

Joint Intent Detection and Slot Filling with Wheel-Graph Attention Networks

Pengfei Wei · Bi Zeng · Wenxiong Liao

Received: date / Accepted: date

Abstract Intent detection and slot filling are two fundamental tasks for building a spoken language understanding (SLU) system. Multiple deep learning-based joint models have demonstrated excellent results on the two tasks. In this paper, we propose a new joint model with a wheel-graph attention network (Wheel-GAT) which is able to model interrelated connections directly for intent detection and slot filling. To construct a graph structure for utterances, we create intent nodes, slot nodes, and directed edges. Intent nodes can provide utterance-level semantic information for slot filling, while slot nodes can also provide local keyword information for intent. Experiments show that our model outperforms multiple baselines on two public datasets. Besides, we also demonstrate that using Bidirectional Encoder Representation from Transformer (BERT) model further boosts the performance in the SLU task.

Keywords Spoken language understanding · Graph neural network · Attention mechanism · Joint learning

1 Introduction

Spoken language understanding (SLU) plays a critical role in the maintenance of goal-oriented dialog systems.

Pengfei Wei
School of Computers, Guangdong University of Technology,
Guangzhou, 510006, China
E-mail: wpf@mail2.gdut.edu.cn

Bi Zeng
School of Computers, Guangdong University of Technology,
Guangzhou, 510006, China
E-mail: zb9215@gdut.edu.cn

Wenxiong Liao
School of Computers, Guangdong University of Technology,
Guangzhou, 510006, China
E-mail: lwx@mail2.gdut.edu.cn

Sentence	play	techno	on	lastfm
Slots	O	B-genre	O	B-service
Intent	PlayMusic			

Table 1: An example with intent and slot annotation (BIO format), which indicates the slot of movie name from an utterance with an intent *PlayMusic*.

The SLU module takes user utterance as input and performs three tasks: domain determination, intent detection, and slot filling [11]. Among them, the first two tasks are often framed as a classification problem, which infers the domain or intent (from a predefined set of candidates) based on the current user utterance [27]. For example, the sentence “*play techno on lastfm*” sampled from the SNIPS corpus is shown in Table 1. It can be seen that each word in the sentence corresponds to one slot label, while a specific intent is assigned for the whole sentence.

In early research, Intent detection and slot filling are usually carried out separately, which is called traditional pipeline methods. Intent detection is regarded as an utterance classification problem to predict an intent label, which can be modeled using conventional classifiers, including regression, support vector machine (SVM) [9] or recurrent neural network (RNN) [19]. The slot filling task can be formulated as a sequence labeling problem, and the most popular approaches with good performances are conditional random field (CRF) [26] and long short-term memory (LSTM) networks [35].

Considering this strong correlation between the two tasks, the tendency is to develop a joint model [8, 21, 22, 37]. However, all these models only applied a joint loss function to link the two tasks implicitly. [11] introduce an RNN-LSTM model where the explicit relationships between the intent and slots are not established.

Subsequently, [7], [1], and [20] proposed the gate/mask mechanism to explore incorporating the intent information for slot filling. [24] adopt the token-level intent detection for the Stack-Propagation framework, which can directly use the intent information as input for slot filling. Recently, some work begins to model the bi-directional interrelated connections for the two tasks. [36] proposed a capsule-based neural network model that accomplishes slot filling and intent detection via a dynamic routing-by-agreement schema. [10] proposed an SF-ID network to establish direct connections for the two tasks to help them promote each other mutually.

We apply the proposed approach to ATIS and SNIPS datasets from [4] and [7], separately. Our experiments show that our approach outperforms multiple baselines. We further demonstrate that using BERT representations [6] boosts the performance a lot. The contributions of this paper can be summarized as follows: (1) Establishing the interrelated mechanism among intent nodes and slot nodes in an utterance by a graph attention neural network (GAT) structure. (2) We establish a novel wheel graph to incorporate better the semantic knowledge and make our joint model more interpretable. (3) Showing the effectiveness of our model on two benchmark datasets. (4) We examine and analyze the effect of incorporating BERT in SLU tasks.

2 Related Works

In this section, we will introduce the related works about SLU and GNN in detail.

2.1 Spoken Language Understanding

Separate Model The intent detection is formulated as a text classification problem. The traditional method is to employ n-grams as features with generic entities, such as locations and dates [37]. This type of approach is restricted to the dimensionality of the input space. Another line of popular approaches is to train machine learning models on labeled training data, such as support vector machine (SVM) and Adaboost [9, 29]. Approaches based on deep neural network technology have shown excellent performance, such as Deep belief networks (DBNs) and RNNs [25, 5]. Slot filling can be treated as a sequence labeling task. The traditional method based on conditional random fields (CRF) architecture, which has a strong ability on sequence labeling tasks [26]. Another line of popular approaches is CRF-free sequential labeling. [35] introduced LSTM architecture for this task and obtained a marginal im-

provement over RNN. [30] and [31] introduce the self-attention mechanism for slot filling.

Implicit Joint Model Recently, there have been some joint models to overcome the error propagation caused by the pipelined approaches, and all these models only applied share parameters a joint loss function to link the two tasks implicitly. [11] proposed an RNN-LSTM architecture for joint modeling of intent detection and slot filling. [37] first proposed the joint work using RNNs for learning the correlation between intent and semantic slots of a sentence. [21] proposed an attention-based neural network model for joint intent detection and slot filling, which further explores different strategies in incorporating this alignment information into the encoder-decoder framework. All these models outperform the pipeline models by mutual enhancement between two tasks. However, these joint models didn't model their correlation.

