

MEDIC: A Multimodal Empathy Dataset in Counseling

Zhou'an Zhu, Xin Li, Jicai Pan, Yufei Xiao, Yanan Chang, Feiyi Zheng, Shangfei Wang

ABSTRACT

Although empathic interaction between counselor and client is fundamental to success in the psychotherapeutic process, there are currently few datasets to aid a computational approach to empathy understanding. In this paper, we construct a multimodal empathy dataset collected from face-to-face psychological counseling sessions. The dataset consists of 771 video clips. We also propose three labels (i.e., expression of experience, emotional reaction, and cognitive reaction) to describe the degree of empathy between the counselors and their clients. Expression of experience describes whether the client has expressed experiences that can trigger empathy, and emotional and cognitive reactions indicate the counselor's empathic reactions. As an elementary assessment of the usability of the constructed multimodal empathy dataset, an interrater reliability analysis of annotators' subjective evaluations for video clips is conducted using the intraclass correlation coefficient and Fleiss' Kappa. Results prove that our data annotation is reliable. Furthermore, we conduct empathy prediction using three typical methods, including the tensor fusion network, the sentimental words aware fusion network, and a simple concatenation model. The experimental results show that empathy can be well predicted on our dataset. Our dataset is available for research purposes.

ACM Reference Format:

Zhou'an Zhu, Xin Li, Jicai Pan, Yufei Xiao, Yanan Chang, Feiyi Zheng, Shangfei Wang. 2018. MEDIC: A Multimodal Empathy Dataset in Counseling. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

In recent years, people have experienced more emotional distress from everyday life [16]. Thus, they need emotional support, and often psychological treatment to recover. The counselors provide professional mental health support via psychological counseling [19, 21]. Empathy is the basis of the humanistic approach to psychotherapy [13] and has long been used effectively. Empathy represents the therapist's ability and willingness to understand the patient's thoughts, feelings, and struggles from the patient's perspective [10]. The ability to empathize is important in promoting positive behavior towards others, and may be the mechanism that drives the desire to help others [30].

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

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

Although empathy has been studied extensively in the field of psychology [3, 4, 10–12, 31], the computational approach to understanding empathy during mental health support has been hindered by the lack of datasets. To the best of our knowledge, there is only one public dataset for empathy understanding [30], collected from text-based, asynchronous conversations on mental health platforms. Since psychological treatment frequently occurs in synchronous, face-to-face settings, a multimodal empathy dataset collected from face-to-face mental health support is urgently required.

In this paper, we build a multimodal empathy dataset (MEDIC) for face-to-face counseling scenarios. To describe empathic communication in multimodal scenarios, we also propose three labels: expression of experience (EE), emotional reaction (ER) and, cognitive reaction (CR). EE describes whether the client has expressed experiences triggering empathy. ER and CR indicate the counselor's empathic reactions on emotion and cognitive dimensions, respectively. The dataset is constructed from counseling case videos. It contains textual modality about the content of the conversations, visual modality about the body and face, and audio modality about the voices.

We evaluated the reliability of annotators' subjective evaluations by calculating the intraclass correlation coefficients and Fleiss' Kappa. Statistical analysis of these annotations showed that professional counselors tend to demonstrate cognitive empathy rather than emotional empathy when establishing empathic connections with their clients, particularly in psychotherapy scenarios. Furthermore, we constructed three baseline multimodal analysis models using the MEDIC dataset and found that each modality plays a significant role.

2 RELATED WORK

Empathy has been researched due to its potential in psychology. There are two data-driven tasks: empathy prediction and empathy conversation generation. For empathy prediction, Alam et al. [1] use speech data from call centers to predict whether agents would show empathy when facing anger or frustration from clients. Buechel et al. [7] predict empathy when people read news reports. Barros et al. [5] asked people to tell stories and then captured the empathy of the listeners for empathy prediction. For empathy conversation generation, the EMPATHETICDIALOGUES dataset contains 25,000 individual conversations based on specific contexts [27]. A dialogue generation model is constructed to generate empathic responses. Although empathy has been applied in the above domains, its computational approach has not been sufficiently studied in the field of psychotherapy due to a lack of datasets. Table 1 shows existing data sets relevant to empathy in a psychotherapeutic context. They are described in detail in the following paragraphs.

In the field of psychotherapy, empathy is generally used in the prediction task. Several empathy datasets are built by collecting texts. For example, a study from a clinical trial used motivational interactions (MI) to construct a dataset from therapy related to

Dataset	Scenes	Annotation	Modal	Data Num	Public
MISC [32]	MI on substance use by college students	MISC 2.1	text	7293	no
MITI [32]	MI on substance use by college students	MITI 2.0	text	88	no
Gibson et al. [15]	Counseling on substance abuse	MITI 2.0, MISC 1.1	text	337	no
CTT [33]	Counseling on substance abuse	MITI 3.0	text, audio	200	no
Pérez-Rosas et al. [26]	Counseling therapy	MITI 4.1	text, audio	276	no
Sharma et al. [30]	Online peer mental health support	Emotional Reactions Interpretations Explorations	text	10,143	yes
MEDIC	Psychological counseling	Expression of experience Emotional reactions Cognitive reaction	text, audio, visual	771	yes

Table 1: Empathy datasets in the field of psychotherapy.

substance use in college students [32]. Researchers manually transcribed the text from the therapist portion of the session. Then, they analyzed 28 MI sessions using the Motivational Interviewing Skill Code (MISC) version 2.1 [22], yielding 854 empathic and 6439 non-empathic utterances. The MISC manual describes the behavior of the counselor and client at the utterance-level and assesses the overall competence of the counselor. This dataset is named MISC. Then they evaluated the empathy level of 88 additional sessions using the Motivational Interviewing Treatment Integrity (MITI) version 2.0 [25]. The MITI is a session-level counselor coding scheme used to give a global score of counselor empathy on the Likert scale. This dataset is named MITI. The data used in the annotation of these two datasets includes audio and original transcripts. James Gibson et al. [15] collected 337 texts from motivational interviews from six separate clinical studies. These studies all focused on addiction counseling related to substance abuse. The researchers manually transcribed and segmented all data at the turn level, and coded them for session-level behavior according to the MITI 2.0 manual. Then talk turns were segmented into utterances, which were assigned utterance-level behavioral codes according to the MISC 1.1 manual [2]. Ashish Sharma et al. [30] used two online mental health support platforms, TalkLife and Mental Health Subreddits, as data sources to obtain 10,143 pairs of communicative texts. They devised a new method to describe empathy that incorporates emotional reaction, interpretation, and exploration. The degree of empathy depends on whether or not these responses are expressed.

