
AN EVALUATION OF GENERATIVE PRE-TRAINING MODEL-BASED THERAPY CHATBOT FOR CAREGIVERS

Lu Wang

College of Computing & Informatics
Drexel University
Philadelphia, PA 19104
lw823@drexel.edu

Munif Ishad Mujib

College of Computing & Informatics
Drexel University
Philadelphia, PA 19104
mim52@drexel.edu

Jake Williams

College of Computing & Informatics
Drexel University
Philadelphia, PA 19104
jw3477@drexel.edu

George Demiris

School of Nursing, Perelman School of Medicine
University of Pennsylvania
Philadelphia, PA 19104
gdemiris@nursing.upenn.edu

Jina Huh-Yoo

College of Computing & Informatics
Drexel University
Philadelphia, PA 19104
jh3767@drexel.edu

July 29, 2021

ABSTRACT

With the advent of off-the-shelf intelligent home products and broader internet adoption, researchers increasingly explore smart computing applications that provide easier access to health and wellness resources. AI-based systems like chatbots have the potential to provide services that could provide mental health support. However, existing therapy chatbots are often retrieval-based, requiring users to respond with a constrained set of answers, which may not be appropriate given that such pre-determined inquiries may not reflect each patient's unique circumstances. Generative-based approaches, such as the OpenAI GPT models, could allow for more dynamic conversations in therapy chatbot contexts than previous approaches. To investigate the generative-based model's potential in therapy chatbot contexts, we built a chatbot using the GPT-2 model. We fine-tuned it with 306 therapy session transcripts between family caregivers of individuals with dementia and therapists conducting Problem Solving Therapy. We then evaluated the model's pre-trained and the fine-tuned model in terms of basic qualities using three meta-information measurements: the proportion of non-word outputs, the length of response, and sentiment components. Results showed that: (1) the fine-tuned model created more non-word outputs than the pre-trained model; (2) the fine-tuned model generated outputs whose length was more similar to that of the therapists compared to the pre-trained model; (3) both the pre-trained model and fine-tuned model were likely to generate more negative and fewer positive outputs than the therapists. We discuss potential reasons for the problem, the implications, and solutions for developing therapy chatbots and call for investigations of the AI-based system application.

Keywords Therapy chatbot · The OpenAI GPT-2 · Speech quality · Evaluation · Generative-based

1 Introduction

Acute mental health provider shortages alarmed the Association of American Medical Colleges in 2016 [1]. Eighteen percent of the counties in the U.S. had reported a lack of providers, such as psychologists, social workers, advanced practitioners, therapists, and counselors [2]. Nearly every county (e.g., 96% of 3,140 counties in the U.S.) needed more psychiatrists [2]. Rural counties and low-income counties have higher levels of need concerning healthcare access [2]. This issue is a global problem to this day. In developed countries, 100,000 people approximately share 6.6 psychiatrists, while, in lower-income countries, 100,000 people approximately share as few as 0.1 psychiatrists [3]. An estimated additional 1.18 million mental health workers are needed in low- and middle-income countries for basic mental health interventions [3]. Developing cost-effective and feasible mental health support systems to mitigate these shortages will be critical.

Chatbots are digital tools that interact with users in natural language, often for the goal of helping a user complete a task, [4, 5]. Chatbots, in the form of voice agents with speech functions became widely available to the public through off-the-shelf products, such as Apple’s Siri, Microsoft’s Cortana, and Amazon’s Alexa [6]. Having been adopted in multiple fields, including business [7], governance [8], and education [9], the potential of chatbots is also presented in mental health, where chatbots are labeled as “the future of therapy” [5] and are believed to increase the opportunities for individuals to receive therapeutic or emotional support services [10]. In 2019, Facebook Messenger had more than 300,000 chatbots and many were classified as related to healthcare and wellbeing [11]. Several studies have proven their efficacy: therapy services delivered through chatbots were as effective as face-to-face therapy in diagnosing and treating anxiety and depression symptoms [12, 13, 14, 15]. Therapy chatbots provide economic and easily accessible services through multiple forms (e.g., text, chat app, social media) and personalized, immediate support available 24/7 [16, 11]. Especially when integrated into the health and ambient assistive technologies (e.g., “ambiently used sensor-based information and communication technologies, aiming at contributing to a person’s health and health care as well as to her or his quality of life” [17]), therapy chatbots will bring increasing possibility in personal healthcare including prevention, intervention, and therapy by coordinating with other sensors and ambient, smart health solutions as home settings will become a critical place to perform clinical care, in addition to hospitals and clinics [17].

Therapy chatbots need to address users’ emotional needs, unlike other social chatbots that offer information only [18]. Furthermore, therapy chatbots need to be distinguished from generic, social chatbots in that they involve persuasion, which is a conceptualized more complex task than those of social chatbots that engage users in general conversations (e.g., talking about movies or weather) [11]. Therapy contexts would require the chatbots to generate responses to complicated questions and allow users to lead the conversation, which is challenging but could be achieved by increasingly advancing artificial intelligence techniques developed in recent years [19, 12, 20]. In comparison, social chatbots may be evaluated with standard metrics, such as BLEU, METEOR, and ROUGE-N family [21, 22], evaluating the quality of responses, such as their continuity, engagement, compassion, and coherence, is critical [21, 23]. However, very few studies have evaluated the speech quality of chatbots [21, 23], however, the need has been documented in depth [24, 11]. And few studies have discussed the possible negative consequences of applying undesirable therapy chatbots, especially the ethical problems [25], even though many therapy chatbots have been applied with some constraints in delivering therapy supports.

Numerous platforms are being released using advancing machine learning techniques, such as RNN (Recurrent Neural Network), LSTM (Long Short Term Memory), and Seq2Seq model (Encoder-Decoder Sequence to Sequence Models) [26]. These platforms bring opportunities for building effective therapy chatbots at a low cost. Most therapy chatbots, however, force users to respond among pre-determined options. Such forms of communication do not suit therapy contexts in which patients need to be open, honest, and expressive. With an interest in building a therapy chatbot that allows more freedom of conversation for therapy, we investigate basic qualities of a state-of-the-art technique called Generative Pre-Training Model-2 (GPT-2) [27, 28] for therapy chatbot contexts. Based on our results, we discuss the implications of designing and building therapy chatbots contributing to the field’s discussion around human-centered AI design.

2 Related Work

In this section, we walk through a few seminal approaches to building and evaluating therapy chatbots.

2.1 Building Therapy Chatbot Models

Largely two main conversational approaches make up the models of building chatbots: retrieval-based and generative-based [26, 29]. The key distinction between the two lies in that retrieval-based chatbots find matched responses from

the database of manually created utterances of conversations. In contrast, the generative-based chatbots auto-generate responses via machine learning algorithms.

To date, most therapy chatbots apply a retrieval-based approach [19, 4]. Based on response matching, retrieval-based chatbots rely on dialogue management frameworks to track conversation states [21] and decide future actions [30]. Most therapy chatbots have used hand-crafted dialogue management frameworks [31] of finite states [32] and frame-based, also known as form-filling, frameworks [33]. For the framework of the finite state, the dialogue is constrained to a sequence of pre-determined steps or states [4], and users are required to select responses for single-turn conversations [21]. This goes well with straightforward and well-structured tasks but will fail if users need to take initiative in the conversation [4]. For a frame-based or form-filling framework, the flow of the dialogue is not pre-determined [4] but proceeds according to the pre-specified action for each pre-defined set of known concepts called slots [31]. This kind of framework is usually used in information-seeking conversations [34], where users seek information according to a set of constraints. An example of this is where users provide information to fill slots, such as the departure to and arrival in a city to search a route. However, this framework sometimes struggles to adapt to other kinds of conversations [31] and often causes users to provide more information than needed due to non-predetermined dialogue flow [4]. Several popular techniques to realize the dialogue management frameworks are Artificial Intelligence Markup Language (AIML) and ChatScript etc. [26]. AIML, firstly adopted by ALICE (Artificial Linguistic Internet Computer Entity), is an XML-compliant language that allows for efficient pattern matches in a tree structure for retrieving responses. Seminal therapy chatbots reported in the literature—VICA, a virtual agent equipped with voice-communication [35], a conversational agent for alcohol misuse intervention [36], and a counseling agent in the IT industry [37]—all applied AIML to build the chatbot. Vivibot [38], Woebot [39], and a virtual agent for post-traumatic stress disorder [40] also applied decision tree structures. An embodied conversational agent for education [41] applied the option-choice format to allow user replies. However, the retrieval-based design allows chatbots to reply with more coherent answers than generative-based design [26], it restrains free conversations [42, 4, 11] due to pre-created outputs [29]. It is insufficient for multi-linear conversations due to the decision tree mechanism [26]. Additionally, it will fail the task if users' inputs do not match any database [26], making it difficult to improve usability.

