

Appraisal Framework for Clinical Empathy: A Novel Application to *Breaking Bad News* Conversations

Allison Lahnala¹, Béla Neuendorf³, Alexander Thomin³, Charles Welch¹
Tina Stibane^{2,3}, Lucie Flek¹

¹Conversational AI and Social Analytics (CAISA) Lab, University of Bonn, Germany

²Marburg Interactive Skills Lab (MARIS), Dr. Reinfried Pohl-Centre for Medical Education

³Philipps-University Marburg, Germany

{alahnala, lflek}@uni-bonn.de

Abstract

Empathy is essential in healthcare communication. We introduce an annotation approach that draws on well-established frameworks for *clinical empathy* and *breaking bad news* (BBN) conversations for considering the interactive dynamics of discourse relations. We construct EMPATHY IN BBNS, a span-relation task dataset of simulated BBN conversations in German, using our annotation scheme, in collaboration with a large medical school to support research on educational tools for medical didactics. The annotation is based on 1) Pounds (2011)'s appraisal framework for clinical empathy, which is grounded in systemic functional linguistics, and 2) the SPIKES protocol for breaking bad news (Baile et al., 2000), commonly taught in medical didactics training. This approach presents novel opportunities to study clinical empathic behavior and enables the training of models to detect causal relations involving empathy, a highly desirable feature of systems that can provide feedback to medical professionals in training. We present illustrative examples, discuss applications of the annotation scheme, and insights we can draw from the framework.

Keywords: clinical empathy, empathy annotation scheme, breaking bad news dialogues

1. Introduction

Empathy in medical encounters is considered a core element to high-quality patient care and an important skill to develop in medical training (Bonvicini et al., 2009; Dey and Girju, 2022). Theoretical models of *clinical empathy* suggest it fosters more open patient-clinician communication for more deeply understanding patients and their conditions, providing valuable information for diagnosis and addressing therapeutic needs. In turn, this leads to better treatment strategies and adherence, therapeutic outcomes, and higher patient satisfaction (Squier, 1990; Neumann et al., 2009; Pounds, 2011), and proper empathy can be intrinsically therapeutic (Suchman et al., 1997).

Empathy is crucial to *breaking bad news* conversations (BBNs), scenarios where a clinician must inform the patient about life-altering circumstances. Clinicians frequently must deliver bad news to patients, a high-stress and complex communication task (Baile et al., 2000). Many medical students must pass formal BBN training, an area where digital tools have the potential to support students via automated feedback and practice with virtual standardized patients (Borish et al., 2014; Lok and Foster, 2019; Reger et al., 2021). Models informed by clinical empathy could be made more transparent and explainable, leading to higher quality feedback on empathic responses or suggestions via an interface that suits their learning needs (Tanana et al., 2019; Girju and Girju, 2022). Natural language

processing (NLP) researchers are currently exploring models for automatic empathy detection and generation, which are essential for such tools.

Objectives: We introduce a new annotation scheme for clinical empathy communication in patient-clinician conversations, guided by three key objectives. First, we aim to identify precise discourse elements and their dynamic interactions that constitute empathic successes and failures during these exchanges. Second, we aim to provide a novel structural representation of empathic conversations that can support addressing existing shortcomings of NLP models for empathy, which struggle in identifying finer-level empathy components and utilizing the broader conversational context necessary for accurate evaluation (Lee et al., 2023). Third, we aim to complement established pedagogical methods for training communication skills for breaking bad news conversations by integrating SPIKES, a protocol for breaking bad news commonly taught in medical didactics training (Baile et al., 2000).

Method Overview: We address these goals through a *span-relation* labeling method. Central to the scheme is identifying interactional sequences formally defined by Suchman et al. (1997)'s *model of empathic communication in medical conversations*. These sequences encompass three elements: 1) *empathic opportunities* (EO) in patient turns, 2) *elicitations* of such opportunities, and 3) *empathic responses* within clinician turns. We label specific spans of patient and clinician turns

containing one of these three elements according to Pounds (2011)'s *appraisal framework for clinical empathy* (AF) which defines types of EOs and responses based on linguistic aspects (see §3.1).

The scheme provides two further innovations. After labeling the spans, we create relations between identified empathic opportunities and the identified elicitations and empathic responses that correspond to them. Furthermore, we incorporate the structure of BBN conversations by labeling the six stages of BBN conversations defined by the SPIKES protocol within the clinician turns (see §3.2). This integration aligns the appraisal framework with the stages of BBN conversations, facilitating the examination of communicative strategies employed at each stage and informing the development of NLP models for digital training tools.

The BBN EMPATHY Dataset: Finally, as part of a collaboration between NLP researchers and medical didactics experts at large medical school to support research on educational tools for medical didactics, we construct a dataset of BBN conversations annotated with the new scheme. The dataset contains practice BBN conversations between medical students and standardized patients and fine-grained annotations of the components of empathic interactions. The BBN EMPATHY Dataset is the first dataset to contain discourse labels and relations for clinical empathy, which we make public for other researchers.

Summary of Contributions: We contribute 1) an innovative annotation scheme for clinical empathy 2) made available on an open-source platform for other researchers (§3), and 3) a public dataset of annotated BBN dialogues (§4).¹

2. Related work

Digital tools have the potential to support training medical professionals in empathetic communication in ways such as offering communication practice with virtual standardized patients (Borish et al., 2014; Lok and Foster, 2019; Reger et al., 2021), assessing the quality of empathetic responses, and providing feedback on empathic responses or suggestions via an interface that suits their learning needs (Tanana et al., 2019; Girju and Girju, 2022). These tools require effective NLP models of empathic language and conversational behaviors. In NLP, there is current momentum toward models for empathy detection and generation. Recognition is the task of determining the presence (Sharma et al., 2020; Hosseini and Caragea, 2021) or degree of empathy (Buechel et al., 2018) or subtypes of empathic behaviors (Welivita and

Pu, 2020; Svikhnushina et al., 2022). Detection models can be designed to evaluate empathic language and identify improvement areas to provide feedback (Wambsganss et al., 2021). Generative models attempt to generate a text response that is empathetic to a conversational partner. Research on generative models has focused on applications to open-domain dialogue (Rashkin et al., 2019), customer support (Firdaus et al., 2020), and counseling (Shen et al., 2020). Generative models can be designed to deliver feedback and provide suggestions for empathetic responses. For example, Sharma et al. (2021) developed a model for "empathic rewriting" to provide suggestions that increase the level of empathy in a given text, an approach that could support students in reflecting on ways to improve their empathic communication.

Despite the progress, current shortcomings in this research include poor operationalization of empathy, tending to employ only abstract notions focused on emotional aspects and overlooking cognitive and behavioral aspects, and a lack of empathic language resources that incorporate these dimensions. These issues are exasperated by lacking measurements with construct validity (Lahnala et al., 2022). In turn, models trained on widely available datasets, such as the empathic dialogues dataset which grounds empathetic engagement in specific emotional situations, could help reveal patterns of emotional understanding, but this is only one facet of the empathy concept (Debnath and Conlan, 2023). Thus, such models are limited in providing detailed and reliable assessments, explanations of the relation between EOs and empathic responses, or validated guidance for developing clinical empathy skills. However, NLP can draw from extensive research in psychology and linguistics, which has empirically studied theoretical models of clinical empathy and measurement approaches.

