Workflow-Guided Response Generation for Task-Oriented Dialogue

Anonymous ACL submission

Abstract

Task-oriented dialogue (TOD) systems aim to achieve specific goals through interactive dialogue. Such tasks usually involve following specific workflows, i.e. executing a sequence of actions in a particular order. While prior work has focused on supervised learning methods to condition on past actions, they do not explicitly optimize for compliance to a desired workflow. In this paper, we propose a novel framework based on reinforcement learning (RL) to generate dialogue responses that are aligned with a given workflow. Our framework consists of ComplianceReward, a metric designed to evaluate how well a generated response executes the specified action, combined with an RL optimization process that utilizes an interactive sampling technique. We evaluate our approach on two TOD datasets, Action-Based Conversations Dataset (ABCD) (Chen et al., 2021a) and MultiWOZ 2.2 (Zang et al., 2020) on a range of automated and human evaluation metrics. Our findings indicate that our RL-based framework outperforms baselines and is effective at generating responses that both comply with the intended workflows while being expressed naturally and fluently.

1 Introduction

011

012

014

021

026

027

037

038

041

Task-oriented dialogue (TOD) focuses on creating conversational systems that assist users in attaining specific objectives. While prior TOD literature has extensively looked at predicting user intents and identifying relevant slots and values (Henderson et al., 2014; Wei et al., 2018; Budzianowski et al., 2018; Byrne et al., 2019; Rastogi et al., 2020; Shalyminov et al., 2020; Balaraman et al., 2021), real-world interactions often involve nuanced workflows and optimizing for such workflows remains underexplored (Chen et al., 2021a; Hattami et al., 2022; Raimondo et al., 2023). Consider a customer support interaction where agents must follow multistep procedures that adhere to company policies.



Figure 1: In this interaction, the customer requests assistance with an expired promo code. The agent must help the customer while following the steps in the agent guideline, consisting of a sequence of actions to be taken to resolve the issue. For example, offering to generate a new promo code without querying the system results in a non-workflow-compliant behavior.

For example, in Figure 1, a customer asks for help with an expired promotional code. A model that accounts for the user intent might respond reactively, offering to generate a new promo code. However, assisting the customer involves not only modeling their intent but also staying consistent with a workflow, in this case, the company policy. This involves the agent executing the necessary actions in the right order, such as pulling up account information and querying the system to make sure the customer qualifies for the promotion.

Many prior approaches in task-oriented dialogue (TOD), such as SimpleTOD (Hosseini-Asl et al., 2020) and PPTOD (Su et al., 2022), have employed supervised learning alongside utterance-level user intents and system dialogue acts (DAs) for system response generation. However, these frameworks lack explicit optimization for compliance in response generation, resulting in responses that may fail to execute the specified action. This prob-

086

880

097

100

101

103

104

105

106

107

108

109

110

062

063

064

lem arises because the response generators neither receive rewards nor penalties based on adherence to the specified actions. Additionally, there is a notable absence of a metric or model to quantitatively assess the degree of compliance, hindering the evaluation and training of response generators.

In this work, we tackle the problem of workflowcompliant response generation in TOD and propose an RL-based approach that addresses the limitations of existing systems. Our approach (COMPLIANCEOPT) employs RL with compliance scoring to construct training data for the Quark (Lu et al., 2022) framework. We evaluate our approach using the Action-Based Conversations Dataset (ABCD) (Chen et al., 2021a), a TOD dataset enriched with policy-based agent behavior constraints in the form of action sequences, and MultiWOZ 2.2 (Zang et al., 2020). Our experiments show that models integrating workflow information surpass baseline models, producing responses that adhere to policies while maintaining a natural and fluent tone. Furthermore, we observe that direct compliance optimization through RL can lead to additional enhancements in the workflow compliance levels of the dialogue system. We validate our results through automated metrics and human evaluations. Our contributions include:

- A reinforcement learning (RL)-based framework for training workflow-compliant response generators, based on an interactive sampling technique to optimize model behavior over multiple dialogue exchanges.¹
- A new compliance metric based on a reward model validated against human evaluations.
- Evaluation on both automated and human evaluation metrics showing that our models, enhanced with workflow information and direct compliance optimization through RL, consistently outperform baselines.

2 Related Work

Task-oriented Dialogue. Recently, there has been an increase in TOD tasks and datasets (Budzianowski et al., 2018; Byrne et al., 2019; Wei et al., 2018; Rastogi et al., 2020), indicating a growing emphasis on advancing natural language processing techniques for practical applications. These datasets encompass diverse domains and enable researchers to tackle a wide spectrum of real-world challenges. However, previous benchmarks have 111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

Workflow Compliance. The problem of workflow compliance is closely related to, but distinct from, dialogue policy management. The primary objective of dialogue policy management is to predict the optimal dialogue action based on the current conversation state (Takanobu et al., 2019; He et al., 2022). In this context, dialogue actions represent intentions or decisions that are isolated to a single user query, such as "book a flight" or "find a nearby restaurant." In contrast, workflow compliance adopts a more holistic approach, considering the sequential workflow from the larger context of the conversation to define success. For example, offering a new promo code is only valid after a system check has been executed first (Figure 1). It emphasizes the fact that user interactions are not isolated actions but rather part of a continuous process with multiple steps. Raimondo et al. (2023) expands upon Chen et al. (2021a)'s work to show that models augmented with workflow-specific information such as workflow names or action plans can boost the generalizability of action prediction models, but does not consider the problem of generating workflow-compliant responses, which is a focus of our work.

SimpleTOD (Hosseini-Asl et al., 2020) is similar to the baselines in our work as both methods involve training an end-to-end model with interleaved actions and utterances as inputs. PPTOD (Su et al., 2022) also uses interleaved actions, but they use a more complex training pipeline that involves multi-task pretraining. While this can help improve performance, such baselines optimize for an objective that is different from our goal of increasing the compliance quality of generated responses.

Reinforcement Learning. RL has been successfully used to improve TOD systems (Pietquin et al., 2011; Gašić et al., 2013; Fatemi et al., 2016; Lewis et al., 2017; Singh et al., 2002). One application is training dialogue managers that maintain dialogue state transitions (Rieser and Lemon, 2011). Another is to use RL in conjunction with supervised learning to improve the quality of language generation, such as in (Lewis et al., 2017). This line of research applies similar techniques used in RL for general-domain dialogue generation, such

predominantly focused on evaluating only some aspects of TOD systems, such as intent recognition and slot filling, with limited focus on aspects like workflow compliance (Chen et al., 2021a).

¹We open-source our code at github.com/ANON.

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

186

187

188

190

191

194

195

196

198

201

203

207

208

210

211

as interleaving supervised learning and RL, offline and online RL, policy gradients, and Q-learning (Li et al., 2016b; Jang et al., 2022; Snell et al., 2023; Sodhi et al., 2023). Our work adopts a similar strategy of supervised learning followed by RL but introduces an interactive sampling step.

3 Problem Formulation

3.1 Workflow-Compliant Response Generation as an MDP

We formalize the problem of workflow-compliant response generation as a Markov Decision Process (MDP). Given a dataset of context-response pairs $\{x^i, y^i\}_{i=1}^N$, where context x is the conversation history, and response $y = \{y_1, \ldots, y_T\}$ is a target sequence of tokens.

