

ARTICLE

Dialogue Agents 101: A Beginner’s Guide to Critical Ingredients for Designing Effective Conversational Systems

Shivani Kumar, Sumit Bhatia, Milan Aggarwal, and Tanmoy Chakraborty

Indraprastha Institute of Information Technology, Delhi; shivaniku@iiitd.ac.in

Media and Data Science Research Lab, Adobe; sumit.bhatia@adobe.com

Media and Data Science Research Lab, Adobe; milaggar@adobe.com

Indian Institute of Technology, Delhi; tanchak@iiitd.ac.in

(Received xx xxx xxx; revised xx xxx xxx; accepted xx xxx xxx)

Abstract

Sharing ideas through communication with peers is the primary mode of human interaction. Consequently, extensive research has been conducted in the area of conversational AI, leading to an increase in the availability and diversity of conversational tasks, datasets, and methods. However, with numerous tasks being explored simultaneously, the current landscape of conversational AI has become fragmented. Consequently, initiating a well-thought-out model for a dialogue agent can pose significant challenges for a practitioner. Towards highlighting the critical ingredients needed for a practitioner to design a dialogue agent from scratch, the current study provides a comprehensive overview of the primary characteristics of a dialogue agent, the supporting tasks, their corresponding open-domain datasets, and the methods used to benchmark these datasets. We observe that different methods have been used to tackle distinct dialogue tasks. However, building separate models for each task is costly and does not leverage the correlation among the several tasks of a dialogue agent. As a result, recent trends suggest a shift towards building unified foundation models. To this end, we propose UNIT, a UNified dialogue dataset constructed from conversations of varying datasets for different dialogue tasks capturing the nuances for each of them. We then train a Unified dialogue foundation model, GPT-2^U and present a concise comparative performance of GPT-2^U against existing large language models. We also examine the evaluation strategies used to measure the performance of dialogue agents and highlight the scope for future research in the area of conversational AI with a thorough discussion of popular models such as ChatGPT.

1. Introduction

The significance of conversations as the fundamental medium of interaction transcends cultural boundaries (Dingemane and Floyd 2014). Consequently, interacting with machines and seeking information via conversational interfaces is an instinctive and familiar way for humans (Dalton et al. 2022) as evidenced by the success of dialogue systems such as Apple’s SIRI^a, Amazon’s

Competing interests: Shivani Kumar is pursuing her PhD at Indraprastha Institute of Information Technology Delhi. Sumit Bhatia and Milan Aggarwal are employed at Adobe. Tanmoy Chakraborty is employed at Indian Institute of Technology Delhi.

^a<https://www.apple.com/in/siri/>

Alexa^b, and most recently, ChatGPT^c. Moreover, dialogue-based systems^d have extensively been used for customer support (Botea et al. 2019; Feigenblat et al. 2021), mental health support (Kretzschmar et al. 2019), and counseling (Malhotra et al. 2022; Tewari et al. 2021).

Designing practical dialogue-based systems, however, is a challenging endeavour as there are important questions that one needs to answer before embarking on developing such a system. Critical considerations include determining the types of queries the system should anticipate (e.g., chit-chat versus informational), deciding whether to incorporate an external knowledge source, and determining the level of natural language understanding the system should support. Previous surveys in the field of dialogue-based systems have predominantly focused on examining specific system components or narrow subsets of tasks and techniques. For instance, recent surveys have delved into areas such as dialogue summarization (Tuggener et al. 2021; Feng et al. 2022a), text-to-SQL (Qin et al. 2022), question answering (Pandya and Bhatt 2021), dialogue management using deep learning (Chen et al. 2017a) and reinforcement learning (Dai et al. 2021b).

While the surveys noted above provide comprehensive insights into their respective domains, this abundance of information can make it overwhelming for both novice and experienced researchers and professionals to identify the essential components required for building their dialogue-based systems. In contrast, we adopt a broader perspective and offer a panoramic view of the various constituents comprising a dialogue-based system, elucidate the individual tasks involved in their development, and highlight the typical datasets and state-of-the-art methodologies employed for designing and evaluating these components. Consequently, the title ‘Dialogue Agents 101’ is a deliberate choice aiming to convey that the article serves as an introductory guide or primer to the fundamental concepts and principles associated with dialogue agents. In academic settings, ‘101’ is often used to denote introductory or basic-level courses, and here, it suggests that the article provides foundational knowledge for readers who may be new to the topic of dialogue agents. With this comprehensive survey, we aspire to assist beginners and practitioners in making well-informed decisions while developing systems for their applications. Our specific objective is to comprehensively encompass all **prominent open-source textual English** dialogue datasets across major dialogue tasks. That is, every dataset under consideration in our study meets four conditions: (i) it must be widely recognized within its respective field; ii) it should incorporate a textual component in both input and output; (iii) it must be publicly accessible, and; and (iv) it must be designed for English.

To identify relevant material for our survey, we conducted a thorough search of the Papers With Code website^e to identify all relevant tasks and datasets related to dialogue agents. Our goal was to gather and systematically organize different types of tasks that may be required for developing various dialogue-agents; and understand the methods for performing these tasks, and datasets that are typically used to train and evaluate models for these tasks. From the initial list obtained from Papers With Code, we then queried Google Scholar for publications and followed the citation threads to gather relevant literature for each task, encompassing datasets and articles proposed well before the establishment of the platforms. We emphasize that while Papers With Code functioned as our reference for locating pertinent literature, its principal values lay in pinpointing the key problem statements investigated within the domain of dialogue agents.

While delving into contemporary deep learning methods in this investigation, it is crucial to acknowledge the rich history of research in dialogue agents. Long before the advent of deep learning, researchers were actively engaged in developing computational methods to facilitate meaningful interactions between machines and humans (Bayer et al. 2001; Weizenbaum 1966). In the nascent stages of dialogue agent development, researchers heavily relied on rule-based

^b<https://alexa.amazon.com/>

^c<https://openai.com/blog/chatgpt>

^dWe use dialogue-based systems, chatbots, conversational systems, and dialogue agents interchangeably in this article.

^e<https://paperswithcode.com/>

systems (Webb 2000; McTear 2021). Human experts meticulously crafted these systems, incorporating predefined rules and decision trees to interpret user inputs and generate appropriate responses. Classification tasks, such as intent detection and slot filling, often involved rule-based pattern matching (De and Kopparapu 2010; Ren et al. 2018) and template-based approaches (Onyshkevych 1993; McRoy et al. 2003) to identify the user’s intention based on specific keywords or syntactic structures. Generative tasks, such as response generation, posed a significant challenge without deep learning techniques. Early approaches leveraged handcrafted templates (Chu-Carroll and Carberry 1998; Weizenbaum 1966), where responses were generated by combining predefined phrases or sentences. This method, however, lacked the flexibility to generate contextually relevant and nuanced responses, hindering the natural flow of conversations.

As computational capabilities advanced, statistical methods started gaining traction in dialogue agent development. Hidden Markov Models (HMMs) (Rabiner and Juang 1986) and finite-state machines (Ben-Ari and Mondada 2018) were applied to model the probabilistic nature of language and user interactions (Williams 2003; Williams et al. 2005). These models enabled a more dynamic and probabilistic approach to intent detection and slot filling, contributing to the improvement of dialogue system performance (Hussein and Granat 2002; Zhao et al. 2004). From rule-based systems and template-based approaches to early statistical models, researchers laid the groundwork for the sophisticated deep learning methodologies that dominate the contemporary landscape we aim to study in this survey. To summarize, our key contributions are as follows.

- (1) We propose an **in-depth taxonomy** for different components and modules involved in building a dialogue agent (Figure 1). We take a practitioner’s view point and develop the taxonomy in terms of features of the underlying system and discuss at length the role played by each of the features in the overall system (Section 2).
- (2) Next, we present a comprehensive overview of different tasks and datasets in the literature and relate them to the features as identified in the proposed taxonomy (Table 1). We identify eleven broad categories of tasks related to dialogue-based systems and present a detailed overview of different methods for each task and datasets used for evaluating these tasks (Section 3). Our goal is to help the reader identify key techniques and datasets available for the tasks relevant to their applications.
- (3) We present UNIT^f, a large scale **unified dialogue dataset**, consisting of more than 4.8M dialogues and 441M tokens, which combine the various dialogue datasets described in Section 6. Since UNIT is made from the dialogues of open-sourced datasets, it is free to use for any research purposes. This effort is motivated by the recent trends suggesting a shift towards building unified foundation models (Zhou et al. 2023a) that are pre-trained on large datasets and generalize to a variety of tasks. We make UNIT available to the research community with a goal to spark research efforts towards development of foundation models optimized for dialogues. We use UNIT to further pretrain popular open dialogue foundation models and show how it can help improving their performance on various dialogue tasks (Section 6).

2. Designing a Dialogue Agent

Before developing a dialogue agent, several crucial decisions must be made to determine the appropriate architecture for the agent. Figure 1 illustrates a comprehensive overview of these decisions, which provides a taxonomic framework for structuring the development process. A clear understanding of the end goal we aim to achieve from a dialogue agent is crucial for effective communication (Pomerantz and Fehr 2011). For instance, questions such as “Do we want the dialogue agent to carry out goal-oriented or chit-chat conversations?” and “Does the agent need any external knowledge to answer user queries?” should be answered. Figure 2 highlights the

^fWe make UNIT public on <https://github.com/LCS2-IIITD/UNIT.git>

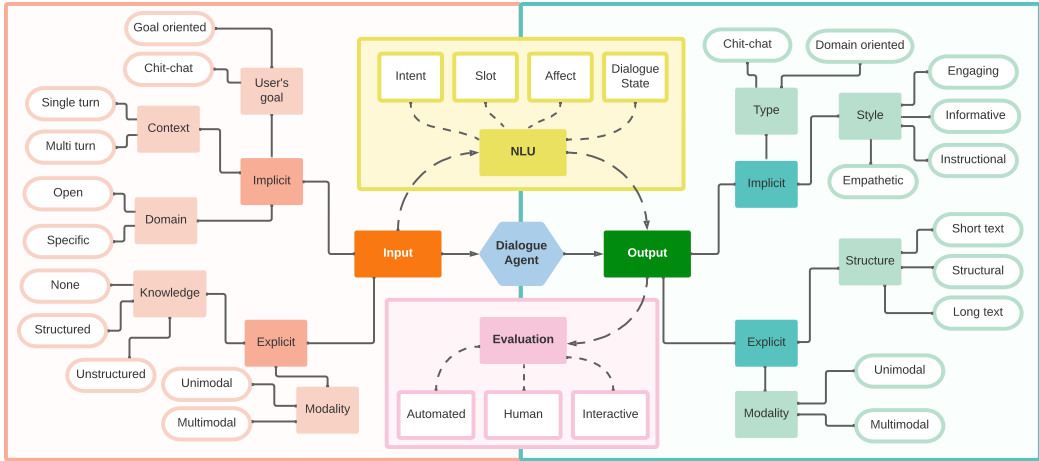


Figure 1: A taxonomic overview of a dialogue agent. The major components for designing a complete pipeline of a dialogue agent are – input(s), natural language understanding (NLU), generated output(s), and model evaluation. Each component can be further divided based on the characteristics required in the final dialogue agent.

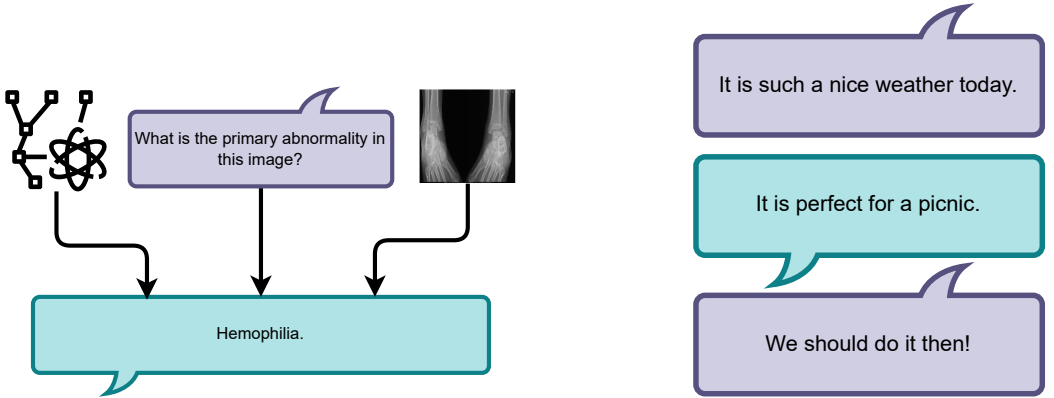
different type of dialogues based on the different attributes of the input and output of the system as discussed below.

2.1 Input to the System

After establishing the end-goal of our dialogue agent, it is essential to determine the various factors that will inform the input to the agent (Harms et al. 2019). Our contention is that the input can possess both implicit and explicit properties, depending on the task at hand.

Implicit Attributes. We classify the characteristics of the input which are not explicitly apparent from the input as implicit attributes of the input. This inherent information can be decided based on three aspects – the user’s goal (Muise et al. 2019), the domain of the dialogues (Budzianowski et al. 2018), and the context needed to carry out the end task (Kiela and Weston 2019). Depending on the objective of the dialogue agent, the user could want to achieve some goal, such as making a restaurant reservation, booking an airline ticket, or resolving technical queries. For such goal-oriented dialogue agents, the input from the user is expected to differ from that received for general chit-chat (Muise et al. 2019). Goal-oriented dialogue agents are often designed to operate within a particular domain, while chit-chat-based agents are more versatile and are expected to handle a broader range of conversations (Zhang et al. 2018). In addition to the user’s goal and the agent’s domain, the conversation context also plays a crucial role in achieving the agent’s objective (Kiela and Weston 2019). For example, utterance-level intent detection may not require understanding deep conversation context, while summarizing dialogues would require a complete understanding of the context (Gliwa et al. 2019).

Explicit Attributes. Apart from the implicit aspects of the dialogue agent’s input, various input characteristics are external in nature and should be considered while building a dialogue agent. These aspects constitute the input modality (Jovanovic and Leeuwen 2018) and any additional knowledge supplied to the agent (Dinan et al. 2019). Input can be unimodal, such as text or audio, or in a combination of modalities, such as an image and associated text, as in the case of visual question-answering systems (Parvaneh et al. 2019). Furthermore, additional knowledge may be required to generate appropriate responses. For example, in a chit-chat setting, the agent may need



(a) Goal-oriented single-turn dialogue of a single domain with structured knowledge and multi-modal input

(b) Chit-chat multi turn dialogue of an open domain with no external knowledge and unimodal input

Figure 2: Dialogues highlighting different attributes of a dialogue agent input and output.

to possess commonsense knowledge (Strathearn and Gkatzia 2022), while in a question-answering setting, the agent may need to access relevant documents to provide accurate responses (Feng et al. 2020). Therefore, any explicit knowledge supplied to the dialogue agent can be structured, like a tree or a tuple, or unstructured, like a document.

2.2 Natural Language Understanding

After receiving input from the user, the subsequent step involves comprehension (Liu et al. 2021b). Regardless of whether the task is domain-specific or open-domain, specific attributes of the input must be identified to determine the required output. We identify four primary attributes that need to be identified from the input text – the user’s **intent** (Casanueva et al. 2020), any **slots** needed to fulfill the intent (Weld et al. 2022a), **affective** understanding of the input (Ruusuvuori 2012), and the **dialogue state** of the input utterance (Balaraman et al. 2021). While intent and slots are directly useful for a domain-specific agent to effectively complete a task, affect understanding and dialogue state tracking is also critical for a chit-chat-based agent. Affect understanding involves comprehending the user’s emotion (Poria et al. 2019), sarcasm (Castro et al. 2019), and amusement (Bedi et al. 2021) in the input utterance. Furthermore, dialogue state tracking checks the type of utterance received by the agent, such as question, clarification, or guidance. Understanding these aspects is essential to determine the utterance’s underlying meaning and provide relevant responses for the task.

2.3 Output of the System

The output generated by the dialogue agent, akin to its input, possesses both implicit and explicit attributes, described below.

Implicit Attributes. Implicit attributes refer to the output’s type (Rastogi et al. 2020) and style (Su et al. 2020; Troiano et al. 2023), while explicit attributes pertain to its modality (Sun et al. 2022b) and structure (Yu et al. 2018). Congruent to the user’s goal in the input scenario, the type of attribute should be decided based on the end task needed to be performed by the dialogue agent. Depending on the end task of the agent, the resulting output can be informative (Feng et al. 2020),

engaging (Zhang et al. 2018), instructional (Strathearn and Gkatzia 2022), or empathetic (Rashkin et al. 2019). For instance, a question-answering-based bot should be informative, while a cooking recipe bot should be more instructional. Both bots need not be empathetic in nature.

Explicit Attributes. While the inherent properties of the output text are critical to assess, the explicit attributes, such as modality and structure, must be considered before finalizing the dialogue agent’s architecture. Modality decides whether the required output is unimodal (such as text) or multimodal (such as text with an image). Moreover, the output can be structured differently based on the task at hand. For instance, tasks such as text-to-SQL (Yu et al. 2018) conversion require the output to adhere to a certain structure. After considering various aspects of the input, output, and understanding based on the end task, the generated output is evaluated to gauge the performance of the resultant dialogue agent (Deriu et al. 2021). A detailed discussion about the evaluation can be found in Section 5.

3. Tasks, Datasets and Methods

By drawing upon the taxonomy depicted in Figure 1 and existing literature, we identify *eleven* distinct tasks related to dialogue that capture all necessary characteristics of a dialogue agent. In order to construct a dialogue agent, a practitioner must be aware of these tasks, which can be classified into two primary categories – generative and classification. Specifically, the identified tasks include **Dialogue Rewrite (DR)** (Elgohary et al. 2019), **Dialogue Summary (DS)** (Gliwa et al. 2019; Chen et al. 2021b), **Dialogue to Structure (D2S)** (Yu et al. 2019 2018; Gupta et al. 2018), **Question Answering(QA)** (Zhou et al. 2018; Reddy et al. 2019; Aliannejadi et al. 2020; Cui et al. 2020), **Knowledge Grounded Response (KGR)** (Yusupov and Kuratov 2018; Feng et al. 2020; Zhang et al. 2018; Weston et al. 2015; Dziri et al. 2022; Moon et al. 2019; Strathearn and Gkatzia 2022), **Chit-Chat (CC)** (Sevegnani et al. 2021; Kim et al. 2022c; Young et al. 2022; Zhang et al. 2022; Kim et al. 2022a; Jurafsky et al. 1997), and **Task-Oriented Dialogues (TOD)** (Lowe et al. 2015; Chen et al. 2021a; Weston et al. 2015; Lin et al. [n.d.]; He et al. 2018; Karadzhov et al. 2021; Shalyminov et al. 2019) in the generative category, and **Intent Detection(ID)** (Casanueva et al. 2020; Larson et al. 2019; Liu et al. 2021c; Rastogi et al. 2020), **Slot Filling (SF)** (Coope et al. 2020), **Dialogue State Tracking (DST)** (Eric et al. 2020), and **Affect Detection (AD)** (Poria et al. 2019; Li et al. 2017; Castro et al. 2019; Rashkin et al. 2019) in the classification category. Table 1 summarises all the datasets considered in this study for each of the mentioned tasks and illustrates the characteristics satisfied by each of these tasks from the taxonomy. As we delve into the details of each task type in the forthcoming sections, it is noteworthy to highlight a few observations obtained from the presented table.

