

A Transformer-based Threshold-Free Framework for Multi-Intent NLU

Lisung Chen¹, Nuo Chen¹, Yuexian Zou¹, Yong Wang², Xinzhong Sun²

¹ADSPLAB, School of ECE, Peking University, Shenzhen, China

²AOTO Electronics Co., Ltd

Abstract

Multi-intent natural language understanding (NLU) has recently gained attention. It detects multiple intents in an utterance, which is better suited to real-world scenarios. However, the state-of-the-art joint NLU models mainly detect multiple intents on threshold-based strategy, resulting in one main issue: the model is extremely sensitive to the threshold settings. In this paper, we propose a transformer-based Threshold-Free Multi-intent NLU model (TFMN) with multi-task learning (MTL). Specifically, we first leverage multiple layers of a transformer-based encoder to generate multi-grain representations. Then we exploit the information of the number of multiple intents in each utterance without additional manual annotations and propose an auxiliary detection task: Intent Number detection (IND). Furthermore, we propose a threshold-free intent multi-intent classifier that utilizes the output of IND task and detects the multiple intents without depending on the threshold. Extensive experiments demonstrate that our proposed model achieves superior results on two public multi-intent datasets.

1 Introduction

Natural language understanding (NLU) consists of two sub-tasks, including intent detection (ID) and slot filling (SF) which allow the dialogue system to create a semantic frame that summarizes the user’s requests. Early works often approach these two tasks separately (McCallum et al., 2000; Sarikaya et al., 2011; Yao et al., 2014; Vu, 2016). Considering intent detection and slot filling are highly related, recent works tend to model these two tasks jointly, where the correlation between the intent and slots are utilized (Goo et al., 2018; E et al., 2019; Qin et al., 2019; Zhou et al., 2021).

The works above only consider the scenario where each utterance has one intent. However, in real-life situations, users may express multi-

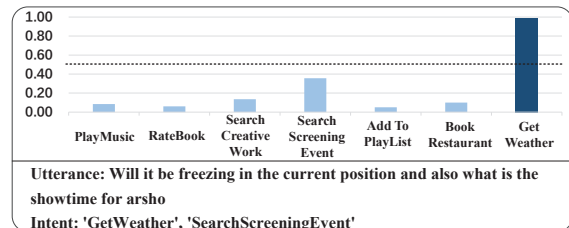


Figure 1: A threshold-based multi-intent detection example in MixSNIPS with given utterance and intent labels. Threshold, which is the dash line, is set to 0.5.

ple intents in an utterance, thus making it difficult to apply single intent NLU models. Recently, several works have studied Multi-intent NLU problem. Gangadharaiyah et al.(2019) investigated an attention-based neural network. Qin et al.(2020) proposed an Adaptive Graph Interactive Framework (AGIF). Qin et al.(2021) explored a non-autoregressive approach to speed up the inference time. Chen et al.(2022a) proposed a Self-distillation Joint NLU model. However, these works all predict multiple intents with threshold, where the common practice is estimating label-instance probabilities and picking the intent labels whose probabilities are higher than the threshold value. We named them threshold-based models. The main issue of threshold-based models is that they are not robust to the threshold settings. As shown in Figure 1, the correct intents for the utterance are 'GetWeather' and 'SearchScreeningEvent'. Although the model can detect that 'GetWeather' and 'SearchScreeningEvent' are the two most probable intents, the threshold-based model only considers 'GetWeather' as the intent due to the threshold which is usually set as 0.5.

In this paper, we propose a transformer-based Threshold-free Multi-NLU model (TFMN) and detect multiple intents without relying on the threshold. Specifically, we leverage the upper layers of a transformer-based encoder to generate multi-grain representations. Next, we fully exploit the annota-

