A Distilled Representation for Zero and Few-Shot Localization of Task-Oriented Dialogue Agents

Anonymous ACL submission

Abstract

Task-oriented Dialogue (ToD) agents are 002 mostly limited to a few widely-spoken languages, mainly due to the high cost of acquiring training data for each language. Existing low-cost approaches that rely on cross-lingual embeddings or naive machine translation sacrifice a lot of accuracy for data efficiency, and largely fail in creating a usable dialogue agent. We propose automatic methods that use ToD training data in a source language to build a high-quality functioning dialogue agent in another target language that has no training data 013 (i.e. zero-shot) or a small training set (i.e. fewshot). Unlike most prior work in cross-lingual ToD that only focus on Dialogue State Tracking (DST), we build an end-to-end agent.

> We show that our approach closes the accuracy gap between few-shot and existing fullshot methods for ToD agents. We achieve this by (1) improving the dialogue data representation, (2) improving entity-aware machine translation, and (3) automatic filtering of noisy translations.

We evaluate our approach on the recent bilingual dialogue dataset BiToD. In Chinese to English transfer, in the zero-shot setting, our method achieves 46.7% and 22.0% in Task Success Rate (TSR) and Dialogue Success Rate (DSR) respectively. In the few-shot setting where 10% of the data in the target language is used, we improve the state-of-the-art by 15.2% and 14.0%, coming within 5% of full-shot training.¹

1 Introduction

018

019

039

While dialogue agents in various forms have become commonplace in parts of the world, their lack of support for most human languages has prevented access to the benefits they provide for much of the world. Commercial virtual assistants for example, only support a handful of languages, as extending their functionality to each new language is extremely costly, partially due to the need for collecting new annotated training data in that language.

In recent years, several non-English taskoriented dialogue (ToD) datasets have been created; they are either collected from scratch (Quan et al., 2020; Zhu et al., 2020), paraphrased from synthetic sentences by crowdworkers (Lin et al., 2021), or manually translated from another language (Li et al., 2021b). All of these approaches are laborintensive, expensive, and time-consuming; such investment is unlikely to be made for less widely spoken languages.

Cross-lingual transfer, i.e. using training data from other languages to build a dialogue agent for a specific language, seems especially appealing. An emerging line of work has employed machine translation of training data, and multilingual pre-trained neural networks to tackle this task (Sherborne et al., 2020; Li et al., 2021a; Moradshahi et al., 2021). However, work in ToD cross-lingual transfer has for the most part, focused on understanding the user input, namely Dialogue State Tracking (DST) and Natural Language Understanding (NLU). Other necessary parts of a dialogue agent like policy and response generation have mostly remained unexplored.

In this paper, we present a methodology for building a fully functional dialogue agent for a new language (e.g. English), by using training data in another language (e.g. Chinese) with little to no additional manual dataset creation effort. We found that despite prior efforts to improve modeling for existing ToD datasets, the dialogue representation used as input to these models, e.g. full dialogue history in natural language (Hosseini-Asl et al., 2020), is sub-optimal, especially when the training data is either scarce or created automatically using noisy machine translation. We propose a new Distilled representation to fix the shortcomings of 041

042

043

044

045

047

049

¹We will release our code and data upon publication.

163

164

165

167

168

169

170

171

172

173

174

175

176

177

178

179

180

current representations. We also found that previously proposed entity-aware translation technique Moradshahi et al. (2021) to be inadequate. Our proposed technique effectively combines entity-aware neural machine translation with text similarity classifiers to automatically create high-quality training data for a new language. This paper explains all the ingredients we found useful, and motivates their use by performing extensive ablation studies.

The contributions of this paper are:

087

097

100

103

104

105

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

- A new state-of-the-art result for the BiToD dataset in both few-shot and full-shot settings according to all of our 6 automatic metrics, including an improvement of 14.0% and 2.9% respectively in Dialogue Success Rate (DSR). In fact, using our Distilled representation, our few-shot model trained on only 10% of the training data, achieves similar results to the previous SOTA model trained on 100% training data.
- 2. The first dialogue agent created in the zero-shot cross-lingual transfer setting, i.e. starting from no training data in the target language. Our agent achieves 71%, 62%, 40%, and 47% of the performance of a full-shot agent in terms of Joint Goal Accuracy (JGA), Task Success Rate (TSR), DSR, and BLEU score, respectively.
- 3. A concise dialogue representation designed for cross-lingual ToD agents. The Distilled dialogue representation works well with our new decomposition of agent's subtasks, making these significant improvements possible.
- 4. An improved methodology for high-quality automatic translation of ToD training data. We adapt and improve an existing entity-aware machine translation system that localizes entities (Moradshahi et al., 2021), extend it to agent response generation, and equip it with a filtering step that increases the quality of the resulting translations.

2 Related Work

2.1 Multilingual Dialogue Datasets

MultiWOZ (Budzianowski et al., 2018; Ramadan 123 et al., 2018; Eric et al., 2019) and CrossWOZ (Zhu 124 et al., 2020) are two monolingual Wizard-Of-Oz 125 dialogue datasets that cover several domains, suitable for building travel dialogue agents in En-127 glish and Chinese respectively. For the 9th Dia-128 log System Technology Challenge (DSTC-9) (Gu-129 nasekara et al., 2020), they were translated to Chi-130 nese and English using Google Translate. Glob-131

alWOZ (Ding et al., 2021) is another translation of MultiWOZ to Spanish, Chinese and Indonesian, with human translators post-editing machine translated dialogue templates, and filling them with newly collected local entities.

Different from these translation approaches, Lin et al. (2021) introduced BiToD, the first bilingual dataset for end-to-end ToD modeling. BiToD uses a dialogue simulator to generate dialogues in 5 tourism domains in English and Chinese, then uses crowdsourcing to paraphrase entire dialogues to be more natural. Unlike WOZ-style datasets which usually suffer from poor annotation quality due to human errors (Moradshahi et al., 2021), BiToD is automatically annotated during synthesis. Since neither manual nor machine translation is used in the creation of BiToD, it does not contain translationese (Eetemadi and Toutanova, 2014) or other artifacts of translated text (Clark et al., 2020), and provides a realistic testbed for cross-lingual transfer of task-oriented dialogue agents.

2.2 Multilingual Dialogue State Tracking

Mrkšić et al. (2017) proposed using cross-lingual word embeddings for zero-shot cross-lingual transfer of DST models. With the advent of large language models, contextual embeddings obtained from pre-trained multilingual language models (Devlin et al., 2018; Xue et al., 2021; Liu et al., 2020) have been used to enable cross-lingual transfer in many natural language tasks, including DST.

Chen et al. (2018) used knowledge distillation (Hinton et al., 2015) to transfer DST capabilities from a teacher DST model in the source language to a student model in the target language.

Schuster et al. (2019) used contextual crosslingual representations obtained from machine translation models and reported that it performs better than training with machine translated training data for single-turn commands. Moradshahi et al. (2021) proposed an entity-aware machine translation method to improve the quality of translated DST data.

