# **Extracting and Inferring Personal Attributes from Dialogue**

**Anonymous ACL submission** 

### Abstract

Personal attributes represent structured information about a person, such as their hobbies, pets, family, likes and dislikes. We introduce the tasks of extracting and inferring personal attributes from human-human dialogue, and analyze the linguistic demands of these tasks. To meet these challenges, we introduce a simple and extensible model that combines an autoregressive language model utilizing constrained attribute generation with a discriminative reranker. Our model outperforms strong baselines on extracting personal attributes as well as inferring personal attributes that are not contained verbatim in utterances and instead requires commonsense reasoning and lexical inferences, which occur frequently in everyday conversation. Finally, we demonstrate the benefit of incorporating personal attributes in social chit-chat and task-oriented dialogue settings.

## 1 Introduction

001

016

017

034

040

041

Personal attributes are structured information about a person, such as what they like, what they have, and what their favorite things are. These attributes are commonly revealed either explicitly or implicitly during social dialogue as shown in Figure 1, allowing people to know more about one another. These personal attributes, represented in the form of knowledge graph triples (*e.g.* I, has\_hobby, volunteer) can represent large numbers of personal attributes in an interpretable manner, facilitating its usage by weakly-coupled downstream dialogue tasks (Li et al., 2014; Qian et al., 2018; Zheng et al., 2020a,b; Hogan et al., 2021).

One such task is to ground open-domain chitchat dialogue agents to minimize inconsistencies in their language use (*e.g.*, **I like cabbage**  $\rightarrow$ (next turn)  $\rightarrow$ **Cabbage is disgusting**) and make them engaging to talk with (Li et al., 2016; Zhang et al., 2018; Mazaré et al., 2018; Qian et al., 2018; Zheng et al., 2020a,b; Li et al., 2020; Majumder et al.,



Figure 1: Overview of obtaining personal attribute triple from utterances using our model GenRe. Attribute values are contained within the utterance in the EXTRACTION task, but not the INFERENCE task.

2020). Thus far, personalization in chit-chat has made use of dense embeddings and natural language sentences. While KG triples have been shown to be capable of grounding Natural Language Generation (Moon et al., 2019; Koncel-Kedziorski et al., 2019), they have yet to be used to personalize chit-chat dialogue agents.

Personal attributes can also help task-oriented dialogue agents to provide personalized recommendations (Mo et al., 2017; Joshi et al., 2017; Luo et al., 2019; Lu et al., 2019; Pei et al., 2021). Such personalized recommendations have only been attempted for single-domain tasks with a small set of one-hot features (< 30). Personalization across a wide range of tasks (recommending food, movies and music by multi-task dialogue agents such as Alexa, Siri and Assistant) however can require orders of magnitude more personal attribute features. This makes KG triples ideal for representing them, given the advantages of this data structure in selecting and utilizing pertinent features (Li et al., 2014;

#### Hogan et al., 2021).

063

064

065

067

086

090

097

100

101

102

103

104

105

Based on these advantages, we investigate how personal attributes can be predicted from dialogue. An important bottleneck for this step lies in the poor coverage of relevant personal attributes in existing labeled datasets. Therefore, we introduce two new tasks for identifying personal attributes in Section 2. As shown in Figure 1, the EXTRACTION task requires determining which phrase in an utterance indicate a personal attribute, while the INFERENCE task adds further challenge by requiring models to predict personal attributes that are not explicitly stated verbatim in utterances. This is common in conversational settings, where people express personal attributes using a variety of semantically related words or imply them using commonsense reasoning. We analyze how these factors allow personal attributes to be linked to utterances that express them.

To tackle these tasks, we propose a simple yet extensible model, **GenRe**, in Section 3. GenRe combines a constrained attribute generation model (that is flexible to accommodate attributes not found verbatim in utterances) with a discriminative reranker (that can contrast between highly similar candidates). Our experiments in Section 4 suggests that such design allows our model to outperform strong baseline models on both the EXTRACTION and INFERENCE tasks. Subsequently in Section 5, detailed ablation studies demonstrate the value of our model components while further analysis identifies future areas for improvement.

Finally in Section 6, we show how personal attributes in the form of KG triples can improve the personalization of open-domain social chit-chat agents as well as task-oriented dialogue agents. In the former case, personal attributes can be utilized to improve chat-bot consistency on the PersonaChat task (Zhang et al., 2018). In the latter case, we suggest how our personal attributes can support personalization in multi-task, task-oriented dialogue settings.

## 2 Personal Attribute Tasks

Having motivated the importance of personal attributes, we propose the task of obtaining them
from natural language sentences. We first explain
how we formulate two complementary tasks from
DialogNLI data and then formally define our tasks.
Finally, we analyze the task datasets to gather insights into the linguistic phenomena that our tasks

involve.