Additionally, each dialogue is associated with a domain d representing the task (e.g., troubleshootsite, subscription-inquiry). Every domain has an associated set of workflow G_d , which is a natural language description of the steps the system must follow to assist the user, as well as a sequence of actions W_d , which represents a flat action list or workflow sequence based on the guidelines. Fully compliant dialogues do not necessarily follow the full sequence W_d and may instead only include a subset of these actions since the guideline includes conditional branching. For example, in Figure 1, the step 3 is dependent on the result of step 2.

Each data instance, denoted as (x, y, G_d) , can be viewed as an episode within an MDP, which we define as follows:

- States, st ∈ S is the context x, workflow Gd, and the partially generated sequence of tokens up to and including time step t, which we denote as ŷ<t := ŷ1,..., ŷt.
- Actions, at ∈ A are the set of possible next tokens ŷt+1 from the vocabulary V.
- Transition function, \$\mathcal{T}(s_{t+1}|s_t, a_t)\$ is deterministic, as each state-action pair (\$\hat{y}_{< t}, \$\hat{y}_{t+1}\$) leads to a unique state \$\hat{y}_{< t+1}\$ for the next step.
- Rewards, r_t : S × A → [0, 1] provide a measure of how well the generated response ŷ executes the provided workflow G_d. It is a terminal reward. Since workflow compliance can be computed only after multiple exchanges, the reward is computed using block evaluation.
- Horizon, T represents the time span of each episode, concluding either when the current time step t exceeds T or when an end-of-

sentence (EOS) token is generated.

The goal is to learn a policy $\pi : s_t \to a_t$ maximizing *return*, *i.e.* the cumulative reward over an episode $\mathbb{E}_{\pi} \sum_{t=0}^{T} \gamma^t r_t$. We assume undiscounted cumulative rewards, *i.e.* $\gamma = 1$.

Block Evaluation. One of our key observations is that compliance is not easily captured in a single dialogue response. For example, in a customer service use case, an agent may need to verify the identity of the user before proceeding to issue resolution. To successfully comply with the next workflow action e.g. verify-identity, the agent needs to take several steps. To better model and leverage this insight, we consider "blocks" of user and agent utterances when evaluating and optimizing for compliance. Blocks refer to the sequence of user and system utterances that occur between two action executions. We define an interaction "block" *b* as a list of user and system utterances between consecutive action executions by the system.

4 Approach

We introduce COMPLIANCEOPT, which directly optimizes compliance with the specified workflow. We define compliance as the extent to which the generated system utterances adhere to the prescribed workflow action at turn t.

Algorithm 1 shows our overall training procedure, which is adapted from the Quark algorithm (Lu et al., 2022). The Quark framework is similar to Decision Transformer in that it treats RL as a sequence modeling problem (Chen et al., 2021b). After interactively sampled (Figure 2-(i)) generations are scored (Figure 2-(ii)), the rewards are quantized to produce reward tokens r_k , which are then used to condition the generations during training (Figure 2-(iii)).

4.1 Interactive Sampling

Diverging from the Quark method, we implement an interactive sampling step, using two distinct models, a system model, and a user simulator. This is because achieving workflow compliance often requires multiple dialogue turns between the participants. Consider a customer service agent who needs to gather a user's name, email, and order id to validate their purchase. This is a multi-turn process where the system needs to gather information over multiple dialogue turns of questions and answers.

We warm-start with a system model trained with standard autoregressive training. The user simula-

235

236

237

238

239

240

241

242

243

244

245

212

213

214

215

216

217

218

219

220

221

222

223

224

246 247

248

249

250

251

252

253

254

255

256

257

258

259



Figure 2: Approach Overview. RL optimizes the model towards better workflow compliance. Interaction-score pairs are processed into RL data in the Quark framework.

tor remains fixed during Quark training and is only used for the interactive sampling procedure. Given a dialogue context $c_0 = [u_0]$, the system model first samples an utterance, which is then concatenated to c_0 , forming $c_1 = [u_0, s_1]$. Then, conditioned on c_1 , the user simulator samples a user turn, forming another context $c_2 = [u_0, s_1, u_2]$. This process is repeated M times, which is a hyperparameter. We denote the generated user and system utterances block as b, which are fed alongside planned workflow actions as inputs to the ComplianceReward.

261

262

263

273

274

276

281

289

291

297

Our interactive sampling technique is independent of Ouark and can be used as a sampling approach for other RL methods, such as proximal policy optimization (PPO) (Schulman et al., 2017). Previous studies, such as those by Zhang et al. (2022), have employed similar multi-utterance sampling techniques to simulate dialogue interactions. However, our approach diverges significantly from these methods. Instead of using a database of dialogue logs between real users and bots to construct RL data, our framework fully simulates responses for both users and systems, recognizing that completing an action typically requires multiple interactions. This introduces a more complex challenge for reinforcement learning (RL) optimization, as it necessitates the inclusion of simulated user responses within the block of utterances.

Compliance Scoring Model 4.2

To quantify compliance and use it as a reward for RL, we developed the ComplianceReward, which measures the alignment between the generated system utterances and the prescribed workflow action.

Reward modeling. We train the ComplianceReward using the reward modeling loss for ranking two responses (Ouyang et al., 2022).

$$l(\theta) = -\sum_{(p,b_w,b_l)\sim D} \log\left(\sigma\left(r_\theta(p,b_w) - r_\theta(p,b_l)\right)\right) \quad (1)$$

Algorithm 1 COMPLIANCEOPT RL Training

- **Input:** Initial Policy l_0 , User Simulator μ , Dialogue Contexts C, reward $r(\cdot)$, KL weight β , number of quantiles K, number of interactions M, number of train iterations N
- 1: Make a copy l_{θ} of initial policy l_0 .
- 2: for iteration = $1, 2, \cdots, N$ do
- for $c_i \in C$ do 3:
- 4: Do **interactive_sample** $(l_{\theta}, \mu, M, c_i)$ to obtain
- b_i . ▷ Interactive Sampling 5:
 - Add $c_i, s_i, r(c_i, s_i)$ into data pool \mathcal{D} ▷ Scoring
- 6: end for
- 7: $\mathcal{D}_i \leftarrow \mathsf{quantize}$
- 8: for step = $1, 2, \cdots, M$ do
- 9: Draw a batch of data (c_i, b_i, r_{k_i}) from quantized data pool \mathcal{D}_i
- 10: Compute the objective in Eqn 2 and update policy θ with gradient descent ⊳ Update 11: end for

12: end for

r(p, b) represents the scalar output generated by ComplianceReward (parametrized by θ), given the planned workflow action p and the generated block b. The term b_w denotes the favored choice among the pair of responses, b_w and b_l , in the comparison dataset D.

298

299

300

301

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

The generated block b can include multiple utterances by both the user and the system. We found that there are advantages in excluding the dialogue context and presenting only b to the model due to the following reasons: (1) The model can focus more effectively on evaluating the text itself rather than being distracted by the typically longer context, (2) Constructing negative instances (b_l) becomes straightforward by replacing the workflow of the positive instance b_w with an alternative.