Unidirectional related Joint Model Recently, some works have explored unidirectional related joint models. These models have exploited the intent information for slot filling. [20] proposed a novel intent-augmented gate mechanism to utilize the semantic correlation between intent and slots fully. [7] proposed a slot gate that focuses on learning the relationship between intent and slot attention vectors to obtain better semantic frame results by global optimization. [2] utilize a mask gating mechanism to model the relationship between intent detection and slot filling. [24] perform the token-level intent detection for the Stack-Propagation framework to better incorporate the intent information.

Interrelated Joint Model Considering this strong correlation between the two tasks, interrelated joint models have been explored recently. [34] introduce their cross-impact to each other using two correlated bidirectional LSTMs (BLSTM) to perform the intent detection and slot filling tasks jointly. [10] introduce an SF-ID network to establish direct connections for two tasks to help them promote each other mutually. [36] proposed a capsule-based neural network that models hierarchical relationships among word, slot, and intent in an utterance via a dynamic routing-by-agreement schema.

2.2 Graph Neural Networks

Applying graph neural networks (GNN) to solve some problems has been a popular approach recently in social network analysis [13], knowledge graphs [12], urban computing, and many other research areas [33, 16].

GNN can model non-Euclidean data, while traditional neural networks can only model regular grid data.

Unlike previously discussed neural network-based methods, our approach explicitly establishes direct connections among intent nodes and slots nodes by GAT [33], which uses weighted neighbor features with feature dependent and structure-free normalization, in the style of attention. Analogous to multiple channels in ConvNet [18], GAT introduces multi-head attention [32] to enrich the model capacity and to stabilize the learning process. Unlike other models [10,36], our model does not need to set the number of iterations during training. We have also established a wheel graph structure to learn context-aware information in an utterance better.

3 Proposed Approaches

In this section, we will introduce our wheel-graph graph attention model for SLU tasks. The architecture of the model is shown in Figure 1. First, we show how to use a text encoder to represent an utterance, which can grasp the shared knowledge between two tasks. Second, we introduce the graph attention network (GAT) user weighted neighbor features with feature dependent and structure-free normalization, in the style of attention. Next, the wheel-graph attention network performs an interrelation connection fusion learning of the intent nodes and slot nodes. Finally, intent detection and slot filling are optimized simultaneously via a joint learning schema.

3.1 Text Encoder

Word Embedding: Given a sequence of words, we first covert each word as embedding vector \mathbf{e}_t , and the sequence is represented as $[\mathbf{e}_1, \dots, \mathbf{e}_T]$, where T is the number of words in the sentence.

Affine Transformation: We perform an affine transformation on the embedding sequence, which is a data standardization method.

$$\mathbf{x}_t = \mathbf{W}\mathbf{e}_t + \mathbf{b} \quad (1)$$

where \mathbf{W} and \mathbf{b} are trainable weights and biases.

Two-Layer BiGRU: As an extension of conventional feed-forward neural networks, it was difficult to train Recurrent neural networks (RNNs) to capture long-term dependencies because the gradients tend to either vanish or explode. Therefore, some more sophisticated activation functions with gating units were designed.

Two revolutionary methods are long short-term memory (LSTM) [15] and gated recurrent unit (GRU) [3]. Similarly to the LSTM unit, the GRU has gating units that modulate the flow of information inside the unit; however, without having a separate memory cells and has less parameters. Based on this, we use GRU in this work.

$$\mathbf{r}_t = \sigma(\mathbf{W}_r\mathbf{x}_t + \mathbf{U}_r\mathbf{h}_{t-1}) \quad (2)$$

$$\mathbf{z}_t = \sigma(\mathbf{W}_z\mathbf{x}_t + \mathbf{U}_z\mathbf{h}_{t-1}) \quad (3)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{r}_t \odot (\mathbf{U}\mathbf{h}_{t-1})) \quad (4)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t \quad (5)$$

where \mathbf{x}_t is the input at time t , \mathbf{r}_t and \mathbf{z}_t are reset gate and update gate respectively, \mathbf{W} and \mathbf{U} are weight matrices, σ is sigmoid function and \odot is an element-wise multiplication. When the reset gate is off (\mathbf{r}_t close to 0), the reset gate effectively makes the unit act as if it is reading the first symbol of an input sequence, allowing it to forget the previously computed state. For simplification, the above equations are abbreviated with $\mathbf{h}_t = GRU(\mathbf{x}_t, \mathbf{h}_{t-1})$.