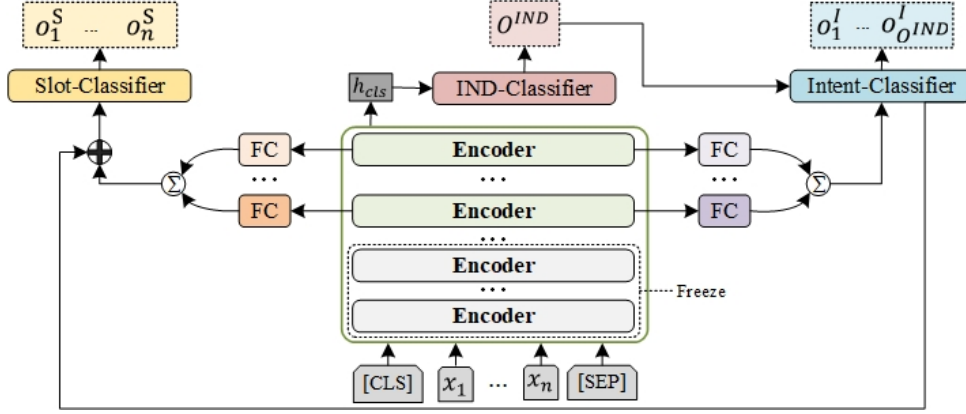


Figure 2: The architecture of transformer-based TFMN model.

tions from original multiple intents data and propose an Intent Number Detection (IND) task. The motivation is to allow the model to detect the intent numbers in a given utterance. Then we propose a threshold-free intent classifier that utilizes the output of IND task to detect the multiple intents.

We validate TFMN on two public datasets (Qin et al., 2020): MixATIS and MixSNIPS, and show that our method outperforms competitive baselines. The contributions of our work are summarized as follows: (1) We propose a novel threshold-free Multi-NLU model based on transformers. (2) We propose IND task, a feasible task to improve the multi-intent NLU without additional manual annotation, and a threshold-free multi-intent classifier that detects multiple intents without relying on threshold. (3) We present extensive experiments demonstrating the effectiveness of our approach.

2 Problem Formulation

Given an input sequence $X = (x_1, \dots, x_n)$, multi-intent detection is defined as a multi-label classification task that outputs $O^I = (o_1^I, \dots, o_m^I)$, where m is the number of predicted intent labels. Slot filling task can be regarded as a sequence labeling task that outputs a slot sequence $O^S = (o_1^S, \dots, o_n^S)$.

3 Approach

In this section, we first introduce the architecture of TFMN model, then detail the proposed IND task and threshold-free intent classifier.

3.1 Threshold-free Multi-intent NLU Model

The architecture of our model is illustrated in Figure 2. TFMN includes a transformer-based encoder with L layers and three task-specific classifiers.

Multiple Intent Detection Following (Qin et al., 2019), we perform a token-level multiple intent detection which can be formalized as a sequence labeling problem (You et al., 2020, 2021b; Chen et al., 2021a, 2022b) that maps the input utterance $X = (x_1, \dots, x_n)$ to sequence of intent label $O^I = (o_1^I, \dots, o_n^I)$. According to (Jawahar et al., 2019; Rogers et al., 2020; Chen et al., 2021b), transformer-based encoder tends to capture syntactic information in the middle and semantic information at the top layers. Therefore, we take the top j layers of the encoder to form multi-grain intent features. First, we map each hidden layer into a different feature space via a fully connected layer, then we combine hidden layers by adding them together:

$$h^I = \sum_{n=L-j}^L w_n^I h_n \quad (1)$$

where w_n^I are trainable parameters and h_n are different hidden layers. We then generate intent logits with the intent feature h^I :

$$l^I = w_i h^I \quad (2)$$

where w_i are trainable parameters. The intent logits will be used to provide token-level intent information for slot filling and detect the final multiple intent labels which we will detail in Section 3.3.

Slot Filling Similar to intent detection, We leverage the top j layers of a transformer-based encoder for slot filling. The slot features h^S are generated by combining hidden layers and concatenating with token-level intent information:

$$h_{temp}^S = \sum_{n=L-j}^L w_n^S h_n \quad (3)$$

$$h^S = h_{temp}^S \oplus l^I \quad (4)$$

then slot classifier computes the slot prediction:

$$p_t^S = \text{softmax}(w_a \text{LeakyReLU}(w_b h_t^S)) \quad (5)$$

where w_a and w_b are trainable parameters.