Data Collection. We first segment each conversation into multiple blocks s, comprising contiguous utterances that are annotated with the same workflow step p. By pairing utterances from different segments with different workflow annotation p_l such that $p \neq p_l$, we generate (p, b_w, b_l) triplets.

322

324

325

330

332

336

338

341

345

347

361

365

4.3 Compliance Optimization

Finally, the system model is updated according to a combination of the standard LM loss and a KL divergence loss between the updated model and the reference model, shown in Equation 2.

$$\max_{\theta} \mathbb{E}_{k \sim \mathcal{U}(1,K)} \mathbb{E}_{(c,b) \sim \mathcal{D}^{k}} [\log l_{\theta}(b|c,r_{k}) -\beta \sum_{t=1}^{T} \mathrm{KL} \left(l_{0}(\cdot|b_{< t},c) \| l_{\theta}(\cdot|b_{< t},c,r_{k}) \right)]$$

$$(2)$$

5 Experimental Setup

Datasets. We evaluate our approach using two TOD datasets, where users aim to accomplish specific tasks through dialogue. Each dataset consists of conversations between two speakers, a system or agent, and a user or customer.

- Action Based Conversations Dataset (ABCD) (Chen et al., 2021a): Contains ~10,000 dialogues between customers and agents and spans 55 intents. The agents have explicit workflows they need to follow according to company guidelines, making it an ideal dataset to evaluate compliance requirements.
 - MultiWOZ 2.2 (Zang et al., 2020): Contains over ~10,000 dialogues spanning multiple domains. We designate annotated user intents as workflow actions to be predicted and include agent dialogue acts (DAs) in the context.

Evaluation. We evaluate the different approaches and baselines on a variety of metrics:

- **LLM compliance:** We automatically evaluated compliance using an LLM (prompt in Appendix A). We used a categorical labeling scheme involving two levels: 0 = 'not compliant,' and 1 = 'fully compliant.'
- Human compliance: For human evaluation, we randomly selected 100 generated outputs from each model (guideline in Appendix A). We used binary labeling (0, 1) for *compliance* and had three annotators rate each example. The annotators had access to the complete policy document containing guidelines for all workflow actions.
- Human coherence: Annotators were asked to also rate each of the same 100 examples on *coherence*, represented as a binary (0,1) label.
- Semantic Similarity: We measure the similarity between generated responses and the corresponding human-annotated ground truth using commonly-used similarity measures such as

BLEU, Meteor, BLEURT, and BERTScore (Papineni et al., 2002; Banerjee and Lavie, 2005; Sellam et al., 2020; Zhang et al., 2020). We report the "Block" version for each, computed by taking the max between each prediction and target utterance pair over all targets and taking the average over predictions.

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

384

387

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

- **Diversity**: We measure generated response diversity using the dist-3 metric (Li et al., 2016a).
- Workflow Accuracy: For the ACTIONPLAN and COMPLIANCEOPT models that predict the next workflow action, we report the exact match accuracy of the predicted action against ground truth.

Methods & Baselines.

- NOACTION: A simple model that only sees user and system utterances without access to completed actions or next workflow steps as done in Sodhi et al. (2023).
- ACTIONAWARE: Action executions are interleaved in the input alongside utterances, allowing the model to understand the history of completed workflow actions in the dialogue context. The model may implicitly learn the relationship between workflow policies and agent utterances, enabling the generation of more contextually relevant responses. This approach applies supervised learning on dialogue context augmented with prior and future actions and is similar to the SimpleTOD model (Hosseini-Asl et al., 2020).
- ACTIONPLAN: The ACTIONPLAN model goes beyond ACTIONAWARE by explicitly modeling future compliance to workflow policy guidelines. It introduces the concept of a "planned" workflow action, representing the next action that must be completed based on the policy. This planned action is incorporated into the dialogue context, and the model generates responses that align with the intended workflow. This approach treats the planned future workflow action as a latent variable in the generation process, resulting in better workflow compliance in responses.
- **GUIDELINE**: Instead of relying on completed actions or predicted future workflow actions, the GUIDELINE approach conditions on a fixed "standard" sequence of actions, referred to as the guideline in (Chen et al., 2021a).
- LLM-PROMPTING: We use prompting and in-context learning with large language models (LLMs) to explore the option of using nat-

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

467

468

469

ural language policy guidelines as a source of 417 workflow information. Similar to prior work 418 (Zhang et al., 2023), our LLM prompt consists 419 of instructions that describe the task and task-420 related text that consists of guidelines, exam-421 ple conversations, and the dialogue context C_t . 422 Our LLM-prompting method assumes ORACLE 423 next workflow and generates corresponding re-424 sponses. We include our prompt in Appendix A. 425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461 462

463

464

465

466

• **PREDICTED/ORACLE** variants: At test time, our ACTIONPLAN and COMPLIANCEOPT models predict the next workflow action and condition system response generation on this action (referred to as PREDICTED). ORACLE models are supplied with ground truth next workflow actions, to gauge the upper bound of performance. Finally, **ACTIONPLANALL FUTURE ORACLE** uses all future remaining workflow steps annotated in the ground truth data.

Training. Our dialogue system models and user simulators are both initialized with pretrained DistilGPT2 (Sanh et al., 2019), which is a condensed variant of GPT-2 (Radford et al., 2019). Tokenization of inputs to the system and user models use pre-trained BPE codes (Sennrich et al., 2016). For the ComplianceReward model, we start with a pretrained RoBERTa model, with its associated BPE tokenizer. (Liu et al., 2019). Training procedure and hyperparameters are included in Appendix A.

6 Results and Analyses

6.1 Overall Results

- Workflow-awareness consistently improves performance: Models incorporating workflow information show higher compliance over the NOACTION baseline (Table 5, Figure 3).
- Direct compliance optimization leads to peak system compliance: Our investigation reveals that COMPLIANCEOPT, which utilizes reinforcement learning to optimize compliance scores, outperforms models trained with teacher-forcing. This approach not only successfully optimizes response for the ComplianceReward model (Table 5), but also leads to high performance in human compliance and LLM-based evaluations (Figure 3).
- Human evaluation validates automated metrics: Human evaluators corroborate the results obtained from automated evaluations, confirming that workflow-aware models consistently outperform baselines. Remarkably, RL optimiza-

tion achieves higher compliance without compromising coherence (Figure 3).

- Consistent performance across datasets: The improved performance of workflow-aware models, particularly ACTIONPLAN and COMPLI-ANCEOPT, extends beyond the primary ABCD dataset. These findings hold even when validated on more general task-oriented dialogue datasets, such as MultiWoz (Table 3).
- *Ablation studies:* We also conduct ablation studies to investigate the effectiveness of explicitly predicting workflow actions compared to directly following standardized workflow guidelines. We also explore the impact of predicting and conditioning on future action sequences as opposed to single actions (Tables 5,2).
- Workflow-aware models maintain high performance on dialogue metrics: While our approach generates responses directly without predicting intermediate slot values, we include Action State Tracking (AST) dialogue metrics in Appendix A.2 by extracting these values from the response. We show that we can perform comparably to prior work on these metrics.

