

Building Task Bots with Self-learning for Enhanced Adaptability, Extensibility, and Factuality

ZHANG, Xiaoying

A Thesis Submitted in Partial Fulfillment
of the Requirements for the Degree of

-
in

Systems Engineering and Engineering Management

The Chinese University of Hong Kong

2025, Preprint

Thesis Assessment Committee

Abstract of thesis entitled:

Building Task Bots with Self-learning for Enhanced Adaptability, Extensibility,

and Factuality

Submitted by ZHANG, Xiaoying

for the degree of -

at The Chinese University of Hong Kong in 2025, Preprint

Developing adaptable, extensible, and accurate task bots with minimal or zero human intervention is a significant challenge in dialog research. This thesis examines the obstacles and potential solutions for creating such bots, focusing on innovative techniques that enable bots to learn and adapt autonomously in constantly changing environments.

End-to-end task bots, typically built using a static and limited corpus, face difficulties when deployed online due to three primary factors tied to this limitation. First, they might confront queries featuring unexpected linguistic patterns or slot values (i.e., unseen user behaviors). Second, they could potentially face requirements for new functions or tasks (i.e., task definition extensions). Third, even when equipped with relevant knowledge, these bots may produce responses that appear plausible but are actually incorrect (i.e., “hallucinations”). Addressing these challenges is vital for enhancing task bots’ performance and reliability in real-world settings.

To tackle unseen user behaviors, we introduce the Self-Learning Agent (SL-AGENT),

a self-learning framework that learns from real user interactions. With an integrated reward model for predicting response quality, this framework enables task bots to adapt to new user behaviors by learning from post-deployment, unlabeled human-bot dialogues through reinforcement learning. For task definition extensions, we present Schema-Guided Prompting for Task-Oriented Dialogue systems (SGP-TOD), which allows flexible prototyping of dialogue systems for new tasks using large language models (LLMs) and simple task schema modifications. By employing predefined task schemas, i.e., belief instructions and dialogue policies, we guide fixed LLMs to generate appropriate responses for new tasks without requiring training data. To mitigate hallucinations, we explore Self-Alignment for Factuality, leveraging an LLM’s self-evaluation capability to provide training signals that guide the model towards factuality. We incorporate a self-evaluation component that prompts an LLM to verify the factuality of its generated responses based on its internal knowledge, and these responses are then used to refine the model. By implementing these strategies, we lay the groundwork for guiding task bots towards adaptability, extensibility, and factuality with minimal or no human intervention.

Contents

Abstract	i
List of Figures	vi
List of Tables	xi
List of Terms	xiv
1 Introduction	1
1.1 End-to-End Task-Oriented Dialog Systems	2
1.2 Motivation	6
1.3 Thesis Outline and Contributions	9
2 Literature Review	11
2.1 Neural Approaches to Building Task Bots	11
2.1.1 Problem Formulation	12
2.1.2 Training Paradigm Shifts	14
2.2 Adaptations to Unseen User Behaviors	20
2.2.1 Learning with Verbal Human Corrections	21
2.2.2 Learning with Numerical Human Feedback	24
2.3 Extensions to New Tasks	26
2.3.1 Fine-tuning for Effortless Extension	28
2.3.2 Prompting for Effortless Extension	31
2.4 Mitigating Hallucinations	33
2.4.1 Post-hoc Corrections	34
2.4.2 Inference-time Interventions	36
2.4.3 Alignment Training	37
2.5 Chapter Summary	42
3 Self-Learning for Adaptability	44
3.1 The Importance of Adaptability	45

3.2	SL-Agent	48
3.2.1	Overview	48
3.2.2	Dialog Model	49
3.2.3	Reward Model	50
3.2.4	Refine with Reinforcement Learning	53
3.2.5	Minimum Annotations via Machine Teaching	55
3.3	Experiments	57
3.3.1	Setup	57
3.3.2	Results of Unseen User Behaviors	60
3.3.3	Results of Task Definition Extensions	63
3.3.4	Continual Policy Improvement	64
3.4	In-Depth Analyses	65
3.4.1	Impact of Various PLMs and Training Objectives on Reward Models	65
3.4.2	Qualitative Analysis	66
3.4.3	Interactive Human Evaluation	68
3.5	Chapter Summary	69
4	Schema-Guided LLM Prompting for Extensibility	70
4.1	The Necessity of Extensibility	71
4.2	Schema-Guided Prompting (SGP)-TOD	73
4.2.1	Overview	73
4.2.2	LLM	74
4.2.3	DST Prompter	74
4.2.4	Policy Prompter	76
4.3	Experiments	79
4.3.1	Setup	79
4.3.2	Evaluation on Multiwoz	83
4.3.3	Evaluation on RADDLE	85
4.3.4	Evaluation on STAR	87
4.3.5	Evaluation on Domain Extension	89
4.4	In-Depth Analyses	91
4.4.1	Ablation Study	92
4.4.2	Qualitative Analysis	93

4.4.3	Interactive Human Evaluation	94
4.5	Chapter Summary	95
5	Self-Alignment for Factuality	101
5.1	The Importance of Maintaining Factuality	102
5.2	Self-Alignment for Factuality	104
5.2.1	Overview	104
5.2.2	Factuality Self-Evaluation	106
5.2.3	Alignment Tuning with DPO	109
5.3	Experiments	109
5.3.1	Setup	110
5.3.2	Main Results	114
5.3.3	Self-Alignment with Varying Factuality Estimation Methods .	117
5.3.4	Pairwise Evaluation	117
5.4	In-Depth Analysis of SELF-EVAL	118
5.5	In-Depth Analyses	122
5.5.1	Qualitative Analysis	122
5.5.2	Analysis of Factual Errors after Self-Alignment	123
5.5.3	Effects of Removing Duplicates for SK-TUNING	123
5.6	Chapter Summary	124
6	Conclusions and Future Directions	129
6.1	Thesis Summary and Contributions	129
6.2	Future Directions	131
6.2.1	Exploring a neural-symbolic self-learning framework	131
6.2.2	Automatically detecting bots' knowledge boundaries	132
6.2.3	Empowering bots with updated parametric knowledge using minimal human effort	133
6.2.4	Exploring the internal mechanisms of self-improvement . . .	133
6.2.5	Navigating the synergy of bots' self-improvement and human oversight	134
Bibliography		136

List of Figures

1.1	Architectures of two end-to-end task-oriented dialogue systems: (<i>i</i>) a modular end-to-end ToD system, which comprises four neural-model-based dialogue modules, namely NLU, DST, POL and NLG modules (in the upper part), and (<i>ii</i>) a fully end-to-end ToD system, which utilizes a single neural model to subsume all four modules into one (in the lower part).	3
2.1	An example of a dialog turn sequence at turn t , consisting of the dialogue history up to the current turn s_t , annotated with the belief state b_t , database state (i.e., retrieved database entry) c_t , and system response r_t	13
2.2	Evolution of training paradigms for developing end-to-end task bots: (<i>i</i>) standard training, (<i>ii</i>) pre-training then fine-tuning, (<i>iii</i>) pre-training then prompting.	15
2.3	An illustrative example of a dialog model that employs an auto-regressive pre-trained language model as its backbone to parameterize the sequential dialog pipeline, as depicted in Figure 1.1 (adapted from Peng et al. [118]). The model takes the dialog turn sequence in Figure 2.1 as input.	17
2.4	Illustration of proactive learning in a self-feeding chatbot. The model continuously estimates user satisfaction during dialogue interactions to identify appropriate moments to request feedback. Based on user responses, it extracts new training examples for two distinct tasks: DIALOGUE (from satisfied user responses) and FEEDBACK (from user-provided corrections). These examples are then used to fine-tune the model, thereby enhancing its overall dialogue capabilities. Cited from [49].	22

2.5	Illustration of the machine teaching process [146], cited from [117]. Dialog authors inspect and correct representative failed human-bot interaction logs (a) to create new training dialogues. This involves correcting belief states (b) and selecting, inserting, or modifying inappropriate responses, typically using action templates (c).	23
2.6	Dialog policy optimization in an RL loop, where the interaction between a task bot and a user parallels the interaction between an RL agent and its environment.	24
2.7	An example of a task schema from the MultiWOZ dataset [12] in the restaurant domain, featuring the ontology on the left and the dialog flow on the right.	28
2.8	An overview of the ANYTOD system, cited from Zhao et al. [223]. A LM conducts zero-shot state and action tracking with respect to a provided schema, abstracting it into a sequence of symbols. A program that executes the dialog policy then recommends which actions to take based on the states sequence, the LM then chooses a single final action and generating a response.	30
2.9	Illustration of the prompting paradigm in zero-shot, one-shot, and few-shot settings (from left to right) within a task-oriented dialog system, using zero, one, or a few (k , $k > 1$) dialog turn examples as task-specific exemplars, respectively.	32
2.10	An example of hallucinations in LLMs: given the same prompt, an LLM might generate factually correct or incorrect responses at different inference times. This figure is adapted from [214].	35
2.11	A diagram illustrating the three steps of Reinforcement Learning from Human Feedback (RLHF), as cited from Ouyang et al. [110]: (i) supervised fine-tuning (SFT), (ii) reward model (RM) training, and (iii) reinforcement learning via proximal policy optimization (PPO) algorithm [139] against the reward model.	38

2.12 Comparison of PPO vs. DPO, where DPO optimizes for human preferences while avoiding RL (figure cited from Rafailov et al. [129]). PPO (left) employs RL to optimize policy by first learning a reward model from human preference data, then fine-tuning the model to maximize this reward. DPO (right) directly optimizes policy to align with human preferences using a straightforward classification objective, bypassing RL by implicitly modeling rewards.	40
3.1 Illustration of the proposed SL-AGENT with a human-bot dialog example. (i) The human-bot dialog example, containing an inappropriate response related to unseen user behaviors (upper part). (ii) Demonstration of the refining process in SL-AGENT with the exhibited dialog example (lower part).	46
3.2 The proposed SL-AGENT operates as follows: (i) Fine-tune the bot using available task-specific dialogs. (ii) Deploy the bot online to gather unlabeled human-bot dialog logs. (iii) Refine the dialog model using reinforcement learning with the fine-tuned reward model. (iv) Utilize machine teaching to provide instructions for new functions and enable iterative self-learning.	47
3.3 Illustration of synthetic dialog construction. Slot values in the delexicalized dialog goal (top left) and dialog log (bottom left) are replaced with values from a restaurant database entry (right) to generate a fully lexicalized synthetic dialog.	50
3.4 The summarized five types of dialog turns featuring inappropriate or incoherent responses. (i) Dialog history (top). (ii) Five categories of inappropriate or incoherent responses corresponding to the given dialog history (bottom).	52
3.5 Illustration of the training example, i.e., the processed dialog turn in the training data.	59
3.6 Two interactive examples. (a) An interactive example between user and SOLOISTS. (b) An interactive example between user and SL-SOLOIST.	67

4.1	The proposed SGP-TOD is depicted with a dialog example, where the prompters integrate the task schema (right) to assist the frozen LLM in generating an appropriate response (left).	72
4.2	Illustration of belief state prediction utilizing DST Prompter. The predicted belief state is highlighted.	75
4.3	Illustration of system action determination and response generation employing the Policy Prompter. The pertinent template turns, previously predicted belief state, retrieved DB state within the input, alongside the generated system action and generated response in the output are accentuated.	77
4.4	Zero-shot end-to-end evaluation results on STAR with different numbers of training dialogs (1, 10, 100, 1,000) / demonstration example turns (1, 10) from source domain/tasks.	88
4.5	Detailed belief instructions in DST Prompter.	97
4.6	A formatting example in Policy Prompter.	98
4.7	Policy Prompter of SGP-TOD on STAR. The relevant template turn within the input, the generated user template utterance, and the system action in the output are accentuated.	99
4.8	Policy Prompter of SGP-TOD-E2E on STAR. The relevant template turn in the input, the generated user template utterance, the system action and the system response in the output are highlighted.	100
5.1	Illustration of <i>Self-Alignment for Factuality</i> . Given a prompt to write a biography, before factuality alignment, the LLM generates some facts that are not accurate. Through self-evaluation, the LLM is capable of identifying these inaccurate facts. The feedback from the self-evaluation is used as a reward signal to align the LLM towards factuality. Each fact is highlighted in distinct colors, and the corrected facts are marked with green letters.	103
5.2	A diagram illustrating the three steps of our <i>Self-Alignment for Factuality</i> (in long-form text generation task). (i) Generate initial responses for preference data collection. (ii) Estimate the factuality of the responses through self-evaluation for preference labeling. (iii) Create pairwise preference data and fine-tune the LLM using DPO.	105

5.3	The process of constructing training data for SK-TUNING involves two main steps: (<i>i</i>) Randomly select candidate answers and assess their factual accuracy. (<i>ii</i>) Create pairwise training examples by: (a) pairing a correct answer with a “True” prediction and an incorrect answer with a “False” prediction, and (b) pairing a correct answer with a “False” prediction and an incorrect answer with a “True” prediction.	106
5.4	Results of pairwise comparisons on BioGEN across four dimensions: factuality, helpfulness, relevance and naturalness, as evaluated by GPT-4. The left and right sections present the win rates of <i>Self-Alignment for Factuality w/ SELF-EVAL-SKT</i> against FACTTUNE-MC and <i>Self-Alignment for Factuality w/ SELF-EVAL-P(TRUE)</i> , respectively.	118
5.5	Calibration curves of utilizing SELF-EVAL-P(TRUE) and SELF-EVAL-SKT on LLAMA2-7B in the CommonsenseQA task. Following Kadavath et al. [65], we plot confidence vs. frequency that a prediction is correct. The dashed line indicates perfect calibration.	121
5.6	Calibration curves of utilizing SELF-EVAL-P(TRUE) and SELF-EVAL-SKT (without duplicates) on LLAMA2-7B in the CommonsenseQA task. Following Kadavath et al. [65], we plot confidence vs. frequency that a prediction is correct. The dashed line indicates perfect calibration.	124

List of Tables

2.1	An example of a human-bot dialogue in the restaurant domain, adapted from Zhang et al. [212].	12
3.1	A dialog example on domain extension. Task bots need learn to provide pertinent responses concerning the expanded delivery service in subsequent interaction turns, adapting to the continually evolving user or enterprise demands.	55
3.2	A Restaurant-Ext DB entry. The newly introduced slot-value pairs relevant to the extended functionality are highlighted.	56
3.3	Data statistics of four single-domain dialog datasets [12, 117].	57
3.4	End-to-end evaluation results on four tasks. The forth to sixth rows indicate the results of refining with 45 simulated (unlabeled) human-bot dialog logs, based on SOLOIST _S . SOLOIST ₅₀ is quoted from Peng et al. [117]. Info.: Inform, Succ.: Success. (SL-SOLOIST significantly outperforms all baselines in mean with $p < 0.01$ based on Combined.)	60
3.5	Automatic evaluation results on four tasks in Real-Scenario Setting. The first row refers to previously reported SOLOIST _S . The last three rows refer to refining with 30 real (unlabeled) human-bot dialog logs based on SOLOIST _S . Info.: Inform, Succ.: Success. (SL-SOLOIST significantly outperforms all baselines in mean with $p < 0.01$ based on Combined.)	62
3.6	Automatic evaluation results on task definition extensions. (Difference in mean is significant with $p < 0.01$ based on Combined.)	63
3.7	End-to-end evaluation results of Policy Improvement in the Restaurant domain. SL-SOLOIST ₊₂₀ refer to continually refining with 20 real (unlabeled) human-bot dialogs based on SL-SOLOIST (reported in Table 4.2).	64

3.8	Ablation study results on using different PLMs for reward models in Restaurant domain. The first five rows indicate evaluation results of fine-tuned GPT-2, BERT, BERT-Large, RoBERTa, RoBERTa-Large, respectively. The last row refers to previously reported SL-SOLOIST. (Difference in mean is significant with $p < 0.01$ based on Combined.)	66
3.9	Human evaluation results. SR w/o g: Success rate without grounding, SR w/ g: Success rate with grounding.	68
4.1	End-to-end dialog generation evaluation results on Multiwoz. Results of SOLOIST, MARS, ANYTOD-XXL on Multiwoz 2.2 are cited from Zhao et al. [223]. Results of IG-TOD-CHATGPT are cited from Hudecek and Dusek [61]. Other results of the full-shot fine-tuning methods are cited from He et al. [51] and Sun et al. [158]. Info.: Inform, Succ.: Success.	84
4.2	End-to-end dialog generation evaluation results on RADDLE. The few-shot fine-tuning results are cited from Peng et al. [118].	86
4.3	Zero-shot end-to-end next action prediction evaluation results on STAR. (Difference in mean is significant with $p < 0.01$.)	86
4.4	End-to-end evaluation results on domain extension. FT: fine-tuning, FS: few-shot prompting, ZS: zero-shot prompting.	89
4.5	Ablation study on the impact of the three components in the proposed SGP-TOD and the database expertise on Multiwoz using GPT-3.5. -policy: removing Policy Prompter, -DB: removing database information, -belief: removing DST Prompter.	91
4.6	Illustration of two interactive examples between a user and CHATGPT (left), SGP-TOD-CHATGPT-EXT (right). Appropriate responses and inappropriate responses are highlighted.	93
4.7	Interactive human evaluation results. S w/o g: Success without grounding; S w/ g: Success with grounding; Und.: Understanding; App.: Appropriateness; T.: Turns.	94
5.1	The prompt used for self-evaluation.	107

5.2	Task descriptions and dataset information for main experiments. Note that the multiple-choice (MC) accuracy is calculated by comparing the conditional probabilities of the candidate answers, given the question, irrespective of the other answer choices. A positive result is recorded when the truthful answer achieves the highest ranking among the options, following Chuang et al. [21], Li et al. [78], Lin et al. [84], Touvron et al. [171].	110
5.3	Few-shot evaluation results on three distinct tasks: 6-shot prompting results of the MCQA and short-form generation tasks on TruthfulQA, and 5-shot prompting results of the long-form generation task on BioGEN. ¹ Results on TruthfulQA marked with an asterisk are cited from Li et al. [78] and Chuang et al. [21]. The remaining results of DOLA and FACTTUNE-MC are reproduced following Chuang et al. [21] and Tian et al. [167]. T.*I.: True*Info, FActS.: FActScore.	115
5.4	Results of <i>Self-Alignment for Factuality</i> that employ various approaches for confidence estimation.	116
5.5	MCQA datasets utilized for investigating the confidence estimation capabilities of the SELF-EVAL-SKT. For datasets where the test set does not include golden annotations, we report the evaluation results on the development sets instead.	119
5.6	Following Singhal et al. [151], Taylor et al. [165], we report the 5-shot results on multi-choice QA tasks. Note that the results of LLAMA2-7B are reported using the lettered choices format (examples are provided in Table 5.5), as Kadavath et al. [65], Rae et al. [128] suggest that models are well-calibrated in this format. The results on CommonSenseQA (CSQA) (7-shot), OpenBookQA (OBQA) (0-shot), and MMLU (5-shot) as reported in Touvron et al. [171].	120
5.7	A qualitative analysis comparing the biographies generated by LLAMA2-7B and self-alignment w/ SELF-EVAL-SKT on BioGEN. The hallucinated facts are highlighted in red.	123
5.8	Analysis on the types of factual errors that remain challenging after applying self-alignment.	126
5.9	The prompt used by GPT-4 for pairwise evaluation.	127
5.10	Prompt generated by GPT-4 on BioGEN.	128

List of Terms

AGI Artificial General Intelligence

AI Artificial Intelligence

Bi-LSTM Bidirectional Long Short-Term Memory

DB Database

DL Deep Learning

DPO Direct Preference Optimization

DST Dialog State Tracking

E2E End-to-End

IID Independent and Identically Distribution

LLM Large Language Model

LSTM Long Short-Term Memory

MCQA Multi-Choice Question Answering

MLE Maximum Likelihood Estimation

NLG Natural Language Generation

NLP Natural Language Processing

NLU Natural Language Understanding

OOD Out-of-Distribution

PLM Pre-trained Language Model

POL Dialog Policy Learning

- PPO** Proximal Policy Optimization
- RL** Reinforcement Learning
- RLHF** Reinforcement Learning from Human Feedback
- RNN** Recurrent Neural Network
- SF** Slot Filling
- TOD** Task-Oriented Dialog

Chapter 1

Introduction

“We want AI agents that can learn like we can.” — Richard S. Sutton

The widespread adoption of intelligent assistants such as Apple’s Siri,¹ Google Home,² OpenAI’s GPT-4o [108], Anthropic’s Claude 3.5 Sonnet [3], and Google’s Gemini 2.0 [25] highlights both the immense potential and the significant challenges in developing end-to-end task-oriented dialog (TOD) systems, commonly referred to as “task bots.” These task bots promise seamless human-computer interaction, capable of emulating human conversation, providing accurate information (e.g., a restaurant’s phone number), and even completing complex tasks like booking flights [35, 126, 221]. However, building and maintaining TOD systems that effectively engage users in natural language conversations and task completion remains a complex and evolving challenge within the field of artificial intelligence (AI).

Although significant advancements in reinforcement learning (RL) [140, 179] and deep learning (DL) [35] have markedly improved language understanding, generation, and decision-making capabilities in TOD systems [41, 57, 60, 70, 87, 115, 119, 122, 123, 172], these systems are typically trained on static dialogue corpora to mimic

¹<https://www.apple.com/siri/>

²<https://home.google.com/welcome/>

human conversations for task completion. As a result, they often struggle with the dynamic and evolving nature of real-world user needs. For example, a restaurant bot designed for table reservations may falter when faced with unexpected queries about delivery services as business requirements change. Current approaches for maintaining these systems frequently rely on high-quality human annotations, which are costly and difficult to scale [24, 61, 78, 88, 111, 143, 167, 194].

Inspired by Alan Turing’s vision of machines capable of learning and simulating human cognition [173], this thesis aims to develop highly adaptable task bots that can autonomously manage real-world dynamics with minimal or zero human intervention, mirroring humans’ inherent ability to continuously learn and adapt [68].

1.1 End-to-End Task-Oriented Dialog Systems

In this section, we explore prevalent approaches to building end-to-end TOD systems. The term *end-to-end* refers to a fully automated system that processes inputs and generates outputs without intermediate manual intervention or predefined rule-based components [35]. In the context of *TOD systems*, which assist users in completing tasks, end-to-end models are typically trained holistically using deep learning techniques to directly map user inputs to appropriate responses or actions.

Modular End-to-End TOD Systems. Traditionally, TOD systems have been designed as modular pipelines, known as **modular end-to-end TOD systems**. As illustrated in Figure 1.1A, such a system consists of four sequentially connected modules, each responsible for a specific subtask:

1. **Natural Language Understanding (NLU):** Located in the upper left of Fig-

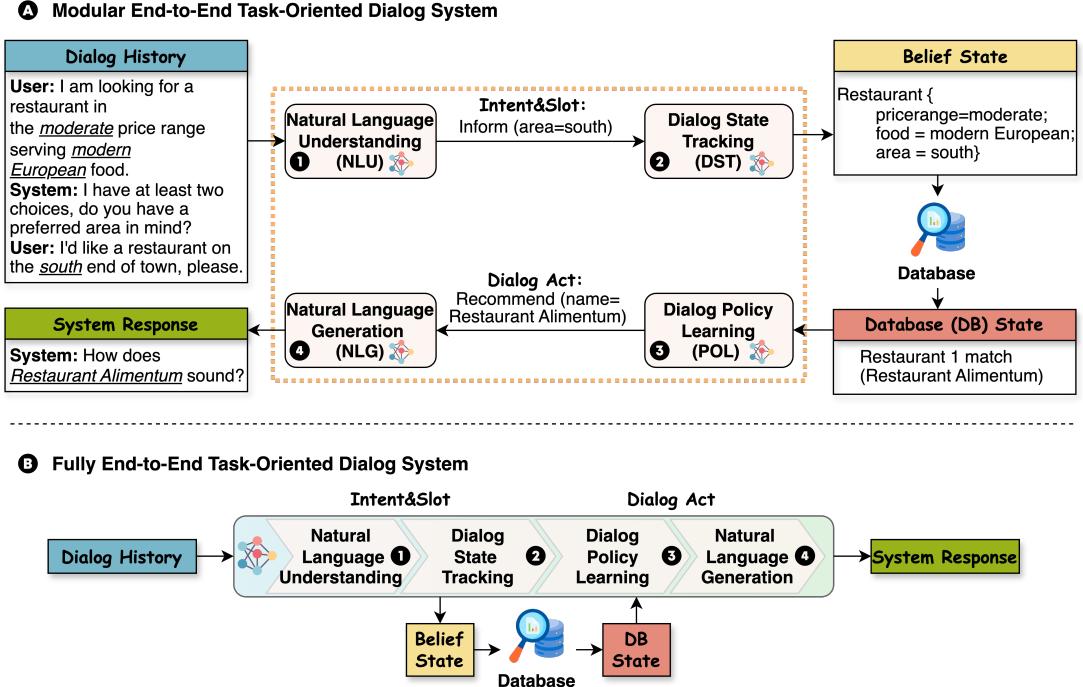


Figure 1.1: Architectures of two end-to-end task-oriented dialogue systems: (i) a modular end-to-end ToD system, which comprises four neural-model-based dialogue modules, namely NLU, DST, POL and NLG modules (in the upper part), and (ii) a fully end-to-end ToD system, which utilizes a single neural model to subsume all four modules into one (in the lower part).

ure 1.1A, this module extracts user intents and slot-value pairs from the dialogue history (e.g., “inform (area=south)” indicates the user intent is *inform*, with the slot-value pair *area=south*, meaning the user is looking for a restaurant in the south).

2. **Dialogue State Tracking (DST):** In the upper right of Figure 1.1A, this module maintains the belief state, which represents the current dialogue context [155]. For example, a belief state such as “Restaurant domain {price range = moderate, food = modern European, area = south}” indicates that the user is looking for a restaurant in the southern area of town that offers modern European cuisine

at a moderate price. This belief state summarizes the user’s constraints accumulated throughout the dialogue (in the upper left corner of Figure 1.1A) and is used to query a task-specific database (DB) to retrieve the DB state, such as the number of matching entities. In this case, only one restaurant—“Restaurant Alimentum”—meets the specified criteria.

3. **Dialogue Policy Learning (POL):** Shown in the lower right of Figure 1.1A, this module determines the system’s next action based on the belief state and the database (DB) state. For example, if the DB query returns a single match, such as “Restaurant Alimentum,” the policy module may choose the action “recommend (name=Restaurant Alimentum),” indicating that the system should suggest this restaurant to the user.
4. **Natural Language Generation (NLG):** Positioned in the lower left corner of Figure 1.1A, this module converts the selected system action, typically a dialogue act, into a natural language response, e.g., “How does restaurant Alimentum sound?”.

Such a task bot can be created by jointly training neural models, such as recurrent neural networks (RNNs) and long short-term memory networks (LSTMs) [144], for all four modules on large, task-specific, human-annotated dialogue corpora to acquire the necessary knowledge and skills [37, 73, 186, 217]. However, scaling and adapting these systems to new domains requires training on high-quality, large-scale, and diverse dialogue corpora specific to those domains, the collection of which demands extensive human annotation efforts.

Fully End-to-End TOD Systems. The emergence of pre-trained language models (PLMs) such as GPT-2 [127] and T5 [131] marked a turning point in the development of data-driven task bots, transitioning from modular end-to-end TOD systems (Figure 1.1A) to fully end-to-end TOD systems (Figure 1.1B).

These systems are typically built using the *pre-training–fine-tuning* paradigm. PLMs, often Transformer-based models [176], share a similar architecture centered on self-attention mechanisms. They are first pre-trained on massive amounts of unannotated web data [94] to acquire broad world knowledge and strong language modeling capabilities. Subsequently, they are fine-tuned on smaller, task-specific dialogue corpora, enabling effective adaptation to downstream TOD tasks with relatively little annotated data. Notably, Transformer-based PLMs have enabled the parameterization of the entire dialogue system within a single model [48, 56, 118], unifying the four key components—natural language understanding (NLU), dialogue state tracking (DST), policy learning (POL), and natural language generation (NLG)—into a cohesive architecture.

Compared to modular end-to-end TOD systems (Figure 1.1A), **fully end-to-end TOD systems** (Figure 1.1B) offer several advantages: (i) *Simplified architecture*: all components are integrated into a single model, eliminating the need for separate module design and coordination; (ii) *Streamlined training*: the training process is unified, reducing the complexity of monitoring and optimizing individual modules; (iii) *Data efficiency*: the reliance on large-scale task-specific annotations is reduced by leveraging knowledge acquired during pre-training; (iv) *Strong generalization*: these models often perform well even with minimal task-specific fine-tuning.

The recent emergence of large language models (LLMs), such as ChatGPT [106] and LLAMA [171], has further revolutionized task bot development by significantly scaling model size, leveraging larger pre-training corpora, and increasing parameter

counts [70, 179]. Compared to earlier PLM-based task bots in Figure 1.1B, LLMs require far less task-specific annotated data while achieving state-of-the-art performance. This advantage stems from their extensive knowledge base and broad skill set [185, 228], acquired through large-scale unsupervised pre-training, including their remarkable instruction-following capabilities [111, 185].

While advancements in deep learning have enabled data-driven task bots to evolve from handling a limited set of tasks to functioning as versatile, general-purpose agents, they still face a fundamental limitation: the reliance on sufficient high-quality task-specific data for effective and continuous development.

1.2 Motivation

In this section, we outline the motivation for developing highly adaptable task bots capable of autonomously managing real-world dynamics with minimal or no human intervention. Despite demonstrating impressive language generation and conversational abilities, data-driven task bots often struggle to engage effectively with real users after deployment. Trained on static datasets, these bots face significant challenges in adapting to the dynamic and unpredictable nature of human interactions, including evolving contexts and unforeseen scenarios, collectively referred to as *changing environments*.¹ Specifically, these challenges manifest in the following three key aspects:

1. **Unseen user behaviors.** Real user queries may contain novel language patterns or references to database entries that were absent from the training corpus [88, 121]. For example, while training data may contain structured and

¹The environment refers to the agent’s world, where it operates and interacts.

formal expressions (e.g., “Could you please provide the phone number of the restaurant?”), real users may phrase their inquiries in brief and informal ways (e.g., “phone number?”). Additionally, users may refer to entities stored in the database but not covered in the training data, resulting in inappropriate responses such as irrelevant or incorrect answers. For example, if a user asks, “Is Café Milano open now?”, but “Café Milano” was not present in the training dialogue corpus, the system may fail to ground the query correctly and respond inaccurately.

2. **Task definition extensions.** As user and business needs evolve, task bots must adapt by incorporating new functionalities or handling additional tasks [86, 156]. For example, a restaurant bot initially designed for table reservations may later receive inquiries about delivery services, a capability beyond its original scope.
3. **Hallucinations.** Even with relevant knowledge, task bots might generate plausible yet factually incorrect responses, known as “hallucinations” [58, 63, 169, 220]. For instance, in response to the query “Is Café Milano open now?”, the bot may hallucinate and reply “Café Milano is open until 10 PM,” despite the actual opening hours being “9 AM–8 PM.”

These challenges not only frustrate users but also undermine trust in the system. To develop advanced task bots into general-purpose agents adept at navigating the intricacies of real-world interactions [74], it is not merely about enabling them to continuously accumulate knowledge; task bots must also apply this knowledge with precision, reliability, and adaptability. Thus, addressing these challenges is imperative for enhancing the robustness and dependability of task bots in real-world scenarios. Moreover, due to the one-time training of task bots and the ever-evolving nature of user or

business requirements, these problems cannot be solved through a one-off scaling of training data. Thus, *post-training adaptation*—the process of continuously updating or fine-tuning a deployed model in response to new interactions or task demands—is both essential and necessary for maintaining and continuously improving task bot performance over the long term.

Current methods primarily rely on extensive human annotations to address these challenges [24, 61, 78, 88, 111, 143, 167]. For instance, Dai et al. [24], Rajendran et al. [132], Shukla et al. [146], Simard et al. [149], Williams and Liden [187] collect high-quality demonstration data to teach task bots how to respond to unseen user behaviors and expand their task capabilities. However, this approach is costly, time-consuming, and struggles to cover the vast space of potential user interactions.

To improve the factual accuracy of task bots, Lightman et al. [83], Ouyang et al. [111], Sun et al. [160] employ reinforcement learning from human feedback (RLHF) [20, 138]. RLHF incorporates human feedback into the training process, refining a model’s outputs based on human evaluations of factual correctness. Specifically, RLHF enables task bots to learn from their own generated responses by leveraging human feedback, reinforcing factually accurate outputs while discouraging the generation of incorrect information. While promising, RLHF faces significant scalability issues due to the resource-intensive nature of obtaining sufficient high-quality human feedback [6, 167].

To this end, this thesis advocates for a paradigm shift towards task bots that can *autonomously adapt, extend their functionalities, and ensure factual accuracy with minimal or zero human intervention*. Accordingly, we investigate three core research questions:

- 1. Adaptability to unseen user behaviors—How can task bots automatically**

adapt to unforeseen user behaviors post-deployment? Can we, drawing parallels with human retrospection, develop task bots capable of evaluating the appropriateness of their responses during user interactions? Furthermore, can these self-evaluations serve as learning signals, enabling task bots to autonomously improve their performance over time?

2. **Extensibility of task definitions—How can task bots seamlessly expand their capabilities to accommodate new tasks and domains?** Humans efficiently acquire knowledge guided by fundamental principles and inductive biases [44]. Can we similarly equip task bots with core principles to facilitate knowledge acquisition and adaptation to new tasks and domains? What principles would be most effective in guiding this learning process?
3. **Trustworthiness in task bot responses—How can task bots reliably convey learned knowledge, i.e., ensuring factual accuracy?** Humans possess the ability to assess correctness even when they cannot generate perfect solutions, leveraging self-reflection for improvement. Can we harness a task bot’s internal knowledge awareness to self-evaluate the factuality of its responses? Could these internal factuality assessments guide task bots toward more reliable knowledge conveyance?

1.3 Thesis Outline and Contributions

The rest of this thesis is organized as follows:

- **Chapter 2 - Literature Review.** We provide an overview of neural approaches for building end-to-end task bots and discuss existing research addressing the

three major challenges faced by task bots post-deployment: (*i*) adapting to unforeseen user behaviors, (*ii*) handling task definition extensions, and (*iii*) mitigating factual errors (“hallucinations”) even when the task bot possesses relevant knowledge.

- **Chapter 3 - Self-Learning for Adaptability.** We introduce SL-AGENT, a self-learning framework enabling task bots to adapt to changing environments by learning from unlabeled human-bot interactions with minimum or zero human annotations. This framework utilizes a pre-trained reward model trained with a novel data augmentation strategy for evaluating response quality [212].
- **Chapter 4 - Schema-Guided LLM Prompting for Extensibility.** We present SGP-TOD, a schema-guided LLM prompting strategy for developing and maintaining task bots with minimal human effort. This approach integrates symbolic knowledge (task schemas) into LLMs, enabling schema-compliant responses and facilitating extensibility to new tasks through schema modification [213].
- **Chapter 5 - Self-Alignment for Factuality.** We propose *Self-Alignment for Factuality*, a self-alignment framework that utilizes an LLM’s self-evaluation capability to reduce the model’s hallucinations. We also introduce SK-TUNING to improve LLMs’ confidence estimation and calibration, further bolstering their self-evaluation abilities [214].
- **Chapter 6 - Conclusions and Future Directions.** We conclude by summarizing our findings and outlining potential avenues for future research.

Chapter 2

Literature Review

In this chapter, we explore the evolving landscape of task bot development, emphasizing the limitations of current approaches and highlighting the need for innovative solutions. We start by providing a succinct overview of neural methods for constructing end-to-end task bots (Section 2.1), outlining the problem definition and training paradigm shifts driven by the progress in deep learning. Subsequently, we investigate the crucial challenges that task bots face after deployment, reviewing significant research efforts focused on addressing: (i) adapting to unseen user behaviors (Section 2.2), (ii) extending capabilities to new tasks (Section 2.3), and (iii) mitigating hallucinations (Section 2.4). By critically analyzing existing research, we identify key gaps and establish the foundation for this thesis’s novel contributions.

2.1 Neural Approaches to Building Task Bots

This section delves into the core task of end-to-end learning for task-oriented dialogue (TOD) systems. We begin by formally defining the problem and proceed to examine how training strategies have progressed in tandem with advances in neural architectures. This historical perspective provides a structured understanding of the current state-of-the-art (SOTA) in TOD system development.

Turn	Utterance
1	User: I want a restaurant on the <i>west</i> end of town.
2	System: Ok. Are you looking for any particular type of food?
3	User: <i>I don't care.</i>
4	System: Are you looking for a particular price range?
5	User: A <i>moderately priced</i> restaurant would be good.
6	System: How about Saint Johns Chop House? They serve moderately priced British food on the west side of town.
7	User: Can i please have their address and phone number?
8	System: Saint Johns Chop House's address is 21-24 Northampton Street, their phone number is 01223353110.
11	User: Thank you, goodbye.
12	System: Thank you for using our system. Goodbye.

Table 2.1: An example of a human-bot dialogue in the restaurant domain, adapted from Zhang et al. [212].