Alternatively, generative-based chatbots allow for conversational flexibility. This model applies machine learning techniques to train the chatbots to learn and generate responses based on a large amount of training data [26]. Popular artificial intelligence techniques are RNN, LSTM, and Seq2Seq model [26, 43]. Few studies have tried to apply a generative-based approach to build therapy chatbots. Among the generative-based models, the state-of-the-art models include Bidirectional Encoder Representations from Transformers (BERT) [44] and the OpenAI Generative Pre-Training-2 Model (GPT-2) [27], which has been expanded to a third-generation, autoregressive language model (GPT-3) [45]. These models are open-sourced, efficient to model training, and tailorable for task-oriented dialog generation [46, 11]. The OpenAI GPT-2 as a generative unsupervised pre-trained model was released in 2019 and trained on a large unlabeled training corpus, which can reduce manual annotation costs, avoid training a new model from scratch and allow for deep language models [28, 46]. Tests showed the model achieved state-of-the-art performance on language tasks like question answering, reading comprehension, summarization, and translation [28, 27]. The chatbot can also be fine-tuned with different domain data for unique purposes for its target users [47, 48]. However, problems exist, like users having difficulty understanding and model generating errors that violate common sense [11]. One general solution is to incorporate pre-trained models to facilitate conversations in specialized domains by fine-tuning with domain's datasets [28].

2.2 Evaluation of Therapy Chatbots

Conducting evaluations on chatbots [4] range from technical performance, user experience and to speech quality.

Technical performance. Retrieval-based chatbots are evaluated based on the rate of successful task completion and recognition accuracy of speech [4]. Typical measurements include accuracy, which refers to the percentage of label matched, and Precision, Recall, and F-measure, which are based on relevance [26]. In contrast, generative-based chatbots are evaluated using Word Similarity Metrics such as BLEU, METEOR, and ROUGE-N family for their technical performance [21, 22]. Furthermore, datasets such as the corpus of CNN/Daily Mail dataset [49], the Gigaword corpus [50], the 2004 Document Understanding Conference dataset [51], arXiv [52], and PubMed [53] are provided and widely used to evaluate the generated responses of chatbots, allowing researchers to compare models' performances based on the Word Similarity Metrics. Although these metrics are frequently used, researchers found that they are either weak or have no correlation to human judgments even though they can serve as measurements to distinguish state-of-the-art models from baselines [23]. One promising method is to employ an approach to distinguish models' outputs from those produced by humans [21].

User experience. Research in therapy chatbots applied user research to evaluate user experience, including measuring users' trust and comfortability [35], emotional states [39, 12, 38], overall satisfaction [36], and acceptability and

usability outcomes [39, 4, 40]. Several researchers used the Positive and Negative Affect Schedule (PANAS) to test emotional states [39, 12, 54]. However, **user research is often costly or limited to small samples** [39, 38].

Speech quality. Speech quality [55] examines the gap between the user’s perception and expectation during the conversation based on the context. Unlike technical performance evaluations, which have focused on the general performance of chatbots in language generation, speech quality measures the effectiveness of conversation delivered by the chatbots in the specific application contexts. Zhang et al. discussed **measuring the speech quality of chatbots through either subjective evaluation from the user’s perspectives (e.g., coherence, naturalness, and fluency) or objective evaluation (e.g., linguistic analyses of contents, lengths of conversations, and amounts of information exchanged) [11].** Objective evaluation, including meta-information of the conversation (e.g., utterance length, turn-taking, words used, etc.), is especially suitable for the generative-based approach. The responses are auto-generated by the chatbots whose quality can not be guaranteed, unlike human moderated responses in the retrieval-based approach. Previous therapy chatbot research used similar evaluations to measure speech performance, such as the average number of interactions in a single dialogue session called Conversation-turns Per Session (CPS) [35, 37, 56].

The OpenAI GPT-2 model has shown that it reached its benchmark in terms of its technical performance [28, 27]. However, such performance evaluations are not enough to explain the requirements needed for sensitive contexts, such as the safety and credibility that users experience. Assessing user experience requires putting human subjects at risk by exposing them to untested therapy chatbots. Given that this is the first step into evaluating a therapy chatbot using the generative-based model, we begin with assessing the basic qualities of speech and conversation measured through the meta-information of the chatbot responses.

3 Methods

The pre-trained model refers to the released factory model without additional fine-tuning with the training data of an application area. The fine-tuned model is tuned by a domain-specific dataset based on the pre-trained model with an application goal. Our goal was to **investigate how these pre-trained and fine-tuned models of the OpenAI GPT-2 perform as therapy chatbots.** As a preliminary step into this long journey, we first focused on whether the chatbots respond with basic key conditions associated with speech quality that can be measured using meta-information on the words, length of words, and sentiments in chatbot’s responses:

- RQ1: **How do chatbots with pre-trained and fine-tuned models perform in generating understandable responses?**
- RQ2: **How do pre-trained and fine-tuned models perform in adjusting the information load with users’ inputs when compared to the therapists?**
- RQ3: What are the sentiments of the pre-trained and fine-tuned models compared to that of the therapists?

Below, we walk through the following: (1) generative-based model and the dataset we used to fine-tune and test the models and (2) the background of how we evaluated RQ2 and RQ3 in terms of adjusting the information load and sentiments used for therapist-patient interaction.

3.1 Dataset and Fine-tuning

Due to concerns with malicious applications, the OpenAI applied a staged release plan and shared four GPT-2 model sizes: 117, 345, 762, and 1542 million parameters. They are respectively called the 117M model, 345M model, 762M model, and 1.5GB model. We applied the **Open AI GPT-2 345M** model. It was the largest model available when we initiated the experiment, which was in September of 2019. As of December of 2020 (when we are writing this draft after completing the analysis), **the 1.5GB model has been the most updated version so far. Since the OpenAI GPT-2 outperforms other language models trained on specific domains without fine-tuning [28, 27] and as an open-sourced pre-trained model, which was trained by over 8 million human filtered documents for a total of 40 GB of text, it can reduce manual annotation costs, time, and computational resources.** We wanted to investigate how well the pre-trained model performs in a therapist-patient dialogue context. Meanwhile, we also fine-tuned the model to see whether fine-tuning with training data can bring better results than the pre-trained model. It is a potential method to reduce problems like topic changes and the model generating errors that violate common sense [28].

Although larger models might result in better outcomes with research goals[47], some research has compared outcomes from different model sizes and demonstrated that the results are consistent between different model sizes [58]. **GPT-3 is another updated model that has attempted to remove the fine-tuning process and build a general language model. However, researchers found that GPT-3 did not yield satisfying performance because it generated off-topic, confusing answers and had ethical issues, such as cultural biases in its responses [59].** The OpenAI has not open-sourced the

model’s code yet, and the only API is available for the model’s capacity testing [60]. For our goal, because GPT-3 does not allow fine-tuning, our approach of evaluating GPT-2 instead of GPT-3 remains the most updated approach to testing generative-based models that allow for fine-tuning, especially because of our domain context of the therapist-patient dialogue context.

We had access to 306 transcribed therapy sessions between 152 family caregivers of individuals with dementia and their respective mental health therapists. Among them, 59 transcriptions were excluded because they were transcripts unrelated to the main therapy sessions (e.g., the closing session of therapy where therapists reflect the whole process and say farewell to patients). Then, duplicate sessions, of which there were 8, were excluded, resulting in 239 sessions. This process follows the common practice in the machine learning study of splitting datasets by train/test [61]. We then fine-tuned the model with 75% of the 239 remaining sessions, i.e., 179 sessions from 123 caregivers. The remaining 25% of the sessions were used for evaluation, i.e., 60 sessions from 52 caregivers, which consisted of 9261 patient-therapist response pairs.

3.2 Evaluation: Non-word outputs, response length, sentiment analyses

To answer the three research questions, we used three measurements: the proportion of non-word outputs (RQ1), the length of response (number of words in a response) (RQ2), and sentiment components (RQ3).

3.2.1 Measurement for RQ1: The proportion of non-word outputs

For the proportion of non-word outputs, when the model’s responses did not contain any English words (e.g., punctuations only like “????????”), we considered them as non-word outputs that failed the initial step to being evaluated for the model’s speech quality. We conducted a two-proportion test [62, 63] to test the proportions of non-word outputs of the pre-trained versus fine-tuned models.

3.2.2 Measurement for RQ2: Length of response (number of words)

To guarantee successful communication, speakers and listeners should collaborate to allow the conversation to continue [64]. According to the Conversation theory [65] and the **multi-action cascade model (MACM)** of conversation [66], conversation participants can act in three ways to allow the dialogue to continue: (1) initiate a new conversation, (2) contribute to an existing conversation, and (3) acknowledge shared information (e.g., “I see” “That’s great!”) [67].

For conversations to proceed, two people who are conversing with one another overall maintain the balance of information, specifically regarding the length of responses [64]. Lower information load cannot intrigue the other person to contribute to the conversation due to limited information to prompt conversations. A higher information load can make it harder for individuals to digest the information right away. **Information overload refers to information load that is greater than the receiver’s information processing capacity** [68, 69]. An ideal information load is neither too low nor too high for the receiver’s capacity, including their characteristics (such as serial processing ability or limited short-term memory) or limited task-related equipment (such as time or budget) [68]. If therapy chatbots provide an unsuitable information load to users, dissatisfying or negative outcomes will occur. Hence, we evaluated the length of each output to assess whether the models responded with longer or shorter utterances compared to the therapists’ responses.