6.2 Compliance

First, we evaluate all approaches on how well they generate compliant responses. We conduct human and LLM evaluations. We also report ComplianceReward on the test dataset and show that it correlates with both LLM and human scores.

LLM evaluation: We first evaluate compliance using the LLM evaluation prompt. The results, shown on the leftmost plot in Figure 3, indicate that ground truth responses achieve the highest compliance and workflow-aware models outperform baselines. COMPLIANCEOPT achieves higher compliance scores compared with ACTIONPLAN, which highlights the benefit of directly optimizing for compliance. Also, the ORACLE model variants perform better than their PREDICTED counterparts, which means that we can generate more compliant responses when we have access the true next workflow action.

	ACTIONAWARE	ACTIONPLAN	COMPLIANCEOPT	Ground Truth
	ACTIONAWARE	PREDICTED	PREDICTED	Ground Truth
Avg Compliance	0.31	0.29	0.35	0.46
At least 1 Compliant (↑)	0.41	0.53	0.47	0.57
All Compliant (†)	0.2	0.19	0.24	0.38
Fleiss Kappa	0.66	0.65	0.67	0.73

Table 1: Human Annotation Metrics & Inter-annotatorAgreement (full table in Appendix)



Figure 3: Left: Sample Simulated Interaction between Agent Models and User Simulator. Right: Evaluation results with human annotators and LLM. We report the average score received for each model.

Human evaluation: Next, we evaluate compli-510 ance using human evaluation. In the middle plot in 511 Figure 3, we show the average compliance score 512 across all annotators for each approach. Results 513 show that COMPLIANCEOPT models are more 514 compliant than their ACTIONPLAN counterparts, in 515 both PREDICTED and ORACLE variants. Addition-516 ally, we observe similar trends for both human and 517 LLM compliance judgments, with ground truth re-518 sponses receiving the highest scores, trailed by the 519 ORACLE models, then by the PREDICTED models. 520

> Table 1 shows further breakdown of the humanannotated scores. We compute the percentage of examples that received at least 1 compliant score and the percentage that received all compliant scores. In this analysis, we observe a similar trend, with ground truth performing best, followed by COM-PLIANCEOPT, which outperforms ACTIONPLAN on "all compliant" and "at least 1 non-compliant". We used Fleiss' Kappa for assessing annotator agreement. We find that the annotators are in "substantial agreement", showing that human compliance judgment is a reliable metric for compliance evaluation (Landis and Koch, 1977).

ComplianceReward While LLM and human evaluation are the most reliable way to evaluate compliance, we also report ComplianceReward, our training reward signal, on the test dataset. We compute the reward values across different approaches (Table 5). We find the compliance reward to be positively correlated with both human (Table 6) and LLM evaluation (Table 7). More details are included in Appendix A.2.

6.3 Coherence

521

522

524

527

535

537

539

540

541

542

543

545

In addition to compliance, we also evaluated whether generated responses were coherent. As

shown in the rightmost plot in Figure 3, ground truth responses have the highest coherence scores, followed closely by both variants of COMPLI-ANCEOPT. This indicates that COMPLIANCEOPT is able to achieve higher compliance scores while also maintaining high fluency and coherence.

6.4 Semantic Similarity and Diversity

Model	Block	Block	Block	Block	dist 3	Workflow
WIGHEI	BertScore	BLEURT	METEOR	BLEU	uist-5	Accuracy
	Ba	aselines &	Ablations			
NOACTION	0.8577	0.2286	0.0549	0.4481	0.7738	N/A
GUIDELINE	0.8670	0 2763	0.0670	0 /028	0 7536	N/A
ORACLE	0.0079	0.2705	0.0079	0.4920	0.7550	11/1
LLM-PROMPTING	0.8676	0 3033	0.0609	0 5493	0 7013	N/A
ORACLE	0.0070	0.5755	0.0007	0.5475	0.7015	10/11
ACTIONPLAN						
ALL FUTURE	0.8560	0.2498	0.0535	0.4606	0.7479	N/A
ORACLE						
		Proposed N	Aethods			
ACTIONAWARE	0.8642	0.2703	0.0726	0.4745	0.7661	N/A
ACTIONPLAN	0.8685	0 2050	0.0808	0 4951	0 7707	0 7011
PREDICTED	0.0005	0.2757	0.0000	0.4751	0.7707	0.7011
COMPLIANCEOPT	0 8740	0 2964	0 0924	0 5075	0.6553	0.6821
Predicted	0.0740	0.2704	0.0724	0.5075	0.0555	0.0021
ACTIONPLAN	0.8683	0 3081	0.0881	0 5021	0 7683	N/A
ORACLE	0.8085	0.5081	0.0881	0.3021	0.7005	19/74
COMPLIANCEOPT	0.8745	0 3312	0 1156	0 5287	0.6501	N/A
ORACLE	0.0743	0.5512	0.1150	0.5207	0.0391	11/1
Ground Truth	N/A	N/A	N/A	N/A	0.7738	N/A

Table 2: Semantic Similarity and Diversity Results

ACTIONPLAN and COMPLIANCEOPT achieve the highest semantic similarity scores when compared with ground truth-compliant responses (Table 2). This result indicates that adding future planned actions can lead to more contextually relevant and compliant system responses. Moreover, ACTIONPLAN and COMPLIANCEOPT ORACLE models outperform their PREDICTED counterparts, which suggests that using the *true* next workflow action results in responses more aligned with the human-annotated compliant behavior.

In addition, the lower dist-3 scores obtained by the COMPLIANCEOPT models, regardless of whether they are in the PREDICTED or ORACLE 546

547

566

553

554

configuration, suggest that these models produce 567 568 less diverse responses. One explanation is that the COMPLIANCEOPT models, as a result of the RL-569 based optimization, learn to focus on generating a narrower range of utterances that are compliant given the context. Since the ground truth responses achieve both higher dist-3 and compliance rewards, 573 this effect seems unique to the RL optimization. 574

6.5 Ablations

576

577

580

582

584

587

590

591

596

597

599

611 612

613

614

616

Effect of including all future actions. Since including future workflow actions results in more compliant responses, we explore if adding all future workflow actions would result in even more compliant behavior (ACTIONPLAN ALL FUTURE ORACLE vs. ACTIONPLAN ORACLE). Table 5 shows that including all future actions can hurt performance, likely because too much future information leads to noise and model confusion. In contrast, simply focusing on the next workflow action leads to compliant localized interactions ("blocks").

Training a model with standardized workflows.

We consider the effect of conditioning on standardized workflows, without dynamically including the next workflow actions in the context. As shown in Table 5, the GUIDELINE ORACLE model performs better than the baseline but worse than workflowaware models because it does not dynamically generate contextually relevant workflow actions and responses. This reinforces the importance of dynamic workflow prediction, which captures the inherent uncertainty in dialogues.

Few-shot LLM prompting with workflow guidelines. The final model variant we considered was directly using an LLM to predict the next workflow action, instead of fine-tuning a separate model. The LLM-PROMPTING ORACLE model achieves the second-highest compliance reward after the COMPLIANCEOPT ORACLE. We see that the COMPLIANCEOPT model, explicitly trained 605 to optimize compliance can outperform or match the LLM with orders of magnitude more parameters (gpt-3.5-turbo). The high text similarity 608 scores achieved by the LLM-PROMPTING ORA-CLE, often outperforming even the best-performing ACTIONPLAN and COMPLIANCEOPT models in terms of metrics like BLEURT and BLEU, validate the value of using guided prompts to improve response compliance. We note that RL optimization of LLMs requires much larger computational resources and remains an interesting future work.