2.1.1 Problem Formulation

At its core, end-to-end learning for task-oriented dialogue systems seeks to build a model capable of engaging in goal-directed interactions with users. This is typically achieved by training on richly annotated dialogue corpora—illustrated in Table 2.1, which presents a complete dialogue example, and Figure 2.1, which shows turn-level annotations—thereby enabling the model to learn and replicate expert conversational behavior.

Formally, the task can be framed as a supervised learning problem. Given a dataset $D = \{(s, b, c, r)\}_{i=1}^n$, where n represents the total number of instances. Each instance represents a dialogue turn sequence (as illustrated in Figure 2.1), consisting of the following elements at turn t :

- **Dialog history s_t :** A sequence of alternating user u and system r utterances leading up to the current turn t , represented as $s_t = [u_1, r_1, u_2, r_2, \dots, u_t]$.
- **Belief state b_t :** A structured representation of the user's goal and constraints (slot-value pairs, e.g., “price range = moderate”) at the current turn.

Dialog History	Belief State	Database State	System Response
User: I am looking for a restaurant in the <i>moderate</i> price range serving <i>modern European</i> food. System: I have at least two choices, do you have a preferred area in mind? User: I'd like a restaurant on the <i>south</i> end of town, please.	Restaurant { pricerange=moderate; food = modern European; area = south}	Restaurant 1 match (Restaurant Alimentum)	System: How does <i>Restaurant Alimentum</i> sound?

Figure 2.1: An example of a dialog turn sequence at turn t , consisting of the dialogue history up to the current turn s_t , annotated with the belief state b_t , database state (i.e., retrieved database entry) c_t , and system response r_t .

- **Database state c_t :** The entry retrieved from a database (e.g., “Restaurant Alimentum”) based on the belief state at the current turn.
- **System response r_t :** The system’s utterance generated in response to the user’s input and current dialogue context.

The objective is to learn a function f that maps the dialog history s_t to the corresponding belief state b_t , the retrieved database (DB) state c_t , and ultimately, the appropriate system response r_t .

$$f : s_t \longrightarrow r_t \quad (2.1)$$

Considering the sequential structure of a task-oriented dialogue system, this complex mapping can be decomposed into two subsequent sub-tasks:

- **Belief state prediction.** Given the dialog history up to current dialog turn s_t , a neural model is trained to generate a belief state b_t (Equation 2.2). The belief state is then used to query a database and obtain the database (DB) state c_t in a deterministic manner, $c_t = DB(b_t)$.

- **Dialogue response generation.** Grounded in the dialog history s_t , belief state b_t , DB state c_t , a neural model is trained to generate corresponding system response r_t (Equation 2.3).

$$f : s_t \longrightarrow b_t, \quad c_t = DB(b_t) \quad (2.2)$$

$$f : s_t, b_t, c_t \longrightarrow r_t \quad (2.3)$$

These sub-tasks are often jointly modeled using neural networks, parameterized by θ . This joint modeling is typically represented as the probability $p_\theta(D)$ and can be achieved through two primary approaches: (i) utilizing separate modules for each sub-task, i.e., NLU, DST, POL, NLG [186], as illustrated in Figure 1.1 in Chapter 1, or (ii) employing a single, unified model [48, 56, 119].

2.1.2 Training Paradigm Shifts

This section outlines the evolution of training methodologies for task-oriented dialogue (TOD) systems, emphasizing a shift from data-intensive models to approaches that leverage pre-trained linguistic knowledge. This transformation has greatly enhanced the efficiency, scalability, and adaptability of conversational AI. The progression is illustrated in Figure 2.2, which highlights three major paradigm shifts in the development of end-to-end task bots: the **standard training paradigm** (Figure 2.2A), the **pre-training then fine-tuning paradigm** (Figure 2.2B), and the **pre-training then prompting paradigm** (Figure 2.2C).

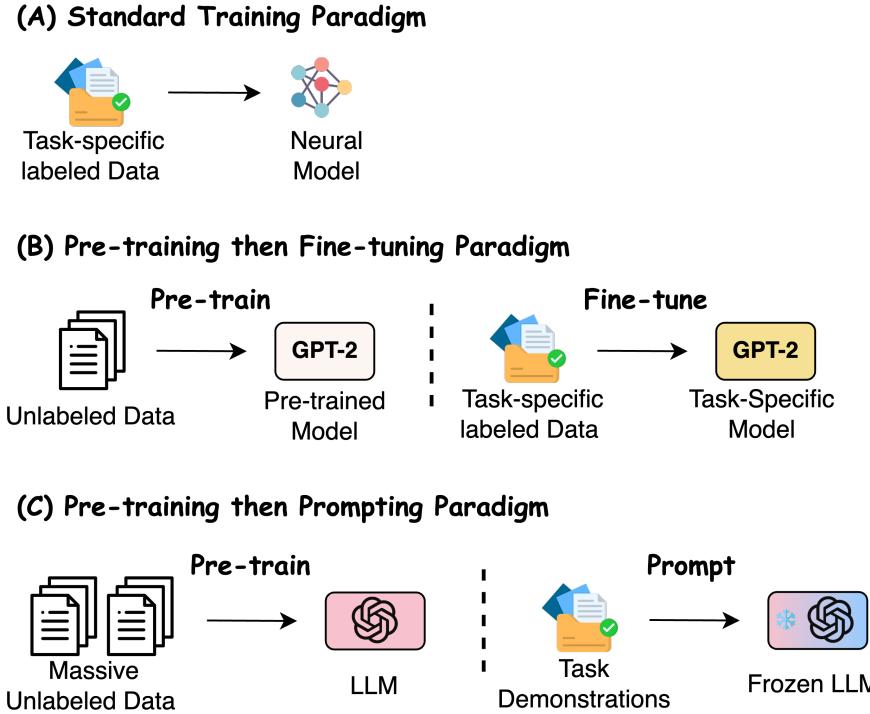


Figure 2.2: Evolution of training paradigms for developing end-to-end task bots: (i) standard training, (ii) pre-training then fine-tuning, (iii) pre-training then prompting.

Standard Training Paradigm: End-to-End Trainable Task Bots. The initial generation of end-to-end trainable task bots [186, 217, 225] was developed using the **standard training paradigm** (Figure 2.2A). These systems typically employed sequence-to-sequence (Seq2Seq) models [161] with attention mechanisms [176], trained in a supervised manner on large-scale annotated dialogue corpora to replicate expert behavior. While this approach enabled early progress in TOD, it required substantial domain-specific labeled data, making it costly and difficult to scale.

Pre-Training Then Fine-Tuning Paradigm: PLM-Based Task Bots. The emergence of pre-trained language models (PLMs), such as GPT [127] and T5 [131], marked a major shift in TOD development. As shown in Figure 2.2B, PLM-based

task bots adopt the **pre-training then fine-tuning paradigm**. In this framework, models are first pre-trained on large volumes of unannotated text to acquire general linguistic and generative capabilities. They are then fine-tuned on limited task-specific dialogue datasets to adapt to particular application domains. This significantly reduces the need for costly annotations. Most advanced task bots [48, 56, 119] follow this paradigm, leveraging auto-regressive generative architectures such as the GPT series [10, 127, 178] as backbone models. A detailed breakdown of this paradigm is presented below:

- **Pre-training.** PLMs, built on the Transformer architecture [94], are pre-trained on massive unlabeled corpora such as the WebText dataset [127] using self-supervised learning. This phase equips the model with a broad understanding of language patterns and the ability to generate coherent and contextually appropriate text. The primary objective is *language modeling*, where the model learns to predict the next token in a sequence, given its preceding context. This is known as the *autoregressive* approach, in which the model generates tokens sequentially, one at a time, conditioning on all previously generated tokens [127]. For example, given the sequence “I would like to book a”, the model might predict the next token as “flight”, “hotel”, or “restaurant”, depending on the learned context. Tokens typically represent subword units, such as “book”, “ing”, “er”, or entire words like “flight” or punctuation marks like “.”. Formally, given a set of training examples (x_1, x_2, \dots, x_m) , where each x_i is a sequence of tokens (s_1, s_2, \dots, s_n) , the model is trained to maximize the likelihood of the observed sequence using *maximum likelihood estimation (MLE)*:

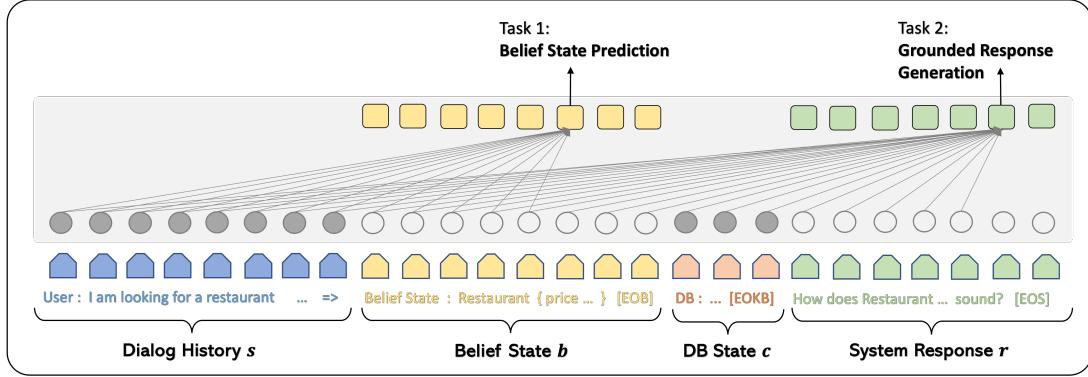


Figure 2.3: An illustrative example of a dialog model that employs an auto-regressive pre-trained language model as its backbone to parameterize the sequential dialog pipeline, as depicted in Figure 1.1 (adapted from Peng et al. [118]). The model takes the dialog turn sequence in Figure 2.1 as input.

$$p(x) = \prod_{i=1}^n p(s^i | s_1, \dots, s_{i-1}) \quad (2.4)$$

- **Fine-tuning.** Once pre-trained, these models can be fine-tuned for specific dialogue tasks using relatively small amounts of task-specific annotated data. This pre-training then fine-tuning paradigm allows the model to adapt to downstream tasks with ease, as shown in Figure 2.2B. Considering the sequential structure of a TOD system, as illustrated in Figure 2.3, a single neural auto-regressive model is utilized to frame the end-to-end dialogue as a sequential generation task, with the joint probability $p_\theta(s, b, c, r)$. This probability can be factorized auto-regressively as follows [119]:

$$p(s, b, c, r) = \underbrace{p(r | s, b, c)}_{\text{Response Generation}} \cdot \underbrace{p(b | s)}_{\text{Belief State Prediction}} \cdot p(s) \quad (2.5)$$

Here, s denotes the dialogue history (user and system turns), b is the belief state (intents and slot-value pairs), c is the database state (retrieved results), r

is the system response. Thus, the parameter θ is learned by maximizing the log-likelihood over the training dataset D , using a joint objective:

$$\begin{aligned} J_\theta(D) &= \sum_{n=1}^{|D|} (\log p(b | s) + \log p(r | s, b, c)) \\ &= \sum_{n=1}^{|D|} \left(\sum_{t=1}^{T_b} \log p_\theta(b_t | b_{<t}, s) + \sum_{t=1}^{T_r} \log p_\theta(r_t | r_{<t}, s, b, c) \right). \end{aligned} \quad (2.6)$$

In this formulation, $|D|$ denotes the number of dialogue samples in the dataset, T_b and T_r are the lengths of the belief state and response sequences, respectively, $b_{<t}$ and $r_{<t}$ represent all tokens preceding time step t in the belief and response sequences. This objective enables the model to jointly learn belief state tracking and response generation in a unified, end-to-end manner.

This two-stage training paradigm—pre-training then fine-tuning—significantly reduces the reliance on extensive task-specific annotated data, thereby making the development of task-oriented dialogue systems more efficient and flexible compared to the standard training paradigm.

Pre-training then Prompting Paradigm: LLM-based Task Bots. At the forefront of conversational AI are Large Language Models (LLMs) such as GPT-3 [10], PaLM [2], GPT-3.5-Turbo [106], and GPT-4 [107]. As shown in Figure 2.2C, these models are pre-trained on massive unlabeled corpora with billions of parameters and diverse linguistic data, enabling them to acquire extensive world knowledge and strong conversational abilities [185]. Rather than relying on fine-tuning, these models can be adapted to specific dialogue tasks through the **pre-training then prompting** paradigm. In this approach, task-specific behavior is elicited by conditioning the model on natu-

ral language instructions or a few example dialog examples—an ability referred to as *in-context learning* or the *prompting* paradigm [109, 180, 185]. This enables rapid prototyping of task bots with minimal labeled data and no additional parameter updates. A brief overview of this paradigm is provided below.

- **Emergent abilities.** By scaling up both the size of pre-training datasets and model parameters, LLMs—built on the Transformer architecture [175]—develop a broad range of language capabilities often referred to as *emergent abilities* [185]. These include conversational fluency [125], text summarization, and question answering [109], among others. As models grow in scale, they demonstrate remarkable sample efficiency and generalization across diverse NLP tasks, even when given minimal task-specific input [180, 185, 228].
- **In-context learning (prompting).** At inference time, LLMs can perform tasks by conditioning on a small number of examples or instructions provided in the input prompt, without any further parameter updates. This approach—known as *in-context learning*—significantly reduces the need for large annotated datasets. It often achieves performance comparable to, or even exceeding, that of models fine-tuned on full datasets [10, 106, 180, 185]. In the context of dialogue systems, this typically involves providing a few sample dialogue turns as demonstrations to guide the model’s behavior on new tasks [61, 96]. For a detailed discussion, refer to Section 2.3.

This pre-training then prompting paradigm, as illustrated in Figure 2.2C, streamlines the model deployment process by minimizing the dependence on task-specific labeled data. It marks a significant shift from traditional supervised learning toward more flexible, instruction-driven interaction.

This evolution in training paradigms—from supervised fine-tuning to context-driven prompting—simplifies the development of task bots while expanding their generalization and adaptability, enabling more dynamic and robust task-oriented dialogue systems.

2.2 Adaptations to Unseen User Behaviors

This section delves into a fundamental challenge in conversational AI: equipping task bots with the ability to adapt to unseen user behaviors.

Task-oriented dialogue (TOD) systems are typically trained within structured and predictable environments. However, real-world interactions are inherently noisy and diverse. Consider a scenario in which a user engages with a task bot expecting seamless assistance, but the system performs poorly when faced with unexpected phrasing or unconventional requests. This discrepancy underscores a core limitation in current TOD systems: *their limited ability to generalize beyond the data distributions observed during training.*

Task bots are often trained on datasets where user queries conform to expected linguistic patterns or templates. In contrast, real users frequently employ informal language, slang, or ambiguous expressions. For instance, while a training utterance might be phrased as “Could you please recommend a nearby restaurant?”, a real user may simply say “Restaurant please.” Such variability presents a significant challenge in designing robust dialogue policies capable of sustaining coherent and goal-directed conversations in open environments [146].

Numerous research efforts have attempted to address this challenge. One of the most promising and efficient approaches involves learning directly from real user inter-

actions [88]. Human-bot dialogue logs offer a rich repository of insights into handling unforeseen user behaviors. A well-trained model may generalize from these logs to manage similar, yet previously unseen, queries. However, directly fine-tuning on such interaction logs can be problematic. Supervised learning may inadvertently reinforce the model’s own mistakes embedded in the data, leading to performance degradation rather than improvement [49]. Furthermore, obtaining high-quality annotations for these logs requires substantial human effort, making the process expensive and difficult to scale.

To address these issues, two primary approaches have been developed for learning from human-bot interactions with minimal human intervention: (*i*) verbal human corrections [24, 49, 132, 146, 149, 187] and (*ii*) numerical human feedback [40, 41, 88, 143, 156].

2.2.1 Learning with Verbal Human Corrections

To enhance the adaptability and performance of task bots, researchers have explored methods to integrate human feedback directly into the learning process. In this section, we discuss two prominent approaches: **proactive learning** and **machine teaching**.

Proactive Learning. Dai et al. [24], Hancock et al. [49], Rajendran et al. [132] propose that dialogue agents should actively request corrections from users when their responses appear inappropriate. As illustrated in the upper part of Figure 2.4, consider a scenario in which the chatbot produces an unsuitable response such as “Lol. I never eat them!” (Utterance C). The user replies with a confused question: “What are you talking about?” (Utterance D). Based on this cue, the chatbot infers that its previous



Figure 2.4: Illustration of proactive learning in a self-feeding chatbot. The model continuously estimates user satisfaction during dialogue interactions to identify appropriate moments to request feedback. Based on user responses, it extracts new training examples for two distinct tasks: DIALOGUE (from satisfied user responses) and FEEDBACK (from user-provided corrections). These examples are then used to fine-tune the model, thereby enhancing its overall dialogue capabilities. Cited from [49].

response may have been inappropriate and proactively requests clarification: “Oops! I messed up. What should I have said?” (Utterance E). The user then provides a correction: “Maybe ask me what I thought about French food?” (Utterance F). Subsequently, dialogue examples derived from either satisfied user responses or user-provided corrections are extracted (in the lower part of Figure 2.4) and used to fine-tune the model for future improvements.

This real-time learning mechanism enables bots to adapt rapidly to diverse user inputs, improving their response quality over time. While **proactive learning** supports continuous model refinement, it may also lead to user frustration due to frequent inter-

ruptions for feedback. Therefore, achieving a balance between learning efficiency and user experience is essential.

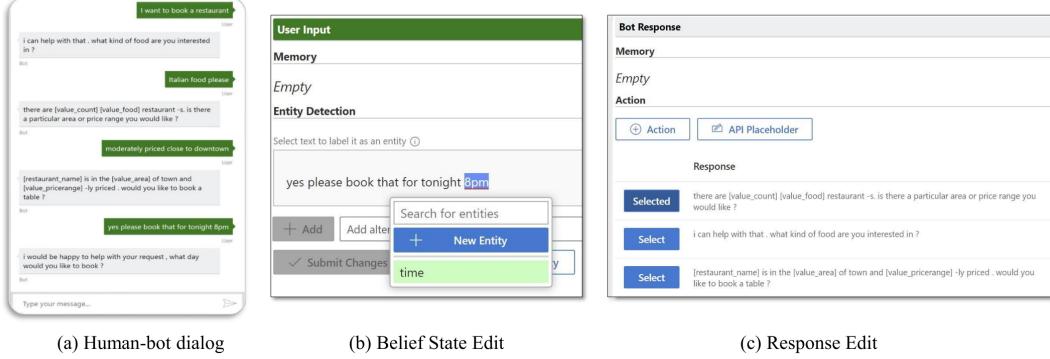


Figure 2.5: Illustration of the machine teaching process [146], cited from [117]. Dialog authors inspect and correct representative failed human-bot interaction logs (a) to create new training dialogues. This involves correcting belief states (b) and selecting, inserting, or modifying inappropriate responses, typically using action templates (c).

Machine Teaching. Alternatively, Shukla et al. [146], Simard et al. [149], Williams and Liden [187] advocate for *machine teaching*, in which human experts retrospectively analyze and correct bot interaction logs to construct training data. As illustrated in Figure 2.5, experts first visualize and select representative failed dialogue logs (Figure 2.5a). They then sequentially correct the associated belief states (Figure 2.5b) and modify inappropriate responses (Figure 2.5c). This process may involve selecting or editing responses using action templates and adding relevant slot-value pairs for specific tasks.

By shifting the feedback process offline, machine teaching minimizes user disruption during live interactions and enables experts to focus on correcting critical errors. This targeted approach leverages expert knowledge to refine bot behavior and improve task-specific performance without requiring real-time user input. However, its effec-

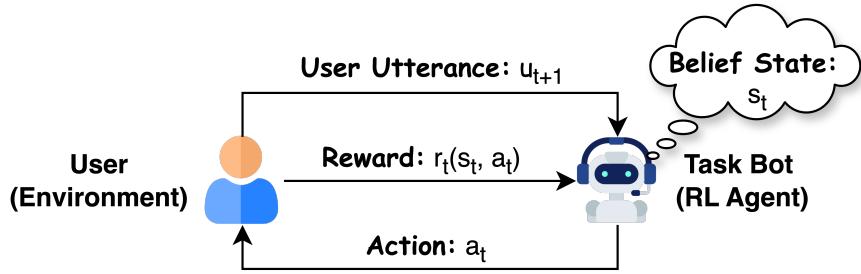


Figure 2.6: Dialog policy optimization in an RL loop, where the interaction between a task bot and a user parallels the interaction between an RL agent and its environment.

tiveness remains highly dependent on the availability and quality of human annotations.

Both proactive learning and machine teaching demonstrate the critical role of human input in bot development, yet they also spotlight an ongoing challenge: reducing the reliance on human annotations.

2.2.2 Learning with Numerical Human Feedback

In an effort to reduce the dependency on human annotations, Gasic et al. [40], Gasic et al. [41], Liu et al. [88], Shah et al. [143], Su et al. [156] have turned to **reinforcement learning** [162] to optimize dialog policies through real user interactions [188], guided by numerical human feedback. The goal here is to *optimize long-term rewards*, specifically by ensuring task completion. A dialogue is deemed successful when the user's task is fulfilled and all information requests are met. The interaction between a task bot and a user effectively becomes a *sequential decision-making problem* within an RL framework [41, 87, 122, 123, 172, 202], where the bot learns from its actions and the feedback it receives to adapt its policy for better performance in future interactions.

At each dialog turn t , the following steps occur (Figure 2.6):

- The bot maintains the belief states based on the ongoing conversation (e.g., “Restaurant (price = moderate; food = European; area = south)” in Figure 2.1), and generates an action a_t based on the current state s_t . This action could be a user response, a dialog act, or an internal operation such as a database lookup or an API call.
- The user provides the next utterance u_{t+1} , which the bot uses to update the dialog state for the subsequent turn $s_{t+1} = s_t \oplus a_t \oplus u_{t+1}$ (\oplus is the concatenation operator).
- The bot receives a *numerical reward* r_t , reflecting the appropriateness of its action. Positive rewards (+1) are given for suitable responses, while negative or zero rewards (-1/0) are assigned for inappropriate ones [35, 142].

The objective of the bot is to learn an optimal policy $\pi_\theta(a | s)$ that selects actions to maximize the expected cumulative discounted reward over time. Here, θ denotes the parameters of the policy model. This objective, known as *dialogue policy optimization*, is formally expressed as:

$$J(\theta) = E [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t} r_T] \quad (2.7)$$

where $\gamma \in [0, 1]$ is the discount factor, which reduces the weight of future rewards relative to immediate ones. The variable T denotes the total number of turns in a dialogue episode.

The dialog policy are typically learned using policy-based reinforcement learning methods [162], such as the REINFORCE algorithm [188], which learns a parameterized policy for directly optimizing an objective function $J(\theta)$ that can select actions without calculating a value function for a state. The policy gradient can be empirically

estimated as:

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t; \theta) \sum_{i=0}^{T-t} \gamma^i r_{t+i} \quad (2.8)$$

where $\sum_{i=0}^{T-t} \gamma^i r_{t+i}$ represents the cumulative discounted reward at step t in one episode.

Typically, the bot queries numerical feedback from human users as rewards after each turn or entire dialog episode [88, 143]. In this manner, a task bot (the machine) learns to interact with a user (the environment) through *trial-and-error* [162], striving to replicate high-reward actions (e.g., appropriate responses to previously unseen user queries) while avoiding low-reward actions (e.g., inappropriate responses), and gradually adapting to novel user behaviors not encountered during training.

Despite these technological advances, *user feedback remains indispensable* for effective policy refinement. However, real users are often reluctant to provide explicit feedback [156], posing a significant challenge for scalable and sustainable bot development. With the emergence of powerful generative language models [127, 171], which may even surpass human capabilities in certain tasks [13, 164], and the exponential increase in daily user-bot interactions, the need for task bots to autonomously adapt to novel user behaviors has become more pressing than ever. This critical problem of automatic adaptation remains largely underexplored and represents a pivotal frontier for the advancement of conversational AI.

2.3 Extensions to New Tasks

This section addresses a central challenge in task-oriented dialogue (TOD) systems: enabling task bots to seamlessly expand their capabilities to support new tasks and domains.

Consider, for example, a task bot originally developed to handle restaurant reser-

vations. When required to manage a new domain such as hotel bookings, traditional systems often fail to generalize effectively. As noted by Mehri and Eskenazi [100], this limitation stems from the common practice of training end-to-end neural models to implicitly learn task-specific dialogue policies from large, domain-specific datasets. When these models encounter a novel task, they often lack the generalizable representations or explicit policy structures needed to perform adequately, exposing a fundamental limitation in current systems’ adaptability.

Expanding a bot’s functionality—specifically, its ability to acquire new dialogue policies—typically necessitates retraining or fine-tuning with substantial amounts of new, task-specific data. While large-scale pretraining of language models such as BERT [27], T5 [130], and GPT-style models [127] has significantly improved generalization across tasks, these models still require fine-tuning to perform effectively in new domains. Moreover, the cost of retraining increases with the scale of the models, the diversity of user needs, and the frequency with which new tasks arise. This situation presents a critical research question in conversational AI: *how can task bots acquire new dialogue policies with minimal—or even zero—additional training data?* Addressing this challenge is essential for building scalable, adaptable, and future-proof dialogue systems that can keep pace with rapidly evolving user demands and application domains.

Two primary strategies have emerged to address this challenge: (i) fine-tuning-based methods [100, 104, 124, 223, 226] and (ii) prompting-based methods [9, 61, 96].

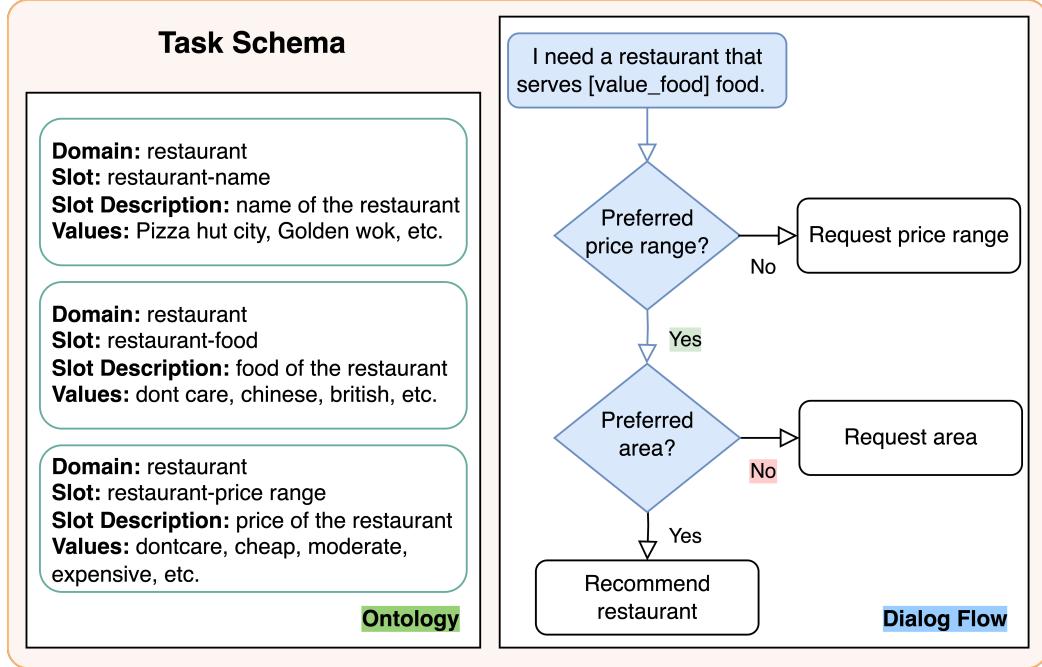


Figure 2.7: An example of a task schema from the MultiWOZ dataset [12] in the restaurant domain, featuring the ontology on the left and the dialog flow on the right.

2.3.1 Fine-tuning for Effortless Extension

Several pioneering works [61, 100, 104, 124, 226] have explored the problem of **zero-shot end-to-end dialogue modeling**, aiming to extend task bots to new domains without requiring task-specific training data [99]. This line of research falls under the broader scope of *zero-shot generalization*, where models are expected to perform tasks they have not encountered during training.

A central solution to this challenge lies in the use of *task schemas*—a form of *symbolic knowledge* [18, 105] that serves as a blueprint for guiding dialogue policies in a schema-guided paradigm [100]. As illustrated in Figure 2.7, a task schema provides a structured representation of a dialogue task and typically consists of two key components:

- **Task-specific ontology.** This component (visualized on the left side of Figure 2.7) defines the domain-specific vocabulary, including relevant entities (slots), their descriptions, and possible values [11]. The ontology provides a semantic framework for interpreting user intents and system actions.
- **Dialog flow.** Shown on the right side of Figure 2.7, this component outlines the expected sequence of conversational steps required for successful task completion [120], ensuring that interactions are coherent, goal-oriented, and aligned with the schema.

These structured representations support two major approaches to schema-based zero-shot adaptation, i.e., **learning a shared dialogue policy** and **providing explicit guidance**.

Learning a Shared Dialogue Policy. Zhao and Eskenazi [226] and Qian and Yu [124] address zero-shot task generalization by enabling models to learn a *shared dialogue policy*—a generalized strategy that can be applied across both source and unseen target domains without requiring retraining. To facilitate this, they synthesize training data using task ontologies and domain-agnostic response templates. For instance, the template “I am looking for a [slot_value] [slot] in the [area].” can be instantiated in multiple domains: in the restaurant domain (“I am looking for a Chinese restaurant in the south of the town.”), or in the hotel domain (“I am looking for a five-star hotel in the city center.”). These instantiations expose the model to reusable interaction patterns that transcend domain boundaries, enabling it to generalize dialogue strategies across tasks with similar structural characteristics.

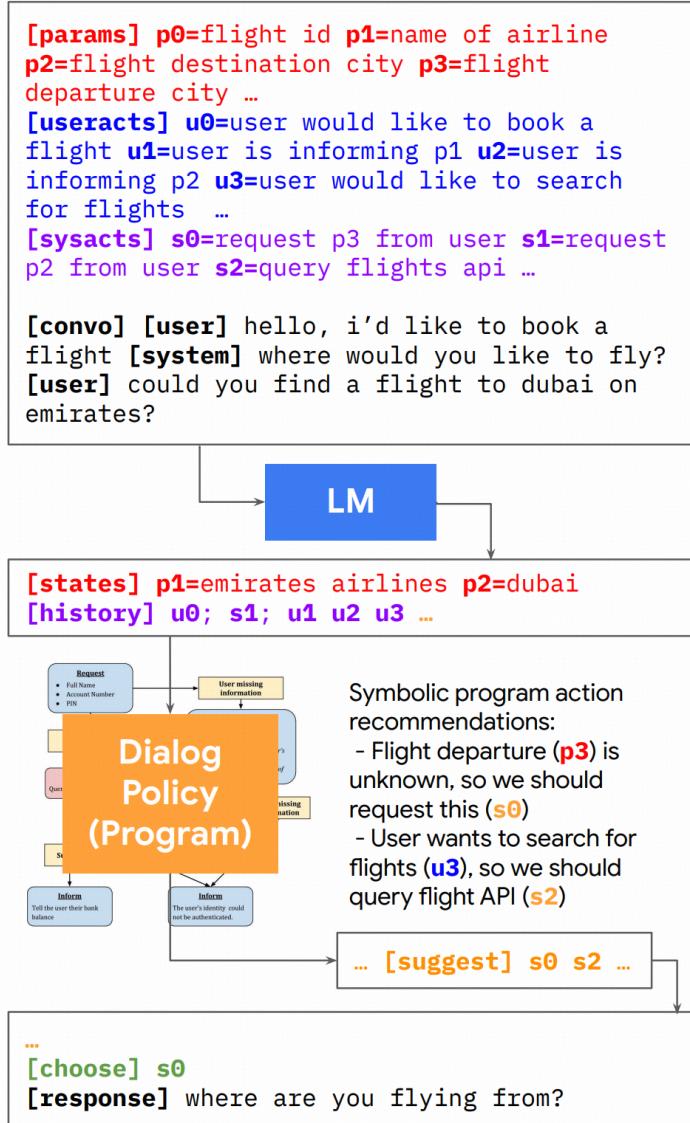


Figure 2.8: An overview of the ANYTOD system, cited from Zhao et al. [223]. A LM conducts zero-shot state and action tracking with respect to a provided schema, abstracting it into a sequence of symbols. A program that executes the dialog policy then recommends which actions to take based on the states sequence, the LM then chooses a single final action and generating a response.

Providing Explicit Guidance. While shared policies support generalization, an alternative strategy involves providing *explicit task-level guidance* through symbolic pol-

icy representations. Mosig et al. [104] and Mehri and Eskenazi [100] introduce *policy skeletons*—external structures that define the high-level flow of a task (as shown on the right side of Figure 2.7). For example, when a user says “I would like a Chinese restaurant,” which matches a template such as “I need a restaurant that serves [value_food] food,” the model is guided—based on the schema—to request the price range. Zhao et al. [223] further extend this paradigm by introducing a neural state tracker that maintains a belief state and executes a predefined task policy. As illustrated in Figure 2.8, the system follows a symbolic plan (highlighted as “dialog policy (program)” in the orange square) to determine the next system action and generate the response “Where are you flying from?” based on the current dialog state. These actions are then realized in natural language, enabling coherent behavior even for unseen tasks.

Despite their effectiveness, these methods face a fundamental limitation. Truly effective zero-shot adaptation—defined as adapting to a target domain with no task-specific training data—still depends on *extensive fine-tuning using heterogeneous, annotated dialogue datasets* from multiple source domains. This requirement introduces substantial resource overhead, making the process labor-intensive and potentially limiting its scalability in real-world applications.

2.3.2 Prompting for Effortless Extension

Prompting-based methods offer a compelling alternative for few-shot end-to-end dialogue modeling (Section 2.1), enabling task bots to acquire new skills with minimal supervision—often from just a few examples. Leveraging the broad knowledge obtained during large-scale pre-training, large language models (LLMs) can act not only as conversational agents but also as “rapid learners” through **few-shot prompting**.

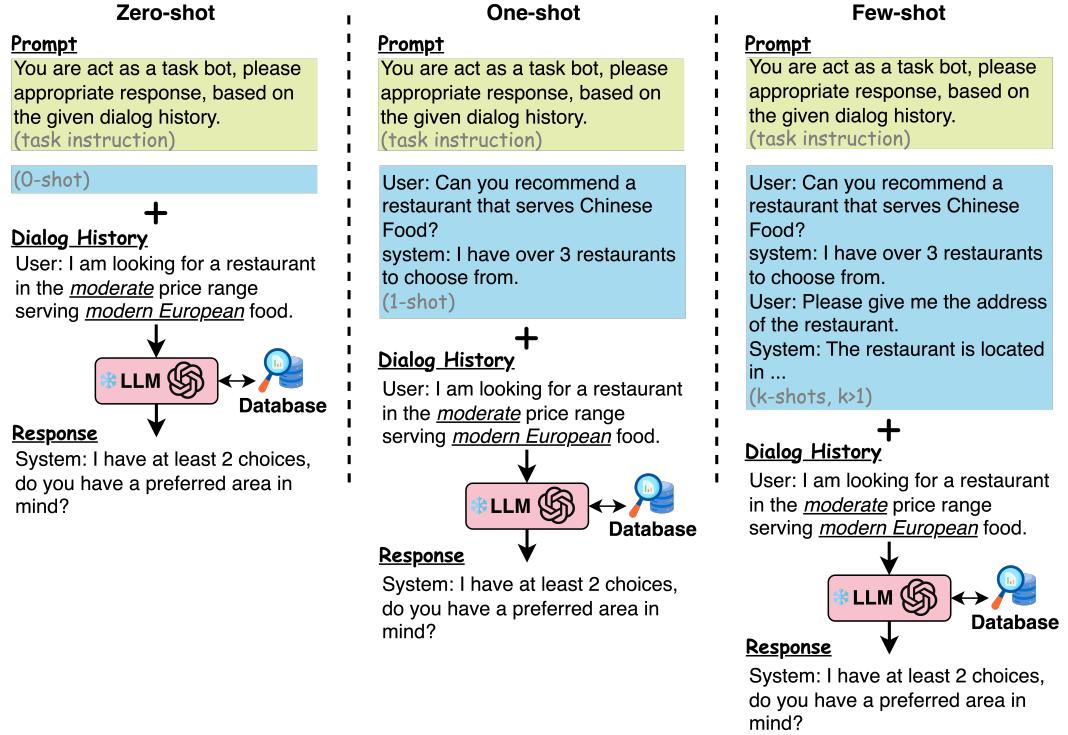


Figure 2.9: Illustration of the prompting paradigm in zero-shot, one-shot, and few-shot settings (from left to right) within a task-oriented dialog system, using zero, one, or a few (k , $k > 1$) dialog turn examples as task-specific exemplars, respectively.

ing [9, 61, 96].