#### 2.1 Source of Personal Attributes

DialogNLI (Welleck et al., 2019) contains samples of PersonaChat utterances (Zhang et al., 2018), each paired with a manually annotated personal attribute triple. Each triple consists of a head entity, a relation, and a tail entity. These triples were initially annotated to identify entailing, contradicting and neutral statements within the PersonaChat corpus. For instance, a statement labelled with (I, [has\_profession], chef) will contradict with another statement labelled with (I, [has\_profession], engineer). The three largest groups of relations are: a.  $has_X$  (where X = hobby, vehicle, pet) b. favourite\_Y (where Y = activity, color, music) c. like\_Z (where Z = read, drink, movie). 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

162

#### 2.2 Extraction and Inference Tasks

By re-purposing the DialogNLI dataset, our tasks seek to extract these personal attribute triples from their paired utterances. We first used a script that obtains pairs of personal triples and utterances. Next, we combined relations with similar meanings such as like\_food and favourite\_food and removed under-specified relations such as favourite, have and others. Finally, we removed invalid samples with triples containing *None* or < blank > and removed prefix numbers of tail entities (*e.g.* 11 dogs), since the quantity is not important for our investigation.

We formulate two tasks by partitioning the DialogNLI dataset into two non-overlapping subsets. Here, each **sample** refers to a sentence paired with an annotated triple. Train/dev/test splits follow DialogNLI, with descriptive statistics shown in Table 1. The dataset for the EXTRACTION task contains samples in which both the head and tail entities are spans inside the paired sentence. An example is (I, [has\_profession], receptionist) from the sentence "I work as a receptionist in my day job". We formulate the EXTRACTION task in a similar way to existing Relation Extraction tasks such as ACE05 (Wadden et al., 2019) and NYT24 (Nayak and Ng, 2020). This allows us to apply modeling lessons learned from Relation Extraction.

The complementary set is the dataset for the IN-FERENCE task, for which the head entity and/or the tail entity cannot be found as spans within the paired sentence. This is important in real-world conversations because people do not always express their personal attributes explicitly and instead use paraphrasing and commonsense reasoning to do so. An example of a paraphrased triple is (I, [physical\_attribute], tall) from the sentence "I am in the 99th height percentile", while one based on commonsense reasoning is (I, [want\_job], umpire) from the sentence "my ultimate goal would be calling a ball game".

163

164

165

166

167

168

169

	EXTRACTION	INFERENCE
Samples		
train	22911	25328
dev.	2676	2658
test	2746	2452
Unique eleme	nts	
head entities	88	109
relations	39	39
tail entities	2381	2522
Avg. words		
head entities	1.03	1.08
relations	1.00	1.00
tail entities	1.20	1.28
sentences	12.9	12.2

Table 1: Statistics of the dataset for the two tasks.

The INFERENCE task is posed as a challeng-170 ing version of the EXTRACTION task that tests 171 models' ability to identify pertinent information 172 in sentences and then make commonsense infer-173 ences/paraphrases based on such information. An 174 existing task has sought to predict personal at-175 tributes that are not always explicitly found within 176 sentences (Wu et al., 2019). However, it did not 177 178 distinguish between personal attributes that can be explicitly found within sentences (i.e. EXTRAC-179 TION) from those that cannot (*i.e.* INFERENCE). 180 We believe that, given that the inherent difficulty of 181 identifying the two types of personal attributes are greatly different, it is helpful to pose them as two separate tasks. In this way, the research commu-184 nity can first aim for an adequate performance on the simpler task before applying lessons to make 186 progress at the more challenging task. This is also 187 the first time that personal attributes that are not explicitly contained in sentences are shown to be 189 derivable from words in the sentence using com-190 monsense/lexical inferences. 191

#### 2.3 Formal Task Definition

192

Given a sentence S, we want to obtain a personalattribute triple in the form of (head entity, **relation, tail entity**). The relation must belong to a set of 39 predefined relations. In the EXTRACTION subset, the head entity and tail entity are spans within S. Conversely, in the INFERENCE subset, the head entity and/or the tail entity cannot be found as spans within S. 195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

### 2.4 Dataset Analysis

We analyze the datasets to obtain insights into how the tasks can be approached. Because the majority of head entities (93.3%) are simply the word "I", our analysis will focus on tail entities.



Figure 2: Bar plot for 10 most common dependency labels of tail entities

**Dataset for the EXTRACTION task** We use dependency parses of sentences to understand the relationship between words within tail entities and the sentence ROOT. Dependency parsing was chosen because it is a well-studied syntactic task (Nivre et al., 2016) and has been shown to contribute to the relation extraction task (Zhang et al., 2017). Dependency parses and labels associated with each dependent word were identified using a pre-trained transformer model from spaCy.<sup>1</sup> The parser was trained on data annotated with the ClearNLP dependency schema that is similar to Universal Dependencies (Nivre et al., 2016).<sup>2</sup>

As shown in Figure 2, objects of prepositions (pobj) and direct objects (dobj) each comprise of 17.5% of tail entities, followed by compound words (compound), attributes (attr) and adjectival complements (acomp), plus 138 other long-tail labels. The range of grammatical roles as well as the fact

<sup>2</sup>https://github.com/clir/clearnlp-

<sup>&</sup>lt;sup>1</sup>https://spacy.io/

guidelines/blob/master/md/specifications/dependency\_labels.md

225that one third of tail entities do not involve nouns226(see Figure in Section A.2) also suggest that the227tail entities in our dataset go beyond proper nouns,228which are what many Relation Extraction datasets229(e.g., ACE05 and NYT24) are mainly concerned230with. Such diversity in grammatical roles played by231tail entities means that approaches based on rule-232based extraction, parsing or named entity recog-233nition alone are unlikely to be successful in the234EXTRACTION task.