Model	Compliance Score	BertScore	BLEURT	METEOR	BLEU	dist-3
		Baselin	e			
NOACTION	0.7446	0.4711	0.2476	0.0959	0.0108	0.4366
	Р	roposed M	ethods			
ACTIONAWARE	0.8451	0.8575	0.4001	0.1959	0.0252	0.8086
ACTIONPLAN PREDICTED	0.8463	0.8496	0.3928	0.1936	0.0249	0.8027
COMPLIANCEOPT PREDICTED	0.8853	0.8616	0.4310	0.1917	0.0265	0.8267
ACTIONPLAN ORACLE	0.8573	0.8506	0.3897	0.1900	0.0242	0.7962
COMPLIANCEOPT ORACLE	0.9153	0.8622	0.4271	0.1951	0.0278	0.8137
Ground Truth	0.8946	N/A	N/A	N/A	N/A	0.8237

Table 3: Automated Evaluation Results on MultiWOZ 2.2. PREDICTED variants of ACTIONPLAN and COM-PLIANCEOPT achieved 69% and 75% workflow accuracy respectively.

MultiWOZ Experiment Results 6.6

In our MultiWOZ experiments, we find consistent support for our approach. Workflow-aware models, particularly ACTIONPLAN and COMPLIANCEOPT, outperform NOACTION and ACTIONAWARE in both PREDICTED and ORACLE settings, showcasing their capacity to generate compliant and contextually relevant responses.

However, there are several differences when compared to the ABCD experiments. MultiWOZ introduces increased response diversity, especially noticeable in the COMPLIANCEOPT models, a departure from the ABCD dataset's behavior. Moreover, workflow-aware models benefit significantly from action annotation, as seen in the NOACTION versus ACTIONAWARE comparison. We conjecture that these disparities may be attributed to differences in action annotation and the nature of actions, which are typically resolved in a single interaction in MultiWOZ, in contrast to the more intricate workflows in the ABCD dataset.

7 Conclusion

In this paper, we propose the problem of workflowguided response generation and introduce a novel RL-based framework to train workflow-compliant models for task-oriented dialogue. By integrating workflow information during training and directly optimizing for compliance, our approach improves upon baseline models and generates responses that are both workflow-compliant and linguistically natural. We evaluate our models on both ABCD and MultiWoz datasets and show empirical improvements in automated and human evaluation metrics. 618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

8 Limitation

650

673

674

675

681

682

683

684

687

697

702

This paper introduced a novel RL-based framework that generates workflow-guided responses for taskoriented dialogue. While this is promising, sev-653 eral limitations warrant discussion. Higher fidelity user simulator models: Fully simulating realworld customer service interactions would provide a more comprehensive evaluation. For instance, our user model was intentionally kept simple to facilitate the development and testing. Using more sophisticated models that incorporate diverse user behaviors can potentially help with better general-661 ization. Using the same models for users and agents can also increase efficiency. Workflow requirements: We consider dialogue settings where workflow information is available, e.g. policy guidelines in customer-service interactions, and indeed more useful to follow compared to slot-value objectives. However, for datasets that do not have explicit workflows, one would have to proxy workflows, e.g. in MultiWoz we had to create workflows 670 671 from from user intent and system act annotations. This can limit the applicability of our method.

References

- Vevake Balaraman, Seyedmostafa Sheikhalishahi, and Bernardo Magnini. 2021. Recent neural methods on dialogue state tracking for task-oriented dialogue systems: A survey. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 239–251, Singapore and Online. Association for Computational Linguistics.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gašić. 2018. MultiWOZ - a largescale multi-domain Wizard-of-Oz dataset for taskoriented dialogue modelling. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 5016–5026, Brussels, Belgium. Association for Computational Linguistics.
- Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Ben Goodrich, Daniel Duckworth, Semih Yavuz, Amit Dubey, Kyu-Young Kim, and Andy Cedilnik. 2019. Taskmaster-1: Toward a realistic and diverse dialog dataset. In *Proceedings of the 2019 Conference on Empirical Meth-*

ods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4516–4525, Hong Kong, China. Association for Computational Linguistics. 703

704

705

706

707

708

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

- Derek Chen, Howard Chen, Yi Yang, Alex Lin, and Zhou Yu. 2021a. Action-based conversations dataset: A corpus for building more in-depth task-oriented dialogue systems. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021*, pages 3002–3017, Online. Association for Computational Linguistics.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. 2021b. Decision transformer: Reinforcement learning via sequence modeling. *CoRR*, abs/2106.01345.
- Mehdi Fatemi, Layla El Asri, Hannes Schulz, Jing He, and Kaheer Suleman. 2016. Policy networks with two-stage training for dialogue systems. In *Proceedings of the 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 101– 110, Los Angeles. Association for Computational Linguistics.
- Milica Gašić, Catherine Breslin, Matthew Henderson, Dongho Kim, Martin Szummer, Blaise Thomson, Pirros Tsiakoulis, and Steve Young. 2013. POMDPbased dialogue manager adaptation to extended domains. In *Proceedings of the SIGDIAL 2013 Conference*, pages 214–222, Metz, France. Association for Computational Linguistics.
- Amine El Hattami, Stefania Raimondo, Issam Hadj Laradji, David Vázquez, Pau Rodríguez López, and Christopher Joseph Pal. 2022. Workflow discovery from dialogues in the low data regime. *Trans. Mach. Learn. Res.*, 2023.
- 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 semisupervised learning and explicit policy injection. *Proceedings of the AAAI Conference on Artificial Intelligence.*
- Matthew Henderson, Blaise Thomson, and Jason D. Williams. 2014. The second dialog state tracking challenge. In *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, pages 263–272, Philadelphia, PA, U.S.A. Association for Computational Linguistics.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. In *Advances in Neural Information Processing Systems*, volume 33, pages 20179–20191. Curran Associates, Inc.

- 761 762 767 772 773 774 775 776 777 778 779 782 790 791 796 797 801
- 806 807

814

815 816

- Youngsoo Jang, Jongmin Lee, and Kee-Eung Kim. 2022. Gpt-critic: Offline reinforcement learning for end-toend task-oriented dialogue systems. In International Conference on Learning Representations (ICLR).
- J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. Biometrics, 33(1):159-174.
- Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021. Dialogue state tracking with a language model using schema-driven prompting. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 4937-4949, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Mike Lewis, Denis Yarats, Yann Dauphin, Devi Parikh, and Dhruv Batra. 2017. Deal or no deal? end-toend learning of negotiation dialogues. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 2443–2453, Copenhagen, Denmark. Association for Computational Linguistics.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. A diversity-promoting objective function for neural conversation models. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 110–119, San Diego, California. Association for Computational Linguistics.
- Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. 2016b. Deep reinforcement learning for dialogue generation. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1192-1202, Austin, Texas. Association for Computational Linguistics.
- 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. CoRR, abs/1907.11692.
- Ximing Lu, Sean Welleck, Jack Hessel, Liwei Jiang, Lianhui Qin, Peter West, Prithviraj Ammanabrolu, and Yejin Choi. 2022. QUARK: controllable text generation with reinforced unlearning. In NeurIPS.
- 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.
- 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.