Though they are still few, NLP studies exploring tools for training and education in empathetic communication have generally integrated insights from these fields, often in collaborations with psychologists. For example, Wambsganss et al. (2021, 2022) investigate the effectiveness of digital empathy training tools in enhancing students' empathetic communication skills when writing peer feedback. Other work focuses on applications for psychotherapy (Imel et al., 2017), examining for instance, linguistic behaviors that signal empathy and session quality in Motivational Interviewing conversations (Pérez-Rosas et al., 2017, 2018, 2019), and crisis counseling conversations (Zhang and Danescu-Niculescu-Mizil, 2020; Zhang et al., 2020).

Some recent NLP studies on clinical empathy started integrating linguistic theories and discourse analysis approaches. Dey and Girju (2022) cre-

¹Dataset and code: <https://github.com/caisa-lab/BBN-Empathy>

ated a corpus of medical students' essays about breaking bad news to patients with sentence-level labels of cognitive, affective, and prosocial empathy. They built a novel system architecture informed by frame semantics (Baker et al., 1998) that outperformed state-of-the-art empathy classification models. On the same dataset, Dey and Girju (2023) showed that incorporating features of Construction Grammar (Michaelis, 2006) and Systemic Functional Grammar (Halliday and Matthiessen, 2013) also improves deep learning models for empathy classification. Shi et al. (2021) also demonstrated the potential of the discourse annotation resources to improve empathy NLP models. Nevertheless, creating clinical dialogue datasets with quality clinical empathy annotations suitable for training NLP models requires expertise, labor, and ethics and privacy protections.

To address the described shortcomings, we contribute a novel dataset of annotated clinical empathy in breaking bad news conversations, an annotation method based on a linguistic framework for structured analysis of clinical empathy, and a framework of communication strategies for BBNs developed by medical experts. We tailor these resources to support the development of NLP models that can be integrated with digital training tools. Related work outside of NLP has also used discourse analysis approaches to identify strategies in clinical BBN conversations (Pun, 2021). Recently, Rey Velasco et al. (2022) applied SFL approaches and Pounds (2011)'s clinical empathy framework to asynchronous health interactions, revealing insights into implicit/explicit EOs and their relationship to trust between healthcare providers and patients. To our knowledge, this work is the first to apply Pounds (2011)'s framework to live conversations and BBN scenarios, which we make public to support NLP research.

3. Annotation Scheme

3.1. Appraisal Framework

Background and Motivation. *Empathic opportunities* (EOs) are expressions or behaviors that reveal a patient's feelings or views, which can be either *explicit* or *implicit* (Suchman et al., 1997). *Explicit empathic opportunities* directly reveal a patient's feelings or attitude via explicit expressions of behavior or emotive behaviors (e.g., crying). *Implicit empathic opportunities* are defined as "patient statements from which a clinician might infer an underlying emotion that has not been explicitly expressed" (Suchman et al., 1997).

The appraisal framework for clinical empathy extends Suchman et al. (1997)'s model of clinical empathic communication (Pounds, 2011), by inte-

grating a finer-grained taxonomy that categorizes EOs and doctor responses based on the linguistic functions of attitudinal expression. It draws on insights from (Wynn and Wynn, 2006)'s linguistic research on interactional sequences that build empathy in psychotherapy settings, and Martin and White (2005)'s appraisal framework, a systemic functional linguistics (SFL) approach to discourse analysis (Halliday and Matthiessen, 2013) that concerns the interpersonal, interactive functions of language in specific social settings. The framework aligns various linguistic functions with components of attitudinal expression within empathic communication in order to gain insights that support the teaching of empathic communication skills. The three dimensions of attitudinal expression are FEELING (i.e., affect), JUDGMENT of oneself or other people, and APPRECIATION; an attitude or perception toward things, events, actions, and behaviors.

Previous research observed that clinicians often overlook empathic opportunities, hindering effective, satisfactory communication with patients (Levinson et al., 2000; Hsu et al., 2012). Moreover, Suchman et al. (1997) finds implicit EOs are more common in a medical interview than explicit empathic opportunities. As implicit EOs are hidden in patients' expressions, they are particularly challenging to identify and infer. Thus, the ability to model explicit and implicit EOs is an important step toward digital tools that support communicative skills training for BBNs. By developing a dataset of EOs and their relations to clinician expressions, not only do we enable training such NLP models, but we also provide a resource for investigating linguistic aspects that could add to the body of knowledge about BBN conversation strategies. This resource contributes to ongoing efforts in broader NLP empathy research seeking multidimensional representations of empathy, including affective and cognitive empathy and empathic behaviors (Lahnala et al., 2022). The segments of inferred implicit EOs provide opportunities, especially to better understand cognitive empathy, and can be leveraged in research on abductive social reasoning (Bhagavatula et al., 2020; Zhao et al., 2023).

As we present the framework, we provide descriptions and examples inspired Pounds (2011).

Categories of Empathic Opportunities

Table 1 contains examples and descriptions of the FEELINGS, JUDGMENTS, and APPRECIATIONS categories of explicit and implicit EOs, yielding six possible labels to apply to EO spans. Explicit EOs are directly observable in the patient's expressions and behaviors. Implicit EOs can be explored by the clinician to more deeply understand patient views and can lead to explicit EOs and a better consensus on empathic accuracy, which informs the clinician's

Patient Role: Empathic Opportunities	
Explicit	Implicit
<p>FEELINGS Describing or exhibiting</p> <ul style="list-style-type: none"> emotion quality: "I'm sad" emotive behavior "I cried" or "I laughed" mental state: "I'm in pain" or "I feel alone" 	<p>May occur through expression of judgement or appreciation. Implicitly expresses feeling or perception by referring to</p> <ul style="list-style-type: none"> negative experiences: "My aunt had the same condition. She was in a lot of pain, and she didn't make it" (fear) critical life stages and experiences: "Everything was going well...I just started my master's thesis" (surprise/disbelief)
<p>APPRECIATION Attitude or perception toward things, events, actions, and behaviors, e.g.</p> <ul style="list-style-type: none"> event: "The MRI was boring." thing: "The medication is not helping me." 	<p>Indirectly conveys attitude or perception toward things, events, actions, and behaviors that a clinician may infer and explore</p> <ul style="list-style-type: none"> thing: "My symptoms don't seem to improve with the medication."
<p>JUDGEMENT Attitude or perception toward</p> <ul style="list-style-type: none"> self: "I'm not good at medication adherence" others: "The nurses weren't helpful" 	<p>Indirectly conveys attitude or perception toward</p> <ul style="list-style-type: none"> themselves: "I'm not very consistent about taking my medication" others: "The nurse had to poke me several times to withdraw blood"

Table 1: Description and examples of explicit and implicit functions (feeling, appreciation, judgement) in patients, representing empathic opportunities.

Clinician Role: Elicitations	
Direct	Indirect
<p>FEELINGS Inquiring directly about the patient's</p> <ul style="list-style-type: none"> emotions or mental state: "How do you feel about that?" emotive behaviors: "What was your reaction?" 	<p>Asking about experiences or emotive behaviors, where the clinician may convey interpretation which invites the patient to confirm, reject, or clarify</p> <ul style="list-style-type: none"> "So you're worried that the treatment won't work."
<p>APPRECIATION Directly asking the patient about appreciation of things, events, actions, or behaviors</p> <ul style="list-style-type: none"> event: "How was the MRI for you?" thing: "Do you find the medication helpful?" 	<p>Explores preferences and statements that convey clinician's interpretation which invites the patient to confirm, reject, or clarify</p> <ul style="list-style-type: none"> "It sounds like the medication isn't helping."
<p>JUDGEMENT Asking the patient about judgement of</p> <ul style="list-style-type: none"> self: "Do you think you are a good father?" others: "Was the nurse helpful?" 	<p>Inquire about behaviors or make statements that convey clinician's interpretation, which invites the patient to confirm, reject, or clarify</p> <ul style="list-style-type: none"> "So the nurse was not very helpful then?"