Dataset for the INFERENCE task A qualitative inspection of the dataset showed that inferences can be made on the basis of semantically-related words and commonsense inferences, as shown in examples discussed in Section 2.2. To better understand 239 how tail entities can be inferred from the sentence 240 in the INFERENCE subset, we analyze the relation-241 ship between words in the tail entity and words 242 in the sentence. 79.2% of tail entities cannot be 243 directly identified in the sentence. We performed 244 a few transformations to identify potential links between the tail entity and the sentence. Concept-Net\_connect refers to words with highest-weighted 247 edges on ConceptNet to sentence words while ConceptNet\_related refers to words that have closest embedding distances to sentence words. Details of their preparation are in Appendix A.3. As in Table 251 2, our analysis shows that a model that can perform well on the INFERENCE task requiring both Word-Net semantic knowledge (Fellbaum, 1998) as well as ConceptNet commonsense knowledge (Speer et al., 2017). 256

Transformation	Example	%
	(sentence→tail entity	y)
ConceptNet_related	mother $\rightarrow$ female	71.3
ConceptNet_connect	wife $\rightarrow$ married	56.8
WordNet_synonym	outside $\rightarrow$ outdoors	39.5
WordNet_hypernym	drum $\rightarrow$ instrument	5.04
WordNet_hyponym	felines $\rightarrow$ cats	4.17
Same_stem	swimming $\rightarrow$ swim	43.3

Table 2: Proportion (%) of tail entities that can be related to sentence words after applying each transformation.

## 3 GenRe

257

This section proposes GenRe, a model that uses a unified architecture for both the EXTRACTION and the INFERENCE tasks. We use a simple and exten-

sible generator-reranker framework to address the needs of the two tasks. On one hand, a generative model is necessary because head and/or tail entities cannot be directly extracted from the sentence for INFERENCE dataset. On the other hand, preliminary experiments using a Generator in isolation showed that a large proportion of correct triples are among the top-k - but not top-1 - outputs (see Table **??**). A Reranker can be used to select the most likely triple among the top-k candidate triples, leading to a large improvement in performance (see Table **4**). 261

262

263

265

266

267

269

270

271

272

273

274

275

276

277

278

279

280

281

283

287

289

290

291

292

293

294

295

296

298

299

300

301

302

303

304

306

307

308

309

310

#### 3.1 Generator

We use an autoregressive language model (GPT-2 small) as our Generator because its extensive pretraining is useful in generating syntactically and semantically coherent entities. The small model was chosen to keep model size similar to baselines. We finetune this model to predict a personal attribute triple occurring in a given input sentence. Specifically, we treat the flattened triples as targets to be predicted using the original sentence as context. The triple is formatted with control tokens to distinguish the head entity, relation, and tail entity as follows:

 $\mathbf{y} = [\text{HEAD}], t_{1:m}^{head}, [\text{RELN}], t^{reln}, [\text{TAIL}], t_{1:k}^{tail}$ where {[HEAD],[RELN], [TAIL]} are control tokens,  $t_{1:m}^{head}$  is the head entity (a sequence of length m),  $t^{reln}$  is a relation, and  $t_{1:k}^{tail}$  is the tail entity.

During evaluation, we are given a sentence as context and seek to generate a personal attribute triple in the flattened format as above. To reduce the search space, we adopt a constrained generation approach. Specifically, after the [RELN] token, only one of 39 predefined relations can be generated, and so the output probability of all other tokens is set to 0. After the [TAIL] token, all output tokens not appearing in the input sentence will have zeroed probabilities in the EXTRACTION task. Conversely for the INFERENCE task, the only allowed output tokens after the [TAIL] token are those which have appeared following the predicted relation in the training data. For example, tail entities that can be generated with a [physical\_attribute] relation include "short", "skinny" or "wears glasses", as these examples occur in the training data. We imposed this restriction to prevent the model from hallucinating attributes that are not associated to the predicted relation (such as "dog" with [physical\_attribute]). Despite limiting the model's ability

405

406

357

311to generate novel but compatible tail entities (and312thereby upper-bounding maximum possible recall313to 75.7%), this approach in balance helped to im-314prove model performance. Implementation details315are in Appendix A.4.

### 3.2 Reranker

317

319

320

322

324

325

326

327

328

329

330

333

335

336

339

340

341

342

343

344

345

We use BERT-base as the Reranker because its bidirectionality allows tail tokens to influence the choice of relation tokens. Furthermore, BERT has demonstrated the best commonsense understanding among pre-trained language models (Petroni et al., 2019; Zhou et al., 2020). These benefits have led to many relation extraction models using BERT as part of the pipeline (Wadden et al., 2019; Yu et al., 2020; Ye et al., 2021).

For each S, we obtain the L most likely sequences using the Generator, including the context sentence. Each sequence is labelled as correct or incorrect based on whether the predicted triple (head entity, relation, tail entity) matches exactly the ground-truth triple. Incorrect sequences serve as challenging negative samples for the Reranker because they are extremely similar to the correct sequence since they were generated together. We finetune a BERT model with a binary cross-entropy loss function to classify whether sequences are correct. During inference, we select the sequence with the highest likelihood of being correct as our predicted sequence. We set L to 10 in all experiments. Implementation details are in Appendix A.5.