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

- Olivier Pietquin, Matthieu Geist, Senthilkumar Chandramohan, and Hervé Frezza-Buet. 2011. Sampleefficient batch reinforcement learning for dialogue management optimization. ACM Trans. Speech Lang. *Process.*, 7(3).
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Stefania Raimondo, Christopher Pal, Xiaotian Liu, David Vazquez, and Hector Palacios. 2023. Improving generalization in task-oriented dialogues with workflows and action plans.
- 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.
- Verena Rieser and Oliver Lemon. 2011. Reinforcement Learning for Adaptive Dialogue Systems: a datadriven methodology for Dialogue Management and Natural Language Generation.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. In NeurIPS EMC2 Workshop.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. CoRR, abs/1707.06347.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7881–7892, Online. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Igor Shalyminov, Alessandro Sordoni, Adam Atkinson, and Hannes Schulz. 2020. Fast domain adaptation for goal-oriented dialogue using a hybrid generativeretrieval transformer. In 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP).

Satinder Singh, Diane Litman, Michael Kearns, and Marilyn Walker. 2002. Optimizing dialogue management with reinforcement learning: Experiments with the njfun system. J. Artif. Int. Res., 16(1):105–133.

874

875

895

897

898

899

900 901

902

903

904 905

906

907

908

909

910

911

912

913

914

915

917

919

920

921

922

925

927

928

929

930

- Charlie Snell, Ilya Kostrikov, Yi Su, Sherry Yang, and Sergey Levine. 2023. Offline RL for natural language generation with implicit language Q learning. In The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net.
- Paloma Sodhi, Felix Wu, Ethan R. Elenberg, Kilian Q
 Weinberger, and Ryan Mcdonald. 2023. On the effectiveness of offline RL for dialogue response generation. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 32088–32104. PMLR.
 - 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. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4661–4676, Dublin, Ireland. Association for Computational Linguistics.
 - Ryuichi Takanobu, Hanlin Zhu, and Minlie Huang. 2019. Guided dialog policy learning: Reward estimation for multi-domain task-oriented dialog. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 100–110, Hong Kong, China. Association for Computational Linguistics.
 - Wei Wei, Quoc Le, Andrew Dai, and Jia Li. 2018. Air-Dialogue: An environment for goal-oriented dialogue research. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 3844–3854, Brussels, Belgium. Association for Computational Linguistics.
 - 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.
 - 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*.
 - Xiaoying Zhang, Baolin Peng, Jianfeng Gao, and Helen Meng. 2022. Toward self-learning end-toend task-oriented dialog systems. *arXiv preprint arXiv:2201.06849*.
 - Xiaoying Zhang, Baolin Peng, Kun Li, Jingyan Zhou, and Helen M. Meng. 2023. Sgp-tod: Building task

bots effortlessly via schema-guided llm prompting.	93
<i>ArXiv</i> , abs/2305.09067v1.	93

A Appendix

A.1 Human Annotation Metrics & 934 Inter-annotator Agreement 935

We include the complete human annotation metrics and inter-annotator agreement in Table 4.

	ACTIONAWARE	ACTIONPLAN PREDICTED	COMPLIANCEOPT PREDICTED	Ground Truth
Avg Compliance	0.31	0.29	0.35	0.46
At least 1 Compliant (↑)	0.41	0.53	0.47	0.57
All Compliant (†)	0.2	0.19	0.24	0.38
At least 1 Non-Compliant (↓)	0.79	0.8	0.75	0.62
All Non-Compliant (↓)	0.58	0.47	0.53	0.42
Fleiss Kappa	0.66	0.65	0.67	0.73

 Table 4: Complete Human Annotation Metrics & Interannotator Agreement

A.2 ComplianceReward Analysis

Model	Compliance Reward
Baselines & Ablation	ns
NOACTION	0.4963
GUIDELINE ORACLE	0.5713
LLM-PROMPTING PREDICTED	0.6421
LLM-PROMPTING ORACLE	0.8410
ACTIONPLAN ALL FUTURE ORACLE	0.6043
Proposed Methods	
ACTIONAWARE	0.6012
ACTIONPLAN PREDICTED	0.6762
COMPLIANCEOPT PREDICTED	0.6742
ACTIONPLAN ORACLE	0.7925
COMPLIANCEOPT ORACLE	0.8670
Ground Truth	0.8676

Table 5: ComplianceReward Results

Table 5 shows that our proposed framework, as expected, outperforms other methods on this metric. We also computed the correlation between ComplianceReward scores and human evaluation (Table 6) on the human-annotated subset of our modelgenerated responses and found that there is a positive correlation. There is a similar positive correlation between ComplianceReward and LLM compliance (Table 7). This validates that ComplianceReward as a training signal is indeed well correlated with human measures of compliance. Moreover, the model performances across human and LLM scores are generally similar to that of the ComplianceReward rankings, with GroundTruth and COMPLIANCEOPT receiving the highest scores. 937

933

936

938

951

952

953

939

	Annotator Avg Score	Model Avg Score	Pearson	Spearman
Ground Truth	0.4631	0.6290	0.2414	0.2437
ACTIONAWARE	0.3100	0.4349	0.4075	0.4146
ActionPlan Predicted	0.2867	0.5416	0.2114	0.1645
ComplianceOpt Predicted	0.3533	0.6394	0.2819	0.2702

Table 6: Comparison of the ComplianceReward with human judgment. All correlations are significant with p < 0.05.

	LLM Score	Model Avg Score	Pearson	Spearman
Ground Truth	0.4552	0.6290	0.3526	0.5030
ACTIONAWARE	0.2553	0.4349	0.3012	0.3549
ACTIONPLAN Predicted	0.2254	0.5416	0.4362	0.5502
ComplianceOpt Predicted	0.2802	0.6394	0.3252	0.4256

Table 7: Comparison of the ComplianceReward with LLM compliance. All correlations are significant with p < 0.05.

A.2.1 Action State Tracking Evaluation

954

955

957

958

960

961

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

982

While our framework focuses on response generation and does not directly predict slot values, we extract these values from the generated responses and compute model performances on Action State Tracking (AST) metrics, which is a set of performance benchmark metrics proposed for the ABCD dataset (Chen et al., 2021a). Specifically, we adopt the approach of Lee et al. (2021) that has competitive performance in extracting slot values from generated responses and train a t5-base model to extract slot values from generated system responses.

We additionally report the workflow accuracy performances of the PREDICTED versions of AC-TIONPLAN and COMPLIANCEOPT models for comparisons. The PREDICTED models first predict the next planned workflow action and condition the next response generation on the predicted action.

From Table 8, we find that ground truth and our models perform similarly on both b-slot and value predictions. We note that our framework does not optimize these metrics, and this evaluation shows that our optimization does not lead to a strong degradation of the standard dialogue metrics.

A.3 User Simulator

We instantiate the user simulator with the NOAC-TION model.

