# HR-Agent: A Task-Oriented Dialogue (TOD) LLM Agent Tailored for HR Applications

Weijie Xu<sup>1</sup>, Jay Desai<sup>1</sup>, Fanyou Wu<sup>1</sup>, Josef Valvoda<sup>2</sup>, Srinivasan H. Sengamedu<sup>1</sup>

Amazon

<sup>2</sup>University of Cambridge

weijiexu@amazon.com

#### **Abstract**

Recent LLM (Large Language Models) advancements benefit many fields such as education and finance, but HR has hundreds of repetitive processes, such as access requests, medical claim filing and time-off submissions, which are unaddressed. We relate these tasks to the LLM agent, which has addressed tasks such as writing assisting and customer support. We introduce an innovative HR LLM Agent (HR-Agent) designed to address these challenges by utilizing compact models trained by synthetic data and LLM. The proposed system has low latency, high accuracy, and versatility in supporting multiple tasks. Since conversation data is not sent to an LLM during inference, it preserves confidentiality required in HR-related tasks.

#### 1 Introduction

Recent advances in natural language processing (NLP) have been used for many domains such as Law (Sargeant et al., 2024), Finance (Masson and Paroubek, 2024) and Education (Zhao et al., 2021b). However, many HR processes, such as requesting time off, scheduling meetings, submitting tickets for IT issues, or filing medical claims, remain highly inefficient. Automating these processes could save a significant amount of time that would otherwise be spent on repetitive work. This paper investigates how LLM Agent could facilitate such automation. For a LLM Agent to be useful in the HR domain, it must satisfy the following five requirements:

- (1) It must have a **fast** response time. Employees are less likely to use a chatbot if it cannot complete the task quickly (Hoxmeier and DiCesare, 2000). Research has shown that user satisfaction decreases as response time increases, and a preferable response time should be less than 2 seconds (Shneiderman and Plaisant, 2010).
- (2) The HR Agent must also be **extractive**. When using a TOD to file a medical claim, users

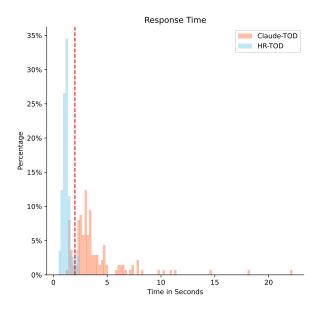


Figure 1: We systematically compared our method to Claude's approach (Zhang et al., 2023) in terms of response time by collecting 40 conversations across four different categories: time off, medical claims, resume creation, and issue ticket handling. The HR-Agent we propose demonstrates significantly faster response times compared to the Claude-based TOD. In fact, our system achieves a response time of less than 2 seconds in 94 percent of cases, while the Claude-based system accomplishes this in only 4 percent of cases. These results highlight the substantial speed advantage of our HR-Agent over the Claude-based solution.

must be able to trust that the system will accurately retrieve the right number.

- (3) Equally important is the **versatility** of the system it must handle the various HR use cases mentioned above.
- (4) Since employee information is highly sensitive, the TOD itself must be **confidential**.
- (5) Finally, the HR Agent must be **HR specific** and perform effectively in HR-relevant tasks as suggested by Xu et al. (2024a)

To complete the task, these systems rely on Dialogue State Tracking (DST) (Rastogi et al., 2020a),

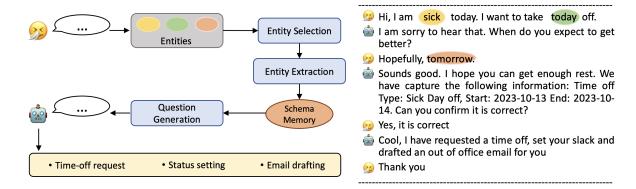


Figure 2: An illustration of the solution. The Entity Selection Model identifies relevant entity. The selected entity is passed to the Entity Extraction model to find the relevant word in the utterance. Based on Schema's memory and previous utterance. The Question Generation model is used to generate the next question. The HR-Agent system then connects to the API to finish the relevant tasks such as drafting email, requesting time off and setting status.

which monitors and predicts user intent and details during conversations.

DST uses schema-guided techniques, categorized as extractive(Rastogi et al., 2020b; Ruan et al., 2020), generative (Feng et al., 2021; Tian et al., 2021), and LLM-based Hudeček and Dušek (2023); Zhang et al. (2023), to track conversation states with natural language explanations. However, extractive and generative methods are limited by training data and poor transfer learning, hindering HR-specific and versatile applications. They therefore do not satisfy points 3 and 5 above. LLM-based approaches are slow and may lack confidentiality, they therefore do not satisfy requirements 1, 2 and 4. In response, we propose an HR-specific LLM Agent (HR-Agent). It is fast in response time, extractive, versatile, confidential, and HR-specific. We summarize our contributions:

- Despite being smaller and faster in speed, as illustrated in Figure 1, modules trained on synthetic datasets achieve superior performance compared to larger language models.
- Based on these modules, we designed an HR-Agent system that significantly improves HR process efficiency.

#### 2 Related Work

Schema-Guided Dialogue (SGD) (Rastogi et al., 2020b) is a dialogue dataset with evolving ontologies, introducing new test set slots and services, emphasizing DST performance and zero-shot generalization. SGD-X (Lee et al., 2022) expands on SGD, presenting five additional schema styles. MultiWOZ (Budzianowski et al., 2020) features

human-human dialogues using a stable ontology. HR-MultiWOZ (Xu et al., 2024a) is aligned with HR related tasks, but it is too small to be used for training. *Collecting data for these datasets is labor-intensive and costly*.

