Inducing Global and Local Knowledge Attention in Multi-turn Dialog Understanding

Anonymous ACL submission

Abstract

In multi-turn dialog understanding, semantic frames are constructed by detecting intents and 002 slots within each user utterance. However, recent works lack the capability of modeling 005 multi-turn dynamics within a dialog where the contexts are mostly adopted for updating di-007 alog states instead of capturing overall intent semantic flows in spoken language understanding (SLU). Moreover, humans rely on commonsense knowledge to better illustrate slot seman-011 tics revealed from word connotations, which many works only considered for end-to-end re-012 sponse generation. In this paper, we propose to amend the research gap by equipping a BERTbased SLU framework with knowledge and context attention modules. We propose three attention mechanisms to induce both global and local attention on knowledge triples. Experimental results in two complicated multi-turn dialog datasets have demonstrated significant improvements of our proposed framework by mutually modeling two SLU tasks with commonsense knowledge and dialog contexts. Attention visualization also provides nice interpretability of how our modules leverage knowledge across the utterance.

1 Introduction

027

034

040

In conventional task oriented dialog systems, spoken language understanding (SLU) modules aim to transform utterances into meaningful semantic representations for dialog management (Weld et al., 2021; Zhang et al., 2020). It mainly detects associated dialog acts or intents and extracts key slot information as so-called '*semantic frames*' (Abbeduto, 1983), shown in Table 1. In order to understand an utterance, besides intra-sentence semantics, humans usually manipulate commonsense knowledge to associate previous contexts with current relevant objects. In Table 1, knowledge triples representing background experiences and act relations such as '*Inform*' may follow '*Request*' acts

Speaker	Utterance
1 Uson	Is there something that's
1. User	maybe a good intelligent comedy ?
Act & Slots:	Request (genre: comedy)
	(<i>intelligent</i> ; related to; well_informed)
Knowledge:	(comedy; related to; comic)
	(comedy; is a; drama)
	Whiskey Tango Foxtrot is the only
2. System	Adult comedy I see playing in your
	area. Would you like to try that?
	Inform (movie: Whiskey Tango Foxtrot)
Act & Slote	Inform (genre: Adult comedy)
Act & Slots.	Inform (distance limits: in your area)
	Confirm_question
	(foxtrot; related to; dance)
Knowladge	(<i>foxtrot</i> ; related to; rhythm)
Knowledge:	(adult; capable of; work)
	(area; is a; region)

Table 1: Excerpt of a single turn within a dialog with corresponding dialog acts, slots and knowledge samples that are related to **keywords** in the utterance.

may benefit the prediction of overall intent semantics and slot values. 042

043

044

045

046

048

051

052

054

056

057

060

061

062

063

However such intuition has usually not been emphasized when automating SLU tasks. In early attempts of SLU systems, utterances were isolated and analyzed separately for user intents and semantic slots (Raymond and Riccardi, 2007; Liu et al., 2017). Models that maximize the joint distribution likelihood were then proposed to allow transitions between two tasks (Liu and Lane, 2016; Wang et al., 2018; Wu et al., 2021a; Li et al., 2018a). Some works also tackled utterances with multiple intents (Qin et al., 2019; Rashmi Gangadharaiah, 2019; Qin et al., 2020). While driven by large pretrained corpus, these methods still fall short of employing complete dynamic interactions within dialogs. Some works have then integrated previous dialog contexts for more robust SLU (Wang et al., 2019; Gupta et al., 2019; Su et al., 2021; Wu et al., 2021c). However, many of them cannot capture dialog flows well with RNN encoders.

Despite considering contexts, relying simply on

training dialog corpus may limit the machine to fully explore the relations between contexts and 065 slots without external commonsense knowledge. 066 Much efforts have pushed forward the progress in knowledge grounded dialog generation (Wang et al., 2021b; Zhao et al., 2020; Zheng et al., 2021), where relevant documents or a knowledge base auxiliarily guide the language autoregressive progress. Term-level denoising (Zheng et al., 2021) or filtering techniques (Wang et al., 2021b) refine the adopted knowledge for better semantic considerations. However, construction of semantic frames 075 may also require knowledge induction in more com-076 plex dialogs. Wang et al. (2019) has proposed to adopt knowledge attention for joint tasks. However, 078 it adopts a single LSTM layer to couple all knowledge and contexts without filtering, which cannot model complex interactions well and is ambiguous in how these two components affect each other.

064

077

079

090

094

095

097

100

101

102

103

104

105

106

107

108

109

110

To solve the above concerns, we propose a Global and Local Knowledge Attention Framework (GLKA) to amend the research gap in joint SLU tasks by effectively incorporating dialog history and external knowledge. We propose three different attention modules that consider local and global awareness of knowledge at token and utterance levels respectively. After obtaining knowledge-enriched vectors, we predict intents and slots coherently with two LSTM decoders with different fused inputs. Experiment results have shown superior performances of our methods in manipulating contexts and knowledge and beat all competitive baselines. Our contributions are as follows:

1. We propose GLKA framework to fill the void of exploring relations between commonsense knowledge and dialog history in recent SLU works. It dynamically selects knowledge with contexts for multiple dialog act and slot filling detection.

