates the same person check-back to one call-out event, see table 25.

b-1/b-2 indicates two check-backs from different speakers to one call-out, see table 26.

The three phases are not necessarily closely next to each other. There might be some other utterances that insert between call-out and check-back, and check-back and closing-of-the-loop.

In our scripts, sometimes, the time span of each utterance might overlap, and starting timestamp may not be ordered properly. We need to pay special attention to the timestamps in order to make sense of the flow of conversations.

The **Checkback_Score** measures the quality of the check-back phases. If the check-back utterance repeated the key information in the call-out utterance, and shows the full understanding of the call-out information with no ambiguity, then the check-back can get a score of 3. But if there's only an acknowledgment like "Okay" or "Alright" but no major information that could clear out the ambiguity, that check-back utterance can only receive a score of 1. If the check-out phase contains some part of the key information in the call-out phase but has some level of ambiguity, the check-back utterance can get a score of 2. Table 3 provides the rubric and example for evaluating the check-back score.

K.3 Example Cases

CLC Phase	Semantic Types
Call-out	request, question, action directive, instruction, commitment, assert, knowledge sharing
Check-back	[another player] acknowledgment, confirm, (key information in call-out)
Closing-of-the-loop	[call-out speaker] acknowledgment, confirm, gratitude

Table 22: Common semantic types of CLC phases

Role	Utterance	CLC_Label	Checkback_Score
Green	where's the management meeting and the transporter here	15a.16a	
	I'm going to go check in there		
Blue	okay	16b	1

Table 23: One sentence contains two events

Role	Utterance	CLC_Label	Checkback_Score
Red	transporter you at M1	42a	
Red	this is medic	42a+	
Green	this is transporter I am almost there	42b	2

Table 24: One sentence is cut off into several utterances

Role	Utterance	CLC_Label	Checkback_Score
Red	okay so E5 we should also be good	7a	
Blue	okay	7b	3
Blue	E5 looks good	7b+	3

Table 25: Two check-backs from one person for the same call-out. The scores should be the same for all "7b" labels because they are considered as one 7b event

Role	Utterance	CLC_Label	Checkback_Score
Red	yeah um can someone come with me to B2	30a	
Green	I'll be back there in a sec	30b-1	2
Blue	B2 yeah	30b-2	2

Table 26: Two check-backs for one call-out

Role	Utterance	CLC_Label	Checkback_Score
Red	I'm heading to A2 medic	12a	
Red	management meeting is in M3	13a	
Blue	B2 okay	12b.13b	1

Table 27: One check-back for two call-outs

Role	Utterance	CLC_Label	Checkback_Score
Green	this is transporter area c as in the hole is there a number associated or am I missing something	13a	
Blue	this is engineer I'm sorry I could not hear what you said could you repeat that for me please	13b	3
Green	B2 this is transporter you said that area C has Rubble	13c	
Green	oh Zone c i see	14a	
Blue	B2 yes on the south Zone C where the critical conditioner it got covered in rubble so I cleared it out I apologize	14b	3

Table 28: Follow-up questions for the call-out. The follow-up question is considered as a 3 scored *b*