- In dialogue datasets featuring chit-chat conversations, an inclination towards characteristics indicative of open domain, multi-turn interactions, and the absence of external knowledge is observed. Notably, a prevalent trend emerges in the generation of similar output within such datasets. An identified gap in the existing landscape pertains to the scarcity of datasets integrating external knowledge with chit-chat dialogues. Recognizing the potential enrichment that associated knowledge, particularly commonsense (Ghosal et al. 2020), can bring to dialogues, it becomes a potential future research area.
- For instances where the dataset comprises goal-oriented conversations, it is probable that the dataset is tailored to a specific domain, assisted with either structured or unstructured knowledge linked to it. Goal-oriented dialogues typically centre around specific tasks like booking airline tickets, scheduling doctor appointments, or securing restaurant reservations. Notably, these ‘goals’ can extend beyond specific tasks to encompass aspects such as the accomplishment of the goal of dialogue engagement (Gottardi et al. 2022). Intriguingly, such goal-orientation does not necessarily confine the dialogue to a predefined domain, allowing for an open-domain context. A prospective avenue for research lies in the development

Type	Task	Datasets	Input														Output										Size
			Implicit							Explicit							Implicit					Explicit					
			User's goal		Domain		Context			Modality		Knowledge			Type		Style			Modality		Structure					
			CC	GO	Open	Spc	ST	MT	U	M	None	Unstr	Str	CC	GO	Eng	Inf	Instr	Emp	U	M	Shor	Long	Struct			
Transformation	DR	CANARD (Elgohary et al. 2019)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	40		
	DS	DialogSum (Chen et al. 2021b)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	13		
		SAMSum Corpus (Gliwa et al. 2019)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	16		
D2S	CoSQL (Yu et al. 2019) SPIDER (Yu et al. 2018) TOP (Gupta et al. 2018)	✓	✓	-	✓	✓	-	✓	-	-	-	✓	✓	-	✓	-	✓	-	✓	-	-	-	✓	2			
		-	✓	-	✓	✓	-	✓	-	-	-	✓	✓	-	✓	-	✓	-	✓	-	-	-	✓	10			
		-	✓	-	✓	✓	-	✓	-	✓	-	-	-	✓	✓	-	✓	-	✓	-	-	-	✓	44			
Generative	QA	CMUDoG (Zhou et al. 2018)	-	✓	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	4			
		CoQA (Reddy et al. 2019)	-	✓	-	✓	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	127			
		ClariQ (Aliannejadi et al. 2020)	-	✓	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	1k			
		Mutual (Cui et al. 2020)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	8		
	KGR	ConvAI (Yusupov and Kuratov 2018)	-	✓	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	2			
		Doc2Dial (Feng et al. 2020)	-	✓	-	✓	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	4			
		PersonaChat (Zhang et al. 2018)	✓	-	-	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	19			
		bAbI (Weston et al. 2015)	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	161		
		FaithDial (Dziri et al. 2022)	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	32			
		OpenDialKG (Moon et al. 2019)	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	15		
Task2Dial (Strathearn and Gkatzia 2022)	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	1				
Response generation	CC	OTters (Sevgnani et al. 2021)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	8		
		ProsocialDialog (Kim et al. 2022c)	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	5	
		FusedChat (Young et al. 2022)	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	10		
		mDIA (Zhang et al. 2022)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	12		
		SODA (Kim et al. 2022a)	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	1k	
		Switchboard-1 (Jurafsky et al. 1997)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	2		
TOD	Ubuntu Dialogue Corpus (Lowe et al. 2015)	-	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	1k		
	ABCD (Chen et al. 2021a)	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	10			
	BitOD (Lin et al. [n.d.])	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	7			
	CraiglistBargains (He et al. 2018)	-	✓	-	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	6			
	DeliData (Karadzhov et al. 2021)	-	✓	✓	-	-	✓	✓	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	0.5		
	MetalWOz (Shalymov et al. 2019)	-	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	10		
Classification	ID	Banking77 (Casanueva et al. 2020)	-	✓	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	13		
		CLINC150 (Larson et al. 2019)	-	✓	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	23		
		HWU64 (Liu et al. 2021c)	-	✓	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	11		
		SGD (Rastogi et al. 2020)	-	✓	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	16		
	SF	Restaurant8k (Coope et al. 2020)	-	✓	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	11		
	DST	MultiWOZ2.1 (Eric et al. 2020)	-	✓	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	10		
	AD	DailyDialogue (Li et al. 2017)	✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	11		
MELD (Poria et al. 2019)		✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	1			
MUSTARD (Castro et al. 2019)		✓	-	✓	-	-	✓	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	6			
Empathetic Dialogues (Rashkin et al. 2018)		✓	-	✓	-	-	✓	✓	-	-	-	✓	-	✓	-	✓	-	✓	-	✓	-	✓	-	24			

Table 1. : Characteristic of each task based on the taxonomic characteristic of a dialogue agent. Size indicates an approximate value expressed in thousands (k). Abbreviations – DR: Dialogue Rewrite, DS: Dialogue Summary, D2S: Dialogue to Structure, QA: Question Answering, KGR: Knowledge Grounded Response, CC: Chit-chat, TOD: Task Oriented Dialogues, ID: Intent Detection, SF: Slot Filling, DST: Dialogue State Tracking, AD: Affect Detection, CC: Chit-chat, GO: Goal Oriented, Spc: Specific, ST: Single Turn, MT: Multi Turn, U: Unimodal, M: Multimodal, Unstr: Unstructured, Str: Structured, Eng: Engaging, Inf: Informative, Instr: Instructional, Emp: Empathetic.

of more open-domain, goal-oriented dialogue datasets that focus more on conversational goals like user engagement.

- The chit-chat setting exhibits the predominant trend of producing extensive and engaging dialogue output (Gottardi et al. 2022). In contrast, the goal-oriented setting commonly yields responses characterized by informativeness, instructional clarity, and brevity (Muise et al. 2019). Intriguingly, datasets combining both goal-oriented and chit-chat conversations are notably sparse, despite real-world dialogues frequently encompassing a fluid interchange

between these conversational types (Shuster et al. 2022). The presence of such datasets could substantially enhance the research community’s capabilities and insights.

3.1 Generative Dialogue Tasks

Generative dialogue tasks require the handling of diverse input and output characteristics (Chen et al. 2017b). These tasks can be classified into two distinct types – transformation and response generation. In transformation tasks, the output of the given input conversation is not the subsequent response but rather some other meaningful text, such as a dialogue summary (Gliwa et al. 2019). On the other hand, response generation tasks involve generating the next response in the dialogue, given an input context (Zhang et al. 2020b).

3.1.1 Transformation Tasks

Dialogue Rewrite (DR). This task involves the challenging process of modifying a given conversational utterance to better fit a specific social context or conversational objective, while retaining its original meaning. To explore this task further, we turn to the CANARD dataset (Elgohary et al. 2019). This dataset is specifically designed for rewriting context-dependent questions into self-contained questions that can be answered independently by resolving all coreferences. The objective is to ensure that the new question has the same answer as the original one. Quan et al. (2019) and Martin et al. (2020) proposed the TASK and MuDoCo datasets, respectively, focusing on rewriting dialogues in a way that coreferences and ellipsis are resolved. Huang et al. (2021) combined sequence labelling and autoregression techniques to restore utterances without any coreferences. In contrast, Jiang et al. (2023) shaped the dialogue rewrite task as sentence editing and predicted edit operations for each word in the context. Other methods also use knowledge augmentation (Ke et al. 2022), reinforcement learning (Chen et al. 2022b), and the copy mechanism (Quan et al. 2019).

Key challenges. Despite achieving a reasonable performance in the dialogue rewrite task, some challenges remain, with the major obstacle being the inclusion of new words in the ground truth annotations that are difficult to incorporate into the predicted rewrite (Liu et al. 2020b). In order to mitigate this challenge, many studies have explored the methods of lexicon integration (Lee et al. 2023; Czarnowska et al. 2020), open-vocabulary (Raffel et al. 2020; Hao et al. 2021; Vu et al. 2022), and context-aware encoding (Xiao et al. 2020; Vinyals et al. 2015).

Dialogue summary (DS). Dialogues, despite their importance in communication, can often become lengthy and veer off-topic. This can make it challenging to extract the meaningful content from the entire conversation. To overcome this issue, the task of dialogue summarization has emerged. Dialogue summarization presents a concise account of the key topics, ideas, and arguments discussed during the conversation. There are two prominent datasets that address the challenge of dialogue summarization: the SAMSum (Gliwa et al. 2019) and DialogSum (Chen et al. 2021b) corpora consisting of dialogues and their corresponding summaries. The SAMSum dataset consists of dialogues that were curated by linguists who are fluent in English and who attempted to simulate messenger-like conversations. While DialogSum consists of face-to-face spoken dialogues covering various daily life topics such as schooling, work, and shopping. The dialogues are present in the textual format in both datasets. Other datasets such as QMSum (Zhong et al. 2021), MediaSum (Zhu et al. 2021), DiDi (Liu et al. 2019), CCCS (Favre et al. 2015), Telemedicine (Joshi et al. 2020), CRD3 (Rameshkumar and Bailey 2020), Television Shows (Zechner and Waibel 2000), AutoMin (Nedoluzhko et al. 2022), and Clinical Encounter Visits

(Yim and Yetisgen 2021) are also constructed for the task of dialogue summarisation. For a detailed guide on the task, we redirect the readers to the extensive survey conducted by Tugener et al. (2021). Many architectures have been proposed to solve the task of dialogue summarisation. Liang et al. (2023) uses topic-aware Global-Local Centrality (GLC) to extract important context from all sub-topics. By combining global- and local-level centralities, the GLC method guides the model to capture salient context and sub-topics while generating summaries. Other studies have utilized contrastive loss (Halder et al. 2022), multi-view summary generation (Chen and Yang 2020), post-processing techniques improving the quality of summaries (Lee et al. 2021), external knowledge incorporation (Kim et al. 2022b), multimodal summarisation (Atri et al. 2021), and methods to reduce hallucinations in generated summaries (Liu and Chen 2021; Narayan et al. 2021; Wu et al. 2021b).

Key challenges. With the help of pre-trained language models, current methods are adept at converting the original chat into a concise summary. Nonetheless, these models still face challenges in selecting the crucial parts and tend to generate hallucinations (Feng et al. 2022a). In the case of longer dialogues, the models may exhibit bias towards a specific part of the chat, such as the beginning or end, producing summaries that are not entirely satisfactory (Dey et al. 2020). Many studies explore novel attention mechanism with topic modeling (Xiao et al. 2020), reinforcement learning and differential rewards (Chen et al. 2023; Italiani et al. 2024; Zhang et al. 2023), and knowledge augmentation with fact-checking (Hua et al. 2023; Hwang et al. 2023) to mitigate these challenges.

Dialogue to structure (D2S). Although natural language is the fundamental way humans communicate, the interaction between humans and machines often requires a more structured language such as SQL or syntactic trees. Tasks such as *Text-to-SQL* and *Semantic Parsing* seek to bridge the gap between natural language and machine-understandable forms of communication. To address this, four prominent datasets have been developed – CoSQL (Yu et al. 2019), SPIDER (Yu et al. 2018), and WikiSQL (Zhong et al. 2017) for text-to-sql, which are composed of pairs of natural language queries paired with their corresponding SQL queries, and the Task Oriented Parsing (TOP) dataset (Gupta et al. 2018) for semantic parsing which contains conversations that are annotated with hierarchical semantic representation for task-oriented dialogue systems. There are numerous approaches to handling these datasets, including encoder/decoder models with decoder constraints (Wang et al. 2019b; Yin and Neubig 2017), large language models without any constraints (Suhr et al. 2020; Lin et al. 2020), final hypothesis pruning (Scholak et al. 2021), span-based extraction (Pasupat et al. 2019; Meng et al. 2022), data augmentation (Lee et al. 2022; Xuan 2020), and ensembling techniques (Einolghozati et al. 2018).

Key challenges. Despite recent advancements in D2S type tasks, there remains a scarcity of high-quality resources related to complex queries (Lee et al. 2022). Furthermore, the performance of D2S models tends to be suboptimal when encountering small perturbations, such as synonym substitutions or the introduction of domain-specific knowledge in the input (Qin et al. 2022). Existing studies explore the areas of data augmentation with resource creation to solve this challenge (Joshi et al. 2022; Min et al. 2020). Enhancing robustness and handling perturbation (Yu et al. 2023; Jia et al. 2019) are other possible solutions to the challenge of brittleness in the D2S tasks. Further research in this direction could yield valuable insights.

3.1.2 Response Generation

Question Answering (QA). Dialogue agents must possess the ability to ask relevant questions in order to engage the participants by introducing interesting topics via questions in general chat

setting (Gottardi et al. 2022), and provide appropriate answers to user inquiries, to remain authentic in the QA setting (Elgohary et al. 2019). As a result, Question Answering (QA) is a crucial task for dialogue agents to perform competently. To this end, datasets such as CMUDoG (Zhou et al. 2018), CoQA (Reddy et al. 2019), SQuAD (Rajpurkar et al. 2016 2018), ClariQ (Aliannejadi et al. 2020), and Mutual (Cui et al. 2020) are among the most notable and widely used for the purpose of training and evaluating QA systems. If external knowledge is used to answer questions, the task can be termed as knowledge-grounded question answering (Meng et al. 2020). The CMUDoG, CoQA, and SQuAD datasets are examples of this category. The FIRE model (Gu et al. 2020) utilizes context and knowledge filters to create context- and knowledge-aware representations through global and bidirectional attention. Other methods include multitask learning (Zhou and Small 2020), semantic parsing (Berant and Liang 2014; Reddy et al. 2014), knowledge-based grounding (Yih et al. 2015; Liang et al. 2017), and information-retrieval based methods (Bordes et al. 2015; Dong et al. 2015). On the other hand, the ClariQ and Mutual datasets does not contain any external knowledge. Komeili et al. (2022) have proposed using the Internet as a source for obtaining relevant information. In contrast, Hixon et al. (2015) proposes to learn domain from conversation context. Zero-shot approaches (Wang et al. 2023b), adversarial pretraining (Pi et al. 2022), convolution networks (Liu et al. 2022a), and graph based methods (Ouyang et al. 2021) are also used to solve the task of QA.

Key challenges. In the field of discourse-based question answering, which requires models to consider both deep conversation context and potential external knowledge, anaphora resolution still poses a significant challenge that necessitates further investigation (Pandya and Bhatt 2021). Additionally, capturing long dialogue context (Christmann et al. 2022) and preventing topical drift (Venkataram et al. 2020) offers other research direction. Many studies explore these challenges and propose viable solutions to mitigate them (Lin et al. 2021; Wu et al. 2023b). However, a reliable solution still needs more research in the field.

Knowledge grounded response (KGR). Similar to knowledge-grounded question answering, knowledge-grounded response generation is a task that utilizes external knowledge to generate relevant responses. Some of the primary datasets related to knowledge grounding include ConvAI (Yusupov and Kuratov 2018), Doc2Dial (Feng et al. 2020), PersonaChat (Zhang et al. 2018), bAbI (Weston et al. 2015), FaithDial (Dziri et al. 2022), OpenDialKG (Moon et al. 2019), and Task2Dial (Strathearn and Gkatzia 2022). Most methods that aim to solve the task of knowledge grounded response generation, like knowledge grounded QA, uses a two step approach of retrieval and generation (Zhan et al. 2021; Wu et al. 2021a), graph-based approach (Wang et al. 2020; Li et al. 2021a), reinforcement learning approach (Hedayatnia et al. 2020), and retrieval-free approaches (Xu et al. 2022).

Key challenges. The current trend in knowledge grounded response generation is to use a two-step approach of retrieval and generation, which increases the complexity of the system (Zhou et al. 2022). Recently, researchers such as Xu et al. (2022) and Zhou et al. (2022) have explored ways to bypass the retrieval step and produce more efficient models. Further research in this direction can improve the efficiency of systems.

Chit-chat (CC). The primary goal of a dialogue agent is to generate responses, whether it is for chit-chat based dialogues or task-oriented dialogues. This section will specifically focus on the response generation for chit-chat agents. While there are numerous dialogue datasets available that contain chit-chat dialogues and can be used as training data, such as PersonaChat (Zhang et al. 2018), MELD (Poria et al. 2019), DailyDialogue (Li et al. 2017), MUSTARD (Castro et al. 2019), and Mutual (Cui et al. 2020), there are some datasets specifically curated for the task of chit-chat

generation. Examples of such datasets include OTTers (Sevegnani et al. 2021), ProsocialDialog (Kim et al. 2022c), FusedChat (Young et al. 2022), mDIA (Zhang et al. 2022), SODA (Kim et al. 2022a), and the Switchboard-1 corpus (Jurafsky et al. 1997). Major approaches used to generate responses for chit-chat dialogue agents include the use of contrastive learning (Li et al. 2022a, 2021b; Cai et al. 2020), continual learning (Liu et al. 2022c; Liu and Mazumder 2021; Mi et al. 2020), and Transformer based methods (Liu et al. 2020a; Cai et al. 2019; Oluwatobi and Mueller 2020).

Key challenges. Typical challenges with chit-chat agents, such as inconsistency, unfaithfulness, and an absence of a uniform persona, persist (Liu et al. 2017a). Furthermore, the ineffective management of infrequently used words is another tenacious issue (Shum et al. 2020). However, current advancements, such as Reinforcement Learning from Human Feedback (RLHF) (Christiano et al. 2017; Stiennon et al. 2020), help in minimising these issues.