To consider both past and future information at the same time. Consequently, we use a two-Layer bidirectional GRU (BiGRU) to learn the utterance representations at each time step. The BiGRU, a modification of the GRU, consists of a forward and a backward GRU. The layer reads the affine transformed output vectors $[\mathbf{x}_1, \dots, \mathbf{x}_T]$ and generates T hidden states by concatenating the forward and backward hidden states of BiGRU:

$$\vec{\mathbf{h}}_t = \overrightarrow{GRU}(\mathbf{x}_t, \vec{\mathbf{h}}_{t-1}) \quad (6)$$

$$\overleftarrow{\mathbf{h}}_t = \overleftarrow{GRU}(\mathbf{x}_t, \overleftarrow{\mathbf{h}}_{t+1}) \quad (7)$$

$$\overleftrightarrow{\mathbf{h}}_t = [\vec{\mathbf{h}}_t, \overleftarrow{\mathbf{h}}_t] \quad (8)$$

where $\vec{\mathbf{h}}_t$ is the hidden state of forward pass in BiGRU, $\overleftarrow{\mathbf{h}}_t$ is the hidden state of backward pass in BiGRU and $\overleftrightarrow{\mathbf{h}}_t$ is the concatenation of the forward and backward hidden states at time t .

In summary, to get more fine-grained sequence information, we use a two-layer BiGRU to encode input information. The representation is defined as:

$$\overleftrightarrow{\mathbf{h}}_t = BiGRU(BiGRU(\mathbf{x}_t)) \quad (9)$$

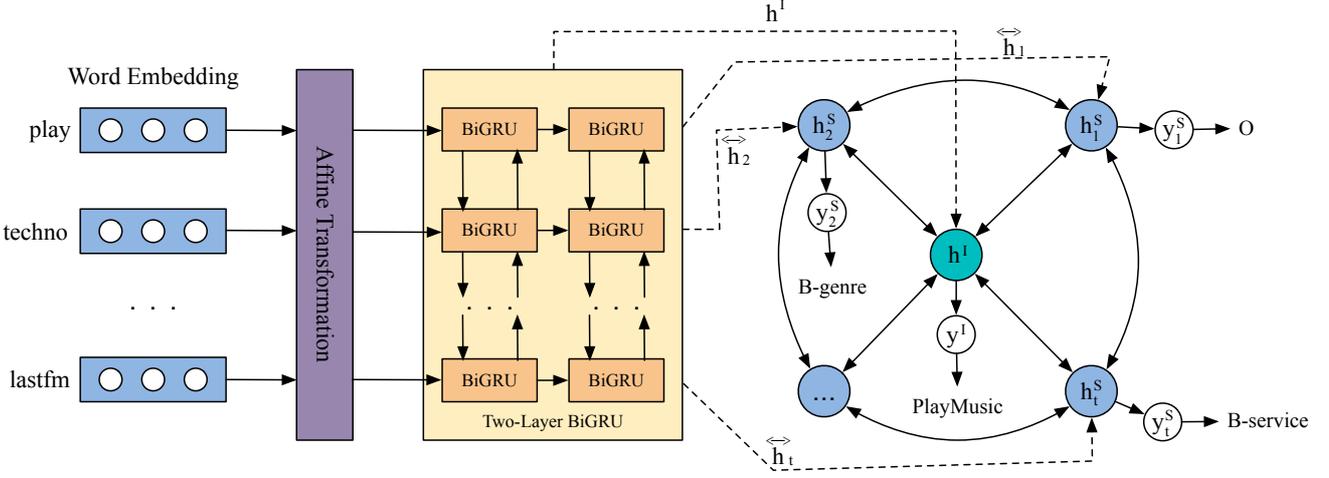


Fig. 1: The overall architecture of the proposed model based on Wheel-Graph attention networks.

3.2 Graph Attention Network

A graph attention network (GAT) [33] is a variant of graph neural network [28] and is an important element in our proposed method. It propagates the intent or slot information from a one-hop neighborhood. Given a dependency graph with N nodes, where each node is associated with a local vector \mathbf{x} , one GAT layer compute node representations by aggregating neighborhood's hidden states.

GAT exploits the attention mechanism as a substitute for the statically normalized convolution operation. Below are the equations to compute the node embedding $\mathbf{h}_i^{(l+1)}$ of layer $l+1$ from the embeddings of layer l .

$$\mathbf{z}_i^{(l)} = \mathbf{W}^{(l)} \mathbf{h}_i^{(l)} \quad (10)$$

$$e_{ij}^{(l)} = f(\vec{\mathbf{a}}^{(l)T} (\mathbf{z}_i^{(l)} \parallel \mathbf{z}_j^{(l)})) \quad (11)$$

$$\alpha_{ij}^{(l)} = \frac{\exp(e_{ij}^{(l)})}{\sum_{k \in N(i)} \exp(e_{ik}^{(l)})} \quad (12)$$

$$\mathbf{h}_i^{(l+1)} = \sigma \left(\sum_{j \in N(i)} \alpha_{ij}^{(l)} \mathbf{z}_j^{(l)} \right) \quad (13)$$

where $\mathbf{W}^{(l)}$ is a linear transformation matrix for input states, \parallel represents vector concatenation, $\vec{\mathbf{a}}^{(l)}$ is an attention context vector learned during training, and \cdot^T represents transposition. $f(\cdot)$ is a LeakyReLU non-linear function [23]. $N(i)$ is the neighbor nodes of node i . σ is the activation function such as \tanh . For simplification, the above equations are abbreviated with $\mathbf{h}_i^{(l+1)} = \text{GAT}(\mathbf{h}^{(l)})$.

3.3 Wheel-Graph Attention Network

In the SLU task, there is a strong correlation between intent detection and slot filling. To make full use of the correlation between intent and slot, we constructed a wheel-graph structure. In Figure 1, this wheel-graph structure contains an intent node and slot nodes.