SGD-baseline (Rastogi et al., 2020b), SGP-DST (Ruan et al., 2020), and DS-DST (Zhang et al., 2020) encode utterances and slot schemas jointly to predict relative slots. Multi-Task BERT adopts slot carryover mechanisms, encoding only the preceding system utterance and the current utterance. LUNA (Wang et al., 2022) separately encodes dialogue history, slots, and slot values, learning to predict the correct utterance to condition slot value predictions. However, these methods lack HR specificity and versatility. Seq2Seq-DU (Feng et al., 2021) and AG-DST (Tian et al., 2021) derive states differently, while DaP (Lee et al., 2021) offers two versions, with the latter being slower. D3ST (Zhao et al., 2021a) deciphers the entire dialogue state at once. Although these generative methods, especially when using T5 models, achieve better Joint Goal Accuracy (JGA) in schema-guided dialogue, they suffer from slow response times because of extensive input.

Despite recent research efforts (Li et al., 2023; Hudeček and Dušek, 2023; Zhang et al., 2023), LLM based TOD's performance metrics like BLEU scores and success rates remain low, even with models like Alpaca-LoRA-7B (Taori et al., 2023) and ChatGPT. LLMs are also non-extractive and suffer from high inference costs and latency issues (Yang et al., 2023), making real-world TOD system deployment challenging. For instance, using the GPT-4 8K context model via OpenAI's API costs \$0.03

for every 1K input tokens and \$0.06 for every 1K output tokens.

Gan et al. (2024) introduced a framework using LLM agents for automated resume screening, improving recruitment efficiency through automated extraction and analysis. It achieved high accuracy and F1 scores in sentence classification and effectively protected privacy by excluding personal information while grading and summarizing resumes. However, unlike HR-Agent, it does not specifically address the variety of repetitive HR tasks, such as access requests, medical claim filing, and time-off submissions (Gan et al., 2024).

Enhancing LLM-based task-oriented dialogue systems by incorporating goal-oriented rewards and reinforcement learning techniques was explored by Singh (Singh, 2023). By aligning dialogue responses with predefined goals, the model improves contextual appropriateness and task-specific outputs. This approach uses Multi-WOZ datasets to measure success and refines dialogue policy models through reinforcement learning. Unlike HR-Agent, Singh's work does not focus on the confidentiality and specific needs of HR-related tasks.

#### 3 Methods

When we have an HR task-specific schema, we first convert each entity into a set of questions with the help of HR domain experts, since it can provide details about the entity (Lee et al., 2021). For each turn of the conversation, we employ an entity selection module to choose the entities that can be addressed in that turn. Since every entity selection is applied on every utterance and constantly updates schema memory. We care less about out of utterance response. Subsequently, we format the question and the utterance as an entity extraction task and select the most relevant entity from the utterance. Afterward, we use the question generation module to produce the utterance for HR-Agent. Once we have gathered sufficient schema information, we employ a few APIs to use this structured information for tasks such as drafting emails, creating tickets, and answering questions as illustrated in Figure 2. Our designed HR-Agent system can be used for tasks such as requesting leaves, inquiring about benefits, applying for internal jobs, navigating the onboarding process, requesting training, reporting workplace issues, participating in surveys, and engaging with HR initiatives. It simplifies benefit enrollment, goal setting, exit interviews, safety guidelines, and compliance training. It is confidential because the external LLM is not used for any inference.

#### 3.1 Dataset Generation

We follow Xu et al. (2024a) to design our synthetic data generation process. After generating our dataset, we need to filter out low-quality data. We enhance the extractive nature of our trained model by removing entities that cannot be extracted from the utterance. We validate sample correctness by dividing them into utterance and question pairs, querying FlanT5 XXL (Chung et al., 2022) and Vicuna-13B (Chiang et al., 2023) to ensure question relevance. Irrelevant questions with a similarity score of less than 0.2, as determined by ST5-XXL (Ni et al., 2022), are removed. To address distribution imbalance, we randomize prompt distribution for equal representation among all choices.

The training set was split into three subsets: a 20,000 data point training set per task, a validation set with 5,000 data points per task, and a test set of 500 data points per task. To ensure test set quality and relevance, we implemented a multi-step validation process. First, three human labellers on Mechanical Turk assessed test data correctness, selecting only unanimously agreed-upon data points. Second, Claude-V3 further validated the data, retaining concurred test data. Finally, domain experts verified data relevance to our specific use cases. This process yielded a test set of 290 data points, with 292 for entity selection tasks and 303 for entity extraction tasks. The detailed steps are documented in Appendix E

After generating synthetic data, our next step involves filtering out low-quality data through a multi-step process. We enhance the extractive nature of our trained model by removing entities that cannot be extracted from the utterance. We validate sample correctness by dividing them into utterance and question pairs, querying FlanT5 XXL (Chung et al., 2022) and Vicuna-13B (Chiang et al., 2023) to ensure question relevance. Irrelevant questions with a similarity score of less than 0.2, as determined by ST5-XXL (Ni et al., 2022), are removed. To address distribution imbalance, we randomize prompt distribution for equal representation among all choices.

	FlanT5 XL	Falcon	MPT	FlanT5-Clean	FlanT5-Raw	Claude V3
Size	3B	7B	7B	220 M	220 M	
Precision	0.581	0.948	0.950	0.910	0.777	0.392
Recall	0.881	0.753	0.754	0.832	0.663	0.852
F1 Score	0.663	0.826	0.828	0.856	0.703	0.519
Response time	1.202	1.085	0.6293	0.366	0.384	1.163

Table 1: The performance for entity selection. From left to right, we document the performance of FlanT5 XL, Falcon 7B, MPT 7B, FlanT5 base trained on filtered data, FlanT5 base trained on unfiltered data, ClaudeV3.