3 Distilled ToD Agent

Our methodology includes a dialogue task decomposition and a Distilled dialogue representation that are tailored to cross-lingual ToD agents. In this section we describe these two components.

We follow the end-to-end task-oriented dialogue (ToD) setting (Hosseini-Asl et al., 2020) where a



Figure 1: Inference-time flow diagram for our dialogue agent. DST, API, DA, and RG share the same neural model. \otimes indicates text concatenation. \oplus refers to the update rule in Equation 1.

user converses freely with an agent over several turns to accomplish his/her goal with all of its constraints (e.g. "book a restaurant that is rated at least 3."). In each turn, the agent must access its database if needed to find the requested information (e.g. find a restaurant that satisfies user constraints), decide on an action (e.g. to present the information to the user or to ask for additional information) and finally respond to the user in natural language based on the action it selects.

3.1 Preliminaries

181

182

183

184

187

188

189

190

192

193

194

196

198

204

Formally, a *dialogue* $D = \{U_1, A_1, ..., U_T, A_T\}$ is a set of alternating user utterances U_t and agent responses A_t for a number of turns T.

A *belief state* at turn t, B_t , consists of a list of $\langle domain, intent \rangle$ tuples and a set of $\langle slot, relation, value \rangle$ tuples. *Intent* is the user intent, either search or book. *Relation* is a comparison or membership operator. *Value* can be one or more entity names or strings from the ontology, or a literal. To see all possible domains, slots and values please refer to Table 4 in Lin et al. (2021).

Levenshtein belief state (Lin et al., 2020) is the difference between belief states in consecutive turns, i.e. $\Delta B_t = B_t - B_{t-1}$. It captures only the relations and values that have changed in the last user utterance, or tuples that have been added or removed.

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

228

229

230

231

232

233

235

236

237

239

249

242

243

244

245

246

247

248

249

251

252

253

258

Agent dialogue acts at turn t, C_t , are a list of $\langle domain, intent \rangle$ tuples and a set of $\langle dialogue_act_name, slot, value \rangle$ tuples indicating the action the agent takes and the information offered to the user, if any.

3.2 Task Decomposition

The task of dialogue agents is usually broken down to several subtasks, which may be performed by a pipelined system (Gao et al., 2018) or by a single neural network (Hosseini-Asl et al., 2020; Lei et al., 2018). Here we describe our subtasks and their inputs and outputs (Figure 1).

After the user speaks at turn t, the agent has access to the belief state up to the previous turn (B_{t-1}) , the history of agent dialogue acts $(C_1, ..., C_{t-1})$, and the history of agent and user utterances so far $(A_1, ..., A_{t-1} \text{ and } U_1, ..., U_t)$. Our agent performs the following four subtasks:

1. Dialogue State Tracking (DST): Generate ΔB_t , the Levenshtein belief state, for the current turn based on the previous belief state, the last two agent dialogue acts, and the current user utterance. ΔB_t is combined with B_{t-1} to produce the current belief state.

$$\Delta B_t = \text{DST}(B_{t-1}, C_{t-2}, C_{t-1}, U_t)$$

$$B_t \leftarrow B_{t-1} + \Delta B_t$$
(1)

2. *API Call Detection(ACD)*: Call an API to query the database, if needed.

$$q_t = \text{ACD}(B_t, C_{t-2}, C_{t-1}, U_t, R_{t-1})$$
 (2)

$$R_t \leftarrow q_t? \operatorname{KB}(B_t) : \emptyset \tag{3}$$

In turn t, ACD determines if an API call is necessary. If so, the result R_t is the top entity in the knowledge base KB, based on some deterministic ranking, that matches the API call constraints in B_t , and is empty otherwise. If no entities match the constraint, we set R_t to the special value NORESULT.

3. Dialogue Act Generation (DAG): Generate C_t , the agent dialogue act for the current turn based on the current belief state, the last two agent dialogue acts, the user utterance, and the result from the API call.

$$C_t = \text{DAG}(B_t, C_{t-2}, C_{t-1}, U_t, R_t)$$
 (4)

4. Response Generation (RG): Convert the agent dialogue act C_t to the new agent utterance A_t .

257

258

Note that C_t contains all the necessary information for this subtask. However, providing U_t improves response fluency and choice of words, translating to a higher BLEU score, partly due to mirroring (Kale and Rastogi, 2020).

$$A_t = \mathrm{RG}(U_t, C_t) \tag{5}$$

3.3 The Distilled Dialogue Representation

The design of Distilled is based on the following principles:

- 1. For cross-lingual agents, it is important to reduce the impact of translation errors. The representation should make minimal use of natural language, and use formal representation where possible.
- 2. Dialogues can get long, but the representation should be succinct, containing only the necessary information, so the neural network need not learn to ignore unnecessary information from copious data. This improves data efficiency as well as training/ inference speed of neural models.

We note that BiToD's original representation (Lin et al., 2021) follows neither of these principles². It makes extended use of natural language: all previous user and agent natural language utterances are included in the input of all subtasks. It has many redundancies: for each subtask, it inputs the concatenation of all previous subtask's inputs and outputs. In the following, we highlight the changes we made to the (Lin et al., 2021) representation.

Replace agent utterances with formal agent dialogue acts. Since agent responses are automatically generated, it is possible to capture all information useful to the different subtasks with formal agent dialogue acts. This way, the neural network needs not to interpret previous natural language utterances.

We take two steps to generate the agent responses: DA (Dialogue Act) first produces the formal act, C_t , which is then fed into RG to generate the natural language response A_t . RG is not part of the dialogue loop, in the sense that A_t only serves to communicate to the user; the C_t from DA is used as input to all the subtasks instead. In contrast, Lin et al. (2021) generates the agent response directly from API results. Hosseini-Asl et al. (2020) also

separates response generation into two steps, but they use A_t instead of C_t as input for the next turn. 304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

340

341

342

343

344

345

347

348

349

351

Note that the agent dialogue acts are languageindependent - this is beneficial to cross-lingual agents as it can learn easier from data available in other languages. Furthermore, DA can be validated on whether the output dialogue acts match the gold answers exactly. This is not possible with natural language results, whose quality is typically estimated with BLEU score.

Shorten user utterance history. Since the belief state formally summarizes what the user has said, we remove previous user utterances $U_1, ..., U_{t-1}$ from input to all subtasks, relying on belief state B_{t-1} instead.

Untangle API Call Detection from Response Generation. After DST is done, depending on whether or not an API call is needed. Lin et al. (2021) either directly generates the agent response, or makes the API call and then generates the response in two steps. Our design is to always take two steps: (1) generate the API call or indicate that there is none, and (2) generate the agent response.

