

Generalized Zero-shot Intent Detection via Commonsense Knowledge

A.B. Siddique, Fuad Jamour, Luxun Xu, and Vagelis Hristidis
msidd005,fuadj,lxu051{@ucr.edu},vagelis@cs.ucr.edu
University of California, Riverside
Riverside, CA, USA

ABSTRACT

Identifying user intents from natural language utterances is a crucial step in conversational systems that has been extensively studied as a supervised classification problem. However, in practice, new intents emerge after deploying an intent detection model. Thus, these models should seamlessly adapt and classify utterances with both seen and unseen intents – unseen intents emerge after deployment and they do not have training data. The few existing models that target this setting rely heavily on the training data of seen intents and consequently overfit to these intents, resulting in a bias to misclassify utterances with unseen intents into seen ones. We propose RIDE: an intent detection model that leverages commonsense knowledge in an unsupervised fashion to overcome the issue of training data scarcity. RIDE computes robust and generalizable *relationship meta-features* that capture deep semantic relationships between utterances and intent labels; these features are linked to those in an intent label via commonsense knowledge. Our extensive experimental analysis on three widely-used intent detection benchmarks shows that relationship meta-features significantly improve the detection of both seen and unseen intents and that RIDE outperforms the state-of-the-art models.

ACM Reference Format:

A.B. Siddique, Fuad Jamour, Luxun Xu, and Vagelis Hristidis. 2021. Generalized Zero-shot Intent Detection via Commonsense Knowledge. In *SIGIR '21: ACM SIGIR Conference on Research and Development in Information Retrieval, July 11–15, 2021, Online*. ACM, New York, NY, USA, 5 pages. <https://doi.org/00.0000/1122445.1122456>

1 INTRODUCTION

Virtual assistants such as Amazon Alexa and Google Assistant allow users to perform a variety of tasks (e.g., Alexa skills) through a natural language interface. For example, a user can set an alarm by simply issuing the utterance “Wake me up tomorrow at 10 AM” to a virtual assistant, and the assistant is expected to understand that the user’s intent (i.e., “AddAlarm”) is to invoke the alarm module, then set the requested alarm accordingly. Intent detection is typically the first step towards performing any task in conversational systems

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGIR '21, July 11–15, 2021, Online

© 2021 Association for Computing Machinery.
ACM ISBN xxx-x-xxxx-XXXX-X/xx/xx...\$15.00
<https://doi.org/00.0000/1122445.1122456>

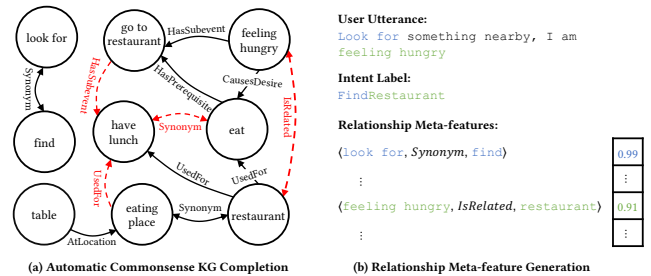


Figure 1: (a) Link Predictor learns from KG (solid lines) to predict missing edges (dashed lines), (b) Example utterance, intent, and computation of relationship meta-features that facilitate GZS intent detection.

and it is a challenging problem due to the vast diversity in user utterances. The challenge is further exacerbated in the more practically relevant setting where intents are added over time. This setting is an instance of the *generalized zero-shot classification problem* [9]: labeled training utterances are available only for seen intents but are unavailable for unseen ones, and at inference time, models do not have prior knowledge on whether the utterances they receive imply seen or unseen intents. This setting is the focus of this paper.

Little research has been conducted on building generalized zero-shot (GZS) models for intent detection, with little success. Earlier works [4, 16, 26] used zero-shot (ZS) learning to train an intent classification model that could classify utterances from unseen intent classes through transferring knowledge from seen classes. The test set in the standard ZS setting is not representative of the real world, as it exclusively includes samples from the unseen classes (as opposed to having samples from both seen and unseen classes as in the GZS setting) at inference time. ZS methods perform poorly in the GZS setting [3, 25], which is primarily caused by their strong bias towards seen classes; ZS intent detection models misclassify almost all test samples from unseen classes into seen ones [20, 32, 33].