### 4 Experiments

We first explain the metrics used in the experiments. Next, we introduce the baseline models. Finally, we examine how GenRe compares to baseline models to understand its advantages.

## 4.1 Metrics

Micro-averaged Precision/Recall/F1 were calculated following Nayak and Ng (2020), in which a sample is considered correct only when all three elements (head\_entity, relation and tail entity) are resolved correctly. We chose these metrics because we are interested in the proportion of all predicted personal attributes that have been correctly identified (precision) and of all ground truth personal attributes (recall). F1 is considered as an aggregate metric for precision and recall.

#### 4.2 Baseline Models

**Generative models** can be used for both the EX-TRACTION and the INFERENCE tasks.

**WDec** is an encoder-decoder model that achieved state-of-the-art performance in the NYT24 and NYT29 tasks (Nayak and Ng, 2020). The encoder is a Bi-LSTM, while the decoder is an LSTM with attention over encoder states. An optional copy mechanism can be used: when used, the decoder will only generate tokens found in the original sentence. The copy mechanism was used on the EXTRACTION dataset but not on the INFER-ENCE dataset (given their better empirical performance).

**GPT2** is an autoregressive language model that we build GenRe on. We use the same configuration as in GenRe.

**Extractive models** can be used only for the EX-TRACTION task, because they select for head and tail entities from the original sentence.

**DyGIE++** is a RoBERTa-based model that achieved state-of-the-art performance in multiple relation extraction tasks including ACE05 (Wadden et al., 2019). It first extracts spans within the original sentence as head and tail entities. Then, it pairs up these entities with a relation and passes them through a graph neural network, with the head and tail entities as the nodes, and relations as the edges. This allows information flow between related entities before passing the triple through a classifier.

**PNDec** is an Encoder-Decoder model that achieved close to SOTA performance in NYT24 and NYT29 (Nayak and Ng, 2020). It uses the same encoder as WDec but uses a pointer network to identify head and tail entities from the original sentence, which it pairs with possible relation tokens to form a triple that is subsequently classified.

All baseline models were trained on our datasets using their suggested hyper-parameters.

#### 4.3 Model Results

The top-performing baseline models on the EX-TRACTION dataset are the extractive models, which select spans within the sentence and classify whether an entire triple is likely to be correct. Because there are only a small number of spans within the sentence, this approach can effectively limit its search space. On the other hand, extractive models cannot solve the INFERENCE task, because the underlying assumption that head and tail entities must

	Ex	TRACT	ION	IN	FEREN	CE
	Р	R	F1	Р	R	F1
GenRe	68.0	52.4	59.2	46.5	35.4	39.2
Generative						
WDec	57.0	49.0	52.7	33.6	34.7	34.1
GPT2	50.9	31.1	38.6	31.3	17.3	22.3
Extractive						
DyGIE++	60.8	50.9	55.3			
PNDec	63.1	49.5	55.5			

Table 3: Performance on the test set. GenRe has significantly higher F1 than all baseline models with 5 runs based on a two-tailed t-test (p < 0.05).

be found within the sentence does not hold. Conversely, generative models perform more poorly on the Extraction task but are capable on the INFER-ENCE task. This is because generation happens in a left-to-right manner, meaning that some elements of the triple have to be generated without knowing what the other elements are. Our approach of combining a Generative model with a BERT-base Reranker that is akin to models used by Extractive approaches combines the best of both worlds. Not only does it perform well on the Extraction task ( $\geq$  3.7 F1 points over baselines), it also excels on the Inference task ( $\geq$  5.1 F1 points over baselines).

## 5 Analysis

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

We first conduct an ablation study to better understand the contribution of constrained generation and the Reranker, by measuring the performance of our model when each component is removed. Then, we seek to understand how errors are made on predicted personal attribute relations to identify future areas of improvement.

#### 5.1 Ablation Study

Table 4 shows that both the Reranker and constrained generation contribute to the performance of GenRe. In particular, the constrained generation plays a larger role on the EXTRACTION dataset while the Reranker plays a greater role on the IN-FERENCE dataset.

**Constrained generation** has a large impact on the EXTRACTION dataset (+13.0% F1), likely because it very much restricts the generation search space to spans from the context sentence. On the INFERENCE dataset, the original search space cannot be effectively limited to tokens in the context sentence. Therefore, applying the heuristic that

	EXTRACTION			INFERENCE		
	Р	R	F1	Р	R	F1
GenRe	68.0	52.4	59.2	46.5	35.4	39.2
- Constr. Gen	53.5	40.7	46.2	37.2	27.1	31.4
- Reranker	67.6	41.0	51.0	31.0	22.3	25.9

Table 4: Ablation study for Reranker and constrained generation.

only tail entities associated with a particular relation (in the training set) can be decoded is useful, even though it upper bounds maximum recall to 75.7%, which is much higher than the achieved 35.4%. Compared to the EXTRACTION dataset, the improvement on the INFERENCE dataset is smaller (+7.8% F1), since the range of tail entities that can be decoded after imposing the constraint is greater. 442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

**The Reranker** is needed because, many times, the correct triple can be generated by the Generator but might not be the triple that is predicted to have the highest likelihood. The maximum possible recall on the EXTRACTION and INFERENCE tasks increases from 41.0% to 59.9% and 22.3% to 41.0% respectively when considering top-10 instead of only top-1 generated candidate. While the achieved recall (52.4% and 35.4% respectively) are still a distance from the maximum possible recall, the achieved recall is much higher than using the Generator alone.