4 **Automatic Dialogue Data Translation**

Given a training dataset for one language, we automatically generate a training set in the target language we are interested in. This problem has been studied in the context of NLU for questions (Moradshahi et al., 2020) and for dialogues (Moradshahi et al., 2021). One challenge is that the translated dataset should refer to entities in the target language. Thus, Moradshahi et al. (2020) proposed to first use cross-attention weights of the neural translation model to align entities in the original and translated sentences, then replace entities in the translated sentences with local entities from a target language knowledge base. Our initial experiments showed that applying this approach directly to end-to-end dialogue datesets does not yield good performance. Thus, we adapted and improved this approach for dialogues as discussed in the following sections.

4.1 **Alignment for Dialogues**

First, we found that while translation with alignment works for NLU, it does not work well for RG. Machine translation introduces two kinds of error: (1) Translated sentences can be ungrammatical, incorrect, or introduce spurious information.

²We found this to be true for several previously-proposed popular representations of MultiWOZ as well (Lei et al., 2018; Chen et al., 2019).

(2) The alignment for entities may be erroneous, which can seriously hurt the factual correctness of the responses. As shown in Moradshahi et al. (2021), these errors are tolerable in NLU since (1) sentences are seen by machines, not shown to users, (2) pre-trained models like mBART are somewhat robust to noisy inputs, since they are pre-trained on perturbed data. However, training with such low-quality data is not acceptable for RG, since the learned responses are shown directly to the user.

353

354

357

370

371

372

374

389

391

395

400

Second, we found alignment recall to be particularly low for an important category: entities that are mostly quantitative. Dates, times, and prices can be easily mapped between different languages using rules. We propose to first try to translate such entities with dictionaries such as those available in dateparser (Scrapinghub, 2015) and num2words (faire Linux, 2017) and to match them in the translated text. Only if no such match is found, do we resort to using neural alignment.

4.2 Filtering Translation Noise for RG

To reduce translation noise for RG, we automatically filter the translated data based on the semantic textual similarity between the source and translated sentences. For this purpose, we use LaBSE (Feng et al., 2020), a multilingual neural sentence encoder based on multilingual BERT (Devlin et al., 2018), trained on translation pairs in various languages with a loss function that encourages encoding pairs to similar vectors. To score a pair of sentences, the model first calculates an embedding for each sentence and computes the cosine distance between those vectors. The lower the distance is, the more semantically similar the sentences are according to the model.

In creating the RG training set, we first translate the source agent utterances to the target language and use LaBSE to remove pairs whose similarity score is below a threshold. We found a threshold of 0.8 to work best empirically. Higher thresholds would inadvertently filter correctly translated utterances. We construct the final training data by pairing aligned translated utterances that pass the filter with their corresponding translated agent dialogue acts.

5 Experiment Setting

5.1 Base Dataset

We perform our experiments on BiToD, a largescale high-quality bilingual dataset created using the Machine-to-Machine (M2M) approach. It is a multi-domain dataset, including restaurants, hotels, attractions, metro, and weather domains. It has a total of 7,232 dialogues (3,689 dialogues in English and 3,543 dialogues in Chinese) with 144,798 utterances in total. The data is split into 5,787 dialogues for training, 542 for validation, and 902 for testing. The training data is from the same distribution as validation and test data. 401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

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

5.2 Evaluation Metrics

We use the following metrics to compare different models. Scores are averaged over all turns unless specified otherwise.

- Joint Goal Accuracy (JGA) (Budzianowski et al., 2018): Is the standard metric for evaluating DST. JGA for a dialogue turn is 1 if all slot-relation-value triplets in the generated belief state match the gold annotation, and is 0 otherwise.
- Task Success Rate (TSR) (Lin et al., 2021): A task, defined as a pair of domain and intent, is completed successfully if the agent correctly provides all the user-requested information and satisfies the user's initial goal for that task. TSR is reported as an average over all tasks.
- **Dialogue Success Rate** (**DSR**) (Lin et al., 2021): DSR is 1 for a dialogue if all user requests are completed successfully, and 0 otherwise. DSR is reported as an average over all dialogues. We use this as the main metric to compare models, since the agent needs to complete all dialogue subtasks correctly to obtain a full score on DSR.
- **API** (Lin et al., 2021): For a dialogue turn, is 1 if the model correctly predicts to make an API call, and all the constraints provided for the call match the gold. It is 0 otherwise.
- **BLEU** (Papineni et al., 2002): Measures the natural language response fluency based on n-gram matching with the human-written gold response. BLUE is calculated at the corpus level.
- Slot Error Rate (SER) (Wen et al., 2015): It complements BLEU as it measures the factual correctness of natural language responses. For each turn, it is 1 if the response contains all entities present in the gold response, and is 0 otherwise.

6 Results and Discussion

We first show how our Distilled representation affects the performance of an agent in a full-shot

Representation	JGA ↑	TSR ↑	DSR ↑	API ↑	BLEU ↑	$SER \downarrow$
Original (Lin et al., 2021)	69.19	69.13	47.51	67.92	38.48	14.93
Distilled (ours)	76.79	75.64	53.39	76.33	42.54	10.61
 Generate full state 	74.30	74.19	50.90	73.93	41.90	11.38
 Natural agent response 	75.62	73.41	49.10	73.93	40.94	11.90
 Only last agent turn 	73.97	74.19	52.71	74.27	41.83	11.81
• Prev. user utterance as state	71.75	61.66	33.94	67.67	39.72	15.97
Remove state	70.84	51.89	24.43	66.47	37.10	19.61

Table 1: Full-shot English monolingual training with ablation. All results are reported on the English test set of BiToD using the same evaluation script. The best result is in bold.

setting. We then evaluate our proposed techniques on cross-lingual settings with varying amounts of available training data.

6.1 Evaluation of Distilled Representation

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

To understand how our design of Distilled representation affects the performance of ToD agents in general, we train an English agent using all the English training data and perform an ablation study (Table 1). We observe that even though the Distilled representation removes a lot of natural language inputs, it improves the best previous English-only results on JGA, TSR, DSR, API, BLEU and SER by 7.6%, 6.5%, 5.9%, 8.4%, 4.1%, and 4.7% respectively. This suggests that natural language utterances carry a lot of redundant information, and the verbosity may even hurt the performance. Note that the improvement in BLEU is also accompanied by an improvement of factuality measured by SER.

Furthermore, using the Distilled representation reduces training time by a factor of 3. See Section A.1 for more details.

Generate full state. Our first ablation study confirms that the proposal by Lin et al. (2020) to predict ΔB_t is indeed better than B_t . Note that the training time per gradient step is more than twice as long in this ablation since the outputs are longer.

Natural agent response. Here we use natural language agent responses as input instead of agent dialogue acts, replacing C_{t-1}, C_{t-2} with A_{t-1}, A_{t-2} . The drop in TSR and DSR shows this is an important design choice - distilling natural language into a concise formal representation improves model ability to understand the important information in the sentence.