Texts are not the only vehicle for perceiving empathy. Pitch et al. [28] carries information about the emotional state of the speaker and has been shown to be related to the perception of empathy in psychotherapy. Bao Xiao et al [33], further provided empathy prediction data, including text and audio, in Context Tailored Training (CTT). The dataset contains 200 motivational interviews, includes an observer rating of a counselor’s empathy using MITI 3.0 [23]. Sessions were transcribed using Automatic Speech Recognition (ASR), and the generated text was used in a text-based empathy prediction model. These speeches and language processing techniques could accurately predict therapists’ performance from the audio recordings. Another study used 276 MI audio sessions from clinical studies [26]. The full set consists of 97.8 hours of audio, with an average session length of 20.8 minutes. The researchers

used MITI 4.1 [24] to label the client’s and counselor’s empathy levels separately.

All of the above datasets contain only text and/or audio modalities. However, empathy can also be expressed visually. For example, facial expressions and gaze can be used to infer empathic behavior [18]. Therefore, we introduce visual modality into empathy research in the field of counseling. In addition, most of the above datasets use different versions of the MISC and MITI scales to measure empathy. These are two classic empathy measurement scales, but neither of them takes into accounts for the complexity of empathy [30]. They assess only cognitive empathy and ignore emotional empathy. More importantly, MISC describes client and counselor behavior at the utterance-level, while MITI describes counselor behavior at the session-level [32]. This does not match the talk turn level sample on our dataset, so we do not use either of these scales. Sharma et al.’s [30] dataset considers both cognitive and emotional empathy. Their description of cognitive empathy in terms of both interpretation and exploration is illuminating, but they only consider the strength of expressions and ignore correctness. They may misinterpret cognitive expressions as highly emotional ones. Therefore, we have combined the above empathy scales for our annotation. We consider both emotional and cognitive empathy. To fit the talk turn level sample, we also introduce an empathy cycle in the labels [3]. Finally, only the dataset from online psychological support platforms with non-professional counseling is publicly available. However, we will make our dataset available for future research.

In contrast to the existing datasets, our proposed dataset, MEDIC, introduces visual modality into empathy research in the field of counseling. It considers both emotional and cognitive empathy and is designed to fit the talk turn level sample. Furthermore, our dataset will be publicly available for easy access, making it a valuable resource for empathy research in the field of psychotherapy.

The major contributions of this paper are:

- To the best of our knowledge, MEDIC is the first empathy dataset for psychotherapy scenarios that considers visual, audio, and textual modalities.
- We devise a new method for describing empathic communication in multimodal scenarios.
- We construct three multimodal analysis baseline models on the MEDIC dataset to demonstrate the contribution of each modality in predicting psychological empathy.

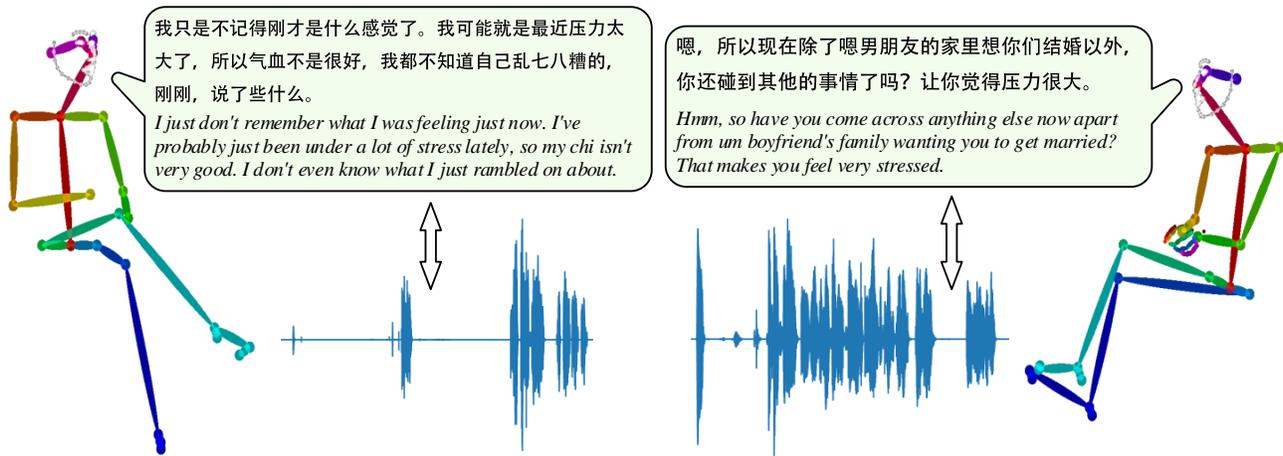


Figure 1: The sample of our dataset holds three modes of data: visual, textual and audio. The italicized text is the English translation of the original. On the left, the client is expressing his or her situation. On the right the counselor is asking about the reasons for the client's emotion.

3 HOW TO MEASURE EMPATHY

Client-centered therapy and psychoanalytics are two therapeutic approaches to empathy in psychotherapy corresponding to the cognitive and emotional aspects of empathy, respectively [29]. Some traditional counselors emphasize a conscious perspective selection process rather than a more automatic body-based emotion simulation process [6]. However, empathy is a complex multidimensional concept, encompassing the individual's perceptions, transpersonal abilities or dispositions, and emotional reactions [10]. Empathy research has also been carried out in emotional [27] and cognitive [33] directions.