Table 2: Description and examples of explicit and implicit functions (feeling, appreciation, judgement) in clinicians in how they elicit empathic opportunities.

empathetic response. Implicit feelings can occur with the expression of judgments or appreciation.

Categories of EO Elicitations

Table 2 provides descriptions and examples of EO elicitation along the three attitudinal dimensions. Elicitations are parts of the empathic interaction in which the clinician elicits the patient's feelings and perspectives. Direct elicitation typically are questions, whereas indirect elicitation may be carried out in various ways and help clinicians avoid imposing their ideas (and potential misunderstandings) on the patient. They may soften or hedge (e.g., *perhaps, it sounds like, I have a feeling*).

Categories of Empathic Responses

The third family of spans, shown in Table 3, are the clinician's empathic response to patient EOs. The framework describes three broad types of responses: 1) Expressing explicit UNDERSTANDING or ACKNOWLEDGEMENT of the patient's feelings/views, 2) SHARING the patient's feelings/views through

expressions of agreement, and 3) expressing ACCEPTANCE in response to patients' explicit, implicit, or potential negative or positive self-judgment.

There are a variety of ways clinicians express their understanding or acknowledgement. This can overlap with indirect elicitation, as one approach is to express an interpretation/inference about the patient's views. We label these as understanding rather than an indirect elicitation because the former is more specific, and we can label the parts of the turn that show the invitation to the patient to confirm, reject, or clarify the interpretation. In the case of an implicit EO, we may observe first an understanding response followed by elicitation of more explicit patient cues.

Clinicians may express shared feelings/views through expressions that agree with the patient's attitudes (e.g., "I would also feel...", "Yes, [agree with judgment/appreciation]"). Acceptance is exhibited through principles of patient-centered care. Pounds (2011) describes two forms of acceptance:

Clinician Role: Responding to Patient Cues	
Acceptance	<p>Unconditional Positive Regard</p> <p><i>Explicit Positive Judgement</i> Expression of positive judgment of the patient as a person</p> <ul style="list-style-type: none"> • “You are a reasonable parent” • “You’re very responsible” <p><i>Implicit Positive Judgement</i> Expression of a judgement of the patients thoughts or feelings</p> <ul style="list-style-type: none"> • “It’s really great that you’ve been taking care of your parents” • “It sounds like you’ve been working hard to improve your health” <p><i>Repetition</i> Repeating or paraphrasing the patients words without countering statements or premature reassurance</p> <ul style="list-style-type: none"> • (patient says they’re worried about cancer spread) “I understand you’re worried about the cancer spreading” <p><i>Allowing Full Expression</i> Allowing patients to express feelings and views fully through minimal responses, nodding, and avoidance of interruption</p>
	<p><i>Explicit Appreciation</i> Appreciation of ideas, feelings, or behaviors regarding the patients normality or acceptability</p> <ul style="list-style-type: none"> • “It’s completely normal to be upset about this” • “It wouldn’t be surprising to feel that way” <p>Neutral Support</p> <p><i>Explicit Judgement</i> Denying negative self-assessment by the patient</p> <ul style="list-style-type: none"> • “You’re not crazy for being worried about that” • “It’s not bad to be thinking about these things”
	<p>Sharing patient views or feelings through expressed agreement</p> <ul style="list-style-type: none"> • “I’m sure I would also feel anxious in this situation if I were going through everything you have going on right now” • “Yes, this medication really does not taste good” • “Oh, no!”
Understanding	<p>Understanding and acknowledgement of the patients views and feelings can be expressed directly as acknowledgement or interpretations and may attempt to elicit explicit expressions from the patient</p> <ul style="list-style-type: none"> • “I get the sense that you found the other doctor unhelpful” • “I see that the medication isn’t working for you” • “I understand that you’re worried about infections and perhaps that makes you anxious”

Table 3: Descriptions and examples of categories of clinician responses to patient cues/empathic opportunities.

unconditional positive regard and neutral support. The former are expressions of positive judgment of the patient. Opportunities for such praise are possible when patients similarly show positive self-judgment. Neutral support can take the form of explicit appreciation or judgement; helping the patient understand their feelings are justified.

Relations form Interactional Sequences

Here, we describe different types of interactional sequences, observable via the relations drawn between empathic opportunities and elicitations or empathic responses. *Empathic response sequences* are cases when an EO is directly linked to an empathic response that explicitly recognizes the attitudes conveyed in the EO. *EO continuer sequences* involve a “potential EO continuer” that facilitates further exploration, which can lead to more explicit empathic opportunities, help the clinician gain more insight, increasing empathic understanding (Suchman et al., 1997). These are identifiable via span relationships that form a sequence of implicit EOs, elicitations, and explicit EOs. *EO terminating sequences* involve an empathic opportunity terminator, when a clinician directs the conversation away from an explicit EO in the patient’s prior turn, or a *potential* EO terminator, which occurs after implicit empathic opportunities when the clinician does not explore the implied feeling via further elicitations, instead directing the conversation away from implied cues. The scenario in §5

demonstrates examples of missed EOs.

3.2. SPIKES Protocol

Background and Motivation. SPIKES is a pedagogical tool commonly employed for training medical students’ communication skills for BBN scenarios. It provides a high-level conversation structure and communication strategies to help manage a BBN conversation compassionately while fulfilling four objectives: 1) gathering information from the patient; 2) transmitting the medical information; 3) providing support to the patient; and 4) eliciting the patient’s collaboration in developing a strategy or treatment plan for the future (Baile et al., 2000). Previous work observed significant increases in confidence in handling aspects of BBN conversations for medical students and faculty trained with the protocol. Mahendiran et al. (2023), similarly, found that it improved learner satisfaction, performance, and knowledge.

Coding SPIKES stages. The protocol contains six steps: (1) Establish a comfortable, private **SETTING**, at a time free from interruptions, and determine the participants (e.g., a patient’s support person). (2) Develop an understanding of the patient’s **PERCEPTION** of their situation (i.e., what knowledge and feelings about it do they already have). (3) Determine the amount of information the patient is ready to hear by seeking their **INVITATION** to provide

it. (4) Deliver the **KNOWLEDGE** (i.e., the information containing the bad news) clearly and compassionately. (5) Respond to the patient’s reactions with **EMPATHY**. (6) **STRATEGY** and **SUMMARY**: Discuss treatment strategies and follow-up steps, which can involve the **INVITATION** step once more to determine how/when these are discussed. Concerning the empathic response of point (5), the **SPIKES** protocol provides a 4-step guide to responding empathetically, which can involve multiple turns exchanged between the patient and clinician. These include 1) observing and 2) identifying the type of reaction or emotion the patient is experiencing (asking open questions as necessary), 3) identifying the emotion reason, and 4) responding in a way that reflects the clinician’s understanding and legitimizes the patient’s feelings.

Except **EMPATHY**, which is covered by the appraisal framework labels, each of the **SPIKES** stages is labeled on clinician turns or segments of the clinician turns when identifiable. We note that **SETTING** is rarely applied given that much of the behaviors involved in this stage occur prior to the BBN conversation.