### 5.2 Misclassification of Relations

Major sources of error on the EXTRACTION dataset 463 came from relation tokens that have close se-464 mantic meanings. They either were related to 465 one another (*e.g.*, [has\_profession] vs [want\_job]) 466 or could be correlated with one another (e.g., 467 [like\_animal] vs [have\_pet] or [like\_music] vs [fa-468 vorite\_music\_artist]), as illustrated in Table 5. 469 Such errors likely arose due to the way that the 470 DialogNLI dataset (Welleck et al., 2019) was anno-471 tated. Specifically, annotators were asked to label a 472 single possible triple given a sentence instead of all 473 applicable triples. Because of this, our evaluation 474 metrics are likely to over-penalize models when 475 they generate reasonable triples that did not match 476 the ground truth. Future work can avoid this prob-477 lem by labelling all possible triples and framing the 478 task as multilabel learning. 479

					Top 3 Most Frequent (n)			
Dataset	True Relation (n)	Р	R	F1	Predicted Relations	True Tail Entities	Predicted Tail Entities	
EXTRACTION	[has_profession] (274)	83.8	62.0	71.3	[has_profession] (189) [employed_by_general] (30) [want_job] (17)	teacher (29) nurse (28) real estate agent (25)	nurse (27) real estate (25) teacher (19)	
	[have_pet] (149)	97.3	55.0	70.3	[have_pet] (88) [have_family] (18) [like_animal] (12)	dog (55) cat (45) pets (22)	cat (32) pets (23) dog (18)	
INFERENCE	[like_food] (77)	46.7	41.6	44.0	[like_food] (62) [like_activity] (5) [like_animal] (4)	pizza (18) onion (9) italian (7)	pizza (19) italian cuisine (10) onion (8)	
	[like_music] (71)	40.8	23.9	30.2	[like_music] (40) [favorite_music_artist] (9) [like_activity] (7)	jazz (10) country (9) rap (6)	the story so far (12) country (8) jazz (7)	

Table 5: Some common relations in EXTRACTION and INFERENCE datasets

## 6 Applications of Personal Attributes

480

481

482

483

484

485

486 487

488 489

490

491

492

493

494

495

496

497

498

499

500

501 502

503

505

506

508

511

512

513

514

515

Personal attributes can make social chit-chat agents more consistent and engaging as well as enable task-oriented agents to make personalized recommendations. In this section, we use personal attributes to improve chit-chat agent consistency and provide information for personalizing task-oriented dialogue agents.

#### 6.1 Consistency in Chit-chat agents

PersonaChat (Zhang et al., 2018) was created to improve the personality consistency of open-domain chit-chat dialogue agents. PersonaChat was constructed by giving pairs of crowdworkers a set of English personal attribute related sentences and asking them to chat in a way that is congruent with those sentences. Models were then trained to generate dialogue responses that are in line with those expressed by crowdworkers using the provided persona information as context.

**Methods** We fine-tune the generative version of Blender 90M (a transformer-based model trained on multiple related tasks) on PersonaChat, which is currently the state-of-the-art generative model on this task (Roller et al., 2020). We prepend a corresponding DialogNLI personal attribute before each utterance (*i.e.* **+Per. Attr.**), in order to better direct the model in generating a suitable response that is consistent with the set persona. This modification is relatively minimal to demonstrate the informativeness of personal attributes, while keeping the model architecture and hyperparameter fine-tuning the same as in the original work (details in Appendix A.1).

Metrics We follow Roller et al. (2020) and Dinan et al. (2019). Metrics for **+Per. Attr.** setting consider both personal attributes and utterances. **Hits@1** uses the hidden states of the generated output to select the most likely utterance amongst 20 candidates (the correct utterance and 19 randomly chosen utterances from the corpus). **Perplexity** reflects the quality of the trained language model. **F1** demonstrates the extent of the overlap between the generated sequence and the ground truth sequence.

	Hits@1↑	Perplexity $\downarrow$	F1 $\uparrow$
Blender	32.3	11.3	20.4
+ Per. Attr.	35.2*	10.4*	20.6*

Table 6: Effects of using personal attributes to augment Blender on Personachat. Higher is better for Hits@1 and F1; lower is better for perplexity. \*Significantly different from Blender with 5 runs based on a two-tailed t-test (p < 0.05).

Fact 1	I love cats and have two cats
Fact 2	I've a hat collection of over 1000 hats.
Blender	My cats names are all the hats i have
+ Per. Attr.	My cats are called kitties
Fact 1	I am a doctor.
Fact 2	My daughter is a child prodigy.
Blender	My daughter is prodigy so she gets a lot of accidents.
+ Per. Attr.	I've seen a lot of accidents.

Table 7: Examples of incorrect utterances generated by Blender by mixing up two facts.