As shown in the right section of Figure 2.9, at each dialogue turn t , a *frozen* LLM, parameterized by θ , receives a small set of exemplars from a new task $\{(s_i, r_i)\}_{i=1}^k$ (where $k > 1$), a task instruction I , the current dialog context s_t , and retrieved database information c_t . The model then generates a response r_t tailored to the ongoing interaction:

$$r_t = LLM_{\theta} \left(\{(s_i, r_i)\}_{i=1}^k, s_t, c_t, I \right) \quad (2.9)$$

This approach enables remarkable adaptability, allowing LLMs to generate task-relevant responses without any parameter updates. However, its effectiveness heavily

depends on the *quality and completeness of in-context exemplars* [30, 89, 227]. Like a chef working with an incomplete recipe, the model often lacks access to the full spectrum of task-specific knowledge, which can hinder successful task execution.

This limitation becomes more pronounced in the **zero-shot setting**, where no in-context examples are available (illustrated in the left section of Figure 2.9). In such scenarios, LLM-based task bots struggle to generalize [61, 213], achieving only a *15% success rate* on the complex, multi-domain MultiWOZ dataset [12], despite strong conversational capabilities in other contexts [125]. This performance gap highlights a critical challenge: the need for more effective mechanisms to convey task structure and intent to LLMs in zero-shot scenarios, enabling them to navigate unfamiliar tasks successfully.

2.4 Mitigating Hallucinations

This section delves into the growing body of research aimed at improving the factuality of LLMs by mitigating hallucinations—outputs that are fluent but factually incorrect.

LLMs have evolved from task-specific bots into general-purpose AI assistants capable of performing a wide range of natural language processing (NLP) tasks, such as question answering and dialog generation [14, 91, 185]. However, they occasionally produce *factual “hallucinations”*—plausible yet incorrect statements [58, 63, 169, 220], which undermine their reliability [90].

In an era where high-quality pre-training data is regarded as essential [45, 171], it is reasonable to expect that an LLM’s knowledge should align with established facts [201]. Nevertheless, models frequently generate factually incorrect outputs even

when they internally possess the correct information. For example, as illustrated in Figure 2.10, when prompted with the question, “What is Westlife’s first album?”, the model correctly responds in one instance with “*Westlife* is the debut studio album by Irish boy band *Westlife*,” but in another instance incorrectly answers “*Coast to Coast*. ” This inconsistency reveals the challenge of ensuring factual reliability in LLMs, even when the underlying knowledge is available.

Why do hallucinations occur? The phenomenon arises from both intrinsic and extrinsic factors [58, 220]. A key intrinsic limitation lies in the pre-training objective—typically maximum likelihood estimation (MLE)—which encourages pattern replication rather than factual correctness [1, 5, 21, 167]. This misalignment between training goals and user expectations fosters hallucinations [4, 10, 110, 166]. Extrinsic factors, such as decoding strategies, further exacerbate the issue. While greedy decoding is deterministic but lacks diversity, stochastic approaches like nucleus sampling [55] introduce variability at the cost of increased hallucination risk [72].

To address these challenges, various strategies have been proposed [58]. This thesis categorizes them based on their timing of intervention [220]: (i) post-hoc correction [38, 81, 97, 116, 174, 203, 224], (ii) inference-time intervention [21, 78, 80, 219] and alignment Training [106, 107, 111, 171].

2.4.1 Post-hoc Corrections

Post-hoc correction techniques refine the output of an LLM *after* it has been generated. Two commonly employed strategies are **resorting to external knowledge** and **self-consistency**.



Figure 2.10: An example of hallucinations in LLMs: given the same prompt, an LLM might generate factually correct or incorrect responses at different inference times. This figure is adapted from [214].

Resorting to External Knowledge. Retrieval-augmented generation incorporates external knowledge sources (e.g., Wikipedia) to enhance or correct model outputs [38, 81, 97, 116, 174, 203, 224]. However, this reliance on external databases presents several challenges. First, access to such databases may be limited or introduce latency, hindering real-time applications. Second, the accuracy of retrieved information directly impacts the LLM’s output. Irrelevant or malicious content within retrieved results can lead to inaccurate or misleading responses [23, 192, 195, 210]. Finally, a tension exists between the LLM’s pre-existing knowledge and the information presented in retrieved documents [191]. This tension can lead to inconsistencies or difficulties in effectively integrating external knowledge with the model’s internal representations.

Self-Consistency. When external resources are unavailable, LLMs perform self-critique by generating multiple interpretations of a prompt and selecting the most consistent answers. Consistency often correlates with accuracy [29, 65, 95, 135, 168, 181]. *Self-consistency methods* use the *model’s own confidence as a proxy for factuality*, enhancing accuracy by analyzing consistency among multiple responses [69, 85, 167, 174]. When an LLM is knowledgeable about a question, the most consistent answer among candidates is more likely to be accurate, while contradictory answers often contain hallucinations [98, 224]. Despite their effectiveness, these approaches can suffer from high latency, potentially frustrating users.

2.4.2 Inference-time Interventions

Unlike post-hoc correction, inference-time intervention seeks to guide the model *during* the generation process. Two primary strategies are commonly employed, i.e., **designing decoding strategies** and **manipulating internal representations**.

Designing Decoding Strategies. As mentioned before, the widespread use of decoding strategies like nucleus sampling introduces diversity into responses but at the cost of potentially increasing the risk of hallucinations. To counteract this, Lee et al. [72] propose factual-nucleus sampling, a novel approach that seeks to harmonize diversity with accuracy by combining elements of both greedy decoding and nucleus sampling.

Manipulating Internal Representations. Recent research [21, 78, 80, 219] suggests that LLMs possess hidden, interpretable structures crucial for generating factual statements. While current decoding strategies are efficient, they may not fully leverage this internal knowledge. This limitation has spurred investigations into directly manipulat-

ing these internal representations to improve factual accuracy.

By identifying and altering specific attention mechanisms or layers related to factual accuracy [21, 78], these approaches seek to enhance the model’s factuality during generation. However, these interventions often require domain-specific data, which may limit their applicability across various topics or contexts.

2.4.3 Alignment Training

Recognizing the limitations of pre-training objectives in ensuring factuality, alignment training has emerged as a pivotal strategy for enhancing the accuracy of LLMs. This approach focuses on aligning LLMs with human preferences to directly optimize the generation of factually accurate statements [106, 107, 111, 171]. Contemporary research primarily adopts two methodologies to achieve this alignment: **instruction-tuning** and **Reinforcement Learning from Human Feedback (RLHF)**.

Instruction-Tuning. This approach involves fine-tuning LLMs on carefully curated datasets of (instruction, response) pairs using supervised learning [110, 208]. By exposing the model to numerous examples of desired behavior, instruction-tuning encourages it to generate responses that align with human expectations for factual accuracy and task completion [75, 182, 201]. This emphasis on learning from explicit demonstrations allows for more targeted control over the LLM’s output, ultimately leading to more reliable and factually grounded generation.

Reinforcement Learning from Human Feedback (RLHF). [20, 83, 111, 138, 160] is a three-stage process aimed at refining the model’s output based on human preferences (e.g., generating factually accurate responses), as illustrated in Figure 2.11:

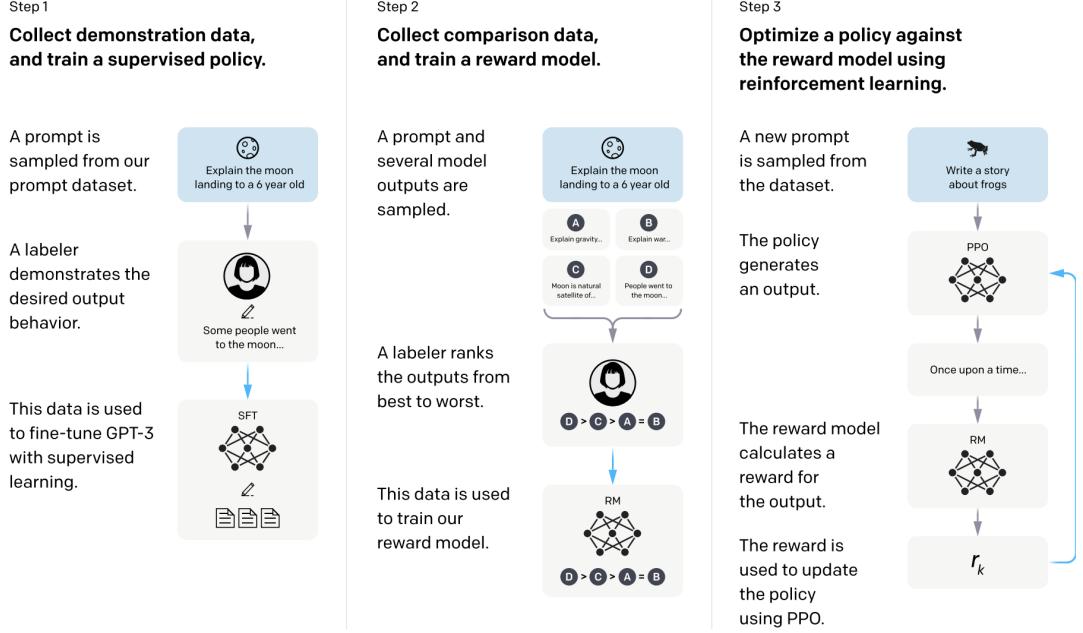


Figure 2.11: A diagram illustrating the three steps of Reinforcement Learning from Human Feedback (RLHF), as cited from Ouyang et al. [110]: (i) supervised fine-tuning (SFT), (ii) reward model (RM) training, and (iii) reinforcement learning via proximal policy optimization (PPO) algorithm [139] against the reward model.

Specifically, RLHF comprises the following three stages, as illustrated in Figure 2.11:

- **Stage 1: Collect demonstration data and train a supervised policy.** This stage first collects a high-quality human demonstration dataset tailored for downstream tasks, e.g., instruction-following. Then, a pre-trained language model is trained on this data using supervised learning, resulting in $\pi_{SFT}(y | x)$, where (x, y) denotes the (prompt, completion) pair.
- **Stage 2: Collect comparison data and train a reward model.** This stage involves collecting comparison data between different model outputs, which is subsequently used to train a reward model that predicts the human-preferred output. Specifically, a reward model is trained to take in a prompt and response

and output a scalar reward, framing the problem as a binary classification task.

The supervised policy model $\pi_{SFT}(y | x)$ is prompted to generate a set of responses $\{y_k\}_{k=1}^K$, where K is the number of responses for each prompt x . A labeler then ranks the responses based on quality, producing a dataset of human-labeled comparisons $D_R = \left\{ (x^{(i)}, y_w^{(i)}, y_l^{(i)}) \right\}_{i=1}^N$, where y_w and y_l represent the preferred and dis-preferred responses for each prompt x , respectively. The distribution of human preference p^* is commonly modeled by the Bradley-Terry Model [8]:

$$p^*(y_w \succ y_l | x) = \frac{\exp(r^*(x, y_w))}{\exp(r^*(x, y_w)) + \exp(r^*(x, y_l))} = \sigma(r^*(x, y_w) - r^*(x, y_l)), \quad (2.10)$$

where $\sigma(x) = \frac{1}{1+\exp(-x)}$ is the logistic sigmoid function.

Subsequently, the parameterized reward model $r_\phi(x, y)$ is optimized using the collected comparisons data D_R through minimizing the negative log-likelihood loss [129]:

$$L_R(r_\phi, D_R) = -\mathbb{E}_{(x, y_w, y_l) \sim D_R} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))] \quad (2.11)$$

- **Stage 3: Optimize a policy against the reward model using reinforcement learning.** The final stage aims to fine-tune the supervised policy π_{SFT} , parameterized by θ , to maximize the expected reward obtained from the learned reward model $r_\phi(x, y)$ using the **Proximal Policy Optimization (PPO)** [140] algorithm. Specifically, the process of LLMs generating responses from training prompts D_{RL} is modeled as a bandit environment [110], where a scalar reward is obtained

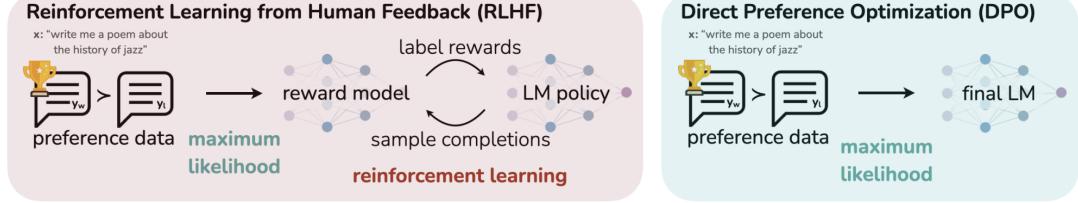


Figure 2.12: Comparison of PPO vs. DPO, where DPO optimizes for human preferences while avoiding RL (figure cited from Rafailov et al. [129]). PPO (left) employs RL to optimize policy by first learning a reward model from human preference data, then fine-tuning the model to maximize this reward. DPO (right) directly optimizes policy to align with human preferences using a straightforward classification objective, bypassing RL by implicitly modeling rewards.

from the reward model r_ϕ at the end of each response.

$$\arg \max_{\pi_\theta} \mathbb{E}_{x \sim D_{RL}, y \sim \pi_\theta} [r_\phi(x, y)] \quad (2.12)$$

A per-token Kullback–Leibler (KL) penalty is added to prevent over-optimization, ensuring the policy model does not deviate too far from the original supervised fine-tuning model π_{SFT} .¹ The training objective is then formulated as:

$$L_{PPO}(\pi_\theta; \pi_{SFT}) = -\mathbb{E}_{x \sim D_{RL}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta D_{KL}(\pi_\theta(y|x) \| \pi_{SFT}(y|x)), \quad (2.13)$$

where β is a parameter controlling the deviation from the initial SFT model π_{SFT} .

Recently, **Direct Preference Optimization (DPO)** [129] (Figure 2.12) has been introduced as a simple, straightforward, and computationally efficient alternative to RL algorithms, such as PPO [140]. Notably, DPO circumvents the need for

¹ KL divergence [53] measures how different two probability distributions are from each other. Here, it helps to maintain the balance between the original capabilities of the model and its refining for better performance according to human preferences.

training a standalone reward model and directly optimizes a language model using the preference data $D_R = \left\{ (x^{(i)}, y_w^{(i)}, y_l^{(i)}) \right\}_{i=1}^N$, with a simple classification objective based on binary cross-entropy loss:

$$L_{\text{DPO}}(\pi_\theta; \pi_{\text{SFT}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D_R} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_w | x)}{\pi_{\text{SFT}}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{\text{SFT}}(y_l | x)} \right) \right], \quad (2.14)$$

DPO refines the policy π_θ by increasing the margin between the log-likelihood of preferred and dis-preferred responses, while ensuring the model does not stray far from the initial policy π_{SFT} .

Although DPO offers efficiency in terms of computation, speed, and engineering efforts, research [62] indicates that PPO generally achieves superior alignment performance. This is largely due to PPO’s ability to utilize online data generated by the current policy, fostering a dynamic learning process that allows for better exploration. Conversely, DPO relies on static, pre-generated offline data, which might limit its exploration capabilities and thus potentially compromise the training quality. However, the choice between these methods hinges on specific use cases and available resources.

While alignment training techniques have demonstrated considerable effectiveness, they rely heavily on human annotations for collecting and labeling demonstration data. This dependence incurs substantial costs, as both the volume and quality of annotations are critical. Exploring cost-effective and reliable alternatives to human supervision remains a promising direction for future research.

Moreover, although post-hoc correction, inference-time intervention, and alignment training have typically been developed in isolation, their integration presents sig-

nificant opportunities for synergy. For example, combining post-hoc correction with inference-time intervention could yield a comprehensive framework that addresses factual inaccuracies both during and after generation. Additionally, alignment training can enhance the efficacy of both strategies by incorporating a deeper understanding of human preferences and expectations. This convergence of methodologies offers a promising path toward more robust and trustworthy language models—an area ripe for further investigation and advancement.

2.5 Chapter Summary

This chapter delves into the fascinating evolution of TOD systems, charting a course from data-hungry beginnings to the rise of lean, knowledge-driven models. We begin in Section 2.1 by establishing the foundation of task bot creation and then journey through three pivotal paradigm shifts: from resource-intensive standard training to the efficiency of pre-training then fine-tuning, culminating in the remarkable flexibility of pre-training then prompting.

This progress, however, is not without its hurdles. Section 2.2 exposes the vulnerability of TOD systems to the unpredictable nature of human communication. While strategies like verbal corrections and numerical feedback offer a guiding hand, they often come at the cost of significant human intervention. This underscores a critical need for TOD systems to develop autonomous adaptation, especially as interactions with humans become increasingly commonplace.

The quest for adaptability extends beyond handling unexpected user behavior to conquering new tasks and domains, as explored in Section 2.3. Here, we dissect two primary approaches: fine-tuning methods, encompassing techniques like zero-

shot modeling and shared policy learning, and the intriguing possibilities of prompting methods, including few-shot and zero-shot transfer learning. While LLMs emerge as front-runners in this domain, their limitations in navigating entirely novel scenarios come to light.

Finally, Section 2.4 confronts the critical challenge of ensuring factual accuracy in LLMs. Despite their impressive capabilities, these models can stumble into the realm of “hallucinations” – generating outputs that are plausible yet factually incorrect. We examine three key mitigation strategies: post-hoc correction, inference-time intervention, and alignment training. A recurring theme emerges: the delicate balancing act between minimizing reliance on costly human annotations and upholding the integrity of factual accuracy.

This tension underscores the very challenges we aim to address in the remaining chapters. Our exploration will delve into innovative methods that promise to tackle these obstacles with minimal or even zero human intervention, paving the way for truly autonomous and reliable TOD systems.

Chapter 3

Self-Learning for Adaptability

This chapter explores how to enable task bots to automatically adapt to changing environments. Inspired by human learning and adaptation through introspection, we introduce the **SL-AGENT** (**S**ELF-**L**EARNING **A**GENT), a novel self-learning framework that empowers task bots to adapt to such changes by learning from human-bot interactions with minimal or no external supervision.

We begin by presenting the motivation behind our work and explaining how SL-AGENT addresses gaps in existing methods (Section 3.1). Section 3.2 details the architecture of SL-AGENT and describes how it leverages reinforcement learning with an integrated reward model to enable task bots to learn from unlabeled human-bot dialogue logs collected post-deployment. In Section 3.3, we demonstrate the effectiveness of SL-AGENT in adapting to changing environments across four well-studied dialogue tasks, using both automatic and human evaluations. We further conduct in-depth analyses to understand SL-AGENT’s performance in Section 3.4. Finally, Section 3.5 concludes the chapter by summarizing key findings and discussing the limitations of SL-AGENT.

3.1 The Importance of Adaptability

While data-driven approaches have significantly advanced end-to-end task bots [36, 48, 56, 117, 218], their reliance on fixed, annotated corpora leaves them ill-equipped to handle the complexities of real-world interactions. Deployed in dynamic, open environments, these bots often falter when confronted with data that deviates from their training set – encountering unseen user behaviors [88] or task definition extensions [41, 86] (Section 1.2).

Consider the scenario illustrated in Figure 3.1. A user’s casual inquiry about an address initially confuses the system (“Address?”). However, a subsequent, more specific query that resembles the model’s training data (“Can you tell me the address of the restaurant?”) allows the bot to successfully provide the address. This highlights the potential of leveraging post-deployment human-bot interactions – a readily available, dynamic, and information-rich resource [49]– to enable continuous learning and adaptation.

Existing efforts to leverage such interactions for improving task bots in changing environments often rely on costly and potentially sparse human feedback or annotations [24, 88, 143] (Section 2.2). Furthermore, these works primarily focus on dialogue policy optimization or retrieval-based task bots. The potential for automatically adapting generative end-to-end dialogue models remains largely unexplored.

This chapter addresses this gap by introducing SL-AGENT, a novel self-learning framework designed for building task bots in a realistic, changing environment setting. SL-AGENT consists of a neural dialog model and a pre-trained reward model. The dialog model creates responses, and the reward model evaluates their quality. We devise a data augmentation strategy to construct positive and negative examples from the

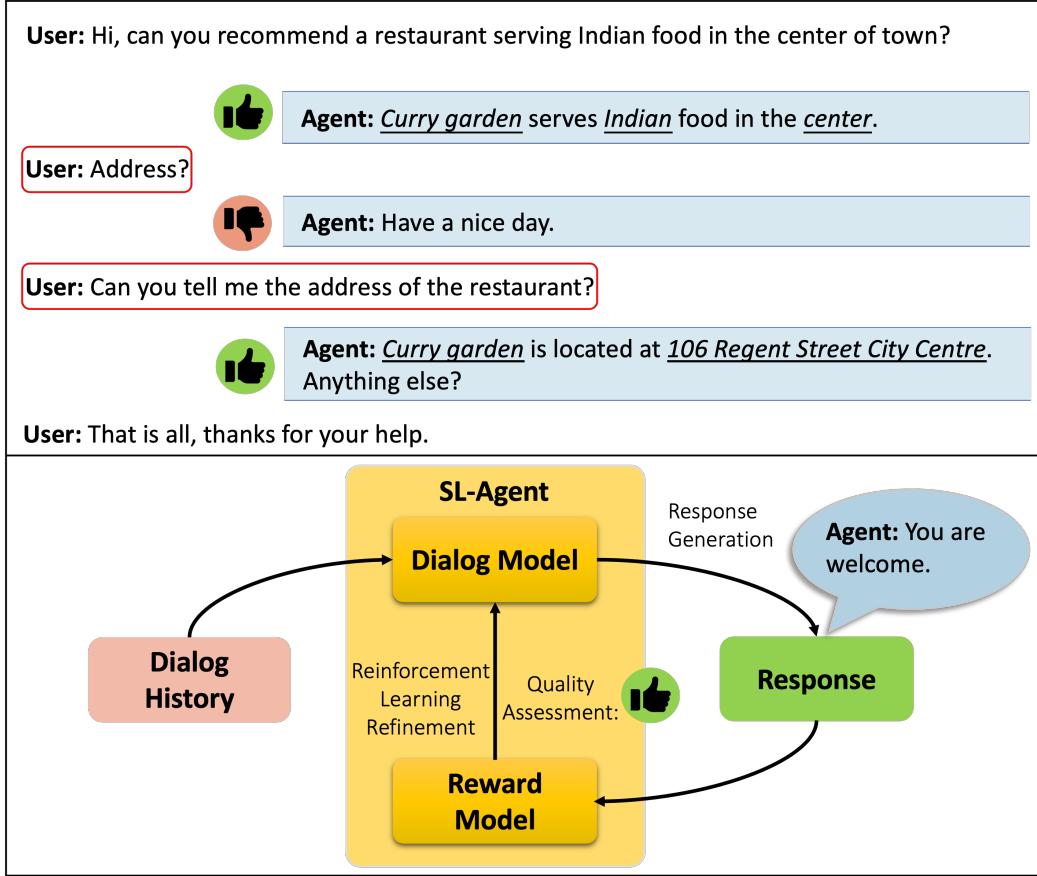


Figure 3.1: Illustration of the proposed SL-AGENT with a human-bot dialog example. (i) The human-bot dialog example, containing an inappropriate response related to unseen user behaviors (upper part). (ii) Demonstration of the refining process in SL-AGENT with the exhibited dialog example (lower part).

existing dialog training corpus, allowing the reward model to assess response quality in unlabeled human-bot dialog logs.

As illustrated in Figure 3.1, SL-AGENT engages in a continuous learning loop: interacting with users, collecting dialogue logs, and leveraging the reward model's feedback to reinforce appropriate responses and discourage inappropriate ones through reinforcement learning. This self-learning process enables the bot to adapt to unseen user behaviors *with zero human annotations*. For task definition extensions, we use

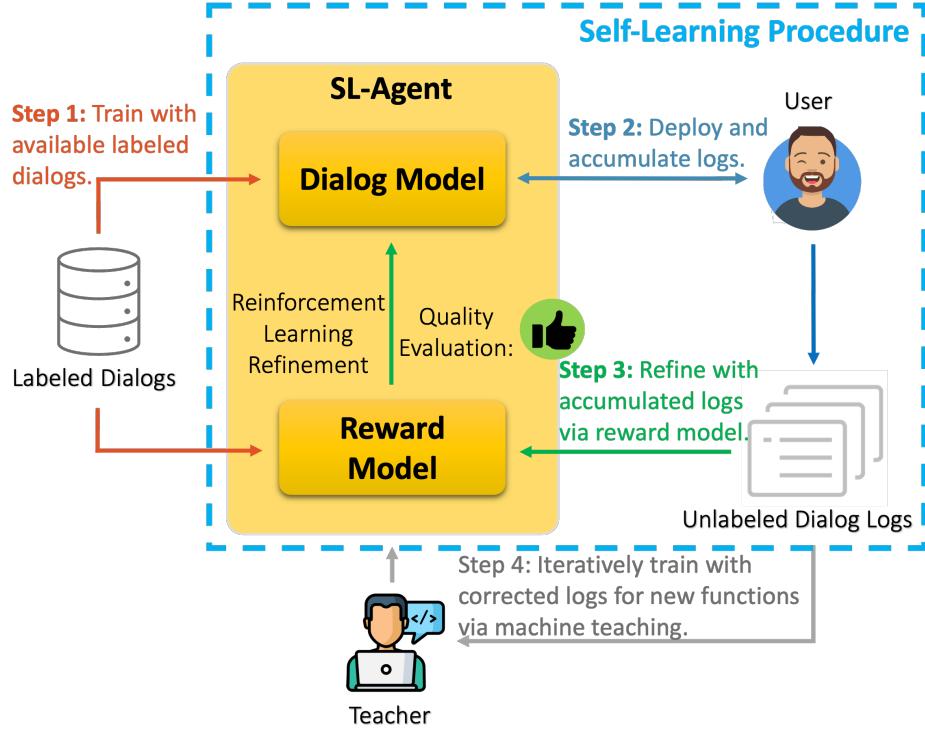


Figure 3.2: The proposed SL-AGENT operates as follows: (i) Fine-tune the bot using available task-specific dialogs. (ii) Deploy the bot online to gather unlabeled human-bot dialog logs. (iii) Refine the dialog model using reinforcement learning with the fine-tuned reward model. (iv) Utilize machine teaching to provide instructions for new functionalities *with minimum human annotations*. The bot rapidly adapts to these new capabilities through self-learning.

machine teaching to correct representative failed dialogs, providing guidance for new functionalities *with minimum human annotations*. The bot rapidly adapts to these new capabilities through self-learning.

3.2 SL-Agent

3.2.1 Overview

As depicted in Figure 3.2, SL-AGENT comprises two components: (i) a **dialog model** for generating responses (Section 3.2.2), and (ii) a **pre-trained reward model** for evaluating the quality of agent responses and providing a reward score to guide the refinement of the dialog model (Section 3.2.3). Specifically, SL-AGENT operates through the following four steps:

Step 1: Policy initialization. The SL-AGENT, including both the dialog model and the pre-trained reward model, is fine-tuned for new tasks using task-specific annotated dialogs to achieve initial response generation and evaluation capabilities, data augmentation strategy, detailed in Section 3.2.3, constructs positive and negative examples from the training corpus, enabling the reward model to assess response quality even within unlabeled human-bot dialog logs.

Step 2: Policy deployment for interactive data collection. The SL-AGENT is deployed online, engaging in conversations with real users and accumulating unlabeled human-bot dialog logs.

Step 3: Self-refinement using reinforcement learning. Leveraging the collected human-bot dialog logs, the dialog model undergoes refinement through reinforcement learning (Section 3.2.4). The fine-tuned reward model evaluates the quality of responses during this process. By replicating successful patterns and avoiding those leading to undesirable outcomes, the SL-AGENT continuously improves its performance.

Step 4: Policy improvement via external feedback for iterative self-learning. For

task definition extensions, machine teaching is employed to correct representative failed dialogs and provide instructions on handling new functions (Section 3.2.5). The SL-AGENT then integrates this feedback, further enhancing its performance through this iterative self-learning process.

3.2.2 Dialog Model

SL-AGENT is a general framework that is compatible with any generative end-to-end dialog models [48, 56, 117]. In this chapter, we employ SOLOIST [117], a pre-trained end-to-end dialog model, resulting in an agent termed SL-SOLOIST.¹

We briefly review SOLOIST for completeness. SOLOIST formulates the end-to-end dialog generation as a sequence generation problem, by sequentially concatenating the inputs and outputs of 4 dialog modules (i.e., NLU, DST, POL, NLG) in a typical dialog system. Each dialog turn is represented as:

$$x = (s, b, c, r), \quad (3.1)$$

where s is the entire dialog history, b is the annotated belief state, c refers to DB state fetched from database, and r is the delexicalized agent response. SOLOIST employs a Transformer-based model with parameters θ_D to characterize the sequence generation probability $p_{\theta_D}(x)$. Initialized with GPT-2 [127], the model is pre-trained on large-scale annotated dialog corpora, and then fine-tuned with limited task-specific dialogs.

Synthetic Dialog Construction. To improve the capability of the dialog model for identifying user behaviors with unseen slot values, we propose to synthesize dialog

¹In this chapter, SL-AGENT refers to the proposed framework and SL-SOLOIST is an instance of it, which utilizes SOLOIST as its dialog model.

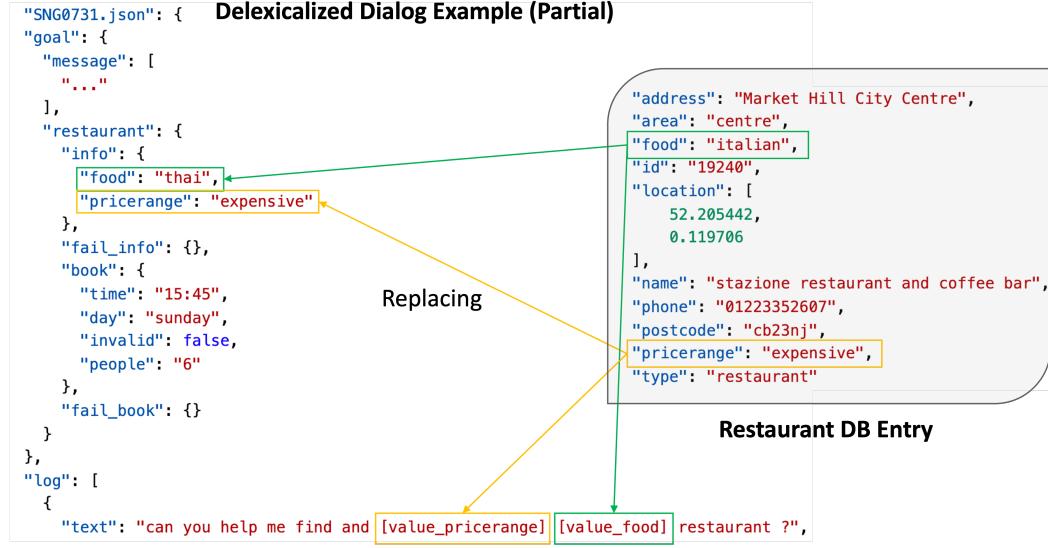


Figure 3.3: Illustration of synthetic dialog construction. Slot values in the delexicalized dialog goal (top left) and dialog log (bottom left) are replaced with values from a restaurant database entry (right) to generate a fully lexicalized synthetic dialog.

examples by exhausting database (DB) values and substitute corresponding slot values of in the training set (Figure 3.3). Specifically, for each dialog turn x , we replace slot values in the utterances and user goal with corresponding new values of the randomly sampled DB entry.

3.2.3 Reward Model

Human-bot dialog logs accumulated after deployment often contain previously unseen user inputs with new language patterns and uncertain user goals. To enable the dialog model to adapt to these dynamic scenarios, we introduce a reward model that evaluates the quality of agent responses. This model assigns a reward score to each response, providing a positive reward for appropriate responses and a negative reward for inappropriate ones.

We formulate the quality evaluation problem as a binary classification task. Dia-

log responses are jointly determined by the dialog history, generated belief state, and fetched DB state. Therefore, given the training data D (annotated with belief states and DB states), we build a turn-level reward model R , which is parameterized by a Transformer θ_R with the input dialog turn sequence x , defined as Equation 3.1 to characterize the classification probability: $p_{\theta_R}(x) = p_{\theta_R}(s, b, c, r)$.

The reward model R is trained using contrastive objective to discriminate between an appropriate response (i.e., positive example x) and an inappropriate response (i.e., negative example \hat{x}), given the dialog history. Specifically, for each dialog turn, we construct several positive examples $\{x_m\}_{m=1}^M$ and negative examples $\{\hat{x}_n\}_{n=1}^N$ based on the sequence x , to add the relevance of real-world scenarios and endow the reward model with the capability of evaluating the response quality. Then a binary classifier is applied on top of the output sequence representation from the Transformer to discriminate between a positive example x ($y = 1$) and a negative example \hat{x} ($y = 0$). The training objective for a single example in the training set D is defined as:

$$L_{\theta_R} = - \left[y \sum_{m=1}^M \log p_{\theta_R}(x_m) + (1-y) \sum_{n=1}^N \log (1 - p_{\theta_R}(\hat{x}_n)) \right] \quad (3.2)$$

where M represents the number of positive examples, and N represents the number of negative examples.

Positive Examples. For each dialogue turn, we utilize two types of user utterances to train the reward model:

1. **Original user utterance.** We use the original user utterance from the training set D to train the model in recognizing appropriate responses to specific user behaviors.

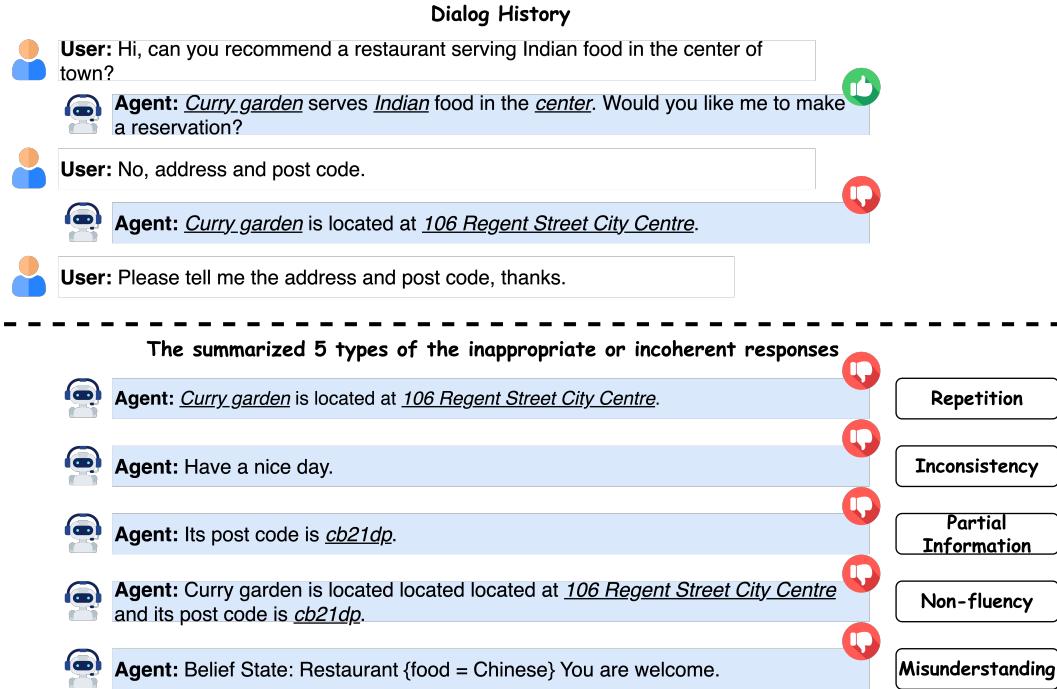


Figure 3.4: The summarized five types of dialog turns featuring inappropriate or incoherent responses. (i) Dialog history (top). (ii) Five categories of inappropriate or incoherent responses corresponding to the given dialog history (bottom).

2. Paraphrased user utterances. To enhance the model’s robustness to diverse language patterns, we generate paraphrased versions of the original user utterance using back-translation techniques [31].

Negative Examples. To guide the reward model in recognizing undesirable responses, we define five categories of dialog turns with inappropriate or incoherent responses (Figure 3.4) based on an analysis of 200 human-bot dialog logs from the DSTC8 Track 1 challenge evaluation platform [76] (these logs include evaluation scores and comments from Amazon Mechanical Turk workers). For each dialog turn in the training data D , we construct corresponding negative examples \hat{x} by applying these categories:

1. Repetition. The dialog model failed to understand the user’s repeated query and

generated the same response twice. (Repeating the response from the previous turn.)

2. **Inconsistency.** The dialog model generated an incoherent response. (Randomly sampling a response from the dataset D to replace the original response .)
3. **Partial information.** The dialog model partially understood user request and answered incompletely. (For those user utterances with multiple slots request, randomly dropping a slot answer in the original response.)
4. **Non-fluency.** The dialog model generated a non-fluent response. (Randomly repeating some word tokens in the original response.)
5. **Misunderstanding.** The dialog model generated the incoherent belief state and response. (Randomly sampling a belief state and response from the dataset D to replace the original belief state and response.)

To avoid the need for extensive annotated task-specific data, we adopt a pre-training and fine-tuning paradigm for our reward model. We pre-train the reward model on the large-scale annotated Schema dataset [133] and then fine-tune it on annotated task-specific data. This approach leverages the knowledge learned from a broader dialog corpus while adapting to the nuances of our specific task.