3.3. Annotation Procedure

Platform Integration. Our scheme is integrated in **INCEPTION** (Klie et al., 2018), an open-source tool for semantic annotations that supports labeling text spans and relations.² We setup custom layers for each. The annotator highlights a text segment and selects a coarse-grained span category (*EOs*, *EO elicitations*, *EO responses*, and *SPIKES*) for the span. This opens tags representing the fine-grained labels for the selected category which the annotator then applies as the span label.³

Annotator Training. Two native German speakers were trained to perform the annotations in three 1-hour sessions. The training included background on SFL, a tutorial on the full coding manual, additional training materials involving real samples for demonstration and practice applying the coding scheme, and a tutorial on the annotation tool.⁴ We explain how to approach the analysis and labeling throughout example scenarios in §5.

After the training sessions, the annotators performed three calibration rounds on dialogues that they coded independently. We met as a group to discuss the source of disagreements between annotators. The primary source in the first round

²<https://inception-project.github.io/>

³We provide a template project for our scheme in our code repository to make it available and simple to deploy on **INCEPTION** for other researchers.

⁴The coding manual and training materials are included in our code repository.

was distinguishing between **JUDGMENTS** and **APPRECIATIONS**, which we clarified and added further examples with explanations to the coding manual. After this, the agreement improved in the subsequent calibration rounds. In those rounds, the main challenge related to identifying the implicit **EOs**. Some disagreements on this aspect are reasonable, as it requires the annotator to make their own inferences, which can vary subjectively. However, during our discussions, we reached a consensus for most implicit **EOs**. We met weekly to review coding conventions, clarify questions, and discuss specific cases throughout the months of the annotations. We present an analysis of interannotator agreement post-training and calibration in the next section.

4. The BBN **EMPATHY** Dataset

The BBN **EMPATHY** Dataset is constructed via a collaboration with medical didactics experts at a large medical school. The dialogues are practice BBN conversations between medical students taking part in a medical didactics seminar and trained standardized patient actors. These simulate BBN scenarios as realistically as possible for student practice. They take place in rooms modeling medical environments, such as a doctor’s office or a hospital bed. During the seminar, the students are trained in BBN communication skills with the **SPIKES** Protocol (§3.2). The standardized patients are provided a role description and background for the scenario, and the students are provided a full scenario description and patient history. The scenarios include delivering cancer diagnoses, failures of treatments for serious diseases, and informing a family member of a severe accident, among others. We collected a total of 63 conversations in German over two semesters of the seminar. Four trained native German speakers transcribe them (see Appendices **A** and **B** for transcription and annotator details). Though these are practice conversations, we anonymize the names of the participants.

4.1. Agreement Study

After training and the calibration rounds, the annotators independently coded eight dialogues. For these, we study the interannotator agreement on the text spans and span labels.

Text span agreement. We measure the interannotator agreement on the text spans labeled by the annotators at the string level based on Wiebe et al. (2005)’s approach shown in Equation 1.

$$agr(a||b) = \frac{|A \text{ matching } B|}{|A|} \quad (1)$$

A and B are the sets of spans highlighted by annotator a and b respectively. $agr(a||b)$ is the proportion of text marked by annotator a that b marked, and $agr(b||a)$ is the proportion of text marked by b that a marked. agr is the mean of both measures.

As with Wiebe et al. (2005)’s span labeling task, agr is a suitable metric for text span agreement for our task because we do not employ nor instruct strict rules about the precision of the text boundaries. The main concern, rather, is whether the annotators mark the same general expression. Example (1) shows a case where for turn t , annotator A marked an additional clause not marked by B , but the expression is generally the same in terms of the appraisals they convey.

(1) A_t : Ähm, okay, Operation? **Okay, also das, das muss raus, oder was?**

B_t : **Okay, also das, das muss raus, oder was?**

English: Um, okay, surgery? Okay, so that, that has to come out, or what?

Here, $agr(a||b) = 0.65$ and $agr(b||a) = 1.0$. In example (2), $agr(a||b) = 0.94$ and $agr(b||a) = 0.86$. Bold indicates the text was marked by both annotators and bold and italic are spans marked by only one annotator.

(2) A_t : **Ähm, okay, kann ich, vielleicht Wasser, oder irgendwie... € (inaudible) [Patient crying]** Oh fuck. (inaudible) Also kann ich dann wieder arbeiten, oder was? (inaudible) **[laughing desperately] Ich weiß nicht, ob Sie das verstehen, aber wenn ich nicht arbeiten kann, dann...**

B_t : Ähm, okay, **kann ich, vielleicht Wasser, oder irgendwie... € (inaudible) [Patient crying]** Oh fuck. (inaudible) Also kann ich dann wieder arbeiten, oder was? (inaudible) **[laughing desperately] Ich weiß nicht, ob Sie das verstehen, aber wenn ich nicht arbeiten kann, dann...**

English: Um, okay, can I have, maybe water, or something... Oh fuck. So can I then go back to work, or what? I don’t know if you understand but if I can’t work, then...

First, we compute agr for each turn. Then, we take the average across all turns to get the agr for a dialogue. Table 4 shows the agr metrics for each dialogue and the means across all eight. The agreement on all spans improved notably between dialogue 0 and 1, which is due to our continued discussions after each dialogue annotation which focused on general clarifications about the scheme and coding conventions rather than resolving disagreements. In our analysis, we observed small differences between the annotator’s choice to include punctuations and short subclauses in the annotations. Overall, the annotators match the same general expression.

Span label agreement. As a first point of reference, we computed Krippendorff’s α for all labels (strict) and report them in Table 4. We also study agreement on the span labels by computing agr

Dialogue	agr	$agr(a b)$	$agr(b a)$	α
Calibration Dialogues				
1	.789	.738	.840	.30
2	.910	.850	.970	.43
3	.902	.844	.960	.47
4	.952	.923	.981	.50
5	.931	.865	.997	.58
6	.912	.935	.890	.63
7	.905	.948	.862	.85
8	.948	.933	.963	.89
mean	.92 ± .02	.90 ± .04	.95 ± .04	
All Dialogues				
mean	.97 ± .03	.96 ± .05	.98 ± .03	

Table 4: Span text agreement for each of the eight calibration dialogues and the mean agreement with standard deviations across all 63 dialogues. The right column shows the span label agreement measured by Krippendorff’s α .

	agr	$agr(a b)$	$agr(b a)$
All labels			
\cup	0.71	0.72	0.69
\cap	0.73	0.74	0.73
Coarse-grained labels			
\cup	0.87	0.90	0.85
\cap	0.90	0.91	0.90
Fine-grained labels			
\cup	0.59	0.59	0.59
\cap	0.61	0.60	0.62
Attitudes			
\cup	0.69	0.70	0.69
\cap	0.72	0.71	0.72
SPIKES			
	0.63	0.67	0.59

Table 5: Mean interannotator agreements on span labels for the eight calibration dialogues. \cup indicates agr over all spans, whereas \cap indicates agr only on spans that have overlap.

for each type of label; 1) All labels include all fine-grained AF labels and SPIKES labels; 2) Coarse-grained labels are the three AF categories (EO, EO elicitation, EO response); 3) Fine-grained labels include the attitudes and explicit/implicit for EOs, direct/indirect for EO elicitations, and each fine-grained category of EO responses; 4) Attitudes only include the *feelings*, *judgement*, and *appreciation* labels (i.e., combining explicit and implicit, and direct and indirect); and 5) only SPIKES labels. $agr(a||b)$ represents annotator b ’s precision evaluated against annotator a ’s labels and $agr(b||a)$ is a ’s with respect to b . We compute these metrics by a *strict* evaluation denoted by \cup , which includes spans where there was no overlap between annotators, and by a *matched* evaluation denoted by \cap ,

which includes only spans with some overlap.