**Results** As shown in Table 6, including personal attributes can improve performance on the PersonaChat task. An inspection of the generated utterances suggests that including personal attributes into Blender can more effectively inform the model which persona statement to focus on during generation. This can prevent Blender from including information in irrelevant persona statements (*e.g.* by mixing up facts from two unrelated persona

523

524

525

526

527

528

529

530

532

535

536

537

541

542

543

544

545

547

549

553

554

555

557

561

565

566

567

statements), as in Table 7.

### 6.2 Personalization in Task-oriented dialogue

While personalization has been incorporated into single-task settings (Joshi et al., 2017; Mo et al., 2017; Luo et al., 2019; Lu et al., 2019; Pei et al., 2021), there has been no attempt for personalization in multi-task settings. This is against the background in which multi-task dialogue is rapidly becoming the standard in task-oriented dialogue evaluation (Byrne et al., 2019; Rastogi et al., 2019; Zang et al., 2020; Shalyminov et al., 2020). To overcome this gap, we show how our dataset can lay a foundational building block for personalization in multi-task dialogue.

**Methods** We used several popular datasets on multi-task task-oriented dialogue (Zang et al., 2020; Shalyminov et al., 2020; Byrne et al., 2019; Rastogi et al., 2019). From each dataset, we manually observed its tasks and categorized them into several overarching domains, as shown in Table 8. We then created a mapping between the various domains and datasets available for personalizing task-oriented dialogue (including ours). Domains that are not supported by any dataset are omitted.

**Results** Compared to existing datasets in Table 8, our dataset is capable of personalizing recommendations in a much larger number of domains. These domains include restaurants and shopping, which have been explored by existing datasets, as well as movies, music, sports and recreation, which have thus far been overlooked. For domains that have been previously explored, such as restaurants, our dataset contains a more diverse set of possible personal attribute values (*e.g.* the foods people like), which can support it to personalize recommendations in more realistic manners.

Dataset	Domains	#Unique features
Ours	Restaurants, Movies,	5583
	Music, Sports,	
	Recreation, Shopping	
Ours	Restaurants only	206
Joshi et al. (2017)	Restaurants	30
Mo et al. (2017)	Restaurants	10
Lu et al. (2019)	Shopping	7

Table 8: Domains covered by various datasets for personalizing task-oriented dialogue. #Uniques features refers to the number of unique attribute-values (*e.g.* the specific food people like) that can be used for personalization.

## 7 Related Work

Personal Attribute Extraction: Most work on extracting personal attributes from natural language (Pappu and Rudnicky, 2014; Mazaré et al., 2018; Wu et al., 2019; Tigunova et al., 2019, 2020) employed distant supervision approaches using heuristics and hand-crafted templates, which have poor recall. In contrast, we use a strong supervision approach in which triples were manually annotated. Li et al. (2014) and Yu et al. (2020) attempted to extract personal information from dialogue using a strongly supervised paradigm. However, they focused on demographic attributes as well as interpersonal relationships, which contrast with our focus on what people own and like. Li et al. (2014) used SVMs to classify relations and CRFs to perform slot filling of entities while Yu et al. (2020) used BERT to identify relations between given entities. Generating KG triple using Language Models: Autoregressive language models have been applied to a wide range of tasks involving the generation of data with similar structures as personal attribute KG triples, including dialogue state tracking (Hosseini-Asl et al., 2020) and commonsense KG completion (Bosselut et al., 2019). The most similar application is Alt et al. (2019), which used the original GPT model (Radford and Narasimhan, 2018) for relation classification. Their task formulation involves identifying a specific relation (out of around 30 possible options) for two given entities. On the other hand, our tasks seek to identify not only the relation, but also the head and tail entities, which have potentially open vocabulary requirements, which makes it much harder.

568

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

### 8 Conclusion

In conclusion, we propose the novel tasks of extracting and inferring personal attributes from dialogue and carefully analyze the linguistic demands of these tasks. To meet the challenges of our tasks, we present GenRe, a model which combines constrained attribute generation and re-ranking on top of pre-trained language models. GenRe achieves the best performance vs. established Relation Extraction baselines on the Extraction task ( $\geq 3.7$ F1 points) as well as the more challenging INFER-ENCE task that involve lexical and commonsense inferences ( $\geq 5.1$  F1 points). Together, our work contributes an important step towards realizing the potential of personal attributes in personalization of social chit-chat and task-oriented dialogue agents.

## 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722

723

724

725

726

669

670

### 618 Ethics and Broader Impact

619Privacy in real world applications Because our620task involves extracting and inferring personal at-621tributes, real-world users should be given the option622to disallow particular types of relations from being623collected and/or used for downstream applications.624Users should also be given the freedom to delete625their collected personal attributes. A further step626might be to restrict the extraction and storage of627personal attributes to only local devices using dif-628ferential privacy and federated learning techniques.