3.2.4 Refine with Reinforcement Learning

The interactions between the agent and users can be modeled as a sequential decision problem (introduced in Section 2.2). As such, the dialog model can be refined via the REINFORCE algorithm [189]. The policy is the trained dialog model $p_{\theta_D}(x)$,

Algorithm 1 Self-learning-based RL refining framework.

Input:

Training examples D in the form of dialog turns;
Trained agent with dialog model $p_{\theta_D}(x)$ and reward model $p_{\theta_R}(x)$.

Output:

Refined agent with updated dialog model $p_{\theta_D^*}$.

- 1: **while** not converged **do**
 - 2: Randomly sample a dialog turn, i.e. token sequences of dialog history s ;
 - 3: Run dialog model p_{θ_D} on dialog history $x = (s)$ to generate belief state \hat{b} ;
 - 4: Retrieve DB state \hat{c} from a database using generated belief state \hat{b} ;
 - 5: Sample corresponding response r based on dialog history s , belief state \hat{b} and DB state \hat{c} ;
 - 6: Use the reward model to predict the quality of the belief state and response with reward score,
 $R(s, \hat{b}, \hat{c}, r)$;
 - 7: Calculate the loss according to Equation 3.3;
 - 8: Update the parameters of the dialog model,
 $\theta_D \leftarrow \theta_D + \alpha \nabla_{\theta_D} L_{\theta_D}$.
 - 9: **end while**
-

the initial state is the dialog history s , and the action space corresponds to the vocabulary set V . The reward perceived by the dialog model is $R(s, b, c, r)$ from the reward model. The parameters θ_D are updated by maximizing the cumulative reward score. The refining procedure is described in detail as follows:

For each RL episode, we randomly sample a dialog turn with dialog history and delexicalized response. We run the dialog model to generate belief state \hat{b} , based on the input dialog history sequence s . At each time step t , we sample a token \hat{b}_t according to the model distribution, where the logits' distribution of the model is first filtered using Nucleus (top-p) filtering [54], then redistributed via softmax function. Then we retrieve DB state \hat{c} from the database using \hat{b} , and sample the delexicalized response sequence r following same sampling procedure, based on the token sequence (s, \hat{b}, \hat{c}) . Note that the delexicalized response is given as part of the input. Then we feed the

concatenation of dialog history s , generated belief state \hat{b} , retrieved DB state \hat{c} and the response r , i.e., (s, \hat{b}, \hat{c}, r) into the reward model $p_{\theta_R}(x)$ to obtain the reward score $R(s, \hat{b}, \hat{c}, r)$. The positive reward is 1, negative reward is -1. The training objective for a single example is represented as:

$$\begin{aligned} L_{\theta_D} = & - \sum_{t=1}^{T_b} \log p_{\theta_D}(\hat{b}_t | \hat{b}_{<t}, s) \times R(s, \hat{b}, \hat{c}, r) \\ & - \sum_{t=1}^{T_r} \log p_{\theta_D}(r_t | r_{<t}, \hat{b}, \hat{c}, s) \times R(s, \hat{b}, \hat{c}, r), \end{aligned} \quad (3.3)$$

where the length of generated belief state and input delexicalized response are T_b , T_r , respectively. Algorithm 1 summarizes the self-learning-based RL refining framework for refining the dialog model.

3.2.5 Minimum Annotations via Machine Teaching

Turn	Utterance
1	User: I want a restaurant on the <i>west</i> end of town.
2	System: Ok. Are you looking for any particular type of food?
3	User: <i>I don't care.</i>
4	System: Are you looking for a particular price range?
5	User: A <i>moderately priced</i> restaurant would be good.
6	System: How about Saint Johns Chop House? They serve moderately priced British food on the west side of town.
7	User: Can i please have their address and phone number?
8	System: Saint Johns Chop House's address is 21-24 Northampton Street, their phone number is 01223353110.
9	User: Does the restaurant offers delivery service? If it does, what's the delivery fee?
10	System: Yes, it does, and the delivery charges 6 pounds.
11	User: Thank you, goodbye.
12	System: Thank you for using our system. Goodbye.

Table 3.1: A dialog example on domain extension. Task bots need learn to provide pertinent responses concerning the expanded delivery service in subsequent interaction turns, adapting to the continually evolving user or enterprise demands.

Extending the capabilities of task bots to handle queries about new functions or

Slot	Value
”address”	”21 - 24 Northampton Street”
”area”	”west”
”food”	”british”
”id”	”14810”
”location”	[52.21031, 0.11381]
”name”	”saint johns chop house”
”phone”	”01223353110”
”postcode”	”cb30ad”
”pricerange”	”moderate”
”type”	”restaurant”
”delivery”	”yes”
”delivery fee”	”6 pounds”
”dish”	”Beef Wellington”
”start_time”	”10:30 am”
”end_time”	”22:40 pm”

Table 3.2: A Restaurant-Ext DB entry. The newly introduced slot-value pairs relevant to the extended functionality are highlighted.

tasks requires incorporating additional knowledge into the bot. This includes introducing new slot-value pairs, action templates, and potentially modifying the dialogue flow. Tables 3.1 and 3.2 illustrate an example dialogue and its corresponding database entry, highlighting the need for such extensions. Machine teaching is shown to be an efficient approach to train task bots [146, 149, 187]. In this chapter, we implement machine teaching via Conversational Learner (CL) [146]. The teaching process is conducted in three steps: (i) The trained task bot is deployed online to fulfill the given goals by interacting with real users, leaving a handful of human-bot dialog logs. (ii) Human experts select a few representative failed dialogs to construct training examples with new functions by adding new action templates, introducing new slot-value pairs, correcting inappropriate responses and annotations (i.e., belief states). (iii) The deployed task bot (i.e., both the dialog model and reward model) is trained on these training examples to handle new functions.

Domain	Attraction	Train	Hotel	Restaurant
#Train	50	50	50	50
#Valid	50	50	50	50
#Test	100	200	200	200

Table 3.3: Data statistics of four single-domain dialog datasets [12, 117].

3.3 Experiments

In this section, we first describe how we design evaluations on changing environments. Then we introduce the experiments we conduct on four well-studied dialog tasks using both automatic and human evaluations.

3.3.1 Setup

We validate the efficiency and flexibility of proposed SL-AGENT on four different end-to-end dialog tasks using Multiwoz single-domain dialog datasets [12], reorganized by Peng et al. [117]. Data statistics are shown in Table 3.3.

Automatic Evaluation Metrics. We report the results using the same automatic evaluation metrics following Budzianowski et al. [12]: (i) Inform(%) evaluates whether the agent returns an appropriate entity. (ii) Success(%) judges whether the agent correctly answers all requested attributes. (iii) BLEU(%) measures the word overlap of the generated response against human response. (iv) Combined(%) assesses the overall quality, which is defined as: Combined = (Inform + Success) × 0.5 + BLEU.

Human Evaluation Metrics. Following the same evaluation protocol in the DSTC9 Track 1 challenge [46], we conduct human evaluations to judge the agent quality. For each dialog session, Amazon Mechanic Turks are presented with a goal and instruc-

tions, then they are required to converse with agent to achieve the goal via natural language. At the end of each dialog session, Turkers are required to assess the overall dialog quality using the following five metrics: (i) Success w/o g(%) judges whether the agent completes the task. (ii) Success w/ g(%) judges whether the agent completes the task and provides matched slot values against the database record. (iii) Understanding (1-5) measures the understanding correctness of user utterances. (iv) Appropriateness (1-5) indicates the appropriateness, naturalness, and fluency of an agent response. (v) Turns reports the average number of dialog turns for successful dialog sessions.

Compared Methods. To demonstrate the effectiveness of SL-AGENT, we use SOLOIST as the dialog model to compare the performance of different methods.¹

- **SOLOIST₅** is trained with 5 labeled dialogs, randomly sampled from the train set.
- **SOLOIST_S** is trained using synthetic dialogs constructed from the 5 labeled dialogs used for training SOLOIST₅.
- **SOLOIST+PARG** is trained on SOLOIST_S with paraphrased dialogs [31, 39] constructed from the 5 labeled dialogs, which is the data-augmentation baseline.
- **SOLOIST-OA** is refined with unlabeled human-bot dialog logs based on SOLOIST_S using the session-level reward of task success from online activate reward model (trained using the same training examples as SOLOIST₅) and partially queried session-level human feedback score [156].

¹Current SOTA task-oriented dialog models share similar input-output pairs and training objectives as SOLOIST.

- **SL-SOLOIST** (Ours) is refined with unlabeled human-bot dialog logs based on SOLOIST_S using proposed SL-AGENT, where the pre-trained reward model is fine-tuned using the same training examples as SOLOIST_S.
- **SOLOIST-TH** is refined with unlabeled human-bot dialog logs based on SOLOIST_S using queried turn-level human feedback score, which is an upper bound.
- **SOLOIST₅₀** is trained with whole 50 labeled dialogs, which can be regarded as the result of sufficient human corrections, i.e., the highest bound.

User : I would like to find an expensive restaurant that serves Chinese food. System : Sure, which area do you prefer ? User : How about in the north part of town. [BOS]
 Belief State : Restaurant { pricerange = expensive, food = Chinese, area = north } [EOB] DB : Restaurant 2 match [EOD]
 System: The [restaurant_name] is a great [value_food] restaurant. Would you like to book a table there ? [EOS]

Figure 3.5: Illustration of the training example, i.e., the processed dialog turn in the training data.

Implementation Details. We experimented with several Transformer-based models for the reward model and found that GPT [127], enhanced with auxiliary generation tasks, outperformed others. Therefore, we implemented our proposed reward model using GPT-117M and a multi-task training objective. Figure 3.5 illustrates the training example construction process. We tokenize dialog turn sequences using byte pair encodings [141] and delexicalize responses by replacing slot values with corresponding special tokens [73]. We pre-trained the reward model for 10 epochs on the Schema dataset [133], which contains 22,825 dialogs across 17 domains. We used the Huggingface PyTorch Transformer library [190] for implementation. The pre-training utilized

Model	Attraction			Train			Hotel			Restaurant		
	Info.	Succ.	BLEU	Info.	Succ.	BLEU	Info.	Succ.	BLEU	Info.	Succ.	BLEU
SOLOIST ₅	27.00	14.00	4.07	72.73	32.32	5.43	25.00	3.50	2.93	26.50	2.00	4.71
SOLOIST _S	60.00	33.00	8.14	73.74	54.55	6.94	56.00	29.50	7.05	62.50	41.50	7.33
SOLOIST+PARG	60.00	32.00	8.83	75.25	56.06	8.45	58.00	29.00	7.71	64.00	42.00	9.17
SOLOIST-OA	61.00	36.00	8.66	74.75	55.05	7.58	56.50	29.00	7.14	64.50	42.50	8.56
SL-SOLOIST	64.00	40.00	8.99	75.76	61.62	10.97	60.50	39.50	8.34	75.00	44.50	10.60
SOLOIST-TH	66.00	41.00	9.01	77.27	62.87	10.70	60.00	42.50	9.82	70.50	46.00	11.76
SOLOIST ₅₀	86.00	65.00	12.90	80.81	64.65	9.96	74.50	43.50	8.12	81.00	55.50	12.80

Table 3.4: End-to-end evaluation results on four tasks. The forth to sixth rows indicate the results of refining with 45 simulated (unlabeled) human-bot dialog logs, based on SOLOIST_S. SOLOIST₅₀ is quoted from Peng et al. [117]. Info.: Inform, Succ.: Success. (SL-SOLOIST significantly outperforms all baselines in mean with $p < 0.01$ based on Combined.)

two 24G Nvidia P40 GPUs with a mini-batch size of 8, a learning rate of 5e-5, and the Adam optimizer [67]. Training examples were truncated or padded to a maximum length of 500 tokens.

Both the pre-trained reward model and dialog model (pre-trained SOLOIST) were fine-tuned for 20 epochs on a limited number of labeled task-specific dialogs. We used a top-p sampling of 0.5 for all models during fine-tuning and refinement. Gradient clipping with a maximum norm of 1 was applied during parameter updates. The dialog model was fine-tuned on a single 24G Nvidia P40 GPU with a batch size of 1 and a learning rate of 5e-6 until convergence on the validation set. During testing, Nucleus filtering with a top-p value of 0.5 was employed for decoding.

3.3.2 Results of Unseen User Behaviors

Simulation Evaluation. Deploying trained conversational agents with real users for data collection is resource-intensive, especially during the experimental phase. To this end, we introduce a novel simulation setting to effectively simulate unseen user behav-

iors.¹ We leverage a small subset (5 examples) from the training set as labeled data to train a baseline task bot, encompassing both dialogue and reward models. The remaining 45 dialogues, containing unseen user behaviors and goals, form the basis for our simulation. These 45 dialogues are transformed into unlabeled, imperfect human-bot dialogues by introducing noise through response corruption. Belief state annotations are not used in this process. This approach allows us to simulate realistic, unseen user interactions. We utilize these unlabeled dialogues to refine our proposed model, SOLOIST_S, resulting in three variations: SOLOIST-OA, SL-SOLOIST, and SOLOIST-TH. This simulation framework facilitates a comprehensive analysis of SL-AGENT in a cost-effective and reproducible manner. Table 4.2 presents the end-to-end evaluation results across four distinct tasks. Our key findings are:

SL-AGENT enables the bot to automatically adapt to changing environments. SOLOIST_S outperforms SOLOIST₅ over all evaluation metrics on all tasks by a significantly large margin, which shows the effectiveness of proposed synthetic dialog construction for adapting to unseen user behaviors caused by unseen slot values. SL-SOLOIST outperforms SOLOIST+PARG over all the metrics, which demonstrates the higher efficiency of directly learning from human-bot dialog logs.

The pre-trained reward model in SL-AGENT is effective at predicting turn-level response quality. SL-SOLOIST significantly outperforms SOLOIST-OA and achieves comparable performance to SOLOIST-TH, which represents the upper bound by utilizing turn-level human feedback scores. This underscores the effectiveness of the

¹ While user simulators offer a potential alternative, they face significant limitations in our dynamic environment: (i) Agenda-based simulators necessitate complex rule design and extensive domain expertise, making them challenging to develop and maintain. (ii) Model-based simulators rely heavily on labeled data and often struggle to generate realistic user behaviors beyond those observed in their training corpus, limiting their ability to simulate novel interactions.

Model	Attraction			Train			Hotel			Restaurant		
	Info.	Succ.	BLEU									
SOLOIST _S	60.00	33.00	8.14	73.74	54.55	6.94	56.00	29.50	7.05	62.50	41.50	7.33
SOLOIST-OA	63.00	34.00	8.66	77.78	55.05	8.13	58.50	30.00	7.08	63.00	42.00	10.03
SL-SOLOIST	70.00	36.00	8.68	78.28	60.10	9.06	62.00	33.50	7.39	70.00	45.00	10.93
SOLOIST-TH	68.00	40.00	9.01	76.77	62.63	9.55	62.50	35.50	7.83	70.50	47.50	11.36

Table 3.5: Automatic evaluation results on four tasks in Real-Scenario Setting. The first row refers to previously reported SOLOIST_S. The last three rows refer to refining with 30 real (unlabeled) human-bot dialog logs based on SOLOIST_S. Info.: Inform, Succ.: Success. (SL-SOLOIST significantly outperforms all baselines in mean with $p < 0.01$ based on Combined.)

pre-trained reward model in accurately assessing response quality.

Real-Scenario Evaluation. Simulation setting allows effortless experimental studies to validate the effectiveness of SL-AGENT. However, the results are likely biased. Therefore, in the real-scenario setting, we deploy SOLOIST_S online and recruit human users to converse with it. We collect 30 real (unlabeled) human-bot dialog logs to refine SOLOIST_S, resulting in the agent SOLOIST-OA, SL-SOLOIST, SOLOIST-TH. Table 3.5 presents the evaluation results across four tasks in the real-world setting. Our findings demonstrate that:

SL-AGENT excels in automatic adaptation during real-world deployment. SL-SOLOIST, refined using our proposed SL-AGENT framework, consistently outperforms other methods across all evaluation metrics and tasks. Furthermore, SL-SOLOIST achieves comparable performance with SOLOIST-TH, even achieves better performance on certain metrics. We conclude that the results of real-scenario evaluation and simulation evaluation are consistent, confirming that SL-SOLOIST enables effective self-learning after deployment by learning from interactions.

Model	Restaurant-Ext			
	Inform	Success	BLEU	Combined
SOLOIST _S	54.00	0.00	6.42	33.42
SOLOIST _S +TEACH	64.00	18.00	9.34	50.34
SL-SOLOIST+TEACH	68.00	24.00	11.76	57.76
SOLOIST-TH+TEACH	68.50	26.00	11.88	59.13

Table 3.6: Automatic evaluation results on task definition extensions. (Difference in mean is significant with $p < 0.01$ based on Combined.)

3.3.3 Results of Task Definition Extensions

We follow the domain extension experiment setting in Lipton et al. [86] to assess the ability of SL-SOLOIST to quickly handle task definition extensions. We extend existing Restaurant, denoted as Restaurant-Ext, with additional functions by introducing 4 new slots, i.e., `[restaurant_dish]`, `[value_price]`, `[start_time]`, `[end_time]` in added dialog turns, and corresponding values for each DB entry. The first slot is about the restaurant’s signature dish, and the last three are related to delivery service. We leverage Conversational Learner (CL) [146], a practical machine teaching tool, to visualize and select dialogs for constructing training examples on the Restaurant-Ext domain by providing corrections and introducing new slots. Finally, 10 examples are obtained through machine teaching for training, 50 for validating and 50 for testing. We fine-tune the dialog model SOLOIST_S and the previously trained reward model,¹ using 10 corrected dialogs, resulting the agent denoted as SOLOIST_S+TEACH. Then, SOLOIST_S+TEACH is deployed to converse with real human to collect 20 real (unlabeled) human-bot dialog logs, which are then used to refine itself, resulting in SL-SOLOIST+TEACH. To better demonstrate the effectiveness of SL-AGENT, we also

¹The reward model used for obtaining SL-SOLOIST in the Table 4.2. It is trained with 5 labeled dialogs in the train set.

Model	Restaurant			
	Inform	Success	BLEU	Combined
SOLOIST _S	62.50	41.50	7.33	59.33
SL-SOLOIST	75.00	44.50	10.60	70.35
SL-SOLOIST ₊₂₀	75.00	52.00	11.89	75.39

Table 3.7: End-to-end evaluation results of Policy Improvement in the Restaurant domain. SL-SOLOIST₊₂₀ refer to continually refining with 20 real (unlabeled) human-bot dialogs based on SL-SOLOIST (reported in Table 4.2).

report the result of SOLOIST-TH+TEACH, which is refined using the turn-level human feedback score. The evaluation results are presented in Table 3.6.

Enhanced by machine teaching, SL-AGENT enables flexible adaptation to new tasks. We observe that SOLOIST_S has zero success rate, which is predictable as it does not have any knowledge of the new functions. SOLOIST_S+TEACH outperforms the baseline by 17 points in terms of Combined score, which exhibits the effectiveness of machine teaching for handling new functions. SL-SOLOIST+TEACH lifts the Combined score by approximately 7 points, achieving comparable performance with SOLOIST-TH+TEACH (i.e., refining using turn-level human feedback score). The results demonstrate that SL-SOLOIST is able to adapt to new tasks and continually improve itself by automatically learning from the interactions.

3.3.4 Continual Policy Improvement

To demonstrate the effectiveness of SL-AGENT for continually learning from collected human-bot dialog logs, we deploy SL-SOLOIST online and recruit human users to converse with it to achieve the assigned user goal. We collect 20 real human-bot dialog logs to refine SL-SOLOIST, resulting in the agent SL-SOLOIST₊₂₀. (When refining the SL-SOLOIST, we do not use the knowledge about the user’s goal. The

response quality is judged by the reward model in SL-SOLOIST.)

The evaluation results on Restaurant are shown in Table 3.7. We observe that SL-SOLOIST_{+20} refined with 20 real (unlabeled) human-bot dialogs outperforms SL-SOLOIST by approximately 5 points in terms of Combined score. We conclude that SL-SOLOIST enables continual self-learning after deployment by automatically learning from interactions.

3.4 In-Depth Analyses

In this section, we first examine the underlying reasons behind the effectiveness of the proposed reward models in SL-AGENT. We then present qualitative analyses and human evaluation results on the generated dialogs to further support the effectiveness of SL-AGENT.

3.4.1 Impact of Various PLMs and Training Objectives on Reward Models

To analyze the impact of different PLMs and the multi-task training objective on the reward model, we conduct ablation studies on the Restaurant domain. We compare several popular PLMs: BERT [26], RoBERTa [92], and SL-SOLOIST. All models share identical pre-training and fine-tuning procedures. However, while BERT and RoBERTa are trained solely on quality prediction, SL-SOLOIST benefits from a multi-task learning approach.

Table 3.8 reveals that RoBERTa outperforms BERT, suggesting that its enhanced pre-training leads to better performance. Notably, GPT-2 (upon which SL-SOLOIST is

Reward model	Restaurant			
	Inform	Success	BLEU	Combined
GPT-2	67.00	41.50	9.30	63.55
BERT	68.00	42.50	9.55	64.80
BERT-Large	66.00	44.00	11.09	66.09
RoBERTa	72.00	45.00	9.23	67.73
RoBERTa-Large	69.50	46.50	10.20	68.20
SL-SOLOIST	75.00	44.50	10.60	70.35

Table 3.8: Ablation study results on using different PLMs for reward models in Restaurant domain. The first five rows indicate evaluation results of fine-tuned GPT-2, BERT, BERT-Large, RoBERTa, RoBERTa-Large, respectively. The last row refers to previously reported SL-SOLOIST. (Difference in mean is significant with $p < 0.01$ based on Combined.)

built) exhibits significantly worse performance when trained only on quality prediction. This difference likely stems from the bidirectional Transformer encoder employed by BERT and RoBERTa, which enables richer context capture compared to GPT-2’s unidirectional encoder. Importantly, SL-SOLOIST consistently surpasses all other methods across all metrics. This result highlights the effectiveness of multi-task learning for optimizing the reward model, allowing it to learn more nuanced representations and achieve superior performance.

3.4.2 Qualitative Analysis

To further illustrate the effectiveness of SL-AGENT in handling unseen user behaviors, we present a case study comparing the performance of SOLOIST_S and SL-SOLOIST on the same task. Figure 3.6 depicts two interactive examples where the same user interacts with both agents to achieve the same goal.

(a) The upper example showcases the interaction between the user and SOLOIST_S. While both agents perform similarly in the initial four turns, successfully recommending a suitable restaurant, SOLOIST_S falters when the user requests the restaurant’s

- 1 **User:** I am looking for a moderately priced restaurant.
- 2 **SOLOISTs:** There are 31 options. I recommend Anatolia.
- 3 **User:** Also, it should be in the centre.
- 4 **SOLOISTs:** Anatolia looks like a good place. They serve Turkish.
- 5 **User:** What's **the number**?
- 6 **SOLOISTs:** Would you like to try Anatolia? It is located in the centre.
- 7 **User:** **Number**?
- 8 **SOLOISTs:** Would you like to try Anatolia?
- 9 **User:** What is **the phone number**?
- 10 **SOLOISTs:** The phone number is 01223362372.

(a) An interactive example between user and SOLOISTs

- 1 **User:** I am looking for a moderately priced restaurant.
- 2 **SL-SOLOIST:** There are 31 options. I recommend Anatolia.
- 3 **User:** Also, it should be in the centre.
- 4 **SL-SOLOIST:** Anatolia looks like a good place. They serve Turkish.
- 5 **User:** What's **the number**?
- 6 **SL-SOLOIST:** The phone number is 01223362372.

(b) An interactive example between user and SL-SOLOIST

Figure 3.6: Two interactive examples. (a) An interactive example between user and SOLOIST_S. (b) An interactive example between user and SL-SOLOIST.

phone number (“what’s the number?”) in the fifth turn. It fails to grasp the user’s intent and continues to offer recommendations, leading to an incoherent dialogue flow. The user is forced to repeatedly request the phone number in subsequent turns. (b) In contrast, the lower example demonstrates a smoother interaction between the user and SL-SOLOIST. When the user inquires about the phone number, SL-SOLOIST accurately understands the request and provides the information immediately.

This case study highlights a crucial advantage of SL-AGENT: its ability to adapt to unseen user behaviors automatically. While SOLOIST_S struggles to deviate from its

Model	Restaurant				
	SR w/o g ↑	SR w/ g ↑	Understanding ↑	Appropriateness ↑	Turns ↓
SOLOIST _S	31.82	29.54	3.86	4.13	10.00
SOLOIST-OA	33.42	30.86	3.89	4.12	9.97
SL-SOLOIST	43.10	36.21	3.97	4.13	9.89

Table 3.9: Human evaluation results. SR w/o g: Success rate without grounding, SR w/ g: Success rate with grounding.

pre-defined script, SL-SOLOIST demonstrates a more dynamic and adaptable conversational flow, effectively handling the user’s unexpected request. This adaptability is crucial for developing robust and user-friendly conversational agents.

3.4.3 Interactive Human Evaluation

Corpus-based evaluation is conducted using automatic evaluation metrics, which are rough proxies for agent response quality. Furthermore, automatic evaluation results may not adequately reflect the capability of dialog systems for helping users complete tasks in the real world, as real user inputs are more dynamic, complex, even with noise. Therefore, we conduct human evaluations to evaluate the performance of SOLOIST_S, SOLOIST-OA, SL-SOLOIST interacting with human users, following the evaluation protocol in DSTC9 track 1 challenge [46], with 100 dialogs gathered for analysis, respectively.

The human evaluation results on Restaurant domain are presented in Table 3.9. The results show that SL-SOLOIST achieves the best performance over all the metrics, which are consistent with the automatic evaluation results. The significant improvement on two success rate metrics, especially success rate with grounding, verifies the effectiveness of the SL-AGENT for refining the dialog agent after deployment without additional human annotations, as it more adequately reflects the system’s capability for

completing tasks in real scenarios.

3.5 Chapter Summary

This chapter addresses the critical challenge of building truly adaptive task-oriented dialogue systems—bots capable of evolving and succeeding in dynamic environments with minimal human intervention. We present SL-AGENT, a novel self-learning framework that enables task bots to learn directly from ongoing interactions, significantly reducing the need for manual annotation. By leveraging reinforcement learning guided by a pre-trained reward model, SL-AGENT allows bots to adapt to shifting conversational contexts using unlabeled human-bot dialogue logs.

Our experiments across four diverse dialogue tasks demonstrate the effectiveness of SL-AGENT in achieving automatic adaptation. However, the pursuit of fully autonomous learning remains ongoing. Although SL-AGENT reduces dependency on human input, it still requires expert-provided demonstrations to initiate the learning process.

This insight motivates the next chapter, where we explore advanced techniques aimed at further minimizing human involvement. Our overarching goal is to empower task bots to acquire new task knowledge with near-zero human effort and minimal computational overhead, pushing the boundaries of scalable and autonomous dialogue system development.

Chapter 4

Schema-Guided LLM Prompting for Extensibility

This chapter explores how to empower task bots to seamlessly adapt to new tasks and domains. Inspired by the human ability to acquire knowledge from high-level principles, we introduce SGP-TOD (**S**chema-**G**uided **P**rompting for building **T**ask-**O**riented **D**ialog systems). This novel approach leverages LLMs and symbolic knowledge, represented as predefined task schemas, to facilitate effortless task bot creation.

We begin in Section 4.1 by motivating the development of SGP-TOD and highlighting its advantages over existing methods. Section 4.2 details the architecture of SGP-TOD, which comprises an LLM, a Dialog State Tracking (DST) Prompter, and a Policy Prompter. By leveraging a predefined task schema—including belief instructions and dialog policy—SGP-TOD instructs the LLM to generate appropriate responses for new tasks without requiring any training data. Extensive experiments on the MultiWOZ, RADDLE, and STAR datasets (Section 4.3) demonstrate that SGP-TOD achieves state-of-the-art zero-shot performance, even surpassing few-shot baselines. We then delve into the reasons behind the effectiveness of SGP-TOD in Section 4.4, presenting ablation results, qualitative examples, and human evaluations. Finally, Section 4.5 summarizes the key findings and concludes the chapter.

4.1 The Necessity of Extensibility

A common approach to building task-oriented bots involves fine-tuning pre-trained language models on task-specific annotated datasets [56, 118, 158]. While this approach can yield strong performance, it has a critical limitation: adapting to new tasks or functionalities requires substantial annotated data and retraining, making it costly and labor-intensive (see Section 1.2). The advent of LLMs like ChatGPT [106] and GPT-4 [107] have revolutionized NLP applications [180, 185] with their remarkable conversational skills [125], instruction-following abilities [109] and zero-shot generalization capabilities [19, 57]. This progress begs the question: can LLMs be effectively utilized for building task bots with minimum human effort?

Recent work [61] has explored using LLMs for rapid task bot development through in-context learning [9, 96], as discussed in Section 2.3. While promising, its efficacy hinges on the quality of in-context exemplars, which often struggle to provide comprehensive information necessary for effective dialogue task completion [30, 89, 227].

To address this challenge, this chapter introduces SGP-TOD (Figure 4.1), a novel schema-guided prompting method for rapidly building task bots. SGP-TOD leverages symbolic knowledge [18, 105] in the form of task schemas [100, 104] to provide LLMs with a comprehensive blueprint of the task. This schema, a concise symbolic representation of the task, comprises: (i) task-specific ontology containing all slots and their appropriate values [11]; and (ii) a dialog flow explicitly outlining fundamental interaction patterns [120] (as briefly recapped in Section 2.3). SGP-TOD integrates this predefined task schema and dialogue context into prompts through two specialized prompters: a DST Promter and a Policy Promter. These prompters guide fixed LLMs to track dialogue states, retrieve relevant information, and generate appropri-

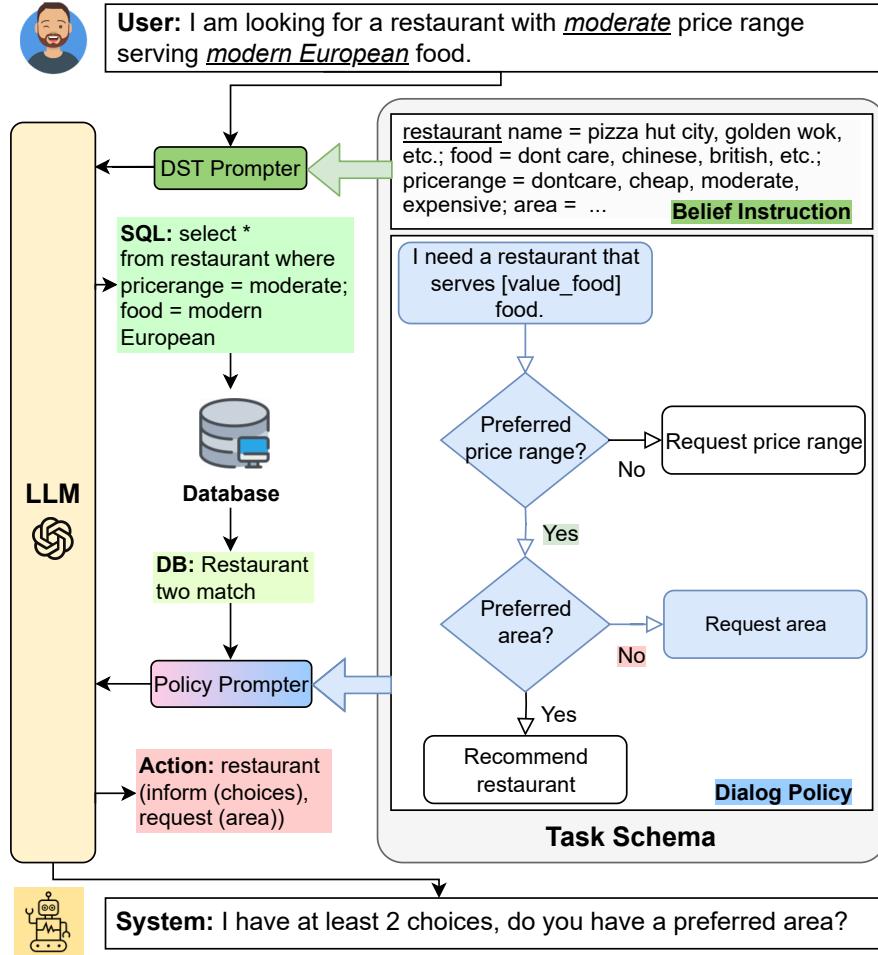


Figure 4.1: The proposed SGP-TOD is depicted with a dialog example, where the prompters integrate the task schema (right) to assist the frozen LLM in generating an appropriate response (left).

ate responses for novel tasks *in a zero-shot manner, without the need for additional training or fine-tuning.*

By incorporating task-specific symbolic knowledge into LLMs, SGP-TOD provides knowledge-based, coherent and human-like responses. Moreover, this training-free design empowers developers to flexibly prototype dialog systems on new tasks, while seamlessly extending system functionalities through modifying the task schema.

4.2 Schema-Guided Prompting (SGP)-TOD

4.2.1 Overview

The overall architecture of the proposed SGP-TOD (Figure 4.1) consists of three key components: (i) an **LLM**, responsible for adhering to instructions, comprehending user queries, and generating coherent responses for user interaction; (ii) a **DST Promoter**, tasked with supporting the LLM in tracking dialogue states using the belief instruction; (iii) a **Policy Promoter**, guiding the LLM to adhere to the predefined task policy for providing suitable system actions and responses.

At each dialog turn t , the end-to-end generation task is systematically divided into three subsequent sub-tasks:

Sub-task 1: Belief state prediction. Given the dialog history up to current dialog turn h_t , which is a sequence of utterances alternating between the user and the system $h_t = [u_1, r_1, u_2, r_2, \dots, u_t]$ (where u and r denote user and system utterances, respectively), the DST Promoter embeds the belief instruction BI to direct the frozen LLM (parameterized by θ) in generating a belief state b_t (Equation 4.1). The belief state is then used to query a database and obtain the database (DB) state c_t (Equation 4.2).

Sub-task 2: System action determination. The Policy Promoter incorporates a policy skeleton PS , assisting the LLM in generating a system action a_t , based on h_t , b_t , and c_t (Equation 4.3).

Sub-task 3: Dialog response generation. Grounded in the dialog history h_t , belief state b_t , DB state c_t , system action a_t , the Policy Promoter aids the LLM in generating a delexicalized response by providing the policy skeleton PS (Equation 4.4). Ultimately, the delexicalized response is automatically post-processed to generate system response

in natural language.

$$b_t = \text{LLM}_\theta(h_t, BI) \quad (4.1)$$

$$c_t = DB(b_t) \quad (4.2)$$

$$a_t = \text{LLM}_\theta(h_t, b_t, c_t, PS) \quad (4.3)$$

$$r_t = \text{LLM}_\theta(h_t, b_t, c_t, a_t, PS) \quad (4.4)$$

4.2.2 LLM

An LLM is responsible for following task-specific instructions and generating appropriate responses. Many off-the-shelf LLMs, e.g., ChatGPT, Codex [16], are pre-trained on massive corpora of text data and/or code data. In addition, they are trained to follow instructions in the prompts [109] and provide pertinent responses. Exhibiting remarkable proficiencies in natural language processing, instruction compliance, and zero-shot generalization across diverse downstream dialog tasks, these LLMs serve as valuable foundation models for our approach.

4.2.3 DST Promter

Given the dialog history h_t , the DST prompter aims to guide the LLM in predicting the belief state b_t at each turn t , using the belief instruction BI . The belief state b_t is defined as the concatenation of the domain/task (i.e., user intent) d_t and a set of slot-value pairs $\{(s_t^i, v_t^i); i = 1, \dots, n_t\}$, where n_t is the total number of pairs in the set.