484 **Only last agent turn.** When we remove C_{t-2} 485 from the input and only use C_{t-1} , we observe a 486 drop across all metrics. This is because some turns 487 in BiToD refer to the agent's states from two turns 488 ago. We experimented with carrying three turns, 489 but there was no improvement. **Previous user utterance as state.** To investigate how much information the previous user utterance U_{t-1} contains, we use it instead of B_t in subtask inputs. Compared to all previous ablations, accuracy drastically decreases across all metrics, especially JGA. This is expected since dialogues contain long-range dependencies and some entities are referenced from earlier turns. This shows that the dataset is highly contextual and therefore a summary of conversation history is necessary.

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

Remove state. We remove B_t without adding back the previous user utterance U_{t-1} . Compared to the previous ablation, TSR and DSR drop by 10.5% and 5.2% respectively. This difference shows U_{t-1} does contain part of the information captured in B_t .

6.2 Evaluation of Cross-Lingual Transfer

The goal of this experiment is to create an agent in a *target* language, given full training data in a source language (\mathcal{D}_{src}), and a varying amount of training data in a target language (\mathcal{D}_{tgt}). We also assume that valuation and test data are available in both source and target languages. We chose Chinese as the source language and English as the target language so we can perform error analysis, and model outputs are understandable for a wider audience.

6.2.1 Varying Target Training Data

Full-Shot. In the full-shot experiments, all of \mathcal{D}_{tgt} is available for training. We train two models on two data sets: (1) on a shuffled mix of \mathcal{D}_{src} and D_{tgt} . (2) on \mathcal{D}_{tgt} alone. The ablation "-*Mixed*" in Table 2 refers to the latter.

Zero-Shot. In our zero-shot experiments, \mathcal{D}_{tgt} is not available. Instead, we automatically create training data from \mathcal{D}_{src} as follows:

Canonicalization: We translate domain names, slot names, agent dialogue acts, and API names in \mathcal{D}_{src} to the target language to match those in the validation and test data of the experiment. BiToD

Setting	JGA ↑	TSR ↑	DSR ↑	API ↑	BLEU ↑	SER \downarrow
Full-Shot						
MinTL(mT5)	72.16	71.18	51.13	71.87	40.71	13.75
- Mixed	69.19	69.13	47.51	67.92	38.48	14.93
MinTL(mBART)	69.37	42.45	17.87	65.35	28.76	-
- Mixed	67.36	56.00	33.71	57.03	35.34	-
Ours	77.52	75.04	54.07	74.44	41.46	11.17
- Mixed	76.79	75.64	53.39	76.33	42.54	10.61
			Zero-Shot			
Ours	55.33	46.74	21.95	63.04	20.01	20.52
- Filtering	54.83	45.03	19.68	60.81	19.11	20.86
- Alignment	47.21	4.72	1.13	52.74	8.26	39.20
- Translation	14.73	3.52	1.58	6.26	0.69	41.30
- Canonicalization	2.13	1.20	0.00	0.26	0.25	42.39
Few-Shot (1%)						
Ours	64.60	57.89	34.16	62.09	28.15	17.94
- Filtering	63.88	57.80	32.35	59.95	28.00	18.57
- Alignment	58.86	51.89	23.76	57.12	26.84	21.56
- Translation	49.58	41.34	19.68	46.05	22.73	24.86
- Canonicalization	44.56	42.97	20.36	46.23	23.08	24.77
- Pre-training	25.08	24.61	11.09	23.67	18.71	32.62
		Fe	w-Shot (10%)			
MinTL(mT5)	58.85	56.43	34.16	57.54	31.20	-
- Translation	48.77	44.94	24.66	47.60	29.53	19.75
- Pre-training	19.86	6.78	1.36	17.75	10.35	-
MinTL(mBART)	37.50	21.61	10.18	27.44	17.86	-
- Translation	42.84	36.19	16.06	41.51	22.50	-
- Pre-training	4.64	1.11	0.23	0.60	3.17	-
Ours	72.70	71.61	48.19	72.56	36.02	12.71
- Filtering	72.45	69.55	44.57	69.55	34.67	13.62
- Alignment	68.40	63.38	38.24	63.38	32.99	16.63
- Translation	67.13	63.12	41.40	63.64	32.86	16.40
- Canonicalization	- Canonicalization 64.51 63.64 40.27 62.69 32.71 16.63					
- Pre-training	57.18	54.80	28.73	55.66	29.61	19.66

Table 2: All results are reported on the original English test set of BiToD using the same evaluation script. The best result in each section is in bold. Each "-" removes one additional component from the previous row. All MinTL results are from Lin et al. (2021). SER numbers are computed only for models that were available.

dataset has a one-to-one mapping for most of those parameters; we added the missing items.

Translation: We use machine translation to convert the user and agent utterances and slot values in \mathcal{D}_{src} to create a training set for the target language.

Alignment: After translating the data, we use alignment (Section 4) to localize entities while ensuring the entities in translated utterances still match the values specified in annotations.

In Table 2, *Ours* refer to our main approach, which combines all three techniques. Each ablation incrementally takes away one of the techniques. Note that any performance we get from training on a dataset with just using canonicalization is from the power of cross-lingual embeddings of mBART model.

546Few-Shot. In the few-shot setting, we start with547our pre-trained zero-shot models (with various ab-548lations) and further fine-tune it on 1% and 10% of549 \mathcal{D}_{tgt} , which include 29 and 284 dialogues respec-550tively. Lin et al. (2021) reported the results only551for the 10% setting. We use their few-shot data

split in that case to be directly comparable. We add one more ablation study where we eliminate pre-training altogether to measure the pure effect of the Distilled representation.

6.2.2 Baseline

We compare our results to the best previously reported result on BiToD from Lin et al. (2021). This SOTA result was obtained using MinTL (Lin et al., 2020) and using a single mT5-small model to perform all dialogue subtasks.

Contrary to what Lin et al. (2021) reported, we found that mBART-large model outperforms mT5-small in all settings. Nevertheless, we have included all the results including MinTL(mBART) in Table 2 for comparison.

6.2.3 Results

The results for our cross-lingual experiment are reported in Table 2. Overall, in the full-shot setting, when training on both source and target language data, we improve the SOTA in JGA by 5.3%, TSR by 3.8%, DSR by 2.9%, API by 2.6%, BLEU by

575

576

577

581

584

585

587

589

590

591

597

598

603

604

610

611

612

613

615

616

617

619

620

621

0.8%, and SER by 2.6%.

Our zero-shot agent achieves 71%, 62%, 40%, and 47% of the performance of a full-shot agent in terms of JGA, TSR, DSR, and BLEU score, respectively. In the 10% few-shot setting, our approach establishes a new SOTA by increasing JGA, TSR, DSR, API, BLEU, and SER absolutely by 13.9%, 9.1%, 3.9%, 2.2%, 0.4%, and 5.1% respectively. Prominently, training with just 10% of the data beats the full-shot baseline which is trained on 100% of the training data, on all metrics except for DSR and BLEU. It also comes within 5% of full training using the Distilled representation confirming building high-quality conversational systems for a new language is possible using translated data alongside a few native training samples.