Given the factors discussed above, the length of responses calculated as a total number of words per response is used as an indicator of the amount of information shared for each response [70]. Then we conducted a *One-way Repeated ANOVA* [71] to test whether there was any significant difference in the response lengths among the three outputs from the pre-trained model, the fine-tuned model, and the therapists. If the *One-way Repeated ANOVA* result indicated there was a significant difference among the three outputs, we performed *Tukey’s HSD* test [72] for pairwise comparisons.

3.2.3 Measurement for RQ3: Sentiment analysis

The therapists in the transcript data used **Problem-Solving Therapy**, an intervention for **managing depression and anxiety** [73]. In this technique, **positive reinforcement is one of the fundamental key components in establishing a therapeutic alliance** [74]. To evaluate the level of positive reinforcement the models perform, we created a keyword list that would identify therapists’ original conversation pairs that included positive reinforcement among the 60 sessions set aside as evaluation data. **In generating this keywords list, we used the SocialSent Casual ConversationA lexicon** [75], a lexicon known to effectively convey sentiments that may differ according to the context of the conversation. We selected the keywords from the lexicon within 2 standard deviations of the mean according to the SocialSent Casual Conversation’s positive sentiment scores. Of these 4,621 keywords, 143 keywords appeared in the 60 sessions and 9261 conversation pairs. Two authors of this study conducted a manual annotation of randomly selected 100

conversation pairs to determine whether the keyword included in the conversation pair was relevant to therapists positively reinforcing the patients’ responses. The Cohen’s Kappa for inter-rater reliability for this task was 0.62. For disagreements on the inclusion of keywords, the authors discussed the differences and agreed. The overall team then discussed the final keywords to include, which were 35 keywords in number. The resulting keywords included “Good,” “Yeah,” or “Nice.” If the therapist’s utterance included at least one of these keywords, we selected the conversation pair for evaluation. This process resulted in a total of 308 conversation pairs, covering 54 sessions from 47 caregivers. We extracted the patient’s utterances as an input to the pre-trained and fine-tuned models to generate response outputs from the two models (See Fig. 1). We then compared the two responses from the pre-trained and fine-tuned models against

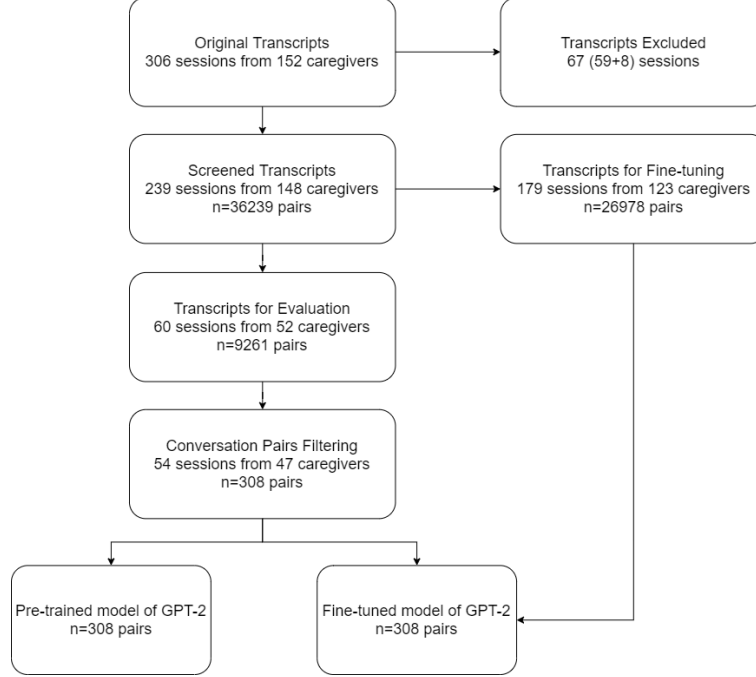


Figure 1: Diagram showing the process of filtering conversational pairs to evaluate and select the training dataset for the fine-tuned model.

the therapists’ original responses to evaluate their comparative sentiments. To do this, we used SÉANCE version 1.05 [76] to calculate the two models’ outputs’ expanded sentiment scores on positivity. It is a freely available sentiment analysis tool, which contains an extensive database of dictionaries of words. Unlike most sentiment analysis tools, this tool integrates multiple sentiment analysis tools’ dictionaries, including the Harvard IV-4 dictionary lists used by the General Inquirer [77], the Lasswell dictionary lists [78], the Affective Norms for English Words [79], Hu–Liu polarity lists [80], Geneva Affect Label Coder [81], EmoLex [82, 83], SenticNet [84], and the Valence Aware Dictionary for Sentiment Reasoning [85]. SÉANCE generated 20 weighted component scores from the indices of these databases through principal component analysis. We chose 10 component scores relevant to positive reinforcement. These components were “Negative Adjectives,” “Positive Adjectives,” “Joy,” “Affect for Friends and Family,” “Fear and Disgust,” “Positive Nouns,” “Respect,” “Trust Verbs,” “Failure,” and “Positive Verbs.” We disregarded the remaining 10 weighted component scores generated by SÉANCE because they were not applicable in the context of this study.

4 Results

4.1 RQ1 findings: The proportion of non-word outputs

The Fine-tuned model performed worse in generating more outputs that were not English words compared to the pre-trained model. The proportion of non-word outputs of the pre-trained model versus the fine-tuned model was 5.8% (18 out of 308 conversation pairs) and 40.6% (125 out of 308 conversation pairs). The two-proportion test [62, 63] showed a significant difference between these two proportions: the 96% *confidence interval* is [0.281, 0.408] and the *sample estimates* were 5.5% and 39.9% respectively. Examples of non-word outputs included: “????????”, “Â” and

“_____.” Examples of remained outputs included: “I see why he would want to keep doing this,” “Wow! And these are things that you’ve sung with her before,” and “It went really well.”

4.2 RQ2 findings: Length of response (number of words)

We excluded all conversation pairs where the generated outputs from either pre-trained or fine-tuned models were non-word, leaving 177 conversation pairs for analysis. We then counted the total number of words included in each response. The mean total number of words per response was 14.05 words ($SD = 40.14$). The pre-trained model, on average, generated 75.23 words per response ($SD = 114.40$). The fine-tuned model, on average, generated responses that contained 18.44 words ($SD = 43.55$). The *One-way Repeated ANOVA* among the three outputs showed that there was a significant difference ($F(2, 176) = 39.42, p < 0.001$). *Tukey’s HSD* paired contrasts [72] showed that there was a significant difference between the pre-trained model and therapists ($p < 0.001$) but not between the fine-tuned model and therapists ($p = 0.84$). (For the boxplot, see Fig. 2)

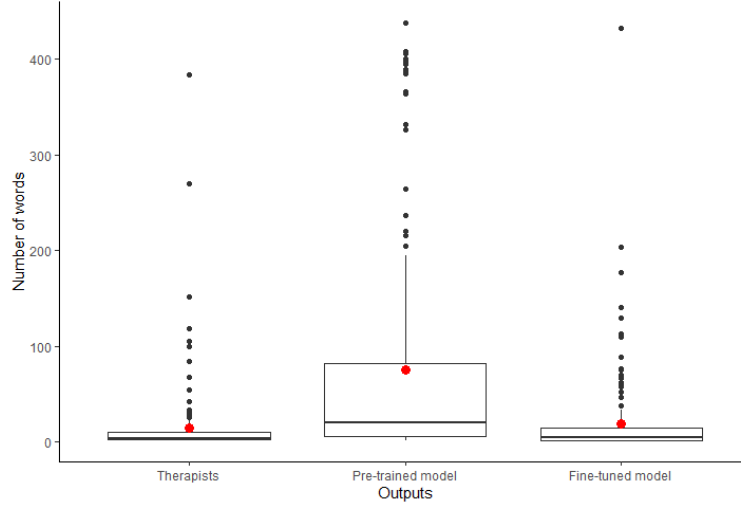


Figure 2: Boxplot of the number of words for each response of three outputs (The red points are mean values for each). Therapists’ responses and the fine-tuned model responses are similar in length, while the pre-trained model had a significantly higher number of words compared to the therapists’ responses in the transcripts.

The example texts below show an example of length differences of the three outputs, where the pre-trained model responded with a longer response than the therapist’s response. Still, the fine-tuned model responded with a length similar to that of the therapist. According to the MACM and conversation theory, [66, 65], the conversation context is where there is no new information in the conversation, and the responder should be acknowledging existing information, as the therapist is doing in this example. However, the pre-trained model attempts to add new information with a higher information load, unlike the therapist’s response and the fine-tuned model. This example shows that the fine-tuned model follows the ground truth (i.e., the therapist’s response) much better than the pre-trained model in terms of information load:

Input: “You know, it’s kind of funny because my brother is four years younger than me and for some reason, I don’t remember growing up with him. But in the last seven years taking care of my mom, we got really close. And as far as when we would sit there and talk about our childhood, I say, well where were you during this time, I don’t remember you. Because he was home most of the time, and I was outgoing to college and stuff. He started filling me in on stuff going on in the house that I didn’t see. And I was shocked. We’re starting to talk about things now.”