#### References

632

633

641

642

643

647

649

650

654

- Christoph Alt, Marc Hübner, and Leonhard Hennig. 2019. Improving relation extraction by pre-trained language representations. In *Automated Knowledge Base Construction (AKBC)*.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Çelikyilmaz, and Yejin Choi. 2019. Comet: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL).*
- Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Daniel Duckworth, Semih Yavuz, Ben Goodrich, Amit Dubey, Kyu-Young Kim, and Andy Cedilnik. 2019. Taskmaster-1: Toward a realistic and diverse dialog dataset.
- Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W Black, Alexander Rudnicky, Jason Williams, Joelle Pineau, Mikhail Burtsev, and Jason Weston. 2019. The second conversational intelligence challenge (convai2).
- Christiane Fellbaum. 1998. WordNet: An Electronic Lexical Database. Bradford Books.
- Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia d'Amato, Gerard de Melo, Claudio Gutierrez, José Emilio Labra Gayo, Sabrina Kirrane, Sebastian Neumaier, Axel Polleres, Roberto Navigli, Axel-Cyrille Ngonga Ngomo, Sabbir M. Rashid, Anisa Rula, Lukas Schmelzeisen, Juan Sequeda, Steffen Staab, and Antoine Zimmermann. 2021. Knowledge graphs.
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue.
- Chaitanya K. Joshi, Fei Mi, and Boi Faltings. 2017. Personalization in goal-oriented dialog.
- Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019.

Text Generation from Knowledge Graphs with Graph Transformers. 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 2284–2293, Minneapolis, Minnesota. Association for Computational Linguistics.

- Aaron W. Li, Veronica Jiang, Steven Y. Feng, Julia Sprague, Wei Zhou, and Jesse Hoey. 2020. Aloha: Artificial learning of human attributes for dialogue agents. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8155–8163.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003, Berlin, Germany. Association for Computational Linguistics.
- X. Li, G. Tur, D. Hakkani-Tür, and Q. Li. 2014. Personal knowledge graph population from user utterances in conversational understanding. In 2014 IEEE Spoken Language Technology Workshop (SLT), pages 224–229.
- Yichao Lu, Manisha Srivastava, Jared Kramer, Heba Elfardy, Andrea Kahn, Song Wang, and Vikas Bhardwaj. 2019. Goal-oriented end-to-end conversational models with profile features in a real-world setting. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Industry Papers), pages 48–55, Minneapolis, Minnesota. Association for Computational Linguistics.
- Liangchen Luo, Wenhao Huang, Qi Zeng, Zaiqing Nie, and Xu Sun. 2019. Learning personalized end-toend goal-oriented dialog. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):6794– 6801.
- Bodhisattwa Prasad Majumder, Harsh Jhamtani, Taylor Berg-Kirkpatrick, and Julian McAuley. 2020. Like hiking? you probably enjoy nature: Personagrounded dialog with commonsense expansions. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 9194–9206, Online. Association for Computational Linguistics.
- Pierre-Emmanuel Mazaré, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. Training millions of personalized dialogue agents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2775–2779, Brussels, Belgium. Association for Computational Linguistics.
- A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Bordes, D. Parikh, and J. Weston. 2017. Parlai: A dialog research software platform. *arXiv preprint arXiv:1705.06476*.

Kaixiang Mo, Shuangyin Li, Yu Zhang, Jiajun Li, and Qiang Yang. 2017. Personalizing a dialogue system with transfer reinforcement learning.

727

728

730

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

770

772

774

775

776

777

778

779 780

781

782

- Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. OpenDialKG: Explainable conversational reasoning with attention-based walks over knowledge graphs. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 845–854, Florence, Italy. Association for Computational Linguistics.
- Tapas Nayak and Hwee Tou Ng. 2020. Effective modeling of encoder-decoder architecture for joint entity and relation extraction. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8528– 8535.
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher D. Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, and Daniel Zeman. 2016. Universal Dependencies v1: A multilingual treebank collection. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 1659–1666, Portorož, Slovenia. European Language Resources Association (ELRA).
  - Aasish Pappu and Alexander Rudnicky. 2014. Knowledge acquisition strategies for goal-oriented dialog systems. In Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL), pages 194–198, Philadelphia, PA, U.S.A. Association for Computational Linguistics.
  - Jiahuan Pei, Pengjie Ren, and Maarten de Rijke. 2021. A cooperative memory network for personalized task-oriented dialogue systems with incomplete user profiles. *arXiv preprint arXiv:2102.08322*.
  - Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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 2463–2473, Hong Kong, China. Association for Computational Linguistics.
- Martin F Porter. 1980. An algorithm for suffix stripping. *Program*, 14(3):130–137.
- Qiao Qian, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Assigning personality/profile to a chatting machine for coherent conversation generation. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 4279–4285. International Joint Conferences on Artificial Intelligence Organization.
  - A. Radford and Karthik Narasimhan. 2018. Improving language understanding by generative pre-training.

Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2019. Towards scalable multi-domain conversational agents: The schema-guided dialogue dataset. *arXiv preprint arXiv:1909.05855*. 784