One empathic process model, the Empathy Cycle [3], is particularly relevant. In the Empathy Cycle, the client and therapist work together to find an accurate expression of the client's experience. The cycle includes four steps. First, the client expresses the experience. The counselor generates empathy, and then expresses it. Lastly, the client accepts empathy. Counselor-generated empathy and client-accepted empathy are inherently unobservable.

Inspired by these studies, we propose a new conceptual method for describing empathy. This method consists of three mechanisms of empathy communication. We find the Empathy Cycle fits perfectly with our talk turn level sample. Therefore, we choose client expression of experience and counselor expression of empathy as important components of our scale. In conjunction with the multidimensional meaning of empathy mentioned above, we subdivide counselor-expressed empathy into cognitive and emotional reactions. In total, this method consists of three mechanisms of empathy communication: expression of experience, emotional reaction, and cognitive reaction. For each mechanism, we adopt a three-category scale: no expression (0), weak expression (1), and strong expression (2). The method is described in detail below.

Expression of Experience (EE). The expression of the client's experience is the first step in the empathy cycle. Empathy is an internal experience; responsive empathy corresponds to the responder [3]. The responder cannot express empathy to an experience that

has not been expressed in some way. The expressing effect depends on the quality of the receiver, the signal, and the sender, so the signal must be measured. The experience here includes not only what clients have done, but also how they feel. Therefore, they should ideally express their feeling or describe an experience. When EE is no expression (0), there are no expressed emotions or described experiences. EE with weak expression expresses a weak emotion or mentions an experience. EE with strong expression corresponds to a strong emotion or a full description of an experience.

Emotional Reactions (ER). The emotional reaction expressed by the counselor is part of the empathic reaction. The counselor first observes the expression of the client's experience, then develops and expresses empathy. The counselor usually expresses emotions such as warmth and compassion. Different modalities have different forms of presentation for emotional reactions. The text contains explicit and/or implicit emotional words [7]. In addition to the verbal information, various acoustic features in audio, such as pitch and loudness, also contain rich emotional reactions [1]. The emotional reactions contained in the videos are primarily reflected in facial expressions and body movements [5]. For generalisability reasons, the emotional reaction labels should correspond to no emotional reaction, weak emotional reaction, and strong emotional reaction. For example, in everyday conversation, when a person shares a very funny experience, the listener may laugh out loud. However, we found that in our dataset, the counselor did not express strong emotional reactions. The final labels in the dataset for emotional reactions include only no emotional reactions (0) and weak emotional reactions (1).

Cognitive Reactions (CR). The cognitive reaction expressed by the counselor is another part of the empathic reaction. Sharma et al. suggest that cognitive reactions can be divided into interpretation and exploration [30]. No cognitive reaction means that the counselor does not refer to or further explore the client's feelings and experiences. A weak cognitive reaction means that the counselor mentions or expresses interest in the client's feelings

	EE	ER	CR	Talker	Message
1)	2	1	2	Client	Did I do a bad job? Did I do something wrong?
				Counselor	It doesn't seem like your daughter has been affected by you telling her about this at all yet. On the contrary, she feels that she can comfort her mother.
2)	1	0	0	Client	Huh? No more? So soon today?
				Counselor	Hmm.
				Client	Well, good.
				Counselor	We can continue next week.
3)	0	0	2	Client	What should I say?
				Counselor	You can start with when you and your boyfriend met.

Table 2: Typical examples of three extreme cases. These samples contain labels with very high and very low values. Due to space limitations, only the translations of the text content and the corresponding labels are provided here.

and experiences. The feelings and experiences may be not accurate. A strong cognitive reaction means that the counselor explicitly explains or explores the client's feelings and experiences. This explanation must accurately correspond to the client's feelings or experiences.

4 DATA COLLECTION

4.1 Data Source

Our data is based on UM Psychology's counseling case videos¹. The case videos contain 11 hours of recorded interactions between clients and counselors. The cases cover counseling issues including marital relationships, professional dilemmas, family education, the meaning of life, and other topics. The cases also contain the reactions of counselors while dealing with clients with a variety of attitudes. Both the client and the counselor in the video are UM Psychology counselors, and they acted out some cases that used to be real. Each case is played out by two different counselors. The client and counselor sit face to face and a camera records the movements and expressions of both. The 10 cases are divided into 38 videos, each averaging 17 minutes and 50 seconds. In total, we record 678 minutes (approximately 11 hours) of video.

4.2 Data Pre-Processing

For privacy reasons, we manually removed personal information and cropped out any mentions of specific people or places. For the image data, we use OpenPose[8] to extract feature points from the face, torso, and arms. The final data set will use feature points to protect user privacy. As our data source consists of a re-enactment of the teaching video, the original patient is absent from the visual content. This approach effectively mitigates potential ethical issues that may arise from the use of identifiable patient data. These ensure that our data set can be public.

Counseling videos are cropped according to talk turns. According to the Empathy Cycle, the dialogue rounds start with the client expressing an experience and end with the counselor expressing an empathic reaction (presumably accepted by the client). However, dialogue is frequently interrupted by vocalized pauses (e.g., uh, um). One of the more common scenarios for this is when the counselor

uses "um" during the conversation to express that he or she is listening intently to what the other person is saying. This is a trick or habit in conversation. However, a problem arises when these pauses are used as segments. In such cases, the counselor's words do not convey any meaningful information, making it difficult for the ER and CR to make accurate judgments due to the lack of information. These vocalized pauses are not used to split the talk turn, so an overall talk turn starts with content expressed by the client and ends with the counselor. Each talk turn corresponds to a sample and has a set of labels. As the talk turn contains both client and counselor sections, it can be further broken down by person. Finally, we separated the audio and image data and used an audio text recognition tool to extract the approximate text modality. The audio-to-text tool we use is a video and audio editing tool from Wondershare². To ensure the accuracy of the text, the video and text are put back together for manual proofreading. Figure 1 shows the form of our sample and contains visual modality about the client and the consultant, the respective audio, and the text of the conversation.