As noted, the annotators improved their text span agreement after dialogue 1 following further clarifications. We observed a stark contrast in the agreements between dialogue 1 and the other dialogues. Dialogue 1 had very low *agr* on the span labels: \cap *agr* for *all labels* was 0.5; AF coarse-grained and SPIKES \cap agreements were 0.85 and 0.60 respectively; the rest ranged from 0.22 (AF fine \cap) to 0.38 (AF attitude \cap). Meanwhile, the label agreements across the rest improved; Table 5 shows the mean values. We find that the annotators generally perform well at matching each other. Disagreement was mostly between implicit feelings and other implicit labels or explicit judgement. Implicit responses are more subjective or difficult to interpret. For SPIKES, the most common disagreement was between the invitation and knowledge steps, as determining how much the patient is ready to hear may be part of the step of delivering that information. After we reached substantial agreement across all categories, the remaining dialogues are annotated by one annotator and revised for quality by the other (see Table 9 in Appendix C for overall agreement).

Patient Role: Empathic Opportunities				
	Explicit		Implicit	
	A	B	A	B
Feelings	248	203	1043	1146
Appreciation	608	618	180	178
Judgement	874	991	135	56

Clinician Role: Empathic Elicitations				
	Direct		Indirect	
	A	B	A	B
Feelings	109	113	55	107
Appreciation	203	172	96	89
Judgement	151	149	68	40

Table 6: Counts of each type of Empathic Opportunity (EO) label and each type of EO Elicitation label from each annotator (*A* and *B*). For EOs, both annotators identified significantly more *implicit* rather than *explicit* feelings, whereas they identified more *explicit* than *implicit* appreciations and judgements. *Implicit feeling* is the most common type of EO. Both annotators identify more *direct* rather than *indirect* EO elicitations.

4.2. Dataset statistics

The label distributions for EOs and EO Elicitations are presented in Table 6 and for EO responses in Table 7. For the Patient EOs, we observe that implicit feelings are much more common than explicit feelings in line with Suchman et al. (1997)’s findings, whereas the opposite is the case for judgment and appreciation. However, this is consistent with Pounds (2011)’s observation that implicit feelings are often identified within explicit judgments and appreciations.

Clinician Role: Empathic Responses		
	ACCEPTANCE	
	A	B
UNCOND. POSITIVE REGARD		
Explicit Positive Judgement	374	454
Implicit Positive Judgement	143	109
Repetition - no counter	31	27
Allowing Full Expression	156	184
NEUTRAL SUPPORT	A	B
Explicit Appreciation	373	405
Explicit Judgement	185	141

SHARING FEELINGS AND VIEWS		
	A	B
Feelings	40	42
Appreciation	50	59
Judgement	44	60

UNDERSTANDING FEELINGS AND VIEWS		
	A	B
Feelings	427	462
Appreciation	75	56
Judgement	167	170

Table 7: Counts of each type of Empathic Response label from each annotator (*A* and *B*). With all sublabels, *unconditional positive regard* is the most common response type, the most frequent among them being *explicit positive judgement* and *explicit appreciation*. *Understanding* rather than *sharing* feelings and views is more common; here, *understanding feelings* is most frequently observed.

EO type	No Response	EO Response
Exp. Feel	35.4	64.6
Imp. Feel	27.8	72.2
Exp. Appreciate	46.9	53.1
Imp. Appreciate	36.5	63.5
Exp. Judge	34.3	65.7
Imp. Judge	27.3	72.7

Table 8: The percentage of EOs by EO type that had *no response* or an *EO response* linked by a relation. The rate of responses is higher for *implicit EOs* than *explicit EOs*. The percentages broken down by response type are shown in Figure 1 in Appendix C.

We quantify the relations between patient EOs and clinician responses, showing the percentage of EOs that are linked to responses in Table 8. Interestingly, we observe that the *implicit EOs* have a higher rate of clinician responses than *explicit EOs*. *Explicit appreciations* and *judgements* are more frequently identified than *implicit* ones. As the annotators remarked on the difficulty with identifying appreciations and judgements compared to feelings, it may be that observing the clinicians’ responses aids the annotator in observing the implicit EOs, thus biasing the relation rates. However, the higher response rate to *implicit* EOs is also the case for *feelings*, for which *implicit* cases are more frequently marked. Future work could investigate

the possible effects further by testing the annotation scheme in a setup in which the identified EOs are locked before observing the rest of the dialogue. Figure 1 in Appendix C shows the percentages broken down by each EO response type. We observe higher proportions of relations between EOs and EO responses of the same attitude type.

5. How do I tell my family?

Here we present an illustrative example with our annotation scheme.

Context. A physician informs a patient that an MRI, performed as part of an anonymous research study, detected a mass in their brain. The patient is in disbelief, as they did not expect a follow-up and had no previous cause for concern, suggesting that the physician had mistaken the results for someone else's.

Scenario. The physician allows the patient to experience shock and express disbelief implicitly by denying that the results could indeed be for them (IMPLICIT FEELING EO). The physician identifies the emotion of disbelief and responds, "You can't believe it quite yet, I have a feeling." The first clause is the physician's interpretation of how the patient feels (UNDERSTANDING/ACKNOWLEDGEMENT). The second serves to soften the delivery of the interpretation and to formulate the statement as an *indirect elicitation*: FEELING and/or APPRECIATION.

The patient speaks more openly, sharing that everything was going great for them. They state, "It is probably not easy. If there's really something there, it won't just go away on its own." This could signal that the patient is accepting the news. One may interpret this as an *implicit EO*, inferring a *negative APPRECIATION* that treatment will be difficult and a *negative FEELING* (anxiety/fear). The physician acknowledges the patient's view, confirming treatment probably will not be easy (*empathic response*: ACCEPTANCE, neutral support). The patient asks if, theoretically, one can die from such a mass; an *implicit EO*: FEELING (fear of severity/uncertainty). The physician confirms this but says further analysis must clarify the type and severity. This balances the SPIKES/BBN principles of *honesty* and *lending hope*, delivering KNOWLEDGE clearly and compassionately.

Later, the patient shares that their aunt had had a brain tumor a few years earlier, describing a quick escalation that was painful and fatal. The loss weighs heavily on the family. We consider this turn to have *implicit EOs*: *negative FEELING* (signaling fear/worry) and APPRECIATION (of this significant event in the family). The physician responds, saying it does not mean the patient's case will be the same. While aspects of this response reflect SPIKES/BBN principles (e.g., *lending hope*,

attempting to *reduce a sense of isolation*), it is also an *EO terminator* since it directs the conversation away from the implicit EOs. The patient's underlying perspectives that this story communicates become clearer as the dialogue continues, suggesting there were indeed missed EOs.

After the physician's response, the patient immediately expresses anxiety about telling their family, asking how to break the news, implying concern for how the family will react. The physician misses this EO, responding "you can figure that out for yourself." The patient asks again what to do, saying that they cannot simply go home and tell their family they might have cancer (*implicit EOs*: *negative FEELING* (anxiety) and APPRECIATION (anticipating significant difficulty in informing their family)). The physician says nothing, perhaps allowing space for the patient's emotive behaviors. As the conversation continues, the patient continues to express the same sentiment about their aunt and not knowing how to tell their family, EOs that are repeatedly missed by the physician.

