

# INTENT INDUCE WITH SCCL AND OUT-DOMAIN DATA AUGMENTATION

Jun Gao

<2022-11-04 >

# Outline

Introduction

our work

dstc11

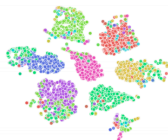
[https://drive.google.com/file/d/  
1itlby2Ypq3sRVt0Y1alr3ygjPZZdB2TT/view](https://drive.google.com/file/d/1itlby2Ypq3sRVt0Y1alr3ygjPZZdB2TT/view)

SCCL(supporting clustering with contrastive learning) is a deep clustering algorithm which can learning more adequate representation of short context. By using Contrastive Learning, same samples will be pull together while different ones will be pushing apart.

Original



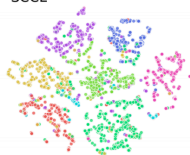
Clustering



Instance-CL

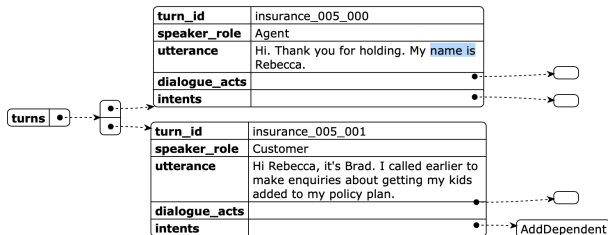


SCCL



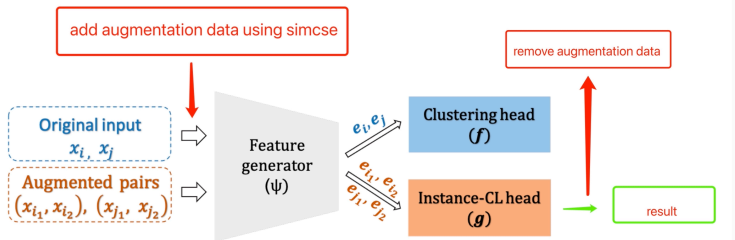
## filtering high quality utterances with statistic rule or hdbscan

In open intent induction, we will not have access to the ground-truth utterances with intents, so need to filter out noise utterances. A tradition way is to use clustering algorithms, like hdbscan, but in this task, we find a more effective way to select high quality utterances, that is statistic rules. Specifically, we count the highest frequency bi-gram phrase before the utterances in development dataset and find it domain-agnostic, thus we migrate the rule to test dataset and obtain high quality utterances. For example, if "insurance<sub>005001</sub>" is a high quality utterance, "name is" is a bi-gram phrase before it.



# data augmentation using out-domain data

After getting high quality utterances for sccl(original input in figure), besides drop-out augmentation in sccl, we can augment the original input by simcse using out-domain data(which has similar utterances in original input). Then we can train sccl normally. After getting clustering results, we have to remove the augmentation data added in previous step.



# results