To mitigate the issue of training data scarcity for unseen intents and ZS models’ inability to effectively handle the GZS setting, we propose incorporating commonsense knowledge into GZS intent detection model. We argue that such knowledge, if incorporated properly, helps overcome training data scarcity and allows detecting intents regardless of whether they are seen or not, given that commonsense knowledge is uniform across intents. We leverage ConceptNet [27] – a rich and widely-used commonsense knowledge graph (KG) that captures a large subset of knowledge in a semi-structured format (i.e., facts in the form $\langle \text{head}, \text{relation}, \text{tail} \rangle$ such as $\langle \text{apple}, \text{IsA}, \text{fruit} \rangle$). Given that ConceptNet is incomplete, similarly to other KGs, we pre-train a link predictor [14] that learns

from an existing KG to infer novel edges (i.e., relationships) among nodes (i.e., head/tail) to overcome the missing information challenge. Figure 1 (a) presents a toy commonsense KG where a link predictor can learn from existing facts such as $\langle \text{feeling hungry}, \text{CausesDesire}, \text{eat} \rangle$ and $\langle \text{restaurant}, \text{UsedFor}, \text{eat} \rangle$, and infer missing facts such as $\langle \text{feeling hungry}, \text{IsRelated}, \text{restaurant} \rangle$. We infuse the knowledge from our link predictor into our model by extracting *relationship meta-features*. These features quantify the level of relevance between an utterance and an intent in the form of relationship weights, where each weight describes the level of relatedness between the phrases in an utterance and an intent label based on a certain relationship type. Figure 1 (b) shows an example utterance, an intent label, and an inference about the relationships between the phrases in an utterance and an intent label in the form of a relationship meta-features. Relationship meta-features augment our embeddings using commonsense knowledge, which significantly reduces our model’s reliance on the scarcely available seen intents training data. Furthermore, these features reduce our model’s bias towards seen intents given that they are similarly computed for both seen and unseen intents; i.e., they are domain-oblivious.

Our model, RIDE¹, combines relationship meta-features with contextual word embeddings [22], and feeds the combined feature vectors into a trainable prediction function. RIDE is able to accurately detect both seen and unseen intents in utterances. Our extensive experimental analysis using the three widely used benchmarks, SNIPS [6], SGD [23], and MultiWOZ [37] show that our model outperforms the state-of-the-art (SOTA) model in F1 scores on unseen intents in the GZS setting by at least 25.66%. The source code of RIDE is available².

A secondary contribution of this paper is that we managed to further improve the performance of GZS intent detection by employing Positive-Unlabeled (PU) learning [8] to predict if a new utterance belongs to a seen or unseen intent. PU learning assists intent detection models by mitigating their bias towards classifying most utterances into seen intents. A PU classifier is able to perform binary classification after being trained using only positive and unlabeled examples. We found out that the PU classifier also improves the performance of existing intent detection works. Our model, however, outperforms existing ones regardless of the PU classifier integration.

2 PRELIMINARIES

GZS Intent Detection. Let $\mathcal{S} = \{I_1, \dots, I_k\}$ be a set of seen intents and $\mathcal{U} = \{I_{k+1}, \dots, I_n\}$ be a set of unseen intents where $\mathcal{S} \cap \mathcal{U} = \emptyset$. Let $\mathcal{X} = \{X_1, X_2, \dots, X_m\}$ be a set of labeled training utterances where each training utterance $X_i \in \mathcal{X}$ is described with a tuple (X_i, I_j) such that $I_j \in \mathcal{S}$. An intent I_j is comprised of an *Action* and an *Object* and takes the form “ActionObject”³ (e.g., “Find-Restaurant”); an Action describes a user’s request or activity and an Object describes the entity pointed to by an Action [5, 28, 29]. Given a test utterance X'_i , the problem is to predict a label $I'_j \in \mathcal{S} \cup \mathcal{U}$.