For the node representation, we use the output of the previous two-layer BiGRU, and the formula is expressed as:

$$\mathbf{h}_0^I = \max_{i=1}^T \overleftarrow{\mathbf{h}}_t \quad (14)$$

where the max function is an element-wise function, and T is the number of words in the utterance. We use \mathbf{h}_0^I as the representation of the intent node and $\overleftarrow{\mathbf{h}}_t$ as the representation of the slot nodes.

For the edge, we created a bidirectional connection between the intent node and the slot nodes. To make better use of the context information of the utterance, we created a bidirectional connection between the slot nodes and connected the head and tail of the utterance to form a loop.

In summary, the feed-forward process of our wheel-graph neural network can be written as:

$$\mathbf{h}_m = [\mathbf{h}_0^I, \overleftarrow{\mathbf{h}}_t] \quad (15)$$

$$\mathbf{h}_m^{(l+1)} = \text{GRU}(\text{GAT}(\mathbf{h}_m^{(l)}), \mathbf{h}_m^{(l)}) \quad (16)$$

$$\mathbf{h}^I, \mathbf{h}_t^S = \mathbf{h}_0^{(l+1)}, \mathbf{h}_{1:m}^{(l+1)} \quad (17)$$

where $m \in 0, 1, \dots, t$, \mathbf{h}^I is the hidden state output of the intent, and \mathbf{h}_t^S is the hidden state output of the slots.

3.4 Joint Intent Detection and Slot Filling

The last layer is the output layer. We adopt a joint learning method. The softmax function is applied to representations with a linear transformation to give the probability distribution \mathbf{y}^I over the intent labels and the distribution \mathbf{y}_t^S over the t -th slot labels. Formally,

$$\mathbf{y}^I = \text{softmax}(\mathbf{W}^I \mathbf{h}^I + \mathbf{b}^I) \quad (18)$$

$$\mathbf{y}_t^S = \text{softmax}(\mathbf{W}^S \mathbf{h}_t^S + \mathbf{b}^S) \quad (19)$$

$$o^I = \text{argmax}(\mathbf{y}^I) \quad (20)$$

$$o_t^S = \text{argmax}(\mathbf{y}_t^S) \quad (21)$$

where \mathbf{W}^I and \mathbf{W}^S are trainable parameters of the model, \mathbf{b}^I and \mathbf{b}^S are bias vectors. o^I and o_t^S are the predicted output labels for intent and slot task respectively.

Then we define loss function for our model. We use \hat{y}^I and \hat{y}^S to denote the ground truth label of intent and slot.

The loss function for intent is a cross-entropy cost function.

$$\mathcal{L}_1 = - \sum_{i=1}^{n_I} \hat{y}^{i,I} \log(y^{i,I}) \quad (22)$$

Similarly, the loss function of a slot label sequence is formulated as:

$$\mathcal{L}_2 = - \sum_{t=1}^T \sum_{i=1}^{n_S} \hat{y}_t^{i,S} \log(y_t^{i,S}) \quad (23)$$

where n_I is the number of intent label types, n_S is the number of slot label types and T is the number of words in an utterance.

The training objective of the model is minimizing a united loss function:

$$\mathcal{L}_\theta = \alpha \mathcal{L}_1 + (1 - \alpha) \mathcal{L}_2 \quad (24)$$

where α is a weight factor to adjust the attention paid to two tasks.

4 Experiments

In this section, we describe our experimental setup and report our experimental results.

Datasets	ATIS	SNIPS
# Train	4,478	13,084
# Validation	500	700
# Test	893	700
# Intents	21	7
# Slots	120	72
Vocab Size	722	11,241
Avg. Length	11.28	9.05

Table 2: Datasets overview.

4.1 Experimental Setup

For experiments, we utilize two datasets, including ATIS [14] and SNIPS [4], which is collected by Snips personal voice assistant in 2018. They are two public benchmark single-intent datasets, which are widely used as benchmark in SLU research. Compared to the single-domain ATIS dataset, SNIPS is more complicated, mainly due to the intent diversity and large vocabulary. Both datasets used in our paper follows the same format and partition as in [24]. The overview of datasets is listed in Table 2.

To validate the effectiveness of our approach, we compare it to the following baseline approaches. It is noted that the results of some models are directly taken from [24].

- **Joint Seq** applies an RNN-LSTM architecture for slot filling, and the last hidden state of LSTM is used to predict the intent of the utterance [11].
- **Attention BiRNN** adopts an attention-based RNN model for joint intent detection and slot filling. Slot label dependencies are modeled in the forward RNN. A max-pooling over time on the hidden states is used to perform the intent classification [22].
- **Slot-Gated Full Atten.** utilizes a slot-gated mechanism that focuses on learning the relationship between intent and slot attention vectors. The intent attention context vector is used for the intent classification [7].
- **Self-Attention Model** first makes use of self-attention to produce a context-aware representation of the embedding. Then a bidirectional recurrent layer takes as input the embeddings and context-aware vectors to produce hidden states. Finally, it exploits the intent-augmented gating mechanism to match the slot label [20].
- **Bi-Model** is a new Bi-model based RNN semantic frame parsing network structure which performs the intent detection and slot filling tasks jointly by considering their cross-impact to each other using two

correlated bidirectional LSTMs [34].