2. We demonstrate the benefits of knowledge and context induction in the low resource setting and non-alphabetic slots.

3. Experimental and attention visualization results show that our model achieves superior performances over several competitive baselines and provides good interpretability of how our model utilizes the knowledge.

Problem Formulation 2

For each utterance $x_n = \{w_1^n, w_2^n, \dots, w_T^n\}$ in a 111 task-oriented dialog \mathbf{X} with N utterances, given 112 the domain ontology of a dialog act set A and a 113

slot set **S**, we aim to find one or more acts $\{a_i^n\}$ 114 ¹ and a sequence of slot tags $\{s_1^n, s_2^n, \ldots, s_T^n\}$ to 115 construct a semantic frame. Namely, we hope to 116 maximize the joint log likelihood of A and S in 117 Eq. 1 given a parametrized model θ , its context 118 $\mathbf{C_n} = \{x_1, \ldots, x_{n-1}\}$ and associated knowledge 119 $\mathbf{K_n} = \phi(K_G, x_n)$ for the current utterance x_n . We 120 deem K_G as an external large knowledge base with 121 knowledge triples and $\phi(\cdot)$ helps to extract related 122 knowledge pairs for x_n . It will be critical to match 123 correct knowledge based on current dialog history 124 and the utterance for better dialog understanding. 125

$$\mathcal{L}(\mathbf{A}, \mathbf{S}) \triangleq \sum_{n} \log P(A_n, S_n \mid x_n, \mathbf{C_n}, \mathbf{K_n}; \theta)$$
(1)

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

Methodology 3

3.1 **Context Attention**

Our overall framework is illustrated in Figure. 1. To allow information flow across the dialog, we first encode the entire dialog with a token-level BERT (Devlin et al., 2019) encoder and a turnlevel context-aware transformer encoder. Instead of concatenating all sentences which may cause an extreme sequence length, we first generate the token-level representations $\mathbf{H} = \{h_1, h_2, \dots, h_N\}$ for each utterance x_n in a dialog X by taking vectors from each [CLS] token. During testing at turn n, we may directly reuse these calculated representations $\{h_1, h_2, \ldots, h_{n-1}\}$ until turn n-1.

In contrast with other contextual SLU (Wang et al., 2019; Gupta et al., 2019) with hierarchical structures, we succinctly introduce a unidirectional transformer encoder with the hidden size H_a to encode $\mathbf{H} \in \mathbb{R}^{N \times H_b}$, which may allow mutual attention flow between dialog contexts. It consists L layers of masked multi-head self-attention (MHA), point-wise feed forward network (FFN), residual sublayer and layer normalization. The future time steps are masked for training since we will not have access to future utterances during testing. We will send H as the first layer input C^1 and iteratively encode it with two sublayers in Eq. 2. Each head $\mathbf{C}_{\mathbf{i}} \in \mathbb{R}^{N \times (H_a/h)}$ will be first mapped into a query $\mathbf{C}^{\mathbf{Q}}$, a key $\mathbf{C}^{\mathbf{K}}$ and a value $\mathbf{C}^{\mathbf{V}}$ which participate in the multi-head self-attention. Here $f(\cdot)$ is softmax function. Finally, we will obtain the final

¹Here we refer the intent detection problem in dialogs as predicting the dialog acts for each utterance.



Figure 1: Illustration of our proposed framework for joint dialog act detection and slot filling in multi-turn dialogs. It consists of context and knowledge attention modules, and two LSTM-based decoders. The utterance-level representations will be encoded with the context attention module and token-level representations will interact with their corresponding knowledge in three proposed awareness submodules.

contextual dialog representations $\mathbf{C}^{\mathbf{L}}$.

158

160

161

162

$$\mathbf{C}^{\mathbf{l}} = FFN(MHA(\mathbf{C}^{\mathbf{l-1}}, \mathbf{C}^{\mathbf{l-1}}, \mathbf{C}^{\mathbf{l-1}}))$$
(2)

$$MHA(\mathbf{C_i^Q}, \mathbf{C_i^K}, \mathbf{C_i^V}) = f(\frac{\mathbf{C_i^Q}(\mathbf{C_i^K})^T}{\sqrt{H_b}})\mathbf{C_i^V}$$
(3)

$$FFN(x) = max(0, x\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2$$
(4)

3.2 Knowledge Attention

Human could naturally associate words and contexts with commonsense knowledge to predict se-164 mantics. Intuitively, some knowledge relations for 165 a particular word may allude its tendency to some slots and intents. Hence we introduce three attention mechanisms to delineate such human heuris-168 tics and obtain the final knowledge-aware vectors 169 V_{K} in Figure. 2. We first purpose Local aware**ness** of knowledge where we select the top |K|171 knowledge relations of each word when predicting its corresponding slot locally. However, semantic 173 slots (BIO scheme) may be expressed as phrases in-174 stead of individual words where knowledge should be possibly shared across words. Therefore, we 176 introduce Global awareness to share all knowl-177 edge gathered from each word for global atten-178 tion directly. Eventually, such treatment is on the 179 contrary opaque on how knowledge is related to slot decision of word individuals and some out-of-181 vocabulary (OOV) words may not have relevant 182 knowledge in data base, which requires knowledge from other words in proximity. We then purpose 184

Global-Local awareness to predict each slot along with entire gathered knowledge.