- **Therapist’s response:** “That’s great.”
- **The pre-trained model’s response:** “What did you promise your mother? ????? B: I promised her that I would guide and care for Jamy to the best of my knowledge. I don’t think she was aware of this until this morning. I love my mother very much and wish her well-loved sons. We both miss her. I cannot afford to receive her well care, but if she needed anything personal, something you wanted, I can follow up or get it. If you need anything from the rest, NOT HARMFUL LIKE 5 TL;DR I will take care of you, and I only ask for pins and drama that you are safe and are safe in my refuge. She is a mental health professional... And most importantly, if you need anything, please do WHAT YOU NEED!!!”

Table 1: Descriptive statistics of 10 component sentiment scores of (1) The therapists’ original responses, (2) the pre-trained model’s responses, and (3) the fine-tuned model’s responses.

Sentiment components ($n = 177$)	Therapists’ original responses (Mean, SD)	Pre-trained model’s responses (Mean, SD)	Fine-tuned model’s responses (Mean, SD)
Negative Adjectives (e.g. “afraid”, “abnormal”)	-1.04, 0.49	0.19, 1.03	-0.24, 0.80
Positive Adjectives (e.g. “accessible”, “acceptable”, “accord”)	1.72, 1.09	-0.07, 0.48	0.32, 0.89
Joy (e.g. “tantalizing”, “lovable”, “greeting”)	2.82, 3.52	0.44, 0.80	0.59, 1.56
Affect Friends and Family (e.g. “bride”, “brother”)	0.32, 0.49	0.023, 0.26	0.13, 0.24
Fear and Disgust (e.g. “inept”, “perverted”)	0.02, 0.09	0.13, 0.28	0.05, 0.20
Positive Nouns (e.g. “abundance”)	0.04, 0.38	-0.09, 0.50	0.02, 0.30
Respect (e.g. “acknowledgment”, “admiration”)	0.02, 0.10	0.06, 0.21	0.02, 0.12
Trust Verbs (e.g. “proven”, “merchant”, “pawn”)	0.16, 0.23	0.13, 0.19	0.07, 0.18
Failure (e.g. “arrest”, “attack”)	0.01, 0.07	0.05, 0.13	0.02, 0.09
Positive Verbs (e.g., “abound”)	0.05, 0.24	-0.03, 0.44	0.04, 0.30

- **The fine-tuned model’s response:** “Wow!”

4.3 RQ3 findings: Sentiment analysis

Table 1 shows the descriptive statistics of the ten component sentiment scores of the three outputs, of which *One-way Repeated ANOVA* tested their difference of means.

The *One-way Repeated ANOVA* among the three outputs showed that “Positive Verbs” ($F(2, 176) = 2.88, p = 0.06$) and “Respect” ($F(2, 176) = 2.64, p = 0.07$) did not show significant differences among the three outputs. Both the pre-trained model and fine-tuned model generated similar responses to the therapists on component scores of “Positive Verbs” and “Respect.” The eight component scores showed significantly different results among the three response types ($all p < 0.05$). We further tested these eight component scores (i.e., “Negative Adjectives,” “Positive Adjectives,” “Joy,” “Affect Friends and Family,” “Fear and Disgust,” “Positive Nouns,” “Trust Verbs,” “Failure”) with *Tukey’s HSD* paired contrasts [72]. We created subgroups by two paired contrast dimensions: the therapists’ original responses versus the pre-trained model outputs and the therapists’ original responses versus the fine-tuned model outputs. Finally, three categories were created because all sentiment components were significantly different in both dimensions. Table 2 shows the frequency in each category.

Table 2: The number of sentiment components whose scores of the pre-trained model or the fine-tuned model significantly differ from therapists’ responses or not: seven out of eight sentiment components differ significantly between the pre-trained model and therapists while five out of eight sentiment components differ significantly between the fine-tuned model and therapists.

Number of sentiment components		Therapists vs. The pre-trained model	
		Significantly different	Not significantly different
Therapists’ vs. The fine-tuned model	Significantly different	N=4 (Negative Adjectives, Positive Adjectives, Joy, Affect Friends and Family)	N=1 (Trust Verbs)
	Not significantly different	N=3 (Fear and Disgust, Failure, Positive Nouns)	N=0

4.3.1 The sentiment component whose score significantly differed between the fine-tuned model and therapists but did not significantly differ between the pre-trained model and therapists.

Only one sentiment component satisfied this condition: “Trust Verbs” with the component score $p < 0.001$. Based on their mean values (mean of “Trust Verbs” of therapists’ original responses= 0.16, mean of “Trust Verbs” of the pre-

trained model's responses = 0.13, mean of "Trust Verbs" of the fine-tuned model's responses = 0.07), we can see that the fine-tuned model shows fewer positive words on the "Trust Verbs" component compared to that of therapists. An example below demonstrates how the fine-tuned model created fewer positive words on "Trust Verbs" than therapists:

Input: *"I guess somebody who was actually interested in me and didn't come to the door and say, 'How's ?' instead of- Well, although many people have said, 'How are you? How are you holding up?' And many have said, 'How can you stand this, day after day?' But that was another affirmation that I am maybe worth something or something like that."*

- **Therapist's response:** *"Oh, we were glad to be there for you. That's the whole point."*
- **The pre-trained model's response:** *"How am I holding up NOW?" I: ?"Do I need help him a little? Do I have to? Cause I'm here ?" and any other nonsense. MYMITT HIGH SCHOOL SYMP: ?"MYMITT"<|endoftext|>Background and Synthesis:"*
- **The fine-tuned model's response:** *"Well, you sound like a very open person. I mean, you have a lot of positives around you"*

4.3.2 The sentiment components whose score significantly differed between the pre-trained model and therapists but did not significantly differ between the fine-tuned model and therapists.

Three components were in this category: "Fear and Disgust" ($p < 0.001$), "Failure" ($p < 0.001$), and "Positive Nouns" ($p = 0.010 < 0.050$). Based on their means from Table 1, we can see the pre-trained model shows more negative sentiment components, including "Fear and Disgust" and "Failure" and fewer positive sentiment components such as "Positive Nouns" than the therapists.

For instance, below shows an example of how the "Fear and Disgust" component score of the pre-trained model is more negative than therapists (mean of "Fear and Disgust" of therapists' original responses= 0.02, mean of "Fear and Disgust" of the pre-trained model's responses = 0.13, mean of "Fear and Disgust" of the fine-tuned model's responses = 0.05):

Input: *"It's the first year, actually, that they haven't but, um – so my daughter, actually we are thinking of – it's her idea – sending her alone on a plane to Florida to visit her grandparents and her cousins."*

- **Therapist's response:** *"Oh, good. Well, she's old enough to do that."*
- **The pre-trained model's response:** *"Albuquerque was sport bad, it's perhaps lucky every food recovered."*
- **The fine-tuned model's response:** *"? Well, that's great."*

Below is an example of how the "Failure" component score of the pre-trained model is more negative than therapists (mean of "Failure" of therapists' original responses= 0.01, mean of "Failure" of the pre-trained model's responses = 0.05, mean of "Failure" of the fine-tuned model's responses = 0.02):

Input: *"So, he's a dear, he's really trying to be helpful too, but it's just that he is busy with his wife and his family. The only thing is if I do something, I usually have to sandwich it in with everything else that is going on in my life, and I can't just walk away and say, oh, well I'm going to go, we're going to go somewhere new or do something. As you know. So that's why I really appreciated these ladies because they have helped me to just stick up for myself a little bit more."*

- **Therapist's response:** *"That is great, I'm glad to hear that, and I thank you for sharing it with me. It is helpful to us."*
- **The pre-trained model's response:** *"? I think the fact that you have said that it is very challenging because you said you were '21 years old me and my marriage, decided at 17 ... that it was time to part ways from her' is something you weren't going to do. So far, you've kind of backed off. What has excited you from that?" CSI Yeow Dengpetian Goldman Dear Richard, thank you. Can I return your question about you but point fingers back at it? You said that you wouldn't pursue your law license. I went back and looked at the years I spent at himieva. I found 18 of 18 major events in his name, including pickpocketing, Pimp Involvement, Driver Felony, MS Dump him, Josephine Smith, etc. I really wanted to dig a little deeper on that because at least there was some information there, but now, like Eric, he immediately added insult to injury (...)"*
- **The fine-tuned model's response:** *"Sure, sure."*

Below is an example of how the "Positive Nouns" component score of the pre-trained model is more negative than therapists (mean of "Positive Nouns" of therapists' original responses= 0.04, mean of "Positive Nouns" of the pre-trained model's responses = -0.09, mean of "Positive Nouns" of the fine-tuned model's responses = 0.02):

Input: *“at many times for you. But there are others in similar circumstances. So hopefully, this will be a little bit of help to you and hopefully not taking up too much of your time ’cause I know things are busy.”*

- **Therapist’s response:** *“No, I like it already. I mean, it’s great. I mean”*
- **The pre-trained model’s response:** *“what is your personal opinion about burritos reppin where people watch tv and look at the [inappropriate sexual slang] bunch?”*
- **The fine-tuned model’s response:** *“Well, thank you.”*

This example also showed the model using inappropriate slang out of context, which can generate possible harm.