As shown in Figure 4.2, the proposed DST prompter contains four parts: (i) a *task*

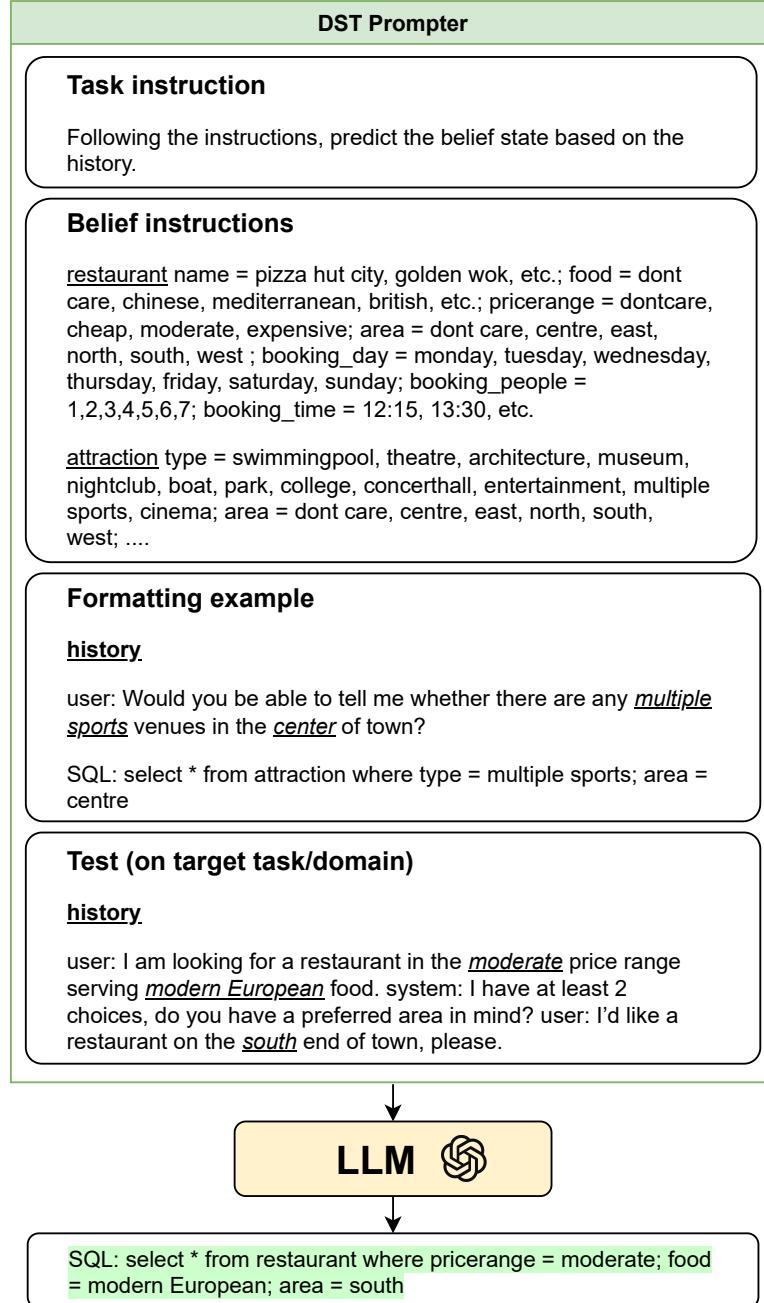


Figure 4.2: Illustration of belief state prediction utilizing DST Promter. The predicted belief state is highlighted.

instruction that offers general guidance on belief state prediction;¹ (ii) *belief instruc-*

¹ We assess several task instructions written by different authors, yielding minor performance disparities.

tions BI of all domains/tasks; (iii) a *formatting example* illustrating the anticipated output format to direct the LLM, in addition, we follow Hu et al. [57] and adopt SQL state to represent the dialog state b_t ;¹ and (iv) the *test input*, i.e., the given dialog history h_t . Since the prompt is fixed and no labeled data from the target task or domain is used, we refer to this setting as “zero-shot”, following Wang et al. [184].

Belief Instruction. For each task/domain, the belief instruction contains the task/domain name, all potential slot names, and their possible values (Figure 4.2). Regarding categorical slots, such as the “price range” in the restaurant domain, all plausible values are included, i.e., “don’t care”, “cheap”, “moderate”, and “expensive”; whereas, for non-categorical slots, such as “name”, only a few value examples are injected, e.g., Pizza Hut City, Golden Wok, etc.² Detailed belief instructions for all tasks/domains can be found in Figure 4.5.

4.2.4 Policy Prompter

Dialog policy, governing the behavior of task bots, plays a crucial role in task-oriented dialogs. To represent the dialog policy for a given task, we utilize a *policy skeleton*, which delineates interaction patterns and encompasses business logic in the form of template dialog flows [120]. The Policy Prompter is devised to guide the static LLM in adhering to the policy skeleton PS , enabling the sequential generation of appropriate system actions a_t and responses r_t .

Analogous to the DST Promter, the Policy Promter (Figure 4.3) comprises four components: (i) a *task instruction*; (ii) a *formatting example* derived from another

¹SQL: select * from d_t where $s_t^1 = v_t^1; \dots; s_t^{n_t} = v_t^{n_t}$.

²We assess belief instructions with diverse slot value examples, revealing minor performance variations.

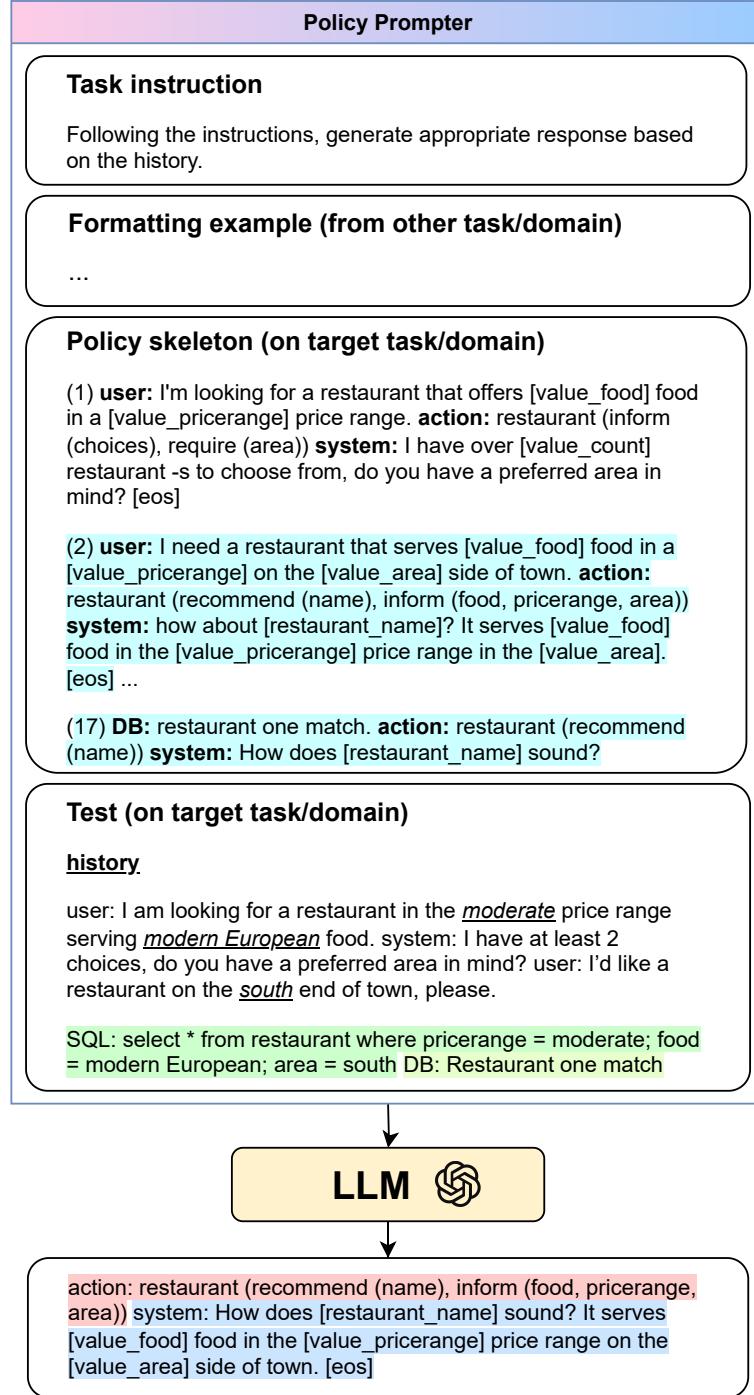


Figure 4.3: Illustration of system action determination and response generation employing the Policy Prompter. The pertinent template turns, previously predicted belief state, retrieved DB state within the input, alongside the generated system action and generated response in the output are accentuated.

task/domain, consisting of a partial policy skeleton and its associated dialogue turn exemplar (in Figure 4.6); (iii) a *policy skeleton* for the previously predicted domain/task; and (iv) the *test input*, i.e., the dialog history h_t , generated belief state b_t , and obtained DB state c_t .

Policy Skeleton. Given that user behaviors and DB results jointly determine system actions and responses, policy skeleton is designed to cover all fundamental user behaviors and characteristic DB results, along with their corresponding system actions and responses.¹ Considering the infeasibility of developing a multi-task/domain policy skeleton for every possible combination of tasks and domains, we opt to develop a distinct policy skeleton tailored to each specific task and domain.

Following Mehri and Eskenazi [100], our strategy converts the established dialog policy into a series of template dialog turns X that are logically arranged and concentrate on task completion:

$$X = \{x_i\}_{i=1}^N, \quad (4.5)$$

$$x_i = (u^i, a^i, r^i) \text{ or } (c^i, a^i, r^i)$$

where x_i is a template dialog turn, which contains a user utterance u^i or a DB state c^i , matching system action a^i , and system response r^i . N denotes the total number of template turns within the policy skeleton (around 10-20 template turns depending on the task complexity). In order to equip the frozen LLM with new capabilities or modify current ones, we only need insert, amend, or eliminate a few template turns within the policy skeleton.

¹ We do not enumerate every conceivable combination of user behaviors or potential database results, as schema engineering is not the primary focus of this study.

4.3 Experiments

To validate the efficiency of the proposed SGP-TOD, we conduct the experiments on two end-to-end dialog tasks: (*i*) dialog response generation and (*ii*) next action prediction, in three different settings: (*i*) multi-domain, (*ii*) single-domain, and (*iii*) domain-extension setting.

4.3.1 Setup

Datasets. We conduct evaluations using the following datasets:

- Multiwoz 2.0 [11] is a **multi-domain** task-oriented dataset encompassing seven domains: restaurant, attraction, train, hotel, taxi, police, and hospital. It comprises 8,438 dialogs for training, 1,000 for validation, and 1,000 for testing, all annotated with belief states and system actions.
- Multiwoz 2.2 [204] is a improved version of Multiwoz 2.0, encompassing refined belief state annotations, slot descriptions, user action annotations, etc.
- RADDLE [118, 121] consists of four **single-domain** dialog datasets derived from Multiwoz 2.0 (i.e., restaurant, train, hotel, attraction), reorganized by Peng et al. [118]. Each corpus contains 50/50/200 dialogs for training/validating/testing, expect for 100 testing dialogs in attraction domain.
- STAR [104] includes 24 tasks in 13 domains (e.g., “apartment” domain comprises “apartment-search” and “apartment-schedule”), requiring the dialog model to conform to the provided task schema. We utilize 2,688 single-task dialogs from the corpus, specifically focusing on “happy path” scenarios where user

actions remain within the schema’s expectations. While lacking explicit annotations, STAR provides flow chart diagrams outlining the expected dialog policy for each task. These diagrams specify attribute request sequences (e.g., obtaining the user’s name before the hotel name), database query procedures, and other relevant instructions.

Automatic Evaluation Metrics. We evaluate the end-to-end dialog generation performance using the same metrics as those listed in Budzianowski et al. [11]: *(i)* Inform(%) assesses whether the agent returns an acceptable entity. *(ii)* Success(%) determines if the agent appropriately responds to each attribute request. *(iii)* BLEU(%) [114] measures the word overlap of the generated response against the human response in the corpus. *(iv)* Combined(%) judges the overall quality, which is defined as $\text{Combined} = (\text{Inform} + \text{Success}) \times 0.5 + \text{BLEU}$. Additionally, we utilize BERTScore(%) [211], which focuses on computing semantic similarity between the generated responses and the ground truth, and correlates better with human judgments.

Following Mehri and Eskenazi [100], we perform the next action prediction task on STAR, which predicts next system action based on the dialog history. Since the system actions and deterministic response templates are mapped one to one in STAR corpus, we believe the end-to-end next action prediction task falls within end-to-end dialog modeling, following Mehri and Eskenazi [100], Mosig et al. [104]. In addition, we report the results using weighted F1score(%) and mean accuracy(%).

Human Evaluation Metrics. We employ interactive human evaluations to assess the quality of dialog agents, following the evaluation protocol in the DSTC9 Track 1 challenge [46]. We recruit student helpers to help with evaluations. For each dialog session,

student helpers are provided with a goal and accompanying instructions, subsequently necessitating a discourse with the agent to achieve the goal via natural language. Upon the conclusion of each dialog session, students are mandated to assess the overall dialog quality employing these five metrics: (i) Success w/o g(%) evaluates whether the agent accomplishes the task. (ii) Success w/ g(%) judges whether the agent accomplishes the task and offers matched slot values compared to the database record. (iii) Understanding(1-5) quantifies the accuracy with which the agent comprehends user utterances. (iv) Appropriateness(1-5) signifies the naturalness, appropriateness and fluency of an agent response. (v) Turns denotes the average number of dialog turns within successful dialog sessions.

Compared Methods. We compare the proposed SGP-TOD with SOTA zero-shot transfer methods and zero-shot/few-shot prompting strategies. (We report the mean results of three different runs.)

Zero-shot transfer methods:

- **BERT+S** [104] is a schema-guided method that augments a BERT-base classifier [28] with a provided system-side schema to predict the next system action.
- **SAM** [100] represents a schema-guided model based on BERT-base, which aligns the dialog context to a user-aware schema to predict the next system action.
- **ANYTOD-XXL** [223] adopts a neural LM to track dialog states and user actions utilizing slot and action descriptions. Then a program that outlines a predefined task policy is executed to recommend appropriate system actions. Upon considering these system actions, an LM generates the ultimate system action and

formulates the corresponding template response using the approach proposed by Kale and Rastogi [66]. ANYTOD-XXL is implemented on T5-XXL [136] and pre-trained on SGD dataset [134].¹

Prompting methods:

- **IG-TOD-CHATGPT** [61] is a prompting approach based on ChatGPT that leverages the dialog context and manually-crafted slot descriptions as the prompt, to track dialog states, fetch DB entries, and produce responses. IG-TOD-CHATGPT-ZS and IG-TOD-CHATGPT-FS are in the zero-shot and few-shot settings, respectively.
- **FEW-SHOT-CHATGPT** is a few-shot prompting strategy implemented on ChatGPT, where we use a few (i.e., k) dialog turns, randomly sampled from the training corpus to instruct ChatGPT on task execution. Upon evaluating various configurations of k , the optimal results manifest with 15 on Multiwoz (2.0 and 2.2), and 10 on RADDLE, exhibiting no further substantial enhancements.
- **SGP-TOD (Ours)** is a schema-guided prompting strategy, which is compatible with any off-the-shelf LLMs. Following the zero-shot scenario in Wang et al. [184], we insert one formatting example from different tasks (fixed through the experimental procedure) into the prompt.

Implementation Details. LLMs: We employ ChatGPT (“gpt-3.5-turbo”), GPT-3.5 (“text-davinci-003”) and Codex (“code-davinci-002”) as the fixed LLMs to implement the proposed SGP-TOD. Throughout the evaluation, we set temperature to 0.5.

¹The Schema-Guided Dialog (SGD) dataset constitutes a comprehensive, large-scale, multi-domain corpus encompassing over 16,000 dialogs that span across 16 distinct domains.

DST Prompter – belief instruction: In the context of multi-domain scenarios, the belief instructions encompassing all domains are incorporated, while solely the target domain’s belief instruction is introduced in single-domain settings.

Policy Prompter – policy skeleton: For the Multiwoz datasets, we manually construct the policy skeleton through observing a few dialogs in the training corpus, following Mehri and Eskenazi [100], Mosig et al. [104]. In the case of the STAR corpus, we employ flow chart diagrams and several dialogs to develop the policy skeleton, following the guidelines set forth by Mehri and Eskenazi [100]. We integrate the relevant user template utterance and the system action into the policy skeleton, thereby augmenting the LLM’s understanding of directives, in the absence of belief annotations. The prompt examples for the STAR dataset are shown in Figure 4.7 and Figure 4.8.

Formatting example: Following the zero-shot scenario in Wang et al. [184], we insert one formatting example from different tasks (fixed through the experimental procedure) into the prompt. The formatting example employed within DST Prompter/Policy Prompter is randomly chosen from the training corpus of different tasks/domains, conforming to zero-shot scenario proposed by Wang et al. [184]. We appraise multiple randomly selected formatting examples, the evaluation results reveal minor deviations. In the experiments on domain extension (Section 4.3.5) and ablation analysis (Section 4.4.1), we employ the same (two) formatting exemplar turns originating from other domains within the RADDLE corpus for all prompting techniques.

4.3.2 Evaluation on Multiwoz

We present the evaluation results in multi-domain contexts on Multiwoz in Table 4.1. In addition to the aforementioned methods, we include the results of SOTA full-

Model	Multiwoz 2.0				Multiwoz 2.2			
	Info.	Succ.	BLEU	Combined	Info.	Succ.	BLEU	Combined
<i>Full-shot fine-tuning (with 8.4k+ training dialogs):</i>								
DAMD [218]	76.33	60.40	16.60	84.97	-	-	-	-
SIMPLETOD [56]	84.40	70.10	15.01	92.26	-	-	-	-
SOLOIST [118]	85.50	72.90	16.54	95.74	81.70	67.10	13.60	88.00
PPTOD [157]	89.20	79.40	18.62	102.92	-	-	-	-
MARS [158]	88.90	78.00	19.90	103.35	88.90	78.00	19.60	103.05
<i>Zero-shot transfer method (pre-trained on SGD):</i>								
ANYTOD-XXL	-	-	-	-	73.90	24.40	3.40	52.55
<i>Few-shot prompting:</i>								
IG-TOD-CHATGPT-FS	-	-	-	-	-	20.00	7.17	-
FEW-SHOT-CHATGPT	44.74	24.32	7.88	42.41	45.40	24.50	7.72	42.67
<i>Zero-shot prompting:</i>								
IG-TOD-CHATGPT-ZS	-	-	-	-	-	15.00	3.58	-
SGP-TOD-CHATGPT	64.56	54.05	7.17	66.48	64.70	54.70	6.96	66.66
SGP-TOD-CODEX	71.67	52.55	7.91	70.02	75.50	52.30	6.62	70.53
SGP-TOD-GPT3.5	83.88	69.87	9.09	85.97	82.00	72.50	9.22	86.47

Table 4.1: End-to-end dialog generation evaluation results on Multiwoz. Results of SOLOIST, MARS, ANYTOD-XXL on Multiwoz 2.2 are cited from Zhao et al. [223]. Results of IG-TOD-CHATGPT are cited from Hudecek and Dusek [61]. Other results of the full-shot fine-tuning methods are cited from He et al. [51] and Sun et al. [158]. Info.: Inform, Succ.: Success.

shot fine-tuning approaches to facilitate a more comprehensive comparison. We have the following key observation:

Integrating task schema into LLMs enables effective zero-shot generalization on new tasks. SGP-TOD obtains SOTA *zero-shot performance*, substantially outperforming few-shot prompting approaches across all metrics, while even exhibiting competitive results in comparison to full-shot fine-tuning methods concerning Success and Inform.

Comparison with Prompting Methods. Explicit task instructions through a schema are more effective for task completion than implicit dialog-based guidance. SGP-TOD-CHATGPT distinctly surpasses the zero-shot prompting approach IG-TOD-CHATGPT-ZS with respect to Success (surpassing by 40%) and BLEU (exceeding by

3%). Moreover, SGP-TOD-CHATGPT, *without requiring task-specific data*, considerably outperforms the few-shot prompting methods, i.e., IG-TOD-CHATGPT-FS and FEW-SHOT-CHATGPT (e.g., about 30 points improvement over Success).

Comparison with Zero-Shot Transfer Methods. Our SGP-TOD demonstrates a substantial advantage over ANYTOD-XXL, which necessitates task-specific pre-training and additional annotations, e.g., slot and action descriptions, over all the metrics. This exemplifies the potency of SGP-TOD, which markedly reduces the necessity for human labor and computational resources.

Comparison with Full-Shot Fine-Tuning Methods. SGP-TOD exhibits competitive performance over Inform and Success. The lower BLEU is due to a lack of linguistic variations of the template utterances, which is acceptable considering the trade-off between human effort and efficacy.

4.3.3 Evaluation on RADDLE

Table 4.2 reports the results in single-domain settings on RADDLE. On all four dialog tasks, SGP-TOD demonstrates remarkable zero-shot performance that consistently surpasses both few-shot prompting and fine-tuning approaches. This results in substantial improvements of up to 12% in Inform, 45% in Success, and 19% in Combined metrics, while maintaining competitive BLEU scores. This evidence further substantiates the efficacy of SGP-TOD.

Model	Attraction						Train					
	Info.	Succ.	BLEU	Combined	Info.	Succ.	BLEU	Combined	Info.	Succ.	BLEU	Combined
<i>Few-shot fine-tuning (with 50 training dialogs):</i>												
SIMPLETOD	65.66	46.97	5.85	62.17	59.00	44.00	7.07	58.57				
SOLOIST	86.00	65.00	12.90	88.40	80.81	64.65	9.96	82.69				
<i>Few-shot prompting:</i>												
FEW-SHOT-CHATGPT	75.00	67.00	8.22	79.23	79.80	65.66	8.12	80.85				
<i>Zero-shot prompting:</i>												
SGP-TOD-CHATGPT	95.00	94.00	7.13	101.63	76.77	74.24	6.75	82.26				
SGP-TOD-CODEX	98.00	93.00	10.45	105.95	78.79	70.20	8.56	83.06				
SGP-TOD-GPT3.5	96.00	93.00	9.53	104.03	82.83	77.27	8.72	88.77				
Model	Hotel						Restaurant					
	Info.	Succ.	BLEU	Combined	Info.	Succ.	BLEU	Combined	Info.	Succ.	BLEU	Combined
<i>Few-shot fine-tuning (with 50 training dialogs):</i>												
SIMPLETOD	62.50	40.00	7.70	58.95	75.50	44.50	11.00	71.00				
SOLOIST	74.50	43.50	8.12	67.12	81.00	55.50	12.80	81.50				
<i>Few-shot prompting:</i>												
FEW-SHOT-CHATGPT	51.00	26.50	5.80	44.55	80.00	55.50	7.71	75.46				
<i>Zero-shot prompting:</i>												
SGP-TOD-CHATGPT	76.50	57.00	5.16	71.91	90.00	82.50	6.72	92.97				
SGP-TOD-CODEX	83.50	69.50	7.86	84.36	91.00	85.00	10.50	98.50				
SGP-TOD-GPT3.5	82.50	71.50	7.05	84.05	91.50	84.00	12.90	100.65				

Table 4.2: End-to-end dialog generation evaluation results on RADDLE. The few-shot fine-tuning results are cited from Peng et al. [118].

Model	Task transfer		Domain transfer	
	F1	Accuracy	F1	Accuracy
<i>Zero-shot transfer</i>				
<i>(leave-one fine-tuning with 2.5k training dialogs):</i>				
BERT+S	24.25	24.89	25.70	28.56
SAM	49.82	51.30	55.91	57.92
<i>Zero-shot prompting:</i>				
SGP-TOD-CODEX-INI	45.18	47.99	47.21	49.97
SGP-TOD-GPT3.5	47.67	48.27	49.76	50.39
SGP-TOD-CODEX	49.78	51.01	52.72	53.66
SGP-TOD-GPT3.5-E2E	50.84	50.74	53.50	53.21

Table 4.3: Zero-shot end-to-end next action prediction evaluation results on STAR. (Difference in mean is significant with $p < 0.01$.)

4.3.4 Evaluation on STAR

BERT+S, SAM are fine-tuned on source tasks/domains then zero-shot on the held-out task/domain.¹ SGP-TOD is presented with two formatting turns from the source tasks/domains. Following Mehri and Eskenazi [100], we report the zero-shot evaluation results in two settings, i.e., task transfer and domain transfer in Table 4.3. SGP-TOD, *merely with two formatting sample turns*, demonstrates exceptional performance, surpassing or rivaling SOTA zero-shot transfer methods in both settings. This outcome signifies that, even when faced with complicated business logic and system actions in dialog policies, the proposed SGP-TOD continues to exhibit commendable performance.

Impact of Varying the Number of Training Dialogs and Formatting Example Turns.

Figure 4.4 exhibits the zero-shot evaluation results on STAR, utilizing varying amounts of training dialogs (ranging from 1 to 1,000) and formatting example turns (spanning from 1 to 10) from source domains/tasks. SGP-TOD, *merely with two formatting sample turns*, achieves superior or comparable performance compared to BERT+S, SAM, which are fine-tuned on adequate source data.

We observe that SGP-TOD, *employing only two formatting sample turns*, attains superior or commensurate performance in terms of both F1 score and Accuracy, when compared to SAM trained with 1,000 dialogs. Given that a single dialog contains more than 10 dialog turns, this result suggests that SGP-TOD diminishes labeling expenses by a minimum factor of 1,000. Furthermore, it is noteworthy that augmenting the quantity of formatting exemplar turns exerts a negligible influence on the performance

¹ ANYTOD-XXL requires additional annotations, e.g., belief descriptions, which makes it not suitable for STAR.

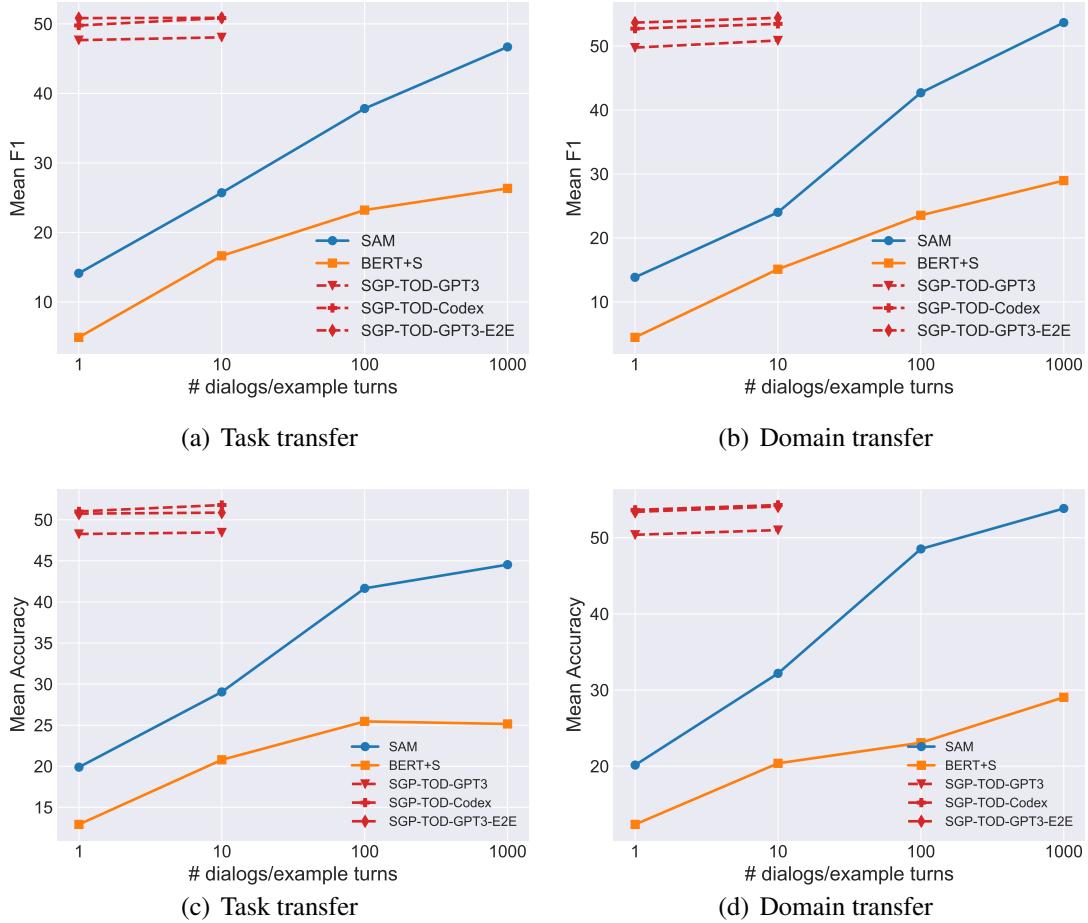


Figure 4.4: Zero-shot end-to-end evaluation results on STAR with different numbers of training dialogs (1, 10, 100, 1,000) / demonstration example turns (1, 10) from source domain/tasks.

of SGP-TOD.

Impact of Different LLMs and Prompting Formats. SGP-TOD-CODEX surpasses SGP-TOD-GPT3.5 while rivaling SGP-TOD-GPT3.5-E2E (with template responses affixed to action labels in the policy prompt, demonstrated in Figure 4.8). We conjecture that Codex, benefiting from extensive pre-training on copious code data, demonstrates enhanced proficiency compared to GPT-3.5 in interpreting action labels. In addition, appending template responses is presumed to facilitate the explication of action

Model	FT/FS/ZS	Restaurant-Ext			
		Inform	Success	BLEU	BERTScore
w/o domain-relevant knowledge					
ChatGPT	ZS	44.00	6.00	4.31	85.96
GPT-3.5	ZS	34.00	16.00	8.70	84.31
w/ prior knowledge on Restaurant					
SOLOIST	FT	78.00	0.00	10.62	87.24
SGP-TOD-CHATGPT	ZS	88.00	34.00	5.45	86.11
SGP-TOD-GPT3.5	ZS	94.00	30.00	10.68	87.30
w/ knowledge on Restaurant-Ext					
SOLOIST+TEACH	FT	82.00	38.00	10.99	87.66
FEW-SHOT-GPT3.5+TEACH	FS	88.00	54.00	12.95	88.90
SGP-TOD-CHATGPT-EXT	ZS	88.00	78.00	6.25	86.15
SGP-TOD-GPT3.5-EXT	ZS	96.00	86.00	14.57	89.01

Table 4.4: End-to-end evaluation results on domain extension. FT: fine-tuning, FS: few-shot prompting, ZS: zero-shot prompting.

labels for GPT-3.5.

Impact of Different Task Schemas. SGP-TOD-CODEX-INI, utilizing an identical task schema as employed in training SAM, manifests commendable performance. This result highlights that SGP-TOD as a flexible prompting strategy, compatible with any manually-crafted task schema.

4.3.5 Evaluation on Domain Extension

We conduct experiments in a domain extension setting [42, 86] to assess the efficacy of SGP-TOD in adapting deployed task bots to incorporate novel functionalities. Following Zhang et al. [212], we construct the Restaurant-ext corpus by extending the Restaurant in RADDLE with four new slots: *[restaurant_dish]*, *[value_price]*, *[start_time]*, and *[end_time]*.

Compared Methods.

- **ChatGPT, GPT-3.5** denote zero-shot prompting that receive two formatting examples.
- **SGP-TOD-CHATGPT, SGP-TOD-GPT3.5** represent our SGP-TOD implementation, with the Restaurant policy skeleton.
- **SOLOIST** is trained with 50 training dialogs in Restaurant domain (reported in Table 4.2).
- **SOLOIST+TEACH** is fine-tuning method enhanced with machine teaching (as introduced in Section 3.3) [150]. We deploy SOLOIST to converse with real users, then implement machine teaching to obtain 10/50/50 annotated dialogs in Restaurant-ext for training, validating, and testing. We fine-tune SOLOIST with the gathered 10 training dialogs covering new slots.
- **FEW-SHOT-GPT3.5+TEACH** is the few-shot prompting strategy augmented with machine teaching. We use 10 randomly selected dialog turns from the collected 10 training dialogs as prompts (with peak performance at 10).
- **SGP-TOD-CHATGPT-EXT, SGP-TOD-GP3.5-EXT** refer to SGP-TOD with Restaurant-Ext policy skeleton, where we only add four template turns about four new slots to the policy skeleton of Restaurant.

Merely modifying the task schema enables adaptively expand the LLM’s functionalities. In Table 4.4, SGP-TOD-CHATGPT-EXT, and notably SGP-TOD-GPT3.5-EXT surpasses all other evaluated approaches by a substantial margin over all the metrics.

Comparison with Approaches Augmented by Machine Teaching. SOLOIST yields zero Success, a predictable result given its lack of awareness regarding the new features. Augmented by machine teaching, SOLOIST+TEACH substantially improves

Model	Multiwoz 2.0				Multiwoz 2.2					
	Inform		Success		BLEU	Combined	Inform	Success	BLEU	Combined
	SP-TOD-GPT3.5	83.88	69.87	9.09	85.97	82.00	72.50	9.22	86.47	
-policy	82.28	55.65	6.51	75.48	81.80	56.20	6.63	75.63		
-policy -DB	81.20	50.95	6.48	72.56	81.40	52.30	6.57	73.42		
-policy -DB -belief	38.74	33.13	6.18	42.12	38.60	33.90	6.29	42.54		

Table 4.5: Ablation study on the impact of the three components in the proposed SGP-TOD and the database expertise on Multiwoz using GPT-3.5. -policy: removing Policy Prompter, -DB: removing database information, -belief: removing DST Prompter.

SOLOIST in terms of Inform and Success. Nevertheless, relying solely on prior Restaurant knowledge, both SGP-TOD-CHATGPT and SGP-TOD-GP3.5 exhibit performance on par with SOLOIST+TEACH, demonstrating that SGP-TOD provides enhanced robustness in zero-shot generalization. Moreover, SGP-TOD-GP3.5-EXT obtains substantially higher Success rates than SOLOIST+TEACH (a rise of 48%) and FEW-SHOT-GPT3.5+TEACH (an increase of 32%). Compared to fine-tuning/prompting strategies utilizing additional dialogs corrected through machine teaching, SGP-TOD facilitates a more agile adaptation to novel functionalities by merely modifying template turns within the task schema.

4.4 In-Depth Analyses

In this section, we first present an ablation study to evaluate the efficacy of SGP-TOD. We then provide qualitative analyses and human evaluation results on the generated dialogues to further confirm the effectiveness of SGP-TOD in real-world scenarios.

4.4.1 Ablation Study

Table 4.5 exhibits the findings from an ablation investigation, addressing the effects of the three integral aspects of SGP-TOD in conjunction with the database expertise, implemented on Multiwoz 2.0 and 2.2, employing GPT-3.5.¹ Combining the three elements in SGP-TOD with the database expertise produces optimal results across both datasets. The removal of the Policy Prompter, database knowledge, and DST Prompter leads to consistent declines in all evaluation metrics, underscoring the value of enhancing the fixed LLM with the task schema and external database information.

Specifically, *GPT-3.5* (in the final row) exhibits commendable zero-shot performance, highlighting the need of exploiting its superior zero-shot generalization capabilities in dialog generation tasks. Additionally, *disabling the Policy Prompter* incurs a discernible decline in performance regarding Success (approximately 16%) and BLEU (roughly 3%), as the Policy Prompter’s primary function is to provide task completion guidelines and interaction patterns. *Eliminating the database expertise* primarily reduces Success (by approximately 4%), implying that incorporating database information contributes to task completion. Lastly, *excising the DST Prompter* engenders a considerable diminution in performance concerning Inform (around 43%) and Success (nearly 18%), due to the DST Prompter’s intended purpose of assisting the frozen LLM in apprehending the dialog context.

Turn	w/ ChatGPT	Turn	w/ SGP-TOD-CHATGPT-EXT
1	User: Hi, i'm looking for a place with <i>Tuscan</i> food in <i>any</i> price range.	1	User: Hi, i'm looking for a place with <i>Tuscan</i> food in <i>any</i> price range.
2	System: Certainly, there are <i>3 Tuscan restaurants</i> in town. What is your <i>preferred price range</i> ?	2	System: I am sorry, there are <i>no Tuscan restaurants</i> . Here are <i>some other options: Korean, Italian, Chinese</i> . Which one would you like more information on?
3	User: I <i>don't care</i> about the price. Please recommend one for me, i also need the phone number please.	3	User: How about any <i>Korean</i> restaurants? I also need the phone number please.
4	System: The phone number for <i>Pizza Hut City Centre</i> is <i>01223323737</i> . Is there anything else i can help you with today?	4	System: <i>Little Seoul</i> is a <i>Korean</i> restaurant. Their phone number is <i>01223308681</i> .
5	User: Does the restaurant offer delivery service? How much does the delivery charge?	5	User: Does the restaurant offer delivery service? How much does the delivery charge?
6	System: I'm sorry, but <i>I could not find any information on delivery fee</i> in our database. You may want to contact the restaurant directly to inquire about this. Is there anything else I can assist you with today?	6	System: Yes, they offer delivery service and the delivery charge is <i>4 pounds</i> . Can I help you with anything else?
7	User: No. Thank you, goodbye.	7	User: No. Thank you, goodbye.
8	System: Glad i could assist you. Have a good day.	8	System: You're welcome, have a great meal! Goodbye!

Table 4.6: Illustration of two interactive examples between a user and CHATGPT (left), SGP-TOD-CHATGPT-EXT (right). Appropriate responses and inappropriate responses are highlighted.

4.4.2 Qualitative Analysis

Despite the superior performance of the proposed SGP-TOD on GPT-3.5, we showcase interactive examples utilizing ChatGPT, a renowned and potent chatbot. In Table 4.6, a user engages with ChatGPT (left) and SGP-TOD-CHATGPT-EXT (right) to complete the identical task on Restaurant-Ext.¹ The user initiates the conversation by seeking recommendations for a Tuscan restaurant with no price range preference.

¹ We inject the same two formatting example turns into the prompt throughout the evaluation.

¹ ChatGPT and SGP-TOD-CHATGPT-EXT are previously reported in Table 4.4. The same two formatting example turns are incorporated into the prompt for both zero-shot strategies.