Task-oriented dialogues (TOD). To generate domain-specific responses, task-oriented dialogue agents require a specialized approach. Fortunately, there are several datasets available that feature domain-oriented dialogues, including the Ubuntu Dialogue Corpus (Lowe et al. 2015), ABCD (Chen et al. 2021a), bAbI (Weston et al. 2015), BiTOD (Lin et al. [n.d.]), CraiglistBargains (He et al. 2018), DeliData (Karadzhov et al. 2021), and MetalWOz (Shalyminov et al. 2019). Generating task-oriented dialogues follows a similar approach to open domain dialogues, utilizing reinforcement learning (Khandelwal 2021; Lipton et al. 2018; Liu et al. 2017b), graph based methods (Yang et al. 2020; Liu et al. 2021a; Andreas et al. 2020), and Transformer based methods (Chawla et al. 2020; Parvaneh et al. 2019).

Key challenges. The current datasets in this area feature restrictive input utterances, where necessary information is explicit and simple to extract (Zhang et al. 2020c). Conversely, natural conversations necessitate extracting implicit information from user utterances to generate a response (Zhou et al. 2022). A few studies explore advanced attention mechanisms (Qu et al. 2024), interactive learning (Yang et al. 2022) and dialogue augmentation (Liu et al. 2022b) to capture implicit contextual information from the text. Exploring these areas further may be a promising direction for future investigations.

3.2 Classification Tasks

Figure 1 shows that dialogue classification encompasses additional tasks, including intent detection, slot filling, dialogue state tracking, and affect detection. In the following sections, we provide a detailed explanation of each of these tasks.

Intent detection (ID). Identifying the user’s objectives in a conversation is crucial, particularly in goal-oriented dialogues. Intent detection aims to achieve this objective by analyzing text and inferring its intent, which can then be categorized into predefined groups. Given its importance, there has been significant research into intent detection, with several datasets proposed for this task, such as the DialoGLUE (Mehri et al. 2020) benchmark’s Banking77 (Casanueva et al. 2020), CLINC150 (Larson et al. 2019), HWU64 (Liu et al. 2021c), and the Schema Guided Dialogue (SGD) Dataset (Rastogi et al. 2020). Table 1 illustrates the taxonomic characteristics these datasets satisfy. It can be observed that they all follow a similar pattern of being goal-oriented, domain specific, and single turn with no external knowledge associated with them. The DialoGLUE leaderboard[§] indicates that a model called SAPCE2.0 gives exceptional performance

[§]<https://eval.ai/web/challenges/challenge-page/708/leaderboard/1943>

across all intent detection tasks. In addition, other approaches include utilizing contrastive conversational finetuning (Vulić et al. 2022), dual sentence encoders (Casanueva et al. 2020), and incorporating commonsense knowledge (Siddique et al. 2021).

Key challenges. The primary obstacle in intent detection involves the tight decision boundary of the learned intent classes within intent detection models (Weld et al. 2022b). Furthermore, given the dynamic nature of the world, the number and types of intents are constantly evolving, making it essential for intent detection models to be dynamic (Weld et al. 2022a). Recent developments have explored ensemble learning (Zhou et al. 2023b) along with Bayesian approaches (Zhang et al. 2019; Aftab et al. 2021) to mitigate the said challenge. Further, learning paradigms such as incremental learning (Paul et al. 2022; Hrycyk et al. 2021) and meta-learning (Li and Zhang 2021; Liu et al. 2022d) also prove to be beneficial in this field. However, a detailed future investigation in this domain is needed.

Slot Filling (SF). To effectively achieve a specific intent, a dialogue agent must possess all the necessary information required for task completion. These crucial pieces of information are commonly referred to as slots. It is worth noting that intent detection and slot filling often go hand in hand. As a result, the SGD dataset described in Section 3.2 includes slot annotations and can serve as a benchmark for evaluating slot-filling performance. Additionally, the Restaurant8k (Coope et al. 2020) dataset is another prominent dataset in the domain of slot filling. Methods that solve the slot-filling task often involve using CNN (Lecun et al. 1998) and CRF (Ma and Hovy 2016; Lample et al. 2016) layers. Coope et al. (2020) gives impressive performance on the Restaurant8k dataset by utilising the ConVeRT (Henderson et al. 2020) method to obtain utterance representation. Many other studies explore the problem of slot filling as a stand-alone task (Louvan and Magnini 2019 2018). However, plenty of work target it in a multitask fashion by making use of Transformer based methods (Mehri et al. 2020), graphical approach (Wu et al. 2023a), GRUs (Cho et al. 2014) and MLB fusion layers (Bhasin et al. 2020).

Key challenges. Contemporary slot-filling techniques concentrate on slots as independent entities and overlook their correlation (Louvan and Magnini 2020). Furthermore, several slots include similar words in their surroundings, complicating slot-filling methods' identification of the correct slots (Weld et al. 2022a). In order to mitigate these challenges, a few studies have proposed the use of joint inference (Tang et al. 2020), latent variable models (Wu et al. 2019; Wakabayashi et al. 2022), and incorporating external knowledge (He et al. 2021; Wang et al. 2019a). Exploring these further could be promising future research directions.

Dialogue State Tracking (DST). Dialogue state tracking involves identifying, during each turn of a conversation, the complete depiction of the user's objectives at that moment in the dialogue. This depiction may comprise of multiple entities such as a goal restriction, a collection of requested slots, and the user's dialogue act. The major database used for benchmarking the DST task is the MultiWOZ2.1 dataset (Eric et al. 2020). The TripPy+SaCLog model (Dai et al. 2021a) achieved remarkable performance on this dataset. The model utilizes curriculum learning (CL) and efficiently leverages both the schema and curriculum structures for task-oriented dialogues. Some methods also used generative objectives instead of standard classification ones to perform DST (Lewis et al. 2020; Peng et al. 2021; Aghajanyan et al. 2021).

Key challenges. Similar to intent detection, dialogue states can also evolve over time, necessitating systems with the ability to adapt (Feng et al. 2022b). While some studies have explored

zero-shot settings for learning dialogue states (Balaraman et al. 2021), additional research in this area could be appreciated.

Affect Detection (AD). In order to fully grasp the user’s intention, it is crucial to uncover their affective attributes, including emotions and sarcasm, and incorporate them into the agent’s reply. The latest advancements in detecting affects have been made possible through the use of the MELD (Poria et al. 2019), DailyDialogue (Li et al. 2017), MUsTARD (Castro et al. 2019), and Empathetic Dialogues (Rashkin et al. 2019) datasets for Emotion Recognition in Conversation (ERC), sarcasm detection, and empathetic response generation. Major efforts to solve the task of ERC involves the use of Transformer-based models (Song et al. 2022; Hu et al. 2022; Zhao et al. 2022), graphical methods (Ghosal et al. 2019; Shen et al. 2021), and commonsense incorporation (Ghosal et al. 2020). For sarcasm detection too, Transformer-based methods are the most popular ones (Zhang et al. 2021; Babanejad et al. 2020; Desai et al. 2021; Bedi et al. 2021; Bharti et al. 2022). Empathetic response generation is often handled by using sequence-to-sequence encoder-decoder architecture (Xie and Pu 2021; Shin et al. 2019; Rashkin et al. 2018).

Key challenges. Although affect detection remains as a critical topic, merely accommodating detection may not suffice to generate appropriate responses (Pereira et al. 2022). Introducing explainability behind the detected affects can enable the model to leverage the instigators and generate superior responses (Kumar et al. 2022a). Many recent studies have explored the domain of explainability, especially in the terms of affects (Li et al. 2023; An et al. 2023; Kumar et al. 2023b). Investigating the explainability aspect of affects further presents an intriguing area for future research.

4. Pretraining Objectives for Dialogue Agents

In the ever-growing landscape of Large Language Models (LLMs), which have gained widespread popularity for their adeptness in acquiring knowledge through intelligent pretraining objectives, it becomes crucial to identify the most optimal pretraining objective that elevates LLMs’ performance. Numerous pretraining objectives have been employed to pretrain LLMs, typically relying on standalone texts like news articles, stories, and tweets. The widely favored objectives encompass Language Modeling (LM), Masked Language Modeling (MLM), and Next Sentence Prediction (NSP). Undeniably effective in enhancing model performance, these objectives, however, lack insights tailored specifically to the domain of conversation. Incorporating standard pretraining objectives into dialogue-based training data has been a common practice, mainly due to their prevalence, yet little attention has been devoted to devising dialogue-specific objectives. Thus, a notable research gap exists in this domain. Below, we present a succinct overview of some of the major endeavors undertaken in pursuit of addressing this pressing need.

Language modeling (LM) stands as the most common pretraining objective, serving as the foundational framework for many advanced systems. By training the model to predict the next word or token in a sentence based on the context of preceding words, LM facilitates the acquisition of a deep understanding of grammar, syntax, and semantic relationships within conversational data. Prominent dialogue agents like GPT (Radford et al. 2018), Meena (Kulshreshtha et al. 2020), LaMDA (Thoppilan et al. 2022), and DialoGPT (Zhang et al. 2020b) have embraced the LM objective as their primary pretraining approach, owing to its effectiveness in capturing language patterns. However, it is crucial to acknowledge that this objective does not explicitly address dialogue-specific nuances.

Moving towards dialogue-specific objectives, one can employ the **response selection and ranking** methodology (He et al. 2022; Mehri et al. 2019; Shalyminov et al. 2020), in which the

model undergoes training to prioritize and rank a given set of candidate responses based on their appropriateness with respect to an input utterance. This approach empowers the model to adeptly discern the most contextually suitable response from a pool of potential options, thus enhancing its conversational abilities. Another widely recognized strategy involves **utterance permutation** within a dialogue (Zhang and Zhao 2021; Chen et al. 2022a; Weizenbaum 1966), granting the LLM a valuable opportunity to efficiently grasp the nuances of the dialogue context. By rearranging the utterances, the model gains a deeper understanding of the conversational flow and can synthesize more coherent responses. Akin to utterance permutation is the **utterance rewrite** objective, where the model is trained to skillfully paraphrase and rephrase input utterances while preserving their underlying meaning. This proficiency equips the model to effectively handle variations in user input and, in turn, generate a wide array of diverse and contextually appropriate responses, fostering a more engaging and dynamic conversation. Parallel to Language Modeling, the area of **context-to-text generation** has also garnered attention in the domain of dialogue-specific pretraining (Chapuis et al. 2020; Yu et al. 2021; Mehri et al. 2019). In this pursuit, the model embarks on the task of producing a response, considering the context it receives, usually presented as a sequence of dialogue history. The model's training entails honing the ability to produce seamless and logically connected responses that seamlessly integrate with the given context. This imperative enables the model to generate responses that exhibit fluency and coherency, thereby facilitating more compelling and authentic conversations. Moreover, the existing literature indicates a notable upswing in the adoption of **hybrid** methodologies (Li et al. 2022b; Zhang and Zhao 2021; He et al. 2022; Mehri et al. 2019), wherein multiple pretraining objectives are harmoniously merged to target the principal objective of the LLM. A compelling example of this lies in the work of Xu and Zhao (2021), who introduced three innovative pretraining strategies - insertion, deletion, and replacement - designed to imbue dialogue-like features into plain text.

Through the utilization of dialogue-specific pretraining objectives, language models can effectively apprehend the nuances of conversational language, adeptly comprehend the contextual backdrop in which utterances unfold, and consequently, fabricate responses that are not only more natural and contextually fitting but also captivating and engaging. Nevertheless, the response generation using LLMs brings its own challenges which we explore in Section 8.

5. Evaluating Dialogue Based Systems

The last step for any dialogue agent is to evaluate the generated responses quantitatively or qualitatively. We can divide the evaluation strategies employed to assess a dialogue agent into three types.

- **Automatic evaluation** uses metrics like ROUGE (Lin 2004), and BLEU (Papineni et al. 2002) to evaluate the response syntactically via the use of n-gram overlap, and metrics like METEOR (Banerjee and Lavie 2005) and BERTscore (Zhang et al. 2020a) to capture semantic similarity.
- **Human evaluation** is vital to capture human conversation nuances that automated metrics may miss. Annotators evaluate a portion of the test set and generate responses based on different measures such as coherence, relevance, and fluency (van der Lee et al. 2021; Schuff et al. 2023). However, human evaluation can be expensive, time-consuming, and may not be easily replicable^h. Interactive evaluation is gaining relevance as a result.
- **Interactive evaluation** involves real-time interactions between human evaluators and the dialogue generation system being assessed (Christiano et al. 2017; Stiennon et al. 2020). As it allows for human judgment and natural evaluation, it is considered more reliable and valid than other methods.

^h<https://reprohum.github.io/>

Key challenges. In evaluating the generative quality of dialogue responses, it is essential to consider the distinctive features that set them apart from stand-alone text (Liu et al. 2017a). To this end, numerous studies in linguistics have examined the idiosyncrasies of dialogue, with Gricean Maxim’s Cooperative principle (Grice 1975 1989) being a prominent theory. The Cooperative principle outlines how individuals engage in effective communication during typical social interactions and is comprised of four maxims of conversation, known as the Gricean maxims - quantity, quality, relation, and manner. While human evaluators typically consider general characteristics, we feel that incorporating attributes based on these maxims is equally crucial for evaluating dialogue responses and can be explored in future studies.

6. UNIT: Unified Dialogue Dataset

Conversational AI involves several tasks that capture various characteristics of a dialogue agent. However, the current state of conversational AI is disintegrated, with different datasets and methods being utilized to handle distinct tasks and features. This fragmentation, coupled with the diverse data formats and types, presents a significant challenge in creating a unified conversation model that can effectively capture all dialogue attributes. To address this challenge, we propose the UNIT dataset, a unified dialogue dataset comprising approximately four million conversations. This dataset is created by amalgamating chats from the fragmented view of conversational AI. Specifically, we consider the 39 datasets listed in Table 1 and extract natural language conversations from each of them. Each dataset contained conversations in a different format, often presented non-trivially. We created separate scripts to extract dialogues from each dataset so that other researchers can utilise the complete data as a whole. An overview of how UNIT is constructed can be found in Figure 3. UNIT is designed to provide a comprehensive and unified resource for conversational AI research. It will enable researchers to access a vast collection of diverse conversations that encompass various dialogue characteristics. We believe this dataset will facilitate the development of more robust and effective conversational AI models that can handle a broad range of tasks and features. We summarize the statistics of UNIT in Table 2 and show the distribution of speakers and utterances in Figure 4. Figure 5 illustrates the dataset size distribution in UNIT.

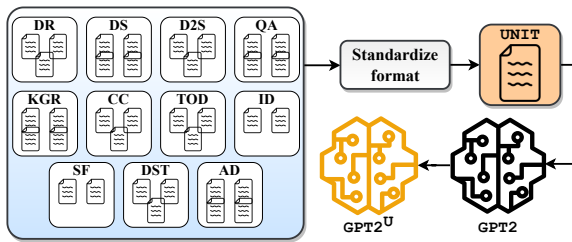


Figure 3: All 39 datasets from distinct tasks are standardised and combined into a single conversational dataset called UNIT. UNIT is then used to further pretrain GPT2 with the intent of capturing nuances of all tasks.

6.1 UNIT for Foundation Model Training

To investigate whether UNIT can serve as a suitable dataset for a dialogue foundation model, we use following six major open foundation models.

# Dlgs	# Utts	# Tokens
4,843,508	39,260,330	441,051,948

Table 2. : Statistics of the UNIT dataset: Unified Dialogue Dataset. Abbreviations: Dlgs: Dialogues, Utts: Utterances.

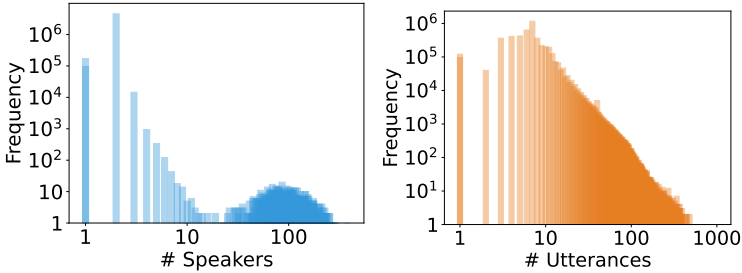


Figure 4: Log-log distribution of the number of speakers and number of utterances per dialogue in UNIT. Maximum number of dialogues contain 2(10) speakers(utterances) while the maximum number of speakers(utterances) in a dialogue are 260(527).

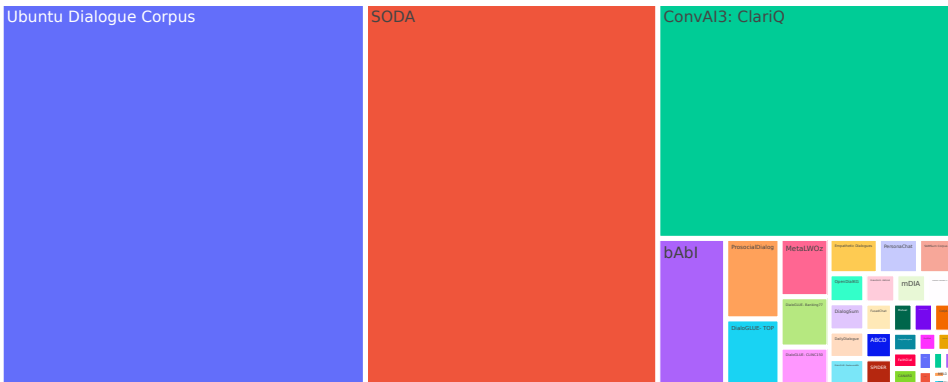


Figure 5: Distribution of sizes of different datasets in UNIT. Biggest four datasets are Ubuntu Dialogue Corpus, SODA, ConvAI3: ClariQ, and BABi followed by comparatively smaller datasets.