Link Prediction in Knowledge Graphs. We pre-train a SOTA link prediction model (LP) [14] on ConceptNet [27] to score novel facts that are not necessarily present in the knowledge graph. Given

Algorithm 1: RMG

Input: $\mathcal{R} = \{r_1, \dots, r_t\}$: relations in KG
 $\mathcal{G}_i = \{g_1, \dots, g_q\}$: utterance n-grams
 $I_j = \{\mathcal{A}, \mathcal{O}\}$: intent’s Action and Object

Output: $e_{relationship}$: X_i - I_j relationship meta-features

Let $e_{X_i}^{\vec{\mathcal{A}}} = \text{RM}(\mathcal{A}, \mathcal{G}_i, \rightarrow)$ // Action to utterance
 Let $e_{X_i}^{\vec{\mathcal{O}}} = \text{RM}(\mathcal{O}, \mathcal{G}_i, \rightarrow)$ // Object to utterance
 Let $e_{X_i}^{\overleftarrow{\mathcal{A}}} = \text{RM}(\mathcal{A}, \mathcal{G}_i, \leftarrow)$ // utterance to Action
 Let $e_{X_i}^{\overleftarrow{\mathcal{O}}} = \text{RM}(\mathcal{O}, \mathcal{G}_i, \leftarrow)$ // utterance to Object

Let $e_{relationship} = [e_{X_i}^{\vec{\mathcal{A}}}, e_{X_i}^{\vec{\mathcal{O}}}, e_{X_i}^{\overleftarrow{\mathcal{A}}}, e_{X_i}^{\overleftarrow{\mathcal{O}}}]$

return $e_{relationship}$

Function $\text{RM}(\text{concept}, \text{phrases}, \text{direction})$:

```

  Let e = []
  foreach r ∈ R do
    if direction = → then
      Let p = Max (LP(concept, r, g)) for g ∈ phrases
    if direction = ← then
      Let p = Max (LP(g, r, concept)) for g ∈ phrases
    e.append(p)
  return e

```

a triple (i.e., fact) in the form $\langle \text{head}, \text{relation}, \text{tail} \rangle$, a link prediction model scores the triple with a value between 0 and 1, which quantifies the level of validity of the given triple.

Positive-Unlabeled Learning Positive-Unlabeled (PU) classifiers learn a standard binary classifier in the unconventional setting where labeled negative training examples are unavailable. The PU classifier [8] learns a probabilistic function $f(X_i)$ that estimates $P(I_j \in \mathcal{S} | X_i)$ as closely as possible. We train a PU classifier using our training set (utterances with only seen intents labeled as positive) and validation set (utterances with both seen and unseen intents as unlabeled). We use 512-dimensions sentence embedding as features when using the PU classifier, generated using a pre-trained universal sentence encoder [2].

3 OUR APPROACH

Given an input utterance X_i , our model first invokes the PU classifier (if it is available) to predict whether X_i ’s intent belongs to set \mathcal{S} or \mathcal{U} . Then, relationship meta-features, utterance embedding, and intent embedding are concatenated and fed into a trainable prediction function that predicts the probability $P(I_j | X_i) \in [0, 1]$. Finally, our model outputs the intent with the highest compatibility probability, i.e., $\text{argmax}_{I_j} P(I_j | X_i)$.

Computing Relationship Meta-features. Relationship meta-features generator (RMG) extracts features by utilizing the “ActionObject” structure of intent labels and commonsense knowledge graphs. RMG takes the following inputs: a set of relations in a knowledge graph (35 in the case of ConceptNet) $\mathcal{R} = \{r_1, r_2, \dots, r_t\}$; the set of n-grams $\mathcal{G}_i = \{g_1, g_2, \dots, g_q\}$ that correspond to the input utterance X_i , where $|\mathcal{G}| = q$; and an intent label $I_j = \{\mathcal{A}, \mathcal{O}\}$, where \mathcal{A} and \mathcal{O} are the Action and Object components of the intent, respectively. RMG computes a relationship meta-features vector in four steps, where each step results in a vector of size $|\mathcal{R}|$. The smaller vectors are: $e_{X_i}^{\vec{\mathcal{A}}}$, $e_{X_i}^{\vec{\mathcal{O}}}$, $e_{X_i}^{\overleftarrow{\mathcal{A}}}$, and $e_{X_i}^{\overleftarrow{\mathcal{O}}}$, where $e_{X_i}^{\vec{\mathcal{A}}}$ captures the weights

¹RIDE: Relationship Meta-features Assisted Intent DEtection.

²<https://github.com/anonymous-sigir-researcher/GZS-IntentDetection>

³If intents are described using a complex textual description, Actions and Objects can be extracted using existing NLP tools such as dependency parsers.

Table 1: Dataset statistics.