185

187

188

189

190

191

192

194

196

197

198

199

200

202

203

204

205

207

208

3.2.1 Knowledge extraction

The first step lies in gathering knowledge for attention at the current utterance $x_n = \{w_1^n, w_2^n, \dots, w_T^n\}$. For each word w_i^n , we first retrieve a list of relations of the exactly same head entity being the word w_i^n from a knowledge base K_G . If no head entities are matched, we instead seek entities that has a substring of w_i^n . Based on the weights of relations provided in K_G , we select top |K| related triples $\gamma = \{h, r, t\}$ as the final knowledge k_i^n for w_i^n . We will finally obtain a T length knowledge sequence $\mathbf{K_n} = \{k_1^n, k_2^n, \dots, k_T^n\}$ gathered from each word w_i^n . In case of non-alphabetic or OOV words with no match in K_G , we instead replace their $\mathbf{K_n}$ as zero vectors to represent agnosticism of knowledge.

3.2.2 Local awareness

1

Locally, we obtain the local knowledge-aware vector v_i^n for each word w_i^n only based on its corresponding knowledge k_i^n (i.e. |K| triples $\gamma = \{h, r, t\}$) with the following attention mechanism.

$$w_i^n = \sum_{j=1}^M \alpha_{ij}[r_{ij}^n; t_{ij}^n]$$
 (5)

$$\alpha_{ij} = \exp(\beta_{ij}) / \sum_{m=1}^{M} \exp(\beta_{im})$$
 (6) 209

$$\beta_{ij} = (h_i^n \mathbf{W}^{\mathbf{H}})(tanh(r_{ij}^n \mathbf{W}^{\mathbf{R}} + t_{ij}^n \mathbf{W}^{\mathbf{T}}))^T$$
(7) 21

 r_{ij}^n , t_{ij}^n are relation and tail entity vectors. 211 $\mathbf{W}^{\mathbf{H}}, \mathbf{W}^{\mathbf{R}}, \mathbf{W}^{\mathbf{T}}$ are learnable matrices during 212



Figure 2: Three submodules to induce knowledge awareness. (a) Local awareness performs attention at token-level with intra-word knowledge. (b) Global awareness takes all knowledge related to the utterance for context-based attention. (c) Global-Local awareness performs attention at token-level but with all inter-word knowledge.

training. M is the number of knowledge triples. [;] is the concatenation of two vectors. Given the token-level representations for each word h_i^n in the utterance x_n , attention weights are assigned to reveal the relevance of each knowledge triple under current contexts.

3.2.3 Global awareness

213

214

215

216

219

221

222

229

230

The above mechanism may be restricted in relying on intra-word knowledge to form token-level knowledge vectors. In real case, some phrases may have continuous words where knowledge should be shared globally. Therefore, instead of attending knowledge locally, we aggregate the knowledge triples from all words $\mathbf{K_n} = \{k_i^n\}$ into a dense matrix and directly find the attention weights for our utterance-level contexts c_n^L to form a global knowledge-aware vector v^n instead.

$$v^{n} = \sum_{t=1}^{T} \sum_{j=1}^{M} \alpha_{tj}[r_{tj}^{n}; t_{tj}^{n}]$$
(8)

$$\alpha_{tj} = \exp(\beta_{tj}) / \sum_{t=1}^{T} \sum_{m=1}^{M} \exp(\beta_{tm}) \quad (9)$$

$$\beta_{tj} = (c_n^L \mathbf{W}^H) (tanh(r_{tj}^n \mathbf{W}^R + t_{tj}^n \mathbf{W}^T))^T$$
(10)

3.2.4 Global-Local awareness

At last, we combine the view of global and local awareness by generating the local knowledgeaware vector v_i^n but with the global knowledge K_n . We could avert the circumstances where some OOV words may not have relevant knowledge by considering knowledge from other word neighbors. Here the knowledge-aware vector v_i^n will be obtained by weighted summing all knowledge in the sentence x_n :

$$v_i^n = \sum_{t=1}^T \sum_{j=1}^M \alpha_{tj}[r_{tj}^n; t_{tj}^n]$$
(11)

241

242

243

244

245

246

247

248

249

250

251

252

254

256

257

259

260

261

262

264

265

266

268

where T is the number of words in the sentence x_n .

3.3 Semantic Decoder

After obtaining the knowledge-enriched representations $\mathbf{V}_{\mathbf{K}} = \{v_i^n\}$ (§3.2.2, 3.2.4) or v^n (§3.2.3) along with contextual dialog representations $\mathbf{C}^{\mathbf{L}}$ and the initial token-level representations \mathbf{H} , we adopt two BiLSTM to predict multiple dialog acts and slot filling.

$$\mathbf{H_{slot}} = \mathbf{BiLSTM}([\mathbf{H}; \mathbf{V_K}], \mathbf{C^L}) \qquad (12)$$

$$\mathbf{H}_{\mathbf{act}} = \mathbf{BiLSTM}([\mathbf{C}^{\mathbf{L}}; \mathbf{V}_{\mathbf{K}}])$$
(13)