	Workflow Accuracy	B-Slot	Value	Action (Joint)
Ground Truth	1	0.5737	0.5789	0.4895
ComplianceOpt Predicted	0.7011	0.5632	0.6000	0.5105
COMPLIANCEOPT ORACLE	N/A	0.5526	0.6211	0.4895
ActionPlan Predicted	0.6821	0.5526	0.5632	0.4684
ActionPlan Oracle	N/A	0.5632	0.5474	0.4579
ACTIONAWARE	N/A	0.5632	0.5789	0.4789

Table 8: Action State Tracking Metrics Results on theABCD dataset.

A.4 Experiment Details

A.4.1 Block Processing

To generate workflow blocks described in Section 3.1 for training the ComplianceReward model, we use the following simple approach: 983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1001

1003

1004

1006

1007

1008

1009

1010

1011

1012

1013

1014

- For each conversation in the dataset, we traverse the utterances in reverse order, starting from the final turn.
- At each moment of the traversal, we maintain an index of the latest completed workflow action. We begin by initializing the latest completed action as "Send-off/Goodbye", which is a placeholder to account for the final moments of the dialogue where no action is completed.
- When a new completed action is seen, the index is updated. Moreover, agent utterances between the old index and the new action index are marked with the old workflow action, and the marked utterances are considered as making up a "block" of exchanges, which then can be used for training the ComplianceReward model.

A.4.2 Input Formatting

For both training and inference, ACTIONPLAN and COMPLIANCEOPT models use a simple input format (Figure 4) which interleaves dialogue history, planned action, and completed actions, similar to that of SimpleTOD (Hosseini-Asl et al., 2020). NOACTION and ACTIONAWARE models use variants of the format, each without planned and completed action, and completed actions, respectively.

A.4.3 Training Procedure

Our method comprises two stages: teacher-forcing1015training and RL training with direct compliance1016





 $\ensuremath{\text{lnput}}$: Dialogue context, complete actions and planned workflow are formatted into a single sequence.

Figure 4: Our framework uses an autoregressive LM with interleaved utterances, actions, and workflow actions as the input.

optimization. ACTIONAWARE, ACTIONPLAN use teacher-forcing training with different workflow 1018 1019 action inputs while the third (COMPLIANCEOPT) uses RL training with direct compliance optimiza-1020 tion. ACTIONAWARE conditions only on completed actions $\{a_t\}$, while ACTIONPLAN additionally conditions on predicted future actions. 1023 In teacher-forcing training (ACTIONAWARE, AC-1024 TIONPLAN), the LM is trained with the standard 1025 negative log-likelihood loss, while in RL training, 1026 it is trained based on a reward signal generated by 1027 the ComplianceReward model. 1028

> ACTIONAWARE. A simple way to include workflow action information is to include the history of past workflow actions in the dialogue context C_t . In this way, an action execution a_t is treated as an utterance, and indicates that the system completed a workflow action at time t. Thus, an example context might be $C_t = [u_0, s_1, u_2, s_3, a_4, \cdots]$. Conditioned on C_t , the LM can then generate the system utterance s_{t+1} . This implicitly models the relationship between workflow policy and agent utterances, since the LM may learn to use patterns between completed actions and the next system response.

1032

1033

1034

1035

1036

1038

1039

1040

ACTIONPLAN. Including only the previously 1041 completed actions does not directly model future 1042 compliance with policy guidelines. To explicitly model this, we introduce a future workflow action, 1044 as the next workflow action that must be completed. 1045 Given a completed ground truth dialogue, we con-1046 struct this input sequence by backpropagating the 1047 1048 action execution a_t to all previous utterances $\{s_k\}$, k < t, before another action occurs. We define this action assignment as a *planned* workflow ac-1050 tion p_t associated with every system utterance s_t . We also include past actions in the context to help 1052

model workflow dynamics. An example context	1053
is $C_t = [u_0, (p_1, s_1), u_2, (p_3, s_3), a_4, \cdots]$, where	1054
$p_1 = p_3 = a_4.$	1055

1056

1057

1059

1060

1061

1062

1063

1064

1065

1067

1068

1069

1071

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

COMPLIANCEOPT. Next, we apply reinforcement learning to directly optimize workflow compliance. Specifically, a teacher-forcing trained AC-TIONPLAN model is used as an initial model for this step and uses the same input formatting. To apply the interactive sampling model to generate a block of utterance, we use a frozen copy of a teacher-forcing trained response generator.

At each training step, the agent and user models simulate a "block" session by generating one turn at a time. A turn is defined as a consecutive sequence of utterances made by the same party. A block session is terminated either when the agent completes the planned action, or the turn count reaches a predefined limit. When a session is terminated, the dialogue history and the simulated exchange are then scored by the ComplianceReward model, which are then quantized and trained according to the modified Quark algorithm (Algorithm 1).

A.4.4 Inference

The inference procedure for generating a block of interaction between the user and agent models is identical to the interactive sampling step used for optimizing the COMPLIANCEOPT model.

A.4.5 Parameters & Hyperparmeters

We list the parameters and hyperparameters we used for our experiments in Table 9. We tune our learning rate, interactive sampling temperature, and the number of quantiles K using grid search.

Teacher-Forcing Setting				
Agent model	distilgpt2			
LLM prompting model	gpt-3.5-turbo			
Scoring model detail	roberta-base			
Training epochs	10 (ABCD) / 1 (MultiWOZ)			
Learning Rate	2e-5			
Special tokens	START_USER, START_WORKFLOW,			
*	END_WORKFLOW, START_AGENT,			
	END_AGENT, START_ACTION,			
	END ACTION, START DIALOG,			
	END_DIALOGUE			
Compliand	ce Optimization RL Setting			
Agent model	distilgpt2 (Warm-start ACTIONPLAN)			
Client model	distilgpt2			
Learning Rate	2e-5			
Sampling Temperature	0.5			
Number of Interactions	3			
Number of Quantiles K	5			
KL weight β	0.05			
Training Steps	80k (ABCD) / 160k (MultiWOZ)			