4.3 Data Annotation

Empathy is not a universal concept, and the meaning of empathy varies by culture and language. To ensure that linguistic meanings were understood accurately, data sources were annotated by five students who spoke the same language (Chinese). Annotators were thoroughly trained in our labeling methods and had to pass a live data annotation trial. Our annotators received full annotation instructions and some classic examples. We then conducted three one-hour offline sessions to explain details about empathy and each label's meaning. During this process, annotators were asked questions about annotations, which largely resolved potential ambiguity issues. Before formal annotation began, we conducted an on-site test to ensure that each annotator was qualified for annotation.

The annotator pairs watched video clips of each sample and identified three mechanisms (the client's expression of experience and the counselor's emotional and cognitive reactions). During the annotation process, the annotator is asked to pay attention not only to the information contained in the spoken words but also to the voice and the intonation. At the same time, they must consider

¹https://appdu96wh3o1781.h5.xiaoeknow.com/v1/goods/goods_detail/p_5fd6221ce4b04db7c094a079?type=3&type=3&jump_from=1_05_37_01

²<https://unicconverter.wondershare.cn/>

	Average Sample	Total
Talk Turns	1	771
Speaking Times	4.29	3306
Number of Words	129.45	99808
Duration	52.76s	678m
Number of Frames	1137	876692

Table 3: Dataset statistics.

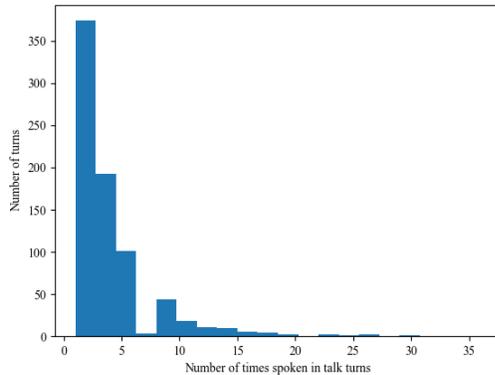


Figure 2: Distribution of times spoken in talk turns.

the expressions and actions of clients and counselors. Each sample was annotated by two annotators. For each label, we verify if the two annotators’ results are identical. If there is a discrepancy in the results, a third annotator will review the video clip again to determine which result is superior.

4.4 Sample Analysis

Table 2 shows a few typical samples on all three labels. This means that the client engages in strong expression, and the counselor gives a relatively strong empathic reaction. The client expresses doubt about her’s own behavior and shows a strong desire for approval. In addition, although no evidence is given in the forms, the client shows an expression of anxiety in the images. The client’s voice sounds anxious. Therefore, in this sample, the EE was marked as expressing strong emotions. Based on the context, the client is concerned that her behavior has a negative effect on her daughter. The counselor uses a soft voice and gentle expression to reassure the client that her daughter is not affected. The ER label corresponds to a weak emotional reaction. Finally, from the cognitive aspect, the counselor clearly explains the client’s previous behavior and further reassures the client. Therefore, the CR label corresponds to a strong cognitive reaction.

The second and third samples received lower scores on the EE, ER, and CR labels. The lowest scores were obtained on the second ER and CR labels. The client was a little surprised at the end of the consultation and expressed weak emotions, so EE was annotated as a weak expression. The counselor had only a simple reaction to the client’s words and did not react emotionally or cognitively. Therefore, ER and CR are marked as having no empathic reaction. The

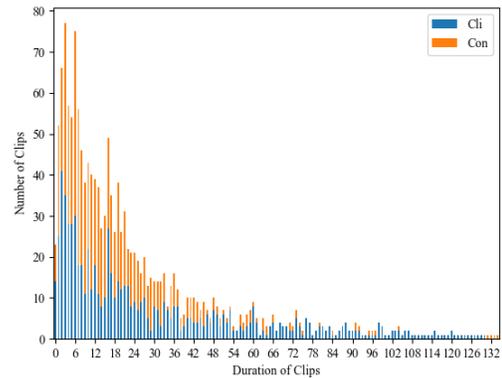


Figure 3: Distribution of duration of single speaking.

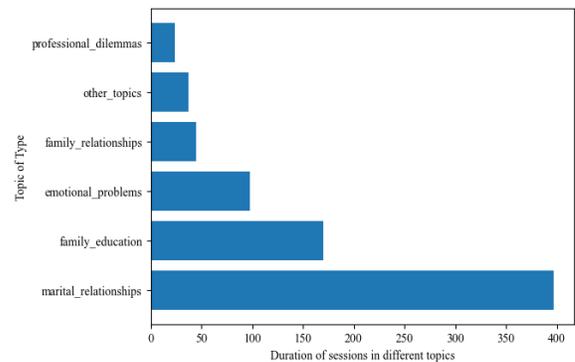


Figure 4: Duration of sessions in different topics

third sample had the lowest EE and ER scores. The client initially asks what she should say, and no emotion or experience is expressed. Therefore, EE corresponds to no expression. The counselor leads the client to discuss past experiences. Therefore, CR corresponds to strong cognitive reactions. ER is marked as no expression because the counselor does not express emotional reaction.

4.5 Statistical Results

Individual statistical results are given in Table 3. The final dataset contains 771 talk turns, each corresponding to one sample. The average length of each talk turn is 53 seconds. Each talk turn contains approximately 4 sentences, 53 seconds of audio, and 1137 frames of image features.