3.2 Intent Number Detection

To achieve threshold-free multi-intent detection, we propose an Intent Number Detection task which trains with the intent detection and slot filling in a multi-task fashion. In IND task, we fully utilize the original intent label annotations by calculating the numbers of intents in each utterance and forming the intent number labels Y^{IND} . Then we train the model to detect how many intents are there in the input utterance with Y^{IND} . Specifically, we take the output of [CLS] token from the last hidden layer h_{cls} as representation for IND task to classify:

$$p^{IND} = \text{softmax}(w_{ind} h_{cls}) \quad (6)$$

$$O^{IND} = \text{argmax}(p^{IND}) \quad (7)$$

We use cross-entropy to optimize IND task:

$$\mathcal{L}_{IND} = - \sum_k y_k^{IND} \log p_k^{IND} \quad (8)$$

3.3 Threshold-free Intent Classifier

Once having the intent logits l^I and being able to predict the intent number with the proposed IND task, we send l^I into a sigmoid activation function:

$$p_t^I = \text{sigmoid}(l_t^I) \quad (9)$$

where p_t^I is the intent probability distribution of t -th token in the utterance. Since the final output should be the utterance-level intent detection, we sum p_t^I up for utterance-level intent probability distribution P^I , and choose the top O^{IND} , which is the predicted intent number of the utterance, most probable intent label as the final result $O^I = (o_1^I, \dots, o_{O^{IND}}^I)$.

3.4 Multi-Task Training

Our model optimizes the parameters jointly. Multiple intent detection is trained with binary cross-entropy and slot filling is trained with cross-entropy. The total loss of TFMN is the weighted sum of three losses:

$$\mathcal{L}_{total} = \alpha \cdot \mathcal{L}_{ID} + \beta \cdot \mathcal{L}_{SF} + \lambda \cdot \mathcal{L}_{IND} \quad (10)$$

with three hyper-parameters α , β , and λ to balance.

4 Experiments

4.1 Datasets

We conduct experiments on two public multi-intent NLU datasets¹. They are MixATIS (Qin et al., 2020) collected from ATIS dataset (Hemphill et al., 1990) with 13162/759/828 utterances for train/validate/test and MixSNIPS (Qin et al., 2020) collected from SNIPS dataset (Coucke et al., 2018) with 39776/2198/2199 utterances for train/validate/test. Both of the datasets have the ratio of sentences with 1~3 intents as [0.3, 0.5, 0.2].

4.2 Experimental Settings

For TFMN, we use the English uncased Bert-Base model (Devlin et al., 2019) which consists of 12 hidden layers, 12 heads, and the hidden size is 768. For fine-tuning, we freeze the bottom half of Bert to save computational memory and empirically choose the top 4 layers to generate representations. The batch size is 128 and the epoch is 80. Adam is used for optimization with learning rate of 2e-5. The hyper-parameters of loss are empirically set as α : β : λ = 0.6: 1: 1 for MixATIS and α : β : λ = 0.7: 0.9: 1 for MixSNIPS. We evaluate the performance of slot filling with F1 score (You et al., 2021a; Chen et al., 2021c), intent detection with accuracy, and the NLU semantic frame parsing with overall accuracy.

4.3 Baselines

We compare our model with both single-intent and multi-intent baselines. For single-intent baselines to handle multi-intent utterances, multiple intent labels are connected with "#" and treated as a single label, named as *concat* version. For multi-intent baselines, they are all threshold-based models, named as *thresh* version. We also obtain our own pre-trained language model (PLM) baseline for comparison, called Bert-baseline. Following (Chen et al., 2019), we obtain the hidden state of the first special token ([CLS]) for detecting multi-intent based on threshold and use hidden states of utterance tokens for slot filling.

4.4 Results

The main results are illustrated in Table 1. We observe that TFMN model outperforms previous state-of-the-art baselines significantly. On slot filling, our model outperforms GL-GIN 1.5% on MixSNIPS. For multiple intent detection, we achieve

¹<https://github.com/LooperXX/AGIF>

Model	MixATIS			MixSNIPS		
	Slot (F1)	Intent (Acc)	Overall (Acc)	Slot (F1)	Intent (Acc)	Overall (Acc)
SF-ID (<i>concat</i>) (2019)	87.4	66.2	34.9	90.6	95.0	59.9
Stack-Propagation (<i>thresh</i> = 0.5) (2019)	87.8	72.1	40.1	94.2	96.0	72.9
Joint Multiple ID-SF (<i>thresh</i> = 0.5) (2019)	84.6	73.4	36.1	90.6	95.1	62.9
AGIF (<i>thresh</i> = 0.5)(2020)	86.7	74.4	40.8	94.2	95.1	74.2
GL-GIN (<i>thresh</i> = 0.5)(2021)	88.3	76.3	43.5	94.9	95.6	75.4
SDJN (<i>thresh</i> = 0.5)(2022a)	88.2	77.1	44.6	94.4	96.5	75.7
SDJN+BERT (<i>thresh</i> = 0.5)(2022a)	87.5	78.0	46.3	95.4	96.7	79.3
Bert-baseline (<i>thresh</i> = 0.3)	83.1	74.8	42.6	95.5	95.7	80.2
Bert-baseline (<i>thresh</i> = 0.5)	86.3	74.5	44.8	95.5	95.6	80.1
Bert-baseline (<i>thresh</i> = 0.8)	85.6	75.8	43.5	95.2	96.7	80.6
TFMN (Bert-base)	88.0	79.8	50.2	96.4	97.7	84.7