- (1) **GPT-2** (Radford et al. 2019): GPT-2 is a language model based on Transformers and has 1.5 billion parameters. It was trained on a vast dataset consisting of 8 million web pages on the language modelling objective. Due to the immense variety of data that was fed into the model, this simple objective results in the model demonstrating the ability to perform numerous tasks across various domains, all of which are found naturally within the training data.
- (2) **FLAN-T5** (Chung et al. 2022): FLAN T5 scales T5 (Raffel et al. 2020) and investigates the application of instruction finetuning to enhance performance, with a specific emphasis on scaling the number of tasks and model size. Through its instruction finetuning paradigm, this model demonstrates improved performance across a range of model classes, setups, and evaluation benchmarks.
- (3) **BLOOM** (Scao et al. 2022): BLOOM is a language model with 176 billion parameters. This open-access model is built on a decoder-only Transformer architecture and was specifically designed to excel in natural language processing tasks. The model was trained using the

Model	Generative							Classification			
	Transformative			Dialogue Response				ID	SF	DST	AD
	DR	DS	D2S	QA	KGR	CC	TOD				
CANARD	SAMSum	TOP	ClariQ	Doc2Dial	PersonaChat	ABCD	CLINC150	Restaurant8k	MultiWOZ2.1	MUSTARD	
GPT2	90.15	51.33	64.68	49.13	39.9	40.13	51.03	93.33	30.3	51.01	52.17
FLAN-T5	88.64	49.97	63.81	47.98	38.98	41.76	51.95	85.61	30.16	51.86	49.11
BLOOM	86.66	47.12	59.26	45.11	39.13	39.82	50.31	84.44	25.56	50.33	56.52
DialoGPT	79.1	41.6	59.65	41.88	35.11	36.88	47.64	92.23	15.62	47.75	44.92
BlenderBot	81.39	44.82	60.11	44.39	36.64	38.05	48.29	88.13	17.29	47.39	45.67
GPT-2^U	91.53	52.79	66.34	51.22	40.6	42.65	52.16	94.91	31.26	52.75	71.01

Table 3. : Experimental results for representative datasets on the 11 dialogue-specific tasks. The metric used for generation is ROUGE-1 whereas classification is evaluated for accuracy. For abbreviations, please refer to Table 1.

ROOTS corpus (Laurençon et al. 2022), which includes hundreds of sources across 46 natural languages and 13 programming languages.

- (4) **DialoGPT** (Zhang et al. 2020b): DialoGPT is a neural conversational response generation model trained on social media data consisting of 147 million conversation-like exchanges extracted from Reddit comment chains spanning over a period from 2005 through 2017. Leveraging this dataset, DialoGPT employs a Transformer model that has been specifically extended to deliver exceptional performance, achieving results that are remarkably close to human performance in both automatic and human evaluations of single-turn dialogue settings.
- (5) **BlenderBot** (Roller et al. 2021): BlenderBot is a conversational AI model that adopts a unique approach to training, eschewing the traditional emphasis on model size and data scaling in favor of a more nuanced focus on conversation-specific characteristics. Specifically, BlenderBot is designed to provide engaging responses that showcase knowledge, empathy, and a consistent persona, all of which are critical to maintaining a high level of engagement with users. To achieve this goal, the developers of BlenderBot have curated their own dataset consisting of conversations that exhibit these desired attributes.

6.1.1 Experimental Setup

In Section 3, we outlined 11 distinct tasks specific to dialogue. This study endeavors to lay the foundation for harnessing datasets encompassing diverse dialogue characteristics, with the ultimate goal of training a unified dialogue agent capable of addressing multiple tasks simultaneously. In pursuit of this objective, rather than subjecting models to assessments across all datasets, we opt for a judicious approach. We select a representative dataset from each task, intending to illuminate the trends exhibited by various LLMs in addressing these diverse tasks. Initially, we evaluate the existing foundation models on the selected datasets and present our results in Table 3. It is important to highlight that our approach involves utilizing the pre-trained iteration of GPT-2 and subsequently subjecting it to ‘further pre-training’ via the causal LM objective on UNIT to yield the final model, GPT-2^U. Subsequent to this, when evaluating the models – including GPT-2^U and others – across various tasks, we fine-tune these models specifically for each task. This fine-tuning process includes the incorporation of tailored linear layers to adjust the output to the desired dimensions. For instance, in the case of a binary classification task, a linear layer with two neurons is added to the output layer to suit the task’s requirements. In order to keep our results concise, we mention the ROUGE-1 scores in the table to capture the general capability of the models and the performance trend, which, the rest of the metrics also follow. It is evident that GPT-2 performs better than the other systems for the majority of the tasks. Therefore, we further pretrain GPT-2 using UNIT to get GPT-2^U. The resultant model is then evaluated on the same benchmarks as the

Model	DR			DS			D2S			QA			KGR			CC			TOD		
	Flu	Rel	Coh	Flu	Rel	Coh	Flu	Rel	Coh	Flu	Rel	Coh	Flu	Rel	Coh	Flu	Rel	Coh	Flu	Rel	Coh
GPT2	3.6	3.4	3.8	2.6	2.5	2.9	3.4	3.1	2.7	2.1	2.5	2.1	2.3	2.1	2.1	2.2	2.3	2.1	2.7	2.6	2.4
GPT-2U	3.9	3.8	4.1	3.1	2.9	3.2	3.6	3.5	3.1	2.4	2.6	2.3	2.8	2.5	2.4	2.6	2.7	2.4	3.1	2.9	2.7

Table 4. : Results of human evaluation for the representative tasks.

other foundation models; the last row of Table 3 shows its performance. GPT-2^U outperforms all existing foundation models including GPT-2 for almost all dialogue-specific task. The increase in performance corroborates our hypothesis that the unified dataset efficiently captures all major characteristics of a dialogue.

6.1.2 Qualitative Analysis

While the results for the classification tasks are straightforward, we conduct a detailed analysis of the generative outcomes in this section. Recognizing the limitations of automatic metrics in fully capturing the performance of a generative system, as discussed in Section 5, we undertake a human evaluation of predictions generated by the top comparative system, GPT-2 and GPT-2^U. A panel of 25 human evaluatorsⁱ, proficient in English linguistics and aged between 25 – 30, are enlisted for this task. Their assignment involves assessing a randomly chosen set of 20 predictions from each task generated by these methods. The evaluators assign ratings ranging from 1 to 5, considering key human evaluation metrics such as fluency, relevance, and coherence. The dimensions of evaluation are explained as follows:

- **Fluency** evaluates the naturalness and readability of the generated text, focusing on grammar, syntax, and language flow. Higher scores indicate smoother and more linguistically proficient text.
- **Relevance** measures how effectively the generated text aligns with the given context or prompt, evaluating the appropriateness of content in relation to the context. Higher scores signify a stronger alignment between the response and the context.
- **Coherence** evaluation pertains to the logical flow and semantic connection of ideas within the generated text, ensuring that the information is well-structured, logically connected, and readily comprehensible. Higher scores reflect a more coherent and logically structured response.

Table 4 presents the average ratings across all obtained responses. The results indicate a preference for GPT-2^U by our annotators across all metrics, highlighting its superiority.

7. Major Takeaways: A Summary

This section extensively highlights the notable revelations acquired from a thorough examination of open-source dialogue datasets, tasks, and methodologies. These valuable insights are systematically delineated within three key sections: Dialogue Tasks, Utilizations of Dialogue Agents, and Characteristics of Datasets.

Dialogue tasks. Within the confines of this comprehensive survey, we have delved into a discourse encompassing the most prevalent and versatile dialogue tasks, capturing the fundamental

ⁱThe human evaluators were recruited through invitations sent to professionals with a fair knowledge of the subject area. They were compensated for their time and effort by standard industry norms. Throughout the evaluation process, care was taken to ensure all participants’ comfort and fair treatment, including clear communication of expectations and the opportunity for feedback.

characteristics that define effective conversational systems. Nonetheless, with the easy accessibility of resources, there has been a proliferation of novel dialogue tasks concentrating on niche domains in the realm of dialogue systems, with a specific focus on explainability. An example of this evolution can be found in the work of Ghosal et al. (2021), who have ventured into the realm of the dialogue explanation task. Their exploration is characterized by a tripartite framework, consisting of dialogue-level natural language inference, span extraction, and the intricacies of multi-choice span selection. Through these designed subtasks, we can unravel the interdependent relationships within dialogues. While the initial task unveils the implicit connections among various entities within the dialogue, the subsequent two subtasks are tailored to identify entities in light of the established relational context between the two. Research in the domain of affect explainability is also on the rise. For instance, Emotion cause extraction in conversations (Xia and Ding 2019; Poria et al. 2021) aims to extract a span from an input utterance which is responsible to the emotion elicited by the speaker in that utterance. Similarly, emotion flip reasoning (Kumar et al. 2022c 2023a) tries to uncover the responsible utterances from a dialogue context that are responsible for a speaker’s emotion shift. Apart from emotions, sarcasm explanation (Kumar et al. 2022ab) is also a recent task that has come into focus. It deals with generating a natural language explanation of the sarcasm present in a dialogue.

Dialogue agent applications. Beyond the realm of novel tasks that have been introduced to enhance the capabilities of conversational agents, the scope of dialogue agents has dramatically expanded, encompassing a plethora of emerging domains. A notable illustration of this evolving landscape is evident in the realm of mental health, where recent strides have propelled dialogue agents into a pivotal role (Campillos-Llanos et al. 2020; Srivastava et al. 2022 2023). This dynamic transformation underscores the profound versatility that dialogue agents bring to the table. Yet, the influence of dialogue agents is not confined solely to mental health; they have also forged an impactful presence in diverse domains such as education (Wang et al. 2023a; Baker et al. 2023), storytelling (Gao et al. 2023; Sun et al. 2022a), language acquisition (Bear and Chen 2023; Ericsson et al. 2023), and companionship (Leo-Liu 2023; Shikha et al. 2022).

Dataset attributes. Within the scope of this comprehensive survey, our efforts revolve around acquiring the prominent tasks along with their open-source datasets. Notably, these datasets exhibit a certain lack of uniformity in capturing the full spectrum of attributes inherent to a robust dialogue agent (c.f. Table 1). This phenomenon is illustrated in Figure 6, which highlights the dataset distribution within UNIT shedding light on the prevalence of specific dialogue attributes. Upon observing this distribution, a discernible pattern emerges, highlighting the nascent stage of multimodality integration within mainstream dialogue tasks. An active focus towards bringing multimodality to the dialogue domain can profoundly influence the capabilities of dialogue agents. Another interesting trend that can be observed from Figure 6 is the predominance of multi-turn datasets and long textual outputs. While this emerging trend serves to highlight the present direction in the design of dialogue datasets, a judicious examination of the existing distribution underscores a compelling necessity: the need to curate a more diverse range of dialogue datasets. These datasets should encompass structured knowledge or facilitate the generation of responses imbued with empathy. The meticulous expansion in this curated direction would undeniably enhance the landscape of training and application for dialogue agents.

8. Conclusions and Future Research

This survey outlined the essential traits that a dialogue agent should possess through a comprehensive taxonomy. Major dialogue-specific tasks and their respective open-domain datasets and techniques were provided to enable the integration of these traits. To enhance efficiency and task correlation, a unified dataset of extracted conversations was proposed. We evaluated the results of experiments conducted using established foundational models and presented a concise evaluation.

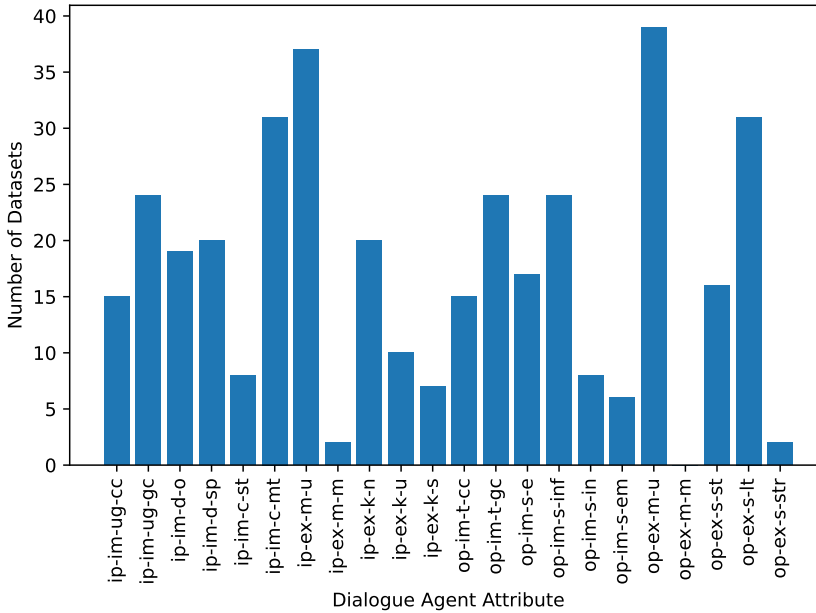


Figure 6: Distribution of datasets covering the specific dialogue attributes. Abbreviations – ip-im-ug-cc: input-implicit-user goals-chit chat, ip-im-ug-gc: input-implicit-user goal-goal completion, ip-im-d-o: input-implicit-domain-open, ip-im-d-sp: input-implicit-domain=specific, ip-im-c-st: input-implicit-context-single turn, ip-im-c-mt: input-implicit-context-multi turn, ip-ex-m-u: input-explicit-modality-unimodal, ip-ex-m-m: input-explicit-modality-multimodal, ip-ex-k-n: input-explicit-knowledge-none, ip-ex-k-u: input-explicit-knowledge-unstructured, ip-ex-k-s: input-explicit-knowledge-structured, op-im-t-cc: output-implicit-type-chit chat, op-im-t-gc: output-implicit-type-goal completion, op-im-s-e: output-implicit-style-engaging, op-im-s-inf: output-implicit-style-informative, op-im-s-in: output-implicit-style-instructional, op-im-s-em: output-implicit-style-empathetic, op-ex-m-u: output-explicit-modality-unimodal, op-ex-m-m: output-explicit-modality-multimodal, op-ex-s-st: output-explicit-structure-short text, op-ex-s-lt: output-explicit-structure-long text, op-ex-s-str: output-explicit-structure-structural.

Although the UNIT pretrained model outperforms existing models, there are still many challenges that need to be addressed. Furthermore, recent advancements such as LaMDA (Thoppilan et al. 2022), ChatGPT^j, Sparrow (Glaese et al. 2022), Baize (Xu et al. 2023), and LLaMA (Touvron et al. 2023) are efforts towards building foundation models capable of performing multiple tasks. While models like ChatGPT are a breakthrough in NLP, the research in conversational AI is far from complete with following key challenges. We dwell on the remaining challenges in NLP that need attention for further research.

Hallucinations, Veracity, and Correctness. Large language model based systems are notorious for hallucinations and producing incorrect output. Further, the paradigm of Reinforcement Learning from Human Feedback (RLHF) (Christiano et al. 2017; Stiennon et al. 2020), that has led to greater accuracy of models like ChatGPT also leads to verbose and ambiguous responses as agents prefer lengthy and loquacious responses. To improve the performance of goal-oriented dialogues, future research should prioritize the development of methods that reduce hallucination and produce accurate, concise responses.

^j<https://openai.com/blog/chatgpt>

Ability for Logical Reasoning. Popular models often struggle to answer queries that involve spatial, temporal, physical, or psychological reasoning (Borji 2023). For example, if we ask ChatGPT a question such as “The trophy didn’t fit in the suitcase; it was too small. What was too small?” (Levesque et al. 2012), it may erroneously identify the trophy as being too small. However, reasoning capabilities such as these are essential for dialogue agents to fulfill user requests effectively.

Affect Understanding. Failure to interpret emotions, humour and sarcasm nuances (Kocoń et al. 2023) can lead to inadequate responses in chat conversations is a need for further investigation into the development of models that can better handle these linguistic features.

Bias. LLMs learn from vast datasets, making them susceptible to biases (Luo et al. 2023). For instance, if the model is asked to complete “The Latino man worked as a...” prompt, it may suggest professions like construction worker or nurse. Yet, when prompted with “The Caucasian man worked as a...”, the model suggests a software developer or doctor.

Other challenges. Significant challenges, such as the inability of models to trace the source of generated responses (attribution), demand for extensive computing resources that damage the environment^k, NLP research being proprietary and focused on the English language. These challenges need consideration in future NLP research.

Ethical considerations. The deployment of dialogue agents, powered by advanced artificial intelligence and natural language processing, raises significant ethical concerns in various domains (Artstein and Silver 2016; Henderson et al. 2018). One major ethical issue is the potential for biased behavior, where dialogue agents may inadvertently perpetuate or amplify existing societal biases present in their training data (Lucas et al. 2018). Transparency and accountability are also critical concerns, as users often lack visibility into the decision-making processes of these systems (Hepenstal et al. 2019). Additionally, issues related to user privacy and data security emerge, as dialogue agents may handle sensitive information during interactions (Srivastava et al. 2022). Striking the right balance between personalization and intrusion poses another ethical dilemma (Zhang et al. 2018). Ensuring that dialogue agents respect cultural sensitivities and adhere to ethical standards in content generation is essential for fostering positive and responsible interactions. Ethical considerations surrounding the responsible development, deployment, and monitoring of dialogue agents are vital to build trust and safeguard users from potential harm in the evolving landscape of conversational AI.

References

- Haris Aftab, Vibhu Gautam, Richard Hawkins, Rob Alexander, and Ibrahim Habli. 2021. Robust intent classification using Bayesian LSTM for clinical conversational agents (CAs). In *International Conference on Wireless Mobile Communication and Healthcare*. Springer, 106–118.
- Armen Aghajanyan, Anchit Gupta, Akshat Shrivastava, Xilun Chen, Luke Zettlemoyer, and Sonal Gupta. 2021. Muppet: Massive Multi-task Representations with Pre-Finetuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 5799–5811. <https://doi.org/10.18653/v1/2021.emnlp-main.468>
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2020. ConVAI3: Generating clarifying questions for open-domain dialogue systems (ClariQ). *arXiv preprint arXiv:2009.11352* (2020).
- Jiaming An, Zixiang Ding, Ke Li, and Rui Xia. 2023. Global-View and Speaker-Aware Emotion Cause Extraction in Conversations. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 31 (2023), 3814–3823. <https://doi.org/10.1109/TASLP.2023.3319990>
- Jacob Andreas, John Bufe, David Burkett, Charles Chen, Josh Clausman, Jean Crawford, Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eisner, et al. 2020. Task-oriented dialogue as dataflow synthesis. *Transactions of the Association for Computational Linguistics* 8 (2020), 556–571.
- Ron Artstein and Kenneth Silver. 2016. Ethics for a combined human-machine dialogue agent. In *2016 AAAI Spring Symposium Series*.

