

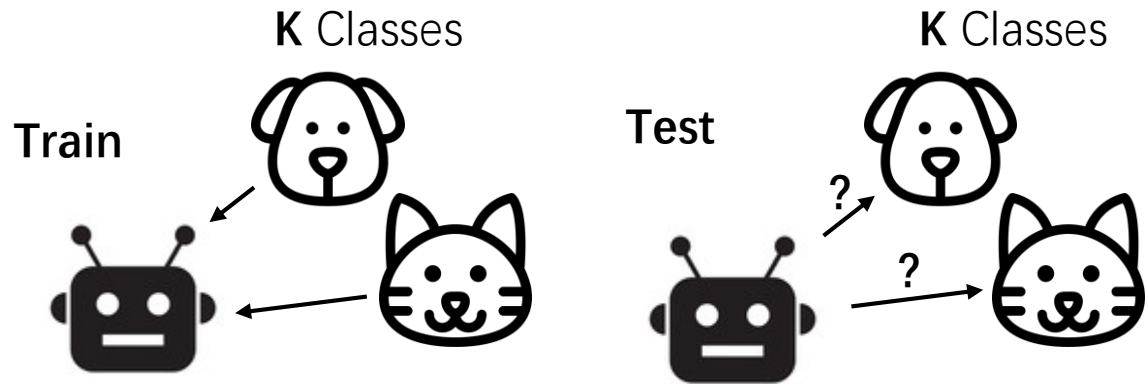
Out-of-Domain Intent Detection Considering Multi-Turn Dialogue Contexts

Hao Lang, Yinhe Zheng, Binyuan Hui, Fei Huang, Yongbin Li

Alibaba Group

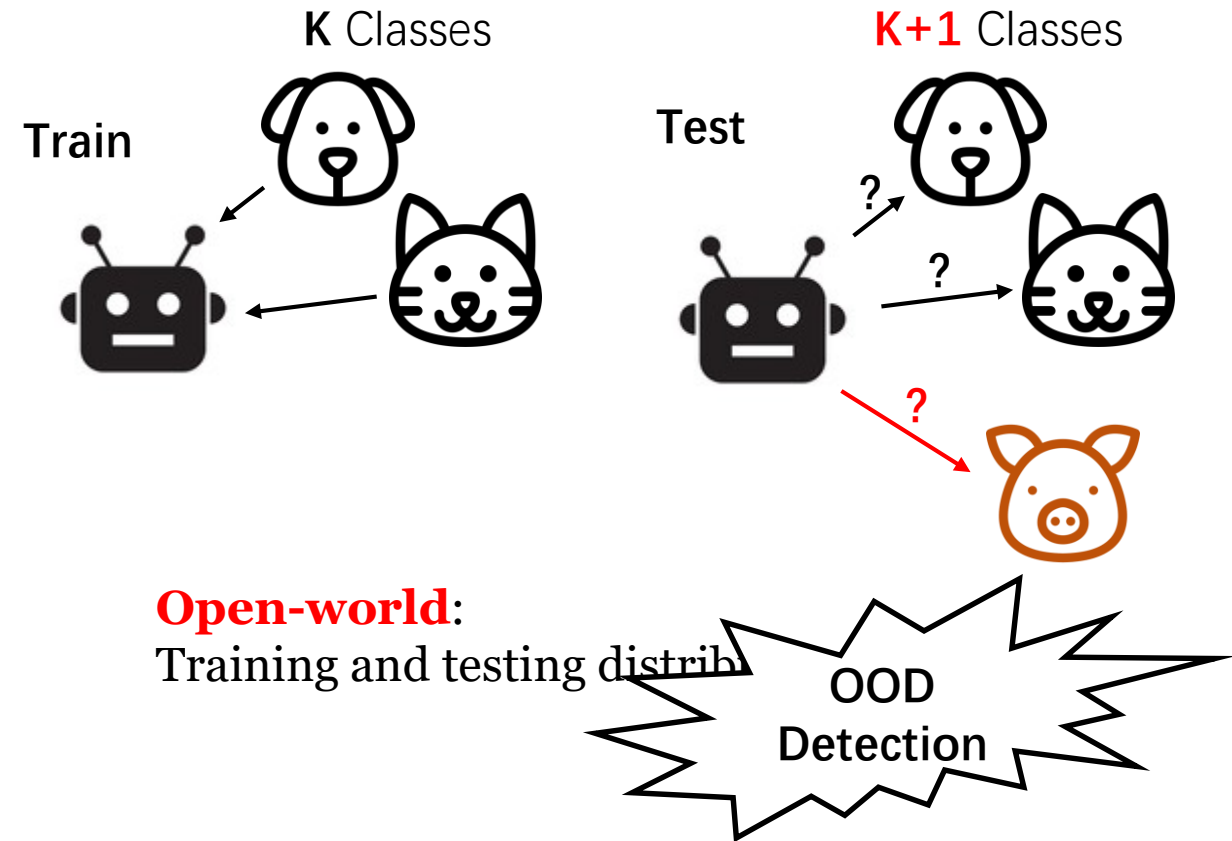
Motivation (1/3)

- Most intent detection model follows the **closed-world assumption**



Closed-world:

Training and testing distributions **match**



Open-world:

Training and testing distrib

Motivation (2/3)

- Most OOD intent detection studies only focus on **single-turn inputs**,
 - i.e., only the most recently issued utterance is taken as the input.
- In real applications, completing a task usually necessitates **multiple turns of conversations**.