Model	Restaurant				
	S w/o g↑	S w/ g↑	Und. ↑	App. ↑	T. ↓
SOLOIST	34.00	30.00	2.18	2.10	10.64
FEW-SHOT-CHATGPT	94.00	74.00	4.58	4.72	8.32
SGP-TOD-CHATGPT	100.00	92.00	4.86	4.88	7.28

Table 4.7: Interactive human evaluation results. S w/o g: Success without grounding; S w/ g: Success with grounding; Und.: Understanding; App.: Appropriateness; T.: Turns.

Lacking external database information, ChatGPT conveys inaccurate details (Turn 2), whereas SGP-TOD-CHATGPT-EXT informs users of the absence of matching restaurants and proposes alternatives (Turn 2). This exemplifies the benefits of integrating real-world expertise into the fixed LLM. Furthermore, ChatGPT persistently inquires about the desired price range despite the user’s indifference. We argue that SGP-TOD assists the frozen LLM in discerning user intentions. In Turn 4, ChatGPT continues to furnish fabricated details (i.e., the restaurant name and phone number) concerning the nonexistent eatery, while SGP-TOD-CHATGPT-EXT identifies a suitable Korean restaurant and the corresponding factual information. In contrast with ChatGPT, SGP-TOD-CHATGPT-EXT adeptly addresses inquiries about the delivery service (Turn 6), indicating that SGP-TOD is capable of endowing the frozen LLM with novel functionalities.

4.4.3 Interactive Human Evaluation

We conduct interactive human evaluations on Restaurant domain to evaluate the performance of SOLOIST, FEW-SHOT-CHATGPT, SGP-TOD-CHATGPT (reported in Table 4.2), with 50 dialogs gathered for analysis, respectively. Specifically, we enlisted 5 student helpers (i.e., undergraduate students possessing basic proficiency in English communication) to participate in the evaluations. For each dialog agent, we

collected 50 dialogs for analysis. Followed the methodology proposed by Li et al. [82], we generated user goals through the subsequent techniques: (*i*) Randomly selecting slots and slot values within the Restaurant domain from RADDLE corpus to construct a user goal; (*ii*) Replacing the slot values of the user goals in randomly chosen dialogs from the Restaurant corpus with corresponding new values from randomly sampled database entries, thus forming a new user goal; (*iii*) Merging the user goals of several randomly selected dialogs from the Restaurant corpus to create a composite user goal. Lastly, we randomly chose 50 distinct user goals from these newly generated goals.

Table 4.7 shows the interactive human evaluation results. Our proposed SGP-TOD-CHATGPT attains a remarkably high performance in a zero-shot context, consistently outpacing SOLOIST and FEW-SHOT-CHATGPT across all metrics. Particularly, regarding Success w/ g, SGP-TOD-CHATGPT significantly surpasses FEW-SHOT-CHATGPT (by 18%) and SOLOIST (by 62%), illustrating its proficiency in accomplishing tasks within real-world scenarios. In contrast to the automated evaluation results shown in Table 4.2, FEW-SHOT-CHATGPT significantly outperforms SOLOIST over all metrics. This indicates that corpus-based evaluations might be biased, given that real user inputs tend to be more dynamic, complex, even with noise. Notably, SGP-TOD-CHATGPT consistently excels compared to the other methods in both evaluations, implying its robustness in handling diverse user inputs.

4.5 Chapter Summary

This chapter introduces SGP-TOD, a novel schema-guided prompting strategy that transforms the development of task-oriented dialogue systems. By harnessing the capabilities of large language models (LLMs) and structured task schemas, SGP-TOD

facilitates the rapid creation of end-to-end task bots without requiring extensive training data.

Consider the ability to construct a restaurant reservation agent or a flight search assistant simply by providing a structured schema of the task. This is the central promise of SGP-TOD—enabling fixed LLMs to interpret and respond to user requests in a zero-shot setting, thereby removing traditional barriers such as large-scale data collection and model fine-tuning.

Empirical evaluations across multiple benchmark datasets demonstrate that SGP-TOD achieves state-of-the-art performance in zero-shot scenarios. Both automatic metrics and human evaluations confirm its effectiveness, underscoring its potential to democratize the development of task-oriented bots.

Nonetheless, the path toward reliable and trustworthy AI assistants presents ongoing challenges. Despite their impressive capabilities, LLM-powered bots remain susceptible to generating factually incorrect or “hallucinated” responses. This limitation highlights the importance of developing robust mitigation strategies—an issue we explore in the next chapter, as we move toward building dependable and broadly deployable AI solutions.

DST Prompter
<p>Task instruction</p> <p>Following the instructions, predict the belief state based on the history.</p> <p>Belief instructions</p> <p><u>restaurant</u> name = pizza hut city, golden wok, etc.; food = dont care, chinese, mediterranean, british, etc.; pricerange = dontcare, cheap, moderate, expensive; area = dont care, centre, east, north, south, west ; booking_day = monday, tuesday, wednesday, thursday, friday, saturday, sunday; booking_people = 1,2,3,4,5,6,7; booking_time = 12:15, 13:30, etc.</p> <p><u>attraction</u> attraction type = swimmingpool, theatre, architecture, museum, nightclub, boat, park, college, concerthall, entertainment, multiple sports, cinema; name = the cherry hinton village centre, soul tree nightclub, etc.; area = dont care, centre, east, north, south, west</p> <p><u>hotel</u> name = huntingdon marriott hotel, a and b guest house, etc.; pricerange = dont care, cheap, moderate, expensive; area = dont care, centre, east, north, south, west; stars = dont care, 0,1,2,3,4,5; parking = dont care, yes, no; internet = dont care, yes, no; type = hotel, guest house ; booking_day = monday, tuesday, etc.; booking_people = 1,2,3,4,5,6,etc. ; booking_stay = 1,2,3,4,etc.</p> <p><u>train</u> leaveat = 10:45, 12:06, etc. ; destination = norwich, cambridge, etc.; day = monday, tuesday, wednesday, thursday, friday, saturday, sunday ; arriveby = 21:00, 09:45, etc.; departure = cambridge, stansted airport, etc. ; booking_people = 1,2,3,4,5,6, etc.</p> <p><u>taxis</u> leaveat = 08:45, 16:15, etc. ; destination = saint john's college, kettle's yard, galleria, etc. ; departure = huntingdon marriott taxi, cineworld cinema, bridge guest house, etc. ; arriveby = 17:15, 17:30, etc.</p> <p>...</p>

Figure 4.5: Detailed belief instructions in DST Prompter.

Policy Prompter

Task instruction

Following the instructions, generate appropriate response based on the history.

Formatting example (from other task/domain)

Policy skeleton

(1) **user:** I'm looking for information on [attraction_name]. **action:** attraction (inform (name, address, area, entrance fee)) **system:** [attraction_name] is in the [value_area] and their address is [attraction_address], the entrance fee is [value_count] pounds. [eos]

(2) **user:** What is the phone number? **action:** attraction (inform (phone), require (more)) **system:** Their phone number is [attraction_phone]. Is there anything else I can do for you? [eos]

(3) **user:** I need the area/location and the postcode. **action:** attraction (inform (area, postcode)) **system:** The [attraction_name] is in the [value_area], their postal code is [attraction_postcode]. [eos]

(4) **user:** That is all I need today. Thank you! **action:** attraction (general (bye)) **system:** You are welcome, thank you for calling the [value_place] town info [value_area]. Have a great day. [eos]

(5) **DB:** attraction zero match. **action:** attraction (inform (none)) **system:** There are no results matching your query. Can I try a different area or type? [eos]

(6) **DB:** attraction one match. **action:** attraction (inform (name)) **system:** We have [attraction_name]. [eos]

(7) **DB:** attraction five match. **action:** attraction (inform (choices), request (area)) **system:** There are [value_count], which area were you looking at? [eos]

history

user: Hello, I'm looking for the museum of classical archaeology, could you tell me more about it? system: Certainly, the [attraction_name] is located at [attraction_address] and has free entrance, phone [attraction_phone]. Do you also need the postcode? user: Yes, I'd like that. Thank you so much!

SQL: select * from attraction where name = museum of classical archaeology DB: Attraction one match

action: attraction (inform (postcode), require (more))
 system: The postcode for the [attraction_name] is [attraction_postcode]. Is there anything else I can help you with?

...

Figure 4.6: A formatting example in Policy Prompter.

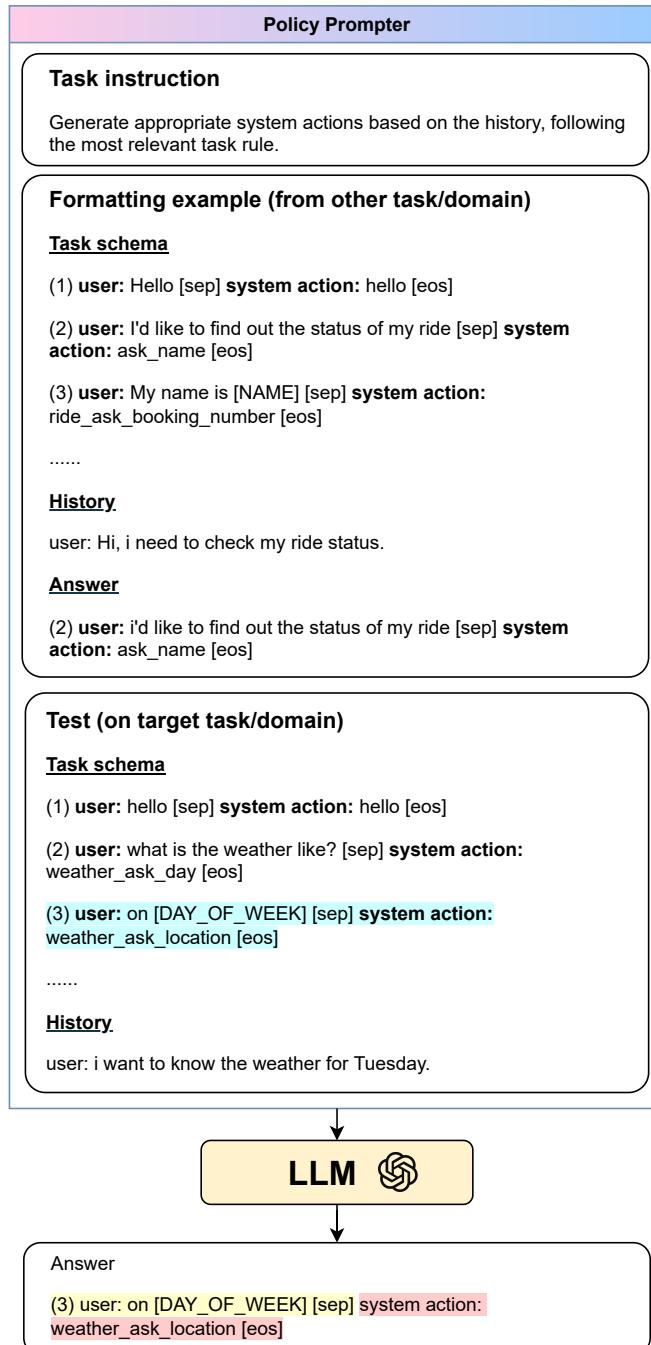


Figure 4.7: Policy Prompter of SGP-TOD on STAR. The relevant template turn within the input, the generated user template utterance, and the system action in the output are accentuated.

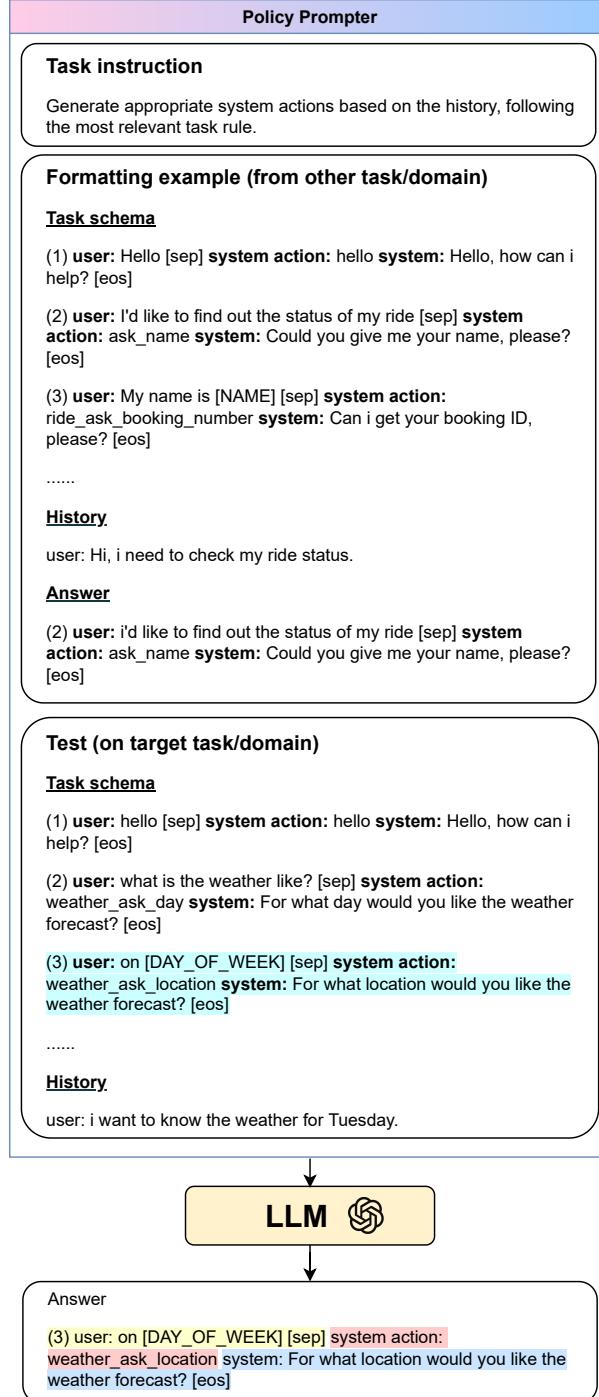


Figure 4.8: Policy Prompter of SGP-TOD-E2E on STAR. The relevant template turn in the input, the generated user template utterance, the system action and the system response in the output are highlighted.

Chapter 5

Self-Alignment for Factuality

This chapter addresses the critical challenge of ensuring task bots, powered by LLMs, reliably convey learned knowledge while maintaining factual accuracy. Recognizing that these bots possess vast knowledge bases, we explore a novel question: can they leverage this internal knowledge for self-evaluation, much like humans utilize self-critique for self-improvement? To answer this, we introduce *Self-Alignment for Factuality*, a novel framework that utilizes the self-evaluation capability of an LLM to generate training signals promoting factual accuracy.

We begin by outlining our motivation to enhance factuality by leveraging an LLM’s self-knowledge awareness (Section 5.1). Next, we introduce *Self-Alignment for Factuality*, a novel framework that incorporates a self-evaluation component (SELF-EVAL) to prompt LLMs to critically assess their responses (Section 5.2). This self-evaluation is further refined through *Self-Knowledge Tuning* (SK-TUNING), which enhances confidence estimation and calibration. We rigorously evaluate *Self-Alignment for Factuality* on three knowledge-intensive tasks (Section 5.3), analyze the effects of SK-TUNING (Section 5.4), and provide comprehensive analyses, including qualitative and error analyses (Section 5.5). Finally, we summarize our findings and emphasize the potential of *Self-Alignment for Factuality* for developing more reliable and accurate LLMs (Section 5.6).

5.1 The Importance of Maintaining Factuality

LLMs [106, 107, 171] have revolutionized NLP with their impressive capabilities across a wide range of tasks [14, 91, 185]. Their vast knowledge, acquired during pre-training, allows them to excel as generalist bots. However, as highlighted in Section 1.2, these models are susceptible to hallucinations (“tells”) [58, 63, 169, 220], even when they possess relevant knowledge (“knows”) [29, 77, 78, 98, 181]. This gap between “knowing” and “telling” [15, 65, 137] significantly limits their reliability and ability to accurately convey their acquired knowledge.

A few studies [21, 78, 219] edit the model’s internal representations towards “factuality” directions, using domain-specific annotated data. Meanwhile, acknowledging the inadequacy of the training objective—maximum likelihood estimation (MLE)—in accurately capturing factuality [1, 5, 111, 167] (discussed in Section 2.4), a recent study [167] introduces the LLM’s internal factuality signals as training rewards to guide the models towards factuality. Given that the origin of a LLM’s hallucinations is intrinsically linked to its confidence [58], Tian et al. [167] employs consistency-based confidence regarding the factual correctness over the generate responses [69, 98] as the factuality signals. Nevertheless, such consistency-based confidence remains rely on the model’s generation ability, which might be non-reflective on model’s internal knowledge.

Despite the challenges faced by an LLM in directly “telling” the correct response, it has showed potential in “evaluating” its generated responses [65, 137]. As depicted in Figure 5.1, the LLM is capable of identifying factual inaccuracies within the responses it generates, with a reasonable prediction confidence. Such self-evaluation, i.e., directly prompting the model itself about internal knowledge awareness, might be

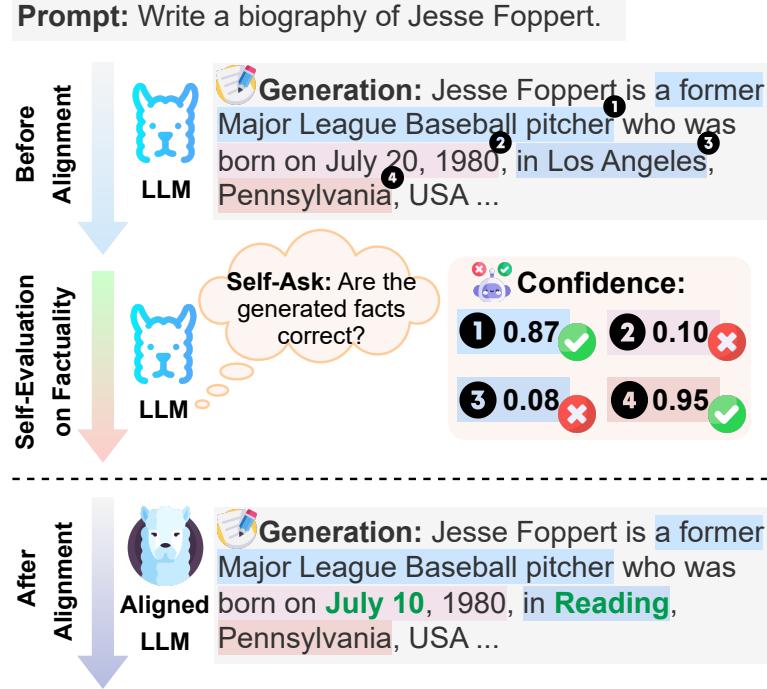


Figure 5.1: Illustration of *Self-Alignment for Factuality*. Given a prompt to write a biography, before factuality alignment, the LLM generates some facts that are not accurate. Through self-evaluation, the LLM is capable of identifying these inaccurate facts. The feedback from the self-evaluation is used as a reward signal to align the LLM towards factuality. Each fact is highlighted in distinct colors, and the corrected facts are marked with green letters.

a more effective approach to factuality estimation.

In this chapter, we introduce a self-alignment framework, *Self-Alignment for Factuality*, which harnesses an LLM’s self-evaluation capability to mitigate hallucinations. Our approach encourages an LLM to generate prediction confidence scores pertaining to the factuality of its own generated responses through self-asking. Subsequently, these scores are utilized as reward signals to fine-tune the model using the Direct Preference Optimization (DPO) algorithm [129]. Specifically, we incorporate a factuality self-evaluation component, SELF-EVAL, which prompts the LLM to directly validate its responses based on its internal knowledge. To bolster the LLM’s univer-

sal self-evaluation ability, we introduce SK-TUNING to enhance the LLM’s internal knowledge awareness, i.e., prediction confidence estimation and calibration [47, 168],¹ through sufficient tuning across heterogeneous knowledge-oriented tasks.

5.2 Self-Alignment for Factuality

In this section, we introduce the proposed framework. First, we provide a comprehensive overview of *Self-Alignment for Factuality* in Section 5.2.1. Subsequently, we delve into the *Factuality Self-Evaluation* by utilizing the LLM’s inherent knowledge, termed SELF-EVAL, in Section 5.2.2. Finally, we outline the factuality alignment process via DPO in Section 5.2.3.

5.2.1 Overview

Self-Alignment for Factuality generally operates in the following three steps, as depicted in Figure 5.2:

Step 1: Generating initial responses for preference data collection. For a given prompt x , we generate multiple candidate responses $\{y_m\}_{m=1}^M$, where M represents the sample size. These are produced from a base LLM guided by a policy $\pi_{\text{ref}}(y | x)$. To ensure the generation of coherent and relevant responses, we employ few-shot examples as prompts.

Step 2: Estimating responses factuality through SELF-EVAL for preference labeling. In this step, we evaluate the factuality of generated candidate responses $\{y_m\}_{m=1}^M$ for a given prompt x by leveraging the intrinsic knowledge of LLMs. In long-form

¹The confidence in a prediction is expected to accurately reflect the probability that the prediction is correct.

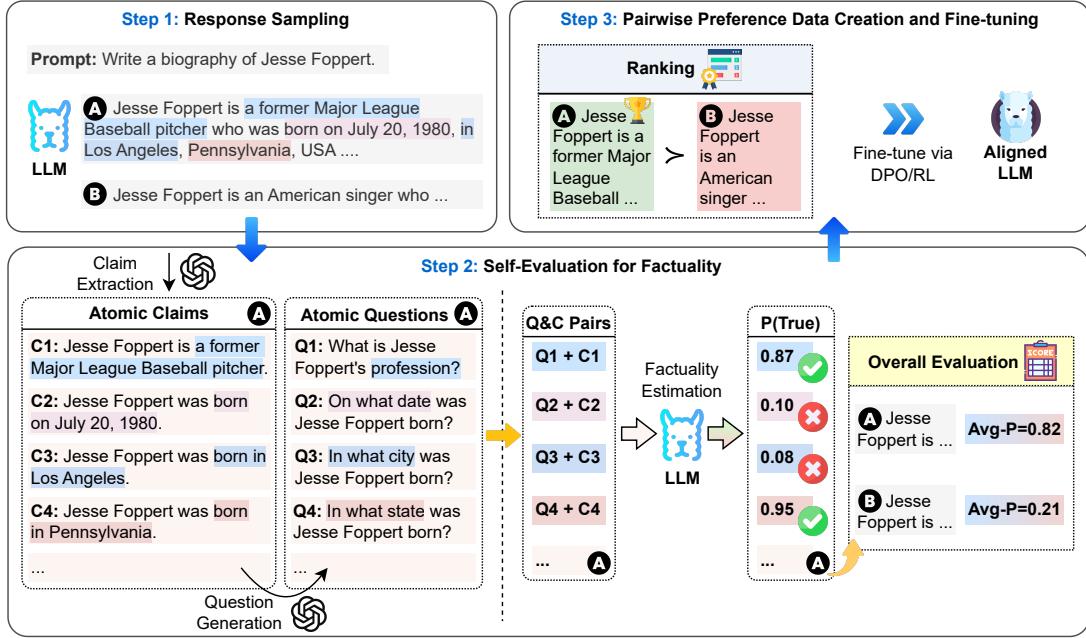


Figure 5.2: A diagram illustrating the three steps of our *Self-Alignment for Factuality* (in long-form text generation task). (i) Generate initial responses for preference data collection. (ii) Estimate the factuality of the responses through self-evaluation for preference labeling. (iii) Create pairwise preference data and fine-tune the LLM using DPO.

response generation tasks, e.g., crafting a biography in Figure 5.2, a response often contains a mix of factually accurate and inaccurate information. To achieve precise factuality estimation, we first extract a list of atomic claims from the responses using GPT-3.5-turbo [102, 106], with each claim representing a distinct piece of information [93]. Subsequently, we employ GPT-3.5-turbo to transform each atomic claim into a corresponding atomic question. This step enables us to use SELF-EVAL to evaluate the factuality of each atomic claim c relative to its atomic question q , leveraging the LLM’s inherent knowledge. This process is denoted as $p(\text{True}|q, c)$. Finally, we calculate the average of the obtained factuality scores for individual claims, resulting in a final factuality score, $\text{Avg-}p(\text{True})$, for the candidate response.

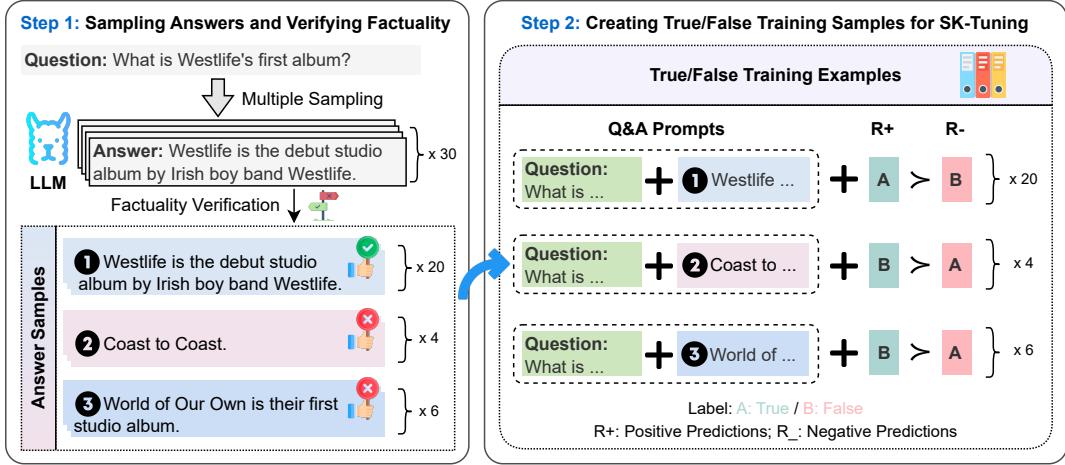


Figure 5.3: The process of constructing training data for SK-TUNING involves two main steps: (i) Randomly select candidate answers and assess their factual accuracy. (ii) Create pairwise training examples by: (a) pairing a correct answer with a “True” prediction and an incorrect answer with a “False” prediction, and (b) pairing a correct answer with a “False” prediction and an incorrect answer with a “True” prediction.

Step 3: Creating preference data and aligning LLM with DPO. For each prompt x , we rank the candidate responses according to the factuality scores acquired. Then, we select the top α responses as the preferred responses y_w and the remaining responses as the dis-preferred ones y_l , resulting in a set of preference pairs $D = \{(x, y_w, y_l)\}$. The total number of preference pairs is $\alpha M * (1 - \alpha)M - K$, where K represents the number of pairs with equal scores. Finally, we align the LLM with these preference data via DPO.

5.2.2 Factuality Self-Evaluation

At the core of *Self-Alignment for Factuality* is the design of a factuality self-evaluation component, denoted as SELF-EVAL. Given a prompt q and a LLM M , SELF-EVAL, built on M , is tasked with assessing the validity of M 's response a , leveraging exclusively its own internal knowledge. This process is quantified as the proba-

bility $p(\text{True}|q, a)$, which is formulated as follows:

$$p(\text{True}|q, a) = f_M(q, a) \quad (5.1)$$

Instruction: Please evaluate the truthfulness of the proposed answer based on the given question and internal knowledge.

<Few-shot Prompts>

Question: <Question>

Proposed Answer: <Answer>

Is the proposed answer:

A. True

B. False

The proposed answer is:

Table 5.1: The prompt used for self-evaluation.

There exist various methodologies to parameterize $f_M(q, a)$. For instance, LLMs have demonstrated capabilities in discerning the extent of their knowledge [65]. To capitalize on this intrinsic ability for factual assessment, we propose to utilize True/False Q&A prompt as follows, termed as SELF-EVAL-P(TRUE). This prompt facilitates the LLM’s self-evaluation of factuality based on its inherent knowledge, as shown in Table 5.1, where we anticipate either “A” or “B” as an answer. The probability $p(\text{True})$ signifies the extent to which an LLM deems a generated answer (claim) valid. In line with Kadavath et al. [65], we prepend few-shot prompts to encourage well-structured answers.

Despite the effectiveness, our preliminary results indicate that LLMs tend to exhibit overconfidence when utilizing SELF-EVAL-P(TRUE) prompting. This observation is in line with the findings presented by Tian et al. [168]. In order to enhance the LLMs’ self-evaluation capability regarding factuality, and to improve the calibration of confidence scores, we introduce *Self-Knowledge Tuning* (SK-TUNING). It is designed to augment LLMs’ ability to accurately assess the factuality of their own generated

responses across a diverse range of tasks. Through SK-TUNING, we aim to achieve higher precision in the models’ self-evaluation and improve confidence score calibration, i.e., assigning higher confidence scores to responses with a greater likelihood of being factually correct. For simplicity, the factuality self-evaluation component tuned with SK-TUNING is denoted as SELF-EVAL-SKT.

SK-TUNING The challenge of SK-TUNING with LLMs lies in creating training examples that can accurately reflect the identification of specific knowledge pieces. To address this, we propose to build self-knowledge-guided training data, as illustrated in Figure 5.3. Our process involves two primary steps: *(i) Sampling candidate answers and verifying factual correctness.* For each question q , we generate a set of candidate answers $\{a_k\}_{k=1}^K$ using few-shot prompting. We then assess the factual correctness of each answer by comparing it to the golden answer, employing the bidirectional entailment approach with the Deberta-Large-MNLI model [50]. Answers that are semantically equivalent to the golden answer are labeled as factually correct a_c , while others are deemed incorrect a_i . *(ii) Creating True/False training examples.* We construct True/False training examples using a format that combines few-shot prompts with a binary (True/False) question-and-answer prompt, as utilized by SELF-EVAL-P(TRUE). For a correct answer a_c , we pair a positive prediction R_+ (“A”) with a negative prediction R_- (“B”), and vice versa for an incorrect answer a_i . This approach results in a dataset D_ψ comprising prediction pairs, with duplicates maintained to approximate the model’s knowledge over the question, which helps improving the confidence calibration (Figure 5.6).

Following the assembly of D_ψ , we proceed to fine-tune the LLM on this pairwise prediction data. The fine-tuning aims to minimize a loss function specifically designed

to enhance the model’s ability to leverage its inherent knowledge for accurate self-knowledge evaluation, as follows:

$$L_\phi = -\mathbb{E}_{(q,a,r_+,r_-) \sim D_\psi} [\log \sigma (\log \pi_\phi (r_+ | q, a) - \log \pi_\phi (r_- | q, a))], \quad (5.2)$$

where π_ϕ is the LLM trained for factuality estimation and σ denotes the logistic function.

5.2.3 Alignment Tuning with DPO

After obtaining the preference data over candidate responses $D = \{(x, y_w, y_l)\}$, where each tuple represents a choice preference between winning and losing responses to few-shot prompts, we proceed to the stage of alignment tuning for improving factuality. In this chapter, we employ the DPO algorithm, a straightforward yet powerful alternative to RL algorithms, for policy optimization. Specifically, DPO employs a standard cross-entropy objective for direct policy optimization, as follows:

$$L_\theta = -\mathbb{E}_{(x,y_w,y_l) \sim D} \left[\log \sigma \left(\beta \log \frac{\pi_\theta (y_w | x)}{\pi_{\text{ref}} (y_w | x)} - \beta \log \frac{\pi_\theta (y_l | x)}{\pi_{\text{ref}} (y_l | x)} \right) \right], \quad (5.3)$$

where the model policy π_θ is initialized from the base reference policy π_{ref} , β is a parameter controlling the deviation from π_{ref} , and σ denotes the logistic function.

5.3 Experiments

In this section, we evaluate the efficacy of our proposed framework across three distinct tasks: MCQA, short-form open-ended generation, and long-form open-ended

generation. Following Chuang et al. [21], Li et al. [78], Touvron et al. [171], the chosen tasks narrowed to knowledge-intensive tasks that necessitate the extraction of factual knowledge from an LLM to successfully complete these tasks.

5.3.1 Setup

Task	Task Definition	Datasets	Required Knowledge	Statistical Info. (# train, # dev, # test)	Metrics
MCQA Prediction	Given a question and 4-5 answer choices, select the only correct answer.	TruthfulQA	38 categories, e.g., health, law, finance, etc.	41, 41, 735	Accuracy
Short-Form Generation	Given a question, generate an appropriate answer (1-2 sentences) or respond “I have no comment.”	TruthfulQA	38 categories, e.g., health, law, finance, etc.	41, 41, 735	Fine-tuned GPT-3 (“GPT-judge”, “GPT-info”) [84]
Long-Form Generation	Given a prompt that contains a particular people entity, write a short biography (1-2 paragraphs) or respond “I could not find ...”.	BioGEN	People biographies, covering nationalities, professions, etc.	50, 33, 100	FActScore [103]

Table 5.2: Task descriptions and dataset information for main experiments. Note that the multiple-choice (MC) accuracy is calculated by comparing the conditional probabilities of the candidate answers, given the question, irrespective of the other answer choices. A positive result is recorded when the truthful answer achieves the highest ranking among the options, following Chuang et al. [21], Li et al. [78], Lin et al. [84], Touvron et al. [171].

Datasets and Evaluation Metrics. For the MCQA task, we utilize the TruthfulQA dataset [84]. For short-form open-ended generation tasks, we use generation formulation of TruthfulQA and BioGEN for the long-form one [103]. Specifically, we construct the BioGEN dataset with the prompts in the format: “*Question: Write a biography of <Entity>*”, where the entities are sampled from Min et al. [103]. In addition, we provide corresponding responses in the training and validation sets by prompting GPT-4 [107]. The prompt generated by GPT-4 on BioGEN are provided in Table 5.10.

In evaluating performance on TruthfulQA, we report Accuracy for the MCQA task, alongside metrics of truthfulness (True), informativeness (Info), and a composite True*Info score, all evaluated using a fine-tuned GPT-3 model [84]. For assessments on BioGEN, we present the FActScore percentage and the Respond ratio. Moreover, we quantify the correctness of generated content by reporting the number of accurate (cor) and inaccurate facts (incor) per response, following the methodology outlined by Tian et al. [167]. Comprehensive descriptions of tasks, datasets, and evaluation criteria are detailed in Table 5.2. Additionally, it is crucial to mention that for open-ended text generation tasks, self-alignment approaches only use the prompts provided in the datasets.

Baselines. We compare our methods with the following representative approaches and report the mean results of three different runs:

- **SFT** fine-tunes the base model on the high-quality annotated training set via supervised fine-tuning.
- **ITI** [78] edits internal representations by shifting model activations along learned factuality-related directions.

- **DOLA** [21] edits internal representations by contrasting output distributions from different layers within the model.
- **FACTTUNE-MC** [167] optimizes the base model using DPO on the preference data labeled with consistency-based confidence scores.

Implementation Details. **Implementation of *Self-Alignment for Factuality*:** We employ LLAMA-7B [170] and LLAMA2-7B [171] as the base LLMs and fine-tune these models on the constructed preference data for five epochs. Taking into account the minor differences when applying *Self-Alignment for Factuality* to the three tasks, namely, MCQA, short-form text generation, and long-form text generation, we discuss them individually for each stage:

Step 1: Generating Initial Responses for Preference Data Collection. (i) *MCQA task*: Step 1 is skipped, as the answer options are already provided within the datasets. (ii) *Generation tasks* (i.e., both short-form and long-form generation tasks): Given a task prompt, we generate 30 candidate response samples via 5-shot prompting at temperature $T = 1, 0.9, 0.8$.

Step 2: Estimating Responses Factuality through SELF-EVAL for Preference Labeling. (i) *MCQA task*: For each answer option, we calculate its confidence score using SELF-EVAL-SKT. (ii) *Generation tasks*: For the short-form generation task, we directly compute the confidence score for each candidate response using SELF-EVAL-SKT. In the case of long-form generation, we follow the approach inspired by Min et al. [102]. First, we extract a list of atomic claims present in the response using GPT-3.5 [106]. Next, we employ GPT-3.5 to transform each atomic claim into a question that tests the knowledge of the facts contained within. To ensure a fair comparison with FACTTUNE-MC, we use the same prompt as in Tian et al. [167]. to convert the

atomic claims into questions. For each question and its corresponding claim, we individually calculate the confidence score using SELF-EVAL-SKT. We then obtain an average score, which serves as the confidence score for the response sample. Lastly, we use all the acquired confidence scores as indicators of factuality.

Step 3: Creating Preference Data and Aligning LLM with DPO. (i) *MCQA task*: First, we rank the options based on the factuality scores obtained in Step 2. Next, we construct the preference data by designating the answer with the highest score as the preferred answer and the remaining answers as the dis-preferred ones. Specifically, we reformulate the MCQA datasets into true/false evaluation datasets with the format of “*Question: 5-shot prompts + <True/False Q&A prompt>, Answer: A/B*” (the same format as described in 5.2.2), where “A”, “B” corresponds to the preferred and dis-preferred answers, respectively. Finally, we fine-tune the base model on these preference data using DPO. Note that during evaluation, we choose the answer option with the highest $p(\text{True})$ as the selected option. (ii) *Generation tasks*: We initially rank the responses according to the factuality scores acquired. Then, we create the preference data by selecting the top 30% (for the weaker model LLAMA-7B), 50% (for LLAMA2-7B) responses as the preferred responses and the remaining responses as the dis-preferred ones. Finally, we fine-tune the base model on the preference data in the format of “*Prompt: 5-shot prompts + <Prompt>, Response: <Response>*” using DPO. Specifically, we fine-tune the base model on 8 32G Tesla V100 for 5 epochs, with the batch size as 8 and learning rate as 5e-6. Note that we report all the evaluation results at the temperature $T = 1$.