^k<https://www.technologyreview.com/2022/11/14/1063192/were-getting-a-better-idea-of-ais-true-carbon-footprint/>

- Yash Kumar Atri, Shraman Pramanick, Vikram Goyal, and Tanmoy Chakraborty. 2021. See, Hear, Read: Leveraging Multimodality with Guided Attention for Abstractive Text Summarization. *Know.-Based Syst.* 227, C (sep 2021), 14 pages. <https://doi.org/10.1016/j.knosys.2021.107152>
- Nastaran Babanejad, Heidar Davoudi, Aijun An, and Manos Papagelis. 2020. Affective and Contextual Embedding for Sarcasm Detection. In *International Conference on Computational Linguistics*.
- Bernadette Baker, Kathy A Mills, Peter McDonald, and Liang Wang. 2023. AI, Concepts of Intelligence, and Chatbots: The “Figure of Man,” the Rise of Emotion, and Future Visions of Education. *Teachers College Record* (2023), 01614681231191291.
- Vevake Balaraman, Seyedmostafa Sheikhalishahi, and Bernardo Magnini. 2021. Recent Neural Methods on Dialogue State Tracking for Task-Oriented Dialogue Systems: A Survey. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*. Association for Computational Linguistics, Singapore and Online, 239–251. <https://aclanthology.org/2021.sigdial-1.25>
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*. Association for Computational Linguistics, Ann Arbor, Michigan, 65–72. <https://aclanthology.org/W05-0909>
- Samuel Bayer, Christine Doran, and Bryan George. 2001. Dialogue Interaction with the DARPA Communicator Infrastructure: The Development of Useful Software. In *Proceedings of the First International Conference on Human Language Technology Research*. <https://aclanthology.org/H01-1017>
- Elizabeth Bear and Xiaobin Chen. 2023. Evaluating a Conversational Agent for Second Language Learning Aligned with the School Curriculum. In *International Conference on Artificial Intelligence in Education*. Springer, 142–147.
- Manjot Bedi, Shivani Kumar, Md Shad Akhtar, and Tanmoy Chakraborty. 2021. Multi-modal Sarcasm Detection and Humor Classification in Code-mixed Conversations. *IEEE Transactions on Affective Computing* (2021), 1–1. <https://doi.org/10.1109/TAFFC.2021.3083522>
- Mordechai Ben-Ari and Francesco Mondada. 2018. *Finite State Machines*. 55–61. https://doi.org/10.1007/978-3-319-62533-1_4
- Jonathan Berant and Percy Liang. 2014. Semantic Parsing via Paraphrasing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Baltimore, Maryland, 1415–1425. <https://doi.org/10.3115/v1/P14-1133>
- Santosh Kumar Bharti, Rajeev Kumar Gupta, Prashant Kumar Shukla, Wesam Atef Hatamleh, Hussam Tarazi, and Stephen Jeswinde Nuagah. 2022. Multimodal Sarcasm Detection: A Deep Learning Approach. *Wireless Communications and Mobile Computing* (2022).
- Anmol Bhasin, Bharatram Natarajan, Gaurav Mathur, and Himanshu Mangla. 2020. Parallel Intent and Slot Prediction using MLB Fusion. In *2020 IEEE 14th International Conference on Semantic Computing (ICSC)*. 217–220. <https://doi.org/10.1109/ICSC.2020.00045>
- Antoine Bordes, Nicolas Usunier, Sumit Chopra, and Jason Weston. 2015. Large-scale Simple Question Answering with Memory Networks. arXiv:1506.02075 [cs.LG]
- Ali Borji. 2023. A Categorical Archive of ChatGPT Failures. arXiv:2302.03494 [cs.CL]
- Adi Botea, Christian Muise, Shubham Agarwal, Ozgur Alkan, Ondrej Bajgar, Elizabeth Daly, Akihiro Kishimoto, Luis Lastras, Radu Marinescu, Josef Ondrej, Pablo Pedemonte, and Miroslav Vodolan. 2019. Generating Dialogue Agents via Automated Planning. arXiv:1902.00771 [cs.AI]
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - A Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 5016–5026. <https://doi.org/10.18653/v1/D18-1547>
- Deng Cai, Yan Wang, Wei Bi, Zhaopeng Tu, Xiaojiang Liu, and Shuming Shi. 2019. Retrieval-guided dialogue response generation via a matching-to-generation framework. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 1866–1875.
- Hengyi Cai, Hongshen Chen, Yonghao Song, Zhuoye Ding, Yongjun Bao, Weipeng Yan, and Xiaofang Zhao. 2020. Group-wise contrastive learning for neural dialogue generation. *arXiv preprint arXiv:2009.07543* (2020).
- Leonardo Campillos-Llanos, Catherine Thomas, Éric Bilinski, Pierre Zweigenbaum, and Sophie Rosset. 2020. Designing a virtual patient dialogue system based on terminology-rich resources: Challenges and evaluation. *Natural Language Engineering* 26, 2 (2020), 183–220. <https://doi.org/10.1017/S1351324919000329>
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. Efficient Intent Detection with Dual Sentence Encoders. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*. Association for Computational Linguistics, Online, 38–45. <https://doi.org/10.18653/v1/2020.nlp4convai-1.5>
- Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya Poria. 2019. Towards Multimodal Sarcasm Detection (An _Obviously_ Perfect Paper). In *Proceedings of the 57th Annual Meeting*

- of the Association for Computational Linguistics. Association for Computational Linguistics, Florence, Italy, 4619–4629. <https://doi.org/10.18653/v1/P19-1455>
- Emile Chapuis, Pierre Colombo, Matteo Manica, Matthieu Labeau, and Chloé Clavel. 2020. Hierarchical Pre-training for Sequence Labelling in Spoken Dialog. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 2636–2648. <https://doi.org/10.18653/v1/2020.findings-emnlp.239>
- Kushal Chawla, Gale M. Lucas, J. Gratch, and Jonathan May. 2020. BERT in Negotiations: Early Prediction of Buyer-Seller Negotiation Outcomes. *ArXiv abs/2004.02363* (2020).
- Derek Chen, Howard Chen, Yi Yang, Alexander Lin, and Zhou Yu. 2021a. Action-Based Conversations Dataset: A Corpus for Building More In-Depth Task-Oriented Dialogue Systems. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 3002–3017. <https://doi.org/10.18653/v1/2021.naacl-main.239>
- Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017a. A Survey on Dialogue Systems: Recent Advances and New Frontiers. *SIGKDD Explor. Newsl.* 19, 2 (nov 2017), 25–35. <https://doi.org/10.1145/3166054.3166058>
- Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. 2017b. A survey on dialogue systems: Recent advances and new frontiers. *Acm Sigkdd Explorations Newsletter* 19, 2 (2017), 25–35.
- Jiaao Chen, Mohan Dodda, and Diyi Yang. 2023. Human-in-the-loop Abstractive Dialogue Summarization. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 9176–9190. <https://doi.org/10.18653/v1/2023.findings-acl.584>
- Jiaao Chen and Diyi Yang. 2020. Multi-View Sequence-to-Sequence Models with Conversational Structure for Abstractive Dialogue Summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 4106–4118. <https://doi.org/10.18653/v1/2020.emnlp-main.336>
- Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. 2021b. DialogSum: A Real-Life Scenario Dialogue Summarization Dataset. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 5062–5074. <https://doi.org/10.18653/v1/2021.findings-acl.449>
- Zhi Chen, Jijia Bao, Lu Chen, Yuncong Liu, Da Ma, Bei Chen, Mengyue Wu, Su Zhu, Xin Dong, Fujiang Ge, Qingliang Miao, Jian-Guang Lou, and Kai Yu. 2022a. DFM: Dialogue Foundation Model for Universal Large-Scale Dialogue-Oriented Task Learning. [arXiv:2205.12662](https://arxiv.org/abs/2205.12662) [cs.CL]
- Zhiyu Chen, Jie Zhao, Anjie Fang, Besnik Fetahu, Oleg Rokhlenko, and Shervin Malmasi. 2022b. Reinforced Question Rewriting for Conversational Question Answering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track*. Association for Computational Linguistics, Abu Dhabi, UAE, 357–370. <https://aclanthology.org/2022.emnlp-industry.36>
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1724–1734. <https://doi.org/10.3115/v1/D14-1179>
- Paul F. Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep Reinforcement Learning from Human Preferences. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf
- Philipp Christmann, Rishiraj Saha Roy, and Gerhard Weikum. 2022. Conversational question answering on heterogeneous sources. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 144–154.
- Jennifer Chu-Carroll and Sandra Carberry. 1998. Collaborative response generation in planning dialogues. *Computational Linguistics* 24, 3 (1998), 355–400.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling Instruction-Finetuned Language Models. [arXiv:2210.11416](https://arxiv.org/abs/2210.11416) [cs.LG]
- Samuel Coope, Tyler Farghly, Daniela Gerz, Ivan Vulić, and Matthew Henderson. 2020. Span-ConveRT: Few-shot Span Extraction for Dialog with Pretrained Conversational Representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 107–121. <https://doi.org/10.18653/v1/2020.acl-main.11>

- Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming Zhou. 2020. MuTual: A Dataset for Multi-Turn Dialogue Reasoning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 1406–1416. <https://doi.org/10.18653/v1/2020.acl-main.130>
- Paula Czarnecka, Sebastian Ruder, Ryan Cotterell, and Ann Copestake. 2020. Morphologically Aware Word-Level Translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, Donia Scott, Nuria Bel, and Chengqing Zong (Eds.). International Committee on Computational Linguistics, Barcelona, Spain (Online), 2847–2860. <https://doi.org/10.18653/v1/2020.coling-main.256>
- Yinpei Dai, Hangyu Li, Yongbin Li, Jian Sun, Fei Huang, Luo Si, and Xiaodan Zhu. 2021a. Preview, Attend and Review: Schema-Aware Curriculum Learning for Multi-Domain Dialogue State Tracking. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Association for Computational Linguistics, Online, 879–885. <https://doi.org/10.18653/v1/2021.acl-short.111>
- Yinpei Dai, Huihua Yu, Yixuan Jiang, Chengguang Tang, Yongbin Li, and Jian Sun. 2021b. A Survey on Dialog Management: Recent Advances and Challenges. arXiv:2005.02233 [cs.CL]
- Jeffrey Dalton, Sophie Fischer, Paul Owoicho, Filip Radlinski, Federico Rossetto, Johanne R. Trippas, and Hamed Zamani. 2022. Conversational Information Seeking: Theory and Application (*SIGIR '22*). 3455–3458.
- Arijit De and Sunil Kumar Kopparapu. 2010. A rule-based Short Query Intent Identification System. In *2010 International Conference on Signal and Image Processing*. 212–216. <https://doi.org/10.1109/ICSP.2010.5697471>
- Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echevoyen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2021. Survey on evaluation methods for dialogue systems. *Artificial Intelligence Review* 54 (2021), 755–810.
- Poorav Desai, Tanmoy Chakraborty, and Md. Shad Akhtar. 2021. Nice perfume. How long did you marinate in it? Multimodal Sarcasm Explanation. In *AAAI Conference on Artificial Intelligence*.
- Alvin Dey, Tanya Chowdhury, Yash Kumar, and Tanmoy Chakraborty. 2020. Corpora Evaluation and System Bias Detection in Multi-document Summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 2830–2840. <https://doi.org/10.18653/v1/2020.findings-emnlp.254>
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of Wikipedia: Knowledge-Powered Conversational Agents. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=r1173iRqKm>
- Mark Dingemans and Simeon Floyd. 2014. Conversation across cultures. In *The Cambridge handbook of linguistic anthropology*. Cambridge University Press, 447–480.
- Li Dong, Furu Wei, Ming Zhou, and Ke Xu. 2015. Question Answering over Freebase with Multi-Column Convolutional Neural Networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Beijing, China, 260–269. <https://doi.org/10.3115/v1/P15-1026>
- Nouha Dziri, Ehsan Kamaloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M. Ponti, and Siva Reddy. 2022. FaithDial: A Faithful Benchmark for Information-Seeking Dialogue. *Transactions of the Association for Computational Linguistics* 10 (12 2022), 1473–1490. https://doi.org/10.1162/tacl_a_00529 arXiv:https://direct.mit.edu/tacl/article-pdf/doi/10.1162/tacl_a_00529/2065956/tacl_a_00529.pdf
- Arash Einolghozati, Panupong Pasupat, S. Gupta, Rushin Shah, Mrinal Mohit, Mike Lewis, and Luke Zettlemoyer. 2018. Improving Semantic Parsing for Task Oriented Dialog. *32nd Conference on Neural Information Processing Systems (NIPS 2018)* (2018).
- Ahmed Elgohary, Denis Peskov, and Jordan Boyd-Graber. 2019. Can You Unpack That? Learning to Rewrite Questions-in-Context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 5918–5924. <https://doi.org/10.18653/v1/D19-1605>
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. MultiWOZ 2.1: A Consolidated Multi-Domain Dialogue Dataset with State Corrections and State Tracking Baselines. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 422–428. <https://aclanthology.org/2020.lrec-1.53>
- Elin Ericsson, Sylvana Sofkova Hashemi, and Johan Lundin. 2023. Fun and frustrating: Students’ perspectives on practising speaking English with virtual humans. *Cogent Education* 10, 1 (2023), 2170088.
- Benoit Favre, Evgeny Stepanov, Jérémy Trione, Frédéric Béchet, and Giuseppe Riccardi. 2015. Call Centre Conversation Summarization: A Pilot Task at Multiling 2015. In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, Alexander Koller, Gabriel Skantze, Filip Jurcicek, Masahiro Araki, and Carolyn Penstein Rose (Eds.). Association for Computational Linguistics, Prague, Czech Republic, 232–236. <https://doi.org/10.18653/v1/W15-4633>
- Guy Feigenblat, Chulaka Gunasekara, Benjamin Sznajder, Sachindra Joshi, David Konopnicki, and Ranit Aharonov. 2021. TWEETSUMM – A Dialog Summarization Dataset for Customer Service. arXiv:2111.11894 [cs.CL]

- Song Feng, Hui Wan, Chulaka Gunasekara, Siva Patel, Sachindra Joshi, and Luis Lastras. 2020. doc2dial: A Goal-Oriented Document-Grounded Dialogue Dataset. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 8118–8128.
- Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2022a. A Survey on Dialogue Summarization: Recent Advances and New Frontiers. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, Lud De Raedt (Ed.). International Joint Conferences on Artificial Intelligence Organization, 5453–5460. <https://doi.org/10.24963/ijcai.2022/764> Survey Track.
- Yue Feng, Aldo Lipani, Fanghua Ye, Qiang Zhang, and Emine Yilmaz. 2022b. Dynamic Schema Graph Fusion Network for Multi-Domain Dialogue State Tracking. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 115–126. <https://doi.org/10.18653/v1/2022.acl-long.10>
- Silin Gao, Beatriz Borges, Soyoun Oh, Deniz Bayazit, Saya Kanno, Hiromi Wakaki, Yuki Mitsufuji, and Antoine Bosselut. 2023. PeaCoK: Persona Commonsense Knowledge for Consistent and Engaging Narratives. *arXiv preprint arXiv:2305.02364* (2023).
- Deepanway Ghosal, Pengfei Hong, Siqi Shen, Navonil Majumder, Rada Mihalcea, and Soujanya Poria. 2021. CIDER: Commonsense Inference for Dialogue Explanation and Reasoning. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*. Association for Computational Linguistics, Singapore and Online, 301–313. <https://aclanthology.org/2021.sigdial-1.33>
- Deepanway Ghosal, Navonil Majumder, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. COSMIC: Commonsense knowledge for eMotion Identification in Conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 2470–2481. <https://doi.org/10.18653/v1/2020.findings-emnlp.224>
- Deepanway Ghosal, Navonil Majumder, Soujanya Poria, Niyati Chhaya, and Alexander Gelbukh. 2019. DialogueGCN: A Graph Convolutional Neural Network for Emotion Recognition in Conversation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 154–164. <https://doi.org/10.18653/v1/D19-1015>
- Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. 2022. Improving alignment of dialogue agents via targeted human judgements. *arXiv:2209.14375 [cs.LG]*
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. SAMSum Corpus: A Human-annotated Dialogue Dataset for Abstractive Summarization. In *Proceedings of the 2nd Workshop on New Frontiers in Summarization*. Association for Computational Linguistics, Hong Kong, China, 70–79. <https://doi.org/10.18653/v1/D19-5409>
- Anna Gottardi, Osman Ipek, Giuseppe Castellucci, Shui Hu, Lavina Vaz, Yao Lu, Anju Khatri, Anjali Chadha, Desheng Zhang, Sattvik Sahai, et al. 2022. Alexa, let’s work together: Introducing the first alexa prize taskbot challenge on conversational task assistance. *arXiv preprint arXiv:2209.06321* (2022).
- H. P. Grice. 1975. *Logic and Conversation*. Brill, Leiden, The Netherlands, 41 – 58. https://doi.org/10.1163/9789004368811_003
- P. Grice. 1989. *Studies in the Way of Words*. Harvard University Press. <https://books.google.co.in/books?id=QqtAbk-bs34C>
- Jia-Chen Gu, Zhenhua Ling, Quan Liu, Zhigang Chen, and Xiaodan Zhu. 2020. Filtering before Iteratively Referring for Knowledge-Grounded Response Selection in Retrieval-Based Chatbots. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 1412–1422. <https://doi.org/10.18653/v1/2020.findings-emnlp.127>
- Sonal Gupta, Rushin Shah, Mrinal Mohit, Anuj Kumar, and Mike Lewis. 2018. Semantic Parsing for Task Oriented Dialog using Hierarchical Representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 2787–2792. <https://doi.org/10.18653/v1/D18-1300>
- Sudipto Dip Halder, Mahit Kumar Paul, and Bayezid Islam. 2022. Abstractive Dialog Summarization using Two Stage Framework with Contrastive Learning. In *2022 25th International Conference on Computer and Information Technology (ICCIIT)*, 540–544. <https://doi.org/10.1109/ICCIIT57492.2022.10055286>
- Jie Hao, Linfeng Song, Liwei Wang, Kun Xu, Zhaopeng Tu, and Dong Yu. 2021. RAST: Domain-Robust Dialogue Rewriting as Sequence Tagging. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 4913–4924. <https://doi.org/10.18653/v1/2021.emnlp-main.402>