Figure 2 shows the distribution of the number of times spoken. The majority of talk turns contain only a few exchanges, with client and counselor typically only speaking once each. In some talk turns, the client or consultant uses a vocalized pause as a reaction. We do not split the talk round at the pause. On average, each client or counselor speaks 4.29 times in a talk turn, for a total of 3306 times. The distribution of the length of each person’s speech is shown in Figure 3. In most cases it is very short. However, it takes longer for the client to describe personal experiences. In addition, the average

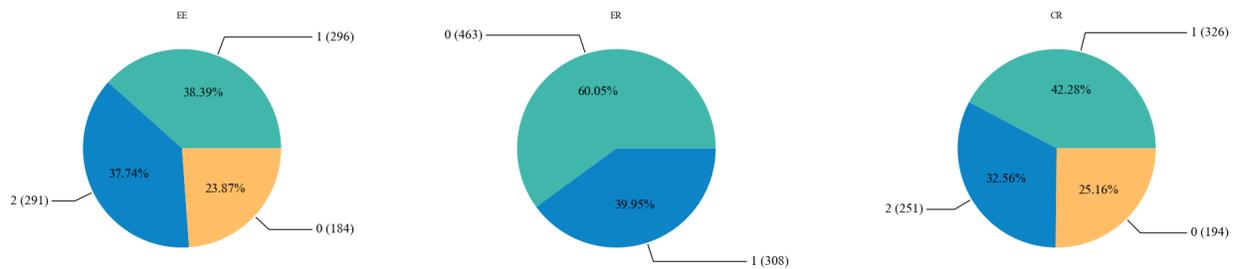


Figure 5: The overall distribution of EE, ER, and CR scores.

length of a single client comment is 11.48s, much more than the length of a single counselor comment of 6.59s. The graph also shows that counselors speak longer when the conversation is very short, and the clients speak more when the conversation is very long. In counseling, clients spend most of the time describing experiences while the counselor listens. Figure 4 presents the total duration of consultations across different themes. The consultations on marital relationships lasted the longest and were significantly longer than the others. This suggests that marital relationships have a major influence on people’s mental health and well-being, and that many people face challenges and hardships in their marriages and require professional assistance and advice. Among all the themes, only professional dilemmas are solely related to their own issues, while most of the other themes stem from relationships with others. This also indicates that interpersonal relationships are one of the key factors affecting mental health.

Figure 5 shows the distribution of labels for the final samples. Overall the distribution of all three labels is relatively balanced. For clients, the percentage of EE labeled 0 is only 24%, indicating that clients are comfortable confiding about their experiences to the counselor. For the counselor, the unrecognized expression of CR is only 25%. This indicates that counselors tend to use cognitive empathy to interact with the client. In contrast, no emotional expressions account for 60% of ER, even when strong emotional expressions are removed beforehand. This indicates that counselors tend not to express empathy through emotions. This situation proves that modern counselors focus on cognitive expression rather than emotional expression, which is consistent with previous research [6].

Figure 6 shows the distribution of labels under different topics. The distribution of Expression of Experience and counselor responses is relatively similar across the different topics, except for family relations and career dilemmas. The topics related to marital relationships, family education, and emotional problems have a more balanced distribution of labels. For family relations, the counselor mostly gave weak responses. For career dilemmas, both the client’s Expression of Experience and the counselor’s responses were strong. This may suggest that the counselor views career difficulties as more personal than situational or interpersonal, and thus tries to help the clients recognize their own problems and resources and improve their problem-solving skills.

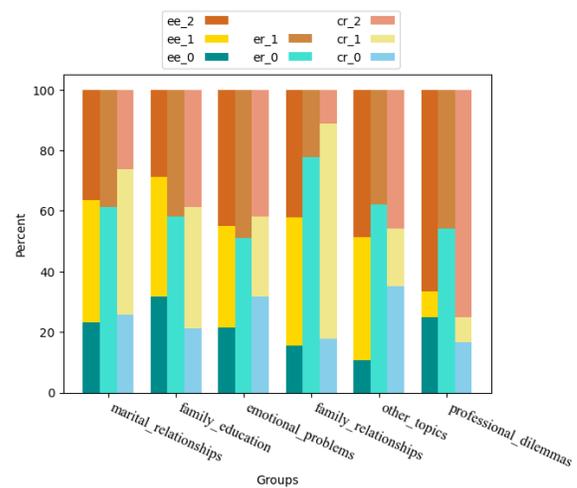


Figure 6: Distribution of labels under different topics

Pearson correlation coefficients of $r = 0.45$ are calculated for the two labels ER and CR. It indicates that the two scales are somewhat correlated but capture different empathy phenomena, emphasizing the importance of using multiple variables to describe the empathic reaction. In terms of label reliability, we analyzed consistency between two annotators for each label. Each sample was scored by two labelers and only given to a third for decision-making if the first two did not agree. We measured annotation consistency between the two individuals using the intraclass correlation coefficient (ICC) [17], Fleiss’ Kappa [14], and percentage agreement. The results are shown in Figure 4. The values of Fleiss’ Kappa range from -1 to 1, where 0.61-0.80 signifies significant agreement. For ICC, values range from 0 to 1, with values closer to 1 indicating greater reliability. Values less than 0.5 indicate poor reliability, between 0.5 and 0.75 indicate moderate reliability, and between 0.75 and 0.9 indicate good reliability. The values of Percent Agreement range from 0 to 1, where 0 means complete disagreement and 1 means complete agreement among evaluators. This means all three labels in the dataset have a good confidence level. For consistency percentages, more than half of the annotated results for each label were free

Measure	EE	ER	CR
Fleiss' Kappa	0.7557	0.6991	0.7102
ICC	0.9021	0.8126	0.8732
Percent Agreement	0.8067	0.8586	0.7846

Table 4: Consistency of annotation across the three labels.

from ambiguity. Our dataset shows high reliability by integrating all results from the three metrics mentioned above.

5 EVALUATION EXPERIMENTS

To evaluate our dataset, features are extracted from three modalities and empathy classification is performed on three different models. TFN [34] is a highly cited classical model, and the SWAFN [9] model focuses on textual modality, which is of paramount importance for consulting. A simple concatenation model is used as a comparison.