For slot filling, $\mathbf{V}_{\mathbf{K}}$ will be first concatenated with \mathbf{H} and serve as the inputs of BiLSTM with $\mathbf{C}^{\mathbf{L}}$ as initial hidden states, where contexts will assist the slot prediction at each knowledge-enhanced time step. At the same time, $\mathbf{V}_{\mathbf{K}}$ will also be concatenated with dialog contexts $\mathbf{C}^{\mathbf{L}}$ to serve as inputs for another BiLSTM. Finally, we can generate logits $\hat{y}_{act} = \sigma(\mathbf{H}_{act}\mathbf{W}_{act})$ by transforming \mathbf{H}_{act} with $\mathbf{W}_{act} \in \mathbb{R}^{H_L \times |\mathcal{Y}^a|}$ and a sigmoid function σ . H_L is LSTM hidden size and $|\mathcal{Y}^a|$ is the size of dialog act set. Likewise, we compute $\hat{y}_{slot} = softmax(\mathbf{H}_{slot}\mathbf{W}_{slot})$. Total loss will be the combination between the binary cross entropy loss based on \hat{y}_{act} and the cross entropy loss based on \hat{y}_{slot} as shown in Eq. 14, 15. Finally, the joint

335

337

339

340

341

342

343

344

345

346

347

348

350

351

352

353

354

355

356

358

359

360

361

362

314

315

269

278

285

290

293

294

295

296

301

objective is formulated as the sum of \mathcal{L}_a and \mathcal{L}_s .

$$\mathcal{L}_{a} \triangleq -\sum_{n=1}^{N} \sum_{a=1}^{|\mathcal{Y}^{a}|} (y_{a}^{n} log(\hat{y}_{a}^{n}) + (1 - \alpha^{n}) log(1 - (\hat{\alpha}^{n})))$$
(1)

 $+(1-y_a^n)log(1-(\hat{y}_a^n))$ (14)

$$\mathcal{L}_{s} \triangleq -\sum_{n=1}^{N} \sum_{t=1}^{T} \sum_{s=1}^{|\mathcal{Y}^{s}|} (y_{s}^{(n,t)} log(\hat{y}_{s}^{(n,t)}))$$
(15)

4 Experiment Setting

4.1 Experimental setup

We evaluate our proposed framework on two largescale dialog datasets, i.e. Microsoft Dialog Challenge dataset (MDC) (Li et al., 2018b) and Schema-Guided Dialog dataset (SGD) (Rastogi et al., 2019). **MDC** contains human-annotated conversations in three task-completion domains (movie, restaurant, taxi) with total 11 dialog acts and 50 slots. **SGD** entails large-scale task-oriented dialogs over 20 domains ranging from travel, weather to banks, etc. It has total 18 dialog acts and 89 slots. We randomly select 1k dialogs for each domain in MDC and two domains (restaurant, flights) from SGD for total 5k dialogs in 7:3 training and testing ratio. Each utterance is labeled with one or more dialog acts and several slots.

4.2 Baselines

We compare our models with several competitive baselines which sequentially include more features for better semantic considerations:

1) MID-SF (Rashmi Gangadharaiah, 2019) which first considers joint multi-intent and slot detection in use of BiLSTMs.

2) ECA (Chauhan A., 2020) which encodes the dialog context with LSTM for joint task prediction.
 3) KASLUM (Wang et al., 2019) which extracts knowledge from the knowledge base and incorporates dialog history for joint tasks.

3024) CASA (Gupta et al., 2019) which encodes the
context with DiSAN and sentence2token where we
replace DiSAN with BERT for better comparison.304replace DiSAN with BERT for better comparison.305We also denote several variations of our proposed
framework with the following detailed descriptions.306framework with the following detailed descriptions.3071) LKA_{AF} (Wang et al., 2021b): we replace only
Local awareness part (§ 3.2.2) with the attention-
based filter (AF) in (Wang et al., 2021b) to compare
different knowledge attention.

- **2) LKA:** local awareness version of our model.
- **312 3) GKA**: global awareness version of our model.
- **4) GLKA** $_T$: we replace the semantic decoder part

(§ 3.3) with a single transformer decoder to both predict dialog acts and slots.

5) **GLKA**: global-local awareness version of our model.

4.3 Implementation details

We adopt the pretrained \mathbf{BERT}_{base} (Devlin et al., 2019) as our utterance encoder. Context attention transformer has L = 6-layer attention blocks with 768 head size and 4 attention heads. The max sequence length is 60. We use ConceptNet knowledge base (Speer et al., 2018) to obtain relevant knowledge for attention. Then, TransE (Bordes et al., 2013) is adopted to represent head, relation and tail as pretrained 100-dim vectors. We retrieve |K| = 5 most related knowledge from each word based on weights assigned on the edges. Both LSTMs have 256 hidden units. We use the batch size of 4 dialogs for MDC and 2 for SGD. In all training, we use Adam optimizer with learning rate as 5e-5. The best performance on validation set is obtained after training 60 epochs on each model. For metrics, we report the dialog act accuracy and slot filling F1 score. Here we only consider a true positive when all BIO values for a slot is correct and forfeit 'O' tags.