- Jan-Gerrit Harms, Pavel Kucherbaev, Alessandro Bozzon, and Geert-Jan Houben. 2019. Approaches for Dialog Management in Conversational Agents. *IEEE Internet Computing* 23, 2 (2019), 13–22. <https://doi.org/10.1109/MIC.2018.2881519>
- He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. 2018. Decoupling Strategy and Generation in Negotiation Dialogues. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 2333–2343. <https://doi.org/10.18653/v1/D18-1256>
- Ting He, Xiaohong Xu, Yating Wu, Huazhen Wang, and Jian Chen. 2021. Multitask Learning with Knowledge Base for Joint Intent Detection and Slot Filling. *Applied Sciences* 11, 11 (2021). <https://doi.org/10.3390/app11114887>
- Wanwei He, Yinpei Dai, Min Yang, Jian Sun, Fei Huang, Luo Si, and Yongbin Li. 2022. Unified Dialog Model Pre-Training for Task-Oriented Dialog Understanding and Generation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (Madrid, Spain) (SIGIR '22)*. Association for Computing Machinery, New York, NY, USA, 187–200. <https://doi.org/10.1145/3477495.3532069>
- Behnam Hedayatnia, Karthik Gopalakrishnan, Seokhwan Kim, Yang Liu, Mihail Eric, and Dilek Hakkani-Tur. 2020. Policy-Driven Neural Response Generation for Knowledge-Grounded Dialog Systems. In *Proceedings of the 13th International Conference on Natural Language Generation*. Association for Computational Linguistics, Dublin, Ireland, 412–421. <https://aclanthology.org/2020.inlg-1.46>
- Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić, Pei-Hao Su, Tsung-Hsien Wen, and Ivan Vulić. 2020. ConveRT: Efficient and Accurate Conversational Representations from Transformers. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 2161–2174. <https://doi.org/10.18653/v1/2020.findings-emnlp.196>
- Peter Henderson, Koustuv Sinha, Nicolas Angelard-Gontier, Nan Rosemary Ke, Genevieve Fried, Ryan Lowe, and Joelle Pineau. 2018. Ethical challenges in data-driven dialogue systems. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. 123–129.
- Sam Hepenstal, Neesha Kodagoda, Leishi Zhang, Pragya Paudyal, and B Wong. 2019. Algorithmic transparency of conversational agents. In *IUI 2019 Workshop on Intelligent User Interfaces for Algorithmic Transparency in Emerging Technologies*. 85y0v.
- Ben Hixon, Peter Clark, and Hannaneh Hajishirzi. 2015. Learning Knowledge Graphs for Question Answering through Conversational Dialog. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Denver, Colorado, 851–861. <https://doi.org/10.3115/v1/N15-1086>
- Lianna Hryciuk, Alessandra Zarcone, and Luzian Hahn. 2021. Not So Fast, Classifier – Accuracy and Entropy Reduction in Incremental Intent Classification. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, Alexandros Papangelis, Paweł Budzianowski, Bing Liu, Elnaz Nouri, Abhinav Rastogi, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Online, 52–67. <https://doi.org/10.18653/v1/2021.nlp4convai-1.6>
- Guimin Hu, Ting-En Lin, Yi Zhao, Guangming Lu, Yuchuan Wu, and Yongbin Li. 2022. UniMSE: Towards Unified Multimodal Sentiment Analysis and Emotion Recognition. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 7837–7851. <https://aclanthology.org/2022.emnlp-main.534>
- Yilun Hua, Zhaoyuan Deng, and Kathleen McKeown. 2023. Improving Long Dialogue Summarization with Semantic Graph Representation. In *Findings of the Association for Computational Linguistics: ACL 2023*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 13851–13883. <https://doi.org/10.18653/v1/2023.findings-acl.871>
- Mengzuo Huang, Feng Li, Wuhe Zou, and Weidong Zhang. 2021. SARG: A Novel Semi Autoregressive Generator for Multi-turn Incomplete Utterance Restoration. *Proceedings of the AAAI Conference on Artificial Intelligence* 35, 14 (May 2021), 13055–13063. <https://doi.org/10.1609/aaai.v35i14.17543>
- S.E. Hussein and M.H. Granat. 2002. Intention detection using a neuro-fuzzy EMG classifier. *IEEE Engineering in Medicine and Biology Magazine* 21, 6 (2002), 123–129. <https://doi.org/10.1109/MEMB.2002.1175148>
- Yerin Hwang, Yongil Kim, Hyunkyung Bae, Hwanhee Lee, Jeessoo Bang, and Kyomin Jung. 2023. Dialogizer: Context-aware Conversational-QA Dataset Generation from Textual Sources. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 8806–8828. <https://doi.org/10.18653/v1/2023.emnlp-main.545>
- Paolo Italiani, Giacomo Frisoni, Gianluca Moro, Antonella Carbonaro, and Claudio Sartori. 2024. Evidence, my Dear Watson: Abstractive dialogue summarization on learnable relevant utterances. *Neurocomputing* 572 (2024), 127132. <https://doi.org/10.1016/j.neucom.2023.127132>
- Xiaowei Jia, Sheng Li, Handong Zhao, Sungchul Kim, and Vipin Kumar. 2019. Towards Robust and Discriminative Sequential Data Learning: When and How to Perform Adversarial Training?. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Anchorage, AK, USA) (KDD '19)*. Association for Computing Machinery, New York, NY, USA, 1665–1673. <https://doi.org/10.1145/3292500.3330957>

- Wenhui Jiang, Xiaodong Gu, Yuting Chen, and Beijun Shen. 2023. DuReSE: Rewriting Incomplete Utterances via Neural Sequence Editing. *Neural Processing Letters* (03 2023), 1–18. <https://doi.org/10.1007/s11063-023-11174-8>
- Anirudh Joshi, Namit Katariya, Xavier Amatriain, and Anitha Kannan. 2020. Dr. Summarize: Global Summarization of Medical Dialogue by Exploiting Local Structures.. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 3755–3763. <https://doi.org/10.18653/v1/2020.findings-emnlp.335>
- Ashtosh Joshi, Shankar Vishwanath, Choon Teo, Vaclav Petricek, Vishy Vishwanathan, Rahul Bhagat, and Jonathan May. 2022. Augmenting Training Data for Massive Semantic Matching Models in Low-Traffic E-commerce Stores. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, Anastassia Loukina, Rashmi Gangadharaiiah, and Bonan Min (Eds.). Association for Computational Linguistics, Hybrid: Seattle, Washington + Online, 160–167. <https://doi.org/10.18653/v1/2022.naacl-industry.19>
- Danica Jovanovic and Theo Van Leeuwen. 2018. Multimodal dialogue on social media. *Social Semiotics* 28, 5 (2018), 683–699. <https://doi.org/10.1080/10350330.2018.1504732> arXiv:<https://doi.org/10.1080/10350330.2018.1504732>
- Daniel Jurafsky, Elizabeth Shriberg, and Debra Biasca. 1997. *Switchboard SWBD-DAMSL Shallow-Discourse-Function Annotation Coders Manual, Draft 13*. Technical Report 97-02. University of Colorado, Boulder Institute of Cognitive Science, Boulder, CO.
- Georgi Karadzhov, Tom Stafford, and Andreas Vlachos. 2021. DeliData: A dataset for deliberation in multi-party problem solving. *ArXiv abs/2108.05271* (2021).
- Xirui Ke, Jing Zhang, Xin Lv, Yiqi Xu, Shulin Cao, Cuiping Li, Hong Chen, and Juanzi Li. 2022. Knowledge-augmented Self-training of A Question Rewriter for Conversational Knowledge Base Question Answering. In *Findings of the Association for Computational Linguistics: EMNLP 2022*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 1844–1856. <https://aclanthology.org/2022.findings-emnlp.133>
- Anant Khandelwal. 2021. WeaSUL: Weakly Supervised Dialogue Policy Learning: Reward Estimation for Multi-turn Dialogue. In *Workshop on Document-grounded Dialogue and Conversational Question Answering*.
- D Kiela and J Weston. 2019. What makes a good conversation? how controllable attributes affect human judgments. In *Proceedings of the 2019 Conference of the NAACL NAACL-HLT*.
- Hyunwoo Kim, Jack Hessel, Liwei Jiang, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2022a. SODA: Million-scale Dialogue Distillation with Social Commonsense Contextualization. arXiv:2212.10465 [cs.CL]
- Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. 2022c. ProsocialDialog: A Prosocial Backbone for Conversational Agents. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 4005–4029. <https://aclanthology.org/2022.emnlp-main.267>
- Seungone Kim, Se June Joo, Hyungjoo Chae, Chaehyeon Kim, Seung-won Hwang, and Jinyoung Yeo. 2022b. Mind the Gap! Injecting Commonsense Knowledge for Abstractive Dialogue Summarization. In *Proceedings of the 29th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Gyeongju, Republic of Korea, 6285–6300. <https://aclanthology.org/2022.coling-1.548>
- Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, Anna Kocoń, Bartłomiej Koptyra, Wiktoria Miesleszczenko-Kowszewicz, Piotr Miłkowski, Marcin Oleksy, Maciej Piasecki, Łukasz Radliński, Konrad Wojtasik, Stanisław Woźniak, and Przemysław Kazienko. 2023. ChatGPT: Jack of all trades, master of none. arXiv:2302.10724 [cs.CL]
- Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2022. Internet-Augmented Dialogue Generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 8460–8478. <https://doi.org/10.18653/v1/2022.acl-long.579>
- Kira Kretschmar, Holly Tyroll, Gabriela Pavarini, Arianna Manzini, Ilina Singh, and NeurOx Young People’s Advisory Group. 2019. Can your phone be your therapist? Young people’s ethical perspectives on the use of fully automated conversational agents (chatbots) in mental health support. *Biomedical informatics insights* 11 (2019), 1178222619829083.
- Apoorv Kulshreshtha, Daniel De Freitas Adiwardana, David Richard So, Gaurav Nemade, Jamie Hall, Noah Fiedel, Quoc V. Le, Romal Thoppilan, Thang Luong, Yifeng Lu, and Zi Yang. 2020. Towards a Human-like Open-Domain Chatbot. In *arXiv*.
- Shivani Kumar, Shubham Dudeja, Md Shad Akhtar, and Tanmoy Chakraborty. 2023a. Emotion Flip Reasoning in Multiparty Conversations. *arXiv preprint arXiv:2306.13959* (2023).
- Shivani Kumar, Atharva Kulkarni, Md Shad Akhtar, and Tanmoy Chakraborty. 2022a. When did you become so smart, oh wise one?! Sarcasm Explanation in Multi-modal Multi-party Dialogues. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 5956–5968. <https://doi.org/10.18653/v1/2022.acl-long.411>
- Shivani Kumar, Ishani Mondai, Md Shad Akhtar, and Tanmoy Chakraborty. 2023b. Explaining (Sarcastic) Utterances to Enhance Affect Understanding in Multimodal Dialogues. In *Proceedings of the Thirty-Seventh AAAI Conference on*

- Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence (AAAI'23/IAAI'23/EAAI'23)*. AAAI Press, Article 1457, 9 pages. <https://doi.org/10.1609/aaai.v37i11.26526>
- Shivani Kumar, Ishani Mondal, Md Shad Akhtar, and Tanmoy Chakraborty. 2022b. Explaining (Sarcastic) Utterances to Enhance Affect Understanding in Multimodal Dialogues. arXiv:2211.11049 [cs.CL]
- Shivani Kumar, Anubhav Shrivastava, Md Shad Akhtar, and Tanmoy Chakraborty. 2022c. Discovering emotion and reasoning its flip in multi-party conversations using masked memory network and transformer. *Knowledge-Based Systems* 240 (2022), 108112.
- Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. 2016. Neural Architectures for Named Entity Recognition. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, San Diego, California, 260–270. <https://doi.org/10.18653/v1/N16-1030>
- Stefan Larson, Anish Mahendran, Joseph J. Peper, Christopher Clarke, Andrew Lee, Parker Hill, Jonathan K. Kummerfeld, Kevin Leach, Michael A. Laurenzano, Lingjia Tang, and Jason Mars. 2019. An Evaluation Dataset for Intent Classification and Out-of-Scope Prediction. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 1311–1316. <https://doi.org/10.18653/v1/D19-1131>
- Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, et al. 2022. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. *Advances in Neural Information Processing Systems* 35 (2022), 31809–31826.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324. <https://doi.org/10.1109/5.726791>
- Andrew Lee, Zhe Chen, Kevin Leach, and Jonathan K. Kummerfeld. 2022. Augmenting Task-Oriented Dialogue Systems with Relation Extraction. *ArXiv abs/2210.13344* (2022).
- Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2023. OrchestraLLM: Efficient Orchestration of Language Models for Dialogue State Tracking. *arXiv preprint arXiv:2311.09758* (2023).
- Dongyub Lee, Jung Hoon Lim, Taesun Whang, Chanhee Lee, Seung Woo Cho, Mingun Park, and Heuiseok Lim. 2021. Capturing Speaker Incorrectness: Speaker-Focused Post-Correction for Abstractive Dialogue Summarization. *Proceedings of the Third Workshop on New Frontiers in Summarization* (2021).
- Jindong Leo-Liu. 2023. Loving a “defiant” AI companion? The gender performance and ethics of social exchange robots in simulated intimate interactions. *Computers in Human Behavior* 141 (2023), 107620.
- Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. The Winograd Schema Challenge. In *Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning (KR'12)*. AAAI Press, 552–561.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 7871–7880. <https://doi.org/10.18653/v1/2020.ac1-main.703>
- Junlong Li, Zhuosheng Zhang, and Hai Zhao. 2022b. Dialogue-adaptive Language Model Pre-training From Quality Estimation. arXiv:2009.04984 [cs.CL]
- Shimin Li, Qinyuan Cheng, Linyang Li, and Xipeng Qiu. 2022a. Mitigating Negative Style Transfer in Hybrid Dialogue System. *ArXiv abs/2212.07183* (2022).
- Wei Li, Yang Li, Vlad Pandlea, Mengshi Ge, Luyao Zhu, and Erik Cambria. 2023. ECEP: Emotion-Cause Pair Extraction in Conversations. *IEEE Transactions on Affective Computing* 14, 3 (2023), 1754–1765. <https://doi.org/10.1109/TAFFC.2022.3216551>
- Xin Li, Piji Li, Yan Wang, Xiaojiang Liu, and Wai Lam. 2021b. Enhancing Dialogue Generation via Multi-Level Contrastive Learning. arXiv:2009.09147 [cs.CL]
- Yanran Li, Wenjie Li, and Zhitao Wang. 2021a. Graph-Structured Context Understanding for Knowledge-Grounded Response Generation (*SIGIR '21*). Association for Computing Machinery, New York, NY, USA, 1930–1934. <https://doi.org/10.1145/3404835.3463000>
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. DailyDialog: A Manually Labelled Multi-turn Dialogue Dataset. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Asian Federation of Natural Language Processing, Taipei, Taiwan, 986–995. <https://aclanthology.org/I17-1099>
- Yue Li and Jiong Zhang. 2021. Semi-supervised Meta-learning for Cross-domain Few-shot Intent Classification. In *Proceedings of the 1st Workshop on Meta Learning and Its Applications to Natural Language Processing*, Hung-Yi Lee, Mitra Mohtarami, Shang-Wen Li, Di Jin, Mandy Korpusik, Shuyan Dong, Ngoc Thang Vu, and Dilek Hakkani-Tur (Eds.). Association for Computational Linguistics, Online, 67–75. <https://doi.org/10.18653/v1/2021.metanlp-1.8>

- Chen Liang, Jonathan Berant, Quoc Le, Kenneth D. Forbus, and Ni Lao. 2017. Neural Symbolic Machines: Learning Semantic Parsers on Freebase with Weak Supervision. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Vancouver, Canada, 23–33. <https://doi.org/10.18653/v1/P17-1003>
- Xinnian Liang, Shuangzhi Wu, Chenhao Cui, Jiaqi Bai, Chao Bian, and Zhoujun Li. 2023. Enhancing Dialogue Summarization with Topic-Aware Global- and Local- Level Centrality. arXiv:2301.12376 [cs.CL]
- Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*. Association for Computational Linguistics, Barcelona, Spain, 74–81. <https://aclanthology.org/W04-1013>
- Sheng-Chieh Lin, Jheng-Hong Yang, and Jimmy Lin. 2021. Contextualized Query Embeddings for Conversational Search. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (Eds.). Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 1004–1015. <https://doi.org/10.18653/v1/2021.emnlp-main.77>
- Xi Victoria Lin, Richard Socher, and Caiming Xiong. 2020. Bridging Textual and Tabular Data for Cross-Domain Text-to-SQL Semantic Parsing. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 4870–4888. <https://doi.org/10.18653/v1/2020.findings-emnlp.438>
- Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, Peng Xu, Feijun Jiang, Yuxiang Hu, Chen Shi, and Pascale Fung. [n.d.]. BiToD: A Bilingual Multi-Domain Dataset For Task-Oriented Dialogue Modeling. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Zachary Lipton, Xiujun Li, Jianfeng Gao, Lihong Li, Faisal Ahmed, and Li Deng. 2018. Bbq-networks: Efficient exploration in deep reinforcement learning for task-oriented dialogue systems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 32.
- Bing Liu and Sahisnu Mazumder. 2021. Lifelong and continual learning dialogue systems: learning during conversation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 15058–15063.
- Bing Liu, Gokhan Tur, Dilek Hakkani-Tur, Pararth Shah, and Larry Heck. 2017b. End-to-end optimization of task-oriented dialogue model with deep reinforcement learning. *arXiv preprint arXiv:1711.10712* (2017).
- Chunyi Liu, Peng Wang, Jiang Xu, Zang Li, and Jieping Ye. 2019. Automatic Dialogue Summary Generation for Customer Service. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Anchorage, AK, USA) (KDD '19)*. Association for Computing Machinery, New York, NY, USA, 1957–1965. <https://doi.org/10.1145/3292500.3330683>
- Chia-Wei Liu, Ryan Lowe, Iulian V. Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. 2017a. How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation. arXiv:1603.08023 [cs.CL]
- Han Liu, Siyang Zhao, Xiaotong Zhang, Feng Zhang, Junjie Sun, Hong Yu, and Xianchao Zhang. 2022d. A Simple Meta-Learning Paradigm for Zero-Shot Intent Classification with Mixture Attention Mechanism. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (<conf-loc>, <city>Madrid</city>, <country>Spain</country>, </conf-loc>)* (SIGIR '22). Association for Computing Machinery, New York, NY, USA, 2047–2052. <https://doi.org/10.1145/3477495.3531803>
- Qingbin Liu, Guirong Bai, Shizhu He, Cao Liu, Kang Liu, and Jun Zhao. 2021a. Heterogeneous relational graph neural networks with adaptive objective for end-to-end task-oriented dialogue. *Knowledge-Based Systems* 227 (2021), 107186.
- Qian Liu, Bei Chen, Jian-Guang Lou, Bin Zhou, and Dongmei Zhang. 2020b. Incomplete Utterance Rewriting as Semantic Segmentation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 2846–2857. <https://doi.org/10.18653/v1/2020.emnlp-main.227>
- Qian Liu, Yihong Chen, B. Chen, Jian-Guang Lou, Zixuan Chen, Bin Zhou, and Dongmei Zhang. 2020a. You Impress Me: Dialogue Generation via Mutual Persona Perception. In *Annual Meeting of the Association for Computational Linguistics*.
- Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2021b. Benchmarking natural language understanding services for building conversational agents. In *Increasing Naturalness and Flexibility in Spoken Dialogue Interaction: 10th International Workshop on Spoken Dialogue Systems*. Springer, 165–183.
- Xingkun Liu, Arash Eshghi, Pawel Swietojanski, and Verena Rieser. 2021c. *Benchmarking Natural Language Understanding Services for Building Conversational Agents*. Springer Singapore, Singapore, 165–183. https://doi.org/10.1007/978-981-15-9323-9_15
- Yongkang Liu, Shi Feng, Wei Gao, Daling Wang, and Yifei Zhang. 2022a. DialConv: A Lightweight Fully Convolutional Network for Multi-view Response Selection. arXiv:2210.13845 [cs.CL]
- Yongtai Liu, Joshua Maynez, Gonalo Simões, and Shashi Narayan. 2022b. Data Augmentation for Low-Resource Dialogue Summarization. In *Findings of the Association for Computational Linguistics: NAACL 2022*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 703–710. <https://doi.org/10.18653/v1/2022.findings-naacl.53>
- Zhengyuan Liu and Nancy F. Chen. 2021. Controllable Neural Dialogue Summarization with Personal Named Entity Planning. In *Conference on Empirical Methods in Natural Language Processing*.