785

789

790

792

793

794

796

800

801

803

804

805

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Kurt Shuster, Eric M. Smith, Y-Lan Boureau, and Jason Weston. 2020. Recipes for building an opendomain chatbot.
- 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).
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI'17, page 4444–4451. AAAI Press.
- Anna Tigunova, Andrew Yates, Paramita Mirza, and Gerhard Weikum. 2019. Listening between the lines: Learning personal attributes from conversations.
- Anna Tigunova, Andrew Yates, Paramita Mirza, and Gerhard Weikum. 2020. CHARM: Inferring personal attributes from conversations. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5391–5404, Online. Association for Computational Linguistics.
- David Wadden, Ulme Wennberg, Yi Luan, and Hannaneh Hajishirzi. 2019. Entity, relation, and event extraction with contextualized span representations. 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 5784–5789, Hong Kong, China. Association for Computational Linguistics.
- Sean Welleck, Jason Weston, Arthur Szlam, and Kyunghyun Cho. 2019. Dialogue natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3731–3741, Florence, Italy. Association for Computational Linguistics.
- Chien-Sheng Wu, Andrea Madotto, Zhaojiang Lin, Peng Xu, and Pascale Fung. 2019. Getting to know you: User attribute extraction from dialogues. *arXiv preprint arXiv:1908.04621*.
- Hongbin Ye, Ningyu Zhang, Shumin Deng, Mosha Chen, Chuanqi Tan, Fei Huang, and Huajun Chen. 2021. Contrastive triple extraction with generative transformer.

Dian Yu, Kai Sun, Claire Cardie, and Dong Yu. 2020. Dialogue-based relation extraction. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4927–4940, Online. Association for Computational Linguistics.

839

842

843

844

845

847

851

852

853

854

855

857

858

859

866

867

870

871

872

873

875

- 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, ACL 2020, pages 109–117.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2204– 2213, Melbourne, Australia. Association for Computational Linguistics.
  - Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017. Positionaware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, pages 35–45.
  - Yinhe Zheng, Guanyi Chen, Minlie Huang, Song Liu, and Xuan Zhu. 2020a. Personalized dialogue generation with diversified traits.
  - Yinhe Zheng, Rongsheng Zhang, Minlie Huang, and Xiaoxi Mao. 2020b. A pre-training based personalized dialogue generation model with persona-sparse data. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):9693–9700.
- Xuhui Zhou, Yue Zhang, Leyang Cui, and Dandan Huang. 2020. Evaluating commonsense in pretrained language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):9733– 9740.

### A Appendix

878

879

882

886

894

895

897

901

902

903

904

905

906

907

908

### A.1 Blender Fine-tuning Details

Finetuning hyperparameters are taken from https://parl.ai/projects/recipes/, with the exception of validation metric changed to Hits@1. Each finetuning epoch takes 1.5 hours on a Nvidia V100 GPU. We only prepend personal attributes before system utterances but not user utterances. Metrics are for the validation set because test set was not available. All experiments were conducted using ParlAI (Miller et al., 2017).

#### A.2 Task Analysis Details



Figure 3: Bar plot for 10 most common POS tags of tail entities.

## A.3 Details of Transformations to Link Tail Entity to Sentence

**ConceptNet\_related**: All words in the tail entity can be found in the 100 most related words to each sentence word based on embedding distance on ConceptNet

**ConceptNet\_connect**: All words in the tail entity can be found in the 100 words that have the highest-weighted edge with each sentence word on ConceptNet.

**WordNet\_synonym**: All words in the tail entity can be found in the synonyms of every synset of each sentence word on WordNet.

WordNet\_hypernym: All words in the tail entity can be found in the hypernyms of every synset of each sentence word on WordNet

**WordNet\_hyponym**: All words in the tail entity can be found in the hyponyms of every synset of each sentence word on WordNet **Same\_stem**: All words in the sentence and tail entity are stemmed using a Porter Stemmer (Porter, 1980) before searching for the tail entity in the sentence 909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

### A.4 Generator Details

GPT-2-small was used. Additional special tokens including the control tokens ([HEAD], [RELN], [TAIL]) as well as relation tokens were added into the tokenizer. Beam search decoding (beam size = 10) was used at inference time. GPT2-small was accessed from HuggingFace Transformers library with 125M parameters, context window 1024, 768-hidden, 768-hidden, 12-heads, dropout = 0.1. AdamW optimizer was used with  $\alpha = 7.5 * 10^{-4}$ for the EXTRACTION dataset and  $\alpha = 2.5 * 10^{-3}$ for the INFERENCE dataset, following a uniform search using F1 as the criterion at intervals of  $\{2.5, 5, 7.5, 10\} * 10^n; -5 \le n \le -3$ . Learning rate was linearly decayed (over a max epoch of 8) with 100 warm-up steps. Each training epoch took around 0.5 hour on an Nvidia V100 GPU with a batch size of 16. Validation was done every 0.25 epochs during training. 5 different seeds (40-44) were set for 5 separate runs.

## A.5 Reranker Details

BERT-base-uncased was used. Additional special tokens including the control tokens ([HEAD], [RELN], [TAIL]) as well as relation tokens were added into the tokenizer. BERT-base-uncased was accessed from HuggingFace Transformers library (with 12-layer, 768-hidden, 12-heads, 110M parameters, dropout = 0.1). The choice of the base model was made to have fairness of comparison with baseline models in terms of model size. AdamW optimizer was used with  $\alpha = 5 * 10^{-6}$ , following a uniform search using F1 as the criterion at intervals of  $\{2.5, 5, 7.5, 10\} * 10^n; -6 \le n \le -3$ . Learning rate was linearly decayed (over a max epoch of 8) with 100 warm-up steps. Each training epoch took around 1 hour on an Nvidia V100 GPU with a batch size of 10. Validation was done every 0.25 epochs during training. 5 different seeds (40-44) were set for 5 separate runs.