Table 9: Experiment Models & Parameters

1085	A.5 LLM Prompts	Evaluate the agent's performance in an
1086		interaction with a customer based
1087	generation_prompt = f"You are a	on a set of categories to gain
1088	cusoumer agent netping a customer	insights into various aspects of
1089	with a issue. Read the dialogue	agent behavior. The primary
1090	context, provided policy	category to consider is workflow
1091	guideline, and generate an agent	compliance, which determines if
1092	utterance to nelp the customer in	the agent has successfully
1093	a way that is compliant to the	achieved the objectives outlined
1094	guideline. The generated agent	in the provided workflow action
1095	turn should be at most 2	during their interaction with the
1096	utterances, and should be similar	client. The other category is
1097	in length to the agent utterances	coherence.
1098	shown in the examples that	
1099	demonsrtate compliant agent	Relevant Documents
1100	<pre>behavior.\n\Custome Situation:</pre>	ABCD Guideline: This document consists
1101	<pre>{s}\n\Policy Action Name:</pre>	of comprehensive descriptions for
1102	<pre>{w}\n\Policy Name Guideline:</pre>	each customer assistance subflow,
1103	<pre>{g}\n\n\{example_str}dialogue</pre>	including specific action steps
1105	Context: {i}\n\n\Agent: "	within each subflow.
1106	ovaluation prompt - f"Poad the provide	To locate a particular workflow
1107	guideline and access the extent to	action, begin by referring to the
8011	yuldeline and assess the extent to	relevant subflow section (e.g.,
1109	which the agent's behavior in the	Initiate Refund), and then
1110	input interaction aligns with the	identify the corresponding
1111	specified workitow action,	workflow action enclosed within
1112	considering the name and a concise	brackets (e.g., [Pull Up Account]).
1113	description of the workflow	
1114	provided. I = Compliant(NO =	In the annotation sheet we will also
1115	Non-compliant\n\nSubtlow:	provide brief policy quidelines
1116	{s}\nworktlow: {w}\nDescription:	alongside examples to aid in the
1117	{g}\n\n\Dialogue	annotation process Ideally
1118	History:\n{i}\n\nInput	policy quidelines should be
1128	Interaction:\n{r}\n\nAnswer:"	sufficient for the annotations.
1101	A.6. Human Evoluation Cuidalinas (Abridged)	Cotomonico
1121	A.o Human Evaluation Guidelines (Abridged)	Lategories
1122	Compliance: Assess if the agent's behavior aligns with the specified workflow action, taking into account the action's	1. Compliance
1123	name and policy guideline. If the agent has already completed	Assess the degree to which the agent's
1125	certain steps or the entire policy guideline behavior in the	behavior aligns with the specified
1126	dialogue history, they should not be penalized for not repeating	workflow action, taking into
1127	mose corresponding steps.	account the action's name and
1128	Coherence: Rate the coherence of the agent's interac-	policy guideline. Please refer to
1129	tion on a binary scale (0=not coherent, 1=coherent). In this	the provided document for more
1130	evaluation, please do not consider repetitive agents as coher-	detailed information. If the agent
1131	ent. Additionally, do not include incoherent or disfluent client	nas already completed certain
1132	behavior in the evaluation (only evaluate agent behavior).	steps or the entire policy guideline behavior in the dialogue

A.7 Human Evaluation Guidelines (Full)

Agent Quality Annotation Task Task

history, they should not be

corresponding steps.

penalized for not repeating those

1189	<pre>1 = Compliant: The agent successfully</pre>	Subf
1190	executes all the steps outlined in	Work
1191	the policy guideline.	Desc
1192	0 = Non-compliant: The agent fails to	
1193	execute any of the steps mentioned	
1194	in the policy guideline.	
1195		
1196		0pti
1197	Examples:	
1198	0 = Non-compliant	
1199	Subflow: out_of_stock_general	Ente
1200	Workflow Action: notify-team	Dial
1201	Policy Guideline: Let the customer	Targ
1202	know that you will write up a	Agen
1203	report and let the Purchasing	
1204	Department know about this, so	
1205	they can do a better job.	Clie
1206		Clie
1207	Enter 'purchasing department' into the	
1208	input box and [Notify Internal	
1209	Team]	Agen
1210	Dialog History: Omitted Here	
1211	Target Generation	Clie
1212	Agent: i m sorry for your	
1213	inconvenience	
1214	Client: ok	
1215	Agent: i can offer you a promo code	2. 0
1216		Plea
1217	1=Compliant:	
1218	Subflow: out_of_stock_general	
1219	Workflow Action: notify-team	
1220	Policy History: Let the customer know	
1221	that you will write up a report	
1222	and let the Purchasing Department	
1223	know about this, so they can do a	
1224	better job.	
1225		
1226	Enter 'purchasing department' into the	
1227	input box and [Notify Internal	Cohe
1228	Team]	
1229	Dialog Context: Omitted Here	
1230	Target Generation	
1231	Agent: okay, i'm sorry for this	
1232	situation, what item is out of	
1233	stock?	Exam
1234	Client: product: iacket brand:	0=no
1235	michael_kors amount: \$49	
1236	Agent: the most i can do for vou is	
1237	write up a report to the	Agen
1238	purchasing department about this.	
1239	would that be alright?	Clie
1240		
-		1

Subflow: shipping status	1241
Norkflow Action: update-order	1242
Description: If the Oracle says No,	1243
then the customer will not be	1244
happy. To resolve, enter how you	1245
will fix the problem	1246
	1247
Options include: 'change date',	1248
'change address', 'change item',	1249
or 'change price'	1250
Enter into [Update Order]	1251
Dialog Context: Omitted Here	1252
Target Generation	1253
Agent: it seems the email was	1254
incorrect. when were you expecting	1255
it to arrive?	1256
Client: tomorrow pm.	1257
Client: i'm just trying to verify if	1258
it is still coming or something	1259
weird happened.	1260
Agent: okay, i can change the date to	1261
be sure you receive by tomorrow.	1262
Client: ok, that's fine, thanks.	1263
	1264
	1265
	1266
2. Coherence	1267
Please rate the coherence of the	1268
agent's interaction on a binary	1269
scale (0=not coherent,	1270
l=coherent). In this evaluation,	1271
please do not consider repetitive	1272
agents as conerent. Additionally,	1273
do not include inconerent or	1274
distigent client behavior in the	1275
evaluation (only evaluate agent	1276
Denavior).	1277
Cabaranca Standard, Ara all of the	1278
system responses in the	1279
system responses in the	1280
to each other?	1201
	1002
	1203
Examples	1204
A=not coherent (agent repeats oneself	1203
contradicts oneself or save	1200
illogical things).	1207
Agent: thank you albert may i also	1200
have your order id?	1209
Client: it was 9940404040	1291

1292	Agent: thank you albert. may i also
1293	have your order id?
1294	Client: sure, it's 9940404040
1295	Agent: thank you albert. may i also
1296	have your account id and order id?
1297	
1298	Agent: how much was the service?
1299	Client: it was \$40.
1300	Agent: how much was the extra price?
1301	Client: i was charged \$40
1302	Agent: how much was the price you were
1303	charged?
1304	
1305	Agent: unfortunately, since the order
1306	has already been shipped i can not
1307	cancel the order. since it has
1308	already been shipped i can only
1309	cancel the order.
1310	Client: ok, thank you
1311	Agent: i apologize for the
1312	inconvenience
1313	Agent: since the order has already
1314	been sent i can not cancel the
1315	shipment.
1316	Client: ok, i understand
1317	Client: that is all i needed
1318	Agent: great. have a nice day!
1319	
1320	
1321	1=coherent:
1322	Agent: ok, i see your refund is in
1323	progress and it looks like it
1324	should be going through to
1325	completion later today or by
1326	tomorrow at the latest
1327	Client: okay, thank you
1328	Agent: you're welcome
1329	Agent: can i help with anything else?
1330	Client: that will be all
1331	Agent: ok, have a good day
1332	
1333	Agent: thanks for your information.
1334	Agent: the system said that your
1335	shipping address is the same as
1336	the one you stated above. the
1337	email was incorrect. you can
1338	ignore it.
1339	Client: thank you that's all i needed
1340	to know
1341	Agent: great, is there anything else
10.10	
1342	that i can help you with?

Agent: have a nice day!