Our Distilled representation improves the performance, especially in few-shot. Comparing our results with that of Lin et al. (2021), in the full-shot monolingual setting (MinTL(mT5) "-Mixed" vs. Ours "-Mixed"), models trained on data with our representation outperform the baseline on all metrics. In the pure few-shot (10%) setting (comparing MinTL(mT5) "-Pre-train" vs Ours "-Pre-train"), our model significantly outperforms the baseline in all metrics. This suggests that our Distilled representation is much more effective in low-data settings.

Canonicalization is useful. Comparing "-Translation" with "-Canonicalization", we see that in all settings, training on data with domain names, slot names, and dialogue acts are translated to the target language significantly improves the results in the zero-shot setting. This is intuitive since canonicalization makes training data closer in vocabulary to the test data in the target language. This improvement comes at almost no cost since translation is done automatically using a dictionary.

Automatic translation of the training set works in zero-shot. By looking at how much performance is lost when we do not use translated data at all (going from "- Alignment" to "- Translation"), we find out the difference between cross-lingual pretraining and translating the train set. The naive translation approach completely fails in the zeroshot setting by achieving only 4.7% in TSR, and 1.1% in DSR, as translated entities might no longer match with ones in the annotation. However, compared to "- Canonicalization", JGA improves by 32.5% as the entities are in English after translation. Adding few-shot data helps significantly as the gap closes between "- Alignment" and "- Translation" ablations.

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

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

Alignment improves translation quality in all settings and metrics. With alignment, the translation approach performs much better in all settings, establishing a new state-of-the-art in zero and fewshot settings according to almost all metrics. As a general trend, the lower data settings benefit more from alignment. In fact, the combination of translation and alignment performs so well that our zeroshot with translated data is better than 1% pure few-shot (without pre-training) on all metrics. We additionally performed an experiment using the alignment proposed by (Moradshahi et al., 2021). There is a 4.0% drop in TSR and 4.5% in DSR, confirming the benefit of our improved alignment.

Filtering noise for RG improves fluency. We perform an ablation by training separate models on filtered and unfiltered translated agent utterances. The filtering process is descried in section 4.2. In 10% fewshot setting, both BLEU and SER improve by 1.4% confirming that automatically removing poor translations from training data improves the agent response quality.

7 Conclusion

In this work, we show how, given a dialogue dataset in one language, we can build a fully functioning dialogue agent in a new language automatically using entity-aware machine translation and our new Distilled dialogue representation.

The performance can be further improved if a few training examples in the target language are available, and we show that our approach outperforms existing ones in this setting as well. Our method achieves 4.5% and 2.9% improvement in TSR and DSR respectively over the previous SOTA in the full-shot setting, and 15.2% and 14.0% in a few-shot setting, showing the effectiveness of a concise data representation in low-resource cross-lingual settings. More importantly, training on translated data and only 10% of original training data comes within 5% of full training, showing building high quality fully functional conversation systems is possible via translation.

We have implemented our methodology as a toolkit for developing multilingual dialogue agents, which will be released open-source upon publication. Our proposed methodology can significantly reduce the cost and time associated with data acquisition for task-oriented dialogue agents.

782

726

727

728

8 Ethical Considerations

674

676

677

678

685

691

697

705

707

708

709

710

711

712

713

714

715

717

718

719 720

721

722

723

724

We do not foresee any harmful or malicious misuses of the technology developed in this work. The data used to train models is about seeking information about domains like restaurants, hotels and tourist attractions, does not contain any offensive content, and is not unfair or biased against any demographic. This work does focus on two widelyspoken languages, English and Chinese, but we think the cross-lingual approach we proposed can improve future dialogue language technologies for a wider range of languages.

We fine-tune multiple medium-sized (several hundred million parameters) neural networks for our experiments. We took several measures to avoid wasted computation, like performing one run instead of averaging multiple runs (since the numerical difference between different models is large enough), and improving batching and representation that improved training speed, and reduced needed GPU time. Please refer to Appendix A.1 for more details about the amount of computation used in this paper.

References

- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Ultes Stefan, Ramadan Osman, 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 (EMNLP).
- Giovanni Campagna, Silei Xu, Mehrad Moradshahi, Richard Socher, and Monica S. Lam. 2019. Genie: A generator of natural language semantic parsers for virtual assistant commands. In *Proceedings of the* 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019, pages 394–410, New York, NY, USA. ACM.
- Wenhu Chen, Jianshu Chen, Pengda Qin, Xifeng Yan, and William Yang Wang. 2019. Semantically conditioned dialog response generation via hierarchical disentangled self-attention. *arXiv preprint arXiv:1905.12866.*
- Wenhu Chen, Jianshu Chen, Yu Su, Xin Wang, Dong Yu, Xifeng Yan, and William Yang Wang. 2018. XL-NBT: A cross-lingual neural belief tracking framework. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 414–424, Brussels, Belgium. Association for Computational Linguistics.
- Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and

Jennimaria Palomaki. 2020. TyDi QA: A benchmark for information-seeking question answering in typologically diverse languages. *Transactions of the Association for Computational Linguistics*, 8:454– 470.

- 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*.
- Bosheng Ding, Junjie Hu, Lidong Bing, Sharifah Mahani Aljunied, Shafiq Joty, Luo Si, and Chunyan Miao. 2021. Globalwoz: Globalizing multiwoz to develop multilingual task-oriented dialogue systems.
- Sauleh Eetemadi and Kristina Toutanova. 2014. Asymmetric features of human generated translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 159–164, Doha, Qatar. Association for Computational Linguistics.
- Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyag Gao, and Dilek Hakkani-Tur. 2019. Multiwoz 2.1: Multi-domain dialogue state corrections and state tracking baselines. *arXiv preprint arXiv:1907.01669*.
- Savoir faire Linux. 2017. num2words. https: //github.com/savoirfairelinux/ num2words.
- Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2020. Language-agnostic bert sentence embedding. *arXiv preprint arXiv:2007.01852*.
- 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.
- Chulaka Gunasekara, Seokhwan Kim, Luis Fernando D'Haro, Abhinav Rastogi, Yun-Nung Chen, Mihail Eric, Behnam Hedayatnia, Karthik Gopalakrishnan, Yang Liu, Chao-Wei Huang, Dilek Hakkani-Tür, Jinchao Li, Qi Zhu, Lingxiao Luo, Lars Liden, Kaili Huang, Shahin Shayandeh, Runze Liang, Baolin Peng, Zheng Zhang, Swadheen Shukla, Minlie Huang, Jianfeng Gao, Shikib Mehri, Yulan Feng, Carla Gordon, Seyed Hossein Alavi, David Traum, Maxine Eskenazi, Ahmad Beirami, Eunjoon, Cho, Paul A. Crook, Ankita De, Alborz Geramifard, Satwik Kottur, Seungwhan Moon, Shivani Poddar, and Rajen Subba. 2020. Overview of the ninth dialog system technology challenge: Dstc9.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network.
- 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*.