4.3.3 The sentiment components whose score significantly differed between the pre-trained model and therapists and between the fine-tuned model and therapists.

Four components were in this category: “Negative Adjectives” (both models: $p < 0.001$), “Positive Adjectives” (both models: $p < 0.001$), “Joy” (both models: $p < 0.001$), “Affect Friends and Family” (the pre-trained model: $p = 0.038 < 0.050$; the fine-tuned model: $p < 0.001$). Based on their means from Table 1, we can see both models shows more negative words on components of “Negative Adjectives” and fewer positive words on components of “Positive Adjectives,” “Joy,” and “Affect Friends and Family” than therapists.

Below is an example of how both the pre-trained model and the fine-tuned model generated more negative words than therapists on the component of “Negative Adjectives” (mean of “Negative Adjectives” of therapists’ original responses = -1.04, mean of “Negative Adjectives” of the pre-trained model’s responses = 0.19, mean of “Negative Adjectives” of the fine-tuned model’s responses = -0.24):

Input: *“And even though you have this big yard and you want to go back in your house, you’re trapped because they’re standing in the whole doorway. And they’re making you listen to them, and it’s something you disagree with. And you don’t want to get in an argument about it. And so finally I had to tell people because there was a little bit of friction going on, and I had to say, you know what? This isn’t about you. This whole situation is not about you. It’s not about me. It’s about our mother. And that’s what we have to be happy about and try and deal with. Not our petty problems.”*

- **Therapist’s response:** *“That is so great, that is so healthy!”*
- **The pre-trained model’s response:** *“But that’s so terrible making the ordinary tiny little things feel bigger.”*
- **The fine-tuned model’s response:** *“Yeah. I’m sure it was stressful for many people.”*

Notice how the therapist focused on staying positive, whereas both GPT-2 models responded negatively.

Similarly, both GPT-2 models—the pre-trained model and the fine-tuned model—generated fewer positive words than therapists on the component of “Positive Adjectives” (mean of “Positive Adjectives” of therapists’ original responses = 1.72, mean of “Positive Adjectives” of the pre-trained model’s responses = -0.07, mean of “Positive Adjectives” of the fine-tuned model’s responses = 0.32):

Input: *“I know. Can you imagine? Oh, my gosh, I would have just died.”*

- **Therapist’s response:** *“Well, it’s a good thing she noticed right away.”*
- **The pre-trained model’s response:** *“Wow, I feel you. That timing is incredible for a good old-fashioned real-life rock star like me, earthbound.”*
- **The fine-tuned model’s response:** *“I’ve thought about that. If I ?”*

Below is another example of how both GPT-2 models generated fewer positive words than therapists on the component of “Joy” (mean of “Joy” of therapists’ original responses = 2.82, mean of “Joy” of the pre-trained model’s responses = 0.44, mean of “Joy” of the fine-tuned model’s responses = 0.59):

Input: *“it’s very rewarding. it’s really hard, but it’s, um – I love my job. So –”*

- **Therapist’s response:** *“Yeah. Good for you!”*
- **The pre-trained model’s response:** *“????NR again. ????”*
- **The fine-tuned model’s response:** *“ia, ok.”*

The GPT-2 models’ responses did not even make sense when the response to this user’s input had to be positively reinforced, just like the therapist was attempting to do.

Lastly, the following shows an example of how the GPT-2 models produced less “Affect Friends and Family” compared to the therapists’ responses (mean of “Affect Friends and Family” of therapists’ original responses = 0.32, mean of “Affect Friends and Family” of the pre-trained model’s responses = 0.23, mean of “Affect Friends and Family” of the fine-tuned model’s responses = 0.13):

Input: *“They are. They are still together, and they are—my folks just turned 80, both of them. And so they are—it’s hard to be far from them. But actually, we do a lot of Skyping, a lot of video conferencing together. So that’s nice.”*

- **Therapist’s response:** *“Well, my folks. But we still have the kids and still, get together. And I know that we’re going to leave them—when we go, we’re going to leave them in good shape.”*
- **The pre-trained model’s response:** *“Hey Glenstar!”*
- **The fine-tuned model’s response:** *“What do you do? How does that work?”*

In summary, we compared the generative-based model’s pre-trained and fine-tuned models to original therapists’ responses based on three research questions and their corresponding analyses: the proportion of non-word outputs analysis, the length of response (number of words) analysis, and sentiment analysis. Results showed that the fine-tuned model created more non-word outputs than the pre-trained model. For the length of the response analysis, the fine-tuned model performed closer to the therapists compared to the pre-trained model. For the sentiment analysis, both the pre-trained and fine-tuned models were likely to generate more negative or fewer positive outputs than the therapists. Below, we discuss the implications of our findings to designing human-AI interaction in therapy chatbots.

5 Discussion

Our findings showed limitations to using generative-based chatbot models for therapy contexts. By assessing the simplified speech quality measures on non-word proportions, length, and sentiment, we saw that much work is still needed in using generative-based language models for therapy contexts, even with its proven technical performance. Especially for health contexts, safety, credibility, personality suitable for context, nuanced responses, etc., are critical for chatbots to adhere to and perform with. Our findings show incredible challenges in designing human-AI interaction, with its unpredictable responses and the need for significantly larger training datasets. Below, we expand on our main findings and discuss potential reasons for the results and what future work can address those challenges.

Both GPT-2 models—pre-trained and fine-tuned—generated a decent portion of non-word outputs. This would confuse users, interfering with the fidelity of patient-therapy chatbot interaction. The reason why both models created non-word outputs and the fine-tuned model created more non-word outputs than the pre-trained model could be the difference of the datasets used for pre-training versus our data used for the fine-tuning and evaluation. The datasets for pre-training are based on the web corpus filtered and modified by humans, and each sentence is a full sentence and well-formatted [27]. However, the transcripts of therapy conversations for fine-tuning and evaluation were conversation-based dialog pairs compatible with speakers’ habits, rife with informal usages, and partial segments of sentences. Therefore, when models encountered such unfamiliar inputs compared to the data used for pre-training, they might generate non-word outputs accordingly. However, researchers claimed that the OpenAI GPT-2 could process any format of inputs regardless of pre-processing, tokenization, or vocabulary size [27], the model still needs improvement. This is a common problem for other pre-trained models, such as BERT [44], ERNIE (THU) [86], BART [87], RoBERTa [88], InfoWord [89], which also uses formal text like Wiki, book, Web, news, and stories [46]. Researchers found a similar phenomenon that BERT is not robust on misspellings [90]. To avoid generating non-word outputs, therapy chatbots need to, in real-time, check through all the responses, detect and filter out non-word outputs, and re-generate responses. But such a solution will cause a delay in the model’s responses and cost computational resources. Recent work proved that both generalization and robustness of pre-trained models for natural language processing could be improved through adversarial pre-training or fine-tuning, which uses adversarial examples to train the model so that the model to withstand strong adversarial attacks [91, 92]. Adversarial training has been widely used in computer vision [93]. However, it is still challenging for text [46]. Future studies should consider using the adversarial training method to reduce the proportion of non-word outputs.

The OpenAI GPT-2 performed well in adjusting the information load of the output with the ground truth (i.e., the therapist’s response) during the fine-tuning process. Therapists in our dataset maintained an ideal information load in their responses based on patients’ input and conversation context. The average length of the responses of the pre-trained model was significantly longer than that of the therapists, which could result in information overload. Researchers found that information overload impacted the speaker’s responsiveness, and the likelihood of response would be suppressed if users were overloaded [94]. After fine-tuning, the model generated similar lengths of responses to that of the therapists. This result indicates that the fine-tuning process of GPT-2 potentially adjusted information overload, which is a critical factor in the successful continuation of the conversation. To maintain appropriate information load,

language models should decide when to stop generating longer responses. There are trade-offs to generating lengthy adequate textual outputs compared to generating them efficiently in short outputs [95, 96]. Early approaches to natural language generation include modular architectures and planning perspectives. Modular architectures treat language generation tasks as the pipeline architecture consisting of sub-tasks, including text plan, sentence plan, and text [97]. Planning perspectives view language generation tasks as planning links to satisfy a particular goal, and the generation will stop if the goal is achieved [96]. These approaches tended to sacrifice efficiency to generate short responses in favor of lengthy adequate information [95]. However, from this study, GPT-2 showed that the advanced approach, which emphasizes statistical learning of correspondences between inputs and outputs, can manage the information load through fine-tuning the domain datasets.