Table 1: Slot filling and multiple intent detection results on two multi-intent datasets.

Model	MixATIS		
	Slot (F1)	Intent (Acc)	Overall (Acc)
TFMN	88.0	79.8	50.2
-w/o <i>T-free Cls</i>	87.1	77.3	47.0
-w/o <i>T-free Cls</i> & <i>IND task</i>	86.3	76.8	46.7

Table 2: Ablation study. *T-free Cls* indicates threshold-free intent classifier.

2.7% and 1.2% improvement compared with SDJN on MixATIS and MixSNIPS respectively. On overall accuracy, our model shows strong performance which surpasses SDJN 5.6% on MixATIS and 9% on MixSNIPS. When comparing PLM baselines, we can first observe that different threshold settings affect the results of Bert-baseline distinctively. Second, TFMN model outperforms PLM baselines in all three metrics on both datasets. The results suggest that our approach brings significant improvements to multi-intent NLU. We believe this is due to the proposed IND task which fully exploits original intent annotations and threshold-free intent classifier that allows our model to detect multiple intents without a threshold and lead to performance gains.

4.5 Ablation Study

We compare TFMN with two simplified versions, *-w/o T-free Cls* and *-w/o T-free Cls & IND task* in Table 2 to analyze the effectiveness of threshold-free intent classifier and IND task. We can see that as the threshold-free intent classifier is removed, the performances drop 0.9%, 2.5%, and 3.2% on slot F1, intent accuracy, and overall accuracy respectively. We attribute this to the fact that the threshold-free approach can better detect the intent number in an utterance compare to threshold

Model	MixATIS			
	Int-1	Int-2	Int-3	Avg.
AGIF	96.5	83.7	76.7	85.6
GL-GIN	96.5	94.6	87.5	92.8
SDJN+BERT	97.2	92.0	84.0	91.2
Bert-baseline (<i>thresh</i> = 0.5)	94.4	87.8	83.5	88.6
TFMN	98.6	99.7	99.3	99.2

Table 3: A comparison of intent number prediction between threshold-based and threshold-free approaches. The evaluation metric is accuracy. **Int-#** means the utterance with the number of “#” intent. Avg. is the average accuracy.

strategy. We further remove the INP task and the performance again drops 0.8%, 0.5%, and 0.3% on slot F1, intent accuracy, and overall accuracy respectively. This indicates the effectiveness of introducing the INP task to multi-intent NLU.

4.6 Threshold-based vs Threshold-free

To compare threshold-based and threshold-free approaches, we evaluate how well a model can detect the number of intents in an utterance. The results are demonstrated in Table 3. We obtained that the threshold-free model, TFMN, significantly outperforms the threshold-based baselines. Our model achieves 2.1%, 5.1%, 11.8%, and 6.4% improvements on one to three intent utterances and average accuracy over GL-GIN. We find it interesting that threshold-based models predict intent number well when there is one intent in the utterance and become worse as the intent number increase while TFMN shows more consistency.

5 Conclusion

In this paper, we propose TFMN model which detects intent numbers in an utterance by a novel IND

task that does not require additional manual annotations. Then we propose a threshold-free intent classifier to detect multiple intents without relying on the threshold. Extensive experiments show that TFMN achieves performance gains over strong baselines.

Acknowledgements

This paper was partially supported by Shenzhen Science & Technology Research Program (No: GXWD20201231165807007-20200814115301001) and NSFC (No: 62176008). Special thanks are given to AOTO-PKUSZ Joint Research Center for AI for its support. Yuexian Zou is the corresponding author of this paper.