	Flan T5 XL	Falcon	MPT	Deberta	Roberta	Flan T5 Trained	Claude V3
Size	3B	7B	7B	135M	125M	220 M	
RougeL	0.786	0.112	0.292	0.729	0.767	0.818	0.793
Response time	0.394	1.083	0.569	0.191	0.088	0.110	1.22

Table 2: The performance for entity extraction. From left to right, we document the performance of FlanT5 XL, Falcon 7B, MPT 7B, Debetra, Roberta, FlanT5 base trained on filtered data, and ClaudeV3.

#### 3.2 Experiments and Model Selection

#### 3.3 Baseline Methods

Falcon (Almazrouei et al., 2023), a Large Language Model (LLM), is designed for task-oriented dialogue systems. It focuses on optimizing efficiency and utility, aiming to provide fast, accurate responses in various conversational scenarios. **MPT** (Team, 2023) is a transformer trained from scratch on 1T tokens of text and code. It is an open source, available for commercial use, and matches the quality of LLaMA-7B **Deberta** (He et al., 2021) Deberta enhances the BERT model with a disentangled attention mechanism for more interpretable attention scores and uses relative position encoding to boost performance in extractive task. Roberta RoBERTa (Liu et al., 2019) is a transformer model pretrained on English data using the Masked Language Modeling technique, masking and predicting 15 percent of input words for bidirectional learning. We use both models finetuned on SQuAD2.0 which performs well on extractive QA. FlanT5 we use the same setup as (Lin et al., 2021b) where we use DDP to train, set up validation loss for early stopping and set max epochs equal to 20.

#### 3.4 Entity Selection

For entity selection, we need to select the relevant entities that could be answered by utterance. We choose FlanT5 as our base model because it contains smaller model size and is explored most in Schema Guided Dialogue literature. Since the entity itself is not informative, we choose to generate questions about the entity because datasets that

are used to train Flan are mostly contain questions. Next, we assess 5 models using the test set. We give the same example for each template and average the result of 10 templates for a pretrained model, such as MPT and Falcon. FlanT5 is trained 5 times and the evaluation result is averaged. The model is trained and evaluated using p3.8xlarge. We do not use an accelerator or ONNX for inference to get a fair comparison of the response time. We benchmark Precision, Recall, F1 Score and Response Time.

Our ideal solution should have high recall but should not output too many entities, which could drastically increase the response time of the TOD system. Decoder-based models perform well in F1 Score but are low on recall. The fine-tuned version of FlanT5 achieves the best F1 score while also achieving the second-best recall. Claude V3 achieves the best recall but tends to select most of the entities and is slow in response. Without data cleaning, FlanT5 tends to choose many first choices, which leads to low recall. We choose FlanT5 trained on filtered data as our entity selection model because it is **fast** and ia **confidential**. Since the training set is **versatile and HR specific**, the trained model also inherits these advantages.

## 3.5 Entity Extraction

Except Falcon and MPT, we also use Deberta and Roberta for the entity extraction benchmark. We also fine-tuned a small FlanT5 model on synthetic dataset(20K). We compare Rouge 1, Rouge L and Response time for each model. Our goal is to select

the model that is extractive and has low response time.

Our ideal solution has a low response time and a high Rouge1 score. Claude does not achieve a good Rouge1 score and takes much longer to train. It fails to use words from the utterance. Based on our selected models, we find that extractive models such as Roberta and Deberta perform well in the extraction task, while decoder-only models such as MPT and Falcon perform poorly. Flan T5 performed the best. Flan T5, trained on a filtered synthetic dataset, further improves the Rouge1 performance by another 5 percent since it is more **extractive**. The trained model is also very **fast** compared to larger models and is comparable to Roberta-based models. Thus, we choose trained FlanT5 as our model for entity extraction.

#### 3.6 Question Generation

Infusing more empathy into the responses is paramount, as it fosters a more human-centric and relatable interaction experience. To achieve that, we use the utterance and the next question to ask as input. We provided this input to Claude-V3 and asked Claude-V3 to rewrite a concise response with empathy, as illustrated in Table 12. To evaluate the performance of Claude3-V3 on rephrasing response with empathy, we want to benchmark against the response generated in the HR-Multiwoz dataset. To achieve that, we first identify response that may reflect negative feelings from the user. We use DistillBert based on a case model finetuned on dataset SST-2 as our sentiment classifier. Since finetuned DistillBert is uncalibrated. Only sentences with a score above 0.998 and classified as negative are selected. We have collected 638 responses from HR-Multiwoz. We then manually rewrite those questions to basic question and leverage Claude3-V3 to rewrite it.

To evaluate the effectiveness and human preference for the responses generated by our HR-Agent, we conducted a preference study involving human labelers using Amazon SageMaker Ground Truth (GT). We compared responses generated by HR-Agent (denoted as Response A) with those from the HR-MultiWOZ dataset (denoted as Response B).

Labelers were presented with dialogue scenarios and asked to choose which response they preferred. The user interface used for this evaluation is shown in Figure 4. Each response was evaluated by three labelers, and the cost per label was \$0.0012. The results of this study are illustrated in Figure 3.

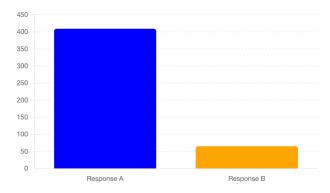


Figure 3: Labeler preferences for responses from HR-Agent (A) and HR-MultiWOZ (B).