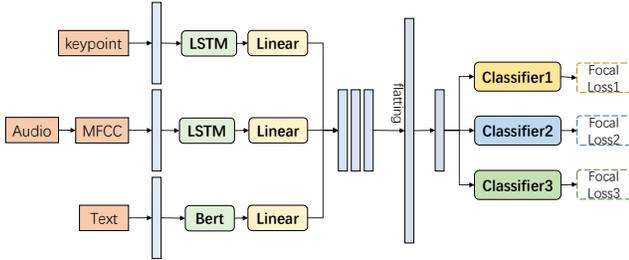


Figure 7: A simple framework for empathy prediction.

5.1 Feature Extraction

Visual features (v). As mentioned previously, OpenPose [8] was used to extract keypoints of the face, hands, and body for each frame of the video for the client and consultant. The keypoint coordinates were then normalized. A distinction was also made between the keypoints of the client and the consultant. The final input image feature dimension was $d_v \in R^{T_v \times 411}$, where T_v is the number of frames in the clip.

Audio features (a). Audio features corresponding to the segments were extracted using librosa [20]. The audio is separate for the client and the consultant. In this experiment, MFCC features are used for training. The final input audio feature dimension is $d_a \in R^{T_a \times 20}$, where T_a is determined by the duration of the audio.

Text features (t). Our text is derived from a transcription. We use the Bert pre-training model to extract text features. The dimension of the final input text feature is $d_t \in R^{T_t \times 768}$, where T_t is determined by the length of the sentence.

5.2 Baseline Model

The TFN model. TFN [34] consists of three components. The sub-network takes single-modal features as input and a multimodal embedding as output. A tensor fusion layer uses a 3-fold Cartesian

	label	EE	ER	CR
train set	0	127	339	145
	1	204	200	234
	2	208	*	160
val set	0	26	39	20
	1	28	38	24
	2	23	*	33
test set	0	31	85	29
	1	64	70	68
	2	60	*	58

Table 5: The number of labels in each category in the divided dataset.

product to explicitly model unimodal, bimodal, and trimodal interactions. A neural network uses the output of the tensor fusion layer as input to predict empathy.

The SWAFN model. SWAFN [9] uses a co-attentive mechanism to learn bidirectional large-scale contextual information between language and other modalities. The sentiment word classification task is then integrated into the model via a multi-task learning mechanism that guides the learning and aggregation of multimodal fusions.

The concatenation model. We also use a very simple multimodal concatenation model as shown in Figure 7. The temporal information of the three modalities are extracted using LSTM. Then they are spliced together and connected to the classifier to predict empathy.

Since the three labels are derived from two different people, in this paper we separately predict corresponding labels using data from the client and the counselor. The client’s features are used to predict EE, and the counselor’s features are used to predict ER and CR.

5.3 Experiment setup

Definition. Let $X_x \in R^{\tau_x \times d_x}$ denote the features of a modality. In our dataset, x can represent visual (v), audio (a), and text (t) modalities. τ_x represents the temporal length of the modality, and d_x represents the feature dimension of the modality. Let $Y = \{y_1, y_2, \dots, y_l\}$ denote the label space with l labels. In our experiments for the three labels EE, ER, and CR, $l = 3$. For the classification task of predicting EE, ER, and CR, the goal is to predict the corresponding label y_i given the input feature of clip X_{x_i} , where $x \in \{a, v, t\}$.

Sample details. We conduct experiments on our dataset, which contains 771 samples. Each sample contains key points of images, audio, and dialogue text from both the client and the consultant. Each sample is annotated with the corresponding EE, ER, and CR. The training, test, and validation sets are split in a ratio of 7:1:2. Details of each category within each label are shown in Table 5. The evaluation metrics are average accuracy(Acc) and macro F1-score.

5.4 Implementation Details

The TFN and SWAFN were trained on the V100 GPU and the concatenation model was trained on the RTX 3090 GPU. The pytorch and Adam optimizer is used on all models. The labels are set with

Modal	TFN		SWAFN		Concatenation	
	Acc	F1	Acc	F1	Acc	F1
v+a+t	0.758	0.758	0.864	0.863	0.819	0.810
v+a	0.555	0.526	0.721	0.718	0.716	0.708
v+t	0.746	0.744	0.805	0.808	0.813	0.804
a+t	0.738	0.736	0.851	0.852	0.800	0.774
v	0.416	0.307	0.636	0.595	0.325	0.301
a	0.527	0.465	0.695	0.686	0.697	0.699
t	0.729	0.726	0.857	0.857	0.768	0.763

Table 6: Baseline and ablation experiments for EE.

Modal	TFN		SWAFN		Concatenation	
	Acc	F1	Acc	F1	Acc	F1
v+a+t	0.729	0.719	0.743	0.743	0.703	0.699
v+a	0.646	0.639	0.743	0.744	0.639	0.631
v+t	0.725	0.720	0.770	0.761	0.697	0.712
a+t	0.743	0.734	0.776	0.777	0.690	0.688
v	0.600	0.375	0.592	0.598	0.568	0.549
a	0.642	0.595	0.724	0.727	0.652	0.647
t	0.734	0.717	0.770	0.764	0.658	0.657

Table 7: Baseline and ablation experiments for ER.

the weights of the classes, which are inversely proportional to the number of classes due to their unbalance. The learning rate in TFN is set to $1e-4$, the batch size is set to 32, and the dropout is set to 0.3. The learning rate in SWAFN is set to $1e-3$, the batch size is set to 32, and the dropout is set to 0.3. The learning rate in the concatenation model is set to $1e-4$, the batch size is set to 16, and the dropout is set to 0.4.