- **SF-ID Network** is a novel bi-directional interrelated model for joint intent detection and slot filling. It contains an entirely new iteration mechanism inside the SF-ID network to enhance the bi-directional interrelated connections [10].
- **CAPSULE-NLU** introduces a capsule-based neural network model with a dynamic routing-by-agreement schema to accomplish intent detection and slot filling tasks. The output representations of IntentCaps and SlotCaps are used to intent detection and slot filling, respectively [36].
- **Stack-Propagation** adopts a Stack-Propagation, which directly uses the intent information as input for slot filling and performs the token-level intent detection to further alleviate the error propagation [24].

4.2 Implementation Details

In our experiments, the dimensionalities of the word embedding are 1024 for the ATIS dataset and SNIPS dataset. All model weights are initialized with uniform distribution. The number of hidden units of the BiGRU encoder is set as 512. The number of layers of the GAT model is set to 1. Graph node representation is set to 1024. The weight factor α is set to 0.1. We use the Adam optimizer [17] with an initial learning rate of 10^{-3} , and L2 weight decay is set to 10^{-6} . The model is trained on all the training data with a mini-batch size of 64. In order to enhance our model to generalize well, the maximum norm for gradient clipping is set to 1.0. We also apply the dropout ratio is 0.2 for reducing overfit.

We implemented our model using PyTorch¹ and DGL² on a Linux machine with Quadro p5000 GPUs. For all the experiments, we select the model which works the best on the validation set and evaluate it on the test set.

4.3 Experimental Results

As with Qin et al [24], we adopt three evaluation metrics in the experiments. For the intent detection task, the accuracy is applied. For the slot filling task, the F1-Score is utilized. Besides, the sentence accuracy is used to indicate the general performance of both tasks,

which refers to the proportion of the sentence whose intent and slot are both correctly-predicted in the whole corpus. Table 3 shows the experimental results of the proposed models on ATIS and SNIPS datasets.

We note that the results of unidirectional related joint models are better than implicit joint models like Joint Seq [11] and Attention BiRNN [22], and the results of interrelated joint models are better than unidirectional related joint models like Slot-Gated Full Attention [7] and Self-Attentive Model [20]. That is likely due to the strong correlation between the two tasks. The intent representations apply slot information to intent detection task while the slot representations use intent information in slot filling task. The bi-directional interrelated model helps the two tasks to promote each other mutually.

We also find that our graph-based Wheel-GAT model performs better than the best prior joint model Stack-Propagation Framework. In ATIS dataset, we achieve 0.6% improvement on Intent (Acc), 0.1% improvement on Slot (F1-score) and 0.7% improvement on Sentence (Acc). In the SNIPS dataset, we achieve 0.4% improvement on Intent (Acc), 0.6% improvement on Slot (F1-score), and 0.5% improvement on Sentence (Acc). This indicates the effectiveness of our Wheel-GAT model. In the previously proposed model, the iteration mechanism used to set the number of iterations is not flexible on training, and the token-level intent detection increases the output load when the utterance is very long. While our model employed graph-based attention network, which uses weighted neighbor features with feature dependent and structure-free normalization, in the style of attention, and directly takes the explicit intent information and slot information further help grasp the relationship between the two tasks and improve the SLU performance.

4.4 Ablation Study

In this section, to further examine the level of benefit that each component of Wheel-GAT brings to the performance, an ablation study is performed on our model. The ablation study is a more general method, which is performed to evaluate whether and how each part of the model contributes to the full model. We ablate four important components and conduct different approaches in this experiment. Note that all the variants are based on joint learning method with joint loss.

- Wheel-GAT w/o intent \rightarrow slot, where no directed edge connection is added from the intent node to the slot node. The intent information is not explicitly applied to the slot filling task on the graph layer.

¹ <https://github.com/pytorch/pytorch>

² <https://github.com/dmlc/dgl>

Model	ATIS Dataset			SNIPS Dataset		
	Slot (F1)	Intent (Acc)	Sentence (Acc)	Slot (F1)	Intent (Acc)	Sentence (Acc)
Joint Seq [11]	94.3	92.6	80.7	87.3	96.9	73.2
Attention BiRNN [22]	94.2	91.1	78.9	87.8	96.7	74.1
Slot-Gated Full Atten. [7]	94.8	93.6	82.2	88.8	97.0	75.5
Self-Attentive Model [20]	95.1	96.8	82.2	90.0	97.5	81.0
Bi-Model [34]	95.5	96.4	85.7	93.5	97.2	83.8
SF-ID Network [10]	95.6	96.6	86.0	90.5	97.0	78.4
CAPSULE-NLU [36]	95.2	95.0	83.4	91.8	97.3	80.9
Stack-Propagation [24]	95.9	96.9	86.5	94.2	98.0	86.9
Wheel-GAT	96.0*	97.5*	87.2*	94.8*	98.4*	87.4*

Table 3: Comparison results of different methods using Wheel-GAN on ATIS and SNIPS datasets. The numbers with * indicate that the improvement of our model over all baselines is statistically significant with $p < 0.05$ under t-test.