Our study revealed a significant preference for responses generated by HR-Agent. Out of the total responses evaluated, Response A was preferred 409 times, while Response B was preferred only 65 times. This strong preference for Response A suggests that the HR-Agent's responses were more aligned with the labelers' expectations and needs.

The HR-Agent's use of empathy in responses contributed to its preference. For instance, when an employee inquired about scheduling a medical appointment, the HR-Agent's response, "I understand that medical emergencies can be stressful. Please share the incident details so we can assist you," was perceived as more empathetic and supportive compared to a more straightforward response from the HR-MultiWOZ dataset.

These findings underscore the value of tailoring dialogue systems to specific domains and incorporating elements such as empathy to improve user satisfaction and engagement.

## 3.7 System Evaluation

To demonstrate that our designed system performs best in HR related tasks. We use HR-Multiwoz as our evaluation set. We compared our method with TransferQA (Lin et al., 2021a) and SGP-TOD (Zhang et al., 2023). ZST adapts a sequence-to-sequence model (T5) pre-trained on large QA datasets to track dialogue states by formulating slot-filling tasks as QA problems. We replace T5 by Deberta and Roberta, suggested by HR Multiwoz benchmark leaderboard. SGP-TOD uses prompt engineer to guide the dialogue and extract the relevant fields.

Joint Goal Accuracy (JGA) and Average Goal Accuracy (AGA) are used to evaluate our models

and baselines. For JGA, the model outputs are only counted as correct when all of the predicted values exactly match the oracle values. AGA is the average accuracy of the active slots in each turn.

As you can see in Table 3, our proposed method achieves the best performance in HR-Multiwoz dataset in terms of JGA and AGA compared to the state-of-the-art methods. This means that the HR Agent works well in HR related SGD tasks.

#### 3.8 Prompt

Prompt we used in the experiment for MPT, Falcon and FlanT5 as you can see in Table 9 and Table 10.

Prompt for Claude-V3 is "The answer is very short and always less than 2 words. Put the answer in <answer></answer> XML tags." for entity extraction and "The answer contains always 2 to 5 choices. Put the answer in <answer></answer> XML tags." for entity selection. We use max tokens to sample be 10, temperature be 0.2 and stop sequence be </answer>

	JGA	AGA
TrasnferQA Roberta	8.65	26.62
TrasnferQA Deberta	18.89	55.61
SGP-TOD	10.74	54.77
HR-Agent	20.47	55.38

Table 3: The overall performance of HR Agent on HR-Multiwoz. It achieves better performance in JGA as well as AZA.

## 3.9 Other considerations

An essential part of this process involves robust fact-checking to validate generated answers, enhancing reliability and credibility. For instance, when identifying a medical provider, the system confirms their details in the database before processing a medical claim. Users can also track the schema throughout the process, ensuring accountability and trustworthiness. Our dataset connects with existing databases to retrieve relevant information and only asks when necessary. After collecting data, users confirm the schema and the task, and we connect to relevant APIs to complete the task. Before putting information into the system, it is important to make information accurate. For example, the time off day collected is today is not informative for the system to track the exact day. To resolve these inaccuracies, we call Claude to auto complete the inaccurate information. For example,

the Claude can change tomorrow to November 1st and 98121 to Seattle WA. This can help the system capture the accurate information. We then confirm it with the user before calling APIs to use the structure information. The system can hugely reduce user's inputs to get required information.

#### 4 Conclusion

In improving the efficiency of corporate employees, we introduce HR-Agent, a rapid, extractive, and confidential dialogue system tailored to HR needs. Our work encompasses several key contributions: identifying challenges in adapting TOD to HR contexts, devising a swift and domain-specific data generation approach, demonstrating the superior performance of smaller, faster-trained modules over larger models, and ultimately delivering HR-Agent as a solution that markedly enhances HR process efficiency.

#### 5 Ethical Statement

This system is a prototype and has not been deployed in production.

The deployment of AI applications in the HR space necessitates careful consideration of ethical issues related to safety, privacy, and bias. There is a possibility that, in attempting to assist, AI may cause more harm than benefit. In response, in collaboration with user experience researcher, security reviewer and HR professionals, we have suggested the following steps for developers who plan to use the HR Agent in order to minimize the risks of harm.

Informed Consent from Users: In the pilot phase, informed consent was obtained from employees using the service. They were made aware that they would be interacting with an AI-based chatbot designed to expedite task completion. It was also communicated that some extracted information might be inaccurate, and users had to verify the correctness of this information before utilizing it for subsequent tasks. To facilitate this, HR Agent should present the collected information at the end of each conversation and request employee confirmation of its accuracy. Developers should ask the employees to participate in a survey after using HR Agent to understand the effectiveness and helpfulness of HR Agent. They should make it clear that HR Agent is used to fill the relevant task and not for anything else. They should provide relevant HR Business Partner contact information as headers to

employees when they interact with HR Agent.

Guardrail:HR agent should pose clarifying questions when unable to extract relevant entities. The conversation should terminate if clarifying questions for the same entity are repeated more than three times. The developers can use sentiment analysis model to monitor the sentiment of employees every four user responses, and the conversation is ended if the negative sentiment score exceeds 0.5. Developers should provide task relevant internal wiki when they end the conversations. For some applications, guardrails are also in place for ranges of money and time that exceed specific thresholds. It's important to note that HR Agent has limited applications at present, and they should update these guardrails accordingly.