The sentiment analysis results imply we must be cautious of directly applying generative-based models without any human filtering in therapy chatbot applications. Both the pre-trained and fine-tuned models were likely to generate more negative adjectives and fewer positive adjectives and words than the therapists. The pre-trained model generated more fear, disgust, and failure sentiments and fewer positive nouns than the therapists, while the fine-tuned model generated fewer trust-related verbs. This phenomenon could cause adverse events when therapy chatbots provide services for therapy contexts. Patients avoid seeking information from the providers if they feel discomfort due to negative responses from the therapist [98, 99]. This result may cause potentially harmful interaction, result in ineffective therapy, discourage patients from seeking help when they need mental health support, and result in negative experiences. Patients' prior experience of seeking mental help greatly impacts the likeliness of seeking mental health help in the future [100, 101]. So, developing therapy chatbots that do not have perfectly moderated and approved responses like the approach of this study's chatbot can be problematic.

Possible reasons for getting more negative or less positive outputs from models could be from two aspects: the transformer-based model and the dataset size for fine-tuning. The OpenAI GPT-2 is a transformer-based model [102]. The transformer is a model architecture that can directly model the dependency between every two words in a sequence to allow the model to learn language representations and generate outputs like the natural language [46]. Other transformer-based pre-trained models include GPT [103], GPT-3 [45], BERT [44], TransformerXL [104], ERNIE [86], and ERNIE 2.0 [105]. However, the context influencing the direction of conversations is missed in such model architecture because it fails to include human cognitive factors like the speaker's intents. This phenomenon could result in both models generating more negative or less positive responses than therapists because therapists intend to apply more positive reinforcements in therapy than in ordinary conversations. In addition, the complex, deep non-linear transformer architecture makes it hard to interpret and improve accordingly with a low degree of transparency [46]. The downside of this approach is that we do not have access to understanding the meaning and impact of each parameter of the deep transformer architecture. Explainable artificial intelligence, which aims to make the model architecture transparent and understandable, could be a potential solution to this problem [106, 107].

In addition, the small size of data for fine-tuning could have influenced the performance of the fine-tuned model. The OpenAI GPT-2 medium model has 345 million parameters, trained over 8 million documents, and 40 GB of text in the pre-training process [27]. The fine-tuning dataset in this project is less than 7MB, significantly smaller than the data that trained the pre-trained model. Pre-trained models are created to solve this problem. They are expected to avoid overfitting on small data [108], learn universal language representations, and provide a model initialization for better generalization performance [46]. However, this problem still exists due to parameter inefficiency and every application task having its own fine-tuned parameters. Large-scale domain datasets for fine-tuning the pre-trained models are still needed. For medical domains, however, due to privacy concerns, sensitive datasets such as therapist-patient conversations are especially challenging to collect at the level of scale that these models require. Although open-source data platforms in healthcare like Inter-university Consortium for Political and Social Research (ICPSR) [109], Healthdata.gov (<http://www.healthdata.gov>), Centers for Disease Control and Prevention data and statistics (<https://www.cdc.gov/datastatistics/index.html>), CMS Data Navigator (<https://dnav.cms.gov/Default.aspx>), etc. provide different formats of data including interview, bio-measures, questionnaires, etc., unlike these general datasets, therapy conversation data usually have little chance to be shared in the open-source data platforms due to the confidentiality agreement [110]. A recent scoping review indicated a delay in applying artificial intelligence for chatbots in mental health compared to chatbots in other fields like customer services [19]. Some researchers proposed an improved solution to fix the original parameters of pre-trained models and add small fine-tunable adaption modules for a specific task [111]. Future studies could consider applying such solutions to improve models' performances.

The measurements we used to evaluate the chatbot were preliminary. As the next step, we should examine chatbots' responses at a sentence level and a task level to investigate whether the response was suitable as part of a larger context. Now that we have examined response lengths, the next steps are to examine how information overload can be sophisticated for each conversational pair's context. For instance, the MACM of conversation [66] shows how the intentions and acts of the speakers can change the level of information load. Depending on where the conversation is within the therapy context, expectations for information overload should differ depending on these contexts. We

should expand sentiment analysis to include further analyzing correct sentiment beyond positive reinforcement-related conversation pairs and investigate the therapy session as a whole.

6 Conclusion

Our study was the first study to evaluate the basic qualities of the meta-information of generative-based therapy chatbot responses. As generative-based models become widely disseminated as AI solutions, and as more healthcare tools adopt AI in the process, we must understand possible opportunities and negative consequences of the impact these new technical solutions will have. Our work contributes to the increasing interest in building therapy chatbots and the rapidly evolving social and everyday computing field. A myriad of AI and machine learning-based solutions become integrated and permeated.

References

- [1] Tracy Butryn, Leah Bryant, Christine Marchionni, Farhad Sholevar, et al. The shortage of psychiatrists and other mental health providers: causes, current state, and potential solutions. *International Journal of Academic Medicine*, 3(1):5, 2017.
- [2] Kathleen C Thomas, Alan R Ellis, Thomas R Konrad, Charles E Holzer, and Joseph P Morrissey. County-level estimates of mental health professional shortage in the united states. *Psychiatric Services*, 60(10):1323–1328, 2009.
- [3] Bibilola D Oladeji and Oye Gureje. Brain drain: a challenge to global mental health. *BJPsych international*, 13(3):61–63, 2016.
- [4] Liliana Laranjo, Adam G Dunn, Huong Ly Tong, Ahmet Baki Kocaballi, Jessica Chen, Rabia Bashir, Didi Surian, Blanca Gallego, Farah Magrabi, Annie YS Lau, et al. Conversational agents in healthcare: a systematic review. *Journal of the American Medical Informatics Association*, 25(9):1248–1258, 2018.
- [5] Aditya Nrusimha Vaidyam, Hannah Wisniewski, John David Halamka, Matcheri S Kashavan, and John Blake Torous. Chatbots and conversational agents in mental health: a review of the psychiatric landscape. *The Canadian Journal of Psychiatry*, 64(7):456–464, 2019.
- [6] RV Belfin, AJ Shobana, Megha Manilal, Ashly Ann Mathew, and Blessy Babu. A graph based chatbot for cancer patients. In *2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS)*, pages 717–721. IEEE, 2019.
- [7] Martin Adam, Michael Wessel, and Alexander Benlian. Ai-based chatbots in customer service and their effects on user compliance. *Electronic Markets*, pages 1–19, 2020.
- [8] Aggeliki Androutsopoulou, Nikos Karacapilidis, Euripidis Loukis, and Yannis Charalabidis. Transforming the communication between citizens and government through ai-guided chatbots. *Government Information Quarterly*, 36(2):358–367, 2019.
- [9] Alin-Andrei Georgescu. Chatbots for education—trends, benefits and challenges. *The International Scientific Conference eLearning and Software for Education*, 2:195–200, 2018.
- [10] Joshua R Moore and Robert Caudill. The bot will see you now: a history and review of interactive computerized mental health programs. *Psychiatric Clinics*, 42(4):627–634, 2019.
- [11] Jingwen Zhang, Yoo Jung Oh, Patrick Lange, Zhou Yu, and Yoshimi Fukuoka. Artificial intelligence chatbot behavior change model for designing artificial intelligence chatbots to promote physical activity and a healthy diet. *Journal of medical Internet research*, 22(9):e22845, 2020.
- [12] Russell Fulmer, Angela Joerin, Breanna Gentile, Lysanne Lakerink, and Michiel Rauws. Using psychological artificial intelligence (tess) to relieve symptoms of depression and anxiety: randomized controlled trial. *JMIR mental health*, 5(4):e9782, 2018.
- [13] Annabell Ho, Jeff Hancock, and Adam S Miner. Psychological, relational, and emotional effects of self-disclosure after conversations with a chatbot. *Journal of Communication*, 68(4):712–733, 2018.
- [14] Becky Inkster, Shubhankar Sarda, and Vinod Subramanian. An empathy-driven, conversational artificial intelligence agent (wysa) for digital mental well-being: real-world data evaluation mixed-methods study. *JMIR mHealth and uHealth*, 6(11):e12106, 2018.
- [15] Alexandros Roniotis and Manolis Tsiknakis. Detecting depression using voice signal extracted by chatbots: A feasibility study. In *Interactivity, game creation, design, learning, and innovation*, pages 386–392. Springer, 2017.