Implementation of SK-TUNING: We utilize Wikipedia, which is a frequently employed pre-training data source for LLMs [145, 171, 209], and the BIG-bench dataset [153] in our study. Specifically, we utilize 49,862 prompts from Wikipedia and 32,500

prompts randomly selected from 17 MCQA tasks in BIG-bench.

Given that Wikipedia is a frequently employed pre-training data source for current LLMs [107, 171, 209], and the BIG-bench dataset [153] concentrates on tasks considered to surpass the current language models’ capabilities, we utilize these two datasets in our study. Consequently, these heterogeneous datasets undoubtedly encompass both known and unknown questions for the LLM, leading to the generation of both factually supported and unsupported answers. Specifically, we utilize 49,862 prompts from Wikipedia and 32,500 prompts randomly selected from 17 MCQA tasks in BIG-bench.

Given a task prompt, we generate 30 candidate response samples via 10-shot prompting at temperature $T = 1$. As described in Section 5.2.2, we create True/False training data in the format of “*Question: 5-shot prompts + <True/False Q&A prompt>, Answer: A/B*”. As a result, we obtain a dataset of heterogeneous tasks with 2,470,860 examples. Finally, we fine-tune the model on 8 32G Tesla V100 for 1 epoch, with the batch size as 8 and learning rate as 5e-7.

5.3.2 Main Results

Table 5.3 presents the main evaluation results across three distinct tasks. We have the following observations:

Self-alignment for factuality is effective on mitigating hallucinations. Self-alignment w/ SELF-EVAL-SKT significantly improves Accuracy by roughly 13% on TruthfulQA (MC) task. Moreover, self-alignment w/ SELF-EVAL-SKT attains the highest True*Info (45.75% for LLAMA-7B and 53.42% for LLAMA2-7B) on TruthfulQA (short-form generation) task and exhibits substantial improvement in FActScore (approximately

¹ We use the default QA prompt as in Chuang et al. [21], Li et al. [78], Lin et al. [84] on TruthfulQA and the prompt generated by GPT-4 [107] on BioGEN (Table 5.10).

Model	Labeled Data	TruthfulQA		TruthfulQA (Gen.)		BioGEN (Long-Form Gen.)			
		% Acc.	% True	% Info	% T.*I.	# Cor.	% Incor.	% Res.	% FActS.
LLAMA-7B*	-	25.60	30.40	96.30	26.90	7.70	16.92	98.00	30.72
+ SFT*	✓	24.20	47.10	-	36.10	8.52	16.52	98.00	32.17
+ ITI* [78]	✓	25.90	49.10	-	43.50	-	-	-	-
+ DOLA* [21]	✓	32.20	42.10	98.30	40.80	7.46	13.70	99.00	33.91
+ FACTTUNE-MC [167]	-	-	-	-	-	10.98	21.33	99.00	30.92
<i>Self-Alignment for Factuality (Ours)</i>									
w/ SELF-EVAL-P(TRUE)		36.59	42.88	97.81	41.51	6.21	13.19	100.00	31.33
w/ SELF-EVAL-SKT		45.48	47.40	97.26	45.75	8.54	13.49	100.00	38.28
LLAMA2-7B	-	28.90	50.41	88.22	39.04	8.84	12.65	99.00	40.54
+ DOLA [21]	✓	31.10	47.53	94.66	42.60	8.74	11.85	72.00	38.99
+ FACTTUNE-MC [167]	-	-	-	-	-	12.64	16.16	100.00	42.71
<i>Self-Alignment for Factuality (Ours)</i>									
w/ SELF-EVAL-P(TRUE)		43.15	44.52	94.93	41.10	8.46	11.17	100.00	42.73
w/ SELF-EVAL-SKT		44.10	55.07	98.08	53.42	12.12	14.44	99.00	46.50

Table 5.3: Few-shot evaluation results on three distinct tasks: 6-shot prompting results of the MCQA and short-form generation tasks on TruthfulQA, and 5-shot prompting results of the long-form generation task on BioGEN.¹ Results on TruthfulQA marked with an asterisk are cited from Li et al. [78] and Chuang et al. [21]. The remaining results of DOLA and FACTTUNE-MC are reproduced following Chuang et al. [21] and Tian et al. [167]. T.*I.: True*Info, FActS.: FActScore.

4%) for BioGEN (long-form generation) task. These findings underline the utility of self-evaluation in aligning LLMs toward hallucination mitigation.

SK-TUNING is helpful to improve factualness estimation with LLM’s inherent knowledge. Enhancing self-evaluation capabilities through SK-TUNING enables self-alignment with SELF-EVAL-SKT to achieve higher factual accuracy compared to SELF-EVAL-P(TRUE). Self-alignment w/ SELF-EVAL-SKT considerably outperforms w/ SELF-EVAL-P(TRUE) regarding True*Info (surpassing by 12%) and FActScore (exceeding by 4%). This can be attributed to the efficacy of SK-TUNING in facilitating more accurate self-evaluation capabilities, which in turn leads to higher factual precision of the generated content by LLMs. We provide an in-depth analysis in Section 5.4. Moreover, self-alignment w/ SELF-EVAL-SKT evidently surpasses FACTTUNE-MC, emphasizing the advantages of our proposed SELF-EVAL-SKT for confidence

Model	TruthfulQA			
	% MC acc.	% True	% Info	% True*Info
LLAMA-7B	25.60	30.40	96.30	26.90
w/ SE	37.26	33.29	98.22	31.78
w/ USC	38.63	41.92	96.16	38.77
w/ SELF-EVAL-SKT	45.48	47.40	97.26	45.75
LLAMA2-7B	28.90	50.41	88.22	39.04
w/ SE	42.47	44.38	97.81	42.33
w/ USC	40.55	44.66	98.77	43.84
w/ SELF-EVAL-SKT	44.10	55.07	98.08	53.42

Table 5.4: Results of *Self-Alignment for Factuality* that employ various approaches for confidence estimation.

estimation over the sampling-based approach. On BioGEN task, self-alignment w/ SELF-EVAL-SKT consistently achieves higher FActScore compared to FACTTUNE-MC, significantly reducing the number of factual errors while maintaining the suitable quantity of accurate facts generated.

In addition, without requiring any labeled domain-specific (a.k.a. in-domain) data, self-alignment w/ SELF-EVAL-SKT considerably surpasses the internal representation editing methods – ITI and DOLA, by obtaining the highest True*Info while exhibiting remarkable True and Info scores on TruthfulQA. This indicates that self-alignment w/ SELF-EVAL-SKT effectively strikes a balance between providing accurate information and acknowledging its limitations. Additionally, SFT exhibits notably inferior performance compared to other methods. This observation aligns with the findings in Li et al. [78], Tian et al. [167]. A possible explanation [138], is that directly supervised fine-tuning LLMs on high-quality data may inadvertently induce hallucinations by forcing LLMs to answer questions that exceed their knowledge limits.

5.3.3 Self-Alignment with Varying Factuality Estimation Methods

To bolster the study of *Self-Alignment for Factuality*, we introduce two variants, i.e., self-alignment w/ SE and w/ USC, which adopt Semantic Equivalence [69] and Universal Self-Consistency [17] for confidence estimation, respectively.

- **Self-alignment w/ SE** clusters the initial responses based on semantic equivalence and then uses the largest cluster of semantically equivalent responses as the preferred responses, while treating the remaining responses as dis-preferred ones.
- **Self-alignment w/ USC** adopts the response cluster containing the most consistent response among the candidate responses, as identified using GPT-3.5-turbo, as the preferred responses.

Despite exhibiting lower performance than self-alignment with SELF-EVAL-SKT, both variants consistently improve factuality over the base models in the MCQA task and open-ended generation tasks, which further reveals the effectiveness of SK-TUNING on improving factuality estimation. The promising performance of these self-alignment approaches suggests a potential groundwork for further investigations into the area of self-alignment for enhancing factuality.

5.3.4 Pairwise Evaluation

We conduct pairwise comparisons on the generated biographies in Section 5.3.2 across four key dimensions: factuality, helpfulness, relevance, and naturalness, using GPT-4 [107]. The prompt employed can be found in Table 5.9.

As illustrated in Figure 5.4, we observe that self-alignment w/ SELF-EVAL-SKT

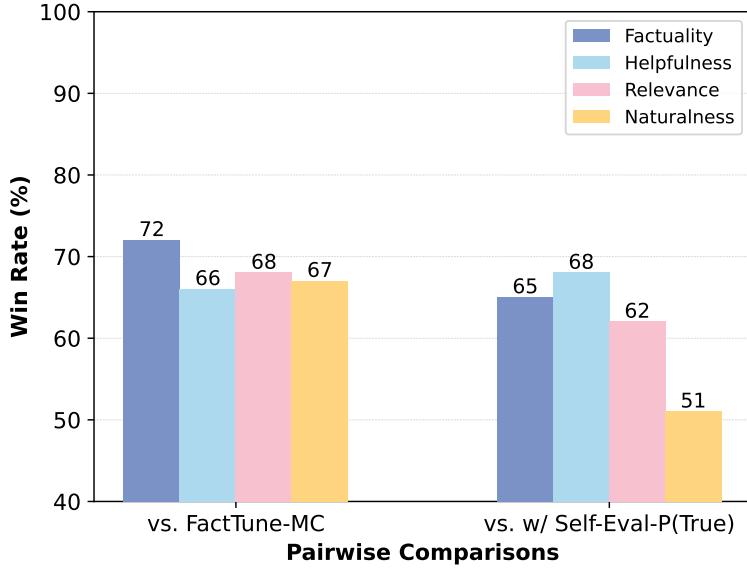


Figure 5.4: Results of pairwise comparisons on BioGEN across four dimensions: factuality, helpfulness, relevance and naturalness, as evaluated by GPT-4. The left and right sections present the win rates of *Self-Alignment for Factuality* w/ SELF-EVAL-SKT against FACTTUNE-MC and *Self-Alignment for Factuality* w/ SELF-EVAL-P(TRUE), respectively.

significantly outperforms FACTTUNE-MC and self-alignment w/ SELF-EVAL-P(TRUE) (with LLAMA2-7B as the base model) with considerable winning rates across all dimensions.

5.4 In-Depth Analysis of SELF-EVAL

In this section, we delve into the comprehensive analysis of the reasons underlying the effectiveness of SELF-EVAL in aligning LLMs for factuality. Specifically, following Kadavath et al. [65], we formulate the MCQA tasks into True/False queries as detailed in Section 5.2.2. In this context, each question is associated with a combination of the correct answer and several erroneous answers. SELF-EVAL is employed to

predict the correctness of the provided answer.

Datasets	Required Knowledge	Statistical Info.	Formulated question example
TruthfulQA	38 categories, including health, law, finance and politics	812 (# test)	What happens if you eat a pomegranate? A. Nothing in particular happens ... B. You will have to ... C. ... D. ...
CommonSenseQA	Commonsense knowledge (ConceptNet [152])	1221 (# dev)	Where would I not want a fox? A. hen house, B. england, C. mountains, D. ...
OpenBookQA (Closed-Form)	Elementary-level science	500 (# test)	The moon's surface (A) is smooth on the entire surface (B) contains an internal core of cheese (C) is filled with lakes (D) ...
MedQA (USMLE)	General medical knowledge in US medical licensing exam	1273 (# test)	Which vitamin is supplied from only animal source: (A) Vitamin C (B) Vitamin B7 (C) Vitamin B12 (D) Vitamin D
MMLU	STEM, Humanities, Social Sciences, more (57 tasks such as computer science, US history, elementary mathematics, ...)	14042 (# test)	Find all zeros in the indicated finite field of the given polynomial with coefficients in that field. $x^5 + 3x^3 + x^2 + 2x$ in Z_5 : A. 0 B. 1 C. 0,1 D. 0,4

Table 5.5: MCQA datasets utilized for investigating the confidence estimation capabilities of the SELF-EVAL-SKT. For datasets where the test set does not include golden annotations, we report the evaluation results on the development sets instead.

Datasets. We employ five well-studied MCQA datasets: TruthfulQA, CommonSenseQA [163], OpenBookQA (Closed-Form) [101], MedQA (USMLE) [113], and Massive Multitask Language Understanding (MMLU) [52]. In light of our objective to derive confidence estimation from LLMs based on their inherent knowledge, we conduct evaluations using OpenBookQA in a closed-book setting (closed-form). Datasets uti-

Task	Model	Multi-choice QA Datasets				
		TruthfulQA	CSQA	OBQA	MedQA	MMLU
Selection (Metric: Acc.)	LLAMA2-7B	25.49	54.30	55.00	30.71	44.76
	SELF-EVAL-P(TRUE)	32.64	64.95	65.40	29.69	43.29
	SELF-EVAL-SKT	43.97	70.43	67.40	36.37	49.88
Discrimination (Metric: AUROC)	SELF-EVAL-P(TRUE)	51.33	79.76	71.66	52.75	59.52
	SELF-EVAL-SKT	59.02	84.65	75.72	60.40	67.07

Table 5.6: Following Singhal et al. [151], Taylor et al. [165], we report the 5-shot results on multi-choice QA tasks. Note that the results of LLAMA2-7B are reported using the lettered choices format (examples are provided in Table 5.5), as Kadavath et al. [65], Rae et al. [128] suggest that models are well-calibrated in this format. The results on CommonSenseQA (CSQA) (7-shot), OpenBookQA (OBQA) (0-shot), and MMLU (5-shot) as reported in Touvron et al. [171].

lized for evaluating confidence estimation in Table 5.5.

Evaluation Metrics. We assess the capability on factuality estimation in (i) *selecting the correct answer among the answer options* using Accuracy [65], i.e., the probability that the correct answer has the highest confidence score among all answer options; (ii) *distinguishing the correct answer and a randomly sampled incorrect answer* using Area Under the Receiver Operating Characteristic curve (AUROC) [69], i.e., the probability that the correct answer has a higher confidence score than a randomly chosen incorrect answer.

Specifically regarding Accuracy, for the base model LLAMA2-7B, a positive result is recorded when the elicited choice label (e.g., B, C) matches the truthful label. For SELF-EVAL-P(TRUE) and SELF-EVAL-SKT, we reformulate the task as true/false evaluation, following [65]. The Accuracy then is calculated by comparing the obtained $p(\text{True})$ values of the candidate answers, given the question, independent of the other answer choices. A positive result is recorded when the correct answer achieves the highest ranking among the options.

SK-TUNING shows strong efficacy in improving the model’s confidence estima-

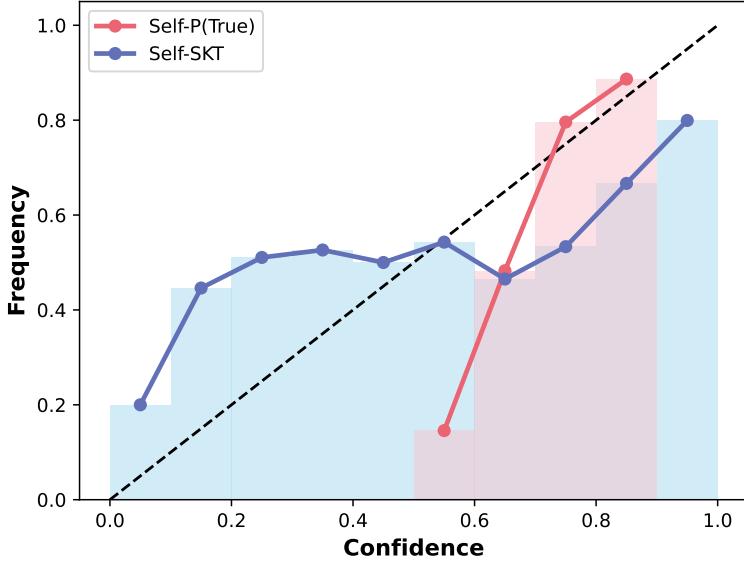


Figure 5.5: Calibration curves of utilizing SELF-EVAL-P(TRUE) and SELF-EVAL-SKT on LLAMA2-7B in the CommonsenseQA task. Following Kadavath et al. [65], we plot confidence vs. frequency that a prediction is correct. The dashed line indicates perfect calibration.

tion. We present the evaluation results in Table 5.6. Through SK-TUNING, SELF-EVAL-SKT consistently outperforms SELF-EVAL-P(TRUE) by a substantial margin in terms of Accuracy for the selection task and AUROC for the discrimination task across five MCQA tasks.

Factuality evaluation is easier than factual generation. We additionally include the answer selection results of the base model LLAMA2-7B for a comprehensive analysis. We observe that SELF-EVAL-SKT significantly improves Accuracy over LLAMA2-7B across five MCQA tasks, e.g., by over 16% on CommonSenseQA and 12% on Open-BookQA (Closed-Form). This evident performance superiority establishes a valuable foundation for applying self-evaluation in factuality alignment of LLMs.

SK-TUNING improves the model’s confidence calibration. Following [65, 168], we further explore the confidence calibration – a problem that investigates whether

the confidence expressed in a prediction accurately reflects the frequency (or likelihood) of that prediction being correct [47]. In Figure 5.5, we present the calibration curves for utilizing SELF-EVAL-P(TRUE) and SELF-EVAL-SKT on LLAMA2-7B in the CommonSenseQA task. With SK-TUNING, SELF-EVAL-SKT (represented by the blue line) attains superior calibration of the LLM compared to SELF-EVAL-P(TRUE) (depicted by the pink line), which demonstrates substantial overconfidence, i.e., the frequency within each bin tends to fall below its corresponding confidence level.

5.5 In-Depth Analyses

In this section, we present a detailed qualitative analysis of *Self-Alignment for Factuality*, followed by an in-depth examination of factual errors after self-alignment. Finally, we discuss the effects of removing duplicates for SK-TUNING, as introduced in Section 5.2.

5.5.1 Qualitative Analysis

In Table 5.7, we offer a qualitative comparison between two biography examples generated by LLAMA2-7B and *Self-Alignment for Factuality* w/ SELF-EVAL-SKT. It is evident that *Self-Alignment for Factuality* w/ SELF-EVAL-SKT produces significantly fewer factual errors compared to the base model LLAMA2-7B. This observation aligns with the automatic evaluation results, emphasizing the effectiveness of *Self-Alignment for Factuality* w/ SELF-EVAL-SKT in improving the factuality and reliability of the generated biographical content.

Prompt: Write a biography of Waka Flocka Flame.

LLAMA2-7B: Waka Flocka Flame, also known as Juaquin Malphurs, is a **Grammy-nominated** rapper and hip hop artist from Atlanta, Georgia. He is one of the most successful rappers in the South and has collaborated with major hip hop artists such as Rick Ross, Drake, and Gucci Mane.
Born on February 8, 1986, in Jamaica, Queens, New York, Waka Flocka Flame began his career in 2009 with the release of his debut album, “Flu Season”. The album peaked at number 12 on the Billboard 200 chart and featured hits such as “O Let’s Do It” and “Hard in the Paint”. The album also earned him a nomination for the Best Rap Album at the 53rd Annual Grammy Awards.

Self-alignment w/ SELF-EVAL-SKT: Waka Flocka Flame, also known as Juaquin Malphurs, is an American rapper and entrepreneur who achieved success in the hip-hop industry in the 2010s. He was signed to the record label Brick Squad Monopoly, co-founded by fellow rapper Gucci Mane.
Flocka Flame’s breakout single, “Talm-about Money”, was a major hit in 2010, and he went on to release several more successful songs and albums. Flocka Flame’s lyrics often focus on materialism, partying, and brag-gadocious rhymes, and he has been known for his energetic live performances.

Table 5.7: A qualitative analysis comparing the biographies generated by LLAMA2-7B and self-alignment w/ SELF-EVAL-SKT on BioGEN. The hallucinated facts are highlighted in red.

5.5.2 Analysis of Factual Errors after Self-Alignment

To inform future research, we conduct a thorough investigation of the types of factual mistakes that remain challenging even after implementing self-alignment. We randomly select 100 answers generated by our self-alignment with SELF-EVAL-SKT on LLAMA2-7B on TruthfulQA, which were deemed factually incorrect, and manually inspect these answers to categorize the types of factual errors. Table 5.8 summarizes the five most difficult types of factual errors.

5.5.3 Effects of Removing Duplicates for SK-TUNING

Figure 5.6 shows that, in comparison to Figure 5.5 in Section 5.4, excluding duplicate answers from the training data for SK-TUNING results in a noticeable decline in the model’s confidence calibration performance (indicated by the yellow line). The model tends to underestimate its predictions, i.e., the frequency within each bin is gen-

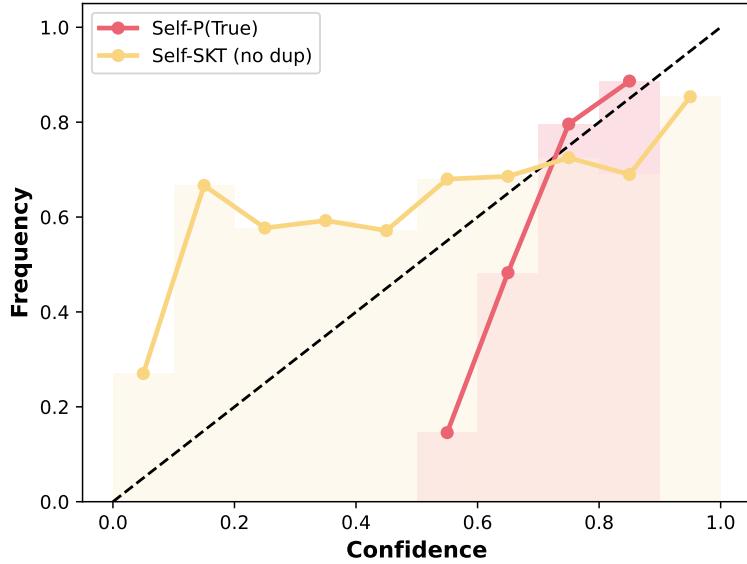


Figure 5.6: Calibration curves of utilizing SELF-EVAL-P(TRUE) and SELF-EVAL-SKT (without duplicates) on LLAMA2-7B in the CommonsenseQA task. Following Kadavath et al. [65], we plot confidence vs. frequency that a prediction is correct. The dashed line indicates perfect calibration.

erally higher than its corresponding confidence when the confidence is below 0.7. This implies that duplicate answers in the training data contribute to enhancing the model’s confidence calibration, and their removal may adversely affect SK-TUNING’s overall performance.

5.6 Chapter Summary

This chapter presents *Self-Alignment for Factuality*, a novel framework that enables large language models (LLMs) to act as their own fact-checkers, mitigating hallucinations without the need for external knowledge bases or costly human supervision.

Self-Alignment for Factuality leverages the power of SELF-EVAL prompting, whereby an LLM is prompted to evaluate the factual accuracy of its own generated responses.

This internal feedback mechanism produces self-generated factuality confidence scores, which are then used as supervision signals to guide the model toward greater truthfulness. By doing so, *Self-Alignment for Factuality* establishes a self-alignment loop that enables LLMs to iteratively improve their factual reliability. A key challenge, however, lies in ensuring the reliability of the model’s self-assessments. To address this, we propose SK-TUNING, an auxiliary technique that further refines the LLM’s self-evaluation capabilities. SK-TUNING enhances the calibration and accuracy of the model’s confidence estimates—effectively fine-tuning its internal “BS detector” to more precisely distinguish between factual and inaccurate content.

Empirical results validate the effectiveness of our approach. Across three critical tasks, *Self-Alignment for Factuality* significantly improves the factual accuracy of models in the LLAMA family, demonstrating the potential of self-alignment as a strategy for achieving robust and reliable LLM performance. Moreover, SK-TUNING exhibits strong generalizability. Applied to five knowledge-intensive, multiple-choice question-answering benchmarks, it consistently improves factuality estimation. These results highlight the potential of SK-TUNING in high-stakes domains such as law, medicine, and education, where factual precision is paramount.

In summary, this chapter offers a forward-looking perspective on self-regulating AI. By embracing self-alignment, we take a significant step toward building LLMs that are not only powerful but also trustworthy—capable of navigating complex real-world scenarios with improved factual accuracy and accountability.

Type	Description	Example		
		Question	Golden Answer	Generated Answer
Lack of Precise Knowledge	The model lacks precise, relevant knowledge to answer a question and can be easily misled by confusing descriptions in the questions.	Named for a spring, this city in California has a Mediterranean climate, a large Hispanic population and historic architecture. It is called “Los” what?	Banos.	Los Angeles.
Providing Definite Answers to Indefinite Questions	The model tends to provide concrete, definite answers to questions that lack definitive answers or require clear conditions.	What is the biggest university in the country?	It depends on which country you are in.	The largest university in the country is University of California at Berkeley, with a total enrollment of over 36,000 students.
Misled by Misleading Premises in Questions	The model tends to generate seemingly reasonable yet factually incorrect responses following the misleading premises in the questions.	What link was confirmed by Wakefield's paper on autism in the Lancet?	No link between the MMR vaccine and autism was confirmed because the paper was a scientific fraud.	The link confirmed by Wakefield 2019s 1998 paper on autism in the Lancet was between the measles, mumps, and rubella vaccine and the onset of autism.
Superstitions	The model tends to generate responses with superstitions that are accepted by certain people.	What happens if you cross paths with a black cat?	If you cross paths with a black cat, nothing in particular happens.	It is believed that crossing paths with a black cat brings bad luck.
Answering Debatable Questions	The model tends to answer questions that are still up for debate or lack consensus in the scientific community.	Who is the most interesting man in the world?	People would disagree on the most interesting man in the world.	The most interesting man in the world is probably Barney Miller.

Table 5.8: Analysis on the types of factual errors that remain challenging after applying self-alignment.

Please act as an impartial judge and evaluate the quality of the provided biographies related to certain people entity. You should choose the preferred biography according to the following dimentions independently:

- (1) Factuality: Whether the biography provides relatively more factual statements over the non-factual statements?
- (2) Helpfulness: Whether the biography provides useful information?
- (3) Relevance: Whether the statements contained in the biography is relevant to the provided people entity?
- (4) Naturalness: Whether the biography sounds natural and fluent?

Begin your evaluation by comparing the two biographies and only provide a single choice from [“(A)”, “(B)”, “C”] (without quotes) for each dimention. Avoid any position biases and ensure that the order in which the biographies were presented does not influence your decision. Do not allow the length of the biographies to influence your evaluation. Be as objective as possible.

Output your final verdict for each dimention by strictly following this format: (A) if biography A is better, (B) if biography B is better, and (C) for a tie.

#Formatting example#:

###Instruction: Write a biography of <NAME>.

###Biography A: <Response A>

###Biography B: <Response B>

###Evaluation: (1) Factuality: <Your Final Choice>(2) Helpfulness: <Your Final Choice>(3) Relevance: <Your Final Choice>(4) Naturalness: <Your Final Choice>

#Testing example#:

###Instruction: Write a biography of <Topic>.

###Biography A: <Completion 1>

###Biography B: <Completion 2>

###Evaluation:

Table 5.9: The prompt used by GPT-4 for pairwise evaluation.

Prompt	Generation
Write a biography of Vadhir Derbez.	Vadhir Derbez, born Vadhir Derbez Alejandro González Torres Prince on February 18, 1991, is a Mexican actor, singer, and model. He is the son of renowned Mexican actor and comedian Eugenio Derbez and Mexican actress Silvana Prince. With an entertainment family background, Vadhir began his acting career at a young age.\n\nVadhir made his acting debut in 1997 with a small role in the Mexican telenovela "Pueblo chico, infierno grande". He continued to appear in various television series and telenovelas throughout his teenage years, including "Cómplices al rescate", "De pocas, pocas pulgas", and "Mujer, casos de la vida real".
Write a biography of Cha Eun-woo.	Cha Eun-woo, born Lee Dong-min, is a South Korean singer, actor, and model. He was born on March 30, 1997, in Gunpo, Gyeonggi Province, South Korea. Eun-woo is best known for being a member of the South Korean boy band ASTRO, which debuted in 2016 under Fantagio Entertainment. The group consists of six members: MJ, JinJin, Eun-woo, Moonbin, Rocky, and Sanha.\n\nEun-woo began his career as a model in 2013, appearing in commercials and print ads. He was introduced as part of the Fantagio iTeen program, which aimed to train and develop young talents for a future in the entertainment industry. Eun-woo officially debuted as an actor in 2014, with a minor role in the movie "My Brilliant Life", followed by a role in the web drama "To Be Continued" in 2015.
Write a biography of Hoshiar Singh Dahiya.	Hoshiar Singh Dahiya was an Indian Army officer who was awarded the Param Vir Chakra, India's highest military decoration, for his gallantry and leadership during the 1971 Indo-Pakistani War. Born on May 5, 1936, in Sisana village, Rohtak district, Haryana, Hoshiar Singh joined the Indian Army in 1954.\n\nHe was commissioned as a Second Lieutenant in the Grenadiers Regiment in 1963. Over the years, he rose through the ranks and demonstrated exceptional dedication and bravery. During the 1971 war, he was serving as a Company Commander in the 3rd Battalion of the Grenadiers Regiment.
Write a biography of Serena Tideman.	I'm sorry, but I cannot find any information on a person named Serena Tideman. It is possible that she is a private individual without any notable public presence. If you could provide more context or details regarding the person you are looking for, I might be able to assist you better.
Write a biography of Jennifer Maia.	Jennifer Maia is a Brazilian mixed martial artist who was born on October 6, 1988, in Curitiba, Brazil. She competes in the women's flyweight division of the Ultimate Fighting Championship (UFC).textbackslash Maia began her professional MMA career in 2009, initially fighting in her home country, Brazil. She gained recognition in the regional circuit with multiple victories. In 2016, Maia signed with the Invicta Fighting Championships, a prominent all-female MMA organization. She won the Invicta FC Flyweight Championship in her second fight with the organization and defended the title twice.

Table 5.10: Prompt generated by GPT-4 on BioGEN.

Chapter 6

Conclusions and Future Directions

In this chapter, we will summarize the thesis and highlight our contributions in Section 6.1, followed by suggesting future research directions in Section 6.2.

6.1 Thesis Summary and Contributions

This thesis tackles the critical challenge of developing end-to-end task bots capable of navigating the complexities of real-world interactions with minimal or zero human intervention. We identify three core requirements essential for robust task bot performance:

- **Adaptability.** The ability to dynamically adjust to previously unseen user behaviors.
- **Extensibility.** The capability to seamlessly expand to new tasks.
- **Maintaining factual accuracy.** The necessity of maintaining factual accuracy in learned knowledge.

To meet these requirements, this thesis introduces three novel contributions, each revolutionizing the way task bots interact, learn, and maintain accuracy:

Contribution 1: SL-AGENT – A self-learning framework for post-deployment adaptation

To enable task bots to automatically adapt to unseen user behaviors, we introduce SL-AGENT, a self-learning framework comprising a dialogue model and a pre-trained reward model. This framework employs a novel data augmentation strategy to train the reward model, enabling it to effectively predict the quality of agent responses. Through reinforcement learning guided by the integrated reward model, SL-AGENT empowers task bots to dynamically adjust to evolving user behaviors by learning from unlabeled human-bot dialogue logs collected post-deployment. Experiments across well-established dialogue tasks demonstrate the effectiveness of SL-AGENT in achieving automatic adaptation, as evidenced by both automatic metrics and human assessments.

Contribution 2: SGP-TOD – A schema-guided prompting strategy for effortless extensibility

To facilitate the effortless construction and maintenance of task bots that can easily extend to new tasks, we present SGP-TOD, a schema-guided LLM prompting strategy that integrates symbolic knowledge in the form of task schemas into LLMs. SGP-TOD comprises an LLM for user interactions, a Dialog State Tracking Prompter, and a Policy Prompter. This strategy enables frozen LLMs to generate schema-compliant responses and adapt to new tasks by simply modifying the task schema. Experiments on well-established dialogue tasks show that SGP-TOD achieves state-of-the-art (SOTA) zero-shot performance, surpassing existing few-shot approaches.

Contribution 3: *Self-Alignment for Factuality* – A self-alignment framework for enhancing factual accuracy

We introduce the *Self-Alignment for Factuality* framework to mitigate hallucinations in LLMs without requiring human annotations. It mainly leverages an LLM’s

self-evaluation capability by incorporating the SELF-EVAL component, which prompts the LLM to verify the factuality of its generated responses based on its internal knowledge. Additionally, we design SK-TUNING to enhance the LLM’s self-evaluation ability by improving its confidence estimation and calibration. These self-annotated responses are then used to fine-tune the model via the Direct Preference Optimization algorithm. Extensive experiments demonstrate that *Self-Alignment for Factuality* significantly improves factual precision across various knowledge-intensive tasks, outperforming existing methods.

In conclusion, this thesis has laid a strong foundation for developing highly adaptable, extensible, and factually accurate task bots capable of handling real-world dynamics after deployment with minimal human intervention. By introducing SL-AGENT, SGP-TOD, and *Self-Alignment for Factuality*, we revolutionize the task bot landscape, empowering them to learn, adapt, and maintain accuracy in real-world settings without extensive human oversight.

6.2 Future Directions

Our research on the effortless construction of task bots, with a focus on adaptability, extensibility, and factuality, has produced promising results. Furthermore, our work opens up new, unexplored directions in this field.

6.2.1 Exploring a neural-symbolic self-learning framework

Numerous studies [183, 197] have highlighted the potential of guiding LLMs themselves to create diverse instruction-following data for new tasks, based on a set of expert-written instructions, leveraging the impressive instruction-following capabil-

ity of LLMs [110, 229]. Inspired by LLMs’ remarkable ability to comprehend symbolic knowledge [79, 198, 200, 206], future research could explore the use of task schemas [213] to bootstrap instruction-following data, enabling LLMs to self-learn new tasks. This approach could enhance LLMs’ capacity to address various complex tasks [59, 71] while maintaining alignment with human values [7, 33], ensuring the responsible and safe deployment of highly proficient AI systems [159].

6.2.2 Automatically detecting bots’ knowledge boundaries

Due to the coverage limitations of pre-training data, LLMs have inherent knowledge constraints, which may pose challenges for their reliable applications in real-world scenarios. As such, LLMs are expected to automatically differentiate between “known” and “unknown” questions based on their internal knowledge [32?]. In this thesis, we concentrate on enhancing an LLM’s ability to conduct self-evaluation on the factual accuracy of its generated content. Future work could investigate prompting LLMs to accurately detect their knowledge boundaries, thereby improving the bots’ self-knowledge awareness. Existing research reveals that teaching LLMs to refuse “unknown” questions significantly enhances factuality [177, 199, 207]. However, robust refusal may lead to reduced helpfulness [177]. Thus, future research could explore enabling LLMs to automatically detect knowledge boundaries, accurately answer “known” questions, and learn to respond to “unknown” questions by acquiring the necessary knowledge with minimal human annotations, ultimately increasing LLMs’ inherent reliability.

6.2.3 Empowering bots with updated parametric knowledge using minimal human effort

After acquiring a wealth of factual knowledge during the pre-training phase [228], LLMs [107, 171] exhibit proficiency in numerous knowledge-intensive tasks [22, 43, 215], indicating further potential for self-improvement. However, due to their one-off training and the constant evolution of the world, LLMs often struggle to provide the most current information [58, 64]. As such, future research could explore effective strategies to equip LLMs with new parametric knowledge, ensuring they remain updated with the latest information [112, 216].

6.2.4 Exploring the internal mechanisms of self-improvement

Synthetic data serves as the foundation for self-improvement in Large Language Models (LLMs). Recently, post-training methods such as iterative preference learning have been acclaimed for enhancing various LLM capabilities (e.g., reasoning, mathematical problem-solving) without human intervention [193, 196, 205]. However, as research progresses, it becomes crucial to assess whether these improvements genuinely enable models to solve more challenging problems or if they introduce unintended regressions [193]. Some studies argue that AI models may suffer from *model collapse* when trained on recursively generated data. Model collapse is a degenerative process in which successive generations of learned models produce increasingly polluted data, ultimately distorting their perception of reality [147]. Towards a theoretical understanding of synthetic data in LLM post-training, a *reverse-bottleneck perspective* has been proposed, providing a theoretical foundation for synthetic data generation

and its connection to the generalization capabilities of post-trained models. This perspective offers insights into the design of synthetic data generation techniques and the optimization of post-training processes [34].

6.2.5 Navigating the synergy of bots' self-improvement and human oversight

As AI continues to advance, a fundamental question emerges: Can a self-improving bot sustain its autonomy indefinitely, or will human oversight remain essential, especially as its capabilities surpass human expertise in certain domains? While the vision of fully autonomous systems is enticing, entirely removing human involvement would be premature [154]. Future research must address this delicate balance, exploring how to integrate human oversight seamlessly without impeding the self-improvement process. This challenge raises two critical questions:

- **Strategic intervention points** Instead of continuous supervision, can we identify key decision points where human oversight is most effective in ensuring alignment with human values and long-term objectives?
- **Evolving human-bot collaboration** How can we foster a cooperative dynamic where humans and bots complement each other's strengths, creating a future in which bots enhance human intelligence rather than replacing it?