Privacy: In the system, information provided by the employees should be kept confidential in a record separate from the employee's general personnel file. Employees have the right to use this record or input this record in other systems or give developers rights to use this record for analysis. This policy follows the Americans with Disabilities Act (ADA and the Genetic Information Nondiscrimination Act (GINA). Please note that all the models in the system should be trained on synthetic data. Developers should not use any real employee data to train the model. Developers also have to make sure the data in the system is in compliance with rigorous internal infoSec policies and standards. For example, security testing includes examining application logs to detect any data leakage into logs.

Negative Examples/Potential Bias: To mitigate potential biases in generative models, developers should employ an extractive approach. None the less, the effectiveness of extractions could vary with the employee's language fluency. This variation could potentially lead to inefficiencies in the Task-Oriented Dialogue (TOD) system for nonnative English speakers. Additionally, the system is prone to errors when extracting multiple entities of the same type, such as time and money from one response. Efforts are underway to understand and address these issues.

Developers should do threat modeling, security testing, penetration test assessment of the system.

#### 6 Limitations and Risks

**Limitations** We have not explore other models for training such as Deberta or LLama2. We have not

evaluated the model on real data.

**Risks** Because of some privacy concerns, we have not discussed the details of the whole architecture. The risk of deploying this architecture is that the performance may drift a bit on real data since test set is synthetic.

#### 7 Future Work

Looking ahead, potential advancements include the capability to generate resumes and emails, interface with various APIs, proficiently answer queries, and identify pertinent tickets, further enhancing the utility and efficiency of the system. To enhance the performance and capabilities of the model, it is crucial to amass a substantial dataset for the fine-tuning of a more extensive model, enabling it to execute multiple tasks concurrently. Essential prerequisites for this advancement encompass the ability to handle longer sequences and the provision of clear, upfront instructions for entity extraction and task execution. Furthermore, illustrating the connections to other agents is vital, offering a comprehensive and interconnected approach to task management and execution, thereby bolstering the overall efficiency and effectiveness of the model. We could leverage topic modeling (Xu et al., 2023b,a, 2024b) to better understand and segment use cases. We could leverage differential private synthetic data generation mechanism (Madl et al., 2023; Xu et al., 2023c) to avoid potential privacy issues.

#### References

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance. *None*.

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2020. Multiwoz – a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality.

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.
- Yue Feng, Yang Wang, and Hang Li. 2021. A sequence-to-sequence approach to 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 1: Long Papers), pages 1714–1725, Online. Association for Computational Linguistics.
- Chengguang Gan, Qinghao Zhang, and Tatsunori Mori. 2024. Application of llm agents in recruitment: A novel framework for resume screening.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2023. Textbooks are all you need.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention.
- John A Hoxmeier and Chris DiCesare. 2000. System response time and user satisfaction: An experimental study of browser-based applications. *AMCIS*.
- Vojtěch Hudeček and Ondřej Dušek. 2023. Are llms all you need for task-oriented dialogue? *arXiv preprint arXiv:2304.06556*.
- Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021. Dialogue state tracking with a language model using schema-driven prompting. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4937–4949, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Harrison Lee, Raghav Gupta, Abhinav Rastogi, Yuan Cao, Bin Zhang, and Yonghui Wu. 2022. SGD-x: A benchmark for robust generalization in schemaguided dialogue systems. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):10938–10946.
- Zekun Li, Baolin Peng, Pengcheng He, Michel Galley, Jianfeng Gao, and Xifeng Yan. 2023. Guiding large language models via directional stimulus prompting. *arXiv preprint arXiv:2302.11520*.

- Zhaojiang Lin, Bing Liu, Andrea Madotto, Seungwhan Moon, Zhenpeng Zhou, Paul Crook, Zhiguang Wang, Zhou Yu, Eunjoon Cho, Rajen Subba, and Pascale Fung. 2021a. Zero-shot dialogue state tracking via cross-task transfer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7890–7900, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zhaojiang Lin, Bing Liu, Seungwhan Moon, Paul A Crook, Zhenpeng Zhou, Zhiguang Wang, Zhou Yu, Andrea Madotto, Eunjoon Cho, and Rajen Subba. 2021b. Leveraging slot descriptions for zero-shot cross-domain dialogue statetracking. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5640–5648.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach.
- Tamas Madl, Weijie Xu, Olivia Choudhury, and Matthew Howard. 2023. Approximate, adapt, anonymize (3a): a framework for privacy preserving training data release for machine learning.
- Corentin Masson and Patrick Paroubek. 2024. Evaluating topic model on asymmetric and multi-domain financial corpus. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 6515–6529, Torino, Italia. ELRA and ICCL.
- Jianmo Ni, Gustavo Hernandez Abrego, Noah Constant, Ji Ma, Keith Hall, Daniel Cer, and Yinfei Yang. 2022. Sentence-t5: Scalable sentence encoders from pretrained text-to-text models. In *Findings of the As*sociation for Computational Linguistics: ACL 2022, pages 1864–1874, Dublin, Ireland. Association for Computational Linguistics.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020a. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. In *Proceedings of the AAAI conference on artificial intelligence*, 05, pages 8689–8696.
- Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020b. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset.
- Yu-Ping Ruan, Zhen-Hua Ling, Jia-Chen Gu, and Quan Liu. 2020. Fine-tuning bert for schema-guided zero-shot dialogue state tracking. *arXiv preprint arXiv:2002.00181*.
- Holli Sargeant, Ahmed Izzidien, and Felix Steffek. 2024. Topic modelling case law using a large language