6. Conclusions and Future Work

We pursued three core objectives toward modeling clinical empathy in patient-clinician conversations: 1) Develop an annotation scheme for labeling precise discourse elements and their relations in clinical empathy encounters; 2) Produce a structural representation of the dynamics of these elements over a full conversation toward addressing current shortcomings of NLP models for empathic interactions; and 3) Tailor the approach for medical didactics research for training empathic communication skills in BBN conversations. We developed a span-relation labeling method based on models of interactional sequences and semantic exchanges in clinical empathy conversations. This establishes links between empathic opportunities and empathic responses, enabling analysis of types of interactional sequences that achieve empathy. In addition, we produce the BBN EMPATHY dataset, the first of its kind curating discourse-level annotations tailored to clinical empathy, which is publicly available for fellow researchers. These contributions will foster open research and interdisciplinary collaboration addressing critical aspects of empathic communication in healthcare contexts.

Future work will explore the linguistic components of the discourse elements through computational linguistic analyses. We also plan to investigate dynamic models of the interactional sequences to study their impacts on empathic understanding. In addition, future work can investigate NLP models trained on the BBN EMPATHY dataset for supporting scaling up such resources in collaboration with trained human annotators.

7. Acknowledgements

We sincerely thank research assistants Lilly Beil and Lilly Metten for their high-quality work on both the transcriptions and annotations, Maceo Kita and Lea Fischbach for their assistance in transcribing the dataset, and Ulvi Shukurzade for his support with the annotation tool. This work has been supported by the German Federal Ministry of Education and Research (BMBF) as a part of the Junior AI Scientists program under reference 01-S20060 and by the Humboldt Foundation through the Fellowship of Dr. Welch.

Limitations

We acknowledge there may be unknown effects that the simulated scenarios have on the dialogues. Further research could investigate artifacts of the simulated scenarios and how this might affect future approaches. Although we see progress in applying the SPIKES protocol and its pedagogical benefits, there is still room for improvement, especially as SPIKES does not specify higher-level aspects of the interaction or the patient's role in the conversation, which appraisal framework for clinical empathy helps address.

The Patient Perspective. Here, we discuss insights from studies on the patient perspective that may imply limitations of SPIKES that may warrant further research. Assessing patients' preferences for BBN communication and their perception and satisfaction of actual BBN disclosures, Seifart et al. (2014) administered the Marburg Breaking Bad News Scale (MABBAN), a questionnaire based on the SPIKES protocol, to 350 cancer patients. They observe that only 46.2% of the patients are fully satisfied by how the bad news was broken to them and that there is a highly significant discrepancy between the patients' preferences for receiving bad news and the actual disclosure. Furthermore, they find that the overall patient's satisfaction with the first BBN disclosure significantly correlates with their emotional state, including depression, anxiety, and sleeplessness, after receiving the bad news. von Blanckenburg et al. (2020) later administered MABBAN to 336 cancer patients. Analyzing its psychometric properties, they observed an accordance between the SPIKES protocol and the MABBAN scale, suggesting that SPIKES meets the preferences of German cancer patients. The study emphasizes that differentiated communication of BBN is highly important due to discrepancies in patient preferences.

Ethics Statement

This research and dataset were evaluated and approved by an IRB at Philipps-University Marburg, Germany, for the purpose of the initial annotation study and NLP experiments. For privacy considerations, the examples are adaptations of observed exchanges of parts of BBN dialogues, and all content is significantly summarized or paraphrased, including quotes (i.e., they do not portray the full context, nor are any part of the scenarios used directly verbatim).

8. Bibliographical References

- Walter F. Baile, Robert Buckman, Renato Lenzi, Gary Gloger, Estela A. Beale, and Andrzej P. Kudelka. 2000. [SPIKES—A Six-Step Protocol for Delivering Bad News: Application to the Patient with Cancer](#). *The Oncologist*, 5(4):302–311.
- Collin F Baker, Charles J Fillmore, and John B Lowe. 1998. The berkeley framenet project. In *COLING 1998 Volume 1: The 17th International Conference on Computational Linguistics*.
- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen tau Yih, and Yejin Choi. 2020. [Abductive commonsense reasoning](#). In *International Conference on Learning Representations*.
- Kathleen A Bonvicini, Michael J Perlin, Carma L Bylund, Gregory Carroll, Ruby A Rouse, and Michael G Goldstein. 2009. Impact of communication training on physician expression of empathy in patient encounters. *Patient education and counseling*, 75(1):3–10.
- Michael Borish, Andrew Cordar, Adriana Foster, Thomas Kim, James Murphy, and Benjamin Lok. 2014. Utilizing real-time human-assisted virtual humans to increase real-world interaction empathy. *Kansei Engineering & Emotion Research (KEER'14)*, 15.
- Sven Buechel, Anneke Buffone, Barry Slaff, Lyle Ungar, and João Sedoc. 2018. [Modeling empathy and distress in reaction to news stories](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4758–4765, Brussels, Belgium. Association for Computational Linguistics.
- Alok Debnath and Owen Conlan. 2023. [A critical analysis of empathetic dialogues as a corpus for empathetic engagement](#). In *Proceedings of*

- the 2nd Empathy-Centric Design Workshop, EMPATHICH '23, New York, NY, USA. Association for Computing Machinery.
- Priyanka Dey and Roxana Girju. 2022. [Enriching deep learning with frame semantics for empathy classification in medical narrative essays](#). In *Proceedings of the 13th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 207–217, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Priyanka Dey and Roxana Girju. 2023. [Investigating stylistic profiles for the task of empathy classification in medical narrative essays](#). In *Proceedings of the First International Workshop on Construction Grammars and NLP (CxGs+NLP, GURT/SyntaxFest 2023)*, pages 63–74, Washington, D.C. Association for Computational Linguistics.
- Mauajama Firdaus, Asif Ekbal, and Pushpak Bhattacharyya. 2020. [Incorporating politeness across languages in customer care responses: Towards building a multi-lingual empathetic dialogue agent](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4172–4182, Marseille, France. European Language Resources Association.
- Roxana Girju and Marina Girju. 2022. [Design considerations for an NLP-driven empathy and emotion interface for clinician training via telemedicine](#). In *Proceedings of the Second Workshop on Bridging Human–Computer Interaction and Natural Language Processing*, pages 21–27, Seattle, Washington. Association for Computational Linguistics.
- Michael Alexander Kirkwood Halliday and Christian MIM Matthiessen. 2013. *Halliday's introduction to functional grammar*. Routledge.
- Mahshid Hosseini and Cornelia Caragea. 2021. [Distilling knowledge for empathy detection](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3713–3724, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ian Hsu, Somnath Saha, Phillip Todd Korthuis, Victoria Sharp, Jonathon Cohn, Richard D. Moore, and Mary Catherine Beach. 2012. [Providing support to patients in emotional encounters: A new perspective on missed empathic opportunities](#). *Patient Education and Counseling*, 88(3):436–442. Patients, providers, and relationships in health care: investigations from the ICCH 2011 conference in Chicago.
- Zac E Imel, Derek D Caperton, Michael Tanana, and David C Atkins. 2017. Technology-enhanced human interaction in psychotherapy. *Journal of counseling psychology*, 64(4):385.
- Jan-Christoph Klie, Michael Bugert, Beto Boulosa, Richard Eckart de Castilho, and Iryna Gurevych. 2018. [The inception platform: Machine-assisted and knowledge-oriented interactive annotation](#). In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 5–9. Association for Computational Linguistics. Event Title: The 27th International Conference on Computational Linguistics (COLING 2018).
- Allison Lahnala, Charles Welch, David Jurgens, and Lucie Flek. 2022. [A critical reflection and forward perspective on empathy and natural language processing](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2139–2158, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Andrew Lee, Jonathan Kummerfeld, Larry An, and Rada Mihalcea. 2023. [Empathy identification systems are not accurately accounting for context](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1686–1695, Dubrovnik, Croatia. Association for Computational Linguistics.
- Wendy Levinson, Rita Gorawara-Bhat, and Jennifer Lamb. 2000. [A Study of Patient Clues and Physician Responses in Primary Care and Surgical Settings](#). *JAMA*, 284(8):1021–1027.
- Benjamin Lok and Adriana E. Foster. 2019. [Can Virtual Humans Teach Empathy?](#), pages 143–163. Springer International Publishing, Cham.
- Meera Mahendiran, Herman Yeung, Samantha Rossi, Houman Khosravani, and Giulia-Anna Perri. 2023. Evaluating the effectiveness of the spikes model to break bad news—a systematic review. *American Journal of Hospice and Palliative Medicine*®.
- J. R. Martin and P. R. R. White. 2005. *The Language of Evaluation*. Palgrave Macmillan UK, London.
- Laura A Michaelis. 2006. Construction grammar. *The encyclopedia of language and linguistics*, 3:73–84.
- Melanie Neumann, Jozién Bensing, Stewart Mercer, Nicole Ernstmann, Oliver Ommen, and Holger Pfaff. 2009. [Analyzing the “nature” and “specific effectiveness” of clinical empathy: A theoretical overview and contribution towards a theory-based research agenda](#). *Patient Education and*