Motivation (3/3)

- However, it is non-trivial to directly extend previous methods to the multi-turn setting.
- We usually experience long distance obstacles when modeling multi-turn dialogue contexts,
 - i.e., some dialogues have extremely long histories filled with irrelevant noises for intent detection.
- Meanwhile, it is expensive to construct OOD samples before training when multi-turn contexts are considered.

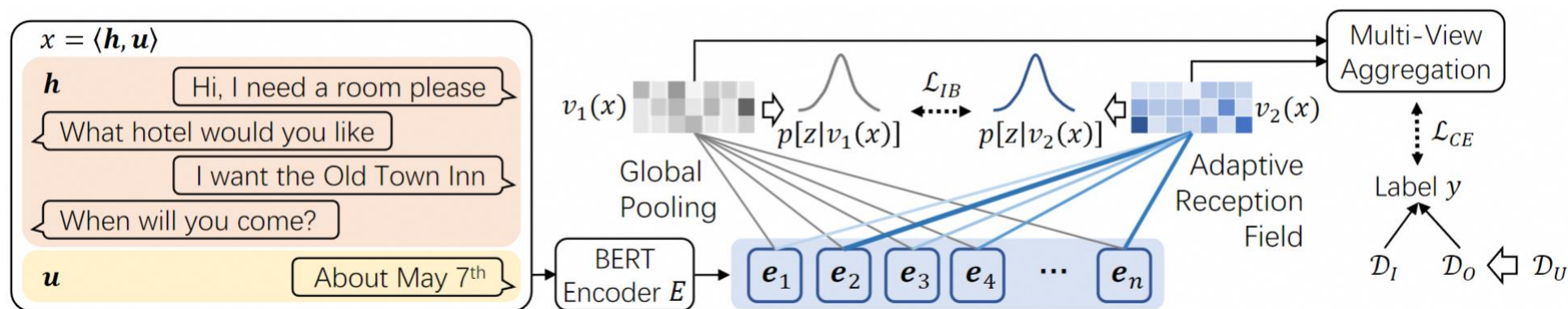
What did we do?

Caro: Context-aware OOD intent detection

Two main challenges to be addressed in Caro:

1. How to alleviate the long distance obstacle and learn robust representations from multi-turn dialogue contexts;
2. How to effectively leverage unlabeled data for OOD intent detection.

How (Overview)



1. Learn robust representations by building diverse views of inputs and optimizing an unsupervised multi-view loss
2. Mine OOD samples from unlabeled data
3. Obtain a $(k + 1)$ -way classifier

How (Representation Learning)

1. Build diverse views of inputs

a) Global Pooling

$$v_1(x) = \sum_{i=1}^n \frac{e_i}{n}$$

b) Adaptive Reception Field

$$v_2(x) = \sum_{i=1}^n \frac{\exp(\alpha_i)}{\sum_{j=1}^n \exp(\alpha_j)} \cdot e_i$$

$$\alpha_i = \sigma(w_i \cdot \text{ReLU}(W_1 \cdot s))$$

2. Optimize information bottleneck loss

$$\mathcal{L}_{IB} = -I(z_1; z_2) + \frac{1}{2} (D_{\text{KL}}[p(z|v_1)||p(z|v_2)] + D_{\text{KL}}[p(z|v_2)||p(z|v_1)])$$

$$p(z|v_i) = \mathcal{N}[\mu(v_i), \Sigma(v_i)]$$

How (OOD Samples Mining)

1. Synthesize pseudo OOD samples \mathcal{D}_P by mixing up IND representations
2. Train a preliminary OOD detector F on $\mathcal{D}_I \cup \mathcal{D}_P$
3. Mine OOD samples \mathcal{D}_O from the unlabeled data \mathcal{D}_U using F

How (OOD Detector Training)

1. Train OOD detector F using \mathcal{L} on \mathcal{D}_I , \mathcal{D}_O , and \mathcal{D}_U

$$\mathcal{L} = \mathbb{E}_{x \in \mathcal{D}_I \cup \mathcal{D}_O} \mathcal{L}_{CE} + \lambda \cdot \mathbb{E}_{x \in \mathcal{D}_U} \mathcal{L}_{IB}$$

Multi-view Aggregation is performed to obtain assembled input representations

$$v(x) = \beta \otimes v_1(x) + (1 - \beta) \otimes v_2(x)$$

$$\beta = \sigma(W_3 \cdot \text{ReLU}(W_2 \cdot (v_1(x) + v_2(x))))$$

Experiments (Datasets)

- We perform experiments on two variants of the STAR dataset (Mosig et al., 2020), i.e., STAR-Full and STAR-Small.
 - STAR is a task-oriented dialogue dataset that has 150 intents.
 - It is designed to model long context dependence, and provides explicit annotations of OOD intents.