Model	ATIS Dataset			SNIPS Dataset		
	Slot (F1)	Intent (Acc)	Sentence (Acc)	Slot (F1)	Intent (Acc)	Sentence (Acc)
Wheel-GAT	96.0	97.5	87.2	94.8	98.4	87.4
Wheel-GAT w/o intent \rightarrow slot	95.5	97.1	86.9	93.5	98.0	85.7
Wheel-GAT w/o slot \rightarrow intent	95.4	96.8	86.6	93.9	97.9	85.8
Wheel-GAT w/o head \leftrightarrow tail	95.6	97.0	86.9	94.0	97.6	85.8
Wheel-GAT w/o GAT	95.0	96.2	84.3	90.8	96.7	77.6

Table 4: Ablation Study on ATIS and SNIPS datasets. \rightarrow indicates that the intent node points to the edge of the slot node. \leftarrow indicates that the slot node points to the edge of the intent node. \leftrightarrow indicates the edge where the head and tail word nodes are connected in an utterance.

- Wheel-GAT w/o slot \rightarrow intent, where no directed edge connection is applied from the slot node to the intent node. The slot information is not explicitly utilized to the intent detection task on the graph layer.
- Wheel-GAT w/o head \leftrightarrow tail, where no bidirectional edge connection is used between the intent node and the slot node. We only use joint loss for joint model, rather than explicitly establishing the transmission of information between the two tasks.
- Wheel-GAT w/o GAT, where no graph attention mechanism is performed in our model. The message propagation is computed via GCN instead of GAT. GCN introduces the statically normalized convolution operation as a substitute for the attention mechanism.

Table 4 shows the joint learning performance of the ablated model on ATIS and SNIPS datasets. We find that all variants of our model perform well based on our graph structure except Wheel-GAT w/o GAT. As listed in the table, all features contribute to both intent detection and slot filling tasks.

If we remove the intent \rightarrow slot edge from the holistic model, the slot performance drops 0.5% and 1.3%

respectively on two datasets. Similarly, we remove the slot \rightarrow intent edge from the holistic model, the intent performance down a lot respectively on two datasets. The result can be interpreted that intent information and slot information are stimulative mutually with each other. We can see that the added edge does improve performance a lot to a certain extent, which is consistent with the findings of previous work [7, 24, 10].

If we remove the head \leftrightarrow tail edge from the holistic model, we see 0.4% drop in terms of F1-score in ATIS and 0.8% drop in terms of F1-score in SNIPS. We attribute it to the fact that head \leftrightarrow tail structure can better model context-aware information in an utterance.

To verify the effectiveness of the attention mechanism, we remove the GAT and use GCN instead. For GCN, a graph convolution operation produces the normalized sum of the node feature of neighbors. The result shows that the intent performance drops 1.3% and 1.7%, the slot performance drops 1.0% and 4.0%, and the sentence accuracy drops 2.9% and 9.8% respectively on ATIS and SNIPS datasets. We attribute it to the fact that GAT uses weighting neighbor features with feature dependent and structure-free normalization, in the style of attention.

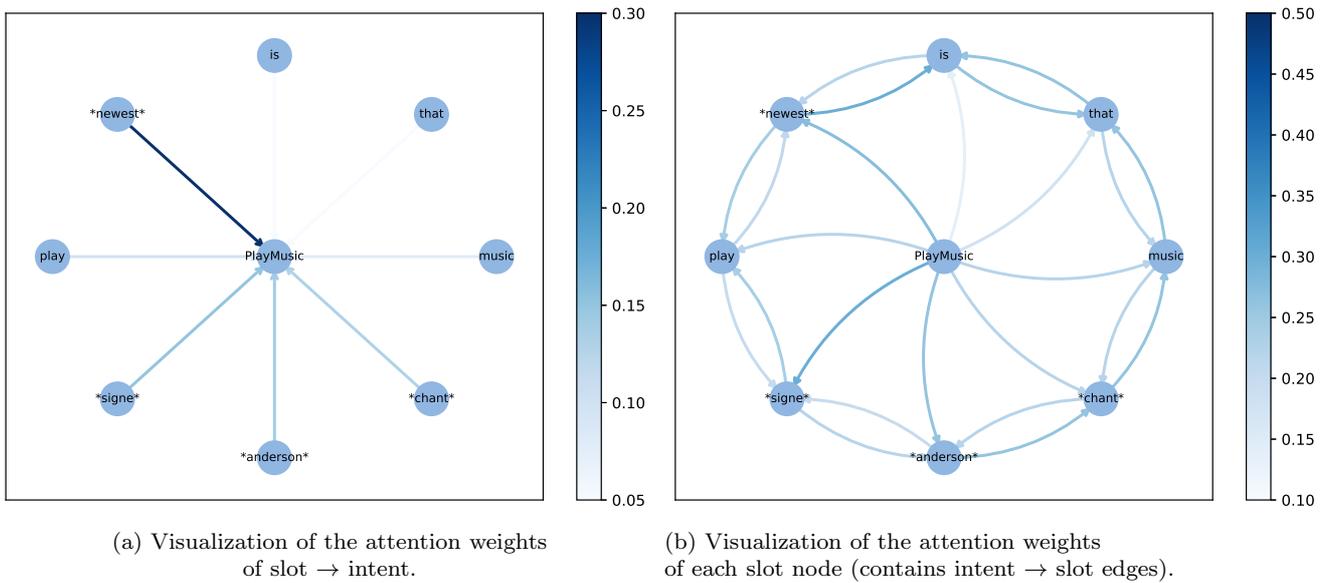


Fig. 2: The central node is intent token and slot tokens are surrounded by *. For each edge, the darker the color, it means that this corresponding of the two nodes is more relevant, so that it integrates more information from this source node features.