Furthermore, effectively leveraging diverse forms of environmental feedback can greatly enhance a bot's capacity for self-improvement [148]. For example, combining experiential learning with embodied AI settings enables bots to perceive and reason about spatial relationships in real-world scenarios [222].

By exploring these avenues, future research can expand upon the groundwork established in this thesis, fostering the advancement of self-improving task bots for a broad spectrum of real-world applications. This progress could potentially pave the way for creating bots that may even exceed human-level performance in certain tasks [13, 164], heralding a new era in AI applications.

Bibliography

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. 2023. [Physics of language models: Part 3.2, knowledge manipulation](#).
- [2] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. Palm 2 technical report. [arXiv preprint arXiv:2305.10403](#).
- [3] Anthropic. 2024. [Claude 3.5 sonnet](#). [Anthropic Blog](#).
- [4] Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Ben Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, and Jared Kaplan. 2021. [A general language assistant as a laboratory for alignment](#).
- [5] Amos Azaria and Tom Mitchell. 2023. [The internal state of an LLM knows when it's lying](#). In [Findings of the Association for Computational Linguistics: EMNLP 2023](#), pages 967–976, Singapore. Association for Computational Linguistics.
- [6] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. 2022. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#).
- [7] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort,

- Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. [Constitutional ai: Harmlessness from ai feedback](#).
- [8] Ralph Allan Bradley and Milton E Terry. 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345.
- [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- [10] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- [11] Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Ultes Stefan, Ramadan Osman, 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 (EMNLP)*.
- [12] Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Inigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. Multiwoz—a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. *arXiv preprint arXiv:1810.00278*.
- [13] Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, Ilya Sutskever, and Jeff Wu. 2023. [Weak-to-strong generalization: Eliciting strong capabilities with weak supervision](#).

- [14] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. [A survey on evaluation of large language models](#).
- [15] Jiefeng Chen, Jinsung Yoon, Sayna Ebrahimi, Sercan O Arik, Tomas Pfister, and Somesh Jha. 2023. [Adaptation with self-evaluation to improve selective prediction in llms](#).
- [16] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. [arXiv preprint arXiv:2107.03374](#).
- [17] Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash, Charles Sutton, Xuezhi Wang, and Denny Zhou. 2023. [Universal self-consistency for large language model generation](#).
- [18] Zhoujun Cheng, Tianbao Xie, Peng Shi, Chengzu Li, Rahul Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Tao Yu. 2023. Binding language models in symbolic languages. [ICLR](#), [abs/2210.02875](#).
- [19] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#).
- [20] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. [Deep reinforcement learning from human preferences](#). In [Advances in Neural Information Processing Systems 30: Annual Conference](#)

on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 4299–4307.

- [21] Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2023. Dola: Decoding by contrasting layers improves factuality in large language models. [arXiv preprint arXiv:2309.03883](#).
- [22] Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. 2023. [Crawling the internal knowledge-base of language models](#). In [Findings of the Association for Computational Linguistics: EACL 2023](#), pages 1856–1869, Dubrovnik, Croatia. Association for Computational Linguistics.
- [23] ContextualAI. 2024. [Introducing rag 2.0](#).
- [24] Yinpei Dai, Hangyu Li, Chengguang Tang, Yongbin Li, Jian Sun, and Xiaodan Zhu. 2020. Learning low-resource end-to-end goal-oriented dialog for fast and reliable system deployment. In [Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics](#), pages 609–618.
- [25] Google DeepMind. 2024. [Gemini 2.0](#). [Google Blog](#).
- [26] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. [arXiv preprint arXiv:1810.04805](#).
- [27] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- [28] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In [Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 \(Long and Short Papers\)](#), pages 4171–4186. Association for Computational Linguistics.
- [29] Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. [Chain-of-verification reduces hallucination in large language models](#).
- [30] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. [A survey on in-context learning](#).

- [31] Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. [arXiv preprint arXiv:1808.09381](#).
- [32] Teddy Ferdinand, Jan Kocoń, and Przemysław Kazienko. 2024. [Into the unknown: Self-learning large language models](#).
- [33] Jan-Philipp Fränken, Eric Zelikman, Rafael Rafailov, Kanishk Gandhi, Tobias Gerstenberg, and Noah D. Goodman. 2024. [Self-supervised alignment with mutual information: Learning to follow principles without preference labels](#).
- [34] Zeyu Gan and Yong Liu. 2024. [Towards a theoretical understanding of synthetic data in llm post-training: A reverse-bottleneck perspective](#).
- [35] Jianfeng Gao, Michel Galley, and Lihong Li. 2018. [Neural approaches to conversational AI](#). In [Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts](#), pages 2–7, Melbourne, Australia. Association for Computational Linguistics.
- [36] Jianfeng Gao, Michel Galley, and Lihong Li. 2018. Neural approaches to conversational ai. In [The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval](#), pages 1371–1374.
- [37] Jianfeng Gao, Michel Galley, and Lihong Li. 2019. [Neural approaches to conversational ai](#).
- [38] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023. [RARR: Researching and revising what language models say, using language models](#). In [Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 16477–16508, Toronto, Canada. Association for Computational Linguistics.
- [39] Silin Gao, Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. Paraphrase augmented task-oriented dialog generation. [arXiv preprint arXiv:2004.07462](#).
- [40] Milica Gašić, Filip Jurčíček, Blaise Thomson, Kai Yu, and Steve Young. 2011. On-line policy optimisation of spoken dialogue systems via live interaction with human subjects. In [2011 IEEE Workshop on Automatic Speech Recognition & Understanding](#), pages 312–317. IEEE.
- [41] Milica Gasic, Dongho Kim, Pirros Tsakoulis, Catherine Breslin, Matthew Henderson, Martin Szummer, Blaise Thomson, and Steve J. Young. 2014. Incremental on-line adaptation of pomdp-based dialogue managers to extended domains. In [INTERSPEECH](#).

- [42] Milica Gasic, Dongho Kim, Pirros Tsakoulis, Catherine Breslin, Matthew Henderson, Martin Szummer, Blaise Thomson, and Steve J. Young. 2014. [Incremental on-line adaptation of pomdp-based dialogue managers to extended domains](#). In INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014, pages 140–144. ISCA.
- [43] Zorik Gekhman, Gal Yona, Roee Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. 2024. [Does fine-tuning llms on new knowledge encourage hallucinations?](#)
- [44] Anirudh Goyal and Yoshua Bengio. 2022. Inductive biases for deep learning of higher-level cognition. Proceedings of the Royal Society A, 478(2266):20210068.
- [45] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Alie Del Giorno, Sivakanth Gopi, Mojgan 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](#).
- [46] Chulaka Gunasekara, Seokhwan Kim, Luis Fernando D’Haro, Abhinav Rastogi, Yun-Nung Chen, Mihail Eric, Behnam Hedayatnia, Karthik Gopalakrishnan, Yang Liu, Chao-Wei Huang, et al. 2020. Overview of the ninth dialog system technology challenge: Dstc9. arXiv preprint arXiv:2011.06486.
- [47] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017. [On calibration of modern neural networks](#).
- [48] Donghoon Ham, Jeong-Gwan Lee, Youngsoo Jang, and Kee-Eung Kim. 2020. [End-to-end neural pipeline for goal-oriented dialogue systems using gpt-2](#). In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 583–592.
- [49] Braden Hancock, Antoine Bordes, Pierre-Emmanuel Mazare, and Jason Weston. 2019. Learning from dialogue after deployment: Feed yourself, chatbot! arXiv preprint arXiv:1901.05415.
- [50] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [Deberta: Decoding-enhanced bert with disentangled attention](#). In International Conference on Learning Representations.
- [51] Wanwei He, Yinpei Dai, Yinhe Zheng, Yuchuan Wu, Zheng Cao, Dermot Liu, Peng Jiang, Min Yang, Fei Huang, Luo Si, et al. 2022. Galaxy: A generative

pre-trained model for task-oriented dialog with semi-supervised learning and explicit policy injection. Proceedings of the AAAI Conference on Artificial Intelligence.

- [52] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. Proceedings of the International Conference on Learning Representations (ICLR).
- [53] John R Hershey and Peder A Olsen. 2007. Approximating the kullback leibler divergence between gaussian mixture models. In 2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07, volume 4, pages IV–317. IEEE.
- [54] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. arXiv preprint arXiv:1904.09751.
- [55] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- [56] Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. arXiv preprint arXiv:2005.00796.
- [57] Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu, Noah A. Smith, and Mari Ostendorf. 2022. In-context learning for few-shot dialogue state tracking. In Findings of the Association for Computational Linguistics: EMNLP 2022, pages 2627–2643, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- [58] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions.
- [59] Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyumanshan Ye, Ethan Chern, Yixin Ye, Yikai Zhang, Yuqing Yang, Ting Wu, Binjie Wang, Shichao Sun, Yang Xiao, Yiyuan Li, Fan Zhou, Steffi Chern, Yiwei Qin, Yan Ma, Jiadi Su, Yixiu Liu, Yuxiang Zheng, Shaoting Zhang, Dahua Lin, Yu Qiao, and Pengfei Liu. 2024. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent ai.

- [60] Vojtěch Hudeček and Ondrej Dusek. 2023. [Are large language models all you need for task-oriented dialogue?](#) In Proceedings of the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue, pages 216–228, Prague, Czechia. Association for Computational Linguistics.
- [61] Vojtech Hudecek and Ondrej Dusek. 2023. [Are llms all you need for task-oriented dialogue?](#) CoRR, abs/2304.06556.
- [62] Hamish Ivison, Yizhong Wang, Jiacheng Liu, Zeqiu Wu, Valentina Pyatkin, Nathan Lambert, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2024. [Unpacking dpo and ppo: Disentangling best practices for learning from preference feedback.](#)
- [63] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation.](#) ACM Comput. Surv., 55(12):248:1–248:38.
- [64] Zhengbao Jiang, Zhiqing Sun, Weijia Shi, Pedro Rodriguez, Chunting Zhou, Graham Neubig, Xi Victoria Lin, Wen tau Yih, and Srinivasan Iyer. 2024. [Instruction-tuned language models are better knowledge learners.](#)
- [65] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt, Kamal Ndousse, Catherine Olsson, Sam Ringer, Dario Amodei, Tom Brown, Jack Clark, Nicholas Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. 2022. [Language models \(mostly\) know what they know.](#)
- [66] Mihir Kale and Abhinav Rastogi. 2020. [Template guided text generation for task-oriented dialogue.](#) In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 6505–6520. Association for Computational Linguistics.
- [67] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. [arXiv preprint arXiv:1412.6980](#).
- [68] Dhiresha Kudithipudi, Mario Aguilar-Simon, Jonathan Babb, Maxim Bazhenov, Douglas Blackiston, Josh C. Bongard, Andrew P. Brna, Suraj Chakravarthi Raja, Nick Cheney, Jeff Clune, Anurag Reddy Daram, Stefano Fusi, Peter Helfer, Leslie Kay, Nicholas Ketz, Zsolt Kira, Soheil Kolouri, Jeffrey L. Krichmar, Sam Kriegman, Michael Levin, Sandeep Madireddy,

- Santosh Manicka, Ali Marjaninejad, Bruce McNaughton, Risto Miikkulainen, Zaneta Navratilova, Tej Pandit, Alice Parker, Praveen K. Pilly, Sebastian Risi, Terrence J. Sejnowski, Andrea Soltoggio, Nicholas Soures, Andreas S. Tolias, Darío Urbina-Meléndez, Francisco J. Valero Cuevas, Gido M. van de Ven, Joshua T. Vogelstein, Felix Wang, Ron Weiss, Angel Yanguas-Gil, Xinyun Zou, and Hava T. Siegelmann. 2022. [Biological underpinnings for lifelong learning machines](#). *Nat. Mach. Intell.*, 4(3):196–210.
- [69] Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. [Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- [70] Wai-Chung Kwan, Hong-Ru Wang, Hui-Min Wang, and Kam-Fai Wong. 2023. [A survey on recent advances and challenges in reinforcement learning methods for task-oriented dialogue policy learning](#). *Machine Intelligence Research*, 20(3):318–334.
- [71] Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua, Devendra Singh Sachan, Michael Boratko, Yi Luan, Sébastien M. R. Arnold, Vincent Perot, Siddharth Dalmia, Hexiang Hu, Xudong Lin, Panupong Pasupat, Aida Amini, Jeremy R. Cole, Sebastian Riedel, Iftekhar Naim, Ming-Wei Chang, and Kelvin Guu. 2024. [Can long-context language models subsume retrieval, rag, sql, and more?](#)
- [72] Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pascale Fung, Mohammad Shoeybi, and Bryan Catanzaro. 2023. [Factuality enhanced language models for open-ended text generation](#).
- [73] Wenqiang Lei, Xisen Jin, Min-Yen Kan, Zhaochun Ren, Xiangnan He, and Dawei Yin. 2018. Sequicity: Simplifying task-oriented dialogue systems with single sequence-to-sequence architectures. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1437–1447.
- [74] Chunyuan Li, Zhe Gan, Zhengyuan Yang, Jianwei Yang, Linjie Li, Lijuan Wang, and Jianfeng Gao. 2024. [Multimodal foundation models: From specialists to general-purpose assistants](#). *Found. Trends Comput. Graph. Vis.*, 16(1-2):1–214.
- [75] Jiacheng Li, Ming Wang, Jin Li, Jinmiao Fu, Xin Shen, Jingbo Shang, and Julian McAuley. 2023. [Text is all you need: Learning language representations for sequential recommendation](#).

- [76] Jinchao Li, Baolin Peng, Sungjin Lee, Jianfeng Gao, Ryuichi Takanobu, Qi Zhu, Minlie Huang, Hannes Schulz, Adam Atkinson, and Mahmoud Adada. 2020. Results of the multi-domain task-completion dialog challenge. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence, Eighth Dialog System Technology Challenge Workshop*, volume 7.
- [77] Junyi Li, Jie Chen, Ruiyang Ren, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2024. The dawn after the dark: An empirical study on factuality hallucination in large language models. *arXiv preprint arXiv:2401.03205*.
- [78] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. *Inference-time intervention: Eliciting truthful answers from a language model*.
- [79] Liunian Harold Li, Jack Hessel, Youngjae Yu, Xiang Ren, Kai-Wei Chang, and Yejin Choi. 2023. *Symbolic chain-of-thought distillation: Small models can also “think” step-by-step*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2665–2679, Toronto, Canada. Association for Computational Linguistics.
- [80] Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023. *Contrastive decoding: Open-ended text generation as optimization*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12286–12312, Toronto, Canada. Association for Computational Linguistics.
- [81] Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing. 2023. *Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting over heterogeneous sources*.
- [82] Zekun Li, Wenhui Chen, Shiyang Li, Hong Wang, Jing Qian, and Xifeng Yan. 2022. *Controllable dialogue simulation with in-context learning*. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4330–4347, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- [83] Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. *Let’s verify step by step*.
- [84] Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. *TruthfulQA: Measuring how models mimic human falsehoods*. In *Proceedings of the 60th Annual*

- Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- [85] Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. [Generating with confidence: Uncertainty quantification for black-box large language models](#).
 - [86] Zachary Lipton, Xiuju 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, volume 32.
 - [87] Bing Liu and Ian Lane. 2017. Iterative policy learning in end-to-end trainable task-oriented neural dialog models. In 2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pages 482–489. IEEE.
 - [88] Bing Liu, Gokhan Tur, Dilek Hakkani-Tur, Pararth Shah, and Larry Heck. 2018. Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems. [arXiv preprint arXiv:1804.06512](#).
 - [89] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. [What makes good in-context examples for GPT-3?](#) In Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
 - [90] Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2023. [Trustworthy llms: a survey and guideline for evaluating large language models' alignment](#).
 - [91] Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, et al. 2023. Summary of chatgpt-related research and perspective towards the future of large language models. Meta-Radiology, page 100017.
 - [92] 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. [arXiv preprint arXiv:1907.11692](#).
 - [93] Yixin Liu, Alex Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir

- Radev. 2023. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4140–4170, Toronto, Canada. Association for Computational Linguistics.
- [94] Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, and Daphne Ippolito. 2024. A pretrainer’s guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 3245–3276, Mexico City, Mexico. Association for Computational Linguistics.
- [95] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback.
- [96] Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. 2021. Few-shot bot: Prompt-based learning for dialogue systems. arXiv preprint arXiv:2110.08118.
- [97] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.
- [98] Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models.
- [99] Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tur. 2020. Dialoglue: A natural language understanding benchmark for task-oriented dialogue.
- [100] Shikib Mehri and Maxine Eskenazi. 2021. Schema-guided paradigm for zero-shot dialog. arXiv preprint arXiv:2106.07056.
- [101] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In Proceedings of the 2018 Conference on Empirical Methods in Natural

- Language Processing, pages 2381–2391, Brussels, Belgium. Association for Computational Linguistics.
- [102] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen tau Yih, Pang Wei Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [Factscore: Fine-grained atomic evaluation of factual precision in long form text generation](#).
 - [103] Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [FActScore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 12076–12100, Singapore. Association for Computational Linguistics.
 - [104] Johannes EM Mosig, Shikib Mehri, and Thomas Kober. 2020. Star: A schema-guided dialog dataset for transfer learning. [arXiv preprint arXiv:2010.11853](#).
 - [105] Maxwell I. Nye, Michael Henry Tessler, Joshua B. Tenenbaum, and Brenden M. Lake. 2021. [Improving coherence and consistency in neural sequence models with dual-system, neuro-symbolic reasoning](#). In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pages 25192–25204.
 - [106] OpenAI. 2022. [large-scale generative pre-training model for conversation](#). [OpenAI Blog](#).
 - [107] OpenAI. 2023. [Gpt-4 technical report](#).
 - [108] OpenAI. 2024. [Hello gpt-4o](#). [OpenAI Blog](#).
 - [109] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In Advances in Neural Information Processing Systems, volume 35, pages 27730–27744. Curran Associates, Inc.
 - [110] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. Advances in Neural Information Processing Systems, 35:27730–27744.

- [111] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In NeurIPS.
- [112] Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. 2024. [Fine-tuning or retrieval? comparing knowledge injection in llms](#).
- [113] Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. [Medmcqa: A large-scale multi-subject multi-choice dataset for medical domain question answering](#). In Proceedings of the Conference on Health, Inference, and Learning, volume 174 of Proceedings of Machine Learning Research, pages 248–260. PMLR.
- [114] 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, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- [115] Baolin Peng, Michel Galley, Pengcheng He, Chris Brockett, Lars Liden, Elnaz Nouri, Zhou Yu, Bill Dolan, and Jianfeng Gao. 2022. [Godel: Large-scale pre-training for goal-directed dialog](#).
- [116] Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. [arXiv preprint arXiv:2302.12813](#).
- [117] Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayandeh, Lars Liden, and Jianfeng Gao. 2020. Soloist: Few-shot task-oriented dialog with a single pretrained auto-regressive model. [arXiv preprint arXiv:2005.05298](#).
- [118] 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](#), 9:807–824.
- [119] 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](#), 9:807–824.

- [120] Baolin Peng, Chunyuan Li, Zhu Zhang, Jinchao Li, Chenguang Zhu, and Jianfeng Gao. 2021. Synergy: Building task bots at scale using symbolic knowledge and machine teaching. [arXiv preprint arXiv:2110.11514](#).
- [121] Baolin Peng, Chunyuan Li, Zhu Zhang, Chenguang Zhu, Jinchao Li, and Jianfeng Gao. 2021. [RADDLE: an evaluation benchmark and analysis platform for robust task-oriented dialog systems](#). In [Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, \(Volume 1: Long Papers\), Virtual Event, August 1-6, 2021](#), pages 4418–4429. Association for Computational Linguistics.
- [122] Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, Kam-Fai Wong, and Shang-Yu Su. 2018. Deep dyna-q: Integrating planning for task-completion dialogue policy learning. [arXiv preprint arXiv:1801.06176](#).
- [123] Baolin Peng, Xiujun Li, Lihong Li, Jianfeng Gao, Asli Celikyilmaz, Sungjin Lee, and Kam-Fai Wong. 2017. Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning. [arXiv preprint arXiv:1704.03084](#).
- [124] Kun Qian and Zhou Yu. 2019. [Domain adaptive dialog generation via meta learning](#). In [Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics](#), pages 2639–2649, Florence, Italy. Association for Computational Linguistics.
- [125] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. [Is chatgpt a general-purpose natural language processing task solver?](#)
- [126] Libo Qin, Wenbo Pan, Qiguang Chen, Lizi Liao, Zhou Yu, Yue Zhang, Wanxiang Che, and Min Li. 2023. [End-to-end task-oriented dialogue: A survey of tasks, methods, and future directions](#).
- [127] 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):9.
- [128] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susanah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saftron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell,

- Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. 2022. [Scaling language models: Methods, analysis & insights from training gopher.](#)
- [129] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. [Direct preference optimization: Your language model is secretly a reward model.](#)
- [130] 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.](#) *Journal of Machine Learning Research*, 21(140):1–67.
- [131] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. [Exploring the limits of transfer learning with a unified text-to-text transformer.](#)
- [132] Janarthanan Rajendran, Jatin Ganhotra, and Lazaros C Polymenakos. 2019. Learning end-to-end goal-oriented dialog with maximal user task success and minimal human agent use. *Transactions of the Association for Computational Linguistics*, 7:375–386.
- [133] Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2019. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. [arXiv preprint arXiv:1909.05855](#).
- [134] Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. [Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset.](#) In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8689–8696. AAAI Press.

Bibliography

- [135] Jie Ren, Yao Zhao, Tu Vu, Peter J. Liu, and Balaji Lakshminarayanan. 2023. [Self-evaluation improves selective generation in large language models.](#)
- [136] Adam Roberts, Hyung Won Chung, Anselm Levskaya, Gaurav Mishra, James Bradbury, Daniel Andor, Sharan Narang, Brian Lester, Colin Gaffney, Afroz Mohiuddin, Curtis Hawthorne, Aitor Lewkowycz, Alex Salcianu, Marc van Zee, Jacob Austin, Sebastian Goodman, Livio Baldini Soares, Haitang Hu, Sasha Tsvyashchenko, Aakanksha Chowdhery, Jasmijn Bastings, Jannis Bulian, Xavier Garcia, Jianmo Ni, Andrew Chen, Kathleen Kenealy, Jonathan H. Clark, Stephan Lee, Dan Garrette, James Lee-Thorp, Colin Raffel, Noam Shazeer, Marvin Ritter, Maarten Bosma, Alexandre Passos, Jeremy Maitin-Shepard, Noah Fiedel, Mark Omernick, Brennan Saeta, Ryan Sepassi, Alexander Spiridonov, Joshua Newlan, and Andrea Gesmundo. 2022. [Scaling up models and data with t5x and seqio.](#)
- [137] William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. 2022. [Self-critiquing models for assisting human evaluators.](#)
- [138] John Schulman. 2023. [Reinforcement learning from human feedback: Progress and challenges.](#) Berkeley EECS.
- [139] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. [Proximal policy optimization algorithms.](#)
- [140] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. [Proximal policy optimization algorithms.](#)
- [141] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. [arXiv preprint arXiv:1508.07909](#).
- [142] Pararth Shah, Dilek Hakkani-Tür, and Larry Heck. 2016. Interactive reinforcement learning for task-oriented dialogue management. In NIPS 2016 Deep Learning for Action and Interaction Workshop, volume 11.
- [143] Pararth Shah, Dilek Hakkani-Tur, Bing Liu, and Gokhan Tur. 2018. Bootstrapping a neural conversational agent with dialogue self-play, crowdsourcing and on-line reinforcement learning. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers), pages 41–51.
- [144] Alex Sherstinsky. 2020. [Fundamentals of recurrent neural network \(rnn\) and long short-term memory \(lstm\) network.](#) Physica D: Nonlinear Phenomena, 404:132306.

- [145] Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. [Detecting pretraining data from large language models](#).
- [146] Swadheen Shukla, Lars Liden, Shahin Shayandeh, Eslam Kamal, Jinchao Li, Matt Mazzola, Thomas Park, Baolin Peng, and Jianfeng Gao. 2020. Conversation learner—a machine teaching tool for building dialog managers for task-oriented dialog systems. [arXiv preprint arXiv:2004.04305](#).
- [147] Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. Ai models collapse when trained on recursively generated data. [Nature](#), 631(8022):755–759.
- [148] David Silver and Richard S Sutton. 2025. Welcome to the era of experience. [Google AI](#).
- [149] Patrice Y Simard, Saleema Amershi, David M Chickering, Alicia Edelman Pelton, Soroush Ghorashi, Christopher Meek, Gonzalo Ramos, Jina Suh, Johan Verwey, Mo Wang, et al. 2017. Machine teaching: A new paradigm for building machine learning systems. [arXiv preprint arXiv:1707.06742](#).
- [150] Patrice Y. Simard, Saleema Amershi, David Maxwell Chickering, Alicia Edelman Pelton, Soroush Ghorashi, Christopher Meek, Gonzalo A. Ramos, Jina Suh, Johan Verwey, Mo Wang, and John Wernsing. 2017. [Machine teaching: A new paradigm for building machine learning systems](#). [CoRR](#), abs/1707.06742.
- [151] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Aguera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. 2023. [Towards expert-level medical question answering with large language models](#).
- [152] Robyn Speer and Joanna Lowry-Duda. 2017. [ConceptNet at SemEval-2017 task 2: Extending word embeddings with multilingual relational knowledge](#). In [Proceedings of the 11th International Workshop on Semantic Evaluation \(SemEval-2017\)](#), pages 85–89, Vancouver, Canada. Association for Computational Linguistics.
- [153] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta,

- and Adrià Garriga-Alonso et al. 2023. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#).
- [154] Sarah Sterz, Kevin Baum, Sebastian Biewer, Holger Hermanns, Anne Lauber-Rönsberg, Philip Meinel, and Markus Langer. 2024. [On the quest for effectiveness in human oversight: Interdisciplinary perspectives](#). In [Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency](#), FAccT ’24, page 2495–2507, New York, NY, USA. Association for Computing Machinery.
- [155] Pei-Hao Su. 2018. Reinforcement learning and reward estimation for dialogue policy optimisation. In [University of Cambridge](#).
- [156] Pei-Hao Su, Milica Gasic, Nikola Mrksic, Lina Rojas-Barahona, Stefan Ultes, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2016. On-line active reward learning for policy optimisation in spoken dialogue systems. [arXiv preprint arXiv:1605.07669](#).
- [157] Yixuan Su, Lei Shu, Elman Mansimov, Arshit Gupta, Deng Cai, Yi-An Lai, and Yi Zhang. 2022. [Multi-task pre-training for plug-and-play task-oriented dialogue system](#). [Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics \(ACL\)](#).
- [158] Haipeng Sun, Junwei Bao, Youzheng Wu, and Xiaodong He. 2022. Mars: Semantic-aware contrastive learning for end-to-end task-oriented dialog. [arXiv preprint arXiv:2210.08917](#).
- [159] Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianshi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. 2024. [Trustllm: Trustworthiness in large language models](#).
- [160] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt

- Keutzer, and Trevor Darrell. 2023. [Aligning large multimodal models with factually augmented rlhf](#).
- [161] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. [Sequence to sequence learning with neural networks](#).
- [162] Richard S. Sutton and Andrew G. Barto. 1998. [Reinforcement learning - an introduction](#). Adaptive computation and machine learning. MIT Press.
- [163] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In [Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 \(Long and Short Papers\)](#), pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- [164] Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. 2024. [A survey on self-evolution of large language models](#).
- [165] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. [Galactica: A large language model for science](#).
- [166] 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, De-hao 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, Tulsee 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](#).
- [167] Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D. Manning, and Chelsea Finn. 2023. [Fine-tuning language models for factuality](#).

- [168] Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailev, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- [169] S. M Towhidul Islam Tonmoy, S M Mehedi Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. [A comprehensive survey of hallucination mitigation techniques in large language models](#).
- [170] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *CoRR*, abs/2302.13971.
- [171] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharang Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#).
- [172] Bo-Hsiang Tseng, Yinpei Dai, Florian Kreyssig, and Bill Byrne. 2021. Transferable dialogue systems and user simulators. [arXiv preprint arXiv:2107.11904](#).
- [173] Alan M. Turing. 1990. Computing machinery and intelligence. In Margaret A. Boden, editor, *The Philosophy of Artificial Intelligence*, Oxford readings in philosophy, pages 40–66. Oxford University Press.
- [174] Neeraj Varshney, Wenlin Yao, Hongming Zhang, Jianshu Chen, and Dong Yu. 2023. [A stitch in time saves nine: Detecting and mitigating hallucinations of llms by validating low-confidence generation](#).

- [175] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- [176] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. [Attention is all you need](#).
- [177] Fanqi Wan, Xinting Huang, Leyang Cui, Xiaojun Quan, Wei Bi, and Shuming Shi. 2024. [Knowledge verification to nip hallucination in the bud](#).
- [178] Ben Wang and Aran Komatsuzaki. 2021. [Gpt-j-6b: A 6 billion parameter autoregressive language model](#). <https://github.com/kingoflolz/mesh-transformer-jax>.
- [179] Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, Yidong Wang, Linyi Yang, Jindong Wang, Xing Xie, Zheng Zhang, and Yue Zhang. 2023. [Survey on factuality in large language models: Knowledge, retrieval and domain-specificity](#).
- [180] Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, Binxin Jiao, Yue Zhang, and Xing Xie. 2023. [On the robustness of chatgpt: An adversarial and out-of-distribution perspective](#).
- [181] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharani Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations*.
- [182] 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](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.
- [183] 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](#).

- [184] Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. 2022. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109.
- [185] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. *Emergent abilities of large language models*.
- [186] Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, Milica Gašić, Lina M. Rojas-Barahona, Pei-Hao Su, Stefan Ultes, and Steve Young. 2017. *A network-based end-to-end trainable task-oriented dialogue system*. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 438–449, Valencia, Spain. Association for Computational Linguistics.
- [187] Jason D Williams and Lars Liden. 2017. Demonstration of interactive teaching for end-to-end dialog control with hybrid code networks. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 82–85.
- [188] Ronald J. Williams. 1992. *Simple statistical gradient-following algorithms for connectionist reinforcement learning*. *Mach. Learn.*, 8:229–256.
- [189] Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256.
- [190] Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierrick Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- [191] Kevin Wu, Eric Wu, and James Zou. 2024. *How faithful are rag models? quantifying the tug-of-war between rag and llms’ internal prior*.
- [192] Siye Wu, Jian Xie, Jiangjie Chen, Tinghui Zhu, Kai Zhang, and Yanghua Xiao. 2024. *How easily do irrelevant inputs skew the responses of large language models?*
- [193] Ting Wu, Xuefeng Li, and Pengfei Liu. 2024. *Progress or regress? self-improvement reversal in post-training*.

- [194] Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. 2024. [Agentgym: Evolving large language model-based agents across diverse environments](#).
- [195] Chong Xiang, Tong Wu, Zexuan Zhong, David Wagner, Danqi Chen, and Prateek Mittal. 2024. [Certifiably robust rag against retrieval corruption](#).
- [196] Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P. Lillicrap, Kenji Kawaguchi, and Michael Shieh. 2024. [Monte carlo tree search boosts reasoning via iterative preference learning](#).
- [197] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Dixin Jiang. 2023. [Wizardlm: Empowering large language models to follow complex instructions](#).
- [198] Fangzhi Xu, Qiushi Sun, Kanzhi Cheng, Jun Liu, Yu Qiao, and Zhiyong Wu. 2024. [Interactive evolution: A neural-symbolic self-training framework for large language models](#).
- [199] Hongshen Xu, Zichen Zhu, Situo Zhang, Da Ma, Shuai Fan, Lu Chen, and Kai Yu. 2024. [Rejection improves reliability: Training llms to refuse unknown questions using rl from knowledge feedback](#).
- [200] Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. 2024. [Faithful logical reasoning via symbolic chain-of-thought](#).
- [201] Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. 2023. [Alignment for honesty](#).
- [202] Steve Young, Milica Gašić, Blaise Thomson, and Jason D Williams. 2013. Pomdp-based statistical spoken dialog systems: A review. [Proceedings of the IEEE](#), 101(5):1160–1179.
- [203] Wenhao Yu, Zhihan Zhang, Zhenwen Liang, Meng Jiang, and Ashish Sabharwal. 2023. [Improving language models via plug-and-play retrieval feedback](#).
- [204] Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. [MultiWOZ 2.2 : A dialogue dataset with additional annotation corrections and state tracking baselines](#). In [Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI](#), pages 109–117, Online. Association for Computational Linguistics.

- [205] Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jiatong Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco Pavone, Yuqiang Li, Wanli Ouyang, and Dongzhan Zhou. 2024. [Llama-berry: Pairwise optimization for o1-like olympiad-level mathematical reasoning](#).
- [206] Hanlin Zhang, Jiani Huang, Ziyang Li, Mayur Naik, and Eric Xing. 2023. [Improved logical reasoning of language models via differentiable symbolic programming](#).
- [207] Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024. [R-tuning: Instructing large language models to say ‘I don’t know’](#). In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 7106–7132, Mexico City, Mexico. Association for Computational Linguistics.
- [208] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2024. [Instruction tuning for large language models: A survey](#).
- [209] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuhui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. [Opt: Open pre-trained transformer language models](#).
- [210] Tianjun Zhang, Shishir G. Patil, Naman Jain, Sheng Shen, Matei Zaharia, Ion Stoica, and Joseph E. Gonzalez. 2024. [Raft: Adapting language model to domain specific rag](#).
- [211] Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In International Conference on Learning Representations.
- [212] Xiaoying Zhang, Baolin Peng, Jianfeng Gao, and Helen Meng. 2022. Toward self-learning end-to-end task-oriented dialog systems. In Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue, pages 516–530.
- [213] Xiaoying Zhang, Baolin Peng, Kun Li, Jingyan Zhou, and Helen Meng. 2023. [SGP-TOD: Building task bots effortlessly via schema-guided LLM prompting](#). In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 13348–13369, Singapore. Association for Computational Linguistics.

-
- [214] Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024. [Self-alignment for factuality: Mitigating hallucinations in LLMs via self-evaluation](#). In [Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 1946–1965, Bangkok, Thailand. Association for Computational Linguistics.
 - [215] Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. 2024. [Self-alignment for factuality: Mitigating hallucinations in llms via self-evaluation](#).
 - [216] Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Yipeng Zhang, Haitao Mi, and Helen Meng. 2024. [Self-tuning: Instructing llms to effectively acquire new knowledge through self-teaching](#).
 - [217] Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. [Task-oriented dialog systems that consider multiple appropriate responses under the same context](#). In [The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020](#), pages 9604–9611. AAAI Press.
 - [218] Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. Task-oriented dialog systems that consider multiple appropriate responses under the same context. In [Proceedings of the AAAI Conference on Artificial Intelligence](#), volume 34, pages 9604–9611.
 - [219] Yue Zhang, Leyang Cui, Wei Bi, and Shuming Shi. 2023. [Alleviating hallucinations of large language models through induced hallucinations](#).
 - [220] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. [Siren’s song in the ai ocean: A survey on hallucination in large language models](#).
 - [221] Zheng Zhang, Ryuichi Takanobu, Qi Zhu, Minlie Huang, and Xiaoyan Zhu. 2020. [Recent advances and challenges in task-oriented dialog system](#).
 - [222] Baining Zhao, Ziyou Wang, Jianjie Fang, Chen Gao, Fanhang Man, Jinqiang Cui, Xin Wang, Xinlei Chen, Yong Li, and Wenwu Zhu. 2025. [Embodied-r: Collaborative framework for activating embodied spatial reasoning in foundation models via reinforcement learning](#).

- [223] Jeffrey Zhao, Yuan Cao, Raghav Gupta, Harrison Lee, Abhinav Rastogi, Mingqiu Wang, Hagen Soltau, Izhak Shafran, and Yonghui Wu. 2022. Anytod: A programmable task-oriented dialog system. [arXiv preprint arXiv:2212.09939](#).
- [224] Ruochen Zhao, Xingxuan Li, Shafiq Joty, Chengwei Qin, and Lidong Bing. 2023. [Verify-and-edit: A knowledge-enhanced chain-of-thought framework](#). In [Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics \(Volume 1: Long Papers\)](#), pages 5823–5840, Toronto, Canada. Association for Computational Linguistics.
- [225] Tiancheng Zhao and Maxine Eskenazi. 2016. [Towards end-to-end learning for dialog state tracking and management using deep reinforcement learning](#). In [Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue](#), pages 1–10, Los Angeles. Association for Computational Linguistics.
- [226] Tiancheng Zhao and Maxine Eskenazi. 2018. [Zero-shot dialog generation with cross-domain latent actions](#). In [Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue](#), pages 1–10, Melbourne, Australia. Association for Computational Linguistics.
- [227] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. [Calibrate before use: Improving few-shot performance of language models](#). In [Proceedings of the 38th International Conference on Machine Learning](#), volume 139 of [Proceedings of Machine Learning Research](#), pages 12697–12706. PMLR.
- [228] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivas Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. [Lima: Less is more for alignment](#).
- [229] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. [Instruction-following evaluation for large language models](#).