5.5 Experimental Results and Analysis

Tables 6, 7, and 8 show the experimental results of three models for the EE, ER, and CR of prediction tasks. The relevant ablation experiment about visual(v), audio (a), and text(t) modalities are also included. For example, $v + a + t$ denotes the fusion of the three modalities, while $v + a$ denotes the fusion of two modalities v and a . From the tables, we can draw the following conclusions.

First, from the three tables, we can find that the F1 scores and accuracy for EE, ER, and CR are all above 69%. This demonstrates that our dataset is highly effective in facilitating empathic prediction.

Secondly, the three modal fusions of $v + a + t$ have achieved the best results. The results of two modal fusions are slightly lower, and the results of one modality alone are the lowest. For example, in Table 8, the best performing model is SWAFN. Its F1-score of $v + a + t$ is 7.6% higher than the results of $a + t$ and 11.5% higher than the results of the text modality alone. Its accuracy of $v + a + t$ is 7.2% higher than $v + t$ and 9.2% higher than the results of the text modality alone. This indicates that the provided multimodal information and fusion are effective.

Focusing on the v modality, we find that the results of $v + a + t$ are higher than $a + t$. For example, the three modal F1-score and accuracy of the TFN model in Table 6 are both 2% higher than $a + t$. The three modal F1-score and accuracy of the Concatenation

Modal	TFN		SWAFN		Concatenation	
	Acc	F1	Acc	F1	Acc	F1
v+a+t	0.722	0.712	0.783	0.785	0.735	0.731
v+a	0.501	0.439	0.658	0.657	0.613	0.581
v+t	0.734	0.725	0.678	0.661	0.697	0.701
a+t	0.716	0.711	0.711	0.709	0.684	0.687
v	0.436	0.259	0.467	0.420	0.490	0.450
a	0.500	0.395	0.605	0.570	0.594	0.556
t	0.690	0.685	0.691	0.670	0.710	0.722

Table 8: Baseline and ablation experiments for CR.

model in Table 8 are both 4.4% higher than $a + t$. This shows the importance of the video modality.

Table 7 shows that $a + t$ achieves the best results on the TFN and SWAFN models, while $v + a + t$ and $v + t$ achieves the best results in concatenation model. This indicates that $v + a + t$ contributes to the training of the ER, but different models focus on different modal information when predicting the ER.

Finally, in the unimodal ablation experiments, t performs the best. In Table 6, for the SWAFN model, the F1-score and accuracy of t are 17.1% and 16.2% higher than a . In Table 7, for the SWAFN model, the F1-score and accuracy of t are 3.7% and 4.6% higher than a . In Table 8, for the concatenation model, the F1-score and accuracy of t are 16.6% and 11.6% higher than a . This is because the text contains the richest verbal information, which is considered to be the most important element in counseling.

5.6 Error Analysis

Our study reveals that all three models incorrectly labeled 1 as 0 for the ER labels in some samples, likely due to imbalanced data distribution. This highlights the need to enhance the models' ability to handle unbalanced samples. Moreover, we discovered that the models tended to rely on sample length, leading to incorrect predictions for some shorter and longer samples across all three labels. This indicates a deficiency in the models' ability to extract semantic meaning.

6 CONCLUSION

Empathy is a critical component in human interaction and plays an important role in facilitating communication between clients and counselors. There is a paucity of multimodal empathy datasets in the field of psychology, and the corresponding empathy evaluation criteria are not well developed. In this paper, we construct a publicly available multimodal dataset for empathy. In order to understand empathy expressed in counseling exchanges, we propose three labels to describe empathic communication in multimodal scenarios. The dataset contains image modality, audio modality, and text modality from counseling scenarios. Collecting multimodal empathy datasets can be challenging due to the need to protect privacy while gathering data from multiple sources. Additionally, accurately measuring empathy, a complex and multidimensional concept, requires carefully designed assessment criteria. As a result, the amount of data we were able to collect was limited. Despite this, the high-quality data in our dataset holds significant value for

empathy research. We hope that more empathy studies in the field of counseling will build on this foundation.