References

- Lisong Chen, Peilin Zhou, and Yuexian Zou. 2022a. [Joint multiple intent detection and slot filling via self-distillation](#). In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7612–7616.
- Nuo Chen, Fenglin Liu, Chenyu You, Peilin Zhou, and Yuexian Zou. 2021a. Adaptive bi-directional attention: Exploring multi-granularity representations for machine reading comprehension. In *ICASSP*, pages 7833–7837. IEEE.
- Nuo Chen, Linjun Shou, Min Gong, Jian Pei, and Daxin Jiang. 2021b. [From good to best: Two-stage training for cross-lingual machine reading comprehension](#).
- Nuo Chen, Linjun Shou, Ming Gong, Jian Pei, and Daxin Jiang. 2022b. [Bridging the gap between language models and cross-lingual sequence labeling](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1909–1923, Seattle, United States. Association for Computational Linguistics.
- Nuo Chen, Chenyu You, and Yuexian Zou. 2021c. Self-supervised dialogue learning for spoken conversational question answering. *arXiv preprint arXiv:2106.02182*.
- Qian Chen, Zhu Zhuo, and Wen Wang. 2019. [BERT for joint intent classification and slot filling](#). *CoRR*, abs/1902.10909.
- A. Coucke, A. Saade, Adrien Ball, Théodore Bluche, A. Caulier, D. Leroy, Clément Doumouro, Thibault Gisselbrecht, F. Caltagirone, Thibaut Lavril, Maël Primet, and J. Dureau. 2018. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *ArXiv*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Haihong E, Peiqing Niu, Zhongfu Chen, and Meina Song. 2019. [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*, pages 5467–5471, Florence, Italy. Association for Computational Linguistics.
- Rashmi Gangadharaiah and Balakrishnan Narayanaswamy. 2019. [Joint multiple intent detection and slot labeling for goal-oriented dialog](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 564–569, Minneapolis, Minnesota. Association for Computational Linguistics.
- Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu, and Yun-Nung Chen. 2018. Slot-gated modeling for joint slot filling and intent prediction. In *Proceedings of The 16th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- C. T. Hemphill, J. Godfrey, and G. Doddington. 1990. The atis spoken language systems pilot corpus. In *HLT*.
- Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics.
- Andrew McCallum, Dayne Freitag, and Fernando CN Pereira. 2000. Maximum entropy markov models for information extraction and segmentation. In *Icml*.
- Libo Qin, W. Che, Yangming Li, Haoyang Wen, and T. Liu. 2019. A stack-propagation framework with token-level intent detection for spoken language understanding. In *EMNLP/IJCNLP*.
- Libo Qin, Fuxuan Wei, Tianbao Xie, Xiao Xu, Wanxiang Che, and Ting Liu. 2021. GI-gin: Fast and accurate non-autoregressive model for joint multiple intent detection and slot filling.
- Libo Qin, Xiao Xu, Wanxiang Che, and Ting Liu. 2020. AGIF: An adaptive graph-interactive framework for joint multiple intent detection and slot filling. In

"Findings of the Association for Computational Linguistics: EMNLP 2020". "Association for Computational Linguistics".

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.

Ruhi Sarikaya, Geoffrey E Hinton, and Bhuvana Ramabhadran. 2011. Deep belief nets for natural language call-routing. In *ICASSP*. IEEE.

Ngoc Thang Vu. 2016. Sequential convolutional neural networks for slot filling in spoken language understanding. In *INTERSPEECH*.

Kaisheng Yao, Baolin Peng, Yu Zhang, Dong Yu, Geoffrey Zweig, and Yangyang Shi. 2014. Spoken language understanding using long short-term memory neural networks. In *2014 IEEE Spoken Language Technology Workshop (SLT)*. IEEE.

Chenyu You, Nuo Chen, Fenglin Liu, Dongchao Yang, and Yuexian Zou. 2020. Towards data distillation for end-to-end spoken conversational question answering. *CoRR*, abs/2010.08923.

Chenyu You, Nuo Chen, and Yuexian Zou. 2021a. Knowledge distillation for improved accuracy in spoken question answering. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7793–7797. IEEE.

Chenyu You, Nuo Chen, and Yuexian Zou. 2021b. Self-supervised contrastive cross-modality representation learning for spoken question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 28–39.

Peilin Zhou, Zhiqi Huang, Fenglin Liu, and Yuexian Zou. 2021. Pin: A novel parallel interactive network for spoken language understanding. *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 2950–2957.