- model and a new taxonomy for uk law: Ai insights into summary judgment.
- Ben Shneiderman and Catherine Plaisant. 2010. *Designing the user interface: strategies for effective human-computer interaction*. Addison-Wesley Computing.
- Vasudha Singh. 2023. Exploring the role of large language model (llm)-based chatbots for human resources. Ph.D. thesis, UT Austin.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models.* https://crfm. stanford. edu/2023/03/13/alpaca. html, 3(6):7.
- MosaicML NLP Team. 2023. Introducing mpt-7b: A new standard for open-source, commercially usable llms. *none*. Accessed: 2023-05-05.
- Xin Tian, Liankai Huang, Yingzhan Lin, Siqi Bao, Huang He, Yunyi Yang, Hua Wu, Fan Wang, and Shuqi Sun. 2021. Amendable generation for dialogue state tracking. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, pages 80–92, Online. Association for Computational Linguistics.
- Yifan Wang, Jing Zhao, Junwei Bao, Chaoqun Duan, Youzheng Wu, and Xiaodong He. 2022. LUNA: Learning slot-turn alignment for dialogue state tracking. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3319–3328, Seattle, United States. Association for Computational Linguistics.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions.
- Weijie Xu, Wenxiang Hu, Fanyou Wu, and Srinivasan Sengamedu. 2023a. DeTiME: Diffusion-enhanced topic modeling using encoder-decoder based LLM. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9040–9057, Singapore. Association for Computational Linguistics.
- Weijie Xu, Zicheng Huang, Wenxiang Hu, Xi Fang, Rajesh Cherukuri, Naumaan Nayyar, Lorenzo Malandri, and Srinivasan Sengamedu. 2024a. HR-MultiWOZ: A task oriented dialogue (TOD) dataset for HR LLM agent. In *Proceedings of the First Workshop on Natural Language Processing for Human Resources (NLP4HR 2024)*, pages 59–72, St. Julian's, Malta. Association for Computational Linguistics.
- Weijie Xu, Xiaoyu Jiang, Jay Desai, Bin Han, Fuqin Yan, and Francis Iannacci. 2024b. Kdstm: Neural semi-supervised topic modeling with knowledge distillation.

- Weijie Xu, Xiaoyu Jiang, Srinivasan Sengamedu Hanumantha Rao, Francis Iannacci, and Jinjin Zhao. 2023b. vONTSS: vMF based semi-supervised neural topic modeling with optimal transport. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4433–4457, Toronto, Canada. Association for Computational Linguistics.
- Weijie Xu, Jinjin Zhao, Francis Iannacci, and Bo Wang. 2023c. Ffpdg: Fast, fair and private data generation.
- Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *arXiv preprint arXiv:2304.13712*.
- Jianguo Zhang, Kazuma Hashimoto, Chien-Sheng Wu, Yao Wang, Philip Yu, Richard Socher, and Caiming Xiong. 2020. Find or classify? dual strategy for slot-value predictions on multi-domain dialog state tracking. In *Proceedings of the Ninth Joint Conference on Lexical and Computational Semantics*, pages 154–167, Barcelona, Spain (Online). Association for Computational Linguistics.
- Xiaoying Zhang, Baolin Peng, Kun Li, Jingyan Zhou, and Helen Meng. 2023. Sgp-tod: Building task bots effortlessly via schema-guided llm prompting. *arXiv* preprint arXiv:2305.09067.
- Jeffrey Zhao, Mahdis Mahdieh, Ye Zhang, Yuan Cao, and Yonghui Wu. 2021a. Effective sequence-to-sequence dialogue state tracking. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jinjin Zhao, Kim Larson, Weijie Xu, Neelesh Gattani, and Candace Thille. 2021b. Targeted feedback generation for constructed-response questions. In *AAAI* 2021 Workshop on AI Education.

#### **A** Synthetic Example

We first try to generate data that is similar to these datasets where we have the conversation and schema generated by the conversation. We can use two claude agent where one acts as an employee and another acts as a TOD system. We then use the third claude agent to track the schema. To generate the conversation, the approach needs multiple api call per sample which is costly.

There are also many other problems. After a few rounds of generation, the result becomes repetitive even with dramatically different prompts. The language model tends to follow the most probable or common paths given their training data and their priors, and they lack the creativity or the incentive to explore alternative. The conversation can also

be irrelevant. Thus, the system trained on these data is less likely to be extractive, versatile and HR specific.

## **B** Designed Prompt

Here is our designed prompt where *schema*, *number1* and *number2* change every time. We use output1 and output2 to capture answers for two tasks. We use Claude-V3 model. We use max toke sample equal to 4096. We keep temperature equal 1, top k equal 1 and top p equal 0.6 to have good output length, capture multiple samples at a time and have enough diversity.

## **C** Output Example

Here is an example of the output using Claude V3: We have also tried cheaper option such as Amazon Titan and Jurassic-2 but we fail to make them follow instruction. Here is an example of titan:

## D Experiment Setup Appendix

Various tables related to Experiment Setup are shown in this section.

#### E Data Validation

#### E.1 LLM use for Validation data

We use the following string to put into Vicuna and FlanT5 XXL to validate the data: Question: *question* Text: *text* Answer: *answer* Does tha Answer answer the Question based on Text? The answer could be yes or no" We only select data with answer contains "yes" For Claude, we input 50 in a group and ask model to find line number with the answer equal no. We save inference cost this way.

#### E.2 Mechanical Turk

We give 0.012 dollar per human labeller per task. We do not enable automated data labeling.