- Counseling*, 74(3):339–346. Theories in Health Communication Research.
- Verónica Pérez-Rosas, Rada Mihalcea, Kenneth Resnicow, Satinder Singh, and Lawrence An. 2017. [Understanding and predicting empathic behavior in counseling therapy](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1426–1435, Vancouver, Canada. Association for Computational Linguistics.
- Verónica Pérez-Rosas, Xuetong Sun, Christy Li, Yuchen Wang, Kenneth Resnicow, and Rada Mihalcea. 2018. [Analyzing the quality of counseling conversations: the tell-tale signs of high-quality counseling](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Verónica Pérez-Rosas, Xinyi Wu, Kenneth Resnicow, and Rada Mihalcea. 2019. [What makes a good counselor? learning to distinguish between high-quality and low-quality counseling conversations](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 926–935, Florence, Italy. Association for Computational Linguistics.
- Gabrina Pounds. 2011. Empathy as "appraisal": developing a new language-based approach to the exploration of clinical empathy. *Journal of Applied Linguistics and Professional Practice*, 7(2):139–162.
- Jack Pun. 2021. A study of chinese medical students' communication pattern in delivering bad news: an ethnographic discourse analysis approach. *BMC Medical Education*, 21(1):286.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. [Towards empathetic open-domain conversation models: A new benchmark and dataset](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy. Association for Computational Linguistics.
- Greg M Reger, Aaron M Norr, Michael A Gramlich, and Jennifer M Buchman. 2021. Virtual standardized patients for mental health education. *Current Psychiatry Reports*, 23:1–7.
- Elena Rey Velasco, Hanne Sæderup Pedersen, Timothy Skinner, et al. 2022. Analysis of patient cues in asynchronous health interactions: Pilot study combining empathy appraisal and systemic functional linguistics. *JMIR Formative Research*, 6(12):e40058.
- Carola Seifart, Mareike Hofmann, Tobias Bär, J Ri-era Knorrenschild, Ulf Seifart, and Winfried Rief. 2014. Breaking bad news—what patients want and what they get: evaluating the spikes protocol in germany. *Annals of Oncology*, 25(3):707–711.
- Ashish Sharma, Inna W Lin, Adam S Miner, David C Atkins, and Tim Althoff. 2021. Towards facilitating empathic conversations in online mental health support: A reinforcement learning approach. In *Proceedings of the Web Conference 2021*, pages 194–205.
- Ashish Sharma, Adam Miner, David Atkins, and Tim Althoff. 2020. [A computational approach to understanding empathy expressed in text-based mental health support](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5263–5276, Online. Association for Computational Linguistics.
- Siqi Shen, Charles Welch, Rada Mihalcea, and Verónica Pérez-Rosas. 2020. [Counseling-style reflection generation using generative pretrained transformers with augmented context](#). In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 10–20, 1st virtual meeting. Association for Computational Linguistics.
- Shuju Shi, Yinglun Sun, Jose Zavala, Jeffrey Moore, and Roxana Girju. 2021. [Modeling clinical empathy in narrative essays](#). In *2021 IEEE 15th International Conference on Semantic Computing (ICSC)*, pages 215–220.
- Roger W Squier. 1990. A model of empathic understanding and adherence to treatment regimens in practitioner-patient relationships. *Social science & medicine*, 30(3):325–339.
- Anthony L. Suchman, Kathryn Markakis, Howard B. Beckman, and Richard Frankel. 1997. [A Model of Empathic Communication in the Medical Interview](#). *JAMA*, 277(8):678–682.
- Ekaterina Svikhushina, Iuliana Voinea, Anuradha Welivita, and Pearl Pu. 2022. [A taxonomy of empathetic questions in social dialogs](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2952–2973, Dublin, Ireland. Association for Computational Linguistics.

- Michael J Tanana, Christina S Soma, Vivek Sriku-mar, David C Atkins, and Zac E Imel. 2019. De-velopment and evaluation of clientbot: Patient-like conversational agent to train basic counseling skills. *Journal of medical Internet research*, 21(7):e12529.
- Pia von Blanckenburg, Mareike Hofmann, Winfried Rief, Ulf Seifart, and Carola Seifart. 2020. [Assessing patients’ preferences for breaking bad news according to the spikes-protocol: the mab-ban scale](#). *Patient Education and Counseling*, 103(8):1623–1629.
- Thiemo Wambsganss, Christina Niklaus, Matthias Söllner, Siegfried Handschuh, and Jan Marco Leimeister. 2021. [Supporting cognitive and emotional empathic writing of students](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4063–4077, Online. Association for Computational Linguistics.
- Thiemo Wambsganss, Matthias Soellner, Kenneth R Koedinger, and Jan Marco Leimeister. 2022. [Adaptive empathy learning support in peer review scenarios](#). In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI ’22, New York, NY, USA. Association for Computing Machinery.
- Anuradha Welivita and Pearl Pu. 2020. [A taxonomy of empathetic response intents in human social conversations](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4886–4899, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39:165–210.
- Rolf Wynn and Michael Wynn. 2006. [Empathy as an interactionally achieved phenomenon in psychotherapy: Characteristics of some conversational resources](#). *Journal of Pragmatics*, 38(9):1385–1397. Focus-on Issue: The Pragmatics of Failure and Success.
- Justine Zhang and Cristian Danescu-Niculescu-Mizil. 2020. [Balancing objectives in counseling conversations: Advancing forwards or looking backwards](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5276–5289, Online. Association for Computational Linguistics.
- Justine Zhang, Sendhil Mullainathan, and Cristian Danescu-Niculescu-Mizil. 2020. Quantifying the causal effects of conversational tendencies. *Proceedings of the ACM on Human-Computer Interaction*, 4(CSCW2).
- Wenting Zhao, Justin Chiu, Claire Cardie, and Alexander Rush. 2023. [Abductive commonsense reasoning exploiting mutually exclusive explanations](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14883–14896, Toronto, Canada. Association for Computational Linguistics.