889

890

891

892

893

894

839

840

841

842

- 783 784
- 786 787
- 789 790
- 79
- 792 793
- 7
- 796 797
- 79
- 79 80
- 8

8

- 8
- 8
- 808
- 810

811 812

813 814 815

816 817

818 819

820 821

8

- 824 825
- 826 827
- 828 829
- 830 831 832

833

834 835 836

837

837 838 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)*, pages 6505–6520, Online. Association for Computational Linguistics.

- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.
- 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.
- Haoran Li, Abhinav Arora, Shuohui Chen, Anchit Gupta, Sonal Gupta, and Yashar Mehdad. 2021a.
 MTOP: A comprehensive multilingual task-oriented semantic parsing benchmark. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 2950–2962, Online. Association for Computational Linguistics.
- Jinchao Li, Qi Zhu, Lingxiao Luo, Lars Liden, Kaili Huang, Shahin Shayandeh, Runze Liang, Baolin Peng, Zheng Zhang, Swadheen Shukla, Ryuichi Takanobu, Minlie Huang, and Jianfeng Gao. 2021b. Multi-domain task-oriented dialog challenge ii at dstc9. In AAAI-2021 Dialog System Technology Challenge 9 Workshop.
- Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, and Pascale Fung. 2020. MinTL: Minimalist transfer learning for task-oriented dialogue systems. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 3391–3405, Online. Association for Computational Linguistics.
- Zhaojiang Lin, Andrea Madotto, Genta Indra Winata, Peng Xu, Feijun Jiang, Yuxiang Hu, Chen Shi, and Pascale Fung. 2021. BiToD: A bilingual multidomain dataset for task-oriented dialogue modeling. Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1 pre-proceedings (NeurIPS Datasets and Benchmarks 2021).
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation.
- Mehrad Moradshahi, Giovanni Campagna, Sina Semnani, Silei Xu, and Monica Lam. 2020. Localizing

open-ontology QA semantic parsers in a day using machine translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5970–5983, Online. Association for Computational Linguistics.

- Mehrad Moradshahi, Victoria Tsai, Giovanni Campagna, and Monica S Lam. 2021. Contextual semantic parsing for multilingual task-oriented dialogues. *arXiv preprint arXiv:2111.02574*.
- Nikola Mrkšić, Ivan Vulić, Diarmuid Ó Séaghdha, Ira Leviant, Roi Reichart, Milica Gašić, Anna Korhonen, and Steve Young. 2017. Semantic specialization of distributional word vector spaces using monolingual and cross-lingual constraints. *Transactions of the Association for Computational Linguistics*, 5:309–324.
- 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 on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026– 8037.
- Ellie Pavlick, Matt Post, Ann Irvine, Dmitry Kachaev, and Chris Callison-Burch. 2014. The language demographics of amazon mechanical turk. *Transactions of the Association for Computational Linguistics*, 2:79–92.
- Jun Quan, Shian Zhang, Qian Cao, Zizhong Li, and Deyi Xiong. 2020. RiSAWOZ: A large-scale multidomain Wizard-of-Oz dataset with rich semantic annotations for task-oriented dialogue modeling. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 930–940, Online. Association for Computational Linguistics.
- Osman Ramadan, Paweł Budzianowski, and Milica Gasic. 2018. Large-scale multi-domain belief tracking with knowledge sharing. In *Proceedings of the* 56th Annual Meeting of the Association for Computational Linguistics, volume 2, pages 432–437.
- Sebastian Schuster, Sonal Gupta, Rushin Shah, and Mike Lewis. 2019. Cross-lingual transfer learning for multilingual task oriented dialog. 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 3795–3805, Minneapolis, Minnesota. Association for Computational Linguistics.

Scrapinghub. 2015. dateparser. https://github. com/scrapinghub/dateparser.

895 896

897

899

900

901

902

903

904

905

906

907

908

909

910 911

912

913

914

915 916

917

918

919 920

921 922

923

924

925

- Tom Sherborne, Yumo Xu, and Mirella Lapata. 2020. Bootstrapping a crosslingual semantic parser. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 499–517, Online. Association for Computational Linguistics.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Semantically conditioned lstm-based natural language generation for spoken dialogue systems. *arXiv preprint arXiv:1508.01745*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R'emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.
- Qi Zhu, Kaili Huang, Zheng Zhang, Xiaoyan Zhu, and Minlie Huang. 2020. CrossWOZ: A large-scale Chinese cross-domain task-oriented dialogue dataset. *Transactions of the Association for Computational Linguistics*, 8:281–295.

A Appendix

926

931

934

935

938

939

945

947

951

952

953

954

957

959

961

962

963

965

966

967

968

970

971

972

973

974

A.1 Implementation details

Our code is implemented in PyTorch (Paszke et al., 2019) using GenieNLP (Campagna et al., 2019) library for training and evaluation metrics. We additionally use BiToD's code ³ for part of data preprocessing and evaluation. We use pre-trained models available through HuggingFace's Transformers library (Wolf et al., 2019). The following model names are from that library. We use *mbart-large-50* as the neural model for our agent in all our experiments. All models use a standard Seq2Seq architecture with a bidirectial encoder and left-to-right autoregressive decoder. mBART is pre-trained to denoise text in 50 languages, while mT5 is trained on 101 languages. mBART uses sentence-piece (Kudo and Richardson, 2018) for tokenization.

In each setting, all four subtasks of DST, API detection, dialogue act generation, and response generation are done in a single model, where we specify the task by prepending a special token to the input. We found mBART to be especially effective in zero-shot settings as the language of its outputs can be controlled by providing a language-specific token at the beginning of decoding. Additionally, its denoising pre-training objective improves its robustness to the remaining translation noise.

For translation, we use the publicly available *mbart-large-50-many-to-one-mmt* (~611M parameters) model which can directly translate text from any of the 50 supported languages to English. It is an mBART model additionally fine-tuned to do translation.

We use greedy decoding and train our models using teacher-forcing and token-level cross-entropy loss. We used Adam (Kingma and Ba, 2014) as our optimizer with a start learning rate of 2×10^{-5} and linear scheduling. These hyperparameters were chosen based on a very limited hyperparameter search on the validation set. For the numbers reported in the paper, due to cost, we performed only a single run for each experiment.

Our models were trained on virtual machines with a single NVIDIA V100 (16GB memory) GPU on the AWS platform. For a fair comparison, all monolingual models were trained for the same number of iterations of 60K, and bilingual models for 120K. In the few-shot setting, we fine-tuned the model for 3K steps on 1% of the data and 6K steps

³https://github.com/HLTCHKUST/BiToD

on 10% of the data. Sentences are batched based on their input and approximate output token count for better GPU utilization. We set the total number of tokens per batch to 800 for mBART. Due to the verbosity and redundancy of the original BiToD representation, they used a batch size of 1 example for training mbart-large. Using our Distilled representation, however, we can fit up to 6 examples in each batch, and still process each larger batch 3 times faster during training. Training and evaluating each model takes about 10 GPU-hours on average.