REFERENCES

- [1] Firoj Alam, Morena Danieli, and Giuseppe Riccardi. 2018. Annotating and modeling empathy in spoken conversations. *Computer Speech & Language* 50 (2018), 40–61.
- [2] Paul C. Amrhein, William R. Miller, Theresa B. Moyers, and Denise Ernst. 2008. Manual for the Motivational Interviewing Skill Code (MISC). *Motivational Interviewing Skill Code v. 2.1* (2008).
- [3] Godfrey T Barrett-Lennard. 1981. The empathy cycle: Refinement of a nuclear concept. *Journal of counseling psychology* 28, 2 (1981), 91.
- [4] Godfrey T. Barrett-Lennard. 2011. The phases and focus of empathy. Barrett-Lennard, G.T. <https://researchrepository.murdoch.edu.au/view/author/Barrett-Lennard_Godfrey.html> (2011) *The phases and focus of empathy. British Journal of Medical Psychology*, 66 (1), pp. 3-14. (2011).
- [5] Pablo Barros, Nikhil Churamani, Angelica Lim, and Stefan Wermtner. 2019. The omg-empathy dataset: Evaluating the impact of affective behavior in storytelling. In *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*. IEEE, 1–7.
- [6] Arthur C Bohart and Leslie S Greenberg. 1997. Empathy: Where are we and where do we go from here? (1997).
- [7] Sven Buechel, Anneke Buffone, Barry Slaff, Lyle Ungar, and Joao Sedoc. 2018. Modeling empathy and distress in reaction to news stories. *arXiv preprint arXiv:1808.10399* (2018).
- [8] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. 2017. Realtime multi-person 2d pose estimation using part affinity fields. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7291–7299.
- [9] Mingping Chen and Xia Li. 2020. SWAFN: Sentimental Words Aware Fusion Network for Multimodal Sentiment Analysis. *international conference on computational linguistics* (2020).
- [10] Matthew Davis, Monique G. Davis, M. Edward Davis, Mark E. Davis, Mark Mitchell Davis, MM Davis, Mdp Davis, Frederick B. Davis, H. Davis, and I. W. Davis. 1980. A multidimensional approach to individual differences in empathy. *JASAS Catalog of Selected Documents in Psychology* (1980).
- [11] Frans B. M. de Waal. 2008. Putting the Altruism Back into Altruism: The Evolution of Empathy. *Annual Review of Psychology* (2008).
- [12] Robert Elliott, Arthur C Bohart, Jeanne C Watson, and Leslie S Greenberg. 2011. Empathy. *Psychotherapy* 48, 1 (2011), 43.
- [13] Robert Elliott, Leslie S Greenberg, Jeanne C Watson, Ladislav Timulak, and Elizabeth Freire. 2013. Research on humanistic-experiential psychotherapies. (2013).
- [14] Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin* 76, 5 (1971), 378.
- [15] James Gibson, Dogan Can, Bo Xiao, Zac E. Imel, David C. Atkins, Panayiotis G. Georgiou, and Shrikanth S. Narayanan. 2016. A Deep Learning Approach to Modeling Empathy in Addiction Counseling. *conference of the international speech communication association* (2016).
- [16] Emily A. Holmes, Ata Ghaderi, Catherine J. Harmer, Paul Ramchandani, Pim Cuijpers, Anthony P. Morrison, Jonathan P. Roiser, Claudi L H Bocking, Rory C. O'Connor, Roz Shafran, Michelle L. Moulds, and Michelle G. Craske. 2018. The Lancet Psychiatry Commission on psychological treatments research in tomorrow's science. *The Lancet Psychiatry* (2018).
- [17] Terry K Koo and Mae Y Li. 2016. A guideline of selecting and reporting intraclass correlation coefficients for reliability research. *Journal of chiropractic medicine* 15, 2 (2016), 155–163.
- [18] Shiro Kumano, Kazuhiro Otsuka, Dan Mikami, and Junji Yamato. 2011. Analyzing empathetic interactions based on the probabilistic modeling of the co-occurrence patterns of facial expressions in group meetings. In *2011 IEEE International Conference on Automatic Face & Gesture Recognition (FG)*. IEEE, 43–50.
- [19] Michael J Lambert, Allen E Bergin, and SL Garfield. 1994. The effectiveness of psychotherapy. *Encyclopedia of psychotherapy* 1 (1994), 709–714.
- [20] Brian McFee, Colin Raffel, Dawen Liang, Daniel P Ellis, Matt McVicar, Eric Battenberg, and Oriol Nieto. 2015. librosa: Audio and music signal analysis in python. In *Proceedings of the 14th python in science conference*, Vol. 8. Citeseer, 18–25.
- [21] John McLeod. 2003. *Doing counselling research*. Sage.
- [22] William R Miller, Theresa B Moyers, Denise Ernst, and Paul Amrhein. 2003. Manual for the motivational interviewing skill code (MISC). *Unpublished manuscript. Albuquerque: Center on Alcoholism, Substance Abuse and Addictions, University of New Mexico* (2003).
- [23] TB Moyers, T Martin, JK Manuel, WR Miller, and D Ernst. 2007. Revised global scales: motivational interviewing treatment integrity 3.0 (MITI 3.0). *University of New Mexico, Center on Alcoholism, Substance Abuse and Addictions (CASAA)* 28 (2007).
- [24] T B Moyers, J K Manuel, & D Ernst, Lisa Hagen Glynn, Christiana Fortini, Ellis Baron, Hans Bertens, Saskia Boom, Angelita Casanovas, Jos Dobber, Jannet De Jonge, Maarten Merckx, Janet Murriss, Hetty De Laet, and Riëtta Oberink. 2022. Motivational Interviewing Treatment Integrity Coding Manual 4.2.1 (MITI 4.2.1).
- [25] Theresa B Moyers, Tim Martin, Jennifer K Manuel, William R Miller, and D Ernst. 2003. The motivational interviewing treatment integrity (MITI) code. *Unpublished manuscript* (2003).
- [26] Verónica Pérez-Rosas, Rada Mihalcea, Ken Resnicow, Satinder Singh, and Lawrence C. An. 2017. Understanding and Predicting Empathic Behavior in Counseling Therapy. *meeting of the association for computational linguistics* (2017).
- [27] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards Empathetic Open-domain Conversation Models: a New Benchmark and Dataset. *meeting of the association for computational linguistics* (2019).
- [28] Catherine M Reich, Jeffrey S Berman, Rick Dale, and Heidi M Levitt. 2014. Vocal synchrony in psychotherapy. *Journal of Social and Clinical Psychology* 33, 5 (2014), 481.
- [29] Robert L Selman. 1980. *The growth of interpersonal understanding: Developmental and clinical analyses*. Academy Press.
- [30] Ashish Sharma, Adam S. Miner, David C. Atkins, and Tim Althoff. 2020. A Computational Approach to Understanding Empathy Expressed in Text-Based Mental Health Support. *empirical methods in natural language processing* (2020).
- [31] Tania Singer and Claus Lamm. 2009. The social neuroscience of empathy. *Annals of the New York Academy of Sciences* (2009).
- [32] Bo Xiao, Dogan Can, Panayiotis G. Georgiou, David C. Atkins, and Shrikanth S. Narayanan. 2012. Analyzing the language of therapist empathy in Motivational Interview based psychotherapy. *asia pacific signal and information processing association annual summit and conference* (2012).
- [33] Bo Xiao, Zac E Imel, Panayiotis G Georgiou, David C Atkins, and Shrikanth S Narayanan. 2015. "Rate my therapist": automated detection of empathy in drug and alcohol counseling via speech and language processing. *PLoS one* 10, 12 (2015), e0143055.
- [34] Amir Zadeh, Minghai Chen, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. 2017. Tensor Fusion Network for Multimodal Sentiment Analysis. *empirical methods in natural language processing* (2017).