A. Transcriptions

The time it takes to transcribe a full conversation is very dependent on the audio quality of the file. With prioritizing simulating a realistic setting for the students’ practice, the microphone position reduces the audio quality. Some transcribers report that in especially low quality circumstances in which the voices are muffled and unclear, the time it takes to complete transcribing the dialogue can amount to 12 hours, since the writing, rewinding, correcting, rewinding, repeat, can take a lot of time.

To aid the transcription labor, we searched for effective ASR tools for German that can be run offline on a private server (to protect the data) and handle low quality. We found setting up Whisper⁵ (Radford et al., 2023) offline was the most effective, and integrated it as a starting point for the transcribers.

In the best-case scenario, with a quality Whisper base script, the task consists mostly of setting the correct timestamps in the right places, correcting the occasional misinterpretation, and filling in the parts that Whisper failed to recognize at all. The task of transcribing the dialogue can be done in two to three hours. Even so, it still depends on the audio quality, which impacts the quality of Whisper transcriptions.

B. Annotator Details

The annotation task is complex and non-trivial, requiring dedicated time and work effort to properly reason through the types of spans. This requires expert annotators who have a comprehensive understanding of the theoretical frameworks, the dialogue setting, and experience applying the framework. The dialogues are privacy protected, so we carefully chose a small group of research assistants

⁵<https://github.com/openai/whisper>

	<i>agr</i>	<i>agr(a b)</i>	<i>agr(b a)</i>
All labels			
U	0.88	0.88	0.87
∩	0.89	0.89	0.90
Coarse-grained labels			
U	0.95	0.96	0.94
∩	0.97	0.96	0.97
Fine-grained labels			
U	0.81	0.81	0.81
∩	0.82	0.81	0.83
Attitudes			
U	0.85	0.85	0.85
∩	0.87	0.86	0.88
SPIKES	0.90	0.91	0.89

Table 9: Mean interannotator agreements on span labels across all 63 dialogues. U indicates a strict evaluation over all spans, including those where the annotators had no overlap, whereas ∩ indicates evaluation only on spans that have overlap.

to whom access is granted. Therefore, the people who perform the dialogue transcriptions also generally perform the annotations. All are L1 German speakers. Across all transcribers and annotators, their expertise and concentrations include psychology, political science, linguistics, informatics, and physics.

C. Additional Reference

Table 9 shows the agreements over all dialogues. Note only the eight dialogues discussed in §4.1 were annotated independently. Otherwise, they were annotated first by one annotator, and reviewed by the other. Table 10 shows the average and standard deviation of the label counts per dialogue.

Label	Avg ±std
<i>Patient EO</i>	33.9 ±14.6
exp.appreciate	5.9 ±5.0
exp.feel	2.4 ±2.3
exp.judge	9.5 ±6.0
imp.appreciate	2.4 ±2.6
imp.feel	12.5 ±5.6
imp.judge	1.2 ±1.7
<i>EO Response</i>	19.1 ±9.3
accept.nt.appreciate	3.8 ±2.5
accept.nt.judge	1.7 ±1.8
accept.pos.allow	1.7 ±1.5
accept.pos.exp.judge	3.5 ±2.8
accept.pos.imp.judge	1.2 ±1.7
accept.pos.repeat	0.4 ±1.0
appreciate.share	0.5 ±1.0
appreciate.understand	0.6 ±1.0
feel.share	0.3 ±0.7
feel.understand	3.4 ±2.4
judge.share	0.5 ±1.1
judge.understand	1.4 ±1.6
<i>EO elicitation</i>	6.8 ±6.4
dir.appreciate	1.8 ±3.0
dir.feel	1.2 ±1.2
dir.judge	1.6 ±1.4
ind.appreciate	0.9 ±1.6
ind.feel	0.8 ±1.7
ind.judge	0.5 ±1.0
<i>SPIKES</i>	15.1 ±5.6
invitation	1.8 ±1.5
knowledge	5.0 ±3.0
perception	2.3 ±2.2
setting	0.7 ±0.6
strategy/summary	5.2 ±3.3

Table 10: Average and standard deviation of label counts per dialogue.

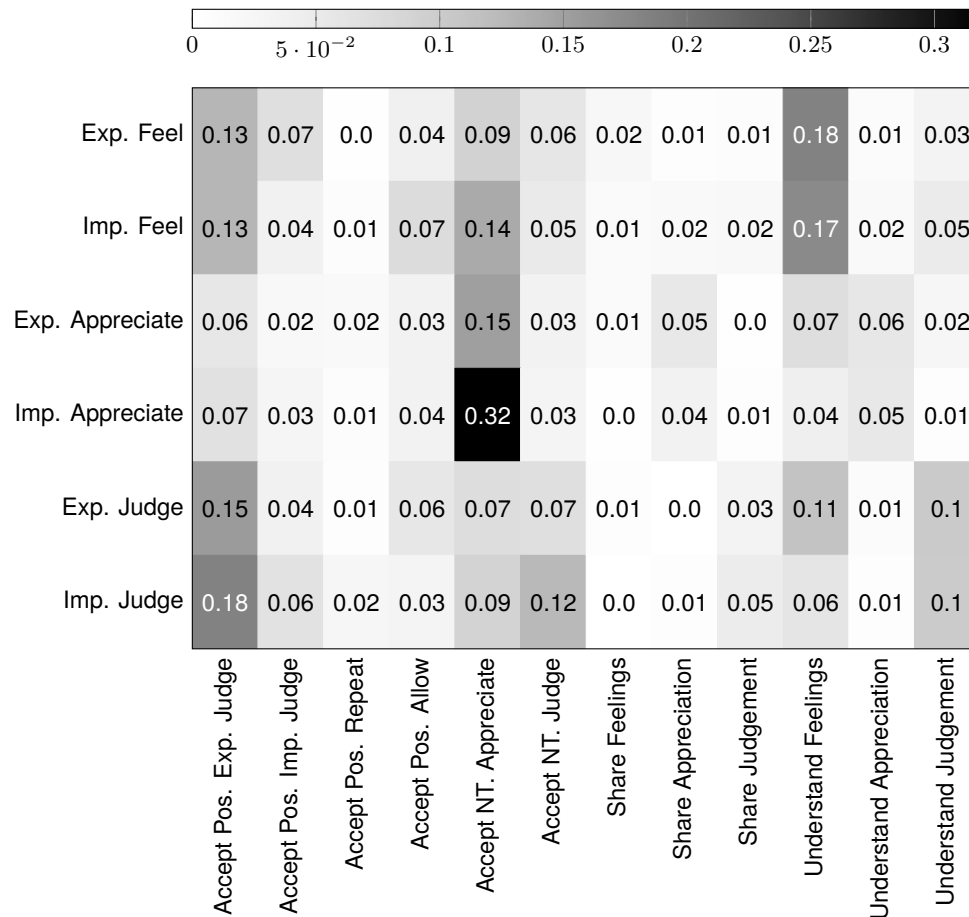


Figure 1: **Relation heatmap.** This heatmap reflects relations between patient EOs (rows) and subsequent EO responses (columns). The cell values reflect the percentage of the EOs of the type specified by the row that was responded to with the EO response type specified by the column.