Dataset	# of Samples	Vocab. Size	Avg. Length	# of Intents
SNIPS	14.2K	10.8K	9.05	7
SGD	57.2K	8.8K	10.62	46
MultiWOZ	30.0K	9.7K	11.07	11

of Action to utterance relationships and $\mathbf{e}_{\mathcal{X}_i}^{\vec{O}}$ captures the weights of Object to utterance relationships. The remaining two vectors capture relationship weights in the other direction; i.e., utterance to Action/Object, respectively. Capturing bi-directional relationships is important because a relationship in one direction does not necessarily imply one in the other direction. The final output of RMG is the concatenation of the four aforementioned vectors.

RMG computes $\mathbf{e}_{\mathcal{X}_i}^{\vec{\mathcal{A}}}$ by considering the strength of each relation in \mathcal{R} between \mathcal{A} and each n-gram in \mathcal{G}_i . That is, $\mathbf{e}_{\mathcal{X}_i}^{\vec{\mathcal{A}}}$ has $|\mathcal{R}|$ cells, where each cell corresponds to a relation $r \in \mathcal{R}$. Each cell is computed by taking $\max(LP(\mathcal{A}, r, g))$ over all $g \in \mathcal{G}_i$. $LP(\text{head}, \text{relation}, \text{tail})$ outputs the probability that the fact represented by the triple $\langle \text{head}, \text{relation}, \text{tail} \rangle$ exists. The vector $\mathbf{e}_{\mathcal{X}_i}^{\vec{O}}$ is computed similarly, but with passing \mathcal{O} instead of \mathcal{A} when invoking the link predictor. The vectors $\mathbf{e}_{\mathcal{X}_i}^{\overleftarrow{\mathcal{A}}}$ and $\mathbf{e}_{\mathcal{X}_i}^{\overleftarrow{O}}$ are computed similarly, but with swapping the head and tail when invoking the link predictor. Algorithm 1 outlines the previous process. Finally, the meta-features are passed through a linear layer with sigmoid activation for normalization.

Utterance and Intent Encoders. We use bi-directional LSTM to produce a d -dimensional representation of the given utterance $\mathcal{X}_i = \{w_1, w_2, \dots, w_u\}$ with u words, where contextual embeddings from a pre-trained ELMo model and parts of speech (POS) tags are employed to embed each word. The concatenation of the last hidden states is used as utterance embedding $\mathbf{e}_{\text{utterance}}$. We encode intent labels similarly to produce an intent embedding $\mathbf{e}_{\text{intent}} \in \mathbb{R}^d$.

Training. The training examples are of the form $((\mathcal{X}_i, \mathcal{I}_j), \mathcal{Y})$, where \mathcal{Y} is a binary label representing whether the utterance-intent pair $(\mathcal{X}_i, \mathcal{I}_j)$ are compatible or not. We prepare our training data by assigning a label of 1 to the available utterance-intent pairs (where intents are seen ones); these constitute positive training examples. We create a negative training example for each positive one by corrupting the example’s intent by modifying their Action, Object, or both. We train the model by minimizing the cross-entropy loss over all the training examples.

4 EXPERIMENTS

4.1 Experimental Setup

Datasets. We used three widely used intent detection benchmarks: SNIPS, SGD, and MultiWOZ; Table 1 summarizes the statistics of these datasets. SNIPS [6] is a crowd-sourced single-turn NLU benchmark with 7 intents across different domains. SGD [23] is a comprehensive and challenging dataset with 46 intents across 16 domains. MultiWOZ [37] is a well-known dataset which has utterances that span 11 intents

Evaluation Methodology. We use F1 scores to evaluate the competing methods and report the per class averages weighted by the respective class support.

Dataset splits. All models are trained on a subset of utterances implying seen intents. At inference time, test utterances are drawn from a set that contains utterances implying a mix of seen and unseen intents (disjoint set from the training set). We decided the train/test splits for each dataset as follows: For SNIPS, we first randomly selected 5 out of 7 intents and designated them as seen intents. We then selected 70% of the utterances that imply any of the 5 seen intents for training. The test set consists of the remaining 30% utterances in addition to all utterances that imply one of the 2 unseen intents. For SGD, we used the standard splits proposed by the dataset authors. Specifically, the test set includes utterances that imply 8 unseen intents and 26 seen intents. For MultiWOZ, we used 70% of the utterance that imply 8 (out of 11) randomly selected intents for training and the rest of the utterances for testing; ; we report average results over 10 runs for all the datasets.