Speaker	Dialogue and Schema
Employee	I'd like to take next Thursday and Friday off.
Chatbot	Okay, which dates would you like off? → Ask repetitive questions
Schema	timeOffStartDate: ? → Not able to capture the right information
	timeOffEndDate: ?
Employee	Next Thursday the 19th and Friday the 20th.
Chatbot	Got it. Can you confirm the reason for the time off?
Schema	timeOffDate: October 19-20, $2023 \rightarrow$ Suddenly change the schema and not extractive
Employee	It's for a doctor's appointment. My name is John Smith.
Chatbot	Okay, I have a request from John Smith for October 19th-20th for a doctor's appointment.
Schema	timeOffDate: October 19-20, 2023
	reason: Doctor's appointment
	name: John Smith

Table 4: This is a generated example by Claude. The schema keeps changing. The schema fails to capture correct information. The chatbot also asked the information that has been mentioned by user already. The chatbot is also not extractive

| Institution | Section |

Figure 4: User interface used for the human preference study.

**Human:** You are asked to come up with a set of 20 diverse scenarios. The input is a user response and a list of questions that the user response could answer about the user. The instruction is to select the right set of questions that could be answered by user input and give an answer for each selected question.

#### **Requirements:**

- 1. Try not to repeat the verb or cases for each input to maximize diversity.
- 2. The language used for the user response should be diverse.
- 3. The type of input should be diverse. The user response should include diverse types of tasks like schema.
- 4. You should generate an appropriate List of Questions for the input. It should involve realistic data and should not contain simple placeholders.
- 5. The list of questions should start with a choice name such as a, b, c, ...
- 6. Output1 should be the choice that could be answered by the user response.
- 7. Output1 should be all correct choice names such as a, b, c.
- 8. The list of questions should contain *number1* questions.
- 9. Output1 contains number2 choices.
- 10. Every question has an equal chance to be the correct answer. The answer should not always contain 'a'.
- 11. Output2 should be the correct answer for each question.
- 12. The answer from Output2 is extracted from the User Response.
- 13. The answer of Output2 cannot be 'yes' or 'no'.
- 14. The scenario is to help corporate extract information from their employees.
- 15. Try not to ask similar questions to maximize diversity.
- 16. Output2 does not contain choice symbols like a, b, c.
- 17. The answer from Output2 can only be extracted from the User Response.
- 18. Output2 uses a comma to separate the answers.
- 19. Two empty lines between each case.

#### **Example:**

User Response: I am taking next Monday off as a vacation day.

#### **List of Questions:**

- a. When is the requested time off?
- b. What action does the user want the recipient to take?
- c. What process has the user completed?
- d. What type of time off is being requested?

Output1: a, d

Output2: next Monday, vacation day

Table 5: This is an example prompt for claude for generation

**User Response:** I would like to schedule a doctor's appointment for next Tuesday at 2pm to get a physical exam.

## **List of Questions:**

- a. What type of appointment does the user want to schedule?
- b. When does the user want to schedule the appointment?
- c. What time does the user want the appointment?
- d. What is the purpose of the appointment?
- e. What action does the user want the recipient to take?
- f. On what date is the user requesting the appointment?
- g. Does the user provide the date for the requested appointment?
- h. Does the user provide the time for the requested appointment?

Output1: a, b, c, d

Output2: schedule, next Tuesday, 2pm, physical exam

Table 6: The example perfectly follow our designed prompts using Claude.

Bot: Great! I can help you with that. What kind of training session would you like to schedule?

User: I need to schedule a meeting with John.

**Bot:** Great! When would you like to have the meeting with John?

## C.1 Scenario: An employee is trying to get a recommendation from a colleague who is leaving the company.

**User:** I want to get a recommendation from my colleague.

**Bot:** Sure, I can help you with that. Would you like a template for a recommendation letter?

#### C.2 Scenario: I want to request a day off

**Bot:** Can you please confirm your name?

## C.3 Scenario: John is a contractor. He is required to submit a timesheet every week.

- 1. Which type of document is John submitting?
  - a. Expense Report
  - b. Meeting Scheduling
  - c. Timesheet Submission

#### 1. John is submitting a Timesheet Submission.

- 2. Which system is John using to submit his timesheet?
  - a. SharePoint
  - b. Web-based
  - c. Email

## 1. John is submitting his timesheet through a Web-based system.

- 3. What is John's purpose for submitting a timesheet?
  - a. To track his hours worked
  - b. To request reimbursement for expenses

#### 1. To track his hours worked.

User: I would like to request a day off.

Table 7: Example generated by Amazon Titan

Here is an example of Jurassic-2 **User Response:** "I would like to inquire about my employee benefits."

## **List of Questions:**

- 1. What type of benefits would you like to inquire about? (e.g., health insurance, retirement plans, etc.)
- 2. Are you currently employed by the company?
- 3. How long have you been with the company?
- 4. Are you eligible for benefits? (e.g., full-time employee, part-time employee, etc.)
- 5. Are you the primary beneficiary of any benefits? (e.g., spouse, children, etc.)
- 6. Do you have any specific questions or concerns about your benefits?
- 7. How would you prefer to receive your benefits information? (e.g., email, letter, etc.)