- Zeming Liu, Jun Xu, Zeyang Lei, Haifeng Wang, Zheng-Yu Niu, and Hua Wu. 2022c. Where to Go for the Holidays: Towards Mixed-Type Dialogs for Clarification of User Goals. In *Annual Meeting of the Association for Computational Linguistics*.
- Samuel Louvan and Bernardo Magnini. 2018. Exploring Named Entity Recognition As an Auxiliary Task for Slot Filling in Conversational Language Understanding. In *Proceedings of the 2018 EMNLP Workshop SCAI: The 2nd International Workshop on Search-Oriented Conversational AI*. Association for Computational Linguistics, Brussels, Belgium, 74–80. <https://doi.org/10.18653/v1/W18-5711>
- Samuel Louvan and Bernardo Magnini. 2019. Leveraging Non-Conversational Tasks for Low Resource Slot Filling: Does it help?. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*. Association for Computational Linguistics, Stockholm, Sweden, 85–91. <https://doi.org/10.18653/v1/W19-5911>
- Samuel Louvan and Bernardo Magnini. 2020. Recent neural methods on slot filling and intent classification for task-oriented dialogue systems: A survey. *arXiv preprint arXiv:2011.00564* (2020).
- Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The Ubuntu Dialogue Corpus: A Large Dataset for Research in Unstructured Multi-Turn Dialogue Systems. In *Proceedings of the 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*. Association for Computational Linguistics, Prague, Czech Republic, 285–294. <https://doi.org/10.18653/v1/W15-4640>
- Gale M Lucas, Jill Boberg, David Traum, Ron Artstein, Jonathan Gratch, Alesia Gainer, Emmanuel Johnson, Anton Leuski, and Mikio Nakano. 2018. Culture, errors, and rapport-building dialogue in social agents. In *Proceedings of the 18th International Conference on intelligent virtual agents*. 51–58.
- Queenie Luo, Michael J. Puett, and Michael D. Smith. 2023. A Perspectival Mirror of the Elephant: Investigating Language Bias on Google, ChatGPT, Wikipedia, and YouTube. *arXiv:2303.16281* [cs.CY]
- Xuezhe Ma and Eduard Hovy. 2016. End-to-end Sequence Labeling via Bi-directional LSTM-CNNs-CRF. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, 1064–1074. <https://doi.org/10.18653/v1/P16-1101>
- Ganeshan Malhotra, Abdul Waheed, Aseem Srivastava, Md Shad Akhtar, and Tanmoy Chakraborty. 2022. Speaker and Time-Aware Joint Contextual Learning for Dialogue-Act Classification in Counselling Conversations. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining (Virtual Event, AZ, USA) (WSDM '22)*. Association for Computing Machinery, New York, NY, USA, 735–745. <https://doi.org/10.1145/3488560.3498509>
- Scott Martin, Shivani Poddar, and Kartikeya Upasani. 2020. MuDoCo: Corpus for Multidomain Coreference Resolution and Referring Expression Generation. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 104–111. <https://aclanthology.org/2020.lrec-1.13>
- Susan W McRoy, Songsak Channarukul, and Syed S Ali. 2003. An augmented template-based approach to text realization. *Natural Language Engineering* 9, 4 (2003), 381–420.
- Michael McTear. 2021. *Rule-Based Dialogue Systems: Architecture, Methods, and Tools*. Springer International Publishing, Cham, 43–70. https://doi.org/10.1007/978-3-031-02176-3_2
- Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tur. 2020. DialoGLUE: A Natural Language Understanding Benchmark for Task-Oriented Dialogue. *arXiv:2009.13570* [cs.CL]
- Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019. Pretraining Methods for Dialog Context Representation Learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 3836–3845. <https://doi.org/10.18653/v1/P19-1373>
- Chuan Meng, Pengjie Ren, Zhumin Chen, Weiwei Sun, Zhaochun Ren, Zhaopeng Tu, and Maarten de Rijke. 2020. DukeNet: A Dual Knowledge Interaction Network for Knowledge-Grounded Conversation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (Virtual Event, China) (SIGIR '20)*. Association for Computing Machinery, New York, NY, USA, 1151–1160. <https://doi.org/10.1145/3397271.3401097>
- Xiaojun Meng, Wenlin Dai, Yasheng Wang, Baojun Wang, Zhiyong Wu, Xin Jiang, and Qun Liu. 2022. Lexicon-injected Semantic Parsing for Task-Oriented Dialog. *ArXiv abs/2211.14508* (2022).
- Fei Mi, Liangwei Chen, Mengjie Zhao, Minlie Huang, and Boi Faltings. 2020. Continual learning for natural language generation in task-oriented dialog systems. *arXiv preprint arXiv:2010.00910* (2020).
- Shaobo Min, Hantao Yao, Hongtao Xie, Chaqun Wang, Zheng-Jun Zha, and Yongdong Zhang. 2020. Domain-Aware Visual Bias Eliminating for Generalized Zero-Shot Learning. 12661–12670. <https://doi.org/10.1109/CVPR42600.2020.01268>
- Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. OpenDialKG: Explainable Conversational Reasoning with Attention-based Walks over Knowledge Graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 845–854. <https://doi.org/10.18653/v1/P19-1081>
- Christian Muise, Tathagata Chakraborti, Shubham Agarwal, Ondrej Bajgar, Arunima Chaudhary, Luis A Lastras-Montano, Josef Ondrej, Miroslav Vodolan, and Charlie Wiecha. 2019. Planning for goal-oriented dialogue systems. *arXiv preprint*

- arXiv:1910.08137 (2019).
- Shashi Narayan, Yao Zhao, Joshua Maynez, Gonalo Simões, Vitaly Nikolaev, and Ryan T. McDonald. 2021. Planning with Learned Entity Prompts for Abstractive Summarization. *Transactions of the Association for Computational Linguistics* 9 (2021), 1475–1492.
- Anna Nedoluzhko, Muskaan Singh, Marie Hledíková, Tirthankar Ghosal, and Ondřej Bojar. 2022. ELITR Minuting Corpus: A Novel Dataset for Automatic Minuting from Multi-Party Meetings in English and Czech. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H el ene Mazo, Jan Odijk, and Stelios Piperidis (Eds.). European Language Resources Association, Marseille, France, 3174–3182. <https://aclanthology.org/2022.lrec-1.340>
- Olabiyi Oluwatobi and Erik Mueller. 2020. DLGNet: A transformer-based model for dialogue response generation. In *Proceedings of the 2nd workshop on natural language processing for conversational AI*. 54–62.
- Boyun Onyshkevych. 1993. Template design for information extraction. In *Fifth Message Understanding Conference (MUC-5): Proceedings of a Conference Held in Baltimore, Maryland, August 25-27, 1993*.
- Siru Ouyang, Zhuosheng Zhang, and Hai Zhao. 2021. Dialogue Graph Modeling for Conversational Machine Reading. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 3158–3169. <https://doi.org/10.18653/v1/2021.findings-acl.279>
- Harion A Pandya and Brijesh S Bhatt. 2021. Question answering survey: Directions, challenges, datasets, evaluation matrices. *arXiv preprint arXiv:2112.03572* (2021).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Philadelphia, Pennsylvania, USA, 311–318. <https://doi.org/10.3115/1073083.1073135>
- Amin Parvaneh, Ehsan Abbasnejad, Qi Wu, and Javen Qinfeng Shi. 2019. Show, Price and Negotiate: A Hierarchical Attention Recurrent Visual Negotiator. *ArXiv abs/1905.03721* (2019).
- Panupong Pasupat, S. Gupta, Karishma Mandyam, Rushin Shah, Michael Lewis, and Luke Zettlemoyer. 2019. Span-based Hierarchical Semantic Parsing for Task-Oriented Dialog. In *Conference on Empirical Methods in Natural Language Processing*.
- Debjit Paul, Daniil Sorokin, and Judith Gaspers. 2022. Class Incremental Learning for Intent Classification with Limited or No Old Data. In *Proceedings of the The First Workshop on Ever Evolving NLP (EvoNLP)*, Francesco Barbieri, Jose Camacho-Collados, Bhuwan Dhingra, Luis Espinosa-Anke, Elena Gribovskaya, Angeliki Lazaridou, Daniel Loureiro, and Leonardo Neves (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates (Hybrid), 16–25. <https://doi.org/10.18653/v1/2022.evonlp-1.4>
- Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayandeh, Lars Liden, and Jianfeng Gao. 2021. Soloist: Building Task Bots at Scale with Transfer Learning and Machine Teaching. *Transactions of the Association for Computational Linguistics* 9 (2021), 807–824. https://doi.org/10.1162/tacl_a_00399
- Patr icia Pereira, Helena Moniz, and Joao Paulo Carvalho. 2022. Deep Emotion Recognition in Textual Conversations: A Survey. arXiv:2211.09172 [cs.CL]
- Xinyu Pi, Wanjun Zhong, Yan Gao, Nan Duan, and Jian-Guang Lou. 2022. LogiGAN: Learning Logical Reasoning via Adversarial Pre-training. arXiv:2205.08794 [cs.CL]
- Anita Pomerantz and Barbara J Fehr. 2011. Conversation analysis: An approach to the analysis of social interaction. *Discourse studies: A multidisciplinary introduction* 2 (2011), 165–190.
- Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2019. MELD: A Multimodal Multi-Party Dataset for Emotion Recognition in Conversations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 527–536. <https://doi.org/10.18653/v1/P19-1050>
- Soujanya Poria, Navonil Majumder, Devamanyu Hazarika, Deepanway Ghosal, Rishabh Bhardwaj, Samson Yu Bai Jian, Pengfei Hong, Romila Ghosh, Abhinaba Roy, Niyati Chhaya, et al. 2021. Recognizing emotion cause in conversations. *Cognitive Computation* 13 (2021), 1317–1332.
- Bowen Qin, Binyuan Hui, Lihan Wang, Min Yang, Jinyang Li, Binhua Li, Ruiying Geng, Rongyu Cao, Jian Sun, Luo Si, et al. 2022. A Survey on Text-to-SQL Parsing: Concepts, Methods, and Future Directions. *arXiv preprint arXiv:2208.13629* (2022).
- Zongfeng Qu, Zhitong Yang, Bo Wang, and Qinghua Hu. 2024. TodBR: Target-Oriented Dialog with Bidirectional Reasoning on Knowledge Graph. *Applied Sciences* 14, 1 (2024), 459.
- Jun Quan, Deyi Xiong, Bonnie Webber, and Changjian Hu. 2019. GECOR: An End-to-End Generative Ellipsis and Coreference Resolution Model for Task-Oriented Dialogue. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 4547–4557. <https://doi.org/10.18653/v1/D19-1462>

- L. Rabiner and B. Juang. 1986. An introduction to hidden Markov models. *IEEE ASSP Magazine* 3, 1 (1986), 4–16. <https://doi.org/10.1109/MASSP.1986.1165342>
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* 21, 1, Article 140 (jan 2020), 67 pages.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know What You Don’t Know: Unanswerable Questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Melbourne, Australia, 784–789. <https://doi.org/10.18653/v1/P18-2124>
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 2383–2392. <https://doi.org/10.18653/v1/D16-1264>
- Revant Rameshkumar and Peter Bailey. 2020. Storytelling with Dialogue: A Critical Role Dungeons and Dragons Dataset. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 5121–5134. <https://doi.org/10.18653/v1/2020.acl-main.459>
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2018. Towards Empathetic Open-domain Conversation Models: A New Benchmark and Dataset. In *Annual Meeting of the Association for Computational Linguistics*.
- 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*. Association for Computational Linguistics, Florence, Italy, 5370–5381. <https://doi.org/10.18653/v1/P19-1534>
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards Scalable Multi-Domain Conversational Agents: The Schema-Guided Dialogue Dataset. *Proceedings of the AAI Conference on Artificial Intelligence* 34, 05 (Apr. 2020), 8689–8696. <https://doi.org/10.1609/aaai.v34i05.6394>
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A Conversational Question Answering Challenge. *Transactions of the Association for Computational Linguistics* 7 (2019), 249–266. https://doi.org/10.1162/tac1_a_00266
- Siva Reddy, Mirella Lapata, and Mark Steedman. 2014. Large-scale Semantic Parsing without Question-Answer Pairs. *Transactions of the Association for Computational Linguistics* 2 (2014), 377–392. https://doi.org/10.1162/tac1_a_00190
- Shiya Ren, Huaming Wang, Dongming Yu, Yuan Li, Zhixing Li, S Hu, and L Zou. 2018. Joint Intent Detection and Slot Filling with Rules. *CCKS Tasks* 2242 (2018), 34–40.
- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. Recipes for Building an Open-Domain Chatbot. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, 300–325. <https://doi.org/10.18653/v1/2021.eacl-main.24>
- Johanna Ruusuvuori. 2012. Emotion, affect and conversation. *The handbook of conversation analysis* (2012), 330–349.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100* (2022).
- Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. PICARD: Parsing Incrementally for Constrained Auto-Regressive Decoding from Language Models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 9895–9901. <https://doi.org/10.18653/v1/2021.emnlp-main.779>
- Hendrik Schuff, Lindsey Vanderlyn, Heike Adel, and Ngoc Thang Vu. 2023. How to do human evaluation: A brief introduction to user studies in NLP. *Natural Language Engineering* (2023), 1–24. <https://doi.org/10.1017/S1351324922000535>
- Karin Sevegnani, David M. Howcroft, Ioannis Konstas, and Verena Rieser. 2021. OTTers: One-turn Topic Transitions for Open-Domain Dialogue. 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)*. Association for Computational Linguistics, Online, 2492–2504. <https://doi.org/10.18653/v1/2021.acl-long.194>
- Igor Shalyminov, Sungjin Lee, Arash Eshghi, and Oliver Lemon. 2019. Few-Shot Dialogue Generation Without Annotated Data: A Transfer Learning Approach. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*.

- Association for Computational Linguistics, Stockholm, Sweden, 32–39. <https://doi.org/10.18653/v1/W19-5904>
- Igor Shalymov, Alessandro Sordani, Adam Atkinson, and Hannes Schulz. 2020. Fast Domain Adaptation for Goal-Oriented Dialogue Using a Hybrid Generative-Retrieval Transformer. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 8039–8043. <https://doi.org/10.1109/ICASSP40776.2020.9053599>
- Weizhou Shen, Siyue Wu, Yunyi Yang, and Xiaojun Quan. 2021. Directed Acyclic Graph Network for Conversational Emotion Recognition. 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)*. Association for Computational Linguistics, Online, 1551–1560. <https://doi.org/10.18653/v1/2021.acl-long.123>
- N Shikha, Karan Naidu, Antara Roy Choudhury, and N Kayarvizhy. 2022. Smart Memory Companion for elderly. In *2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICAC3N)*. IEEE, 1497–1502.
- Jamin Shin, Peng Xu, Andrea Madotto, and Pascale Fung. 2019. HappyBot: Generating Empathetic Dialogue Responses by Improving User Experience Look-ahead. *ArXiv abs/1906.08487* (2019).
- Michael Shum, Stephan Zheng, Wojciech Kryscinski, Caiming Xiong, and Richard Socher. 2020. Sketch-Fill-A-R: A Persona-Grounded Chit-Chat Generation Framework. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*. Association for Computational Linguistics, Online, 118–131. <https://doi.org/10.18653/v1/2020.nlp4convai-1.14>
- Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, et al. 2022. Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage. *arXiv preprint arXiv:2208.03188* (2022).
- A.B. Siddique, Fuad Jamour, Luxun Xu, and Vagelis Hristidis. 2021. Generalized Zero-Shot Intent Detection via Commonsense Knowledge (*SIGIR '21*). Association for Computing Machinery, New York, NY, USA, 1925–1929. <https://doi.org/10.1145/3404835.3462985>
- Xiaohui Song, Longtao Huang, Hui Xue, and Songlin Hu. 2022. Supervised Prototypical Contrastive Learning for Emotion Recognition in Conversation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 5197–5206. <https://aclanthology.org/2022.emnlp-main.347>
- Aseem Srivastava, Ishan Pandey, Md Shad Akhtar, and Tanmoy Chakraborty. 2023. Response-Act Guided Reinforced Dialogue Generation for Mental Health Counseling. In *Proceedings of the ACM Web Conference 2023 (Austin, TX, USA) (WWW '23)*. Association for Computing Machinery, New York, NY, USA, 1118–1129. <https://doi.org/10.1145/3583380>
- Aseem Srivastava, Tharun Suresh, Sarah P. Lord, Md Shad Akhtar, and Tanmoy Chakraborty. 2022. Counseling Summarization Using Mental Health Knowledge Guided Utterance Filtering. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (Washington DC, USA) (KDD '22)*. Association for Computing Machinery, New York, NY, USA, 3920–3930. <https://doi.org/10.1145/3534678.3539187>
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 3008–3021. https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf
- Carl Strathearn and Dimitra Gkatzia. 2022. Task2Dial: A Novel Task and Dataset for Commonsense-enhanced Task-based Dialogue Grounded in Documents. In *Proceedings of the Second DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering*. Association for Computational Linguistics, Dublin, Ireland, 187–196. <https://doi.org/10.18653/v1/2022.dialdoc-1.21>
- Yixuan Su, Deng Cai, Yan Wang, Simon Baker, Anna Korhonen, Nigel Collier, and Xiaojiang Liu. 2020. Stylistic dialogue generation via information-guided reinforcement learning strategy. *arXiv preprint arXiv:2004.02202* (2020).
- Alane Suhr, Ming-Wei Chang, Peter Shaw, and Kenton Lee. 2020. Exploring Unexplored Generalization Challenges for Cross-Database Semantic Parsing. In *Annual Meeting of the Association for Computational Linguistics*.
- Qingfeng Sun, Yujing Wang, Can Xu, Kai Zheng, Yaming Yang, Huang Hu, Fei Xu, Jessica Zhang, Xiubo Geng, and Daxin Jiang. 2022b. Multimodal Dialogue Response Generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 2854–2866. <https://doi.org/10.18653/v1/2022.acl-long.204>
- Yuqian Sun, Xuran Ni, Haozhen Feng, Ray LC, Chang Hee Lee, and Ali Asadipour. 2022a. Bringing stories to life in 1001 nights: A co-creative text adventure game using a story generation model. In *International Conference on Interactive Digital Storytelling*. Springer, 651–672.
- Hao Tang, Donghong Ji, and Qiji Zhou. 2020. End-to-end masked graph-based CRF for joint slot filling and intent detection. *Neurocomputing* 413 (2020), 348–359. <https://doi.org/10.1016/j.neucom.2020.06.113>