Competing Methods. We compare our model RIDE against the following SOTA models and several strong baselines:

SEG [35]: semantic-enhanced gaussian mixture model coupled with a density-based outlier detection algorithm LOF.

ReCapsNet-ZS [19]: employs a capsule neural network (CapsNet) and a dimensional attention module to learn generalizable transformational metrics from seen intents.

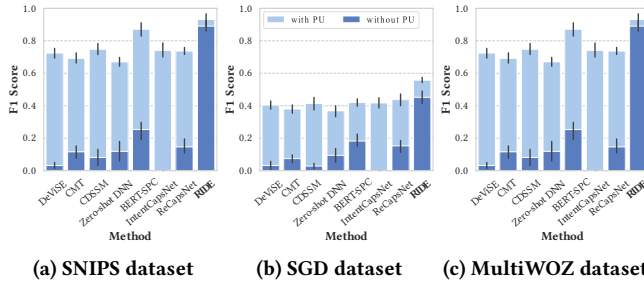
IntentCapsNet [31]: utilizes CapsNet and routing-by-agreement to adapt to unseen intents. This model was originally proposed for detecting intents in the standard ZS setting, so we extended it to support the GZS setting with the help of its authors.

Other Baseline Models. (i) Zero-shot DDN [16]: achieves zero-shot capabilities by projecting utterances and intent labels into the same high dimensional embedding space. (ii) CDSSM [4]: utilizes a convolutional deep structured semantic model to generate embeddings for unseen intents. (iii) CMT [26]: employs non-linearity in the compatibility function between utterances and intents to find the most compatible unseen intents. (iv) DeViSE [11]: was originally proposed for zero-shot image classification that learns a linear compatibility function. (v) BERT Sentence Pair Classifier (BERT-SPC): BERT is pre-trained on the sentence-pair classification task and fine-tuned on the utterances from the seen intents. *Note that baseline ZS models have been extended to support GZS setting.*

Implementation Details. We trained our link predictor on the lemmatized version of ConceptNet KG (1 million nodes, 2.7 million edges, and 35 relation types). The link predictor has two 200-dimensional embedding layers and a negative sampling ratio of 10; it is trained for 1,000 epochs using Adam optimizer with a learning rate of 0.05, L2 regularization value of 0.1, and batch size of 4800. Our relationship meta-features generator takes in an utterance’s n-grams with $n \leq 4$ and an intent label, and uses the pre-trained link predictor to produce relationship meta-features with 140 dimensions. Our utterance and intent encoders use pre-trained ELMo contextual word embeddings with 1024 dimension and POS tags embeddings with 300 dimension, and two-layer bidirectional LSTM with 300-dimensions. Our prediction function has two dense layers with relu and softmax activation. Our model is trained for up to 200 epochs with early stopping using Adam optimizer and a cross entropy loss with initial learning rate of 0.001 and ReduceLROnPlateau scheduler with 20 patience epochs. It uses a dropout rate of 0.3 and a batch size of 32. A negative sampling ratio of up to 6 is

Table 2: Main results: F1 scores for competing models.

Method	SNIPS		SGD		MultiWOZ	
	Unseen	Seen	Unseen	Seen	Unseen	Seen
DeViSE	0.0439	0.6521	0.0177	0.5451	0.0270	0.5770
CMT	0.0910	0.6639	0.0621	0.5803	0.0679	0.6216
CDSSM	0.0484	0.7028	0.0284	0.6379	0.0244	0.6515
Zero-shot DNN	0.1273	0.6687	0.1168	0.6098	0.1149	0.6012
BERT-SPC	0.2761	0.7152	0.1872	0.6401	0.1932	0.6413
IntentCapsNet	0.0000	0.6532	0.0000	0.5508	0.0000	0.6038
ReCapsNet	0.1601	0.6783	0.1331	0.5751	0.1467	0.6170
SEG	0.6991	0.8651	0.4032	0.6356	0.4143	0.6456
RIDE <i>w/o</i> PU	0.9103	0.8799	0.4634	0.8295	0.4645	0.8816
RIDE <i>/w</i> PU	0.9254	0.9080	0.5734	0.8298	0.5206	0.8847

**Figure 2: F1 scores for unseen intents for the competing models after integrating a PU classifier.**

used. We use the same embedding and training mechanism for all competing models.