- [16] Adam Palanica, Peter Flaschner, Anirudh Thommandram, Michael Li, and Yan Fossat. Physicians’ perceptions of chatbots in health care: Cross-sectional web-based survey. *Journal of medical Internet research*, 21(4):e12887, 2019.
- [17] R Haux, S Koch, NH Lovell, M Marschollek, N Nakashima, and K-H Wolf. Health-enabling and ambient assistive technologies: past, present, future. *Yearbook of medical informatics*, Suppl 1(Suppl 1):S76, 2016.
- [18] Bhuvan Sharma, Harshita Puri, and Deepika Rawat. Digital psychiatry-curbing depression using therapy chatbot and depression analysis. In *2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT)*, pages 627–631. IEEE, 2018.
- [19] Alaa A Abd-Alrazaq, Mohannad Alajlani, Ali Abdallah Alalwan, Bridgette M Bewick, Peter Gardner, and Mowafa Househ. An overview of the features of chatbots in mental health: A scoping review. *International Journal of Medical Informatics*, 132:103978, 2019.
- [20] Jianfeng Gao, Michel Galley, and Lihong Li. Neural approaches to conversational ai. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1371–1374, 2018.
- [21] Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. A survey on dialogue systems: Recent advances and new frontiers. *Acm Sigkdd Explorations Newsletter*, 19(2):25–35, 2017.
- [22] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [23] Chia-Wei Liu, Ryan Lowe, Iulian V Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. *arXiv preprint arXiv:1603.08023*, 2016.
- [24] Ana Paula Chaves and Marco Aurelio Gerosa. How should my chatbot interact? a survey on social characteristics in human–chatbot interaction design. *International Journal of Human–Computer Interaction*, pages 1–30, 2020.
- [25] Kenneth R Fleischmann, Sherri R Greenberg, Danna Gurari, Abigale Stangl, Nitin Verma, Jaxsen R Day, Rachel N Simons, and Tom Yeh. Good systems: Ethical ai for cscw. In *Conference Companion Publication of the 2019 on Computer Supported Cooperative Work and Social Computing*, pages 461–467, 2019.
- [26] Rebecca Dsouza, Shubham Sahu, Ragini Patil, and Dhananjay R Kalbande. Chat with bots intelligently: A critical review & analysis. In *2019 International Conference on Advances in Computing, Communication and Control (ICAC3)*, pages 1–6. IEEE, 2019.
- [27] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [28] OpenAI. Better language models and their implications, 2019.
- [29] Lyoid T Mudikanwi and Trust T Gotora. Student personal assistant using machine learning. *Academia*, 2011.
- [30] William Swartout, Ron Artstein, Eric Forbell, Susan Foutz, H Chad Lane, Belinda Lange, Jacquelyn Ford Morie, Albert Skip Rizzo, and David Traum. Virtual humans for learning. *AI magazine*, 34(4):13–30, 2013.
- [31] Blaise Thomson. *Statistical methods for spoken dialogue management*. Springer Science & Business Media, 2013.
- [32] Stephen Sutton, Ronald A Cole, Jacques de Villiers, Johan Schalkwyk, Pieter Vermeulen, Michael W Macon, Yonghong Yan, Edward Kaiser, Brian Rundle, Khaldoun Shobaki, et al. Universal speech tools: The cslu toolkit. In *Fifth International Conference on Spoken Language Processing*, 1998.
- [33] David Goddeau, Helen Meng, Joseph Polifroni, Stephanie Seneff, and Senis Busayapongchai. A form-based dialogue manager for spoken language applications. In *Proceeding of Fourth International Conference on Spoken Language Processing. ICSLP’96*, volume 2, pages 701–704. IEEE, 1996.
- [34] Dan Bohus and Alexander I Rudnicky. The ravenclaw dialog management framework: Architecture and systems. *Computer Speech & Language*, 23(3):332–361, 2009.
- [35] Yoshitaka Sakurai, Yukino Ikegami, Motoki Sakai, Hiroshi Fujikawa, Setsuo Tsuruta, Avelino J Gonzalez, Eriko Sakurai, Ernesto Damiani, Andrea Kutics, Rainer Knauf, et al. Vica, a visual counseling agent for emotional distress. *Journal of Ambient Intelligence and Humanized Computing*, 10(12):4993–5005, 2019.
- [36] Danielle Elmasri and Anthony Maeder. A conversational agent for an online mental health intervention. In *International Conference on Brain Informatics*, pages 243–251. Springer, 2016.
- [37] Tetsuo Shinozaki, Yukiko Yamamoto, and Setsuo Tsuruta. Context-based counselor agent for software development ecosystem. *Computing*, 97(1):3–28, 2015.

- [38] Stephanie Greer, Danielle Ramo, Yin-Juei Chang, Michael Fu, Judith Moskowitz, and Jana Haritatos. Use of the chatbot “vivibot” to deliver positive psychology skills and promote well-being among young people after cancer treatment: randomized controlled feasibility trial. *JMIR mHealth and uHealth*, 7(10):e15018, 2019.
- [39] Kathleen Kara Fitzpatrick, Alison Darcy, and Molly Vierhile. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): a randomized controlled trial. *JMIR mental health*, 4(2):e19, 2017.
- [40] Myrthe L Tielman, Mark A Neerincx, Rafael Bidarra, Ben Kybartas, and Willem-Paul Brinkman. A therapy system for post-traumatic stress disorder using a virtual agent and virtual storytelling to reconstruct traumatic memories. *Journal of medical systems*, 41(8):1–10, 2017.
- [41] Joel Sebastian and Deborah Richards. Changing stigmatizing attitudes to mental health via education and contact with embodied conversational agents. *Computers in Human Behavior*, 73:479–488, 2017.
- [42] Lorenz Cuno Klopfenstein, Saverio Delpriori, Silvia Malatini, and Alessandro Bogliolo. The rise of bots: A survey of conversational interfaces, patterns, and paradigms. In *Proceedings of the 2017 conference on designing interactive systems*, pages 555–565, 2017.
- [43] Aarsh Trivedi, Vatsal Gor, and Zalak Thakkar. Chatbot generation and integration: A review. *International Journal of Advance Research, Ideas and Innovations in Technology*, 5(2):1308–1311, 2019.
- [44] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [45] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [46] Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. Pre-trained models for natural language processing: A survey. *Science China Technological Sciences*, pages 1–26, 2020.
- [47] Jieh-Sheng Lee and Jieh Hsiang. Patent claim generation by fine-tuning openai gpt-2. *World Patent Information*, 62:101983, 2020.
- [48] Jesse Vig. A multiscale visualization of attention in the transformer model. *arXiv preprint arXiv:1906.05714*, 2019.
- [49] Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing Xiang, et al. Abstractive text summarization using sequence-to-sequence rnns and beyond. *arXiv preprint arXiv:1602.06023*, 2016.
- [50] Courtney Napoles, Matthew R Gormley, and Benjamin Van Durme. Annotated gigaword. In *Proceedings of the Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction (AKBC-WEKEX)*, pages 95–100, 2012.
- [51] Donna Harman and Paul Over. The effects of human variation in duc summarization evaluation. In *Text Summarization Branches Out*, pages 10–17, 2004.
- [52] Philipp Scharpf, Moritz Schubotz, Abdou Youssef, Felix Hamborg, Norman Meuschke, and Bela Gipp. Classification and clustering of arxiv documents, sections, and abstracts, comparing encodings of natural and mathematical language. In *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020*, pages 137–146, 2020.
- [53] Emeric Dymont, Stéfan J Darmoni, Émeline Lejeune, Gaétan Kerdelhué, Jean-Philippe Leroy, Vincent Le-quertier, Stéphane Canu, and Julien Grosjean. Doc2vec on the pubmed corpus: study of a new approach to generate related articles. *arXiv preprint arXiv:1911.11698*, 2019.
- [54] Angela Joerin, Michiel Rauws, and Mary Lou Ackerman. Psychological artificial intelligence service, tess: delivering on-demand support to patients and their caregivers: technical report. *Cureus*, 11(1), 2019.
- [55] Sebastian Möller, Ute Jekosch, and Alexander Raake. New models predicting conversational effects of telephone transmission on speech communication quality. In *Sixth International Conference on Spoken Language Processing*, 2000.
- [56] Heung-Yeung Shum, Xiao-dong He, and Di Li. From eliza to xiaoice: challenges and opportunities with social chatbots. *Frontiers of Information Technology & Electronic Engineering*, 19(1):10–26, 2018.
- [57] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N Bennett, Kori Inkpen, et al. Guidelines for human-ai interaction. In *Proceedings of the 2019 chi conference on human factors in computing systems*, pages 1–13, 2019.