4.5 Visualization of Wheel-Graph Attention Layer

In this section, with attempt to better understand what the wheel-graph attention structure has learnt, we visualize the attention weights of slot \rightarrow intent and each slot node, which is shown in Figure 2.

Based on the utterance “*play signe anderson chant music that is newest*”, the intent “PlayMusic” and the slot “O B-artist I-artist B-music_item O O O B-sort”, we can clearly see the attention weights successfully focus on the correct slot, which means our wheel-graph attention layer can learn to incorporate the specific slot information on intent node in Figure 2a. In addition, more specific intent token information is also passed into the slot node in Figure 2b, which achieves a fine-grained intent information integration for guiding the token-level slot prediction. Therefore, the node information of intent and slots can be transmitted more effectively through attention weights in our proposed wheel-graph attention interaction layer, and promote the performance of the two tasks at the same time.

4.6 Effect of BERT

In this section, we also experiment with a pre-trained BERT-based [6] model instead of the Embedding layer, and use the fine-tuning approach to boost SLU task performance and keep other components the same as with our model.

As can be seen from Table 5, Stack-Propagation + BERT [24] joint model achieves a new state-of-the-art performance than another without a BERT-based model, which indicates the effectiveness of a strong pre-trained model in SLU tasks. We attribute this to the fact that pre-trained models can provide rich semantic features, which can help to improve the performance on SLU tasks. Wheel-GAT + BERT outperforms the Stack-Propagation + BERT. That is likely due to we adopt explicit interaction between intent detection and slot filling in two datasets. It demonstrates that our proposed model is effective with BERT.

5 Conclusion and Future Work

In this paper, we first applied the graph network to the SLU tasks. And we proposed a new wheel-graph attention network (Wheel-GAT) model, which provides a bidirectional interrelated mechanism for intent detection and slot filling tasks. The intent node and the slot node construct a explicit two-way associated edge. This graph interaction mechanism can provide fine-grained information integration for token-level slot filling to predict the slot label correctly, and it can also provide specific slot information integration for sentence-level intent detection to predict the intent label correctly. The bidirectional interrelated model helps the two tasks promote performance each other mutually.

Model	ATIS Dataset			SNIPS Dataset		
	Slot (F1)	Intent (Acc)	Sentence (Acc)	Slot (F1)	Intent (Acc)	Sentence (Acc)
Wheel-GAT	96.0	97.5	87.2	94.8	98.4	87.4
BERT SLU [2]	96.1	97.5	88.2	97.0	98.6	92.8
Stack-Propagation + BERT [24]	96.1	97.5	88.6	97.0	99.0	92.9
Wheel-GAT + BERT	96.5	98.0	90.2	97.4	99.3	93.6

Table 5: The SLU performance on BERT-based model on ATIS and SNIPS datasets.

We discuss the details of the prototype of the proposed model and introduced some experimental studies that can be used to explore the effectiveness of the proposed method. We first conduct experiments on two datasets ATIS and SNIPS. Experimental results show that our approach outperforms the baselines and can be generalized to different datasets. Then, to investigate the effectiveness of each component of Wheel-GAT in joint intent detection and slot filling, we also report ablation test results in Table 4. In addition, We visualize and analyze the attention weights of slot \rightarrow intent and each slot node. Besides, we also explore and analyze the effect of incorporating a strong pre-trained BERT model in SLU tasks. Our proposed model achieves the state-of-the-art performance.

In future works, our plan can be summarized as follows: (1) We plan to increase the scale of our dataset and explore the efficacy of combining external knowledge with our proposed model. (2) Collecting multi-intent datasets and expanding our proposed model to multi-intent datasets to explore its adaptive capabilities. (3) We plan to introduce reinforcement learning on the basis of our proposed model, and use the reward mechanism of reinforcement learning to improve the performance of the model. (4) Intent detection and slot filling are usually used together, and any task prediction error will have a great impact on subsequent dialog state tracking (DST). How to improve the accuracy of the two tasks while ensuring the stable improvement of the overall evaluation metrics (Sentence accuracy) still needs to be further explored.

Acknowledgements This work is supported by the National Natural Science Foundation of China under Grant No.61876043, 5467-5471 (2019) National Natural Science Foundation of Guangdong Province under Grant No.2018A030313868 and Major Industry-University Research Project of Guangdong Province under Grant No.2016B010108004. The corresponding author of this paper is Bi Zeng.