5 Main Results

5.1 Main results

Table. 2 shows our main results on the joint task performances. MID-SF with only LSTMs has relatively inferior performances on both datasets especially in SGD. ECA by taking dialog contexts into consideration has much greater increase in SGD than in MDC and further knowledge induction gives 3.5 % increase in KASLUM. Leveraging BERT-based encoder seems to substantially increase semantic visibility in CASA and our proposed frameworks. Eventually, all of our knowledge-enhanced models beat all baselines both in MDC and substantially in SGD, by more efficiently incorporating external knowledge and dialog contexts with the proposed mutual attention mechanism. We first see our purposed knowledge attention in LKA has better effectiveness than LKA $_{AF}$. And GLKA almost beats every baseline to demonstrate the advantage of sharing knowledge globally while maintaining local attention on each word. Finally, we see a single transformer decoder may still entangle the act and slot information while updating gradients simultaneously, where separate

Dataset			M	DC				SC	GD	
Domain	Mo	ovie	Resta	urant	Ta	axi	Resta	urant	Flig	ghts
Model	MDA	SL								
MID-SF (Rashmi Gangadharaiah, 2019)	76.56	67.56	77.35	65.77	85.03	70.03	74.26	81.38	84.74	84.48
ECA (Chauhan A., 2020)	77.10	69.72	77.56	66.85	86.61	71.28	87.98	84.87	95.16	87.91
KASLUM (Wang et al., 2019)	81.86	73.32	80.76	68.36	88.31	74.07	86.81	87.82	92.87	90.05
CASA (Gupta et al., 2019)	84.22	79.59	83.17	74.89	90.00	78.54	92.54	94.20	95.00	91.79
LKA_{AF}^{\dagger} (Wang et al., 2021b)	85.25	79.46	83.27	74.89	90.05	79.59	96.84	94.61	97.17	91.14
LKA^{\dagger} (ours)	85.59	80.21	83.48	75.30	90.01	79.14	98.25	94.57	98.00	92.31
\mathbf{GKA}^{\dagger} (ours)	85.94	80.56	83.64	75.94	90.28	79.08	98.44	94.75	98.74	91.71
GLKA_T^{\dagger} (ours)	85.98	79.94	83.27	75.19	90.40	78.33	97.35	94.34	98.20	91.95
GLKA [†] (ours)	86.09	80.58	84.01	75.27	90.80	79.60	98.47	94.86	99.22	92.67

Table 2: Experimental Results on several SLU models including our proposed frameworks which are specified in percentage (%). MDA indicates the dialog act detection accuracy by counting corrects when all acts are predicted correctly. SL indicates the slot filling F1 score. † denotes models related to our proposed structures.

Dataset	MDC					SGD				
Domain	Mo	ovie	Resta	urant	Ta	nxi	Resta	urant	Flig	ghts
Model	MDA	SL								
GLKA	86.09	80.58	84.01	75.27	90.80	79.60	98.47	94.86	99.22	92.67
w/ KG_V	85.63	80.26	83.43	75.76	89.77	80.03	98.38	94.31	98.93	91.99
w/o KG	86.01	79.92	83.53	74.76	90.56	78.29	97.53	94.83	97.73	92.23
w/o CA	84.87	79.79	81.33	74.68	89.00	78.50	95.88	94.36	97.17	91.94
w/o LSTM	84.57	79.14	82.70	74.35	89.65	79.00	90.96	93.64	94.80	91.33

Table 3: Ablation Results of joint tasks (%) by removing some key components of our proposed frameworks GLKA.

LSTMs perform better in our case. To note, the word matching accuracies in the knowledge base are 78.12% (MDC) and 80.97% (SGD), which indicate that there is still about 20% of zero vectors introduced as redundant noises.

5.2 Ablation analysis

363

365

366

To better estimate the effectiveness of each module of our best model: GLKA, we conduct ablation 370 experiments in Table. 3. We sequentially ablate 371 each component from GLKA to observe the performance drops. We first replace the top |K| knowledge vectors with those ranked within $|K| \sim 2|K|$ 374 behind to compare the effect of knowledge qual-375 ity (w/ KG_V). We could see an overall performance drop except slot accuracy may increase in 378 some domains, which indicates that the selection of knowledge may play a critical role in how model leverages the relations. By removing the entire 380 knowledge attention module, we can see more obvious reduction in slot filling tasks denoting the necessity of external knowledge in enriching the current word representations. By substituting a unidirectional LSTM on top of BERT for our context attention module (CA), we obtain poorer performance in dialog act detection instead. Finally, we 387 see dialog contexts are more crucial in SGD where 388 drop seems significant by removing all context fusion modules. Overall, we observe dialog act detection relies more on contexts while slot filling tasks may concentrate on inter-utterance relations where external knowledge benefits more instead. 391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

5.3 Further Discussion

Could knowledge amend the data scarcity? We also study how knowledge could contribute to the joint tasks when resources are scarce. Figure. 3 shows the performance changes with different numbers of training data. We found that inducing the knowledge will have the positive effect on both tasks. In the few-shot setting, we see the performance difference enlarges where knowledge becomes beneficial to enrich the external information aside from data itself. However, knowledge becomes less useful when we have extreme low dataset particularly for slot detection.