4.2 Results

Main Results. Table 2 shows F1 scores averaged over 10 runs for all competing models. For both seen and unseen intents, our model RIDE outperforms all other competing models with a large margin. Specifically, RIDE achieves 32.37%, 42.21%, and 25.66% better F1 scores than the SOTA model SEG on SNIPS, SGD, and MultiWOZ for unseen intents, respectively. Moreover, our model consistently achieves the highest F1 score on seen intents, which confirms its generalizability. We highlight that RIDE outperforms the SOTA model SEG regardless of whether a PU classifier is incorporated or not. For SNIPS, the role of the PU classifier is negligible as it causes only a slight improvement in F1 score. For SGD and MultiWOZ, which are more challenging datasets, the PU classifier causes significant improvements in F1 scores. Specifically, it results in 23.74% and 12.08% improvement for SGD and MultiWOZ, respectively, on unseen intents.

Effect of PU Classifier on Other Models. We observed that one of the main sources of error for most models in the GZS setting is their tendency to misclassify utterances with unseen intents into seen ones due to overfitting to seen intents. We investigated whether existing models can be adapted to accurately classify utterances with unseen intents by partially eliminating their bias towards seen intents. Figure 2 presents F1 scores of all models with and without PU classifier. A PU classifier significantly improves the results of all the competing models. For instance, the IntentCapsNet model with a PU classifier achieves an F1 score of 74% for unseen intents on SNIPS dataset compared to an F1 score of less than 0.01% without the PU classifier. Note that the PU classifier has an accuracy (i.e., correctly predicting whether the utterance implies a seen or an unseen intent) of 93.69, 86.13, and 87.32 for SNIPS, SGD, and MultiWOZ datasets, respectively. Interestingly, our model RIDE

Table 3: Ablation study: F1 scores for unseen intents.

Configuration	SNIPS	SGD	MultiWOZ
UI-Embed <i>w/o</i> PU	0.2367	0.1578	0.1723
Rel-M <i>w/o</i> PU	0.7103	0.3593	0.3321
RIDE <i>w/o</i> PU	0.9103	0.4634	0.4645
UI-Embed <i>/w</i> PU	0.7245	0.4202	0.4124
Rel-M <i>/w</i> PU	0.8463	<u>0.5167</u>	<u>0.4781</u>
RIDE <i>/w</i> PU	0.9254	0.5734	0.5206

without PU classifier outperforms all the competing models even when a PU classifier is incorporated into them, which highlights that the PU classifier is not the component that does the heavy lifting in our model. We did not incorporate the PU classifier into SEG model because it already incorporates an equivalent mechanism to distinguish seen intents from unseen ones (i.e., outlier detection).

Ablation Study. To quantify the effectiveness of each component in our model, we present the results of our ablation study in Table 3. Utilizing utterance and intent embeddings only (i.e., UI-Embed) results in very low F1 score, i.e., 23.67% on SNIPS dataset. Employing relationship meta-features only (i.e., Rel-M) results in significantly better results: an F1 score of 71.03% on SNIPS dataset. When utterance and intent embeddings are used in conjunction with relationship meta-features (i.e., RIDE *w/o* PU), it achieves a better F1 score compared to the Rel-M or UI-Embed configurations. A similar trend can be observed for the other datasets as well. Finally, when our entire model is deployed (i.e., including utterance and intent embeddings, relationship meta-features, and the PU classifier, i.e., RIDE */w* PU), it achieves the best results on all datasets.

5 RELATED WORK

Supervised intent detection works [15, 18, 24, 34, 38] assume the availability of a large amount of labeled training data for all intents to learn discriminative features. Whereas standard zero-shot intent detection models [1, 7, 10, 12, 16, 30, 36] assume that all utterances faced at inference time imply unseen intents only. Extending such works to handle the generalized zero-shot intent detection setting (i.e., removing the aforementioned assumptions) is nontrivial. Our model is specifically designed for the generalized zero-shot intent detection setting. The authors in [19] attempted to accommodate GZS setting by adding a dimensional attention module to a capsule network that learns generalizable transformation matrices from seen intents. Recently, the authors in [35] proposed using a density-based outlier detection algorithm LOF [1] and semantic-enhanced gaussian mixture model with large margin loss to learn class-concentrated embeddings to detect unseen intents. In contrast, we leverage a rich commonsense knowledge graph to capture deep semantic and discriminative relationships between utterances and intents, which significantly reduces the bias towards classifying unseen intents into seen ones. In a related, but orthogonal, line of research, the authors in [13, 17, 21] addressed the problem of intent detection in the context of dialog state tracking where an annotated dialog state and conversation history are available in addition to an input utterance. In contrast, this work and the SOTA models we compare against in our experiments only consider an utterance without having access to any dialog state elements.