References

- Chen, M., Zeng, J., Lou, J.: A self-attention joint model for spoken language understanding in situational dialog applications. arXiv preprint arXiv:1905.11393 (2019)
- Chen, Q., Zhuo, Z., Wang, W.: Bert for joint intent classification and slot filling. arXiv preprint arXiv:1902.10909 (2019)
- Cho, K., van Merriënboer, B., Bahdanau, D., Bengio, Y.: On the properties of neural machine translation: Encoder-decoder approaches. In: Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation, pp. 103-111 (2014)
- Coucke, A., Saade, A., Ball, A., Bluche, T., Caulier, A., Leroy, D., Doumouro, C., Gisselbrecht, T., Caltagirone, F., Lavril, T., et al.: Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. arXiv preprint arXiv:1805.10190 (2018)
- Deoras, A., Sarikaya, R.: Deep belief network based semantic taggers for spoken language understanding. In: Interspeech, pp. 2713-2717 (2013)
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171-4186 (2019)
- Goo, C.W., Gao, G., Hsu, Y.K., Huo, C.L., Chen, T.C., Hsu, K.W., Chen, Y.N.: Slot-gated modeling for joint slot filling and intent prediction. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pp. 753-757 (2018)
- Guo, D., Tur, G., Yih, W.t., Zweig, G.: Joint semantic utterance classification and slot filling with recursive neural networks. In: 2014 IEEE Spoken Language Technology Workshop (SLT), pp. 554-559. IEEE (2014)
- Haffner, P., Tur, G., Wright, J.H.: Optimizing svms for complex call classification. In: 2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03)., vol. 1, pp. I-I. IEEE (2003)
- Haihong, E., Niu, P., Chen, Z., Song, M.: A novel bi-directional interrelated model for joint intent detection and slot filling. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 800-804 (2019)
- Hakkani-Tür, D., Tür, G., Celikyilmaz, A., Chen, Y.N., Gao, J., Deng, L., Wang, Y.Y.: Multi-domain joint semantic frame parsing using bi-directional rnn-lstm. In: Interspeech, pp. 715-719 (2016)
- Hamaguchi, T., Oiwa, H., Shimbo, M., Matsumoto, Y.: Knowledge transfer for out-of-knowledge-base entities: a graph neural network approach. In: Proceedings of the 26th International Joint Conference on Artificial Intelligence, pp. 1802-1808 (2017)
- Hamilton, W., Ying, Z., Leskovec, J.: Inductive representation learning on large graphs. In: Advances in neural information processing systems, pp. 1024-1034 (2017)
- Hemphill, C.T., Godfrey, J.J., Doddington, G.R.: The atis spoken language systems pilot corpus. In: Speech

- and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990 (1990)
15. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8), 1735–1780 (1997)
 16. Huang, L., Ma, D., Li, S., Zhang, X., Houfeng, W.: Text level graph neural network for text classification. 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), pp. 3435–3441 (2019)
 17. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
 18. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp. 1097–1105 (2012)
 19. Lai, S., Xu, L., Liu, K., Zhao, J.: Recurrent convolutional neural networks for text classification. In: Twenty-ninth AAAI conference on artificial intelligence (2015)
 20. Li, C., Li, L., Qi, J.: A self-attentive model with gate mechanism for spoken language understanding. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pp. 3824–3833 (2018)
 21. Liu, B., Lane, I.: Attention-based recurrent neural network models for joint intent detection and slot filling. *arXiv preprint arXiv:1609.01454* (2016)
 22. Liu, B., Lane, I.: Joint online spoken language understanding and language modeling with recurrent neural networks. *arXiv preprint arXiv:1609.01462* (2016)
 23. Maas, A.L., Hannun, A.Y., Ng, A.Y.: Rectifier nonlinearities improve neural network acoustic models. In: Proc. icml, vol. 30, p. 3 (2013)
 24. Qin, L., Che, W., Li, Y., Wen, H., Liu, T.: A stack-propagation framework with token-level intent detection for spoken language understanding. 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), pp. 2078–2087 (2019)
 25. Ravuri, S., Stolcke, A.: Recurrent neural network and lstm models for lexical utterance classification. In: Sixteenth Annual Conference of the International Speech Communication Association (2015)
 26. Raymond, C., Riccardi, G.: Generative and discriminative algorithms for spoken language understanding. In: Eighth Annual Conference of the International Speech Communication Association (2007)
 27. Sarikaya, R., Hinton, G.E., Deoras, A.: Application of deep belief networks for natural language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* **22**(4), 778–784 (2014)
 28. Scarselli, F., Gori, M., Tsoi, A.C., Hagenbuchner, M., Monfardini, G.: The graph neural network model. *IEEE Transactions on Neural Networks* **20**(1), 61–80 (2008)
 29. Schapire, R.E., Singer, Y.: Boostexter: A boosting-based system for text categorization. *Machine learning* **39**(2-3), 135–168 (2000)
 30. Shen, T., Jiang, J., Zhou, T., Pan, S., Long, G., Zhang, C.: Disan: Directional self-attention network for rnn/cnn-free language understanding (2018)
 31. Tan, Z., Wang, M., Xie, J., Chen, Y., Shi, X.: Deep semantic role labeling with self-attention. In: AAAI (2018)
 32. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems, pp. 5998–6008 (2017)
 33. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: International Conference on Learning Representations (2018)
 34. Wang, Y., Shen, Y., Jin, H.: A bi-model based rnn semantic frame parsing model for intent detection and slot filling. In: Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pp. 309–314 (2018)
 35. Yao, K., Peng, B., Zhang, Y., Yu, D., Zweig, G., Shi, Y.: Spoken language understanding using long short-term memory neural networks. In: 2014 IEEE Spoken Language Technology Workshop (SLT), pp. 189–194. IEEE (2014)
 36. Zhang, C., Li, Y., Du, N., Fan, W., Philip, S.Y.: Joint slot filling and intent detection via capsule neural networks. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 5259–5267 (2019)
 37. Zhang, X., Wang, H.: A joint model of intent determination and slot filling for spoken language understanding. In: IJCAI, vol. 16, pp. 2993–2999 (2016)