Does global knowledge helps non-alphabetic slots? We are interested if knowledge for other words would also help with the slot prediction of the non-alphabetic words. Table. 4 shows the results for each non-alphabetic slot for our local and global attention models. Since there is no knowledge for the non-alphabetic words, we observe an overall 2% increase by inducing global attention. Contexts are beneficial especially for slots associated with rating, money and address, which should be likely inferred by other keywords near them. However, time and zip code are rather independent



Figure 3: SLU performance by training GLKA with a subsample (%) of the original training data of two datasets: MDC and SGD. We show the results with or without the knowledge induced.

to contexts which may be disturbed by introducing more irrelevant noises.

5.4 Knowledge Attention

419

420

421

In Figure. 4, we visualize the attention heatmap 422 of tokens with their slot labels vs. all knowledge 423 triples from each token. First, we focus on the 424 rows of the heat map. Without attached knowledge 425 for the words like numbers or punctuations, their 426 attention weights are perceived blank across all to-427 kens in the utterance. Second, for valid attention 428 weights, we found the knowledge corresponding 429 to keywords like 'you', 'with', 'restaurant' and 430 'antioch' are most adopted for overall knowledge 431 representations across all the utterance. It reck-432 ons that the model will mostly grasp knowledge in 433 words especially tagged as valued slots (non-O tag) 434 for overall semantic understanding. Interestingly, 435 this collection of knowledge is more emphasized 436 on predicting a word to be non-valued than those 437 words with valued slots. For the columns, we could 438 see for non-valued words, they will accentuate on 439 knowledge of valued words like 'restaurant' and 440 'antioch', than the knowledge related to itself. It 441 substantiates the belief that the overall semantics of 442 the utterance may be driven by these valued words. 443

Slot	GLKA (%)	LKA (%)	$\Delta(\%)$
address	17.39	0.00	+17.39
price	66.67	50.00	+16.67
critic_rating	34.48	23.08	+11.41
dress_code	50.00	44.44	+5.56
rating	52.17	49.32	+2.86
cost	95.54	95.29	+0.26
numberofpeople	95.63	95.51	+0.12
date	86.96	86.99	-0.02
pricing	42.55	43.14	-0.58
starttime	76.80	77.68	-0.88
numberofkids	73.68	77.78	-4.09
mpaa_rating	76.92	83.33	-6.41
zip_code	77.65	84.44	-6.80
pickup_time	75.19	82.29	-7.09
total	65.83	63.80	+2.03

Table 4: F1 scores for GLKA and LKA of nonalphabetic slots in overall MGD dataset.

444

445

446

447

448

449

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

In Table. 5, we further show an utterance example with some highlighted words including 'you', 'restaurant' and 'Antioch' with their extracted knowledge and weights for semantic detection. We take the average of all attention weights across all tokens for that knowledge triple; then normalized across the knowledge triples in the same word (head). We could see 'you' as an object is most adopted to clarify the user being offered and informed counts. Then we observe that the knowledge triple (restaurant, atl, city) where restaurant is at a location of the city is most recognized to illustrate the relations of restaurant and city tags. Finally, knowledge for 'Antioch' keyword is mostly relevant to a country which is conducive when the system seldom sees this word during training. But without further contexts, our model believes 'Antioch' is more of a part of Turkey.

6 Related Work

Intent detection and slot filling are two main SLU tasks (Weld et al., 2021). Many classification-based approaches (Sarikaya et al., 2011; Raymond and Riccardi, 2007; Liu et al., 2017) had been proposed to solve single intent detection problems. However, treating two tasks separately may experience error propagation. Liu and Lane (2016) first proposed an attention-based LSTM network to model the correlations between intents and slots. Li et al. (2018a) proposed the gating mechanism for better self-attention on joint tasks. However simply relying on the gate function is not ideal for long sequences. Wang et al. (2018) instead proposed the bi-model to directly model the cross impacts and Zhang et al. (2019) utilized capsule neural net-



Figure 4: Attention visualization of a single utterance example with respect to all knowledge related to each word. We denote an utterance with tokens followed by their predicted tag in x-axis. For y-axis, each word will have five knowledge triples with each as a single tick. The blank area is where attention weights are zero.

works. Memory networks are also popular choices to model long-range dependency (Wu et al., 2021a). However, a single utterance may have many intents. Rychalska et al. (2018) first proposed hierarchical structures to explore multiple intents. Qin et al. (2019) proposed a stack-propagation networks to predict intents on each token. Rashmi Gangadhara-iah (2019) and (Qin et al., 2020) considered the dynamic interactions between two tasks by jointly detecting multiple intents. Wu et al. (2021b) extended the multiple intent scenario with zero-shot cases. These methods nevertheless restrict their resources to current utterances for prediction.

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

Contexts and knowledge With respect to dialogs, contexts are also critical for semantic understanding. Bertomeu et al. (2006) first studied the contextual phenomena in words. Bhargava et al. (2013) and Shi et al. (2015) then introduced contextual signals to the joint intent-slot tasks. Advanced hierarchical structures are also emphasized to encode multi-turn dialog contexts efficiently (Chauhan A., 2020; Wang et al., 2019; Gupta et al., 2019; Wu et al., 2021c). Knowledge is also another