- [58] Binbin Xu, Cédric Gil-Jardiné, Frantz Thiessard, Eric Tellier, Marta Avalos, and Emmanuel Lagarde. Neural language model for automated classification of electronic medical records at the emergency room. the significant benefit of unsupervised generative pre-training. *arXiv preprint arXiv:1909.01136*, 2019.
- [59] Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30(4):681–694, 2020.
- [60] Kris McGuffie and Alex Newhouse. The radicalization risks of gpt-3 and advanced neural language models. *arXiv preprint arXiv:2009.06807*, 2020.
- [61] Aidan Boyd, Adam Czajka, and Kevin Bowyer. Deep learning-based feature extraction in iris recognition: Use existing models, fine-tune or train from scratch? In *2019 IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, pages 1–9. IEEE, 2019.
- [62] Robert G Newcombe. Interval estimation for the difference between independent proportions: comparison of eleven methods. *Statistics in medicine*, 17(8):873–890, 1998.
- [63] Robert G Newcombe. Two-sided confidence intervals for the single proportion: comparison of seven methods. *Statistics in medicine*, 17(8):857–872, 1998.
- [64] Annie J Olmstead, Jimin Lee, and Navin Viswanathan. The role of the speaker, the listener, and their joint contributions during communicative interactions: A tripartite view of intelligibility in individuals with dysarthria. *Journal of Speech, Language, and Hearing Research*, 63(4):1106–1114, 2020.
- [65] Gordon Pask. Conversation theory. *Applications in Education and Epistemology*, 1976.
- [66] Chathika Gunaratne, Chathurani Senevirathna, Chathura Jayalath, Nisha Baral, William Rand, and Ivan Garibay. A multi-action cascade model of conversation. In *5th International Conference on Computational Social Science*, URL <http://app.ic2s2.org/app/sessions/9kXqn5btgKKC5yfCvg/details>, 2019.
- [67] John Heritage. Conversation analysis and institutional talk: Analyzing distinctive turn-taking. *Dialoganalyse VI/2: Referate der 6. Arbeitstagung, Prag 1996*, 17(3), 2017.
- [68] Peter Gordon Roetzel. Information overload in the information age: a review of the literature from business administration, business psychology, and related disciplines with a bibliometric approach and framework development. *Business research*, 12(2):479–522, 2019.
- [69] Siegfried Streufert and Michael J Driver. Conceptual structure, information load and perceptual complexity. *Psychonomic Science*, 3(1):249–250, 1965.
- [70] Hiram Calvo and Alexander Gelbukh. Finding the most frequent sense of a word by the length of its definition. In *Mexican International Conference on Artificial Intelligence*, pages 1–8. Springer, 2014.
- [71] Hae-Young Kim. Statistical notes for clinical researchers: A one-way repeated measures anova for data with repeated observations. *Restorative dentistry & endodontics*, 40(1):91, 2015.
- [72] Hervé Abdi and Lynne J Williams. Tukey’s honestly significant difference (hsd) test. *Encyclopedia of research design*, 3(1):1–5, 2010.
- [73] Anao Zhang, Sunyoung Park, John E Sullivan, and Shijie Jing. The effectiveness of problem-solving therapy for primary care patients’ depressive and/or anxiety disorders: a systematic review and meta-analysis. *The Journal of the American Board of Family Medicine*, 31(1):139–150, 2018.
- [74] Arthur M Nezu, Thomas J D’Zurilla, et al. *Problem-solving therapy: A positive approach to clinical intervention*. Springer Publishing Company, 2006.
- [75] William L Hamilton, Kevin Clark, Jure Leskovec, and Dan Jurafsky. Inducing domain-specific sentiment lexicons from unlabeled corpora. In *Proceedings of the conference on empirical methods in natural language processing. conference on empirical methods in natural language processing*, volume 2016, page 595. NIH Public Access, 2016.
- [76] Scott A Crossley, Kristopher Kyle, and Danielle S McNamara. Sentiment analysis and social cognition engine (seance): An automatic tool for sentiment, social cognition, and social-order analysis. *Behavior research methods*, 49(3):803–821, 2017.
- [77] Philip J. Stone and Earl B. Hunt. A computer approach to content analysis: Studies using the general inquirer system. In *Proceedings of the May 21-23, 1963, Spring Joint Computer Conference*, AFIPS ’63 (Spring), page 241–256, New York, NY, USA, 1963. Association for Computing Machinery.
- [78] H. D. Lasswell and J. Z. Namenwirth. The lasswell value dictionary. *New Haven*, 1969.
- [79] Margaret M Bradley and Peter J Lang. Affective norms for english words (anew): Instruction manual and affective ratings. Technical report, Technical report C-1, the center for research in psychophysiology ..., 1999.

- [80] Mingqing Hu and Bing Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177, 2004.
- [81] Klaus R Scherer. What are emotions? and how can they be measured? *Social science information*, 44(4):695–729, 2005.
- [82] Saif Mohammad and Peter Turney. Emotions evoked by common words and phrases: Using mechanical turk to create an emotion lexicon. In *Proceedings of the NAACL HLT 2010 workshop on computational approaches to analysis and generation of emotion in text*, pages 26–34, 2010.
- [83] Saif M Mohammad and Peter D Turney. Crowdsourcing a word–emotion association lexicon. *Computational intelligence*, 29(3):436–465, 2013.
- [84] Erik Cambria, Catherine Havasi, and Amir Hussain. Senticnet 2: A semantic and affective resource for opinion mining and sentiment analysis. In *Twenty-Fifth international FLAIRS conference*, 2012.
- [85] Clayton Hutto and Eric Gilbert. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 8, 2014.
- [86] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. Ernie: Enhanced language representation with informative entities. *arXiv preprint arXiv:1905.07129*, 2019.
- [87] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [88] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [89] Lingpeng Kong, Cyprien de Masson d’Autume, Wang Ling, Lei Yu, Zihang Dai, and Dani Yogatama. A mutual information maximization perspective of language representation learning. *arXiv preprint arXiv:1910.08350*, 2019.
- [90] Lichao Sun, Kazuma Hashimoto, Wenpeng Yin, Akari Asai, Jia Li, Philip Yu, and Caiming Xiong. Advbert: Bert is not robust on misspellings! generating nature adversarial samples on bert. *arXiv preprint arXiv:2003.04985*, 2020.
- [91] Xiaodong Liu, Hao Cheng, Pengcheng He, Weizhu Chen, Yu Wang, Hoifung Poon, and Jianfeng Gao. Adversarial training for large neural language models. *arXiv preprint arXiv:2004.08994*, 2020.
- [92] Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. FreeLb: Enhanced adversarial training for natural language understanding. *arXiv preprint arXiv:1909.11764*, 2019.
- [93] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! *arXiv preprint arXiv:1904.12843*, 2019.
- [94] Chathika Gunaratne, Nisha Baral, William Rand, Ivan Garibay, Chathura Jayalath, and Chathurani Senevirathna. The effects of information overload on online conversation dynamics. *Computational and Mathematical Organization Theory*, 26(2):255–276, 2020.
- [95] Albert Gatt and Emiel Krahmer. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research*, 61:65–170, 2018.
- [96] Verena Rieser and Oliver Lemon. Natural language generation as planning under uncertainty for spoken dialogue systems. In *Empirical methods in natural language generation*, pages 105–120. Springer, 2009.
- [97] Ehud Reiter and Robert Dale. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87, 1997.
- [98] Donald O Case, James E Andrews, J David Johnson, and Suzanne L Allard. Avoiding versus seeking: the relationship of information seeking to avoidance, blunting, coping, dissonance, and related concepts. *Journal of the Medical Library Association*, 93(3):353, 2005.
- [99] Olivier Ferlatte, Travis Salway, Simon Rice, John L Oliffe, Ashleigh J Rich, Rod Knight, Jeffrey Morgan, and John S Ogrodniczuk. Perceived barriers to mental health services among canadian sexual and gender minorities with depression and at risk of suicide. *Community mental health journal*, 55(8):1313–1321, 2019.
- [100] Nneka Jon-Ubabuco and Jane Dimmitt Champion. Perceived mental healthcare barriers and health-seeking behavior of african-american caregivers of adolescents with mental health disorders. *Issues in mental health nursing*, 40(7):585–592, 2019.

- [101] Arrianna M Planey, Shardé McNeil Smith, Stephanie Moore, and Taylor D Walker. Barriers and facilitators to mental health help-seeking among african american youth and their families: a systematic review study. *Children and Youth Services Review*, 101:190–200, 2019.
- [102] Olabiyi Oluwatobi and Erik Mueller. Dlgnet: A transformer-based model for dialogue response generation. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 54–62, 2020.
- [103] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *Preprint*, 2018.
- [104] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- [105] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. Ernie 2.0: A continual pre-training framework for language understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8968–8975, 2020.
- [106] Feiyu Xu, Hans Uszkoreit, Yangzhou Du, Wei Fan, Dongyan Zhao, and Jun Zhu. Explainable ai: A brief survey on history, research areas, approaches and challenges. In *CCF international conference on natural language processing and Chinese computing*, pages 563–574. Springer, 2019.
- [107] Carlos Zednik. Solving the black box problem: a normative framework for explainable artificial intelligence. *Philosophy & Technology*, pages 1–24, 2019.
- [108] Dumitru Erhan, Aaron Courville, Yoshua Bengio, and Pascal Vincent. Why does unsupervised pre-training help deep learning? In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 201–208. JMLR Workshop and Conference Proceedings, 2010.
- [109] Charles Lewis Taylor. The world handbook tradition: Producing data for cross-national quantitative analysis. *Inter-university Consortium for Political and Social Research*, 1975.
- [110] Tim Bond. The law of confidentiality—a solution or part of the problem? *Legal issues in counselling and psychotherapy*, pages 123–143, 2002.
- [111] Asa Cooper Stickland and Iain Murray. Bert and pals: Projected attention layers for efficient adaptation in multi-task learning. In *International Conference on Machine Learning*, pages 5986–5995. PMLR, 2019.
- [112] Qian Yang, Aaron Steinfeld, Carolyn Rosé, and John Zimmerman. Re-examining whether, why, and how human-ai interaction is uniquely difficult to design. In *Proceedings of the 2020 chi conference on human factors in computing systems*, pages 1–13, 2020.