6 CONCLUSION

We have presented an accurate generalized zero-shot intent detection model. Our extensive experimental analysis on three intent detection benchmarks show that our model achieves 25.66% to 42.21% better F1 score than the SOTA model for unseen intents. The main novelty of our model is its utilization of relationship meta-features and limited reliance on training data. Furthermore, our idea of integrating Positive-Unlabeled learning in GZS intent detection models further improves our models' performance, and significantly improves the accuracy of existing models as well.

REFERENCES

- [1] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. LOF: identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*. 93–104.
- [2] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175* (2018).
- [3] Wei-Lun Chao, Soravit Changpinyo, Boqing Gong, and Fei Sha. 2016. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *European conference on computer vision*. Springer, 52–68.
- [4] Yun-Nung Chen, Dilek Hakkani-Tür, and Xiaodong He. 2016. Zero-shot learning of intent embeddings for expansion by convolutional deep structured semantic models. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6045–6049.
- [5] Zhiyuan Chen, Bing Liu, Meichun Hsu, Malu Castellanos, and Riddhiman Ghosh. 2013. Identifying intention posts in discussion forums. In *Proceedings of the 2013 conference of the North American chapter of the association for computational linguistics: human language technologies*. 1041–1050.
- [6] Alice Coucke, Alaa Saade, Adrien Ball, Théodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, et al. 2018. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *arXiv preprint arXiv:1805.10190* (2018).
- [7] Yann N Dauphin, Gokhan Tur, Dilek Hakkani-Tur, and Larry Heck. 2013. Zero-shot learning for semantic utterance classification. *arXiv preprint arXiv:1401.0509* (2013).
- [8] Charles Elkan and Keith Noto. 2008. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*. 213–220.
- [9] Rafael Felix, Vijay BG Kumar, Ian Reid, and Gustavo Carneiro. 2018. Multi-modal cycle-consistent generalized zero-shot learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 21–37.
- [10] Emmanuel Ferreira, Bassam Jabaian, and Fabrice Lefevre. 2015. Online adaptative zero-shot learning spoken language understanding using word-embedding. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 5321–5325.
- [11] Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Marc' Aurelio Ranzato, and Tomas Mikolov. 2013. Devise: A deep visual-semantic embedding model. In *Advances in neural information processing systems*. 2121–2129.
- [12] Varun Gangal, Abhinav Arora, Arash Einolghozati, and Sonal Gupta. 2019. Likelihood Ratios and Generative Classifiers for Unsupervised Out-of-Domain Detection In Task Oriented Dialog. *arXiv preprint arXiv:1912.12800* (2019).
- [13] Pavel Gulyaev, Eugenia Elistratova, Vasily Konovalov, Yuri Kuratov, Leonid Pugachev, and Mikhail Burtsev. 2020. Goal-oriented multi-task bert-based dialogue state tracker. *arXiv preprint arXiv:2002.02450* (2020).
- [14] Seyed Mehran Kazemi and David Poole. 2018. Simple embedding for link prediction in knowledge graphs. In *Advances in neural information processing systems*. 4284–4295.
- [15] Joo-Kyung Kim, Gokhan Tur, Asli Celikyilmaz, Bin Cao, and Ye-Yi Wang. 2016. Intent detection using semantically enriched word embeddings. In *2016 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 414–419.
- [16] Anjishnu Kumar, Pavankumar Reddy Muddireddy, Markus Dreyer, and Björn Hoffmeister. 2017. Zero-Shot Learning Across Heterogeneous Overlapping Domains.. In *INTERSPEECH*. 2914–2918.
- [17] Miao Li, Haoqi Xiong, and Yunbo Cao. 2020. The spdd system for schema guided dialogue state tracking challenge. *arXiv preprint arXiv:2006.09035* (2020).
- [18] Bing Liu and Ian Lane. 2016. Attention-based recurrent neural network models for joint intent detection and slot filling. *arXiv preprint arXiv:1609.01454* (2016).