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

During error analysis, we noticed that although certain slots (max_temp and min_temp slots in Metro domain, and time and price_range slots in Weather domain) are present in the retrieved knowledge base values, the model does not learn to output them in the agent dialogue act generation subtask. To mitigate this, during evaluation, we automatically check if these slots are present in the input and append them to the generated agent dialogue acts. Additionally, we found that in dialogues that contain the Metro domain, it helps to accumulate API results from the previous turn to make the correct prediction. Thus, in the data prepossessing step, we do keep the API result history for the metro domain.

At inference time, we use the predicted belief state as input to subsequent turns instead of ground truth. However, to avoid the conversation from diverging from its original direction, we use groundtruth agent acts as input for the next turn. Similarly, Lin et al. (2021) use ground-truth natural language agent response as input for the next turn. We made sure the settings are equivalent for a fair comparison.

A.2 Limitations, Risks, and Future Work

As discussed in Section 2.1, organic (i.e. with-1012 out the use of translation) multilingual dialogue 1013 datasets are scarce, which has limited the scope 1014 of our experiments. Our guidelines to improve di-1015 alogue representation mentioned in Section 4 are 1016 general and applicable to any Human-to-Human 1017 or Machine-to-Machine dialogues annotated with 1018 slot-values, but we would have liked to evaluate 1019 the generalization of our cross-lingual approach on 1020 multiple datasets and more languages. For instance, 1021 we partially rely on machine translation models for 1022 Chinese-to-English translation. Available transla-1023 tion models for other language pairs, especially 1024 from/to low-resource languages have much lower quality, and it would be desirable to measure the effect of that in our experiments.

Another limitation is the lack of human evaluation for agent responses. BLEU score does not correlate well with human judgment, and SER only accounts for the factuality of the response but not grammaticality or fluency. This problem is also reported in many prior works (see Section 5). Although finding native speaker evaluators for different languages is a challenge (Pavlick et al., 2014), in future work, we wish to address this by conducting human evaluations.

A.3 Dialogue Examples

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038 1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1050

1051

1052

1053

1054

1055

1056

1057 1058

1059

1060

1062

1063

1064

1065

1067 1068

1069

We include the same example from BiToD's English validation set both in our Distilled representation (Table 3) and in the original (Table 4) representation, along with model predictions in the full-shot setting. For brevity, only the first 3 turns are shown.

In Table 4, we observe that the model fails to ask for the hotel price-range in the second turn and makes an API call instead. Since the API call results are carried over between turns in this representation, in the third turn, the model sees those results in the input and falsely assumes it does not need to make an API call anymore, ultimately resulting in an incorrect response. Compare this to our representation in Table 3. This example shows the importance of separation between API call detection and response generation.

Another phenomenon we often observe is that the model asks for more information than it should according to the gold agent dialogue act. As shown in Table 3, in the second turn, the agent requests user to provide the desired location for the hotel as well as the price range. We believe the main reason for this behavior is the randomness in the agent policy of the BiToD's dialogue simulator. For example, if the agent needs to fill out two slots to make an API call, it can do so by requesting both in the same turn, or one turn at a time. This behavior, though reasonable, is penalized during evaluation, and predictions are considered incorrect if they contain extraneous slots.

5.07		Input	DST: <state> null <endofstate> <history> USER: I'd like hotel recommendations. <endofhistory></endofhistory></history></endofstate></state>
	DST	Target	(hotels search)
		Prediction	(hotels search)
		Input	API: <knowledge> null <endofknowledge> <state> (hotels search) <endofstate> <history> USER: I'd like hotel recommendations. <endofhistory></endofhistory></history></endofstate></state></endofknowledge></knowledge>
	API	Target	no
		Prediction	no
Turn 1		Input	ACTS: <knowledge> null <endofknowledge> <state> (hotels search) <endofstate> <history> USER: I'd like hotel recommendations. <endofhistory></endofhistory></history></endofstate></state></endofknowledge></knowledge>
	DA	Target	(hotels search) request rating , request stars
		Prediction	(hotels search) request rating, request stars
		Input	RG: <actions> (hotels search) request rating , request stars <endofactions> <history> USER: I'd like hotel recommendations. <endofhistory></endofhistory></history></endofactions></actions>
	RG	Target	Certainly. Do you have any requirements for the hotel's rating or the number of stars of the hotel?
		Prediction	Do you have a preference on how many stars and what rating the hotel should have?
	DST	Input	DST: <state> (hotels search) <endofstate> <history> AGENT_ACTS: (hotels search) request rating , request stars USER: The rating doesn't matter, but should be at least 5 stars. <endofhistory></endofhistory></history></endofstate></state>
		Target	(hotels search) rating equal_to " don't care ", stars at_least " 5 "
		Prediction	(hotels search) rating equal_to " don't care ", stars at_least " 5 "
	API Input		API: <knowledge> null <endofknowledge> <state> (hotels search) rating equal_to " don't care ", stars at_least " 5 " <endofstate> <history> AGENT_ACTS: (hotels search) request rating , request stars USER: The rating doesn't matter, but should be at least 5 stars. <endofhistory></endofhistory></history></endofstate></state></endofknowledge></knowledge>
Turn 2		Target	no
		Prediction	no
	DA	Input	ACTS: <knowledge> null <endofknowledge> <state> (hotels search) rating equal_to " don't care ", stars at_least " 5 " <endofstate> <history> AGENT_ACTS: (hotels search) request rating , request stars USER: The rating doesn't matter, but should be at least 5 stars. <endofhistory></endofhistory></history></endofstate></state></endofknowledge></knowledge>
		Target	(hotels search) request price_level
		Prediction	(hotels search) request location , request price_level
	DC	Input	RG: <actions> (hotels search) request price_level <endofactions> <history> USER: The rating doesn't matter, but should be at least 5 stars. <endofhistory></endofhistory></history></endofactions></actions>
	RG	Target	Do you have a price range for the hotel?
		Prediction	And what about location? Do you have a price range for the hotel?