Utterance Example in Figure 4				
Utterance	I found 2 places that may interest you. Starting with Celia's Mexican restaurant located in Antioch .			
Dialog acts	Offer, Inform Count			
Slots	000000000000 B-res I-res I-res			
	O O B-city			
Keyword	Knowledge			
Keyword	Knowledge (hc, noun) (0.29), (hc, object) (0.7)			
Keyword you	Knowledge (hc, noun) (0.29), (hc, object) (0.7) (rel, guys) (6e-4), (hc, object) (8e-5)			
Keyword you	Knowledge (hc, noun) (0.29), (hc, object) (0.7) (rel, guys) (6e-4), (hc, object) (8e-5) (isa, establishment) (8e-9), (atl, hotel) (0.2)			
Keyword you restaurant	Knowledge (hc, noun) (0.29), (hc, object) (0.7) (rel, guys) (6e-4), (hc, object) (8e-5) (isa, establishment) (8e-9), (atl, hotel) (0.2) (atl, town) (0.14), (atl, city) (0.65)			
Keyword you restaurant	Knowledge (hc, noun) (0.29), (hc, object) (0.7) (rel, guys) (6e-4), (hc, object) (8e-5) (isa, establishment) (8e-9), (atl, hotel) (0.2) (atl, town) (0.14), (atl, city) (0.65) (rel, orontes) (4e-5), (rel, swiss) (2e-2)			

Table 5: The utterance example in Figure 4 for joint task prediction. Knowledge (Relation, Tail) related to three **keywords** as head are presented with their attention weights (number after the knowledge). We only show the top four knowledge adopted for each keyword based on the attention weights. 'hc' represents 'has context', 'rel' represents 'related to', 'atl' represents 'at location' and 'ptof' represents 'part of'.

important resource to induce commonsense for understanding. In task-oriented dialogs, Main emphasis lies in the interaction with task-related knowledge bases (Madotto et al., 2020; Yang et al., 2020). Most of works also focus on open-domain dialog response generation (Zhao et al., 2020; Wang et al., 2021b; Rashkin et al., 2021; Zheng et al., 2021) or task-specific responses (Wang et al., 2021a). Wang et al. (2019) also tried to apply knowledge in SLU but it is not suitable for complex dialog modeling. To amend the gap in modeling knowledge and context interactions of SLU, we follow these previous works' paradigms and explore the mechanisms of characterizing their mutual effects in details. 501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

7 Conclusion

In this paper, we propose a novel BERT-based knowledge augmented network to consider dialog history and external knowledge in the joint SLU tasks. We propose three approaches of inducing knowledge awareness, which are capable of selecting relevant knowledge triples for useful knowledge representation. We found that our best model (GLKA) combines the benefits both from local and global awareness, whose effectiveness is verified in two complex multi-turn dialog datasets. We visualize how our models adopt word knowledge spreading in an utterance to provide better interpretability for decision making. These knowledge fusion vectors could be easily applied to downstream dialog state tracking or management tasks.

References

531

533

534

535

536

537

538

539

540

541

542

544

545

547

548

550

551

552

553

554

556

562

571

572

573

574

575

576

577

578

579

580

584

- Leonard Abbeduto. 1983. Linguistic communication and speech acts. kent bach, robert m. harnish. cambridge: M.i.t. press, 1979, pp. xvii 327. *Applied Psycholinguistics*, 4(4):397–407.
 - Núria Bertomeu, Hans Uszkoreit, Anette Frank, Hans-Ulrich Krieger, and Brigitte Jörg. 2006. Contextual phenomena and thematic relations in database QA dialogues: results from a Wizard-of-Oz experiment. In *Proceedings of the Interactive Question Answering Workshop at HLT-NAACL 2006*, pages 1–8, New York, NY, USA. Association for Computational Linguistics.
 - A. Bhargava, A. Celikyilmaz, D. Hakkani-Tür, and R. Sarikaya. 2013. Easy contextual intent prediction and slot detection. In 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, pages 8337–8341.
 - A Bordes, N Usunier, A Garcia-Duran, J Weston, and O Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
 - Singh A. Arora J. Shukla S. Chauhan A., Malhotra A. 2020. Encoding context in task-oriented dialogue systems using intent, dialogue acts, and slots. In Saini H., Sayal R., Buyya R., Aliseri G. (eds) Innovations in Computer Science and Engineering. Lecture Notes in Networks and Systems, vol 103. Springer, Singapore.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding.
 - Arshit Gupta, Peng Zhang, Garima Lalwani, and Mona Diab. 2019. Casa-nlu: Context-aware self-attentive natural language understanding for task-oriented chatbots.
 - Changliang Li, Liang Li, and Ji Qi. 2018a. A selfattentive model with gate mechanism for spoken language understanding. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 3824–3833, Brussels, Belgium. Association for Computational Linguistics.
 - Xiujun Li, Sarah Panda, Jingjing Liu, and Jianfeng Gao. 2018b. Microsoft dialogue challenge: Building end-to-end task-completion dialogue systems. *arXiv preprint arXiv:1807.11125*.
 - Bing Liu and Ian Lane. 2016. Attention-based recurrent neural network models for joint intent detection and slot filling.
 - Ting Liu, Xiao DING, Yue QIAN, and Yiheng CHEN. 2017. Identification method of user's travel consumption intention in chatting robot. *SCIENTIA SINICA Informationis*, 47:997.