- [19] Han Liu, Xiaotong Zhang, Lu Fan, Xuandi Fu, Qimai Li, Xiao-Ming Wu, and Albert YS Lam. 2019. Reconstructing Capsule Networks for Zero-shot Intent 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)*. 4801–4811.
- [20] Kun Liu, Wu Liu, Huadong Ma, Wenbing Huang, and Xiongxiang Dong. 2019. Generalized zero-shot learning for action recognition with web-scale video data. *World Wide Web* 22, 2 (2019), 807–824.
- [21] Yue Ma, Zengfeng Zeng, Dawei Zhu, Xuan Li, Yiyang Yang, Xiaoyuan Yao, Kaijie Zhou, and Jianping Shen. 2019. An end-to-end dialogue state tracking system with machine reading comprehension and wide & deep classification. *arXiv preprint arXiv:1912.09297* (2019).
- [22] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- [23] 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* (2019).
- [24] Suman Ravuri and Andreas Stolcke. 2015. Recurrent neural network and LSTM models for lexical utterance classification. In *Sixteenth Annual Conference of the International Speech Communication Association*.
- [25] Walter J Scheirer, Anderson de Rezende Rocha, Archana Sapkota, and Terrance E Boult. 2012. Toward open set recognition. *IEEE transactions on pattern analysis and machine intelligence* 35, 7 (2012), 1757–1772.
- [26] Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. 2013. Zero-shot learning through cross-modal transfer. In *Advances in neural information processing systems*. 935–943.
- [27] Robyn Speer, Joshua Chin, and Catherine Havasi. 2016. Conceptnet 5.5: An open multilingual graph of general knowledge. *arXiv preprint arXiv:1612.03975* (2016).
- [28] Nikhita Vedula, Nedim Lipka, Pranav Maneriker, and Srinivasan Parthasarathy. 2020. Open Intent Extraction from Natural Language Interactions. In *Proceedings of The Web Conference 2020*. 2009–2020.
- [29] Jinpeng Wang, Gao Cong, Wayne Xin Zhao, and Xiaoming Li. 2015. Mining User Intents in Twitter: A Semi-Supervised Approach to Inferring Intent Categories for Tweets. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (Austin, Texas) (AAAI'15)*. AAAI Press, 318–324.
- [30] Kyle Williams. 2019. Zero Shot Intent Classification Using Long-Short Term Memory Networks. *Proc. Interspeech 2019* (2019), 844–848.
- [31] Congying Xia, Chenwei Zhang, Xiaohui Yan, Yi Chang, and Philip S Yu. 2018. Zero-shot user intent detection via capsule neural networks. *arXiv preprint arXiv:1809.00385* (2018).
- [32] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. 2018. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence* 41, 9 (2018), 2251–2265.
- [33] Yongqin Xian, Bernt Schiele, and Zeynep Akata. 2017. Zero-shot learning—the good, the bad and the ugly. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4582–4591.
- [34] Puyang Xu and Ruhi Sarikaya. 2013. Convolutional neural network based triangular crf for joint intent detection and slot filling. In *2013 IEEE workshop on automatic speech recognition and understanding*. IEEE, 78–83.
- [35] Guangfeng Yan, Lu Fan, Qimai Li, Han Liu, Xiaotong Zhang, Xiao-Ming Wu, and Albert YS Lam. 2020. Unknown Intent Detection Using Gaussian Mixture Model with an Application to Zero-shot Intent Classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 1050–1060.
- [36] Majid Yazdani and James Henderson. 2015. A model of zero-shot learning of spoken language understanding. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 244–249.
- [37] Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. MultiWOZ 2.2 : A Dialogue Dataset with Additional Annotation Corrections and State Tracking Baselines. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*. Association for Computational Linguistics, Online, 109–117. <https://doi.org/10.18653/v1/2020.nlp4convai-1.13>
- [38] Chenwei Zhang, Yaliang Li, Nan Du, Wei Fan, and Philip S Yu. 2018. Joint slot filling and intent detection via capsule neural networks. *arXiv preprint arXiv:1812.09471* (2018).