	DST Input DST: <state> (hotels s <endofstate> <history> A stars AGENT_ACTS: (h</history></endofstate></state>		DST: <state> (hotels search) rating equal_to " don't care " , stars at_least " 5 " <endofstate> <history> AGENT_ACTS_PREV: (hotels search) request rating , request stars AGENT_ACTS: (hotels search) request price_level USER: cheap <endofhistory></endofhistory></history></endofstate></state>
		Target	(hotels search) price_level equal_to " cheap "
	Prediction		(hotels search) price_level equal_to " cheap "
	API	Input	API: <knowledge> null <endofknowledge> <state> (hotels search) price_level equal_to " cheap " , rating equal_to " don't care " , stars at_least " 5 " <endofs- tate> <history> AGENT_ACTS_PREV: (hotels search) request rating , request stars AGENT_ACTS: (hotels search) request price_level USER: cheap <endofhistory></endofhistory></history></endofs- </state></endofknowledge></knowledge>
		Target	yes
		Prediction	yes
Turn 3	DA	Input	ACTS: <knowledge> (hotels search) available_options " 4 " , location " Mong Kok Kowloon Yau Tsim Mong District " , name " Royal Plaza Hotel " , price_level " cheap " , price_per_night " 793 HKD " , rating " 9 " , stars " 5 " <endofknowledge> <state> (hotels search) price_level equal_to " cheap " , rating equal_to " don't care " , stars at_least " 5 " <endofstate> <history> AGENT_ACTS_PREV: (hotels search) request rating , request stars AGENT_ACTS: (hotels search) request price_level USER: cheap <endofhistory></endofhistory></history></endofstate></state></endofknowledge></knowledge>
	Target (hotels search) offer av Plaza Hotel ", offer ratin		(hotels search) offer available_options equal_to " 4 " , offer name equal_to " Royal Plaza Hotel " , offer rating equal_to " 9 "
		Prediction	(hotels search) offer available_options equal_to " 4 " , offer name equal_to " Royal Plaza Hotel " , offer rating equal_to " 9 "
	RG Input RG: <actions> (hotels search) offer available_options e equal_to " Royal Plaza Hotel ", offer rating equal_to " 9 ' USER: cheap <endofhistory> Target Okay. There are 4 hotels available. I recommend the Royal rating.</endofhistory></actions>		RG: <actions> (hotels search) offer available_options equal_to " 4 " , offer name equal_to " Royal Plaza Hotel " , offer rating equal_to " 9 " <endofactions> <history> USER: cheap <endofhistory></endofhistory></history></endofactions></actions>
			Okay. There are 4 hotels available. I recommend the Royal Plaza Hotel, which has a 9 rating.
		Prediction	There are 4 available hotels. I recommend Royal Plaza Hotel. Its rating is 9.

Table 3: An example from BiToD English validation set in Distilled representation, along with our mBART model predictions. For brevity, only the first 3 turns are shown.

		Input	Track Dialogue State: <knowledge><dialogue_state> <user> I'd like hotel</user></dialogue_state></knowledge>	
	DST	Torgot	recommendations.	
	_ ~ _	Dradiation	<api> hotels search</api>	
Turn 1		Flediction	<api> notels search</api>	
	A DI / D osponso	Input	search <user> I'd like hotel recommendations.</user>	
	Ar I/ Kesponse	Target	Certainly. Do you have any requirements for the hotel's rating or the number of stars of the hotel?	
		Prediction	What rating would you like the hotel to have?	
	DST	Input	Track Dialogue State: <knowledge><dialogue_state> <api> hotels search<user> I'd like hotel recommendations.<system> Certainly. Do you have any requirements for the hotel's rating or the number of stars of the hotel?<user> The rating doesn't matter, but should be at least 5 stars.</user></system></user></api></dialogue_state></knowledge>	
		Target	<api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5</value></relation></slot></value></relation></slot></api>	
Turn 2		Prediction	<api> hotels search<slot> rating<relation> equal_to<value> do care<slot> stars<relation> at_least<value> 5</value></relation></slot></value></relation></slot></api>	
	API/ Response	Input	Generate Response: <knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<user> I'd like hotel recommenda- tions.<system> Certainly. Do you have any requirements for the hotel's rating or the number of stars of the hotel?<user> The rating doesn't matter, but should be at least 5 stars.</user></system></user></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge>	
		Target	Do you have a price range for the hotel?	
		Prediction	I found 43 hotels. I recommend Cordis, Hong Kong, which has a rating of 10.	
	DST	Input	Track Dialogue State: <knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<user> The rating doesn't matter, but should be at least 5 stars.<system> Do you have a price range for the hotel?<user> chean</user></system></user></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge>	
			noter (user) encap	
		Target	<pre><api> hotels search<slot> price_level<relation> equal_to<value> cheap</value></relation></slot></api></pre>	
		Target Prediction	<pre><api> hotels search<slot> price_level<relation> equal_to<value> cheap <api> hotels search<slot> price_level<relation> equal_to<value> cheap</value></relation></slot></api></value></relation></slot></api></pre>	
Turn 3	API / Response	Target Prediction Input	<pre><api> hotels search<slot> price_level<relation> equal_to<value> cheap <api> hotels search<slot> price_level<relation> equal_to<value> cheap Generate Response:<knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<slot> price_level<relation> equal_to<value> cheap<user> The rating doesn't matter, but should be at least 5 stars.<system> Do you have a price range for the hotel?<user> cheap</user></system></user></value></relation></slot></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge></value></relation></slot></api></value></relation></slot></api></pre>	
Turn 3	API / Response	Target Prediction Input Target	<pre></pre> <pre><pre><pre><pre><pre>API> hotels search<slot> price_level<relation> equal_to<value> cheap</value></relation></slot></pre><pre><pre><pre>API> hotels search<slot> price_level<relation> equal_to<value> cheap</value></relation></slot></pre><pre>Generate Response:<knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<slot> price_level<relation> equal_to<value> cheap</value></relation></slot></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge></pre><pre>cheap</pre><pre><pre><pre><pre><pre><pre><pre><</pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre>	
Turn 3	API / Response	Target Prediction Input Target Prediction	<pre>API> hotels search<slot> price_level<relation> equal_to<value> cheap <api> hotels search<slot> price_level<relation> equal_to<value> cheap Generate Response:<knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<slot> price_level<relation> equal_to<value> cheap<user> The rating doesn't matter, but should be at least 5 stars.<system> Do you have a price range for the hotel?<user> cheap </user></system></user></value></relation></slot></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge></value></relation></slot></api></value></relation></slot></pre>	
Turn 3	API/ Response API/ Response	Target Prediction Input Target Prediction Input	<pre>API> hotels search<slot> price_level<relation> equal_to<value> cheap </value></relation></slot></pre> <pre> <api> hotels search<slot> price_level<relation> equal_to<value> cheap Generate Response:<knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<slot> price_level<relation> equal_to<value> cheap<user> The rating doesn't matter, but should be at least 5 stars.<system> Do you have a price range for the hotel?<user> cheap </user></system></user></value></relation></slot></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge></value></relation></slot></api></pre> <pre> </pre> <	
Turn 3	API/ Response API/ Response	Target Prediction Input Target Prediction Input	<pre>API> hotels search<slot> price_level<relation> equal_to<value> cheap </value></relation></slot></pre> <pre> <api> hotels search<slot> price_level<relation> equal_to<value> cheap Generate Response:<knowledge><dialogue_state> <api> hotels search<slot> rating<relation> equal_to<value> don't care<slot> stars<relation> at_least<value> 5<slot> price_level<relation> equal_to<value> cheap<user> The rating doesn't matter, but should be at least 5 stars.<system> Do you have a price range for the hotel?<user> cheap </user></system></user></value></relation></slot></value></relation></slot></value></relation></slot></api></dialogue_state></knowledge></value></relation></slot></api></pre> <pre> </pre> <pr< td=""></pr<>	

Table 4: Same example as in Table 3 but in the original representation from Lin et al. (2021), along with MinTL(mT5) model predictions. For brevity, only the first 3 turns are shown.