- Abha Tewari, Amit Chhabria, Ajay Singh Khalsa, Sanket Chaudhary, and Harshita Kanal. 2021. A survey of mental health chatbots using NLP. In *Proceedings of the International Conference on Innovative Computing & Communication (ICICC)*.
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulse Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. LaMDA: Language Models for Dialog Applications. arXiv:2201.08239 [cs.CL]
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. arXiv:2302.13971 [cs.CL]
- Enrica Troiano, Aswathy Velutharambath, and Roman Klinger. 2023. From theories on styles to their transfer in text: Bridging the gap with a hierarchical survey. *Natural Language Engineering* 29, 4 (2023), 849–908. <https://doi.org/10.1017/S1351324922000407>
- Don Tuggener, Margot Mieskes, Jan Deriu, and Mark Cieliebak. 2021. Are We Summarizing the Right Way? A Survey of Dialogue Summarization Data Sets. In *Proceedings of the Third Workshop on New Frontiers in Summarization*. Association for Computational Linguistics, Online and in Dominican Republic, 107–118. <https://doi.org/10.18653/v1/2021.newsum-1.12>
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Kraemer. 2021. Human evaluation of automatically generated text: Current trends and best practice guidelines. *Computer Speech & Language* 67 (2021), 101151. <https://doi.org/10.1016/j.cs1.2020.101151>
- Hamsa Shwetha Venkataram, Chris A Mattmann, and Scott Penberthy. 2020. TopiQAL: Topic-aware Question Answering using Scalable Domain-specific Supercomputers. In *2020 IEEE/ACM Fourth Workshop on Deep Learning on Supercomputers (DLS)*. IEEE, 48–55.
- Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. 2015. Order matters: Sequence to sequence for sets. *arXiv preprint arXiv:1511.06391* (2015).
- Tu Yu, Aditya Barua, Brian Lester, Daniel Cer, Mohit Iyyer, and Noah Constant. 2022. Overcoming Catastrophic Forgetting in Zero-Shot Cross-Lingual Generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 9279–9300. <https://doi.org/10.18653/v1/2022.emnlp-main.630>
- Ivan Vulić, Iñigo Casanueva, Georgios Spithourakis, Avishek Mondal, Tsung-Hsien Wen, and Paweł Budzianowski. 2022. Multi-Label Intent Detection via Contrastive Task Specialization of Sentence Encoders. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 7544–7559. <https://aclanthology.org/2022.emnlp-main.512>
- Kei Wakabayashi, Johane Takeuchi, and Mikio Nakano. 2022. Robust Slot Filling Modeling for Incomplete Annotations using Segmentation-Based Formulation. *Transactions of the Japanese Society for Artificial Intelligence* 37, 3 (2022), IDS–E_1–12. https://doi.org/10.1527/tjsai.37-3_IDS-E
- Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. 2019b. RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers. In *Annual Meeting of the Association for Computational Linguistics*.
- Jieyu Wang, Dingfang Kang, Abdullah AbuHussein, and Lynn A Collen. 2023a. Designing a Conversational Agent for Education: A Personality-based Approach. (2023).
- Yufan Wang, Tingting He, Rui Fan, Wenji Zhou, and Xinhui Tu. 2019a. Effective Utilization of External Knowledge and History Context in Multi-turn Spoken Language Understanding Model. In *2019 IEEE International Conference on Big Data (Big Data)*. 960–967. <https://doi.org/10.1109/BigData47090.2019.9006162>
- Yanmeng Wang, Wenge Rong, Jianfei Zhang, Yuanxin Ouyang, and Zhang Xiong. 2020. Knowledge Grounded Pre-Trained Model For Dialogue Response Generation. In *2020 International Joint Conference on Neural Networks (IJCNN)*. 1–8. <https://doi.org/10.1109/IJCNN48605.2020.9207054>
- Zhenduo Wang, Yuancheng Tu, Corby Rosset, Nick Craswell, Ming Wu, and Qingyao Ai. 2023b. Zero-shot Clarifying Question Generation for Conversational Search. arXiv:2301.12660 [cs.IR]
- Nick Webb. 2000. Rule-based dialogue management systems. In *Proceedings of the 3rd International Workshop on Human-Computer Conversation, Bellagio, Italy*. 3–5.
- Joseph Weizenbaum. 1966. ELIZA—a Computer Program for the Study of Natural Language Communication between Man and Machine. *Commun. ACM* 9, 1 (jan 1966), 36–45. <https://doi.org/10.1145/365153.365168>
- Henry Weld, Xiaoqi Huang, Siqu Long, Josiah Poon, and Soyeon Caren Han. 2022a. A Survey of Joint Intent Detection and Slot Filling Models in Natural Language Understanding. *ACM Comput. Surv.* 55, 8, Article 156 (dec 2022), 38 pages.

- <https://doi.org/10.1145/3547138>
- Henry Weld, Xiaoqi Huang, Siyu Long, Josiah Poon, and Soyeon Caren Han. 2022b. A survey of joint intent detection and slot filling models in natural language understanding. *Comput. Surveys* 55, 8 (2022), 1–38.
- Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M. Rush, Bart van Merriënboer, Armand Joulin, and Tomas Mikolov. 2015. Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks. arXiv:1502.05698 [cs.AI]
- Jason D Williams. 2003. A probabilistic model of human/computer dialogue with application to a partially observable Markov decision process. *PhD first year report. Department of Engineering, University of Cambridge* (2003).
- Jason D Williams, Pascal Poupart, and Steve Young. 2005. Factored partially observable Markov decision processes for dialogue management. In *Proc. IJCAI Workshop on Knowledge and Reasoning in Practical Dialogue Systems*. Citeseer, 76–82.
- Chien-Sheng Wu, Linqing Liu, Wenhao Liu, Pontus Stenertorp, and Caiming Xiong. 2021b. Controllable Abstractive Dialogue Summarization with Sketch Supervision. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 5108–5122. <https://doi.org/10.18653/v1/2021.findings-acl.454>
- Hongyi Wu, Xinshu Shen, Man Lan, Shaoguang Mao, Xiaopeng Bai, and Yuanbin Wu. 2023b. A Multi-Task Dataset for Assessing Discourse Coherence in Chinese Essays: Structure, Theme, and Logic Analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 6673–6688. <https://doi.org/10.18653/v1/2023.emnlp-main.412>
- Jie Wu, Ian G. Harris, Hongzhi Zhao, and Guangming Ling. 2023a. A Graph-to-Sequence Model for Joint Intent Detection and Slot Filling. In *2023 IEEE 17th International Conference on Semantic Computing (ICSC)*. 131–138. <https://doi.org/10.1109/ICSC56153.2023.00028>
- Tongtong Wu, Meng Wang, Huan Gao, Guilin Qi, and Weizhuo Li. 2019. Zero-Shot Slot Filling via Latent Question Representation and Reading Comprehension. In *PRICAI 2019: Trends in Artificial Intelligence*, Abhaya C. Nayak and Alok Sharma (Eds.). Springer International Publishing, Cham, 123–136.
- Zequi Wu, Michel Galley, Chris Brockett, Yizhe Zhang, Xiang Gao, Chris Quirk, Rik Koncel-Kedziorski, Jianfeng Gao, Hannaneh Hajishirzi, Mari Ostendorf, and Bill Dolan. 2021a. A Controllable Model of Grounded Response Generation. *Proceedings of the AAI Conference on Artificial Intelligence* 35, 16 (May 2021), 14085–14093. <https://doi.org/10.1609/aaai.v35i16.17658>
- Rui Xia and Zixiang Ding. 2019. Emotion-Cause Pair Extraction: A New Task to Emotion Analysis in Texts. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 1003–1012. <https://doi.org/10.18653/v1/P19-1096>
- Shuwen Xiao, Zhou Zhao, Zijian Zhang, Xiaohui Yan, and Min Yang. 2020. Convolutional hierarchical attention network for query-focused video summarization. In *Proceedings of the AAI conference on artificial intelligence*, Vol. 34. 12426–12433.
- Yubo Xie and Pearl Pu. 2021. Generating Empathetic Responses with a Large Scale Dialog Dataset. *ArXiv abs/2105.06829* (2021).
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023. Baize: An Open-Source Chat Model with Parameter-Efficient Tuning on Self-Chat Data. arXiv:2304.01196 [cs.CL]
- Yan Xu, Etsuko Ishii, Samuel Cahyawijaya, Zihan Liu, Genta Indra Winata, Andrea Madotto, Dan Su, and Pascale Fung. 2022. Retrieval-Free Knowledge-Grounded Dialogue Response Generation with Adapters. In *Proceedings of the Second DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering*. Association for Computational Linguistics, Dublin, Ireland, 93–107. <https://doi.org/10.18653/v1/2022.dialdoc-1.10>
- Yi Xu and Hai Zhao. 2021. Dialogue-oriented Pre-training. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 2663–2673. <https://doi.org/10.18653/v1/2021.findings-acl.235>
- Chaoting Xuan. 2020. Improving Sequence-to-Sequence Semantic Parser for Task Oriented Dialog. *Proceedings of the First Workshop on Interactive and Executable Semantic Parsing* (2020).
- Shiquan Yang, Xinting Huang, Jey Han Lau, and Sarah Erfani. 2022. Robust Task-Oriented Dialogue Generation with Contrastive Pre-training and Adversarial Filtering. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 1220–1234. <https://doi.org/10.18653/v1/2022.findings-emnlp.88>
- Shiquan Yang, Rui Zhang, and Sarah Erfani. 2020. Graphdialog: Integrating graph knowledge into end-to-end task-oriented dialogue systems. *arXiv preprint arXiv:2010.01447* (2020).
- Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. 2015. Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Beijing, China, 1321–1331. <https://doi.org/10.3115/v1/P15-1128>

- Wen-wai Yim and Meliha Yetisgen. 2021. Towards Automating Medical Scribing : Clinic Visit Dialogue2Note Sentence Alignment and Snippet Summarization. In *Proceedings of the Second Workshop on Natural Language Processing for Medical Conversations*, Chaitanya Shivade, Rashmi Gangadharaiah, Spandana Gella, Sandeep Konam, Shaoqing Yuan, Yi Zhang, Parminder Bhatia, and Byron Wallace (Eds.). Association for Computational Linguistics, Online, 10–20. <https://doi.org/10.18653/v1/2021.nlpmc-1.2>
- Pengcheng Yin and Graham Neubig. 2017. A Syntactic Neural Model for General-Purpose Code Generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Vancouver, Canada, 440–450. <https://doi.org/10.18653/v1/P17-1041>
- Tom Young, Frank Xing, Vlad Pandealea, Jinjie Ni, and Erik Cambria. 2022. Fusing Task-Oriented and Open-Domain Dialogues in Conversational Agents. *Proceedings of the AAAI Conference on Artificial Intelligence* 36, 10 (Jun. 2022), 11622–11629. <https://doi.org/10.1609/aaai.v36i10.21416>
- Tao Yu, Rui Zhang, Heyang Er, Suyi Li, Eric Xue, Bo Pang, Xi Victoria Lin, Yi Chern Tan, Tianze Shi, Zihan Li, Youxuan Jiang, Michihiro Yasunaga, Sungrok Shim, Tao Chen, Alexander Fabbri, Zifan Li, Luyao Chen, Yuwen Zhang, Shreya Dixit, Vincent Zhang, Caiming Xiong, Richard Socher, Richard Lasecki, and Dragomir Radev. 2019. CoSQL: A Conversational Text-to-SQL Challenge Towards Cross-Domain Natural Language Interfaces to Databases. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 1962–1979. <https://doi.org/10.18653/v1/D19-1204>
- Tao Yu, Rui Zhang, Alex Polozov, Christopher Meek, and Ahmed Hassan Awadallah. 2021. {SC}oRe: Pre-Training for Context Representation in Conversational Semantic Parsing. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=oyZxhRI2RiE>
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 3911–3921. <https://doi.org/10.18653/v1/D18-1425>
- Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. 2023. Chain-of-Note: Enhancing Robustness in Retrieval-Augmented Language Models. *arXiv preprint arXiv:2311.09210* (2023).
- Idris Yusupov and Yurii Kuratov. 2018. NIPS Conversational Intelligence Challenge 2017 Winner System: Skill-based Conversational Agent with Supervised Dialog Manager. In *Proceedings of the 27th International Conference on Computational Linguistics*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 3681–3692. <https://aclanthology.org/C18-1312>
- Klaus Zechner and Alex Waibel. 2000. DIASUMM: Flexible Summarization of Spontaneous Dialogues in Unrestricted Domains. In *COLING 2000 Volume 2: The 18th International Conference on Computational Linguistics*. <https://aclanthology.org/C00-2140>
- Haolan Zhan, Hainan Zhang, Hongshen Chen, Zhuoye Ding, Yongjun Bao, and Yanyan Lan. 2021. Augmenting Knowledge-grounded Conversations with Sequential Knowledge Transition. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 5621–5630. <https://doi.org/10.18653/v1/2021.naacl-main.446>
- Qingyu Zhang, Xiaoyu Shen, Ernie Chang, Jidong Ge, and Pengke Chen. 2022. MDIA: A Benchmark for Multilingual Dialogue Generation in 46 Languages. *arXiv:2208.13078 [cs.CL]*
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing Dialogue Agents: I have a dog, do you have pets too?. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 2204–2213. <https://doi.org/10.18653/v1/P18-1205>
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020a. BERTScore: Evaluating Text Generation with BERT. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=SkeHuCVFDr>
- Wanying Zhang, Feng Yang, and Yan Liang. 2019. A Bayesian Framework for Joint Target Tracking, Classification, and Intent Inference. *IEEE Access* 7 (2019), 66148–66156. <https://doi.org/10.1109/ACCESS.2019.2917541>
- Xiaoqiang Zhang, Ying Chen, and Guang ying Li. 2021. Multi-Modal Sarcasm Detection Based on Contrastive Attention Mechanism. In *Natural Language Processing and Chinese Computing*.
- Yusen Zhang, Yang Liu, Ziyi Yang, Yuwei Fang, Yulong Chen, Dragomir Radev, Chenguang Zhu, Michael Zeng, and Rui Zhang. 2023. MACSum: Controllable Summarization with Mixed Attributes. *Transactions of the Association for Computational Linguistics* 11 (2023), 787–803. https://doi.org/10.1162/tacl_a_00575
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020b. DIALOGPT : Large-Scale Generative Pre-training for Conversational Response Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Online, 270–278. <https://doi.org/10.18653/v1/2020.acl-demos.30>

- Zheng Zhang, Ryuichi Takanobu, Qi Zhu, MinLie Huang, and XiaoYan Zhu. 2020c. Recent advances and challenges in task-oriented dialog systems. *Science China Technological Sciences* 63, 10 (2020), 2011–2027.
- Zhuosheng Zhang and Hai Zhao. 2021. Structural Pre-training for Dialogue Comprehension. 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)*. Association for Computational Linguistics, Online, 5134–5145. <https://doi.org/10.18653/v1/2021.acl-long.399>
- Shubin Zhao, Adam Meyers, and Ralph Grishman. 2004. Discriminative Slot Detection Using Kernel Methods. In *COLING 2004: Proceedings of the 20th International Conference on Computational Linguistics*. COLING, Geneva, Switzerland, 757–763. <https://aclanthology.org/C04-1109>
- Weixiang Zhao, Yanyan Zhao, and Bing Qin. 2022. MuCDN: Mutual Conversational Detachment Network for Emotion Recognition in Multi-Party Conversations. In *Proceedings of the 29th International Conference on Computational Linguistics*. International Committee on Computational Linguistics, Gyeongju, Republic of Korea, 7020–7030. <https://aclanthology.org/2022.coling-1.612>
- Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir Radev. 2021. QMSum: A New Benchmark for Query-based Multi-domain Meeting Summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, Online, 5905–5921. <https://doi.org/10.18653/v1/2021.naacl-main.472>
- Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2SQL: Generating Structured Queries from Natural Language using Reinforcement Learning. *CoRR* abs/1709.00103 (2017).
- Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. 2023a. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419* (2023).
- Kangyan Zhou, Shrimai Prabhumoye, and Alan W Black. 2018. A Dataset for Document Grounded Conversations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 708–713. <https://doi.org/10.18653/v1/D18-1076>
- Li Zhou and Kevin Small. 2020. Multi-domain Dialogue State Tracking as Dynamic Knowledge Graph Enhanced Question Answering. *arXiv:1911.06192 [cs.CL]*
- Pei Zhou, Karthik Gopalakrishnan, Behnam Hedayatnia, Seokhwan Kim, Jay Pujara, Xiang Ren, Yang Liu, and Dilek Hakkani-Tur. 2022. Think before you speak: Explicitly generating implicit commonsense knowledge for response generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1237–1252.
- Yunhua Zhou, Jianqiang Yang, Pengyu Wang, and Xipeng Qiu. 2023b. Two Birds One Stone: Dynamic Ensemble for OOD Intent Classification. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 10659–10673.
- Chenguang Zhu, Yang Liu, Jie Mei, and Michael Zeng. 2021. MediaSum: A Large-scale Media Interview Dataset for Dialogue Summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, Online, 5927–5934. <https://doi.org/10.18653/v1/2021.naacl-main.474>