## Output 1:

- a) health insurance
- b) retirement plans
- c) life insurance
- d) dental insurance
- e) disability insurance
- f) vacation and sick leave

## Output 2:

- 1. health insurance
- 2. yes
- 3. full-time employee
- 4. yes
- 5. yes
- 6. no
- 7. email

Table 8: Example generated by Jurassic-2

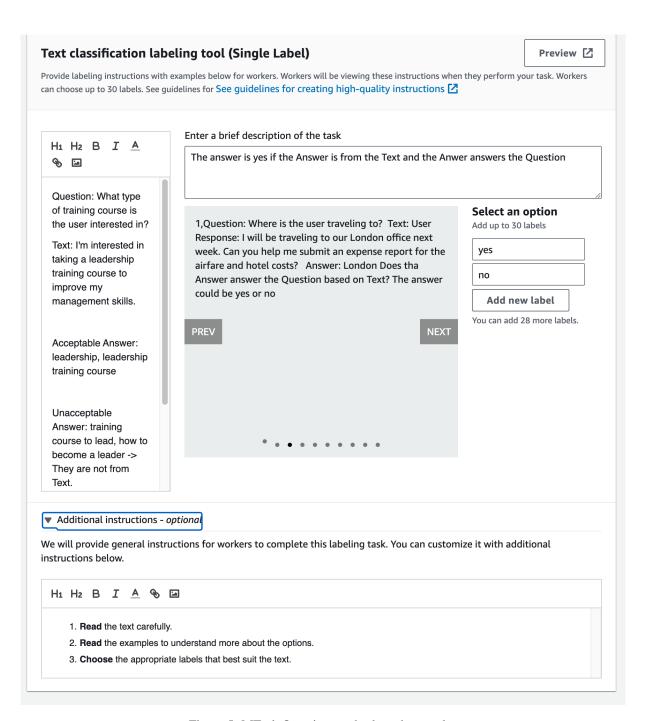


Figure 5: MTurk Questions and selected examples

No.	Sentence
1	Could you retrieve the answer to the Question from the Text?
2	Can you pull out the response to the Question within the Text?
3	Would you mind extracting the reply to the Question from the Text?
4	I'd like you to get the answer to the Question from the Text.
5	Could you find the solution to the Question in the Text?
6	Please identify the answer to the Question in the Text.
7	Can you locate the response to the Question from the Text?
8	I'd appreciate if you could extract the answer to the Question from the Text.
9	Would it be possible to get the reply to the Question from the Text?
10	Please search for the answer to the Question within the Text.

Table 9: Table format prompt for entity extraction

No.	Sentence
1	Could you identify the appropriate question that Text can answer?
2	Please find the correct question for which the Text provides an answer.
3	Can you determine the suitable question that can be resolved using the Text?
4	I'd like you to pinpoint the right question that the Text can address.
5	Please locate the question that aligns with the Text's answer.
6	Could you discern the fitting question that the Text can respond to?
7	I'd appreciate if you could determine the exact question that can be answered using the Text.
8	Can you select the question that the Text can satisfactorily answer?
9	Would it be possible to identify the question that matches the Text's answer?
10	Please deduce the right question that aligns with the Text's response.

Table 10: Table format prompt for entity selection

## F Question Generation Data

The following gives an example to illustrate the synthetic sample.

Table 11: User Input, Next Question, and Empathetic Response

User Input	Next Question	<b>Empathetic Response</b>
"I had a medical emergency last week and need to file a claim."	"Can you please provide details about the incident?"	"I understand that medical emergencies can be stressful. Please share the incident details so we can assist you."
"I've received a medical bill that seems incorrect."	" Could you please provide me with the details of the bill?"	"Dealing with medical bills can be confusing. Please share the bill details, and we'll look into it."
"I'm getting a '404 Not Found' error on my website. What should I do?"	"Can you provide more context about when and where it occurs?"	"Website errors can be frustrat- ing. Please tell me more, so I can assist you effectively."
"My code won't compile, and I don't understand the error message."	"Could you share the error message and a snippet of your code for better assistance?"	"Coding errors can be challenging. Please provide the error message and your code, and we'll figure it out together."

Table 12: Generated examples for training a question generation model with empathy

#### **G** HR Use Case

HR-Agent systems are used to request leaves, inquire about benefits, retrieve payroll details, apply for internal jobs, navigate the onboarding process, schedule performance reviews, request training, report workplace issues, access policies, participate in surveys, and engage with HR initiatives. It simplifies benefit enrollment, goal setting, exit interviews, safety guidelines, and compliance training.

## **H** Data Generation Process

While there are many Schema-Guided Dialogue datasets discussed in the previous section, these datasets are not HR-specific. Therefore, we need to generate synthetic data for creating a TOD system. We employ Claude for generating synthetic data, which is more cost-effective than GPT-4 in terms of API costs and easily integrates with the AWS Ecosystem. In our search for a cost-effective data generation option, we attempted to use Claude to generate both conversations and schemas. However, after trying various prompts, we found that the generated schema does not always contain entities from the conversation (non-extractive), fails to capture correct information, and remains inconsistent. The chatbot's questions can also become redundant. We provide an example and discuss other options and problems in Appendix A. These issues collectively make the synthetic data challenging to

As mentioned in the previous section, most models have performed two tasks. The first task is to select the relevant entity, and the second task is to identify the correct answer for that entity (Zhang et al., 2020). Therefore, we considered generating data that can be used for both of these cases and using separate models for each case. We transformed the synthetic data into a different format, where we asked the model to provide the following: utterance, some questions from the same domain, relevant questions that could be asked based on the utterance, and the extracted entity for each answer.

We follow (Taori et al., 2023; Wang et al., 2023; Gunasekar et al., 2023) to create prompts. We use a batch process by explicitly asking the model to produce multiple samples. We increase the diversity of papers by randomly selecting a domain from the list of the following domains (See Appendix G for all HR use cases), choosing a random number of questions and a random number of answers. Thus, our proposed model can cover all these use

cases and are more versatile. We also provide a few examples and randomly select one to put into prompt for each generation. Our example prompt is shared in Table 5, and we present an example of generated data in Table 6. We have also tried other model options but none of them work as illustrated in Appendix C