Andrea Madotto, Samuel Cahyawijaya, Genta Indra Winata, Yan Xu, Zihan Liu, Zhaojiang Lin, and Pascale Fung. 2020. Learning knowledge bases with parameters for task-oriented dialogue systems. 586

587

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

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

- Libo Qin, Wanxiang Che, Yangming Li, Haoyang Wen, and Ting Liu. 2019. A stack-propagation framework with token-level intent detection for spoken language understanding.
- Libo Qin, Xiao Xu, Wanxiang Che, and Ting Liu. 2020. Agif: An adaptive graph-interactive framework for joint multiple intent detection and slot filling.
- Hannah Rashkin, David Reitter, Gaurav Singh Tomar, and Dipanjan Das. 2021. Increasing faithfulness in knowledge-grounded dialogue with controllable features.
- Balakrishnan Rashmi Gangadharaiah. 2019. Joint multiple intent detection and slot labeling for goaloriented dialog. Proc. of NAACL.
- 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*.
- Christian Raymond and Giuseppe Riccardi. 2007. Generative and discriminative algorithms for spoken language understanding. In *Proc. Interspeech* 2007, pages 1605–1608.
- B. Rychalska, H. Glabska, and A. Wroblewska. 2018. Multi-intent hierarchical natural language understanding for chatbots. In 2018 Fifth International Conference on Social Networks Analysis, Management and Security (SNAMS), pages 256–259.
- Ruhi Sarikaya, Geoffrey E. Hinton, and Bhuvana Ramabhadran. 2011. Deep belief nets for natural language call-routing. In 2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 5680–5683.
- Yangyang Shi, Kaisheng Yao, Hu Chen, Yi-Cheng Pan, Mei-Yuh Hwang, and Baolin Peng. 2015. Contextual spoken language understanding using recurrent neural networks.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2018. Conceptnet 5.5: An open multilingual graph of general knowledge.
- Ruolin Su, Ting-Wei Wu, and Biing-Hwang Juang. 2021. Act-Aware Slot-Value Predicting in Multi-Domain Dialogue State Tracking. In *Proc. Interspeech* 2021, pages 236–240.
- Qingyue Wang, Yanan Cao, Junyan Jiang, Yafang Wang, Lingling Tong, and Li Guo. 2021a. Incorporating specific knowledge into end-to-end task-oriented dialogue systems. In 2021 International Joint Conference on Neural Networks (IJCNN), pages 1–8.

639 640 Yanmeng Wang, Ye Wang, Xingyu Lou, Wenge Rong,

Zhenghong Hao, and Shaojun Wang. 2021b. Im-

proving dialogue response generation via knowledge

graph filter. In ICASSP 2021 - 2021 IEEE Interna-

tional Conference on Acoustics, Speech and Signal

Yu Wang, Yilin Shen, and Hongxia Jin. 2018. A bimodel based rnn semantic frame parsing model for

Yufan Wang, Tingting He, Rui Fan, Wenji Zhou, and

H. Weld, X. Huang, S. Long, J. Poon, and S. C. Han. 2021. A survey of joint intent detection and slotfilling models in natural language understanding.

Jie Wu, Ian Harris, and Hongzhi Zhao. 2021a. Spoken language understanding for task-oriented dialogue systems with augmented memory networks. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 797–806, Online. Association for Computational Lin-

Ting-Wei Wu, Ruolin Su, and Biing Juang. 2021b. A label-aware BERT attention network for zero-shot multi-intent detection in spoken language understanding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4884–4896, Online and Punta Cana, Dominican Republic. Association for Computational Lin-

Ting-Wei Wu, Ruolin Su, and Biing-Hwang Juang.

Shiquan Yang, Rui Zhang, and Sarah Erfani. 2020.

GraphDialog: Integrating graph knowledge into endto-end task-oriented dialogue systems. In *Proceed*-

ings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 1878–1888, Online. Association for Computational

Chenwei Zhang, Yaliang Li, Nan Du, Wei Fan, and

Zheng Zhang, Ryuichi Takanobu, Qi Zhu, Minlie Huang, and Xiaoyan Zhu. 2020. Recent advances and challenges in task-oriented dialog system.

Xueliang Zhao, Wei Wu, Can Xu, Chongyang Tao, Dongyan Zhao, and Rui Yan. 2020. Knowledgegrounded dialogue generation with pre-trained lan-

tion via capsule neural networks.

Philip S. Yu. 2019. Joint slot filling and intent detec-

2021c. A Context-Aware Hierarchical BERT Fusion Network for Multi-Turn Dialog Act Detection. In *Proc. Interspeech 2021*, pages 1239–1243.

Xinhui Tu. 2019. Effective utilization of external knowledge and history context in multi-turn spoken language understanding model. In 2019 IEEE International Conference on Big Data (Big Data), pages

Processing (ICASSP), pages 7423-7427.

intent detection and slot filling.

960-967.

guistics.

guistics.

Linguistics.

guage models.

- 64 64
- 64
- 64
- 648
- 649 650 651
- 6

6

- 654 655 656
- 6
- 6
- 6
- 6
- 6 6
- 60 60
- 6
- 6
- 673 674
- 676

678

679 680

- 6
- 684
- (
- 6

- 6
- 6

69 69 Wen Zheng, Natasa Milic-Frayling, and Ke Zhou. 2021. Knowledge-grounded dialogue generation with termlevel de-noising. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2972–2983, Online. Association for Computational Linguistics.

694

695

697

698

699