	Train		Valid	Test	# Avg. Context Turns
	\mathcal{D}_I	\mathcal{D}_U	\mathcal{D}_V	\mathcal{D}_T	
STAR-Full	15.4K	7.9K	2.8K	2.9K	6.13
STAR-Small	7.7K	3.9K	2.8K	2.9K	6.12

Table 1: Dataset statistics.

Experiments (Main Results)

Model		STAR-Full			STAR-Small		
		F1-All	F1-OOD	F1-IND	F1-All	F1-OOD	F1-IND
Oracle		50.1	64.46	50	46.54	58.23	46.46
\mathcal{D}_I	MSP	40.83	19.74	40.97	37.17	18.1	37.31
	MSP w/o h	17.29	14.12	17.31	17.12	13.49	17.14
	SEG	17.45	6.85	17.53	11.66	7.39	11.69
	SEG w/o h	0.06	2.77	0.04	0.05	2.27	0.04
	DOC	26.53	16.80	26.60	3.47	11.78	3.41
	DOC w/o h	11.31	14.16	11.29	0.08	11.04	0
	ADB	44.64	20.56	44.80	41.36	18.23	41.51
	ADB w/o h	23.27	17.63	23.30	20.08	21.27	20.07
	DAADB	37.27	22.87	37.37	34.81	20.43	34.91
	DAADB w/o h	17.87	15.15	17.88	16.34	17.03	16.33
	Outlier	43.84	19.53	44.01	39.51	19.92	39.64
	Outlier w/o h	23.35	16.75	23.39	19.56	15.42	19.59
	CDA	43.76	5.26	44.03	40.02	10.48	40.22
$\mathcal{D}_I + \mathcal{D}_U$	ASS+MSP	41.97	25.15	42.08	40.85	19.47	40.99
	ASS+LOF	39.87	17.65	40.02	39.54	18.49	39.68
	ASS+GDA	43.73	21.24	43.88	40.86	16.72	41.02
	Caro (ours)	48.75(±1.0)	54.75(±3.2)	48.71(±1.0)	45.02(±1.1)	46.78(±1.8)	45.01(±1.1)

Performance improvement by 5-9% absolutely

Table 2: Performance of Caro and baselines. All results are averages of three runs and the best results are bolded. The standard deviation of the performance of Caro is provided in parentheses.

Experiments (Ablation Study)

Model	STAR-Full			STAR-Small		
	F1-All	F1-OOD	F1-IND	F1-All	F1-OOD	F1-IND
Caro	48.75	54.75	48.71	45.02	46.78	45.01
w/o \mathcal{D}_U	45.97	21.45	46.14	42.24	23.23	42.37
w/o MV	47.71	53.35	47.67	44.42	38.89	44.46
w/o VA	47.34	50.85	47.32	44.14	43.88	44.15
w/o IB	48.23	49.37	48.22	44.14	37.06	44.19

Table 3: Ablation on different components of Caro.

Model	STAR-Full			STAR-Small		
	F1-All	F1-OOD	F1-IND	F1-All	F1-OOD	F1-IND
Caro	48.75	54.75	48.71	45.02	46.78	45.01
InfoMax	47.27	49.92	47.25	44.27	36.66	44.32
MVI	48.46	51.99	48.44	44.70	36.16	44.76
CL	48.18	52.54	48.15	44.59	35.31	44.65
SimCSE	47.73	47.74	47.73	44.30	27.02	44.42

Table 4: Ablation on the representation learning loss.

Experiments (Further Analysis)

Context Len		F1-All	F1-OOD	F1-IND
Long	w/o IB	44.14	37.06	44.19
	w IB	45.02 (+0.88)	46.78 (+9.72)	45.01 (+0.82)
Short	w/o IB	43.61	40.68	43.63
	w IB	43.70 (+0.09)	43.32 (+2.64)	43.70 (+0.07)

Table 5: Benefit of \mathcal{L}_{IB} under different context lengths on the STAR-Small dataset. Long context means retaining all the original dialogue contexts (6 turns on average), and short context means truncating contexts longer than 3 turns. Scores in parentheses is the performance improvement brought by \mathcal{L}_{IB}

Thanks for Your Attention

- If you are interested in our work, feel free to reach out. We are always open to collaboration.
- Email: hao.lang@alibaba-inc.com

**WE'RE
HIRING!**